

Behaviorally Informed Machine Learning for Human Mobility

Ekin Uğurel

A dissertation proposal
submitted in partial fulfillment of the
requirements for the degree of

Doctor of Philosophy

University of Washington

2025

Reading Committee:

Cynthia Chen, Chair

Shuai Huang, Co-Chair

Qi "Ryan" Wang, Member

Program Authorized to Offer Degree:

Civil & Environmental Engineering

©Copyright 2025
Ekin Uğurel

University of Washington

Abstract

Behaviorally Informed Machine Learning for Human Mobility

Ekin Uğurel

Co-Chairs of the Supervisory Committee:

Professor Cynthia Chen

Civil & Environmental Engineering

Professor Shuai Huang

Industrial & Systems Engineering

Large-scale digital trace datasets hold considerable promise for long-range transportation planning, offering the potential to observe mobility at metropolitan scales with far greater temporal and spatial resolution than traditional household travel surveys while capturing the regularities in daily travel behavior that underlie trip-making. Passively-collected mobile (PCM) data (e.g., location signals from smartphones and in-vehicle GPS) are central to this promise. Their usefulness, however, is limited by discontinuities in individual trajectories, privacy constraints that restrict data sharing and integration, and representativeness biases that distort inferred patterns of regional travel demand and travel behavior across population groups. This dissertation addresses these limitations through four contributions. First, it develops a multi-task Gaussian Process-based imputation method (grounded in recurring daily, weekly, and seasonal travel behavior patterns) capable of handling both short- and long-duration gaps in GPS traces, significantly improving the completeness and usability of mobility data. Second, it introduces an individualized, physics-regularized learning framework that produces high-fidelity mobility traces reflective of observed movement patterns. These generated trajectories can be scaled to build richer, more diverse mobility datasets for

developing and validating activity-based models. Third, it investigates the predictive signal linking mobility patterns as expressions of travel behavior to sociodemographic attributes that shape those behaviors, and develops imputation strategies for enriching PCM datasets with these inferred labels. This enrichment supports both more detailed planning analyses and a clearer diagnosis of representativeness biases in passively collected data. Finally, through a qualitative study of long-range transportation planners, this dissertation investigates barriers to the adoption of big data products and provides recommendations for their effective integration into planning processes. Together, these contributions bridge methodological advances in machine learning with insights from travel behavior research and the practical needs of public agencies, offering a more transparent and behaviorally coherent foundation for data-driven planning and travel behavior analysis.

TABLE OF CONTENTS

	Page
List of Figures	1
List of Tables	5
Chapter 1: Introduction	8
1.1 Motivation	8
1.2 State of the Art	11
1.3 Research Objectives	15
1.4 Organization of the Dissertation	15
Chapter 2: Correcting Missingness in Passively-generated Mobile Data with Multi- Task Gaussian Processes	17
2.1 Introduction	17
2.2 Related Work	22
2.3 Methodological Framework	25
2.4 Implementation	33
2.5 Dataset	34
2.6 Experiments	36
2.7 Discussion	46
Chapter 3: Learning to Generate Synthetic Human Mobility Data: a Physics- regularized Gaussian Process approach based on Multiple Kernel Learning	51
3.1 Introduction	51
3.2 Related Work	55
3.3 Methodology	58
3.4 Numerical Experiments	71
3.5 Discussion	81

Chapter 4: On Predicting Sociodemographics from Mobility Signals	84
4.1 Introduction	84
4.2 Literature Review	88
4.3 Datasets	89
4.4 Methodology	91
4.5 Experiments	101
4.6 Conclusion	106
Chapter 5: MPO’s uses of and needs for big data	111
5.1 Introduction	111
5.2 Literature Review	112
5.3 Methods	116
5.4 Findings	117
Chapter 6: Conclusion, Discussion, and Future Work	124
6.1 Discussion	125
6.2 Future Work	127
Appendices	153
Appendix A: Supplementary Material for Chapter 2	154
A.1 Initialization Strategy & Parameter Optimization	154
A.2 Algorithm to Determine Training Data for Experiments in 2.6.2	155
A.3 Detailed Imputation Accuracy Metrics for Robustness Experiments	156
A.4 Benchmark Methods and Model Parameters	156
A.5 Max Speed Threshold Sensitivity Analysis	156
Appendix B: Supplementary Material for Chapter 3	165
B.1 List of Notation Used	165
B.2 Dataset and Software Implementation	165
B.3 Descriptive Statistics	167
B.4 Data Preprocessing and Modeling Considerations	168
B.5 Predictor Variables	170
B.6 Algorithms	171

B.7 Handling Nonstationarity with GPs	174
Appendix C: Supplementary Material for Chapter 4	176
C.1 Experimental Details and Full Model Results	176

LIST OF FIGURES

Figure Number	Page
1.1 The relationships among chapters in this dissertation	16
2.1 Overview of imputation workflow for passively-generated mobile data with varying levels of missingness (long and short gaps). Data Preprocessing reduces the computational complexity of the imputation task while also improving the accuracy of training points. Model Development specifies the kernel and mean functions and learns the kernel parameters through marginal log-likelihood maximization. Gap Imputation “corrects” the raw data by imputing the missing locations.	26
2.2 Descriptive analysis of Spectus data for all 2,000 users. (Top Left) Cumulative distribution of location accuracy. (Top Right) Distribution of observations within each day of the week. (Bottom Left) Histogram of gaps categorized by size. (Bottom Right) Boxplots of temporal occupancy with three temporal resolutions.	37
2.3 Examples of trips fit with an underestimated lengthscale (left) and appropriate lengthscale (right)	39
2.4 (Left) Boxplot of total RMSE for trips in different mobility metric clusters (Right) Boxplot of optimized lengthscales for each trip cluster.	40
2.5 Scatterplots of DTW and temporal occupancy for various gap lengths. (a) $\tau = 1$ week, (b) $\tau = 1$ day, (c) $\tau = 6$ hours, (d) $\tau = 5$ minutes.	43
2.6 Predictions (red) and ground truth test data (blue) for a user tested at $\tau = 6$ hours. MTGP with the proposed kernel (left) outperforms exponential smoothing (right).	46
3.1 (top) Proposed algorithm for greedy multiple kernel learning; (bottom) example greedy learning tree with three steps (checkmark denotes lowest BIC option)	68
3.2 Physics-regularized GP inference for passively-generated mobile data. Circles indicate variables, while diamonds indicate estimation (kernel function) or inference (kernel matrix) processes.	70
3.3 Components of the single-task temporal kernel in a 3-D space	72

3.4	(a) Loss function progression and convergence behavior of different functional forms for a composite kernel; (b) Generating one week of trips using discovered kernel and the temporal GP; (c) Derived mobility metrics by kernel. Bold denotes the optimal kernel. r_g denotes radius of gyration, B_s denotes spatial burstiness, and ΔH denotes distance between user’s imputed and actual home location.	73
3.5	Speed and Bearing Constraints in King County, Washington	75
3.6	Speed and Bearing Predictions and True Observations for a Spectus User	76
3.7	Evaluation metrics for all models being compared. Y-axis is shown in log-scale. The lower the better for all metrics.	81
4.1	Example daily mobility graph where the edges are chronologically numbered. Width of trip arrows corresponds to frequency of visits over observation period. Dashed arrows denote known trips that are not observed on this day.	92
4.2	Illustration of daily travel motifs (Schneider et al., 2013; Wu et al., 2019); (a) Out-and-back; (b) Chain; (c) Cycle-chain; (d, h) Double-cycle; (e) Single-no-return; (f, g) Single-Cycle	95
4.3	Illustration of selected metrics; (left) example travel day. In this case, $n_{\text{trips}} = 7$ and $f_{\text{comp}} = 3/7$; (right) each color denotes a tour to/from the anchors. In this case, $n_{\text{tour}} = 4$ and $f_{\text{mm}} = 2/4$	96
4.4	Shared-trunk multitask architecture. Mobility descriptors are mapped to a shared representation $h = g(X; \theta)$ by a two-layer feed-forward network with ReLU activations and dropout. Four task-specific heads f_t produce class probabilities via softmax. Training minimizes the cross-entropy loss $\sum_t w_t \cdot \ell_t$ where we set the weights to be equal.	100
4.5	Largest magnitude Spearman rank correlations (all $p < 0.001$) between mobility descriptors and demographics. (top left): age; (top right): household income; (bottom left): gender; (bottom right): number of children. Bars to the left indicate negative associations; to the right, positive.	102
4.6	Marginal changes in top-1 accuracy (higher = better), AUROC (higher = better), NLL (lower = better), ECE (lower = better) as more features are added. (top four rows) Overall split; (bottom four rows) training with the 2017/2019 data and testing on the 2023 data.	108
4.7	Reliability diagrams for representative settings. (TOP) Overall split, Age task with C+ST+D features (left) and All features (right). (BOTTOM) Cross-temporal split, Number of Children task with C (left) and C+ST+D (right). Bars show empirical accuracy within 15 equal-width confidence bins; the dashed line is the identity (perfect calibration). Grey histograms (right axes) give the number of samples per bin. Points above (below) the diagonal indicate under- (over-) confidence.	109

4.8	Performance of the MT variant (in blue) compared to ST variants (in orange) at different fractions of training data. (top four rows) Overall split; (bottom four rows) training with the 2017/2019 data and testing on the 2023 data.	110
5.1	Bar plot of answers mentioned in response to “In which area does your MPO operate?”	117
5.2	Word cloud of answers to ”What word comes to mind when you think of “big data”? (98 responses)	119
5.3	Bar plot of factors mentioned in response to “What would help you gain confidence or have trust in using big data?”	120
A.1	Mobility error metrics for $\tau = 1$ week	156
A.2	Mobility error metrics for $\tau = 1$ day	157
A.3	Mobility error metrics for $\tau = 6$ hours	157
A.4	Mobility error metrics for $\tau = 1$ hour	158
A.5	Mobility error metrics for $\tau = 30$ minutes	158
A.6	Mobility error metrics for $\tau = 15$ minutes	159
A.7	Mobility error metrics for $\tau = 5$ minutes	159
A.8	Classical error metrics for $\tau = 1$ week	160
A.9	Classical error metrics for $\tau = 1$ day	160
A.10	Classical error metrics for $\tau = 6$ hours	160
A.11	Classical error metrics for $\tau = 1$ hour	161
A.12	Classical error metrics for $\tau = 30$ minutes	161
A.13	Classical error metrics for $\tau = 15$ minutes	161
A.14	Classical error metrics for $\tau = 5$ minutes	162
A.15	The number of observations eliminated by varying speed limits across the users’ data analyzed in Section 2.6	162
B.1	Violin plots of training and testing set sizes across our passively-collected mobile data experiments	167
B.2	Descriptive analysis of Spectus data for all 2,000 users. (a) Cumulative distribution of location accuracy. (b) Distribution of observations within each day of the week.	168
B.3	Cumulative distribution of location accuracy for Spectus users in the Experiments sample (blue) and the whole dataset (grey)	169
B.4	Normalized observation counts by time of day for Spectus users in the Experiments sample (blue) and the whole dataset (orange)	169

B.5	Evolution of time interval between two consecutive points Week 1 to 27. Lines show the percentiles for the experiment sample, while the green and blue fillings shows the 70-80 and 15-35 percentile range for the entire dataset, respectively.	170
B.6	GP regression using cyclical features and a composite kernel accurately captures nonstationary behavior in CO2 concentration of the Mauna Loa volcano . . .	175
C.1	Comparison of training times between the unified multi-task learning (MT) model and four separate single-task variants (STVs) across varying data fractions (log-scaled on x-axis). Despite its larger architecture, the MT model trains faster overall because its shared trunk amortizes computation across tasks, whereas STVs require four independent forward-backward passes. Error bars show the standard deviation across cross-validation folds.	177

LIST OF TABLES

Table Number	Page
2.1 Temporal dimensions used in the robustness experiments.	34
2.2 Summary of trip clusters.	40
2.3 Median error with respect to the testing sets.	47
3.1 Optimal kernel parameters for long gap imputation example. Asterisk* denotes the highest lengthscale among categorical/binary variable.	74
3.2 Models and median test set metrics	78
4.1 Descriptive statistics of the PSRC HTS in the three waves used in this study (post-processing; values in parentheses denote the strata percentage associated with wave)	90
4.2 Selected linear regression models	103
5.1 Distribution of answers to exploratory questions about the audience	118
5.2 Types of questions MPO staffers have solved or plan on solving, categorized	122
A.1 Parameters of benchmark methods.	163
A.2 Parameters used during robustness experiments.	164
B.1 Variable dictionary	166
B.2 Temporal and physical dimensions used in our experiments. Asterisk* denotes a dimension that was not employed in the GeoLife experiments.	171
C.1 Performance across feature sets and models for Age (overall split). Values are mean±sd across folds; best per model in bold ; best in metric in red	178
C.2 Performance across feature sets and models for Gender (overall split). Values are mean±sd across folds; best per model in bold ; best in metric in red	179
C.3 Performance across feature sets and models for HH Income (overall split). Values are mean±sd across folds; best per model in bold ; best in metric in red	180
C.4 Performance across feature sets and models for Number of Children (overall split). Values are mean±sd across folds; best per model in bold ; best in metric in red	181

C.5	Performance across feature sets and models for Age (2017–2019 train, 2023 test split). Values are mean±sd across folds; best per model in bold ; best in metric in red	182
C.6	Performance across feature sets and models for Gender (2017–2019 train, 2023 test split). Values are mean±sd across folds; best per model in bold ; best in metric in red	183
C.7	Performance across feature sets and models for HH Income (2017–2019 train, 2023 test split). Values are mean±sd across folds; best per model in bold ; best in metric in red	184
C.8	Performance across feature sets and models for Number of Children (2017–2019 train, 2023 test split). Values are mean±sd across folds; best per model in bold ; best in metric in red	185

ACKNOWLEDGMENTS

I would like to express my sincere appreciation to my primary advisor, Dr. Cynthia Chen, for her mentorship throughout my time at the University of Washington. Her academic rigor and continued support were the driving forces that allowed this work to reach completion. I also thank my co-advisor and ISE mentor, Dr. Shuai Huang, and my committee members, Dr. Qi "Ryan" Wang and Dr. Yen-Chi Chen, for their invaluable guidance and the thoughtful questions that pushed me to sharpen my thinking. Additionally, the methodological ideas in Chapter 4 benefited greatly from discussions with Dr. Filipe Rodrigues, whom I also thank for hosting me as a visiting Ph.D. student at the Technical University of Denmark.

This journey would have been far more difficult without the friends and colleagues who provided a vibrant community, making the Ph.D. experience richer than I could have imagined. I am grateful to my former housemates, Michele Martino and J. Felipe Montano-Campos, for their friendship and the memorable times we shared. Thanks also to the wonderful lab mates in Wilcox Hall for their camaraderie and support.

Last, but certainly not least, I am grateful for the endless support of my family. To my parents, Necdet and Aynur Ugurel, thank you for believing in me. To my wonderful partner, Rubina Singh, and our dog, Nia—thank you for your patience and love; I surely would not have made it across the finish line without you.

This research was generously supported by the National Science Foundation on the project entitled "A Whole-Community Effort to Understand Biases and Uncertainties in Using Emerging Big Data for Mobility Analysis" (Award no.: 2114260). The author was also supported by the Center for Teaching Old Models New Tricks (TOMNET), a Tier-1 University Transportation Center sponsored by the U.S. Department of Transportation under grant 69A3551747116, as well as the Center for Understanding Future Travel Behavior and Demand (TBD), a National University Transportation Center sponsored by the US Department of Transportation under grant 69A3552344815 and 69A3552348320.

Chapter 1

INTRODUCTION

1.1 Motivation

For decades, the transportation planning industry has tried to answer some of the most difficult questions that surround the future of urban areas. These include, for example, how to set investment priorities to accommodate population growth while maintaining level of service standards, how to balance competing demands between different transportation modes in limited urban space, how to ensure equitable access to transportation services across diverse socioeconomic groups, and how to develop resilient infrastructure that can adapt to changing climate conditions and emerging technologies like autonomous vehicles. As Rittel and Webber (1973) noted five decades prior, these are rather difficult questions that tend not to have a single correct answer; they represent inherently normative dilemmas where potential solutions are judged not as true or false, but as good or bad depending on stakeholder perspective.

Traditionally, long-range transportation planners (LRTPs) have relied on household travel surveys (HTS) to illuminate the behavioral mechanisms underlying daily travel. These surveys ask residents to record detailed information about their trips, including purposes, modes, timings, and destinations. Yet despite their richness, HTS datasets have well-known limitations: they are costly to administer, achieve modest response rates, and depend on respondents accurately recalling their travel (Li et al., 2025; Richardson et al., 1996; Stopher et al., 2008). Moreover, they capture only a small sample of the population and are collected infrequently (e.g., every eight years in the U.S.; Alberini et al., 2021), and thus cannot fully represent the spatial or temporal heterogeneity of traveler behavior.

The emergence of GPS-based devices has fundamentally transformed our ability to collect mobility data at scale. Unlike the actively-collected HTS, these devices passively generate location data as a by-product of their operation and billing purposes, which planners can then repurpose to infer mobility patterns (Chen et al., 2016a, 2026). This passively-collected data offers unprecedented insights into travel behavior at fine spatial and temporal resolutions, capturing longitudinal information on (1) when people travel, (2) where they go, (3) how long they stay at those destinations, and (4) what routes they take to get there. The past 15

30 years has seen a surge of studies leveraging this rich data source to analyze collective and
31 individual mobility patterns (Alessandretti et al., 2020; Alexander et al., 2015b; Hao et al.,
32 2020; Sulis et al., 2018).

33 However, the promise of GPS-based data comes with significant methodological and
34 data-related challenges that must be carefully addressed before it can be used to solve
35 transportation planning problems. First, despite its continuous nature, GPS data often
36 suffers from both short- and long-gaps in observations that can distort the derived mobility
37 metrics (McCool et al., 2022). These gaps in data collection can occur for multiple reasons:
38 user-driven choices (such as restricting location access to when apps are actively in use; Kim et
39 al., 2019), technical limitations (including aggressive background app management by mobile
40 operating systems; Zhou et al. 2020), and environmental factors (like signal loss in urban
41 canyons or tunnels; Ben-Moshe et al., 2011). The temporal distribution of observations is
42 often uneven, with characteristic patterns showing peaks during commute hours on weekdays
43 and fewer observations during nights and weekends. This missingness is not random—it
44 systematically varies across users, time periods, and geographic contexts, creating complex
45 patterns of data sparsity. Such gaps can significantly impact our understanding of mobility
46 patterns. When observations are missing, derived metrics like stay durations, travel distances,
47 and activity patterns may be systematically underestimated. While brief gaps might have
48 minimal impact, longer periods of missingness can substantially bias our understanding of
49 movement patterns, potentially leading to incorrect conclusions about travel behavior and
50 infrastructure needs. This sparsity issue is the motivation behind the imputation method
51 proposed in Chapter 2.

52 Second, the fine spatiotemporal nature of location traces raises serious privacy concerns.
53 Data aggregators and providers increasingly face strict privacy regulations and user-centric
54 demands for data protection, leading to more restrictive data sharing agreements and easier
55 opt-out options for users. This shift is driven by heightened consumer awareness of data
56 privacy implications and industry-led initiatives to enhance user control over personal data
57 sharing (Apple, 2021). The resulting decrease in available data has profound implications
58 for transportation planning applications, particularly in developing synthetic populations for
59 large-scale mobility simulations. Traditional approaches, such as survey-estimated activity-
60 based models (ABMs) or probabilistic models derived from aggregated mobile data, struggle
61 to reproduce the fine-grained dynamics needed for realistic simulation (Bhat and Koppelman,
62 1999; Pendyala et al., 1997). Their reliance on coarse aggregates or sparse survey samples leaves

63 them unable to represent how personal characteristics, situational factors, and environmental
64 context jointly shape movement. These limitations highlight the need for synthetic data
65 methods that both preserve privacy and faithfully reconstruct the underlying data-generating
66 process at the individual level. This is the motivation for Chapter 3, which prescribes a
67 method for generating synthetic data that preserves statistical properties while protecting
68 individual privacy.

69 Finally, even when privacy-protected GPS data can be collected and processed, the
70 resulting mobility metrics often exhibit systematic biases that must be addressed (Li et al.,
71 2024; Wang et al., 2025). These biases stem from multiple sources: the data may over-
72 represent certain demographic groups while under-sampling others or disproportionately reflect
73 particular travel patterns based on when and where data collection occurs. Unfortunately,
74 the extent of such bias is typically unknown because demographic information is not directly
75 available in passively collected datasets. This missing dimension fundamentally limits the kinds
76 of analyses transportation planners seek to conduct. Metropolitan planning organizations
77 (MPOs) increasingly rely on mobile data to evaluate network investments or assess the
78 distributional impacts of new projects. Many of these applications require linking mobility
79 traces to travelers' demographic attributes in order to gauge, for example, whether tolling
80 policies differentially affect low-income households or whether expanded transit lines serve
81 underserved communities. Without demographic context, agencies risk making decisions
82 based on travel patterns that reflect only the behaviors of overrepresented subpopulations.
83 Consequently, to quantify and ultimately correct for these biases, it becomes essential to infer
84 travelers' underlying demographic attributes from their observed mobility behavior alone. To
85 this end, in Chapter 4, we introduce a standardized set of higher-order mobility descriptors
86 to strengthen demographic signal extraction, calibrate predictive models to characterize
87 uncertainty, and apply multitask learning to improve model generalization across contexts.

88 Besides the domain context, the common denominator across the technical chapters 2-4
89 is a shared methodological backbone in behavioral integration. Each technical contribution
90 is guided by empirical regularities and theoretical insights long established in the travel
91 behavior literature. For example, daily and weekly mobility patterns motivate the structured
92 kernels and periodic priors used for imputing missing data. Principles of activity scheduling
93 and space-time constraints inform the individualized, physics-guided mechanisms that shape
94 our synthetic data generation framework. Likewise, decades of work documenting how
95 sociodemographic attributes manifest in activity patterns motivate the higher-order mobility

96 descriptors and uncertainty-aware predictive models developed for demographic inference.
97 Across all chapters, behavioral theory serves primarily as a structuring device: it constrains
98 model classes to reflect known regularities in human movement, guiding the choice of priors and
99 feature representations. These behavioral constraints help ensure that the resulting models
100 remain grounded in plausible mobility patterns rather than purely data-driven artifacts.

101 Building models around behavioral principles addresses only part of the challenge; equally
102 important is understanding how such methods are received and evaluated within the insti-
103 tutions that rely on them. Despite the growing availability of passively generated mobile
104 data, many planning agencies remain cautious about integrating these products into core
105 decision-making processes. Their hesitancy is not solely due to the technical limitations
106 discussed above (i.e., sparsity, privacy constraints, and demographic bias/uncertainty). It
107 also reflects deeper questions about how these datasets fit within existing institutional, pro-
108 fessional, and regulatory structures. Planners must judge whether new data sources are
109 trustworthy, whether they align with established norms of practice, and how they should be
110 used in the absence of clear federal or state guidance. These questions become even more
111 acute for smaller Metropolitan Planning Organizations (MPOs), which often lack the staff
112 capacity or analytical resources to independently assess complex data products. Chapter 5
113 examines these institutional and organizational dimensions through an in-person workshop at
114 the Association of Metropolitan Planning Organizations (AMPO) Planning Tools & Training
115 Symposium in Albuquerque, New Mexico, in May 2024, using this setting to understand
116 what planners need in order to responsibly and confidently adopt big data products.

117 **1.2 State of the Art**

118 This dissertation develops behaviorally informed machine learning approaches to address
119 issues in passively-collected GPS traces. The goal is to make these data usable and reliable
120 for transportation planning applications. Specifically, with respect to these datasets, this
121 dissertation aims to (1) fix sparsity (both short and long gaps); (2) overcome privacy
122 concerns by generating synthetic data; (3) characterize representativeness bias by imputing
123 sociodemographics from mobility signals; and (4) understand how MPOs use big data, what
124 they need from it, and what barriers shape its adoption. Correspondingly, I will review
125 several existing research areas that intersect with these challenges.

126 1.2.1 *Correcting Missingness*

127 Missing data imputation for mobile datasets generally involves three strategies: (1) time-
 128 series smoothing and interpolation, (2) leveraging external datasets, and (3) machine learning
 129 methods, including kernel- and deep learning-based approaches. Traditional time-series
 130 models (e.g., SARIMAX) offer interpretability but struggle with the non-linear, irregular
 131 nature of mobility data (Huo et al., 2010; Kohn and Ansley, 1986; McCool et al., 2022).
 132 External datasets, including transit smart card records or mobile data from other providers
 133 (Gong et al., 2020), offer a promising way to bolster sparse mobile traces. Yet they also pose
 134 substantial obstacles: access to these data can be limited, privacy rules may restrict sharing,
 135 and cross-source similarity matching tends to be non-trivial (Cao et al., 2016; Kondor et al.,
 136 2020). Machine learning approaches like RNNs or GNNs effectively capture spatiotemporal
 137 dependencies but often operate as black-boxes, with limited interpretability (Ren et al., 2021;
 138 Sun et al., 2021). In Chapter 2, we proposed a novel Gaussian Process (GP)-based framework
 139 to impute both short- and long-gaps, capturing correlations across user coordinates via multi-
 140 task learning and leveraging a distance-based compression algorithm to manage scalability.
 141 We find that our Bayesian framework outperforms a host of competing approaches based on
 142 how closely we recover the underlying mobility statistics of partially-observed trajectories.

143 1.2.2 *Generating Synthetic Data*

144 While imputation addresses sparsity, many planning applications still require privacy-
 145 preserving representations of mobility that can be shared, simulated, and analyzed without
 146 exposing individual traces. This has motivated a growing literature on generating synthetic
 147 trajectories that reproduce heterogeneous travel patterns. Early work relied on structured
 148 stochastic models that separate temporal routines from spatial movement, learning diary-like
 149 Markov chains over activity states and then assigning locations via preferential exploration
 150 and return rules (Pappalardo and Simini, 2018). More recent studies rely on deep genera-
 151 tive models. Recurrent neural networks (RNNs) have been trained to simulate individual
 152 trajectories and, in extensions, full synthetic populations over multiple days (Berke et al.,
 153 2022; Kulkarni and Garbinato, 2017). Generative adversarial networks (GANs) augment
 154 these sequence models with discriminators and mobility-regularization terms to better match
 155 empirical distributions and support downstream tasks such as epidemic simulation (Feng
 156 et al., 2020). Diffusion models push fidelity further by learning denoisers that map noise to
 157 realistic GPS trajectories (Zhu et al., 2023). Finally, transformer architectures have also been

158 adapted to treat activities, locations, and modes as tokens (conditioned on sociodemographics)
 159 to generate rich mobility profiles (Zhang et al., 2024b).

160 Despite this progress, two gaps remain for applications of methods datasets in trans-
 161 portation contexts. First, most existing generators are global models that do not explicitly
 162 tailor their representation to each traveler’s idiosyncratic temporal and spatial patterns, even
 163 though heterogeneity in routines is central to travel behavior (Kroesen, 2014; Lee et al.,
 164 2007; McGuckin and Murakami, 1999a). Second, these models seldom embed simple physical
 165 or built-environment constraints and rarely quantify uncertainty, which limits their use for
 166 risk-sensitive planning applications. Chapter 3 addresses these gaps by recasting synthetic
 167 trajectory generation as a Bayesian learning problem. We learn individualized Gaussian
 168 Process kernels that adapt to each person’s observed mobility, and we incorporate physical
 169 and temporal attributes of the built environment when shaping these kernels. Sampling from
 170 the resulting posterior yields synthetic traces that reflect traveler-specific structure, respect
 171 basic physical constraints, and come with principled uncertainty quantification.

172 1.2.3 *Inferring Demographics from Mobility Behavior*

173 Complementing efforts to generate privacy-preserving synthetic mobility data, a growing
 174 body of research has explored inferring sociodemographic attributes from digital behavioral
 175 traces, motivated by the need to link mobility data with population characteristics for
 176 planning and policy analysis. Within this literature, mobility-focused approaches rely on
 177 *where and when* people travel, transforming raw coordinates into spatial, temporal, and
 178 semantic descriptors that proxy for demographic traits. Spatial metrics (e.g., radius of
 179 gyration, location diversity), temporal regularities (e.g., commuting rhythms), and semantic
 180 signals (e.g., POI visit frequencies, tour structures) have each been shown to carry predictive
 181 power for attributes such as age, gender, household structure, and work status (Auld et al.,
 182 2015; Ding et al., 2019; Doi et al., 2021; Wu et al., 2019). However, these associations are
 183 often context-specific and ad hoc, and few studies systematically evaluate the robustness or
 184 incremental value of each feature family across settings.

185 Methodologically, the field has evolved from classical statistical models using hand-crafted
 186 features (as in Auld et al., 2015) to representation learning techniques that encode daily
 187 trajectories as structured sequences or graphs (as in Solomon et al., 2018). Neural embedding
 188 methods and recurrent models have improved predictive performance by capturing temporal
 189 and contextual regularities (Ding et al., 2019; Xu et al., 2020), but most prior work trains

190 separate models for each demographic attribute, limiting data efficiency and transferability
191 under distributional shifts. Moreover, uncertainty is often handled superficially: while most
192 studies report aggregate accuracy, few quantify how confident the model is in its predictions
193 for individual travelers. Recent advances in calibration and theoretical separability provide
194 tools to address this gap, enabling more uncertainty-aware inference. Building on these
195 developments, Chapter 4 evaluates the predictive contribution of different mobility feature
196 families, leverages multitask learning to exploit shared structure across demographic targets,
197 and incorporates calibration methods to produce more reliable, interpretable uncertainty
198 estimates.

199 *1.2.4 Uses and Adoption of Big Data by MPOs*

200 MPOs increasingly view big data as a way to strengthen core planning tools, but most
201 documented uses remain targeted and cautious. Agencies commonly deploy probe-based
202 speed and volume data to calibrate and validate regional travel demand models (Bauer et al.,
203 2014; FHWA, 2020). Some MPOs have used vendor products to estimate AADT or construct
204 OD matrices for corridor studies, subarea analyses, or to cross-check regional commuting
205 patterns, often treating these datasets as a “second source” rather than a primary input
206 (KimleyHorn, 2021; StreetLight, 2025). Planners also see potential to use smartphone traces
207 to study ride-hailing and to evaluate accessibility across communities (Boyd et al., 2024).
208 In practice, however, many of these aspirations remain exploratory and are constrained by
209 concerns about representativeness bias (Richardson, 2021; Wang et al., 2025).

210 Despite the proliferation of big-data products, adoption within MPOs is limited by
211 intertwined technical, organizational, and institutional barriers. Planners frequently cite
212 opaque vendor methods, uncertain data quality, and shifting proprietary algorithms as
213 major obstacles to trusting and integrating these datasets into long-range planning workflows
214 (KimleyHorn, 2021; Singh et al., 2022). High and often opaque cost structures make it difficult
215 for smaller agencies to justify long-term subscriptions (Bauer et al., 2014; KimleyHorn, 2021).
216 Many MPOs also lack in-house data engineering or advanced analytics capacity, facing
217 difficulties in recruiting and retaining staff with big-data expertise (Pecheux et al., 2020;
218 Richardson, 2021). These constraints interact with organizational culture and leadership;
219 skepticism among executives about the value of big data, siloed data practices, and low
220 psychological safety around experimentation all reduce agencies’ absorptive capacity and
221 confine the use of big data to pilot projects (Pecheux et al., 2020). Chapter 5 engages directly

222 with these issues by examining how planners articulate the uses they see for big data, how
223 they evaluate vendor products, and which institutional, technical, and organizational factors
224 most strongly shape adoption.

225 **1.3 Research Objectives**

226 The objectives of this dissertation are:

- 227 1. To develop and validate a multi-task Gaussian process-based imputation method for
228 short- and long-gaps in passively-collected GPS traces;
- 229 2. To create an individualized, physics-regularized framework for generating high-fidelity
230 synthetic traces that replicate observed patterns;
- 231 3. To develop and evaluate methods for inferring travelers' sociodemographic attributes
232 solely from mobility behavior, using higher-order mobility descriptors, uncertainty
233 calibration, and multitask learning to improve predictive accuracy and generalization
234 across contexts, and;
- 235 4. To investigate and document the barriers to adoption of big data products among
236 transportation planning organizations, providing insights into how these tools can be
237 more effectively integrated into planning processes.

238 **1.4 Organization of the Dissertation**

239 This dissertation is presented in a multiple manuscript format. Chapters 2-5 are written
240 as individual research papers, including an abstract, a main body, and references. The
241 relationships among these chapters are depicted in Figure 1.1. Chapter 2, proposes a Bayesian
242 ML method to correct (i.e., impute) missing data in passively-generated GPS traces, taking
243 advantage of multiple periodicities that exist in humans' travel trends. Chapter 3 extends this
244 work with individualized kernel learning and physics-regularization, facilitating the generation
245 of synthetic traces that replicate users' observed travel patterns. This allows the research
246 community to avoid compromising individual privacy while working with passively-generated
247 GPS traces and share their synthetic counterparts. Chapter 4 investigates the extent to which
248 key sociodemographic attributes (such as age, gender, income, and household structure) can
249 be inferred solely from individuals' mobility patterns. The results highlight the potential

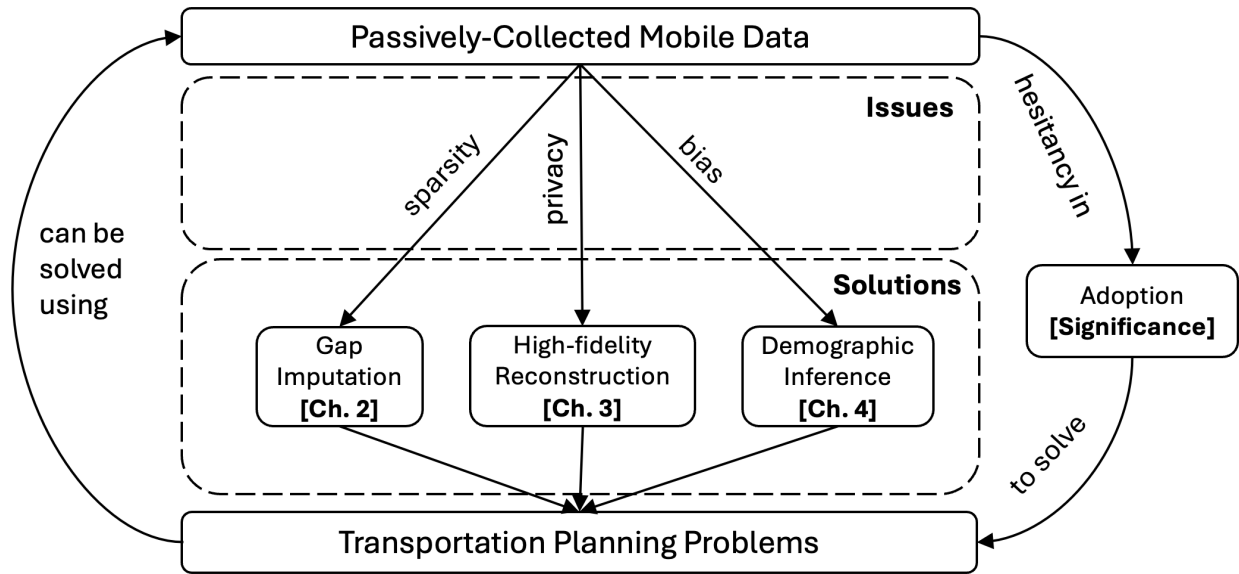


Figure 1.1: The relationships among chapters in this dissertation

250 to enrich passively collected mobility datasets with inferred sociodemographic information,
 251 enabling more comprehensive behavioral analyses without relying on direct survey linkage.
 252 Finally, Chapter 5 presents a qualitative study of long-range transportation planners' uses of
 253 and needs for big data products, tackling the limited adoption of these datasets among the
 254 planning community.

Chapter 2

**CORRECTING MISSINGNESS IN PASSIVELY-GENERATED
MOBILE DATA WITH MULTI-TASK GAUSSIAN PROCESSES**

The prevalence of mobile devices and the ubiquity of network connectivity have generated a massive amount of temporally- and spatially-stamped data. A key characteristic of this mobile data is the prevalence of sparsity—for every recorded user, there is a significant amount of missing data. Sparsity leads to bias in the inferred mobility patterns and thus, correcting bias by imputing the missing data potentially creates opportunities for directly using the corrected mobile trajectory data for large-scale simulations in various applications. We propose a multi-task Gaussian process regression model to correct missingness in mobile data. Gaussian processes (GPs) allow for flexible modeling of diverse (i.e., non-linear, locally periodic) data patterns and quantify prediction uncertainty in an interpretable manner. We develop a methodological framework for applying GPs to mobile data. In doing so, we consider the correlations between users’ coordinates (latitudes and longitudes) through multi-task learning and adjust for individual-level differences in data characteristics through parameter initialization and optimization. We introduce and demonstrate the use of smooth (i.e., rational quadratic) and periodic kernels in modeling human mobility data. Relatedly, we analyze our model’s parameters and imputation accuracy with respect to different types of trips (e.g. walking trips vs. highway trips). We also demonstrate our model’s performance in three experiments with real app-based data, in which it outperforms alternative imputation methods.

2.1 Introduction

The past two decades have brought a range of technological advancements that have made it easy to gather large sets of time- and location-stamped mobile data. Coupled with the ubiquity of network connectivity, this has led to a massive increase in individual trace sightings. Those datasets include tracks generated by Global Positioning Systems (GPS) or other position-signaling devices, geo-tagged posts from social media platforms, and call detail records (CDRs) derived from triangulation of cellular tower positions. Human mobility

283 patterns inferred from such datasets have been used for many applications, including for
284 example, to quantify urban vitality (Sulis et al., 2018), understand commuting patterns
285 (Frias-Martinez et al., 2012), and predict pandemic spreading (Hao et al., 2020).

286 Unlike household travel surveys, which design and implement a set of pre-defined questions
287 to answer (thus called “actively-solicited data”), mobile data is passively-generated as a
288 byproduct of processes that have little to do with answering research questions (Chen et al.,
289 2016b). As a result, while mobile data offers sample sizes orders of magnitudes greater than
290 traditional travel surveys, there exist many issues relating to its representativeness. Chief
291 among these issues is sparsity—for a substantial portion of users, there is a significant amount
292 of missing data. In other words, sparsity occurs when the observation frequency of a device
293 is low.

294 The observation frequency of mobile data varies greatly across users, periods, and ge-
295 ographies for various reasons. Users commonly have the choice of only allowing location
296 transmissions while apps are in use (rather than continuously in the background). Perceived
297 benefits in network externalities (e.g., increased accuracy for Google Maps travel time estima-
298 tion) have been linked to greater willingness to provide personal information, such as location
299 data (Kim et al., 2019). Therefore, third-party apps that do not provide any perceived benefits
300 to the data-solicited user may receive limited access to location data. Certain versions of the
301 Android operating system (OS) have also been found to aggressively shut down apps running
302 in the background, preventing continuous data collection (Zhou et al., 2020). In addition to
303 stochastic events like battery drain, user- and OS-based decisions are likely to lead to long
304 observation gaps for users, ranging anywhere from a few hours to multiple days or weeks.

305 Additionally, app-based datasets have large variations in their intra- and inter-day spar-
306 sity—observations tend to cluster temporally, resulting in frequency ‘peaks’ and ‘valleys’.
307 Ban et al. (2018) offer two insights regarding the temporal distribution of observations in
308 app-based datasets: First, during weekdays, observations peak in the morning (7-9 AM)
309 and in the evening (4-6 PM), while during weekends there is only one mid-day peak (12-3
310 PM). There are also significantly fewer observations overnight than during the day. Second,
311 weekdays have more observations than weekends—Fridays tend to have the most observations,
312 while Sundays tend to have the fewest, with the exception to the rule being holidays. This
313 clustering tendency can be caused by users’ mobile device usage patterns or the update
314 frequency of the data provider. The problem is further shrouded by the opaqueness of those
315 providers in disclosing neither the source of the data nor the reason for its generation.

316 Finally, missingness may also be caused by a range of geographical factors. Shorter
317 gaps in continuity are often the result of temporary signal loss, which may occur in dense
318 urban areas, while traveling through tunnels and other enclosed infrastructure, or due to
319 the “cold-start problem”—when signals are dropped due to a lack of clear line of sight from
320 GPS satellites. Additionally, the “urban canyon effect” describes the tendency of receivers in
321 Global Navigation Satellite Systems (GNSS) to output erroneous location estimations while
322 operating in areas with a high density of tall buildings, which block the satellite’s direct line
323 of sight to the receiving device (Ben-Moshe et al., 2011).

324 When a mobile dataset has missing data, derived mobility patterns are vulnerable to
325 bias, meaning that they may misrepresent individuals’ actual movements. Several outputs
326 may be altered due to missing data, such as the duration of the inferred stays (Ban et al.,
327 2018), the extent of one’s travel (McCool et al., 2022), and users’ interactions with the
328 built environment (e.g., the types of locations visited; Merrill et al., 2020). While short
329 and infrequent observation gaps have minimal influence on derived mobility metrics, longer
330 and more frequent gaps will result in a downward bias (McCool et al., 2022). On the other
331 hand, improvements in data quality, such as decreases in rates of missingness or increases
332 in location data accuracy (i.e., the closeness of observation to the real location), have been
333 linked to more accurate calculations of related mobility metrics, such as home census tracts
334 and trip rates (Ban et al., 2018). Correcting missingness in mobile data is thus important to
335 researchers and practitioners alike.

336 Before prescribing a method, we briefly outline some of the challenges that are present
337 while modeling individual mobile data. In general, people tend to use various modes of
338 transportation to get around. Different modes have different underlying physics (i.e., average
339 velocity, acceleration behavior), meaning that they possess varying data-generating mecha-
340 nisms. Additionally, people tend to change modes, either within the same trip or between
341 different trips. These realities suggest that any method to correct missingness needs to be
342 flexible enough to capture different data-generating mechanisms and the transition between
343 these states.

344 A further complexity involves individuals’ heterogeneous travel behavior. Previous works
345 have explored the predictability of human mobility, proposing various models and outlining
346 their properties (Barbosa et al., 2018; Gonzalez et al., 2008; Song et al., 2010a). In general,
347 however, no single model has outperformed every other model in every context—in contrast,
348 each model has had success in different domains. Thus, a key takeaway is that there may not

349 exist universal “laws” of human mobility—different people have different tendencies to explore
350 and exploit their environments and therefore the distributions of their data are different.
351 Oftentimes, the data distribution depends heavily on context and time (Hills et al., 2015).
352 For example, an individual who has recently moved to a new city may be more interested
353 in trying out a range of restaurants, since they have had no prior experience to guide their
354 preferences. On the other hand, a 15-year resident of a neighborhood may have only one or
355 two restaurants they regularly visit.

356 These individual-level complexities suggest that making rigid assumptions about the
357 shape or form of the underlying distribution of an individual’s mobile data may not be a
358 successful approach. Rather, a non-parametric data-driven approach may be a promising
359 one—such models are less restrictive, can accommodate complex patterns, and can better
360 capture underlying relationships in the data, which may involve both temporal and spatial
361 correlations. One example of a family of non-parametric methods includes Gaussian processes
362 (GP). GPs can be understood as a generalization of multivariate Gaussian distributions to
363 infinitely many variables. Two components define a GP—the mean function specifies the
364 mean at any point in the input space, and the covariance function (or kernel) embeds a
365 measure of similarity between any pair of points in a multi-dimensional space. In practice,
366 specifying a kernel that can accurately capture relationships in the data is a rather difficult
367 task. Simultaneously, a fitting kernel can seamlessly capture complex dependencies in mobile
368 data.

369 In the context of mobile phone trajectory data, which often exhibits highly non-linear and
370 context-dependent human mobility patterns, Gaussian Processes (GPs) offer an appropriate
371 framework for capturing these complexities. For instance, various travel modes exhibit distinct
372 underlying physics—cars move faster than bikes, and whereas a walker can move freely in any
373 direction in a park, a car is limited to paved roads. GPs can manage these differences with
374 kernel functions, facilitating context-dependent modeling of mobile data points. Moreover,
375 GPs provide not only point estimates but also uncertainty estimates, which are invaluable for
376 downstream analysis. The precision of the mobility patterns derived from trajectory data can
377 be significantly influenced by these uncertainty measures. Additionally, there is an increasing
378 necessity to integrate diverse data types, such as categorical, continuous, or even textual
379 data, and GPs are well-equipped to handle this complexity effectively.

380 In this study, we propose a generalizable multi-task GP framework to correct missingness
381 in mobile data at the individual level. Our framework imputes a user’s latitudes and longitudes

382 simultaneously based on a set of predictor variables (primarily time), hence the word multi-
 383 task (the modeling of latitudes and longitudes as two tasks). We train the framework’s
 384 parameters on longitudinal trajectory data. We use the word “generalizable” to refer to two of
 385 our framework’s properties: (1) its ability to accurately correct varying levels of missingness
 386 (i.e. short and long gaps), (2) its propensity to capture changing data-generating mechanisms
 387 (i.e. due to mode changes). Both properties are achieved through the specification of an
 388 appropriate covariance function and the ensuing parameter optimization process. In this way,
 389 the model is kept general enough to apply to any user and unique enough to be accurate.

390 No model is perfect and GPs are no exception. One notable trade-off of using GPs is its
 391 demand on computational resources: though GPs have become more scalable, they can still
 392 be computationally demanding for extensive datasets. Our implementation addresses this
 393 challenge by leveraging GPyTorch (Gardner et al., 2018), an efficient GP implementation
 394 that reduces the asymptotic complexity of exact GP inference from $\mathcal{O}(n^3)$ to $\mathcal{O}(n^2)$. Another
 395 trade-off involves model complexity. Setting up GPs requires meticulous attention to kernel
 396 functions and hyperparameters, contributing to their perceived complexity. Our study
 397 demonstrates that kernel choices are directly related to the underlying behavioral patterns
 398 and the estimated hyperparameters exhibit systematic patterns that can be explained by
 399 human mobility behaviors (Sections 2.3.3 and 2.6.1). Additionally, compared to parametric
 400 regression methods, GPs are often seen as “black-box” models compared to parametric
 401 models, which might limit the interpretability of the results for some users. However, GPs
 402 are considered less of a black-box in comparison to deep neural networks, as choices of basic
 403 kernels are based on the nature of the phenomenon of interest (in this context, human mobility
 404 patterns) and their hyperparameters such as lengthscale are explainable (they directly tell us
 405 the periodicity of the underlying function, which characterizes the periodicity of individual
 406 mobility patterns). Additionally, in our context which is imputing missing data in the raw
 407 mobility trajectories, interpretability is not a primary concern. This is because the outputs we
 408 seek to obtain are repaired trajectories with the missingness filled; those repaired trajectories
 409 still need further processing to obtain meaningful mobility metrics such as the number of
 410 stays.

411 We summarize our contributions below:

- 412 • We develop a methodological framework to correct missingness in heterogeneous mobile
 413 data using a multi-task GP model. The resulting (repaired) data can be used for
 414 many transportation and mobility applications downstream. Our model considers the

415 correlations between users’ coordinates and adjusts for individual-level differences in
 416 data characteristics;

- 417 • We analyze our model’s key hyperparameters (e.g., lengthscale parameters) and impu-
 418 tation accuracy concerning different types of trips (e.g. walking trips vs. highway trips).
 419 Specifically, we show the linkage between trip types and lengthscale parameters;
- 420 • We introduce and demonstrate the combined use of smooth kernels (i.e., the rational
 421 quadratic kernel) and kernels with periodicities (i.e., the periodic kernel) in modeling
 422 human mobility data; and
- 423 • We benchmark our model against other missing data imputation methods, demonstrating
 424 its effectiveness under a range of missingness scenarios (short and long gaps).

425 The rest of this chapter is organized as follows. Section 2.2 summarizes related work on
 426 missing mobile data imputation. Section 2.3 outlines our methodology. Section 2.4 presents
 427 our implementation while Section 2.5 describes the dataset we used in our experiments.
 428 Section 2.6 shows our results. We discuss the implications of our work in Section 2.7.

429 **2.2 Related Work**

430 *2.2.1 Methods for imputing missingness in the data*

431 In general, missing data imputation methods for mobile data fall into one of three categories:
 432 (1) using time-series smoothing and interpolation methods, (2) leveraging external data
 433 sources to facilitate co-learning or (3) employing machine learning methods such as kernels
 434 or deep learning architectures. The first approach uses existing methods for smoothing and
 435 time-series interpolation to impute missing data points. These range from simpler methods
 436 like linear interpolation to more complex ones like seasonal auto-regressive integrated moving
 437 averages (SARIMAX). These models have a long history of being used to impute missing
 438 data in a range of contexts (i.e., in hydrology, finance, etc.) and therefore are more mature
 439 (Cipra et al., 1995; Huo et al., 2010; Kohn and Ansley, 1986; McCool et al., 2022). Though
 440 these methods are relatively interpretable, their interpretability depends on the validity of
 441 their assumptions, which tend to relate to the regularity of the data (i.e., seasonality or
 442 trend), linear dependence over time, and independent noises. Relatedly, the drawback of these

443 models is they are often not flexible enough to handle the non-linear nature of individual-level
444 trajectory data.

445 The second approach leverages external data sources to make inferences about the original
446 dataset. This strategy is especially common with CDRs, which provide limited spatial and
447 temporal precision and thus struggle to capture fine-grained movement patterns. External
448 datasets such as transit smart card records or mobile data from other providers can help
449 enrich CDR or GPS trace datasets when the overlap between sources is strong. Gong et al.
450 (2020) developed two indicators to measure the similarity of spatiotemporal trajectories across
451 multiple data streams, and related work has explored the *matchability* of individuals and
452 trajectories across sources using probabilistic techniques (Cao et al., 2016; Kondor et al., 2020).
453 Despite this promise, multi-source fusion faces several practical obstacles: external datasets
454 may not align spatially or temporally with the primary data, may be costly to obtain, or
455 may be restricted by privacy regulations. Even when available, integrating and preprocessing
456 multiple data sources can be resource-intensive, and robust cross-source matching remains
457 technically challenging.

458 Finally, kernel- and deep architecture-based learning methods have also gained traction in
459 the context of missing data. These methods often boil down to one of two approaches. The
460 first involves using kernels or activation functions to predict the similarity between two spatial
461 trajectories, then using the more complete trajectory to fill in gaps in the sparser trajectory.
462 For example, Wang et al. (2020) utilize the head-direction information of trajectories together
463 with the displacement attributed to an attention mechanism to learn from past trajectory
464 points with different priorities. Similarly, Liu and Onnela (2021) bidirectionally sample
465 discrete displacements for missing segments based on similar trajectories in linear time
466 complexity. The drawback of this approach is that even if two trips have the same ground-
467 truth trajectory, varying data sampling rates (due to the data provider) may distort the
468 shape of the trajectory observed in the data. This could lead one to conclude that these
469 two trips are not similar. Unless access to consistently high-sampling rate training data is
470 available, using trajectory shapes as a measure of similarity is liable to produce errors.

471 The third approach involves using recurrent neural nets (RNN), attention mechanisms,
472 or other related methods to capture long- and short-term temporal dependencies in users'
473 locations. This is the approach we take in this paper, but rather than using time- and
474 resource-intensive deep neural networks, we opt for multi-task GPs. We review two examples
475 of the deep neural network approach. Sun et al. (2021) propose a model that captures

476 complex location transition patterns with graph neural networks and uses two attention
 477 mechanisms to capture the multi-level and shifting periodicity of human mobility. However,
 478 this model is less adept at fine-grained trajectory recovery, which involves accurately tracking
 479 and predicting very specific, detailed movements at a high temporal resolution. On the other
 480 hand, Ren et al. (2021) propose a framework that can recover and map match the fine-grained
 481 points in trajectories from coarse-grained GPS data. It uses a multi-task sequence-to-sequence
 482 learning model to capture the spatial and temporal dependencies of trajectories. However,
 483 their model requires knowledge of the underlying street network, which may not be available
 484 in a lot of scenarios (relating to the above paragraph on external data sources).

485 *2.2.2 GPs in transportation and mobility applications*

486 The use of Gaussian processes in transportation and mobility applications is not new. As
 487 early as 2009 (a few years after Rasmussen & Williams (2006) released their pivotal book),
 488 GPs were being proposed to predict travel times for arbitrary origin-destination pairs (Idé
 489 and Kato, 2009). Since then, GPs have been leveraged in a variety of domains, including in
 490 traffic operations and forecasting (Xie et al., 2010), transportation system estimation (Liu
 491 et al., 2022c), and pedestrian behavior modeling (Nasernejad et al., 2021).

492 Within the field of traffic operations and forecasting, Gaussian Processes (GPs) have
 493 emerged as a powerful tool for modeling and predicting traffic dynamics. Yuan et al. (2021)
 494 proposed a physics-regularized GP framework, pushing the boundaries of macroscopic traffic
 495 flow modeling. Beyond macroscopic flows, GPs have been used to predict fine-grained traffic
 496 speeds, as demonstrated by Le et al. (2016). Similar to the method of our paper, Rodrigues
 497 et al. (2018) leveraged multi-task GPs to impute missing traffic speeds while considering
 498 the spatial dependencies with nearby road segments. They found that using GPs to capture
 499 spatial correlations with nearby road segments led to substantial improvements in imputation
 500 performance over the benchmark methods.

501 Concurrently, GPs have demonstrated efficacy in handling uncertainties within trans-
 502 portation systems. Storm et al. (2022), for example, introduced an efficient method for
 503 evaluating stochastic traffic flow models using GP-based approximations. This not only
 504 enhanced computational efficiency but also refined the accuracy of stochastic traffic models.
 505 Expanding on uncertainty management, Steentoft et al. (2024) utilized GPs to provide
 506 uncertainty estimates for mobility flows derived from large-scale taxi data. Their approach
 507 enhanced the reliability of mobility flow predictions in addition to addressing the critical

508 issue of variable selection in complex transportation datasets.

509 GPs have also found their use in shared mobility and city-scale travel demand modeling.
510 In addressing the challenges of micro-mobility planning, (Gammelli et al., 2020) focused on
511 improving demand forecasting accuracy with censored data, which inherently contains a biased
512 representation of the true demand due to supply constraints. They innovatively employed
513 Gaussian Processes (GPs) to replace the traditional linear functional form in specifying the
514 likelihood of data being censored. This approach enables the model to avoid bias due to
515 censoring and improve the accuracy of demand predictions. On the other hand, Batista et
516 al. (2022) employed GPs for estimating city-scale OD matrices, focusing on supply-related
517 characteristics of urban networks. Their methodology provides an efficient alternative to
518 Monte Carlo simulations, offering computational efficiency through iterative OD pair selection
519 and shortest trip determination. Significantly, their method creates smaller synthetic sets,
520 substantially reducing computational demands.

521 In our study, we also leverage Gaussian Processes, extending their application to a different
522 facet of travel demand modeling—specifically, addressing missing data in mobile trajectory
523 datasets. There is a thematic resonance between our work and that of Batista et al. (2022),
524 particularly in the use of GPs to handle complex urban mobility data. Both studies underscore
525 the flexibility of GPs in urban mobility contexts, whether it be in handling synthesizing
526 trip sets or managing missing trajectory data. Uniquely, our method introduces multi-task
527 learning to exploit correlations in user coordinates while also incorporating a distance-based
528 compression algorithm to reduce the size of our training data, acknowledging the bounded
529 nature of individual mobility patterns as suggested by González et al. (2008). Furthermore,
530 our model utilizes a unique kernel designed to align with the characteristics of our dataset as
531 well as findings from the human mobility literature, which suggest varying levels of regularity
532 in individual travel behavior patterns (Kitamura and Van Der Hoorn, 1987; Song et al., 2010b;
533 Teixeira et al., 2021). This tailored kernel selection echoes the importance of kernel choice in
534 GP models, as also highlighted by Batista et al. (2022) in their work.

535 **2.3 Methodological Framework**

536 Our imputation workflow has three modules: Data Preprocessing, Model Development, and
537 Gap Imputation, described in detail in Figure 2.1. The first module filters out erroneous and
538 noisy data points, compresses a user’s trajectory using a pairwise distance-based algorithm, and
539 normalizes a user’s coordinates. Model Development identifies appropriate kernel functions and

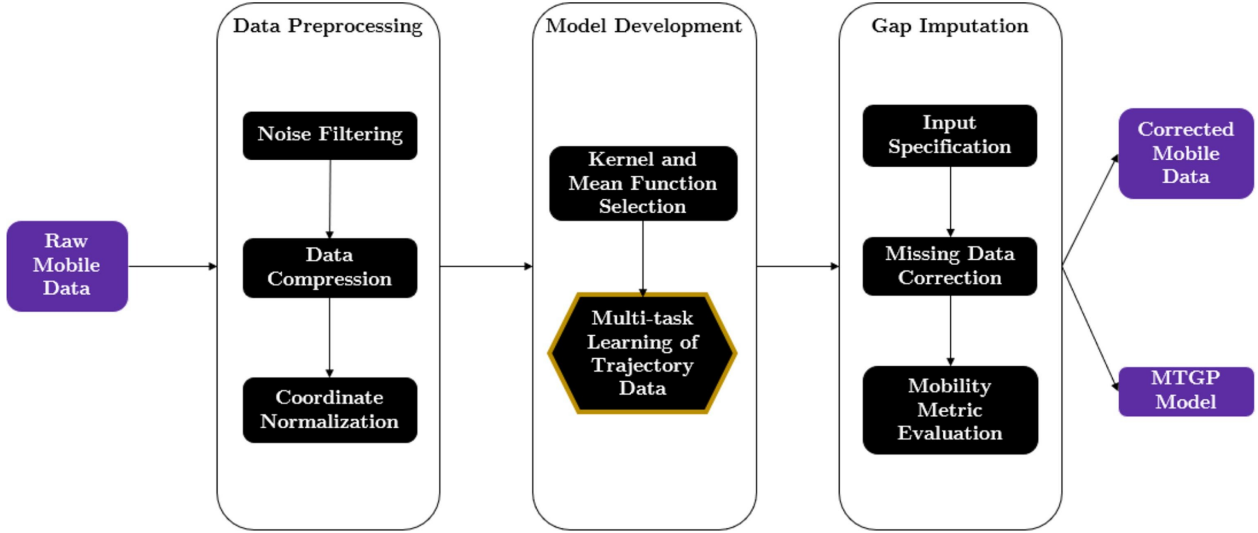


Figure 2.1: Overview of imputation workflow for passively-generated mobile data with varying levels of missingness (long and short gaps). Data Preprocessing reduces the computational complexity of the imputation task while also improving the accuracy of training points. Model Development specifies the kernel and mean functions and learns the kernel parameters through marginal log-likelihood maximization. Gap Imputation “corrects” the raw data by imputing the missing locations.

540 learns related parameters based on the marginal log-likelihood (the loss function) concerning
 541 the training data. Given a set of missing time inputs, Gap Imputation predicts the most
 542 likely location for those inputs.

543 2.3.1 Spatiotemporal Modeling of Human Mobility with Multi-Task GP Learning

544 We adopt a multi-task GP framework to capture the correlations between two highly correlated
 545 tasks in mobile data: predicting latitudes (ϕ) and longitudes (λ). Given a set \mathbf{X} of n inputs

546 $\mathbf{x}_1, \dots, \mathbf{x}_n$, we define $\mathbf{Y} = \begin{bmatrix} \mathbf{y}_\phi \\ \mathbf{y}_\lambda \end{bmatrix}^T = \begin{bmatrix} y_{1\phi}, \dots, y_{n\phi} \\ y_{1\lambda}, \dots, y_{n\lambda} \end{bmatrix}^T$, where y_{ij} is the response for the
 547 j^{th} task at i^{th} observation, and j refers to either latitudes ϕ or longitudes λ .

548 We assume the underlying data generation process for y_{ij} as follows:

$$y_{ij} = f_j(\mathbf{x}_i) + \epsilon_{ij} \quad (2.1)$$

549 where f_j is the latent function mapping inputs \mathbf{x}_i to outputs y_{ij} and ϵ_{ij} is a white noise

550 process associated with the j^{th} task and $\epsilon_{ij} \sim \mathcal{N}(0, \delta_j^2)$ are independent random variables.
 551 y_{ij} is assumed to be normally distributed, or $y_{ij} \sim \mathcal{N}(f_j(\mathbf{x}_i), \delta_j^2)$.

552 In multi-task Gaussian process learning, a multivariate normal distribution is used as
 553 the prior for modeling multi-outputs (in our case, there are two: latitudes and longitudes).
 554 In this setting, we can consider both the temporal correlation within a task but also the
 555 correlation between tasks. Thus, we denote the covariance matrix for all n observations and
 556 all m tasks as \mathbf{K} , which can be expressed as follows:

$$\mathbf{K} = \mathbf{K}^f(\mathbf{y}_\phi, \mathbf{y}_\lambda) \otimes \mathbf{K}^x(\mathbf{X}, \mathbf{X}) \quad (2.2)$$

557 where \otimes is the Kronecker product, \mathbf{K}^x is the covariance matrix of the training data for
 558 temporal correlations (n by n), \mathbf{K}^f is the inter-task covariance matrix (Bonilla et al., 2007). In
 559 our context where $m = 2$, the resulting covariance matrix \mathbf{K} is $2n$ by $2n$. The hyperparameters
 560 associated with \mathbf{K}^x and \mathbf{K}^f can be estimated by minimizing the negative marginal log-
 561 likelihood of the dataset (see Section 2.3.4).

562 In inference, for a new input \mathbf{x}_* , the predictive distribution of the output \mathbf{y}_* is also
 563 Gaussian and can be computed as¹

$$p(\mathbf{y}_* | \mathbf{x}_*, \mathbf{X}, \mathbf{Y}) \sim \mathcal{N}(\boldsymbol{\mu}_*, \boldsymbol{\sigma}_*^2) \quad (2.3)$$

564 where

$$\boldsymbol{\mu}_* = (\mathbf{k}_j^f \otimes \mathbf{k}_*) \left(\mathbf{K}^f \otimes \mathbf{K}^x + \mathbf{D} \otimes \mathbf{I} \right)^{-1} \text{vec}(\mathbf{Y}), \quad (2.4)$$

$$\boldsymbol{\sigma}_*^2 = \left(\mathbf{k}_j^f \otimes \mathbf{k}_{**} \right) - \left(\mathbf{k}_j^f \otimes \mathbf{k}_* \right) \left(\mathbf{K}^f \otimes \mathbf{K}^x + \mathbf{D} \otimes \mathbf{I} \right)^{-1} \left(\mathbf{k}_j^f \otimes \mathbf{k}_* \right). \quad (2.5)$$

565 Here, \mathbf{k}_j^f selects the j^{th} column of \mathbf{K}^f , $\mathbf{k}_* = \mathbf{k}(\mathbf{x}_*, \mathbf{X})$ is the covariance vector between the
 566 test point and the training inputs, $\mathbf{k}_{**} = \mathbf{k}(\mathbf{x}_*, \mathbf{x}_*)$ is the auto-covariance of the test point, \mathbf{D} is
 567 a 2×2 diagonal matrix with the variances of the noise processes for latitude and longitude,
 568 and \mathbf{I} is the identity matrix.

¹For more background on the derivation of these equations, refer to Section 2.2 on Rasmussen and Williams (2006)

569 *2.3.2 Using Kernels to Model Human Mobility*

570 Human mobility observations at nearby time inputs tend to be spatially close to one another.
 571 This is an expanded application of Tobler’s first law of geography—that near things are more
 572 related than distant things. To model this phenomenon, we use the rational quadratic kernel:

$$k_{RQ}(x, x') = \Omega^2 \left(1 + \frac{(x - x')^2}{2\alpha l^2} \right)^{-\alpha} \quad (2.6)$$

573 where x and x' are two inputs, Ω is a scale parameter, determining the average distance of the
 574 function away from its mean, $\alpha > 0$ is the scale mixture parameter, determining the relative
 575 weighting of large-scale and small-scale variations in the data, and $l > 0$ is the lengthscale,
 576 specifying the smoothness of the function (i.e., the frequency of the gradient of the function to
 577 change sign). The lengthscale is commonly interpreted as the maximum distance the model
 578 can extrapolate beyond the training data without reverting to the mean. Note that as x and
 579 x' approach each other, the term with the exponent $-\alpha$ goes to 1, signifying the increasing
 580 likelihood of nearby points to covary.

581 This works fine for smooth trajectories with relatively low levels of sparsity. For prediction
 582 regions in which we have little-to-no observations, data from other days, weeks, and months
 583 for the same individual can be leveraged to improve model performance. This is possible
 584 because human mobility data exhibit patterns of periodicity.

585 Previous literature suggests that people’s mobility patterns exhibit history dependency,
 586 meaning that where one is at time t depends upon their location at $t - 1$ (Kitamura and
 587 Kermanshah, 1983). More specifically, activity locations center around frequently- and
 588 recently-visited locations such as home and work (Barbosa et al., 2018; Song et al., 2010a;
 589 Teixeira et al., 2021). There are rhythms of activity and travel patterns—people usually
 590 go to their workplaces and return home every day, inducing correlations between the days
 591 of the week. Other behavior like grocery shopping, going out with friends, or attending
 592 family gatherings happens less frequently, and hence has longer intervals between consecutive
 593 instances. These activities tend to be correlated between the weeks of the month and
 594 sometimes even the months of the year (i.e., Christmas, Thanksgiving, etc.).

595 To reflect these patterns, we incorporate categorical variables like days and weeks by
 596 representing them as sets of binary variables, using a one-of-k encoding. For example, as
 597 the days of the week x_d can take one of seven values, $x_d \in \{M, Tu, W, Th, F, Sa, Su\}$,

598 then a one-of-k encoding of x_d will correspond to seven binary inputs and one-of-k(Tu) =
 599 $[0, 1, 0, 0, 0, 0, 0]$. One approach to embed multiple inputs in a GP framework is to
 600 multiply kernels defined on each individual input. This family of kernels is called Automatic
 601 Relevance Determination (or ARD), so named due to the existence of a different lengthscale
 602 parameter for each input dimension d (Duvenaud, 2014). In this case, the kernel function
 603 \mathbf{K}_{RQ} in Equation 2.6 becomes (for notational simplicity we drop the subscript i in \mathbf{x}_i):

$$k_{RQ-ARD}(\mathbf{x}, \mathbf{x}') = \prod_{s=1}^S \Omega_s^2 \left(1 + \frac{(x_s - x'_s)^2}{2\alpha l_s^2} \right)^{-\alpha} \quad (2.7)$$

604 where l_s is the lengthscale parameter for input dimension d . Thus, if the optimal lengthscale
 605 for a given categorical variable (i.e., day of the week) is small, the model has determined
 606 relatively low correlation between data in that category. On the other hand, categorical
 607 input dimensions with large lengthscales imply relatively little variation (in other words, high
 608 correlation) along those dimensions in the output variable (i.e., locations).

609 In addition to a one-of-k encoding strategy, we also use periodic kernels, which allow GPs
 610 to model functions that repeat themselves:

$$k_{PER}(x, x') = \Omega^2 \exp\left(-\frac{2 \sin^2(\pi |x - x'|/p)}{l^2}\right) \quad (2.8)$$

611 where the period length p determines the distance between repetitions of the function and the
 612 lengthscale l works in the same way as in the RQ kernel. The incorporation of the periodic
 613 kernel captures the various rhythmic patterns revealed by human mobility (people tend to
 614 conduct different activities and trips at varying frequencies). We note that the periodic kernel
 615 works best for continuous input dimensions that are unbounded in the input space.

616 2.3.3 Kernel and Mean Function Selection

617 Before prescribing a covariance function, we give some intuition on algebraic operations with
 618 kernels. The set of positive semi-definite kernels is closed under sum and product operations².
 619 Multiplying two kernels can be interpreted as an AND operation while adding two kernels can
 620 be interpreted as an OR operation (Duvenaud, 2014). Let k_1 and k_2 be kernels which each

²For details on why this is the case, we refer the reader to Section 4.2.4 of Rasmussen and Williams (2006).

depend on a single input vector, \mathbf{x} and \mathbf{y} , respectively. The product of k_1 and k_2 will result in a prior over functions that vary across both \mathbf{x} and \mathbf{y} and hence the function value $f(x_i, y_i)$ is only expected to be similar to some other function value $f(x_g, y_g)$ if x_i is close to x_g AND y_i is close to y_g . On the other hand, the sum $k_1 + k_2$ will result in a prior over functions which are a sum of one-dimensional functions, and hence the function $f(\mathbf{x}, \mathbf{y}) = f_x(\mathbf{x}) + f_y(\mathbf{y})$.

We design a composite kernel that is the sum of two nonlinear product kernels. We multiply the RQ-ARD kernel (introduced earlier) with a periodic kernel and then add a second identical component:

$$\mathbf{K}(\mathbf{x}, \mathbf{x}') = \eta_1(\mathbf{K}_{RQ-ARD} \times \mathbf{K}_{PER}) + \eta_2(\mathbf{K}_{RQ-ARD} \times \mathbf{K}_{PER}) \quad (2.9)$$

where $\eta = \Omega^2$ denotes the weight of a kernel component.

As described by Duvenaud (2014), the product of a smooth kernel (like the RQ kernel) and the periodic kernel results in functions that are periodic but can slowly vary over the input space—that is, the shape of the repeating part of the function can change. This is an appropriate kernel to use as human mobility (as expressed in the way one’s latitude and longitude changes over time) tends to be periodic (i.e., home-work-home routine), but with slight variations over time (Kitamura and Van Der Hoorn, 1987; Song et al., 2010b; Teixeira et al., 2021).

The two identical product terms in Equation 2.9 are to capture periodic patterns with different scales: the first periodic kernel is initialized with $p = 24$ hours (aiming to capture the home-to-work-to-home mobility pattern), while the second is initialized with $p = 7$ days (aiming to capture the weekly ebbs and flows of mobility). η_1 and η_2 are weights associated with the two product terms, respectively. When $\eta_1 = \eta_2$, the model places equal emphasis on the daily and weekly components. When $\eta_1 \gg \eta_2$, the model identifies that the daily pattern is much more pronounced one than the weekly one and vice versa for $\eta_1 \ll \eta_2$. We constrain the weights such that their sum is always equal to one. We initialize the two components with different period lengths. In estimation, these two weights are treated as hyperparameters, like the lengthscale of a kernel component. Therefore, they are estimated in the same optimization procedure as described in Appendix A.1.

Furthermore, we use a constant mean function, for which the baseline value is set as the median value of the training set. The median latitude and longitude tend to correspond to (or are close to) a person’s home location, as one’s activity pattern tends to evolve around home. Thus, that becomes our model’s baseline prediction in the absence of any information captured

652 in the covariance function. Compared to other mean functions (e.g., linear, multivariate
 653 orthogonal polynomial), the constant mean is a safe choice for data-driven models—in intra-
 654 trip inference, a majority of the prediction task is related to points temporally close to
 655 the training set and hence the covariance function explains most of the variation in the
 656 predictions. Our experiments with other mean functions showed that they tend to produce
 657 irrational predictions once the distance between the testing and training sets exceeds a certain
 658 threshold.

659 2.3.4 Model Training

660 In training our model, the goal is to optimally determine the set of hyperparameters Θ that
 661 minimize error with respect to the training data. The most straightforward method is to do
 662 so by minimizing the negative marginal log-likelihood (Rasmussen and Williams, 2006).

$$-\log p(\mathbf{Y} \mid \mathbf{X}, \Theta) = -\frac{1}{2} \text{vec}(\mathbf{Y})^T \boldsymbol{\Sigma}^{-1} \text{vec}(\mathbf{Y}) - \frac{1}{2} \log |\boldsymbol{\Sigma}| - n \log(2\pi), \quad (2.10)$$

663 where Θ is the set of model parameters and $\boldsymbol{\Sigma} = \mathbf{K}^f \otimes \mathbf{K}^x + \mathbf{D} \otimes \mathbf{I}$ and $|\boldsymbol{\Sigma}|$ is the
 664 determinant of the covariance matrix. Equation 2.10 provides a target function for kernel
 665 learning. The first component here estimates the model fit, while the second and third terms
 666 act as the regularization (Rasmussen and Ghahramani, 2000). Practically, the inverse term
 667 in the first component of this equation may pose a computational barrier when the size of
 668 the training data (and hence the kernel matrix) becomes too large. Traditional exact GP
 669 inference scales in $\mathcal{O}(n^3)$, while newer implementations take advantage of batched linear
 670 conjugate gradients to reduce this to $\mathcal{O}(n^2)$ (Gardner et al., 2018).

671 Various approximation methods have been proposed to overcome the computational
 672 limits of GPs on large datasets. Sparse Gaussian Processes (SGPs), for example, ease the
 673 computational burden by employing a subset of the training data (“inducing points”) to
 674 reduce the dimensionality of the covariance matrix (Titsias, 2009). The key to the success
 675 of these methods is the careful selection of inducing points, which must be informative
 676 enough to explain most of the variation in the training data. Popular SGP methods use
 677 gradient-based optimization algorithms to choose optimal inducing points—these approaches
 678 run into problems when considering a large inducing point set size and/or high dimensional
 679 input spaces (Snelson and Ghahramani, 2005).

680 For passively-generated mobile data, most of the variation is generated from users’ taking
 681 trips that differ from each other in a number of dimensions. Specifically, destination and

Algorithm 1: Trajectory Data Compression

Input: Raw trajectory data of one or multiple users; spatial radius parameter r
Output: Compressed trajectory data
 LenTraj \leftarrow Count the total number of observations in the raw trajectory;
 Sorting: Sort by ‘User ID’ and ‘Datetime’;
 Initialization: Create lists $[\phi]$ and $[\lambda]$, initializing them with ϕ_i and λ_i . Initialize $i = 0$
 and $j = 1$;
for $i < LenTraj$ **do**
 for $j < LenTraj + 1$ **do**
 $d_{ij} \leftarrow$ Measure Haversine distance between (ϕ_i, λ_i) and (ϕ_j, λ_j) ;
 if $d_{ij} > r$ **then**
 | break // End current for loop;
 end
 Add ϕ_j and λ_j to lists $[\phi]$ and $[\lambda]$;
 $j \leftarrow j + 1$;
 end
 $(\phi, \lambda) \leftarrow$ np.median($[\phi]$), np.median($[\lambda]$) // Replace each point in $[\phi]$ and $[\lambda]$ with
 the median point of each list;
 $i \leftarrow i + j$;
 $j \leftarrow i + 1$ // Update indices so that the compressed points are skipped in the next
 iteration;
end
return *The compressed trajectory data;*

682 mode choice, as well as departure time, are the three primary dimensions in which individuals’
 683 travel patterns vary. A wide range of literature has proposed parametric models to estimate
 684 these variables (Abkowitz, 1981; Daisy et al., 2018; Kitamura, 1988). In our model, we relate
 685 the variance of these variables to changes in a high-dimensional time matrix through the
 686 covariance function. However, we also do not need each available data point while training
 687 the model—many observations are close to each other in space, particularly if the sampling
 688 rate for a given period is high.

689 We reduce the size of trajectory data through a batched compression algorithm that aims
 690 to generate an approximated trajectory largely retaining the shape of the original trajectory.
 691 This algorithm is a variation of the well-known Douglas-Peucker algorithm (Douglas and
 692 Peucker, 1973), though instead of specifying an error requirement, we specify a Haversine
 693 distance—if multiple points are within a certain distance of one another, we replace them

694 using the median of those points. We chose this variation as it was faster to implement
695 and used a more interpretable parameter based on distance. Algorithm 1 describes the full
696 algorithm, which achieves the same outcome as a gradient-based optimization algorithm that
697 carefully chooses training points which capture most of the variability. This reduces the
698 computational burden of our model.

699 **2.4 Implementation**

700 We implement our algorithm and data analysis in Python. Specifically, we leverage GPyTorch,
701 an efficient and modular implementation of GPs, as well as scikit-mobility, a data processing
702 framework for GPS traces in the context of human mobility (Gardner et al., 2018; Pappalardo
703 et al., 2019). The entire library of functions and classes we develop can be found at
704 <https://github.com/ekinugurel/GPSImpute>. The requirements.txt file in the repository lists
705 the required packages and their versions to run our program.

706 *2.4.1 Data Preprocessing: Oscillation Correction, Noise Filtering, and Coordinate Normal-* 707 *ization*

708 We begin by preprocessing raw GPS traces to remove noisy data points—specifically, we filter
709 by maximum velocity, using 200 km/h as the upper limit. Because segment velocities are
710 calculated “as the crow flies” in our analysis (and hence expected to be smaller than the true
711 value), a more conservative limit was appropriate (we provide more thoughts on this choice in
712 Appendix A.5). This erases many of the oscillations or physically infeasible jumps that may
713 be observed in mobile data due to the urban canyon effect. We then remove noisy points
714 from the raw data by excluding observations more than 300 meters in precision (this is the
715 location radius for which the data provider has 95% confidence). In practice, the potential of
716 these noisy points to provide previously unobserved location information is shadowed by the
717 problems they cause for model calibration, disrupting the continuity of smooth trajectories.
718 Finally, prior to training the model, we normalize each user’s training coordinates such that
719 they have a mean of 0 and a variance of 1.

720 *2.4.2 Gap Imputation: Input Specification*

721 As described under Section 2.3.2, we represent time in multiple ways to encode the complexities
722 of human mobility patterns (e.g., history dependency, periodicity) into the multi-task GP

Table 2.1: Temporal dimensions used in the robustness experiments.

Variable	Notation	Type	Model Inputs	Code Notation
Unix Time (normalized)	t_u	Continuous	$[0, 1, \dots, n]$	<code>unix_min</code>
Second of the Day Sine	t_{ss}	Continuous	$[0, \dots, 1]$	<code>sam_sin</code>
Second of the Day Cosine	t_{sc}	Continuous	$[0, \dots, 1]$	<code>sam_cos</code>
Day of Week	t_d	Categorical	$[0, 1, 2, 3, 4, 5, 6]$	<code>dow</code>
Public Holiday	t_h	Binary	$[0, 1]$	<code>holiday</code>
Weekend	t_{we}	Binary	$[0, 1]$	<code>weekend</code>
AM Peak	t_{am}	Binary	$[0, 1]$	<code>am_peak</code>
PM Peak	t_{pm}	Binary	$[0, 1]$	<code>pm_peak</code>

723 model. We use a combination of continuous, categorical, and binary variables. The continuous
724 variables include a monotonically-increasing Unix time, normalized to be 0 at users’ first data
725 point, as well as the number of seconds elapsed after midnight in a given day. Furthermore,
726 we represent the days of the week using a one-of-k encoding as previously described, and
727 they make up our categorical variables. Finally, we use binary variables to denote whether
728 the given day is a public holiday, whether it is a weekend, and whether it is included in the
729 AM or PM peaks (defined as 7-10am and 3-6pm, respectively). Table 2.1 summarizes the
730 temporal dimensions used in our experiments.

731 With these in mind, we briefly remark on the previously-introduced kernel function (in
732 Equation 2.9). The RQ-ARD kernels are specified to fit every input dimension (i.e. produce
733 a unique lengthscale for each input vector), while the PER kernels fit only a monotonically-
734 increasing Unix time variable t_u , the only continuous and unbounded input in our input
735 space. Also note that the PER kernels are multiplied by the RQ-ARD kernels—this can be
736 loosely thought of as an AND operation, rather than an OR operation (which corresponds to
737 addition). Therefore, the periodic kernels are only meant to supplement the hidden structures
738 captured through the specification of additional temporal variables, not replace them.

739 2.5 Dataset

740 We employ privacy-protected, passively-generated mobile data from Spectus, a U.S.-based
741 data solution provider specializing in geospatial analytics. The dataset contains time-stamped
742 location traces of 2,000 anonymous, opted-in individuals in the Greater Seattle Area between
743 December 2019 and July 2020. More specifically, the data includes timestamps, unique device

744 identifiers, latitude and longitude coordinates, and a measure of data precision (i.e., a location
 745 radius for which we have 95% confidence). The dataset does not include any demographic or
 746 socioeconomic information. Locational data is sent to Spectus servers in encrypted form, via
 747 an HTTPS protocol through three options: an Application Protocol Interface (API), through
 748 a publisher that has licensed Spectus’ software development kit, or via direct server-to-server
 749 integration (Spectus, 2022).

750 Spectus enhances users’ privacy through two methods. First, they remove data from
 751 locations that do not meet privacy standards determined by their Sensitive Points of Interest
 752 (SPOI) policy. Though not an exhaustive list, the following is an example of data Spectus
 753 does not use for attribution: health-related, locations with vulnerable populations, sensitive
 754 lawful businesses, military-related, locations with first responders, correctional facilities,
 755 locations with firearms, churches/religious facilities, Native American reservations, sexual
 756 orientation-related, adult-oriented entertainment, and social demonstrations (Spectus, 2022).
 757 Second, Spectus anonymizes data in around home locations using patent-pending technology.
 758 What this means is that traces near identified home locations are replaced by points at
 759 the centroid of the address’ census block group (CBG), thereby not revealing one’s actual
 760 residence. CBGs are the second smallest geographical unit for which the U.S. Census Bureau
 761 publishes sample data, and they contain a nationwide average of 51 blocks.

762 2.5.1 Defining Missingness: Temporal Occupancy

Because Spectus’ datasets are passively-generated, the gaps between any two adjacent
 observations are rarely equal in length. Thus, we need a convention to mathematically denote
 varying levels of missingness in mobile data. We discretize a user’s total available data time
 \mathcal{T} into P intervals of length τ , which we refer to as the “temporal resolution.” The choice
 of τ is important—it decides the sparseness of a user’s observed trajectory, in which each
 interval is assigned an indicator variable

$$I_p = \begin{cases} 1 & \text{if } p \text{ has at least one observation} \\ 0 & \text{otherwise} \end{cases}$$

763 We thus define temporal occupancy (or the inverse of sparsity) as

$$q_\tau = \frac{1}{P} \sum_{p=1}^P I_p \tag{2.11}$$

764 Note that $0 \leq q_\tau \leq 1$. A period with high q has fewer gaps in data and vice versa. For
 765 example, a period with $q_{30} = 0.9$ suggests that there exists at least one observation in 90% of
 766 the 30-minute intervals.

767 2.5.2 Descriptive Analysis: Quantifying Missingness

768 We provide some brief statistics on variables in the Spectus data, shown on Figure 2.2. Most
 769 data points are highly precise—95 percent of all traces are within a 65-meter radius of their
 770 true location (top left). This is a significant change from a 2018 study looking at a similar
 771 dataset, in which about 7 percent of all observations had a precision less than 1,000 meters
 772 (Ban et al., 2018). In terms of the distribution of observations throughout the day (top right),
 773 our findings are consistent with that of Ban et al. (2018). We observe two peaks on weekdays,
 774 one in the morning (around 7am) and one in the evening (around 6pm). For weekends, there
 775 is one mid-day peak. For every day, there is a positive trend on the number of observations
 776 as the hours progress.

777 As expected, a mobile dataset’s level of “missingness” depends on the temporal resolution
 778 one chooses. The bottom right figure in Figure 2.2 shows a boxplot of individual temporal
 779 occupancies of 2,000 users from the Spectus dataset for three different temporal resolutions.
 780 As the temporal interval τ increases, so does the average temporal occupancy. The bottom
 781 left figure in Figure 2.2 shows that the gaps are not always short either—roughly 40 percent
 782 of all users recorded in Spectus data have at least one week that is continuously missing
 783 (i.e., there are no observations for a whole week). Meanwhile, 96 percent of all users have at
 784 least one day of missingness, and 99 percent have at least one six-hour period of missingness.
 785 Thus, we anticipate the need to prescribe a method to correct missingness in the general
 786 sense, whether the gaps are short or long.

787 2.6 Experiments

788 We do two sets of experiments to showcase our model. The first set explores the parameter
 789 space of multi-task GPs in the context of human mobility. Our goal here is to understand
 790 what model parameters work well for mobile data with certain characteristics (i.e., highway
 791 drives vs. urban walks), a benefit of which is developing a greater intuition on how to initialize
 792 GP model parameters in different contexts. The second set of experiments benchmarks our
 793 model against alternative missing data imputation methods, including simple exponential
 794 smoothing, exponential smoothing (Huo et al., 2010), Holt-Winters (Cipra et al., 1995),

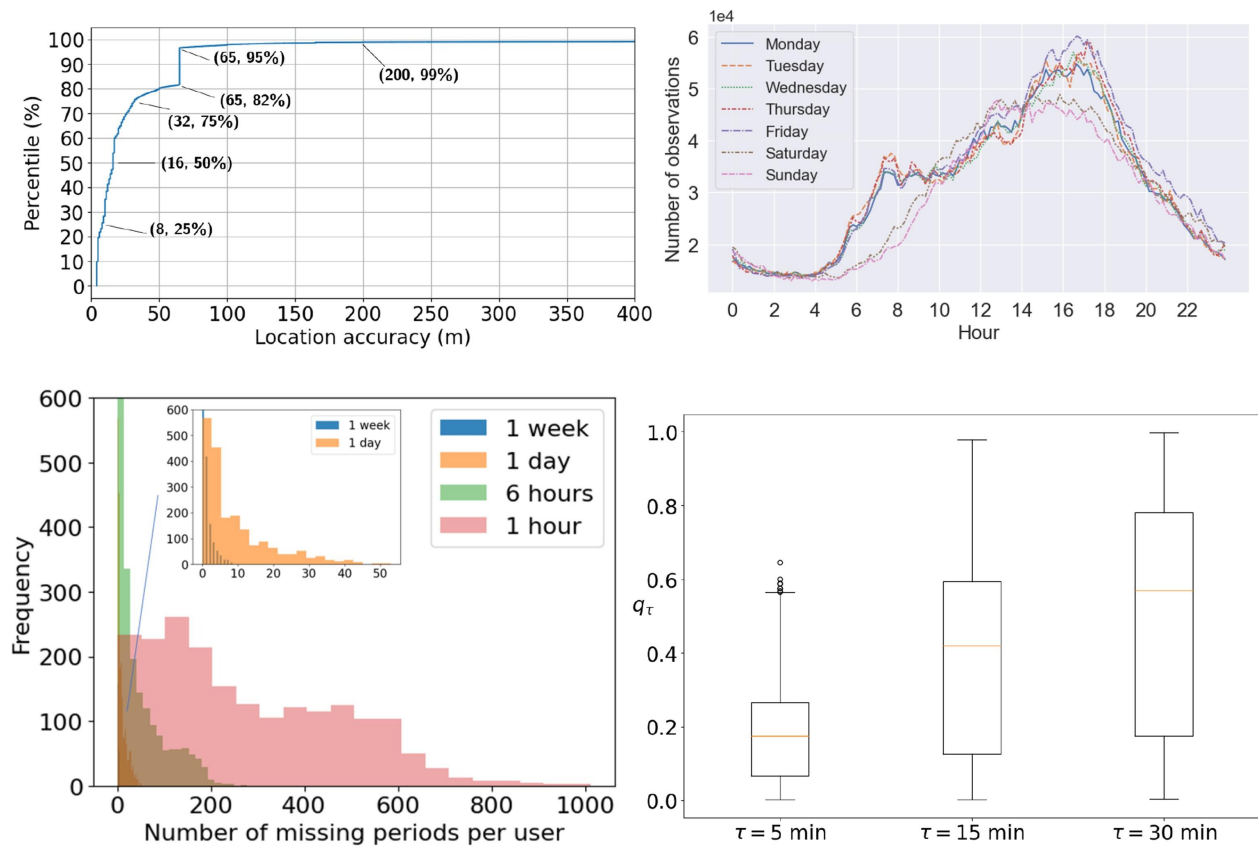


Figure 2.2: Descriptive analysis of Spectus data for all 2,000 users. (Top Left) Cumulative distribution of location accuracy. (Top Right) Distribution of observations within each day of the week. (Bottom Left) Histogram of gaps categorized by size. (Bottom Right) Boxplots of temporal occupancy with three temporal resolutions.

795 ARIMA and SARIMAX (Kohn and Ansley, 1986), as well as a multi-task GP using the
 796 standard RBF-ARD kernel on each input dimension (Duvenaud, 2014).

797 For the first set of experiments (Section 2.6.1), we first preprocess raw GPS traces to obtain
 798 two datasets for 10 randomly chosen users (1920 trips total): (1) a set showing all trip-related
 799 data points for a user, and (2) a set containing various trip information, including average
 800 velocity, total distance, total duration, among other metrics. We fit GPs to trip-related data
 801 points to predict missing locations. Then, we relate the optimized model parameter values
 802 with respect to the compressed trip information, which allows us to infer the type of trip (i.e.,
 803 highway drive or a walk). For the second set (Section 2.6.2), we randomly select a subset of
 804 50 users with at least 10,000 observations and at least 3 months of data (from the first data
 805 point until the last). Within the selected 50, we only retain data from January and February
 806 2020 to reduce the number of data points. We use this longitudinal data to train the model
 807 and make predictions on varying gap lengths.

808 2.6.1 Lengthscale Analysis

809 In this experiment, we determine well-fitting model parameters (i.e., ones that minimize the
 810 loss function) and relate them to a range of data characteristics, including average speed,
 811 total distance, and trip duration among other metrics. Identifying optimal kernel parameters
 812 is important for model accuracy. Figure 2.3 illustrates how a suboptimal lengthscale can
 813 deteriorate model fit. The left subplot shows a trip with noisy training data and a lengthscale
 814 value that is too low. The prediction’s mean reverts to the constant mean (likely near zero)
 815 only a few seconds after a training observation.

816 Furthermore, the optimization process deduces high levels of noise, producing uncertainty
 817 even on training inputs. Both processes result in low prediction accuracy as the prediction
 818 mean remains distant from the testing set. The right subplot, on the other hand, shows
 819 a model with low levels of noise and an appropriate lengthscale for the same trip. The
 820 prediction mean shifts smoothly between the different trip segments and the confidence
 821 interval is narrow near training points.

822 Using preprocessed Spectus data, we leverage k-means clustering to group together similar
 823 trips using trip characteristics (shown in the columns of Table 2.2). We filter a subset of
 824 mobile data based on three criteria, which help avoid anomalies like airplane-based travel
 825 and out-of-state roadtrips: average velocity does not exceed 80 mph, total travel time does 6
 826 hours, and the number of observations (in a trip) is greater than 4. Using the elbow method,

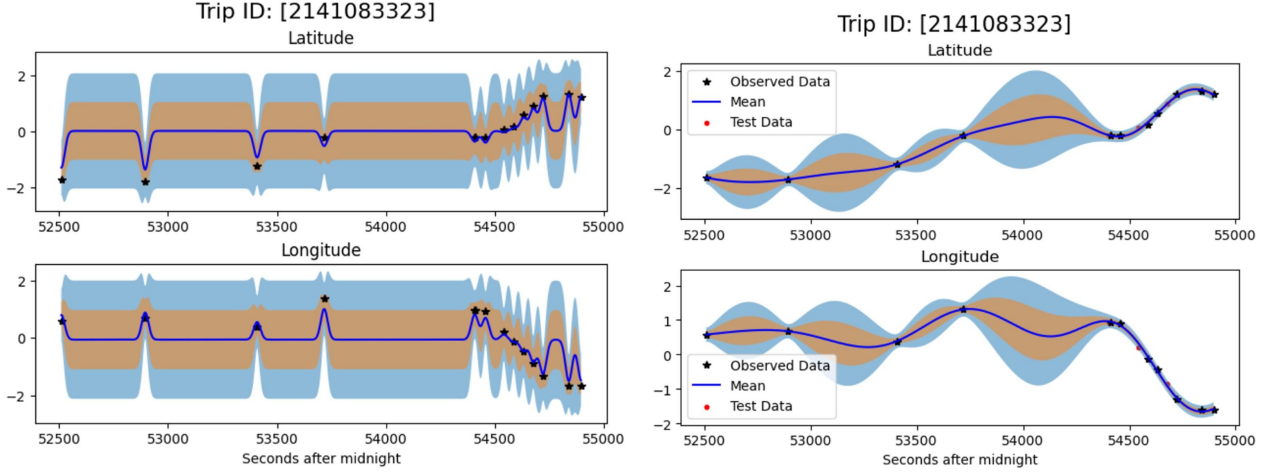


Figure 2.3: Examples of trips fit with an underestimated lengthscale (left) and appropriate lengthscale (right)

827 we identify three clusters from the filtered data as sufficient to explain most of the variation
 828 in trip-level data. Table 2.2 summarizes the average mobility metrics associated with the
 829 trips in the different clusters. The heading change rate shows the ratio of consecutive data
 830 points where a user changes direction with an angle exceeding a threshold (we use 0.33 rad).
 831 The velocity change rate shows the ratio of consecutive data points where the user exceeds a
 832 speed variation threshold (we use 26%). Finally, the stop rate represents the ratio of data
 833 points with an inferred velocity lower than a threshold (we use $0.89 \frac{m}{s}$). These threshold
 834 values are chosen in accordance with findings from Zheng et al.³ (2008) and are meant to
 835 discriminate between moving and being stationary.

836 We only used local data (i.e., observations from the same trip) in fitting GP models for
 837 these trips and thus our only input is \mathbf{t}_u (normalized Unix time). Accordingly, we use the
 838 RBF kernel, a simpler form of the earlier-defined RQ kernel:

$$k_{RBF}(x, x') = \Omega^2 \exp\left(-\frac{(x - x')^2}{2l^2}\right) \quad (2.12)$$

839 Figure 2.4a reveals that walking trips exhibit the lowest average RMSE, indicating higher
 840 prediction accuracy, whereas vehicle trips present the highest RMSE. Correspondingly, Figure
 841 2.4b's boxplot of optimized lengthscale parameters for each trip cluster indicates that faster

³See "Parameter Selection" on Zheng et al. (2008)

Table 2.2: Summary of trip clusters.

Cluster	Avg. Vel. [m/s]	Distance [m]	Duration [s]	Heading Change Rate	Velocity Change Rate	Obs.	Stop Rate
Slower, shorter trips	9.29	8k	1062	2e-3	2e-3	22.8	7e-4
Medium speed/distance trips	13.9	30k	2362	7e-4	8e-4	49.9	2e-4
Fast, distant trips	17.9	60k	3449	5e-4	6e-4	142	1e-4

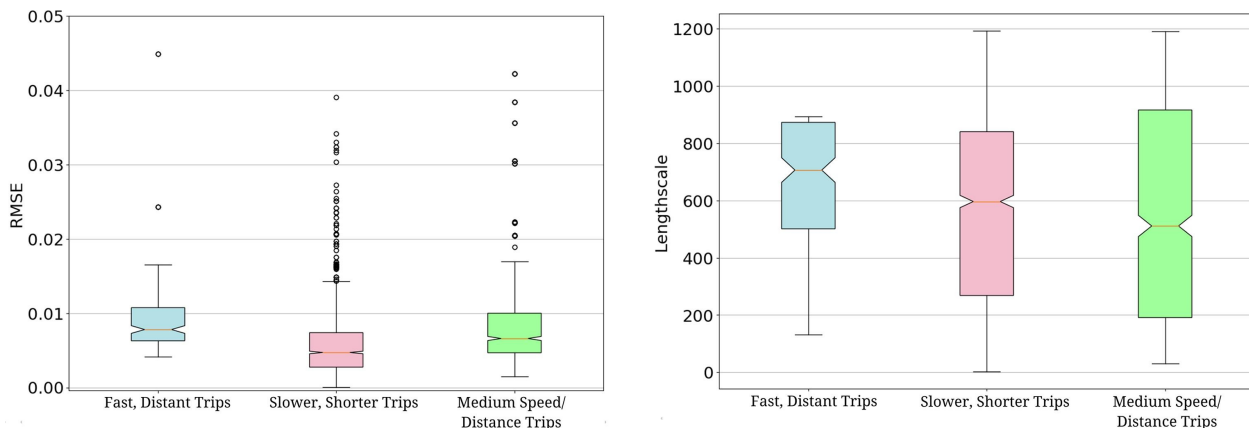


Figure 2.4: (Left) Boxplot of total RMSE for trips in different mobility metric clusters (Right) Boxplot of optimized lengthscales for each trip cluster.

842 trips, typically associated with vehicles, necessitate larger lengthscales, reflecting less frequent
843 changes in mobility patterns. In contrast, walking trips and mixed trips generally require
844 shorter lengthscales, suggestive of more variable and less predictable movements. Notably,
845 mixed trips display the widest range of lengthscales, underscoring a significant variability in
846 temporal correlation. This can be ascribed to the heterogeneity of mixed trips, which often
847 involve a combination of transportation modes, each with distinct velocities and movement
848 patterns. Such diversity introduces numerous non-smooth transitions, or 'kinks,' into the
849 data, which the Gaussian Process (GP) must adeptly capture.

850 The variability in lengthscales observed for walking trips may be attributed to the wide

851 spectrum of walking behaviors. On one hand, some individuals walk with a clear destination
 852 and route in mind, resulting in relatively predictable paths. On the other, walking can be
 853 exploratory, with no pre-determined endpoint, leading to more stochastic patterns featuring
 854 frequent stops and directional changes. This dichotomy in walking behaviors is corroborated
 855 by the high heading change and stop rates for walking trips as detailed in Table 2.2. Despite
 856 the complexities introduced by such diverse activities, both walking and mixed trips achieve
 857 lower median RMSE values, suggesting that the GP model is effectively modeling these trip
 858 types. This indicates that the GP framework is robust enough to handle the stochasticity
 859 inherent in walking trips and the complexity of mixed trips, accurately reflecting the nuanced
 860 shifts in mobility without a significant loss in predictive performance.

861 2.6.2 Robustness Checks: Gap Imputation

862 The goal of this experiment is to assess the performance of our model against other imputation
 863 methods in a variety of missingness conditions. In doing so, we conduct a robustness analysis
 864 on the gap length. Our benchmark methods include simple exponential smoothing (SES),
 865 exponential smoothing (ES), the Holt-Winters (Holt) method, ARIMA, SARIMAX, and an
 866 RBF-ARD kernel. The first three (ES, SES, and Holt) work similarly—they weigh the average
 867 of past observations in producing forecasts, giving exponentially lower weights as the time gap
 868 increases. SES is most suitable for forecasting data with no clear trend or seasonal pattern,
 869 while ES allows for forecasting with a trend. The Holt-Winters method uses three types of
 870 exponential smoothing to model the level (the typical value), the trend (the slope), and the
 871 seasonality (repeating patterns) of data. ARIMA and SARIMAX are both linear regression
 872 models fit for forecasting univariate time series data. While ARIMA can handle data with a
 873 trend, its extension SARIMAX can also handle exogenous variables and seasonal components.
 874 Each of these methods is regarded as a reliable, well-studied alternative commonly used to
 875 impute missing values in time series, and therefore make a useful comparison. The inclusion
 876 of the RBF-ARD kernel is to demonstrate that an off-the-shelf kernel implementation does
 877 not perform as well as our composite formulation. Table A.1 in Appendix A.4 describes and
 878 quantifies the parameters of each method.

879 Our results demonstrate that our model outperforms existing methods in all scenarios
 880 with varying levels of missingness. Specifically, we simulate gaps in the selected 50 users by
 881 reserving a subset of their data for testing, and we choose this subset such that the temporal
 882 occupancy of their training data meets a lower target temporal occupancy (i.e., a decimal

Algorithm 2: Gap Simulation in Trajectory Data

Input: Trajectory data; Bin length τ
Output: Trajectory data with gaps; New q_τ
 $\text{bins} \leftarrow$ Create an array of integers spanning from $t_{u,1}$ to $t_{u,N}$ with step size τ ;
 $\text{bins_dict} \leftarrow$ Use dictionary comprehension to map each data point to the relevant bin;
 $\text{non_empty_bins_dict} \leftarrow$ Make a similar dictionary with only the non-empty bins;
 $\text{target_ocp} \leftarrow \text{np.random.uniform}(0, q_{\text{curr}})$ // Take decimal floor of randomly chosen value between 0 and the current temporal occupancy as the target occupancy;
while $q_\tau > \text{target_ocp}$ **do**
 $\text{np.random.choice}(\text{non_empty_bins_dict.keys}()) \leftarrow$ Randomly choose a bin to remove and remove all values in this bin from original data;
 $\text{bins_dict} \leftarrow$ Update the original dictionary;
 $\text{non_empty_bins_dict} \leftarrow$ Update the non-empty bins dictionary;
 $\text{new_ocp} \leftarrow$ Calculate the new temporal occupancy with the gapped data;
end
return *The gapped trajectory data as well as the new q_τ ;*

883 between 0 and the current temporal occupancy) according to the temporal resolution being
 884 tested. We test six different temporal resolutions (τ): one week, one day, six hours, one hour,
 885 thirty minutes, and fifteen minutes. Algorithm 2 describes the full process to remove points
 886 from users' data⁴.

887 In the context of plotting and analyzing the results of various models with differing levels
 888 of data sparsity, we utilized Dynamic Time Warping (DTW) as an ancillary metric. DTW, a
 889 method established by (Müller, 2007), is suitable for measuring the similarity between two
 890 temporal sequences, accommodating variations in speed and timing. DTW values give us
 891 an indication of the dissimilarity between sequences; a lower DTW value indicates a closer
 892 match, whereas a higher value suggests a greater disparity between the compared sequences.

893 From the analysis depicted in Figure 2.5, we discern an inverse-linear relationship between
 894 temporal occupancy and DTW values, indicating that as temporal occupancy diminishes,
 895 DTW values—reflecting the disparity between testing set outcomes and model predictions—
 896 correspondingly escalate. This trend illustrates a progressive, rather than precipitous, decline
 897 in the model's predictive accuracy with sparser datasets. Notably, this inverse-linear rela-
 898 tionship appears less distinct at smaller gap lengths, as illustrated when comparing Figure
 899 2.5a for $\tau = 1$ week and Figure 2.5d for $\tau = 5$ minutes, suggesting a more stable model

⁴We also provide a step-by-step breakdown of this algorithm in Appendix A.2

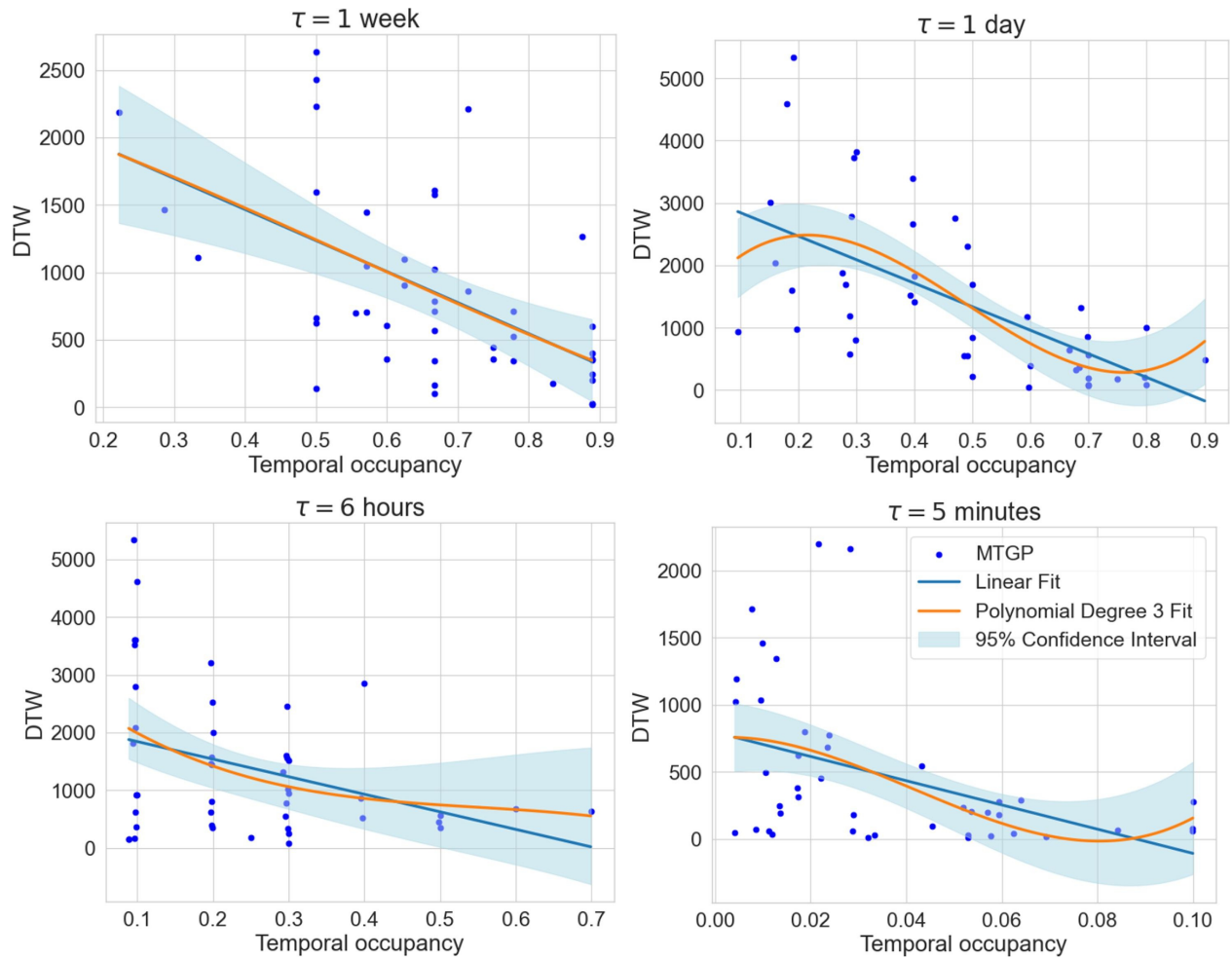


Figure 2.5: Scatterplots of DTW and temporal occupancy for various gap lengths. (a) $\tau = 1$ week, (b) $\tau = 1$ day, (c) $\tau = 6$ hours, (d) $\tau = 5$ minutes.

900 performance within these intervals. Despite the non-linear trends at the extremities of the
 901 temporal occupancies, potentially stemming from limited sample sizes, we do not observe a
 902 specific threshold beyond which the model’s performance deteriorates drastically. Instead,
 903 the model’s accuracy gradually tapers off with reduced temporal occupancy. This analysis
 904 emphasizes the model’s consistent performance across different scenarios while indicating
 905 that ensuring a minimum level of temporal occupancy—adjusted according to the specific
 906 temporal resolution of the data—is crucial for achieving the best predictive accuracy from
 907 the model.

908 For the benchmark algorithms, we used only the monotonically increasing \mathbf{t}_u as input
 909 since they do not have a straightforward way to deal with categorical or binary variables. For
 910 each individual, we trained an optimal benchmark model by maximum likelihood estimation
 911 (MLE) whenever applicable. In cases when we could not use MLE (i.e., determining the order
 912 parameters for ARIMA and SARIMAX), we used the grid search method to identify the
 913 optimal parameters (see Table A.2 in Appendix A.4).

914 We analyzed the accuracy of derived mobility metrics after the imputation of gapped
 915 training data compared to before simulating any gaps. Table 2.3 presents the median error
 916 results in each bin length. The boxplots of these error results can be found in Appendix A.4
 917 (Figures A.1-A.7). The same 50 users’ data was evaluated across the different gap lengths,
 918 allowing for an apples-to-apples comparison. In all tests, our model outperformed the other
 919 models: The average error was always the lowest, and the 5th and 95th percentile bounds
 920 tended to be closer. The metrics we compared across the different methods included the
 921 number of distinct locations, the radius of gyration, the straight-line distance traveled, real
 922 entropy, random entropy, and uncorrelated entropy. The number of distinct locations simply
 923 counted how many unique combinations of coordinate pairs existed in the imputed and
 924 original datasets. The straight-line travel distance traveled by an individual is computed as
 925 the sum of the distances traveled (Williams et al., 2015) and can be computed as

$$d_{SL} = \sum_{j=2}^m dist(r_{j-1}, r_j), \quad (2.13)$$

926 The radius of gyration indicates the characteristic travel distance of a user during a period

927 (Song et al., 2010b) and is given by the equation

$$r_g = \sqrt{\frac{1}{n} \sum_{i=1}^n (r_i - r_{cm})^2}, \quad (2.14)$$

928 where r_i represents the $i = 1, \dots, n$ locations recorded for the user and $r_{cm} = \frac{1}{n} \sum_{i=1}^n r_i$ is
 929 the center of mass of the period’s trajectory. Finally, the concept of entropy is used to assess
 930 a user’s predictability. The real entropy is given by the equation (Song et al., 2010b)

$$E_{real} = - \sum_T P(T') \log_2(P(T')) \quad (2.15)$$

931 where $P(T')$ is the probability of finding a particular time-ordered subsequence T' in the
 932 trajectory T . Therefore, the real entropy depends not only on the frequency of visitation, but
 933 also the order in which the locations were visited, and the time spent at each location, thus
 934 capturing the entire spatiotemporal order present in a user’s mobility. The random entropy
 935 captures a user’s degree of predictability if each of their distinct locations were visited with
 936 equal probability. The (temporal) uncorrelated entropy characterizes the heterogeneity of a
 937 user’s visitation patterns (Eagle and Pentland, 2009). The two are given by the equations

$$E_{rand} = \log_2(L) \quad (2.16)$$

938

$$E_{unc} = - \sum_{j=1}^L P(j) \log_2(P(j)) \quad (2.17)$$

939 where L is the number of distinct locations a user visits and $P(j)$ is the historical probability
 940 that a location was visited by the user.

941 We see that, on average, our model outperforms the rest of the benchmarks by a significant
 942 margin in estimating the number of locations visited, the radius of gyration, the straight-line
 943 travel distance, real entropy, random entropy, and uncorrelated entropy. Moreover, on these
 944 metrics, the distance between our model and the second-best model is significant. The
 945 standard RBF kernel multiplied across all input dimensions does not capture the variability
 946 of trips in sparse regions, performing worse than simpler smoothing algorithms even in low
 947 sparsity regimes. This highlights the importance of using appropriate kernel specifications
 948 while dealing with high-dimensional, mixed-type datasets.

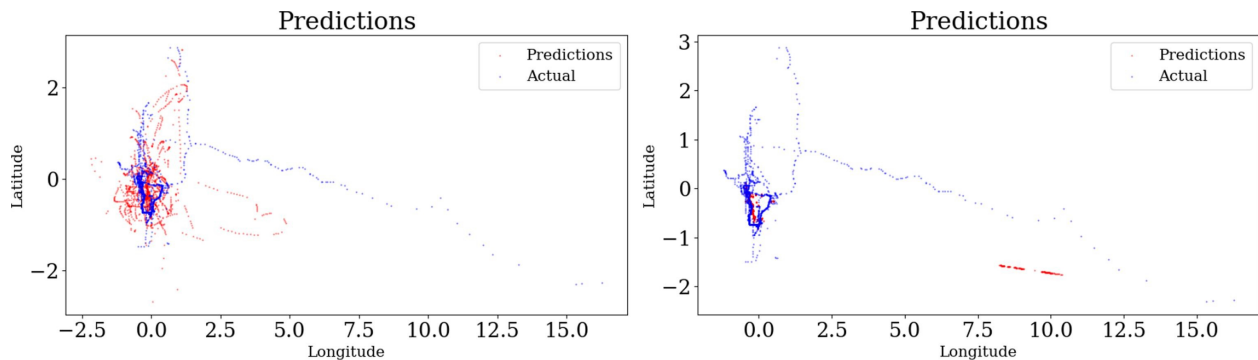


Figure 2.6: Predictions (red) and ground truth test data (blue) for a user tested at $\tau = 6$ hours. MTGP with the proposed kernel (left) outperforms exponential smoothing (right).

949 Figure 2.6 offers some perspective on the results presented here. In this example, a perfect
 950 predictor would have the red points align exactly with the blue points. The MTGP method
 951 (with the proposed kernel) better captures the variability of a user’s location over time than
 952 Exponential Smoothing—hence why it is more accurate in metrics like straight-line travel
 953 distance (Equation 2.13) and radius of gyration (Equation 2.14). We also note that, in
 954 general, MTGP seems less accurate in capturing route choice of users, but more accurate in
 955 learning destination choice. This is what we attribute to the performance we see on metrics
 956 like real and uncorrelated entropy (Equations 2.15 and 2.17, respectively).

957 2.7 Discussion

958 The use of passively-generated mobile data for transportation applications and decision-
 959 making is an irreversible trend. And yet, such data exhibits significant and varying levels of
 960 missingness compared to actively solicited data such as household travel surveys with GPS
 961 loggers or smartphone data collection components. Left untreated, the resulting mobility
 962 metrics are biased (McCool et al., 2022). To correct missingness in the mobile data, we have
 963 developed a novel framework using multi-task Gaussian processes. More specifically, our
 964 framework leveraged the correlations between users’ coordinates (thus multi-task), and allowed
 965 individual-level differences in data characteristics resulting from fundamentally different trip
 966 generation processes (e.g., different modes of transportation used). We introduced and
 967 demonstrated the effectiveness of RQ and PER kernels in the context of mobile data for
 968 human mobility modeling, and highlighted the ARD extension of these kernels for fitting

Table 2.3: Median error with respect to the testing sets.

Time Gap	Method	Number of Locations	Radius of Gyration	Straight Line Travel Distance	Random Entropy	Real Entropy	Uncorr. Entropy
1 week	MTGP	26	-0.07	205.029	0.045	0.278	0.153
	RBF	-801	-0.835	-888.441	-9.647	-9.323	-9.527
	SES	-801	-0.835	-888.441	-9.574	-9.323	-9.527
	Holt	-801	-0.835	-888.441	-9.574	-9.323	-9.527
	ES	-778	-0.621	-643.909	-4.989	-7.726	-4.955
	ARIMA	-801	-0.835	-888.441	-9.276	-9.323	-9.527
	SARIMAX	-801	-0.835	-888.441	-9.647	-9.323	-9.527
1 day	MTGP	33	-0.245	236.805	0.036	0.227	0.117
	RBF	-1050	-0.909	-1303.06	-10.038	-9.612	-9.806
	SES	-1050	-0.871	-1303.06	-9.309	-9.612	-9.806
	Holt	-1027	-0.718	-768.678	-4.875	-7.907	-5.225
	ES	-1050	-0.834	-1303.06	-8.506	-9.609	-9.803
	ARIMA	-1050	-0.834	-1303.06	-8.506	-9.609	-9.803
	SARIMAX	-1050	-0.909	-1303.06	-10.038	-9.612	-9.806
6 h	MTGP	34	-0.187	-13.641	0.042	0.237	0.155
	RBF	-956.5	-0.645	-1223.47	-9.809	-9.493	-9.751
	SES	-954	-0.645	-1177.23	-9.139	-9.493	-9.608
	Holt	-929	-0.645	-1177.19	-4.895	-7.869	-5.066
	ES	-929	-0.718	-768.56	-4.895	-7.869	-5.066
	ARIMA	-954	-0.645	-1178.65	-8.317	-9.493	-9.608
	SARIMAX	-956.5	-0.645	-1223.47	-9.901	-9.493	-9.751
1 h	MTGP	38	-0.074	389.303	0.049	0.257	0.157
	RBF	-990	-0.994	-1319.47	-9.818	-9.548	-9.761
	SES	-901.5	-0.761	-1262.95	-7.989	-9.546	-9.642
	Holt	-878	-0.627	-1161.74	-6.875	-9.149	-5.174
	ES	-898.5	-0.761	-1262.76	-6.801	-9.546	-9.613
	ARIMA	-898.5	-0.761	-1262.76	-6.801	-9.546	-9.613
	SARIMAX	-830.5	-0.644	-1161.14	-6.435	-9.455	-9.449
15 min	MTGP	22	-0.299	-7.116	0.048	0.323	0.161
	RBF	-670	-2.15	-1112.56	-8.925	-8.871	-9.312
	SES	-660	-2.17	-1162.11	-7.734	-8.871	-8.994
	Holt	-660	-2.099	-1162.07	-6.435	-8.871	-8.849
	ES	-637	-2.056	-513.035	-3.906	-7.711	-4.497
	ARIMA	-591	-1.931	-1162.07	-6.435	-8.871	-8.851
	SARIMAX	-670	-2.15	-1119.07	-9.39	-8.871	-9.416

969 high-dimensional data (whether continuous or binary). Our experiments highlighted the
970 importance of specifying a fitting composite kernel for our domain, as using the generic
971 RBF kernel resulted in sub-optimal predictions due to not being able to capture multi-level
972 periodicities. In the results, we showed that trips made by different modes are associated
973 with different, optimally estimated kernel parameters (i.e., the lengthscale) and provided
974 guidelines on different parameter initializations for different kinds of trips. More importantly,
975 we showed that our model outperformed six existing methods for correcting missingness in
976 mobile data and assessed the temporal sensitivity of our model, confirming that our model
977 achieves enhanced predictive performance when trained on a diverse temporal dataset.

978 While the current study is motivated by the need to directly address the significant sparsity
979 issue exhibited in passively-generated mobile data (McCool et al., 2022), the implications
980 of this work can be significant in several aspects. For one, it directly addresses the bias
981 issue induced by data sparsity and the resulting, corrected trajectory data can be viewed as
982 pseudo-ground truth from which mobility metrics can be derived. In another, it opens us the
983 opportunity to create city-wide simulations of human mobility patterns using the generated
984 mobile data directly. This is in contrast with the current process where parametric models are
985 first estimated using the household travel survey data (which often represents less than 1% of
986 a region’s population) and then those models are extrapolated to a synthetic population. As
987 noted earlier, the latter method is not scalable as travel surveys must be frequently collected
988 and models must be updated. An added disadvantage is that those models do a poor job of
989 capturing the nonlinearities embedded in the data.

990 Transportation behaviors, whether they refer to individual travel behaviors as studied
991 in this paper or community- or region-level traffic phenomena, are highly complex and
992 non-linear. Thus, expecting a universal model to capture all is likely unrealistic. Instead,
993 context-dependent modeling is much needed for the field of transportation. Context-dependent
994 modeling requires the development of flexible modeling frameworks that can capture not only
995 heterogeneity but also adapt to changing contexts. We demonstrate that GP-based methods
996 are one suitable framework for this purpose in the context of modeling individual mobility
997 patterns.

998 Our work underscores the critical role of kernel specification in capturing the nuances
999 exhibited by individual mobility patterns. In our study, we used a domain-based approach,
1000 i.e., using our domain knowledge in the travel behavior literature to guide the selection and
1001 the composition of the kernels. Through our experiments, we show that standard, off-the-

shelf kernels, such as the Radial Basis Function (RBF) kernel, frequently underperform in comparison to more complex kernels across a range of imputation scenarios. The process of specifying kernels is akin to a traditional, parametric modeling process where one would have to start with a hypothesized model structure. One may argue that even in data-driven modeling, one would still have to start with a modeling framework or architecture that define how input data are taken in and evolves over different layers, and whether there are feedback loops or not. Just as these models often require rigorous cross-validation to determine the most suitable hyperparameters, Gaussian Processes similarly benefit from a structured approach to kernel selection. Addressing this issue, notable advancements have been made in the development of automated methods for kernel determination. Duvenaud et al. (2013) have proposed a method for automating the construction of kernel expressions, thereby facilitating the identification of suitable kernel structures directly from the data. In a similar vein, Wilson and Adams (2013) introduced a framework that learns expressive covariance functions for GPs through spectral mixture kernels capable of automatically adapting to the structure of the data. These innovations are significant steps forward in simplifying the kernel selection process, thereby enhancing both the performance and the broader applicability of the model across various domains.

Nevertheless, our study also recognizes certain limitations. One limitation is the challenge in selecting kernels that accurately emulate the data’s inherent features, which is compounded by the issue of model interpretability. Although our model offers improved interpretability over deep learning (DL) models, which are often criticized for their ‘black box’ nature, it lacks the transparency of simpler, more traditional statistical methods. The inherent complexity of Gaussian processes, particularly when employing sophisticated kernels such as Rational Quadratic (RQ) and Periodic (PER), introduces significant challenges in terms of interpretability. While model interpretability is not of paramount significance to our study (since the study objective is simply fixing missingness in raw mobile trajectories), potential applications of GPs in other transportation applications (such as those reviewed earlier in Section 2.2 for forecasting travel demand or traffic flows) would require higher levels of model interpretability. Thus, balancing model complexity with interpretability is an essential area for future research that involve GPs or other big data methods.

Another future direction of this work is to tackle computational complexity—the time it takes to train a model scales cubically with the size of the covariance matrix, which is a function of the size of the data. We outline a few ideas to explore this front. One is

1035 to accelerate the process of GP parameter optimization through CUDA (Compute Unified
1036 Device Architecture), a parallel computing platform and API created by NVIDIA that allows
1037 the software to use certain types of graphical processing units (GPUs) for general-purpose
1038 processing. Doing so would allow future models to add additional layers in a GP framework
1039 (i.e., deep GPs), predicting covariates like velocity and bearing before producing a location
1040 output. We can also pre-process the input data to drastically reduce the number of times we
1041 have to re-evaluate Equation 2.10 above. An example is provided by Lee et al. (2017).

1042 There is also the need to incorporate the underlying built environment including multi-
1043 modal road networks into the model, such that bodies of water, buildings, or locations of
1044 different types can be recognized. Methods like map-matching may be employed to post-
1045 process model results such that the predictions can only be made within allowable regions.
1046 However, this approach would have a trade-off with computational complexity—though
1047 scalable map-matching algorithms have been proposed in recent years (Fiedler et al., 2019;
1048 Zeidan et al., 2020), convenient implementations are not yet widely available. Alternatively,
1049 one may also explore additional selection and identification of composite kernels (like Equation
1050 2.9) to potentially account for those contextual effects.

1051 Application-wise, there are numerous ways that the corrected human trajectory data
1052 can be used. There are numerous applications of passively-generated mobile data in the
1053 existing literature, which are exclusively based on uncorrected trajectory data. It would
1054 thus be interesting to see whether the existing findings still stand if the corrected trajectory
1055 data is used. In a previous paper by one of the authors, we showed that not removing
1056 the oscillation phenomenon in the data leads to the non-negligible overestimation of the
1057 regularity of individuals' mobility (Wang and Chen, 2018). We similarly expect improvements
1058 in accuracy for downstream studies that use sparsity-corrected passively-generated mobile
1059 data. Future efforts could also be made to unify existing literature on applications of GPs
1060 to transportation problems, such as the works of Batista et al. (2022) and Gammelli et al.
1061 (2020), highlighting shared model properties, cautioning on common modeling mistakes, and
1062 establishing standards of benchmarking.

Chapter 3

**LEARNING TO GENERATE SYNTHETIC HUMAN MOBILITY
DATA: A PHYSICS-REGULARIZED GAUSSIAN PROCESS
APPROACH BASED ON MULTIPLE KERNEL LEARNING**

Passively-generated mobile data has grown increasingly popular in the travel behavior (or human mobility) literature. A relatively untapped potential for passively-generated mobile data is synthetic population generation, which is the basis for any large-scale simulations for purposes ranging from state monitoring, policy evaluation, and digital twins. And yet, this significant potential may be hindered by the growing sparsity or rate of missingness in the data, which stems from heightened privacy concerns among both data vendors and consumers (users of service platforms generating individual mobile data). To both fulfill the great potential and to address sparsity in the data, there is a need to develop a flexible and scalable model that can capture individual heterogeneity and adapt to changes in mobility patterns. We propose a conditional-generative Gaussian process framework that learns kernel structures characterizing individual mobile data and can provably replicate observed patterns. Our approach integrates physical knowledge to regularize the framework such that the generated data obeys constraints imposed by the built and natural environments (such as those on velocity and bearing). To capture travel behavior heterogeneity at the individual level, we propose a data-driven multiple kernel learning approach to determine the optimal composite kernel for every user. Our experiments demonstrate that: (1) the impact of kernel choice on mobility metrics derived from synthetic data is non-negligible; (2) physics-regularization not only reduces model bias but also improves uncertainty estimates associated with the predicted locations; and (3) the proposed method is robust and generalizes well to varying individuals and modes of travel.

3.1 Introduction

The ubiquity of sensor-equipped mobile devices has enabled the generation of passively-generated large-scale time- and location-stamped data (hereafter called "mobile data"), including GPS traces, geo-tagged posts from social media platforms, and call detail records

1091 (Chen et al., 2016b). These datasets are often massive in size (covering millions of residents in
1092 and travelers passing through a region) and longitudinal that can last from days to months and
1093 even years. These two attributes are fundamentally different from household travel surveys
1094 which are snapshots of a region’s travel patterns with a small sample (typically less than 1%
1095 of a region’s population) and have been traditionally used for travel behavior and demand
1096 forecasting studies. These big mobile datasets offer great promises in moving travel behavior
1097 or human mobility studies forward in various applications. Existing studies of the past 15
1098 years using these passively-generated mobile data have mostly focused on identifying models
1099 that can capture both regularities and variations of human mobility patterns (Gonzalez et al.,
1100 2008; Yuan and Raubal, 2012). The emergence of the COVID pandemic as well as AI in the
1101 last several years have stimulated additional studies that seek to identify cause and effect
1102 relationships (e.g., whether a change in the built environment leads to changes in mobility
1103 patterns; Wang et al., 2021a) and to understand changes in mobility patterns to support
1104 timely policy making (e.g., social distancing policy-making during pandemic requires rapid
1105 assessment of changes in people’s mobility patterns; Morris et al., 2021).

1106 An emerging application of using passively-generated mobile data is synthetic population
1107 generation for large-scale simulation of individual mobility patterns (Deng et al., 2021),
1108 whose applications are wide-ranging from state monitoring, to whole-region or targeted
1109 policy making, and long-term investment decisions. Using passively-generated mobile data,
1110 which are not only massive but also longitudinal, for this purpose has the potential to
1111 transform the current state of the art in the field that are primarily of two kinds. The
1112 first kind is the long-standing Activity-based Models (ABM) which are essentially a set of
1113 linked econometric models whose parameters are estimated using household travel survey
1114 data (Bhat and Koppelman, 1999; Pendyala et al., 1997; Timmermans and Arentze, 2011).
1115 Estimated econometric models reflect average trends and thus direct application of those
1116 models at the individual level result in large amounts of errors. Furthermore, ABM models
1117 are often estimated on small cross-sectional datasets and thus are static and cannot capture
1118 how behaviors evolve over time. The second kind are data-driven models that leverage
1119 population-level probabilistic distributions estimated from the big mobile data. These models
1120 capture individual heterogeneity through the expression of a distribution and fundamentally
1121 they cannot account for or adapt to nuances that are unique to an individual. These nuances
1122 may include factors unique to the individual (e.g., socio-demographics), the trip itself (e.g.,
1123 the specific modes of transportation being used), and the underlying physical environment

1124 at the moment (e.g., road networks regulate average velocity and bearing). Fundamentally,
1125 neither of the two kinds of the existing studies touches upon uncovering the underlying data
1126 generation process at the individual level, which is a critical task for passively-generated
1127 mobile data to be used for synthetic population in large-scale simulations. This constitutes
1128 the central aim of our study.

1129 The great potential of the passively-generated mobile data as stated above can be hampered
1130 by the fact that the consumer data privacy protection landscape has become increasingly
1131 user-centric, bestowing individuals with much greater flexibility in how their data is handled
1132 (i.e., easily allowing users to opt out of location-based data sharing agreements). This is
1133 due to a combination of factors both from the supply and demand sides: From the latter,
1134 consumers are now much more aware of the implications of unlimited data collection (i.e.,
1135 hyper-targeted advertisements) and, as a result, wary of opting into agreements from which
1136 they do not receive any benefits (Kim et al., 2019). Simultaneously, companies like Apple have
1137 spearheaded efforts to promote user-centric privacy, implementing a range of new features in
1138 2021’s iOS 15 to help users better control and manage access to their data (Apple, 2021).
1139 Though the jury is still out on whether the goal of Apple’s new initiatives were to maximize
1140 stakeholder returns or avoid anti-trust legislation (or both), the end result is the same: there
1141 has been a significant decrease in the amount of data generated, or a significant increase in
1142 sparsity. This emphasizes, from another perspective, the need to develop models that can
1143 uncover the underlying data generation process of passively-generated mobile data, as such
1144 models can generate synthetic data where consumer privacy will be less of an issue.

1145 In this paper, we propose a framework that can be used to potentially generate synthetic
1146 individual mobile data that replicates real travel behavior. This approach can achieve two
1147 ends: It can augment existing datasets with sparsity, and alternatively, it can be substituted
1148 for real raw data while conducting analyses without compromising user privacy. In modeling
1149 travel behavior or mobility patterns, individual heterogeneity has been well-documented in
1150 numerous studies, as affected by many systematic factors such as demographic, socioeconomic,
1151 temporal, cultural, and built environment-related factors (Bayarma et al., 2007; Kitamura
1152 and Van Der Hoorn, 1987; McGuckin and Murakami, 1999a; Nishii et al., 1988; Wallace
1153 et al., 2000), as well as spur of the moment factors (Lee and McNally, 2006, 2003). Thus, the
1154 challenge is to develop a flexible model framework that can adapt to individual uniqueness
1155 and easily scale up as more and new data come in. This model cannot take a parametric form
1156 that requires the pre-specification of a functional form. To this end, we propose a multi-task,

1157 multiple kernel Gaussian process framework that discovers optimal kernel specification for each
 1158 user. Gaussian Processes (GPs) are a generalization of multivariate Gaussian distributions to
 1159 infinitely many variables. A GP usually consists of two main components—a mean function
 1160 that specifies the mean at any point in the output space, and a covariance (kernel) function
 1161 that embeds a measure of similarity between any pair of observations in a multi-dimensional
 1162 space. The most crucial component of a GP model is its kernel, as it determines the class
 1163 of functions the GP relies upon (Rasmussen and Williams, 2006). To learn the form of the
 1164 covariance function, we combine candidate kernels in a domain-inspired and data-driven
 1165 manner to optimize model performance and use the Kronecker product to reduce the time it
 1166 takes to invert a complex covariance matrix.

1167 While powerful, unconstrained GPs can generate predictions that do not align with
 1168 real-world conditions. Discrete observations of mobile data tend to satisfy certain physical
 1169 bounds, which stem from the underlying street network and the broader built and natural
 1170 environment. Specifically, we deal with two dimensions of physics: (1) Land-based human
 1171 mobility has inherent limitations on the speed of travel, mostly related to the underlying
 1172 street network as well as the time of day¹. For example, a driver cannot travel 100 kilometers
 1173 per hour in a residential street. Similarly, trips at night may have higher average velocities
 1174 due to reduced traffic. (2) The direction of a user between consecutive data points is limited
 1175 by the mode of travel and the built environment. If there is no paved road going east at a
 1176 juncture, a person traveling by car would not be able to traverse that direction, just as a
 1177 person cannot walk into a lake. Thus, our approach integrates physical knowledge, such as
 1178 segment velocities and bearings, into the GP inference process. Our conditional-generative
 1179 framework, inspired by (but also differs from) the likes of Raissi et al. (2019), Lasserre et
 1180 al. (2006), Wang et al. (2022b), works by first sampling a set of locations as a function of
 1181 time (conditional model) and then using a generative process to simulate physical variable
 1182 observations at a certain time and location. However, our framework differs from existing
 1183 methods in that it does not generate new training data before posterior estimation. Instead,
 1184 it samples training data during estimation and generates synthetic physical inputs for model
 1185 inference.

1186 Our Contributions

- 1187 • We propose a conditional-generative GP framework to generate synthetic individual

¹This applies to water-based and air-based human mobility as well. However, because these trips are not captured in the dataset we use, we do not consider them in our methodology.

- 1188 mobile data that provably replicates observed travel patterns (Theorem 1);
- 1189 • We formulate a flexible physics-regularization framework that permits the learning of
1190 composite kernels that can accommodate nonsmooth and nonstationary patterns in
1191 human mobility modeling;
 - 1192 • We develop a data-driven multiple kernel learning (MKL) approach to determine the
1193 best spatiotemporal kernel for each individual;
 - 1194 • Our experiments lead us to three findings: First, the impact of kernel choice on mobility
1195 metrics derived from synthetic data is non-negligible. Second, incorporating physics-
1196 based regularization not only diminishes model bias but also enhances the precision of
1197 uncertainty estimates in predicting locations. Third, the proposed method is robust
1198 and generalizes well to varying individuals and modes of travel.

1199 **3.2 Related Work**

1200 Machine learning techniques like artificial neural networks (ANN) and kernel-based approaches
1201 have been particularly promising in modeling human mobility, with notable recent develop-
1202 ments in map-matching (Jiang et al., 2023) and trajectory completion (Ugurel et al., 2024a;
1203 Wang et al., 2020). The basic idea with many ANN-based models is to replace unknown
1204 functions with a neural network and then minimize a loss function with respect to the param-
1205 eters of the network. Although ANN-based models have demonstrated empirical success in
1206 tackling complex problems, their performance can be inadequate for simpler problems if the
1207 ANN architecture is not tailored to the underlying domain (Wang et al., 2021b, 2022a). This
1208 goes hand-in-hand with the fact that the theory of ANNs has lagged behind its empirical
1209 success, with many ANN-based PDE solvers lacking rigorous theoretical guarantees, and
1210 convergence analyses often being limited to linear instances with constraints (Shin et al.,
1211 2020; Wang et al., 2022a). As Chen et al. (2021) points out, when compared to kernel
1212 methods, ANN methods suffer from limited theoretical foundations and present less favorable
1213 trade-offs between computational complexity and accuracy estimates. These limitations are
1214 especially salient in generative settings, where complex ANN-based architectures are now
1215 widely used to synthesize human mobility traces but often remain black boxes with limited
1216 interpretability or uncertainty quantification. The generative nature of kernel methods can

1217 provide a more elucidated mathematical form of its model and can provide uncertainty
1218 quantification, probabilistic inference, and simulation.

1219 Beyond tasks like map-matching and trajectory completion, a growing line of work uses
1220 data-driven models to generate synthetic human mobility data as a privacy-preserving proxy
1221 for raw trajectories. Early approaches build structured stochastic models that decouple
1222 temporal routines from spatial movement, for example diary-based frameworks that learn
1223 Markovian patterns of routine versus non-routine activity and then assign locations via
1224 preferential exploration and return mechanisms (Pappalardo and Simini, 2018). More recent
1225 studies rely on deep generative models. Recurrent neural networks have been used to simulate
1226 individual-level traces by learning next-location sequences (Kulkarni and Garbinato, 2017),
1227 and later extended to generate synthetic populations over multiple days while explicitly
1228 tuning randomness to balance statistical fidelity with person-level dissimilarity for privacy
1229 (Berke et al., 2022). Generative adversarial networks with mobility-aware regularization
1230 further integrate spatial relations and daily periodicity to reproduce a wide range of empirical
1231 mobility statistics and support downstream simulations such as epidemic spread (Feng et al.,
1232 2020). Diffusion-based models push this idea further by training powerful denoisers that
1233 produce high-fidelity synthetic GPS datasets for urban analysis (Zhu et al., 2023). Most
1234 recently, large language model style transformers have been adapted to mobility generation,
1235 treating sequences of activities, locations, and modes as tokens and conditioning on socio-
1236 demographic attributes to generate rich traveler profiles (Zhang et al., 2024b). A recent
1237 systematic review highlights that, despite their empirical realism, these ANN-based generators
1238 generally prioritize utility, offer limited formal privacy guarantees, and provide little in the
1239 way of calibrated uncertainty or interpretable structure (Kapp et al., 2023).

1240 In contrast to these predominantly ANN-based generative models, Gaussian Process
1241 priors offer a Bayesian alternative for modeling and simulating complex spatio-temporal
1242 processes. Rasmussen & Williams (2006) wrote extensively on modeling and predicting
1243 complex interactions in space and time while accounting for spatial non-stationarity. The
1244 flexibility of Gaussian Processes in capturing nonlinearity and incorporating prior knowledge
1245 through kernel selection has demonstrated their efficacy in uncovering hidden patterns and
1246 predicting in dynamic geospatial contexts. The major limitation of using GPs is that the
1247 expert must choose a kernel to describe the underlying functional relationships in the data.
1248 An ill-fitting kernel can degrade GP model performance, leading to poor generalization and
1249 inadequate feature representation (Gönen and Alpaydm, 2011). In the human mobility

1250 domain, only a few previous works have explored preliminary methods on this kernel selection
1251 problem, while many are limited to choosing a kernel from an established off-the-shelf list
1252 of common kernels. Liu and Onnela (2021) proposed an additive kernel with a daily and a
1253 weekly temporal component as well as a distance-based spatial component to impute missing
1254 values in longitudinal mobile data, while Batista et al. (2022) focused primarily on the RBF
1255 and Matern kernels. The literature lacks a systematic framework established to account for
1256 heterogeneity in travel behavior at the individual level.

1257 As a data-driven method, the performance of GP models can suffer when the training
1258 data does not sufficiently reflect the complexity of the system (Wang et al., 2022b). Left
1259 unconstrained, a GP’s behavior may not satisfy expected system outcomes (Ugurel, 2023).
1260 Furthermore, GPs may produce predictions that are inconsistent with real-world conditions,
1261 defying the laws of physics or the constraints imposed by the built and natural environments.
1262 Physical models, built using physics knowledge expressed through differential equations, can
1263 characterize the underlying principles governing a system without needing large amounts
1264 of training data (Lapidus and Pinder, 1999). Thus, marrying data-driven methods with
1265 physical models not only provides insights into the system’s mechanics but also enhances
1266 prediction accuracy by overcoming data limitations, which addresses the issues arising from
1267 unconstrained GPs.

1268 For synthetic mobility generation, the lack of explicit physical constraints in existing
1269 ANN-based models means that generated traces can violate simple laws of motion or built-
1270 environment constraints, even when they match aggregate statistics. Physics-regularization
1271 has been extremely popular in the context of ANNs (Cuomo et al., 2022; Raissi et al., 2019).
1272 Though several studies have proposed analogous approaches using GPs and other kernel
1273 methods (Chen et al., 2021; Nevin et al., 2021; Wang et al., 2022b; Yang et al., 2019), the
1274 methodology has not been developed as extensively. Specifically, the studies focusing on
1275 physics-regularized GPs have proposed methods for deriving physical knowledge for the
1276 estimation phase (i.e., to obtain the posterior). In this study, we use the human mobility
1277 domain to our advantage in overcoming this estimation problem, specifically by reducing
1278 the noise associated with direct estimation through proper data labeling (see Section 3.3.3).
1279 Conversely, the challenge of inference using this posterior, particularly when it involves
1280 physical variables that are dependent on both the inputs and the outputs, has not been
1281 extensively explored.

1282 Within the transportation domain, Yuan et al. (2021) proposed a physics-regularized GP

1283 analogous to that of Wang et al. (2022b) to model macroscopic traffic flow, while Zhu et al.
 1284 (2022) proposed a physics-regularized multi-output GP to model the numbers of pickups,
 1285 returns, and idle bikes of a large-scale bike-sharing system. Recently, Wu et al. (2024a)
 1286 used anisotropic Gaussian Processes to model congestion propagation in traffic flow data.
 1287 These studies highlight the growing interest in applying physics-regularization in diverse
 1288 contexts, yet there remains a gap in its application to individual-level mobility data where
 1289 each individual has their unique spatial coverage and the data can extend for days, weeks or
 1290 months. To our knowledge, this study is one of the first to propose a method for generating
 1291 GPS-based human mobility traces that (i) models individual-level heterogeneity through
 1292 tailored kernels and (ii) explicitly regularizes the latent functions with physical attributes
 1293 associated with the built environment, thereby complementing existing ANN-based synthetic
 1294 mobility generators.

1295 **3.3 Methodology**

1296 *3.3.1 Problem Definition*

1297 Mathematically, the problem of uncovering the underlying data generation process for synthetic
 1298 individual mobile data generation translates into two related research questions:

- 1299 1. Given time, how do we infer (predict) spatial locations?
- 1300 2. How do we infuse physics (i.e., velocity and bearing) into the inference problem from
 1301 time to location, as stated above?

1302 The input time can be represented by a set of variables, in addition to being simply a
 1303 time index (e.g., the Unix time). Those additional time-related variables include periodic
 1304 elements like the day of the week and weeks of the month. They capture varying rhythms that
 1305 have long been identified in travel behavior literature (Bayarma et al., 2007; Gonzalez et al.,
 1306 2008; Hanson and Huff, 1988): e.g., people usually go to their workplaces and return home
 1307 every day and other trips like grocery shopping, going out with friends, or attending family
 1308 gatherings happen less frequently. Additionally, there are behavioral patterns correlated
 1309 with certain weeks of a month or months of a year (i.e., Christmas, Thanksgiving, etc.). We
 1310 denote with \mathbf{T} the matrix of temporal variables, where $t_{u,i}$ is the u^{th} temporal dimension at
 1311 the i^{th} observation. The physics-related variables such as average velocity and bearings are
 1312 to capture the constraints that the built environment has on one’s movement pattern. We

1313 denote with \mathbf{P} the matrix of physics-related variables, where v_i and β_i is the average velocity
 1314 and bearing at i^{th} observation respectively. The output variable location \mathbf{Y} is represented
 1315 with latitude λ and longitude ϕ ²:

$$\mathbf{T} = \begin{bmatrix} t_{1,1} & \dots & t_{d,1} \\ \vdots & \ddots & \vdots \\ t_{1,n} & \dots & t_{d,n} \end{bmatrix} = \begin{bmatrix} \mathbf{t}_1 \\ \vdots \\ \mathbf{t}_n \end{bmatrix}, \quad \mathbf{P} = \begin{bmatrix} v_1 & \beta_1 \\ \vdots & \vdots \\ v_n & \beta_n \end{bmatrix}, \quad \mathbf{Y} = \begin{bmatrix} y_{\lambda,1} & y_{\phi,1} \\ \vdots & \vdots \\ y_{\lambda,n} & y_{\phi,n} \end{bmatrix}. \quad (3.1)$$

1316 Human movements, as reflected in individual-level mobile phone trajectory data, are highly
 1317 non-linear, context-dependent, and heterogeneous. Underlying the observed mobile data is
 1318 the wide range of transportation modes used, each with different underlying physics (i.e.,
 1319 average velocity, acceleration behavior). When more than one mode of transportation is used,
 1320 there is also a change of modes that need to be captured. Individual movement behaviors
 1321 are also heterogeneous. While previous works have explored "universal" models on human
 1322 movement patterns (Gonzalez et al., 2008; Song et al., 2010b), later studies have revealed
 1323 many other models that can capture human movement patterns and no single model will
 1324 outperform others across contexts. These individual-level complexities suggest that making
 1325 rigid assumptions about the shape or form of the underlying distribution of an individual's
 1326 mobile data is unlikely a fruitful approach. A non-parametric, data-driven method holds
 1327 greater promise, being more flexible and less constrained by predefined assumptions. Such an
 1328 approach is capable of accommodating complex patterns and more adeptly capturing the
 1329 inherent relationships present within the data. Thus, to answer the first research question
 1330 stated above, we propose a multi-task Gaussian Process that learns individualized kernels to
 1331 infer a set of locations \mathbf{Y} from a set of temporal variables \mathbf{T} . As noted earlier, identifying
 1332 kernels, or more specifically the correct structure for a combination of kernels, is critical as a
 1333 fitting kernel can seamlessly capture complex dependencies in mobile data while a non-fitting
 1334 one will be inadequate. In this paper, we address this multiple kernel learning problem in the
 1335 context of mobile trajectory data (see Section 3.3.4).

1336 To answer the second research problem stated above, we need information on the physics-
 1337 related variables, or \mathbf{P} . And yet, \mathbf{P} is not directly observed. More specifically, this issue of
 1338 non-observability is different in the estimation problem where the kernels are learned from

²For simplicity, we project λ and ϕ to a two-dimensional grid such that a user's location at time t is $\mathbf{y}_t = [y_\lambda, y_\phi]$. As travel distances between consecutive data points in mobile datasets tend to be short, this does not make a significant difference compared with extensions to spherical geometry.

1339 \mathbf{T} to \mathbf{Y} versus the inference problem where \mathbf{Y} is inferred based on \mathbf{T} . In the estimation
 1340 problem, since \mathbf{T} and \mathbf{Y} are both observed in the training data, \mathbf{P} can be calculated. This
 1341 calculation, of course, is unavoidably tampered with noise, especially when the time gap
 1342 between two consecutive observations is large. We address this issue by a sample labeling
 1343 procedure that reduces bias in the estimation of physical variables (see Section 3.3.3). In the
 1344 inference problem, \mathbf{Y} is not available, since it is to be inferred. Thus, we must first estimate
 1345 \mathbf{P} by considering its dependency on the yet-to-be-predicted output variable \mathbf{Y} , as well as on
 1346 \mathbf{T} . We achieve this by formulating a model to predict the unobserved physical constraints
 1347 \mathbf{P} as a function of \mathbf{T} and \mathbf{Y} . We then incorporate these generated constraints as observed
 1348 inputs to aid in the inference of locations. We detail this strategy in Section 3.3.5.

1349 The rest of Section 3.3 is organized as such: Section 3.3.2 translates the defined problem
 1350 into one that can be addressed with Gaussian processes; Section 3.3.3 describes how to
 1351 incorporate $\mathbf{P} \rightarrow \mathbf{Y}$ in the modeling framework; Section 3.3.4 discusses our greedy kernel
 1352 learning approach to learn individual kernel structures; Section 3.3.5 ties together Sections
 1353 3.3.1 through 3.3.4 and describes our model inference strategy.

1354 3.3.2 Multi-task Gaussian Process for Mobile Data

1355 We first focus on modeling the relationship $\mathbf{T} \rightarrow \mathbf{Y}$. Consider the task of learning a function
 1356 $f_j : \mathbb{R}^d \rightarrow \mathbb{R}$ from a training set $\mathcal{D} = (\mathbf{T}, \mathbf{Y})$ where j refers to either latitudes ϕ or longitudes
 1357 λ . The basic form of our learning problem is

$$y_{ji} = f_j(\mathbf{t}_i) + \epsilon_{ji}, \quad (3.2)$$

1358 where f_j is a systematic function mapping inputs \mathbf{t}_i to output y_{ji} , and $\epsilon_{ji} \sim \mathcal{N}(0, \delta_j^2)$ are
 1359 independent random variables for noise associated with the j^{th} task. We place a GP prior on
 1360 f_j such that $f_j \sim \mathcal{GP}(m(\cdot), k(\cdot, \cdot))$, where $m(\cdot) = \mathbb{E}[f_j(\cdot)]$ is the mean function (hereafter set
 1361 to $\mathbf{0}$ assuming proper normalization), and $k(\cdot, \cdot)$ is the covariance (or kernel) function.

1362 **Assumption 1.** *The set $\mathbf{f}_j = [f_j(\mathbf{t}_1), \dots, f_j(\mathbf{t}_n)]^\top$ follows a multivariate normal distribution*
 1363 *such that $p(\mathbf{f}_j | \mathbf{T}) = \mathcal{N}(\mathbf{0}, \mathbf{K}_t)$, where \mathbf{K}_t is the covariance matrix such that $[\mathbf{K}_t]_{i,g} = k_t(\mathbf{t}_i, \mathbf{t}_g)$,*
 1364 *where subscript t is used to indicate temporal correlations and subscript g is an index like*
 1365 *subscript i .*

1366 In our approach, we relate multiple tasks (each multivariate normal distributed) by
 1367 leveraging the correlations between them. Thus, we specify the covariance matrix for all n

1368 observations and two tasks as

$$\mathbf{K}_t = \mathbf{K}^t(\mathbf{T}, \mathbf{T}) \otimes \mathbf{K}^f(\mathbf{y}_\lambda, \mathbf{y}_\phi), \quad (3.3)$$

1369 where \mathbf{K}^t is the $n \times n$ covariance matrix of the training times using any valid PSD kernel, \otimes
 1370 is the Kronecker product, and \mathbf{K}^f is a 2×2 PSD matrix containing the inter-task covariances,
 1371 learned by jointly minimizing the negative marginal log-likelihood (Bonilla et al., 2007). The
 1372 dimension of \mathbf{K}_t for two tasks is then $2n \times 2n$. A key property of this model is that the joint
 1373 distribution over the entire output space is not block-diagonal with respect to tasks. That is,
 1374 observations of one task can affect the predictions on another task.

1375 **Assumption 2.** *The inferred location $f_j(\mathbf{t}_*)$ of a new input \mathbf{t}_* conditioned on the training*
 1376 *data is assumed to be distributed with the following form:*

$$p(f_j(\mathbf{t}_*)|\mathbf{t}_*, \mathbf{T}, \mathbf{Y}, \delta_j^2) \sim N(\boldsymbol{\mu}_*, \boldsymbol{\sigma}_*^2), \quad (3.4)$$

1377 where

$$\begin{aligned} \boldsymbol{\mu}_* &= (\mathbf{k}_j^f \otimes \mathbf{k}_*) (\mathbf{K}^f \otimes \mathbf{K}^t + \mathbf{D} \otimes \mathbf{I})^{-1} \text{vec}(\mathbf{Y}) \\ \boldsymbol{\sigma}_*^2 &= (\mathbf{k}_j^f \otimes \mathbf{k}_{**}) - (\mathbf{k}_j^f \otimes \mathbf{k}_*) (\mathbf{K}^f \otimes \mathbf{K}^t + \mathbf{D} \otimes \mathbf{I})^{-1} (\mathbf{k}_j^f \otimes \mathbf{k}_*). \end{aligned} \quad (3.5)$$

1378 Here, \mathbf{k}_j^f selects the j^{th} column of \mathbf{K}^f , $\mathbf{k}_* = k(\mathbf{t}_*, \mathbf{T})$ is the vector of covariances between
 1379 the test point and the training set, \mathbf{D} is a 2×2 diagonal matrix with the variances of the
 1380 noise processes for latitude and longitude δ_j^2 , and $\mathbf{k}_{**} = k(\mathbf{t}_*, \mathbf{t}_*)$. Finally, we minimize the
 1381 negative marginal log-likelihood (MLL) of the output vectors with respect to the training
 1382 data in determining the optimal hyperparameters.

$$-\log(p(\mathbf{Y}|\mathbf{T}, \Theta)) = \frac{1}{2} [\text{vec}(\mathbf{Y})^\top \boldsymbol{\Sigma}^{-1} \text{vec}(\mathbf{Y}) + \log(\det(\mathbf{K}_t)) + |\Theta| \log(2\pi)], \quad (3.6)$$

1383 where Θ is the set of model parameters, $|\Theta|$ denotes the cardinality of the set of model pa-
 1384 rameters, which include both kernel-specific parameters as well as the weights associated with
 1385 each kernel component (discussed further in Section 3.3.4), $\boldsymbol{\Sigma} = \mathbf{K}^f \otimes \mathbf{K}^t + \mathbf{D} \otimes \mathbf{I}$, and $\det(\mathbf{K}_t)$
 1386 is the determinant of the \mathbf{K}_t matrix. We emphasize that the elements of \mathbf{K}^f are treated as
 1387 free parameters which are learned through Equation 3.6. The negative marginal log-likelihood
 1388 provides a target function for kernel learning. The first component, $\text{vec}(\mathbf{Y})^\top \boldsymbol{\Sigma}^{-1} \text{vec}(\mathbf{Y})$,
 1389 estimates the model fit, while the second and the third terms, $\log(\det(\mathbf{K}_t)) + |\Theta| \log(2\pi)$, act
 1390 as the regularization (Rasmussen and Ghahramani, 2000).

1391 *3.3.3 Physics-regularized GP*

1392 Only using the relationship $\mathbf{T} \rightarrow \mathbf{Y}$ to generate synthetic data can be problematic, as
 1393 unconstrained GPs may produce predictions that are inconsistent with real-world conditions.
 1394 We can indeed regularize the framework we have thus described through measurements of
 1395 velocity and bearing, which are realizations of physics-based constraints acting on every
 1396 user. Incorporating physical knowledge helps reduce the risk of overfitting and improves the
 1397 likelihood of model generalization since the physical knowledge provides data augmentation.
 1398 We derive observations of these variables with

$$v_\lambda = \frac{dy_\lambda}{dt}, \quad v_\phi = \frac{dy_\phi}{dt}, \quad v = \sqrt{v_\lambda^2 + v_\phi^2} \quad (3.7)$$

$$\beta = \arctan\left(\frac{v_\phi}{v_\lambda}\right) = \arctan\left(\frac{dy_\phi}{dy_\lambda}\right) \quad (3.8)$$

1399 where v_λ and v_ϕ denote the velocity in the horizontal and vertical direction, respectively, v
 1400 denotes the overall velocity, and β denotes the bearing.

1401 One issue with using physical variables is the noise that may be introduced by irregular
 1402 sampling of training data. That is, if time gaps between data points are large, the estimated
 1403 segment velocities and bearings may not be representative of the distribution of physics
 1404 within that period. To get around this, Wang et al. (2022b) proposed to sample the physical
 1405 observations from the posterior of the GP modeling $\mathbf{T} \rightarrow \mathbf{Y}$, thereby retaining control over
 1406 where \mathbf{P} is observed by utilizing a continuous function. A drawback of this approach is
 1407 that the resulting likelihood function is difficult to solve, necessitating the calculation of
 1408 the evidence lower bound (ELBO). Additionally, it significantly increases the computational
 1409 demand, raising the complexity of the conditional GP from $O(n^3)$ to $O(n^3 + m^3)$, where m
 1410 represents the number of sampled locations for \mathbf{P} .

1411 Alternatively, Alvarez and colleagues (2013; 2009) propose a different strategy that
 1412 sidesteps the problem by embedding physical knowledge directly into the kernel’s structure,
 1413 rather than incorporating it as part of the training data. While this method simplifies the
 1414 model, it also imposes limitations on the selection of kernels. Specifically, it restricts the
 1415 choice to simpler and smoother kernel functions, which in turn constrains the ability to
 1416 incorporate complex physical knowledge into the model through more sophisticated and
 1417 adaptable kernels.

1418 To mitigate the introduction of noise from data points with extensive time gaps, we simply

1419 assign each data point the same label until the time gap between the next observation exceeds
 1420 a threshold (in our case, we used 1 hour). This essentially labels continuous trajectories
 1421 with the same unique identifier. For the first point of each unique trajectory, we assign a
 1422 velocity and bearing of zero as we cannot observe the physical variables from the most recent
 1423 trajectory (which would be > 1 hour ago). This avoids having to drop training data with
 1424 noisy physical variables without complicating the estimation procedure.

1425 After ensuring the training data is efficiently sampled, adding physical variables to the
 1426 posterior of the temporal GP translates to sampling from a larger covariance matrix that
 1427 also considers physical correlations between the physical variables and their interactions with
 1428 the temporal variables.

1429 **Assumption 3.** Let $\mathbf{X} = \begin{bmatrix} \mathbf{T} & \mathbf{P} \end{bmatrix}$ such that a data vector \mathbf{x}_i has both temporal and physical
 1430 variables. Given the set of physics observations \mathbf{P} , the set $\mathbf{f}_j = [f_j(\mathbf{x}_1), \dots, f_j(\mathbf{x}_n)]$ follows a
 1431 normal distribution such that $p(\mathbf{f}_j|\mathbf{X}) = \mathcal{N}(\mathbf{0}, \mathbf{K}_{comp})$, where \mathbf{K}_{comp} is the covariance matrix
 1432 with entries $[\mathbf{K}_{comp}]_{i,g} = k(\mathbf{x}_i, \mathbf{x}_g)$.

1433 The inclusion of the physical variables is compatible with the multi-task kernel structure,
 1434 as the kernel in Equation 3.3 becomes $\mathbf{K}_{comp} = \mathbf{K}^x(\mathbf{X}, \mathbf{X}) \otimes \mathbf{K}^f(\mathbf{y}_\lambda, \mathbf{y}_\phi)$. This then changes
 1435 the probability of observing a location $f_j(\mathbf{x}_*)$ at unobserved input \mathbf{x}_* to be conditional on
 1436 not only time but physics as well.

1437 **Assumption 4.** The inferred location $f_j(\mathbf{x}_*)$ of a new input \mathbf{x}_* conditioned on the training
 1438 data is assumed to be distributed with the following form:

$$p(f_j(\mathbf{x}_*)|\mathbf{x}_*, \mathbf{T}, \mathbf{P}, \mathbf{Y}, \delta_j^2) \sim \mathcal{N}(\boldsymbol{\mu}_{**}, \boldsymbol{\sigma}_{**}^2), \quad (3.9)$$

1439 where $\boldsymbol{\mu}_{**}$ and $\boldsymbol{\sigma}_{**}^2$ are defined as before but using the larger covariance matrix \mathbf{K}_{comp} .

1440 This process transforms the MLL shown in Equation 3.6 to

$$-\log(p(\mathbf{Y}|\mathbf{T}, \mathbf{P}, \Theta)) = \frac{1}{2}[\text{vec}(\mathbf{Y})^\top \boldsymbol{\Sigma}^{-1} \text{vec}(\mathbf{Y}) + \log(\det(\mathbf{K}_{comp})) + |\Theta| \log(2\pi)]. \quad (3.10)$$

1441 The impact of including physical variables is non-negligible, both for estimation and
 1442 inference (the latter of which we discuss in Section 3.3.5). Only looking at the temporal
 1443 learning problem $\mathbf{T} \rightarrow \mathbf{Y}$ risks letting the GP remain too unconstrained, leading to illogical
 1444 travel behavior predictions or unfeasible routes. For example, a well-fit GP with no physics-
 1445 regularization tends to predict a user hovering around an activity location (rather than staying

stationary while not on a trip), potentially leading to a false "walking trip" interpretation. On the other hand, a well-fit GP with physics-regularization largely solves this issue, as the model is statistically informed on whether the user is moving or remaining stationary.

3.3.4 Greedy Algorithm for Multiple Kernel Learning

Multiple kernel learning comprises two steps: the first step is to select a set of base kernels that are suitable to the context of the study (in this case, it is human mobility patterns) and the second step is to identify the optimal combination of the base kernels.

3.3.4.1 Selecting base kernels

The selection of base kernels is based on two key features of human mobility patterns. The first one relates to that observations of human mobility at nearby time intervals tend to be geographically close to each other. This concept extends Tobler's first law of geography (Tobler, 1970), which states that things in close proximity are more closely related than those far apart. Such dependencies can be captured rather straightforwardly using smooth functions derived from a variety of kernels—in this work, we consider the squared exponential (SE), rational quadratic (RQ), and Matern (MAT) 5/2 kernels:

$$k_{SE} = \eta \exp\left(-\frac{|x_i - x_g|^2}{2l^2}\right), \quad (3.11)$$

$$k_{RQ} = \eta \left(1 + \frac{|x_i - x_g|^2}{2\alpha l^2}\right)^{-\alpha}, \quad (3.12)$$

$$k_{MAT_{5/2}} = \eta \left(1 + \frac{\sqrt{5}|x_i - x_g|}{l} + \frac{5|x_i - x_g|^2}{3l^2}\right) \exp\left(-\frac{\sqrt{5}|x_i - x_g|}{l}\right), \quad (3.13)$$

where $|x_i - x_g|$ represents the Euclidian distance between any pair of inputs x_i and x_g ; η is a scale parameter³; the lengthscale l determines the smoothness of the function; and the scale mixture α determines the relative weight of large- and small-scale variations in the data.

The SE kernel is characterized by infinite differentiability and is defined by a single lengthscale parameter. It is suitable for capturing smooth and gradual variations in the data. The RQ kernel introduces an additional scale mixture parameter that controls the smoothness

³The scaling parameter can also be thought of as a weight in the composite kernel, determining the influence of each base kernel.

1467 of the kernel function and therefore can model both short- and long-range dependencies in
 1468 the data. It can capture various patterns, from smooth variations to abrupt changes, making
 1469 it suitable for datasets with diverse characteristics. Finally, the Matern 5/2 kernel exhibits
 1470 differentiability up to the second order and is often favored when there is prior knowledge
 1471 that the underlying process has a roughness level between that of the SE and the Matern 3/2
 1472 kernels. The Matern 5/2 kernel can capture moderate levels of non-smoothness in the data
 1473 while still maintaining a certain degree of smoothness.

1474 The second feature as noted in Section 3.3.1 is that individuals' travel behavior tends
 1475 to center around a few important locations, such as home and work and there tend to be
 1476 rhythms or cyclical patterns of different lengths with respect to how individuals conduct
 1477 activities. To account for these, we further introduce categorical variables to capture cyclical
 1478 patterns of different lengths. We represent calendar-based structures like days, weeks, and
 1479 months as binary variables using a one-of-k encoding. For example, as the days of the week
 1480 can take one of seven values $\{Mo, Tu, We, Th, Fr, Sa, Su\}$, a one-of-k encoding of *We* will
 1481 correspond to $\{0, 0, 1, 0, 0, 0, 0\}$. We embed multiple categorical inputs in a GP framework by
 1482 multiplying the same kernel across the one-hot encodings:

$$k_{RQ_{cat}} = \prod_{u=1}^d k_{RQ_u}. \quad (3.14)$$

1483 Equation 3.14 is a product kernel containing a unique lengthscale parameter for each input
 1484 dimension being multiplied across the function. A small optimal lengthscale for a categorical
 1485 variable implies a relatively low correlation among data in that category, while a large optimal
 1486 lengthscale implies a high correlation. In addition to the one-hot encoding strategy, we also
 1487 use the periodic kernel (Equation 3.15), which allows GPs to model functions that repeat
 1488 themselves:

$$k_{PER} = \exp\left(-\frac{2 \sin^2(\pi|x_i - x_g|/p)}{l^2}\right), \quad (3.15)$$

1489 where the period length p determines the distance between repetitions of the function. Due
 1490 to its sinusoidal nature, the periodic kernel is a better fit for continuous input dimensions
 1491 that are unbounded in the input space.

1492 3.3.4.2 Learning multiple kernel structures

1493 The basic rationale of MKL builds on the possibility to create composite kernels comprising
 1494 base kernels. So, the task at hand is kernel function discovery, i.e., finding the optimal
 1495 structure for the composite kernel function that comprises a set of base ones. Kernel function
 1496 discoveries have been done in other domains that aim to find concrete analytic forms of the
 1497 kernels (Barla et al., 2003; Gibbs, 1998; Ong et al., 2002; Wilson, 2014), but there have
 1498 been no specialized kernel functions developed for individual mobile data. In this paper, we
 1499 develop a multiple kernel learning algorithm that is data-driven.

1500 A GP’s kernel is only valid if it has a PSD covariance matrix (Rasmussen and Williams,
 1501 2006). The direct sum or the tensor product of two valid kernels is also a valid kernel (for
 1502 proof, we refer the reader to Rasmussen & Williams (2006) section 4.2.4). To put it another
 1503 way, the set of PSD kernels is closed under sum and product operations. Multiplying two
 1504 kernels can be interpreted as an **AND** operation while adding two kernels can be interpreted
 1505 as an **OR** operation (Duvenaud, 2014). Let k_1 and k_2 be kernels which each depend on a
 1506 single input vector, \mathbf{x} and \mathbf{y} , respectively. The product of k_1 and k_2 will result in a prior
 1507 over functions that vary across both \mathbf{x} and \mathbf{y} and hence the function value $f(x_i, y_i)$ is only
 1508 expected to be similar to some other function value $f(x_g, y_g)$ if x_i is close to x_g **AND** y_i is
 1509 close to y_g . On the other hand, the sum $k_1 + k_2$ will result in a prior over functions which are
 1510 a sum of one-dimensional functions, and hence the function $f(\mathbf{x}, \mathbf{y}) = f_x(\mathbf{x}) + f_y(\mathbf{y})$.

1511 Our MKL learning is essentially learning about six components: target function, the base
 1512 learner, the learning method, the functional form, the training method, and the computational
 1513 complexity (Gönen and Alpaydm, 2011). We use the negative marginal log-likelihood
 1514 (Equation 3.6) as our target function to assess the model’s goodness of fit. However, simply
 1515 comparing MLLs after optimizing kernel parameters would result in a bias in favor of more
 1516 complex models. We avoid this overfitting issue by approximating the integral of the marginal
 1517 likelihood over all free parameters using the Bayesian information criterion (BIC) (Schwarz,
 1518 1978):

$$\text{BIC}(n, |\Theta|) = |\Theta| \log(n) - \log(p(\mathbf{Y}|\mathbf{X}, \Theta)). \quad (3.16)$$

1519 Our greedy multi-kernel learning is essentially iteratively adding or multiplying new kernels
 1520 until the BIC is minimized. The functional form of the derived kernel expressions can be both
 1521 linear (i.e., additive) or non-linear (i.e., product)—each component has a weight parameter
 1522 that determines its relevance and proportion of importance. We use the Adaptive Moment

1523 Estimation algorithm (Kingma and Ba, 2014) as the training method, which optimizes
 1524 both the kernel weight parameters $\boldsymbol{\eta}$ and the base kernel parameters simultaneously. The
 1525 computational complexity is a function of the algebraic operations involved, which determines
 1526 the number of times the composite covariance matrix has to be inverted. We initialize the
 1527 algorithm with a set of base kernels B , the maximum number of steps M , and a set of
 1528 algebraic operations A ⁴. This optimization procedure can be formally stated as:

$$\begin{aligned}
 & \underset{\boldsymbol{\eta}, \mathbf{K}, \Theta}{\operatorname{argmin}} && -\log(p(\mathbf{Y}|\mathbf{X}, \Theta, \boldsymbol{\eta}, k_{curr})) + |\Theta|\log(n), \\
 & \text{s.t.} && k_{curr} = \sum_{h=1}^M \eta_h k_h \quad \text{where} \quad k_h = \begin{cases} k_{prev} + k_{base}, & \text{if addition at step } h \\ k_{prev} \times k_{base}, & \text{if multiplication at step } h \end{cases} \\
 & && \sum_{h=1}^M \eta_h = 1, \quad \eta_h \geq 0 \quad \forall h \in 1, \dots, M.
 \end{aligned} \tag{3.17}$$

1529 As illustrated in Figure 3.1, we first estimate the BIC of each kernel in B and choose
 1530 the lowest one. We then begin the greedy procedure: we apply each operation in the set of
 1531 algebraic operations A to the chosen kernel and every available kernel in B (that is, they are
 1532 added or multiplied) and choose the combination with the lowest BIC. This procedure is
 1533 repeated until no new kernels result in a lower BIC than that of the existing kernel or until
 1534 we reach M . The convergence properties of this algorithm depend largely on the complexity
 1535 of the training data—if the data includes a large timeframe and the user displays multiple
 1536 periodicities, trends, or random variations, the algorithm may not converge until we reach M .
 1537 Otherwise, only one or two steps can be sufficient in describing the user’s behavior.

1538 3.3.5 Model Inference

1539 So far, we have dealt with the estimation of the posterior for the GP model. However, inference
 1540 using physical variables presents a significant obstacle as it is a random variable determined
 1541 by both temporal and spatial factors (locations), the latter of which are unobserved. Thus,
 1542 prior to generating a set of synthetic locations \mathbf{Y}_{gen} using the conditional model, we must
 1543 specify a model to predict the unobserved physical constraints \mathbf{P}_{gen} as a function of time
 1544 and space. For the rest of the paper, \mathbf{T} , \mathbf{P} , and \mathbf{Y} will be represented with \mathbf{T}_{cond} , \mathbf{P}_{cond} , and

⁴For now, we only allow addition and multiplication as other operations do not always result in PSD kernels.

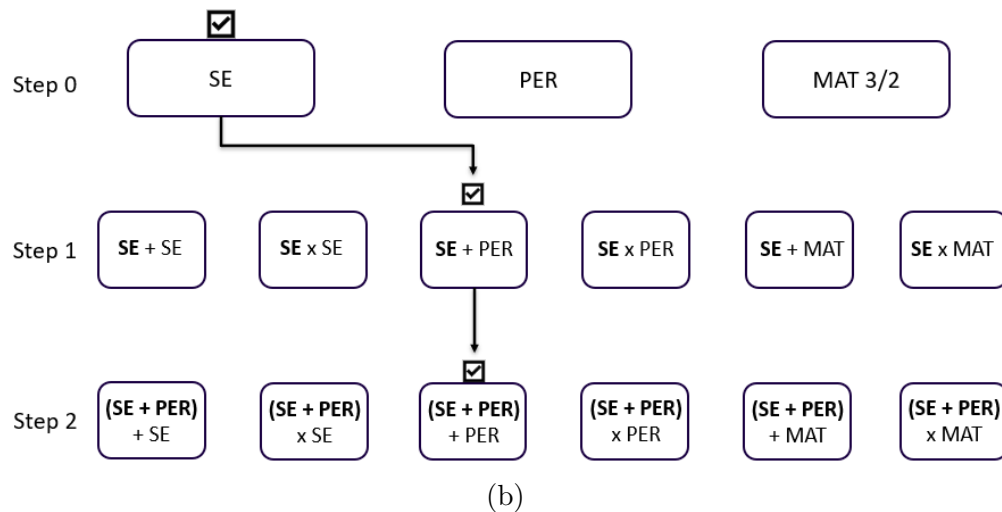
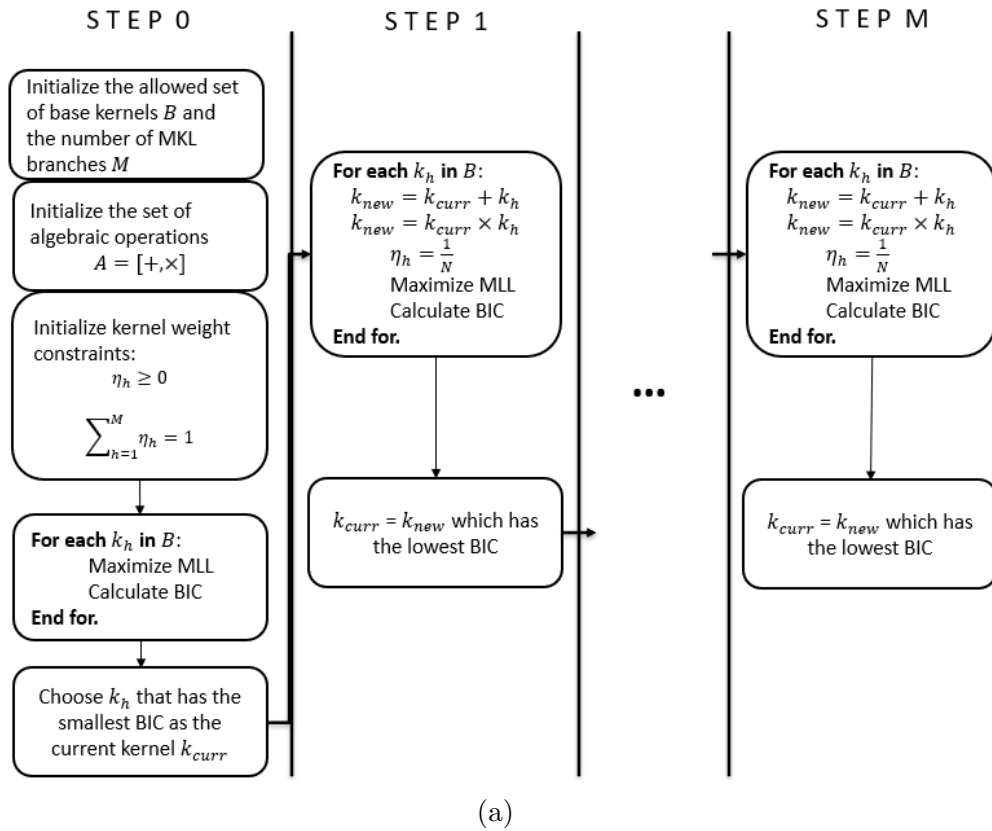


Figure 3.1: (top) Proposed algorithm for greedy multiple kernel learning; (bottom) example greedy learning tree with three steps (checkmark denotes lowest BIC option)

1545 \mathbf{Y}_{cond} to distinguish the fact that they belong to the conditional model, while sets relating to
 1546 the generative model will have the subscript "gen".

1547 Inspired by Lasserre et al. (2006), we propose a conditional-generative inference model.
 1548 The generative component is for realizations of physical variables beyond the training data.
 1549 Unlike LFM and the physics-regularized GP proposed by Wang et al. (2022b), the physical
 1550 variables in our method are determined partially by the output of the conditional model
 1551 (coordinates). Therefore, we first generate a set of noisy locations $\bar{\mathbf{Y}}_{gen} = [\bar{\mathbf{y}}_{gen,1}, \dots, \bar{\mathbf{y}}_{gen,m}]^\top$
 1552 induced by a set of times \mathbf{T}_{gen} using the conditional multi-task GP $f^{\bar{y}} \sim \mathcal{GP}(\mathbf{0}, \mathbf{K}_t)$ via
 1553 Equations 3.4 and 3.5, where $[\mathbf{K}_t]_{i,g} = k_t(\mathbf{t}_i, \mathbf{t}_g)$.

1554 To estimate \mathbf{P}_{gen} , we require a second posterior that takes in spatial and temporal obser-
 1555 vations $\mathbf{Z} = \begin{bmatrix} \mathbf{Y}_{cond} & \mathbf{T}_{cond} \end{bmatrix}$ and approximates a function $f^p : \mathbf{Z} \rightarrow \mathbf{P}_{cond}$. This is achieved
 1556 by defining another multi-task GP ("hidden GP" in Figure 3.2) for the physical variables
 1557 $f^p \sim \mathcal{GP}(\mathbf{0}, \mathbf{K}_p)$ where $[\mathbf{K}_p]_{i,g} = k_{comp}(\mathbf{z}_i, \mathbf{z}_g)$. Sampling from the posterior distribution

$$\mathbf{P}_{gen} \sim p(\mathbf{T}_{gen})p(f^p | \bar{\mathbf{Y}}_{gen}, \mathbf{Y}_{cond}, \mathbf{T}_{cond}, \mathbf{P}_{cond}) \quad (3.18)$$

1558 where we use the marginal independence of \mathbf{T}_{gen} to simplify the joint conditional distribution.
 1559 The set of generated physical variables are then incorporated as inputs to the physics-
 1560 regularized GP model described in Section 3.3.3. Finally, we sample from the physics-
 1561 regularized GP $f^y \sim \mathcal{GP}(\mathbf{0}, \mathbf{K}_{comp})$ where $[\mathbf{K}_{comp}]_{i,g} = k_{comp}(\mathbf{x}_i, \mathbf{x}_g)$.

1562 Figure 3.2 shows the workflow for the proposed model, which works in two passes. On
 1563 the first pass, the physical knowledge is embedded as a function of time and space using a
 1564 GP ("Physics Embedding"). On the second pass, synthetic data points are generated given
 1565 \mathbf{T}_{gen} and the estimated physical variables \mathbf{P}_{gen} .

1566 **Theorem 1.** *The empirical distribution functions used for sampling the generated synthetic*
 1567 *data \mathbf{Y}_{gen} and \mathbf{P}_{gen} are constructed only using the basis vectors generated from the conditional*
 1568 *and hidden GPs, respectively.*

1569 *Proof.* Let $F_{\mathbf{Y}_{gen}}$ and $F_{\mathbf{P}_{gen}}$ denote the empirical distribution functions for \mathbf{Y}_{gen} and \mathbf{P}_{gen} ,
 1570 respectively. Due to Mercer's theorem⁵, each sample of \mathbf{P}_{gen} and \mathbf{Y}_{gen} is a linear combination
 1571 of the basis functions determined by the covariance matrices \mathbf{K}_p and \mathbf{K}_{comp} . Because $F_{\mathbf{Y}_{gen}}$
 1572 and $F_{\mathbf{P}_{gen}}$ are constructed from these samples, they will also be linear combinations of these
 1573 basis vectors generated from the conditional and hidden GPs. \square

⁵For details, see Section 4.3 of Rasmussen & Williams (2006)

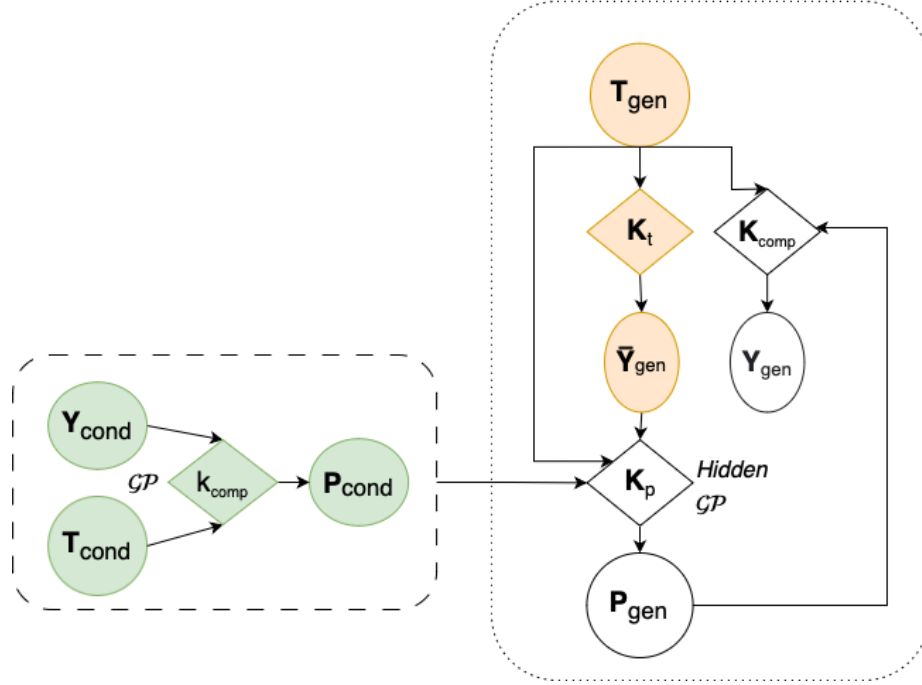


Figure 3.2: Physics-regularized GP inference for passively-generated mobile data. Circles indicate variables, while diamonds indicate estimation (kernel function) or inference (kernel matrix) processes.

1574 The implication of Theorem 1 is three-fold. First, the generated \mathbf{P}_{gen} and \mathbf{Y}_{gen} reflect
 1575 the patterns, trends, and noise characteristics present in the observed data \mathbf{P}_{cond} and \mathbf{Y}_{cond} .
 1576 Second, the model is not just statistically informed but also physically informed. This means
 1577 that any generated data point respects the underlying physical laws or constraints that
 1578 have been encoded through the model framework. These constraints might represent known
 1579 relationships between variables that must be maintained, such as the laws of motion in a
 1580 physical system. Third, by adhering to the learned relationships and physical constraints, the
 1581 synthetic data generation process is not just mimicking the observed data but generalizing
 1582 from it. This enables the model to create data points that are plausible under the model's
 1583 understanding of the domain, even if they were not explicitly present in the training data.

1584 **3.4 Numerical Experiments**

1585 In this section, we document the results of four experiments that evaluate the formulation
 1586 we have thus presented. In Sections 3.4.1 and 3.4.2, we discover the best-fitting covariance
 1587 structure for temporal and physical features to create a composite kernel, respectively. In
 1588 Section 3.4.3, we compare the performance of this composite kernel against alternative
 1589 formulations over many users in a passively-collected mobile dataset. Finally, to analyze
 1590 our framework’s generalizability, we conduct another set of experiments in Section 3.4.4 on
 1591 Microsoft Asia’s GeoLife dataset using trips of varying modes. All class and function objects
 1592 relevant to this experiment can be found on GitHub: [https://github.com/ekinugurel/](https://github.com/ekinugurel/physics-regularized-MTGP)
 1593 [physics-regularized-MTGP](https://github.com/ekinugurel/physics-regularized-MTGP). Additionally, Appendix B.6 outlines a set of algorithms we
 1594 use throughout the experiments.

1595 Before all experiments, we perform data standardization (z-score normalization) such that
 1596 each feature is centered at 0 with a standard deviation of 1. For example, $v_i \in \mathbf{v}$ (the i -th
 1597 observation of velocity in \mathbf{v}) can be standardized with

$$\tilde{v}_i = \frac{v_i - \bar{v}}{s_{\mathbf{v}}} \quad (3.19)$$

1598 where \bar{v} is the sample’s average velocity and $s_{\mathbf{v}}$ is the standard deviation of the sample’s velocity.
 1599 As different variables have different units, doing so simplifies the process of optimization via
 1600 gradient descent—allowing the algorithm to converge more efficiently by ensuring that the
 1601 scales of all features are comparable. This standardization also aids in preventing certain
 1602 variables with larger magnitudes from disproportionately influencing the optimization process,
 1603 promoting a more balanced and effective model training.

1604 *3.4.1 Temporal Structure Discovery using Multiple Kernel Learning*

1605 We show the process of discovering the temporal structure associated with an individual’s
 1606 mobile data. Specifically, we conduct a test on data over three weeks, using the first two
 1607 weeks for training the model, and measure its accuracy using the last week. We run the
 1608 MKL algorithm on the training set with the following base kernels: SE (Equation 3.11), PER
 1609 (Equation 3.15), and RQ (Equation 3.12). We initiate two instances of the PER kernel: one
 1610 initialized with $p = 1,440$ (representing a day) and the other initialized with $p = 10,080$
 1611 (representing a week). After a sensitivity analysis, we settle on $M = 3$ (as more steps do not
 1612 result in lower BICs) and initialize the kernel weight constraints. The resulting kernel is a

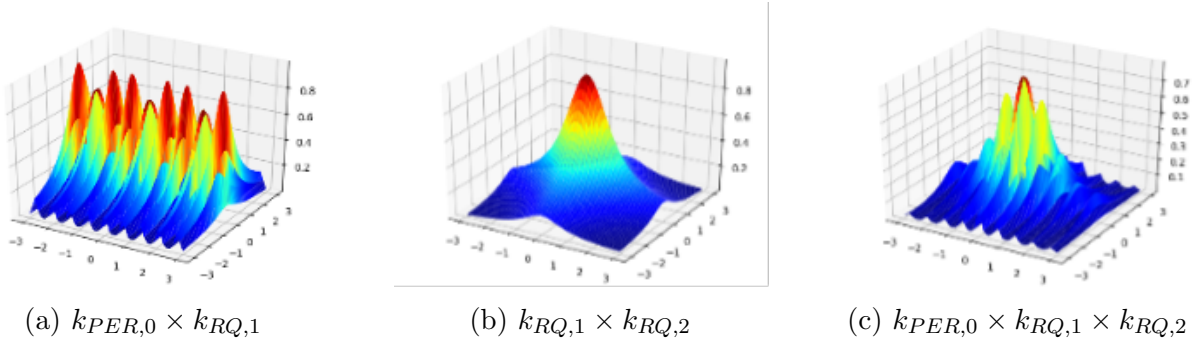


Figure 3.3: Components of the single-task temporal kernel in a 3-D space

1613 set of rational quadratic kernels multiplied on all input dimensions and one periodic kernel
 1614 on the Unix time dimension. Figure 3.3 visualizes components of the discovered kernel in
 1615 a single-task space. To evaluate the success of our algorithm, we compare this kernel with
 1616 similar kernels that can be constructed using the above base kernel set.

1617 We measure model accuracy with three mobility metrics: spatial burstiness B_s , represent-
 1618 ing the distribution pattern of travel distances between consecutive stay points (Kim and
 1619 MacEachren, 2014); radius of gyration r_g , indicating the characteristic travel distance of a
 1620 user during a period (Song et al., 2010b); and home location inaccuracy ΔH , measuring the
 1621 haversine distance between user’s imputed and actual home location. Our results are detailed
 1622 in Figure 3.4, while the optimal kernel parameters are shown in Table 3.1. From Figure
 1623 3.4a, the convergence behavior of each kernel is discernible, with most of them stabilizing as
 1624 the number of iterations increases, indicating that the models are reaching an optimal point
 1625 in terms of the parameters being estimated. A few kernels exhibit some volatility in their
 1626 loss values, even in later iterations, which could suggest a less stable model fit. Figure 3.4b
 1627 shows the generated traces using the discovered kernel and the temporal GP based on the
 1628 training data. Sigma 1 and Sigma 2 represent 1 and 2 standard deviations in the normalized
 1629 distribution. The coefficients of variation using de-normalized generated traces are 0.127 for
 1630 the latitudes and 0.026 for the longitudes, respectively, suggesting that GP does a good job
 1631 in estimating the variance of the predictions.

1632 Figure 3.4c suggests that there is a trade-off in the performance of a composite kernel
 1633 depending on the mobility metric being prioritized and that the impact of the chosen kernel
 1634 is non-negligible in this context. Notably, while certain kernels may not be optimal for all

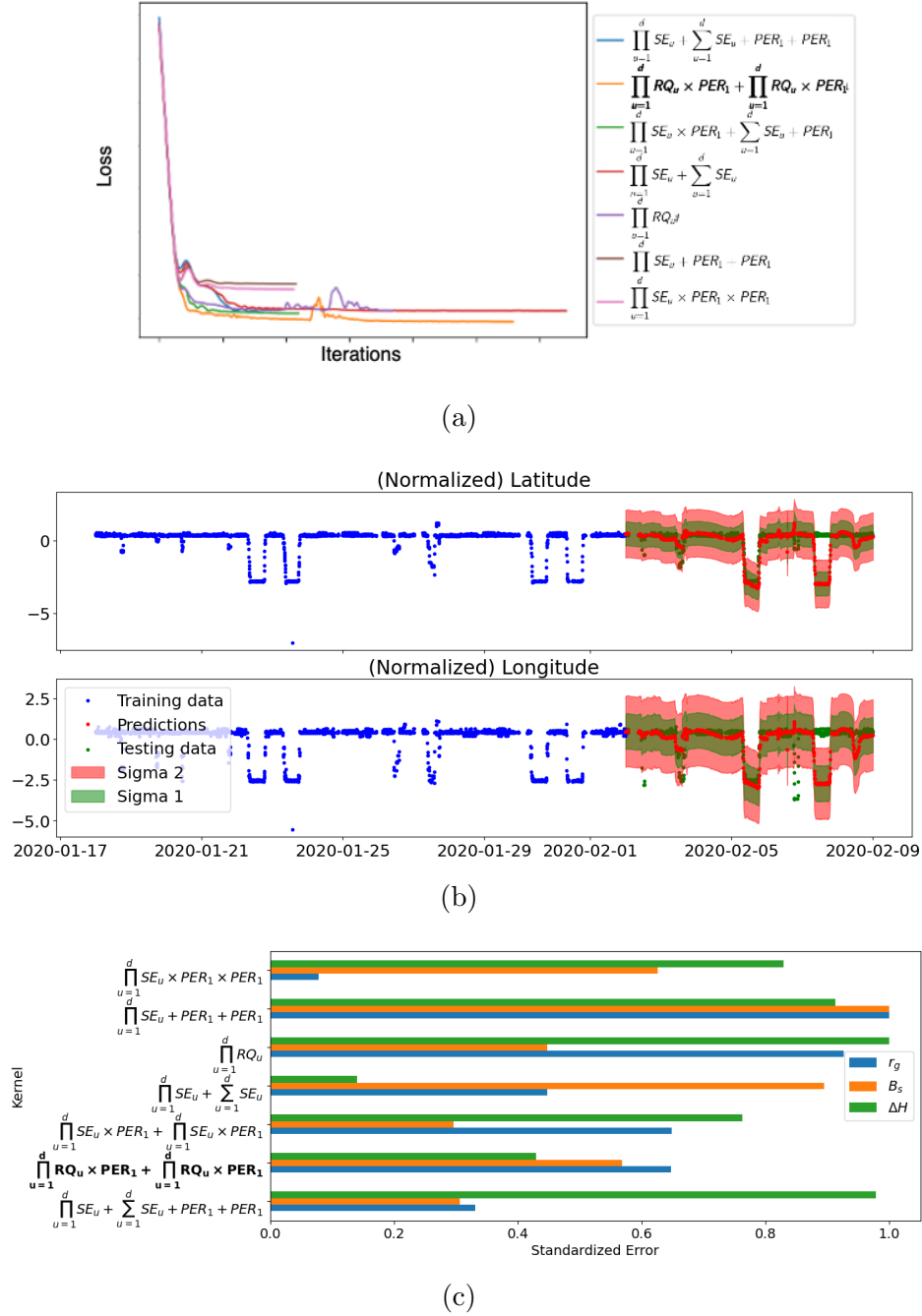


Figure 3.4: (a) Loss function progression and convergence behavior of different functional forms for a composite kernel; (b) Generating one week of trips using discovered kernel and the temporal GP; (c) Derived mobility metrics by kernel. Bold denotes the optimal kernel. r_g denotes radius of gyration, B_s denotes spatial burstiness, and ΔH denotes distance between user's imputed and actual home location.

Table 3.1: Optimal kernel parameters for long gap imputation example. Asterisk* denotes the highest lengthscale among categorical/binary variable.

Parameter	Value	Explanation
Initial $p_{a,1}$	1,440 mins = 1 day	Initial value for period length in PER_a
Initial $p_{b,1}$	10,080 mins = 1 week	Initial value for period length in PER_b
η_a	0.27	Optimal weight of kernel component a
η_b	0.73	Optimal weight of kernel component b
$p_{a,1}$	629 mins = 0.44 days	Optimal period length of PER_a
$p_{b,1}$	5,290 mins = 3.67 days	Optimal period length of PER_b
$l_{a,RQ,t_d=3}$	39.4*	Lengthscale of RQ_a for $t_d = 3$
$l_{a,RQ,t_{am}=1}$	18.4*	Lengthscale of RQ_a for $t_{am} = 1$
$l_{b,RQ,t_d=6}$	55.6*	Lengthscale of RQ_b for $t_d = 6$
$l_{b,RQ,t_{am}=1}$	18.7*	Lengthscale of RQ_b for $t_{am} = 1$

1635 metrics, they still provide some degree of insight into the mobility patterns. A kernel that
 1636 may have a higher loss or error for one metric might still be useful for understanding certain
 1637 aspects of mobility through other metrics. This may be because the mobility metrics analyzed
 1638 here are not necessarily aligned with the objective function of the MKL process. We discuss
 1639 ideas for future studies on this topic in Section 3.5.

1640 In practical applications, composing a kernel would then depend on the specific aspect of
 1641 mobility one aims to capture. If the goal is to have a balanced model across various metrics,
 1642 then kernels that show moderate performance across the board might be preferred. However,
 1643 if the objective is to excel in one particular metric, such as predicting home location with high
 1644 accuracy, then the kernel with the lowest error for ΔH would be chosen, despite potential
 1645 compromises in other areas.

1646 Analyzing optimal kernel parameters highlights the capacity of our model to uncover
 1647 individual mobile data generation processes. The last four rows in Table 3.1 show the highest
 1648 characteristic lengthscale along each categorical dimension—a high smoothing parameter
 1649 for an input dimension suggests low variation along that dimension in the function being
 1650 modeled, and vice versa. We notice that the example user has high correlations in their
 1651 behavior on Thursdays ($t_d = 3$), Sundays ($t_d = 6$) and in the AM peaks ($t_{am} = 1$) compared
 1652 to the other days of the week and the PM peak, respectively. The optimal period lengths also
 1653 suggest that the user’s periodicities occur at shorter intervals than initially assumed—the
 1654 first periodic kernel is optimized at 0.44 days and the second at 3.67 days. Relatedly, the

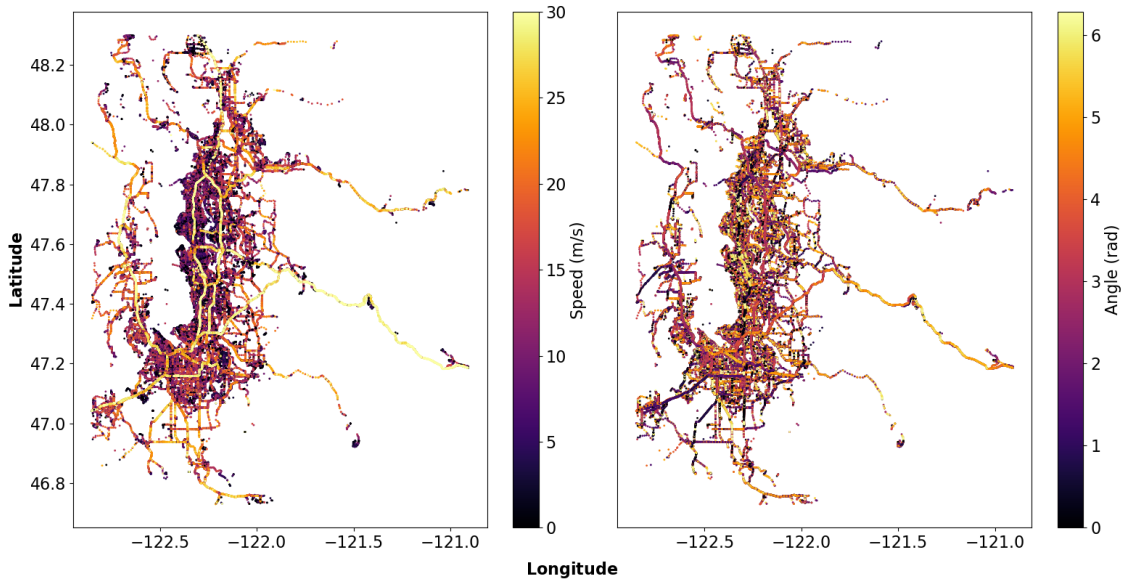


Figure 3.5: Speed and Bearing Constraints in King County, Washington

1655 daily kernel component influences the model less than the weekly component judging by the
 1656 weights η_a and η_b .

1657 3.4.2 Learning the Physical Kernel

1658 In this section, we discuss the process of learning the optimal physical kernel k_p . Figure
 1659 3.5 shows discrete observations of average segment speeds and angular directions observed
 1660 from 10 users' trips over six months. We note that one can deduce roadway types by the
 1661 average speed observed on them (e.g., the clear vertical lines on the left are SR-99 and I-5
 1662 respectively).

1663 To learn the form of the covariance function for physics-based variables, we apply the
 1664 greedy MKL algorithm on the training set with the following base kernels: SE (Equation 3.11),
 1665 and RQ (Equation 3.12), MAT (Equation 3.13). This is one realization of the $\mathbf{P}_{cond} \rightarrow \mathbf{Y}_{cond}$
 1666 GP shown on the left side of Figure 3.2. The resulting kernel (Equation 3.20) has an additive
 1667 form and uses the RQ (Equation 3.12) and MAT 5/2 (Equation 3.13) kernels. Figure 3.6
 1668 shows predictions on one week of velocities and bearings after the covariance matrix associated
 1669 with Equation 3.20 was estimated for two weeks of training data.

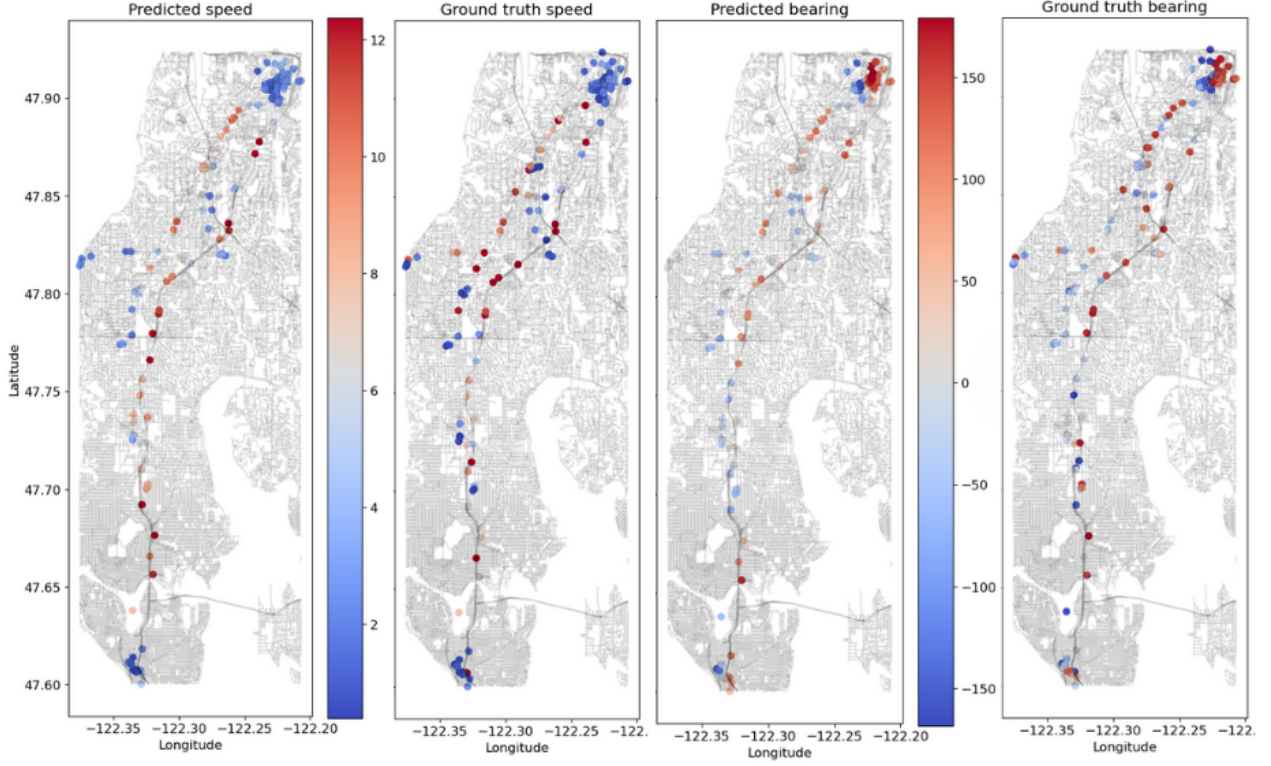


Figure 3.6: Speed and Bearing Predictions and True Observations for a Spectus User

$$\begin{aligned}
 k_p = & \eta_{RQ,1} \left(1 + \frac{|\mathbf{x}_i - \mathbf{x}_g|^2}{2\alpha_1 l_{RQ,1}^2} \right)^{-\alpha_1} \\
 & + \eta_{MAT,1} \left(1 + \frac{\sqrt{5}|\mathbf{x}_i - \mathbf{x}_g|}{l_{MAT,1}} + \frac{5|\mathbf{x}_i - \mathbf{x}_g|^2}{3l_{MAT,1}^2} \right) \exp \left(-\frac{\sqrt{5}|\mathbf{x}_i - \mathbf{x}_g|}{l_{MAT,1}} \right) + \eta_{RQ,2} \left(1 + \frac{|\mathbf{x}_i - \mathbf{x}_g|^2}{2\alpha_2 l_{RQ,2}^2} \right)^{-\alpha_2} \\
 & + \eta_{MAT,2} \left(1 + \frac{\sqrt{5}|\mathbf{x}_i - \mathbf{x}_g|}{l_{MAT,2}} + \frac{5|\mathbf{x}_i - \mathbf{x}_g|^2}{3l_{MAT,2}^2} \right) \exp \left(-\frac{\sqrt{5}|\mathbf{x}_i - \mathbf{x}_g|}{l_{MAT,2}} \right).
 \end{aligned} \tag{3.20}$$

1670 3.4.3 Passively-collected Mobile Dataset

1671 We randomly select a subset of 33 users with at least 10,000 observations and 3 months of
 1672 data (from first data point until last). Within the selected 33, we only retain data from
 1673 the month with the most observations. We minimize Equation 3.10 using the mean and
 1674 covariance structures in Equation 3.9 to evaluate the performance of the physics-regularized
 1675 GP formulation against two alternative formulations: One using only physical variables, and

1676 another using only temporal variables. We note that the former is unrealistic in practice, as
 1677 variables like speed and bearing have to be derived from coordinates (which are assumed
 1678 missing in our tests). However, we show the results as they illuminate the motivation for
 1679 integrating the two types of variables.

1680 We use the root mean square error (RMSE) and the mean standardized log loss (MSLL)
 1681 as our performance metrics—the former is the root of the squared residual between the mean
 1682 prediction and the target, averaged over the test set; the latter is a probabilistic measure that
 1683 takes into account the uncertainty of the predictions. The MSLL incorporates the logarithmic
 1684 transformation of the predicted and observed values, normalized by the associated uncertainty.
 1685 This allows for a more nuanced assessment of the model’s performance, as it captures the
 1686 relative differences between predicted and observed values in a probabilistic manner. We use
 1687 these metrics to compare our predictions to various testing sets, allowing us to comment on
 1688 the model’s generalizability and the impact of the physics-regularization on bias and variance.

$$RMSE = \sqrt{\frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2}, \quad (3.21)$$

$$MSLL = \frac{1}{n} \sum_{i=1}^n \left(\frac{\log(\hat{y}_i + 1) - \log(y_i + 1)}{\sigma_i} \right)^2. \quad (3.22)$$

1689 where y_i and \hat{y}_i are respectively the true and predicted values of the dependent variable at
 1690 index i and σ_i is the predicted standard deviation.

1691 We remove consecutive chunks of data to test the model’s accuracy in imputing time
 1692 series, progressively increasing the size of the training set in each iteration while keeping the
 1693 size of the test set constant. Specifically, we divided the dataset into three sequential folds for
 1694 training and testing, ensuring chronological order is maintained. The first split (S1) uses the
 1695 first and the second weeks of the month as training and testing sets, respectively; the second
 1696 set (S2) then uses the first two weeks to train and the third week to test. Finally, the third
 1697 set (S3) uses the first three weeks to train and the last week to test. This approach simulates
 1698 a real-world scenario where a model is trained on an increasing amount of historical data and
 1699 tested on the subsequent period. It also allows us to observe the performance of the model as
 1700 it is exposed to more data over time.

1701 Table 3.2 shows the median performance of each model across the three splits. The
 1702 Physics-regularized GP outperforms the individual physics-only and temporal-only models

Table 3.2: Models and median test set metrics

Models	Metrics								
	RMSE			MSLL			RT		
	S1	S2	S3	S1	S2	S3	S1	S2	S3
Physical GP	0.105	0.128	0.130	0.767	0.816	0.855	6.69	13.9	30.3
Temporal GP	0.107	0.119	0.116	0.885	0.808	0.874	7.99	17.4	38.5
Physics-reg. GP	0.103	0.119	0.113	0.882	0.779	0.849	11.5	22.8	50.2

Note: S1, S2, and S3 represent three different time-series splits.

1703 in each testing chunk with respect to RMSE and MSLL. The increase in RMSE from S1 to
 1704 S2 for all models and from S2 to S3 for the Temporal GP indicates a decrease in predictive
 1705 accuracy as the training set size increases, possibly due to overfitting or increased model
 1706 complexity. However, the Physics-regularized GP maintains consistently low RMSE values
 1707 across all splits, demonstrating better generalization as more data is included. The steady
 1708 performance improvement in this model suggests that the integration of both physical and
 1709 temporal inputs helps to counteract overfitting and improve predictive accuracy as the dataset
 1710 grows.

1711 The MSLL is lower for the Physics-regularized GP in two out of the three splits, indicating
 1712 not only better prediction accuracy but also better calibrated uncertainty estimates. The
 1713 consistent performance of the Physics-regularized GP across the splits in terms of MSLL
 1714 suggests that it remains reliable and well-calibrated regardless of the amount of training
 1715 data. Although the Physics-regularized GP estimates more input dimensions simultaneously
 1716 than the other two models, its Kronecker product structure keeps the increase in training
 1717 time relatively modest. The physical-only model has the fastest runtime, but the added
 1718 computational cost of the Physics-regularized GP is justified by its improvements in both
 1719 prediction and uncertainty calibration.

1720 3.4.4 *GeoLife Dataset*

1721 To test the generalizability of our framework across multiple modes and populations, we
 1722 conduct an experiment on Microsoft Asia’s GeoLife dataset (Zheng et al., 2008, 2010, 2009),
 1723 specifically using the subset of all GeoLife users who had mode labels available. The essence
 1724 of this experiment is to assess how well our model can generate unobserved trips based on
 1725 a training set of similar trips. Rather than a time-series-based experimental design where

1726 the training and testing sets are temporally sequential, we use a stratified sampling strategy
 1727 (shown within Algorithm 7 in Appendix B.6) applied to all of one user’s trips. That is, we
 1728 cluster trips based on their start and end locations, and sample out of the set of similar trips
 1729 to build the training set. We do this for two reasons: (1) in GeoLife, users may appear in
 1730 the dataset over the course of years and data continuity is low (i.e., we do not see trips from
 1731 the same user in consecutive days), making a time-series split more difficult; (2) Ensuring
 1732 that the training set of trips is similar to the testing set allows us to better isolate the effect
 1733 of physics-regularization, as a pair of start/end locations can have multiple routes (with
 1734 presumably different underlying physics).

1735 We assess the proposed exact GP formulation (PIMTGP) against three alternatives: a
 1736 sparse GP model (SparsePIMTGP), an exact temporal GP using multiple kernel learning
 1737 but without physics-regularization (TempMTGP), and a Long Short-Term Memory (LSTM)
 1738 network (Hochreiter and Schmidhuber, 1997). Unlike probabilistic models, LSTM networks
 1739 do not inherently provide measures of uncertainty. To facilitate a comparison between
 1740 LSTM predictions and that of the GPs, we conduct Monte Carlo (MC) simulations with
 1741 a dropout rate of 0.5 across 70 simulations during inference. This approach introduces
 1742 randomness during inference to approximate the model estimation uncertainty of the LSTM.
 1743 It is important to note that this does not equate to the prediction uncertainty provided by the
 1744 GP models, which is derived from their probabilistic nature. Therefore, while MC dropout
 1745 helps in assessing the robustness of the LSTM predictions, the comparison in probabilistic
 1746 measures between LSTMs and GPs should be interpreted with caution as they reflect different
 1747 aspects of uncertainty. We use the Torch implementation of LSTM with a mean squared
 1748 error loss, optimized via Adam (Kingma and Ba, 2014) as with all other models. We choose
 1749 hyperparameters via cross-validation and limit the number of MC simulations to 70 to
 1750 maintain scalability.

1751 The key idea of the sparse GP model is that we can learn a memory-efficient representation
 1752 of the $n \times n$ covariance matrix using c inducing points, which capture most of the variation
 1753 in the dataset (where $c \ll n$) (Snelson and Ghahramani, 2005). The benchmark formulation
 1754 uses a variational approximation where the predictive distribution is given by

$$p(\mathbf{f}(\mathbf{x}_*)) = \int_{\mathbf{o}} p(f(\mathbf{x}_*|\mathbf{o}))q(\mathbf{o})d\mathbf{o} \quad (3.23)$$

1755 where \mathbf{o} represents the function values at the c inducing points, $q(\mathbf{o}) = \mathcal{N}(\mathbf{m}, \mathbf{S})$. Here,

1756 $\mathbf{m} \in \mathbb{R}^c$ and $\mathbf{S} \in \mathbb{R}^{c \times c}$ are learnable parameters. Its multitask extension leverages the linear
 1757 model of coregionalization (LMC) which assumes that each task j is the linear combination
 1758 of some latent functions $\mathbf{g}(\cdot) = [g^{(1)}(\cdot), \dots, g^{(Q)}(\cdot)]$ and

$$f_j(\mathbf{X}) = \sum_{v=1}^Q a^{(v)} g^{(v)}(\mathbf{X}), \quad (3.24)$$

1759 where $a^{(v)}$ are learnable parameters.

1760 Figure 3.7 shows the performance of each model across seven evaluation metrics and three
 1761 modes of travel—bike, walk, and bus. RMSE and MSLL are defined as before. The mean
 1762 absolute error (MAE) and the median absolute deviation (MAD) are defined as follows,

$$MAE = \frac{1}{n} \sum_{i=1}^n |y_i - \hat{y}_i|, \quad (3.25)$$

$$MAD = \text{median}(|y_i - \hat{y}_i|). \quad (3.26)$$

1763 We also assess model performance with metrics specific to curves (or trajectories). These
 1764 metrics better capture differences in predicted and actual route choice, an important variable
 1765 in travel behavior research. The curve length (CL) measures the arc-length distance along
 1766 the curve from the origin, capturing total travel distance in a given trip. Meanwhile, dynamic
 1767 time warping (DTW) incorporates timestamps of trajectories when considering sequence
 1768 alignment, capturing differences in speed of travel (Ugurel et al., 2024a). A lower DTW value
 1769 indicates a closer match whereas a higher value suggests a greater disparity between the
 1770 compared sequences. Finally, we compare the training runtime (RT) of all methods.

1771 From Figure 3.7, we observe a few things: PIMTGP tends to outperform the other
 1772 models in the classical error metrics (RMSE, MAE, and MAD) across all modes. In the
 1773 probabilistic MSLL metric, we observe a clear distinction between GP-based models and the
 1774 LSTM network. The MKL-equipped TempMTGP seems to be competitive with PIMTGP in
 1775 this probabilistic measure, highlighting the importance of using Multiple Kernel Learning.
 1776 In the trajectory comparison metrics (CL and DTW), the outcomes are fairly similar, with
 1777 TempMTGP and PIMTGP tending to perform slightly better. Finally, taking a look at
 1778 the runtimes paints a clear picture—the more complex models take longer to train but also
 1779 tend to perform better at individual data generation. It is also worthy to note that the
 1780 SparsePIMTGP has similar runtimes to the LSTM model while maintaining the core of our

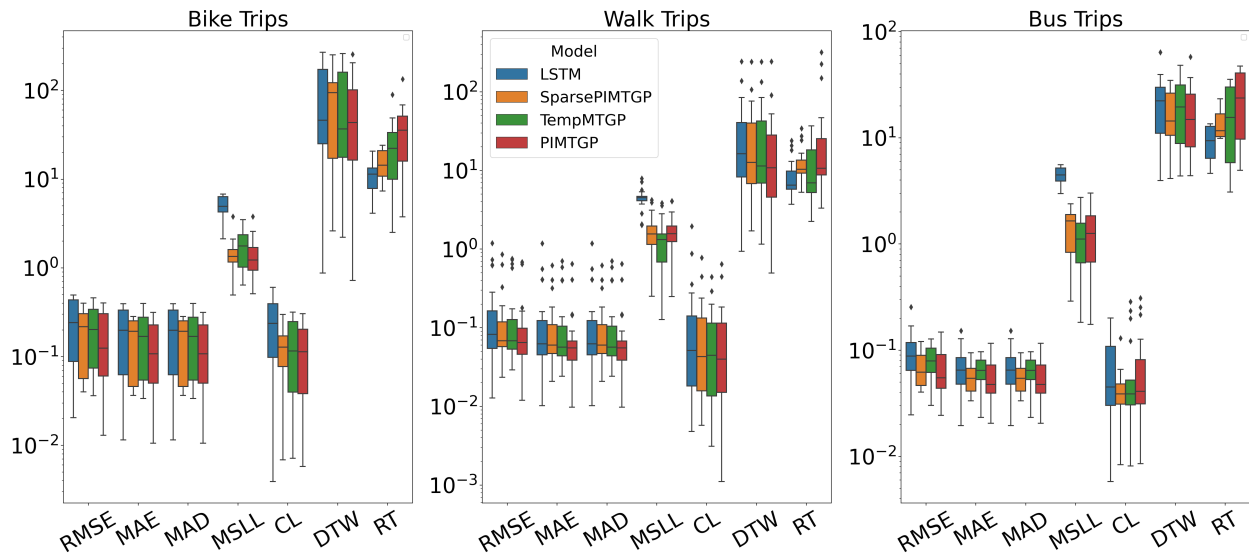


Figure 3.7: Evaluation metrics for all models being compared. Y-axis is shown in log-scale. The lower the better for all metrics.

1781 framework, opening a path to future scalable implementations. We believe this is a promising
 1782 research avenue.

1783 3.5 Discussion

1784 To realize the potential of creating synthetic population for large scale simulations of individual
 1785 mobility patterns, we need to develop a flexible and scalable model that can capture individual
 1786 heterogeneity and also adapt to changes over time. These two properties: flexible and scalable,
 1787 require nonparametric models whose forms can easily adapt. In this paper, we propose a
 1788 physics-regularized multi-task Gaussian Process (GP) model for learning kernel structures that
 1789 characterize individual movement patterns. Though there are other choices for nonparametric
 1790 models such as deep neural networks, GP models are another promising method that offers
 1791 several benefits over deep learning-based models. First, GPs are generally more interpretable
 1792 than deep learning architectures, as choices of basic kernels are based on the nature of the
 1793 phenomenon of interest (in this context, human mobility patterns) and their hyperparameters
 1794 such as lengthscale are explainable (they directly tell us the periodicity of the underlying
 1795 function). A related second point is that kernel methods are more mathematically mature
 1796 than their deep learning counterparts. Especially in exact GPs, all calculations are done in

1797 closed-form—this signals a potential for theoretical bounds on generalization. Additionally,
1798 the probabilistic framework of GPs facilitates robust uncertainty quantification, motivating
1799 the use of different metrics to capture this added benefit (Mohamed et al., 2022).

1800 As noted earlier, the study is motivated by an emerging application of passively-generated
1801 mobile data that will have great potential when realized, a problem that can hamper the
1802 great potential, and the fact that few have looked into the problem at hand. That emerging
1803 application is synthetic population generation for large-scale simulations of cities for various
1804 needs such as state monitoring, policy formulation and evaluations, and building digital twins.
1805 Existing methods such as activity-based models or probability models are inflexible, unable
1806 to capture heterogeneity exhibited in individual mobility patterns and adapt to changing
1807 patterns over time. The reality now and in the future is that the quality of data is degrading
1808 (substantial sparsity or missingness associated with the data) with rising awareness to preserve
1809 privacy on the part of both data vendors and consumers. To fulfill the great potential that
1810 passively-generated mobile data has for synthetic population generation and to address the
1811 increasing issue of sparsity, an understanding of the underlying data generation process
1812 behind the observed individual mobile data is needed. This is the key theme of our paper:
1813 the development of a flexible modeling framework that can adapt to heterogeneity at the
1814 individual level.

1815 We summarize our contributions as follows: We propose a conditional-generative GP
1816 framework to generate synthetic individual mobile data that integrates physical knowledge and
1817 provably replicates observed travel patterns. To account for heterogeneity at the individual
1818 level, we develop a data-driven multiple kernel learning approach to determine the most
1819 suitable spatiotemporal kernel for each user. The learned multiple kernel GP models can be
1820 viewed as mechanisms to uncover underlying data generation processes that result in the
1821 observed big mobile data, which in turn can be used to generate synthetic data mimicking
1822 the properties of collected data.

1823 Moving forward, our research opens several avenues for advancement. Future studies can
1824 aim to explore alternative formulations for the MKL objective function that allow for the
1825 prioritization of different mobility metrics in the synthetic data generation process, tailoring
1826 the model to specific applications. Furthermore, new methods in one-step ahead forecasting
1827 strategies (i.e., predicting the next data point in a time series based on the observed values
1828 up to the current time point) infused with our framework could provide real-time insights
1829 and predictions, enhancing the model’s applicability in dynamic environments. Finally, to

1830 address the computational burden often associated with GP models, studies should further
1831 investigate approximation techniques (as we do with the SparsePIMTGP) as well as algebraic
1832 methods, such as exploiting the grid structure in inputs, to streamline the model's efficiency
1833 without sacrificing accuracy. These approaches will be key to unlocking the production-level
1834 (i.e., scalable) potential of such methods.

Chapter 4

ON PREDICTING SOCIODEMOGRAPHICS FROM MOBILITY SIGNALS

Inferring sociodemographic attributes from mobility data could help transportation planners better leverage passively collected datasets, but this task remains difficult due to weak and inconsistent relationships between mobility patterns and sociodemographic traits, as well as limited generalization across contexts. We address these challenges from three angles. First, to improve predictive accuracy while retaining interpretability, we introduce a behaviorally grounded set of higher-order mobility descriptors based on directed mobility graphs. These features capture structured patterns in trip sequences, travel modes, and social co-travel, and significantly improve prediction of age, gender, income, and household structure over baselines features. Second, we introduce metrics and visual diagnostic tools that encourage evenness between model confidence and accuracy, enabling planners to quantify uncertainty. Third, to improve generalization and sample efficiency, we develop a multitask learning framework that jointly predicts multiple sociodemographic attributes from a shared representation. This approach outperforms single-task models, particularly when training data are limited or when applying models across different time periods (i.e., when the test set distribution differs from the training set).

4.1 Introduction

Over several decades, travel behavior scholarship has revolved around the interplay between *who* a person is and *how* they move. Household-travel surveys (HTSs) have long provided the empirical backbone for this work by coupling rich trip diaries with respondent characteristics such as age, gender, income, and household composition. Analyses drawing on these surveys consistently show that, after accounting for the built environment, sociodemographic traits still correlate with car ownership, mode choice, trip frequency, and trip-chaining behavior (Bhat and Koppelman, 1994; Lee et al., 2007; Lu and Pas, 1999; McGuckin and Murakami, 1999b; Mokhtarian and Chen, 2004)

In the past dozen years, the ubiquity of GPS-enabled smartphones has spawned a parallel,

1863 industry-scale source of mobility evidence in passively-generated mobile data, which includes
1864 call-detail records (CDR), location-based service (LBS) pings, connected-vehicle traces, and
1865 the like (Chen et al., 2016b). These datasets dwarf HTSs in both sample size and temporal
1866 length, are refreshed continuously, and can often be licensed at a fraction of the cost of
1867 running a tailored survey. Their content, however, is almost exclusively spatial temporal; they
1868 record where and when a device was observed but remain agnostic about *who* was holding it.
1869 This missing dimension limits many distributional and behavioral analyses, including those
1870 that require understanding how travel patterns vary across population subgroups.

1871 Despite this blind spot, public agencies have been keen on experimenting with mobile data
1872 products (Ugurel et al., 2024b). Metropolitan planning organizations (MPOs) see potential in
1873 using them to stitch origin-destination matrices (Alexander et al., 2015a; Iqbal et al., 2014a),
1874 site electric-vehicle chargers (Yang et al., 2017a), and evaluate complete-street retrofits (Bian
1875 et al., 2023). Yet the lack of respondent attributes imposes two related hazards. First,
1876 representativeness: smartphone datasets systematically under sample certain population
1877 subgroups (Li et al., 2024; Wang et al., 2025; Wesolowski et al., 2013; Wu et al., 2024c). As a
1878 result, decisions based solely on such data may reflect the travel patterns of overrepresented
1879 groups while ignoring others. Second, for many analyses MPOs would like to conduct, such as
1880 gauging the effects of new toll roads or assessing if expanded transit lines reach underserved
1881 communities, linking mobility traces to travelers' sociodemographic profiles is required.

1882 To unlock the full value of passively-generated traces, researchers need a way to infer
1883 or impute sociodemographic variables from the mobility patterns embedded in the devices.
1884 This issue has been studied from a variety of angles, including travel behavior (Auld et al.,
1885 2015; Zhang et al., 2024a), data mining of destination choice (Doi et al., 2021; Solomon et al.,
1886 2018; Wu et al., 2019; Zhong et al., 2015), communication metadata (Jahani et al., 2017;
1887 Razavi et al., 2024), transit smart card use (Ding et al., 2019), and information theory (Zhao
1888 et al., 2022). As our goal is to enable the use of mobile data for transportation planning, we
1889 primarily focus on studies that solely use mobility behavior to achieve the aim of imputing
1890 sociodemographic attributes from mobility traces.

1891 Investigating this sociodemographic-inference problem presents two principal obstacles.
1892 First, most of the foundational literature in travel behavior frames sociodemographics as
1893 shaping mobility, not the reverse. This perspective aligns with behavioral models in which
1894 age, income, and household responsibilities influence when, where and how people travel. As
1895 a result, conceptualizing mobility traces as predictive signals for inferring sociodemographic

1896 attributes remains a relatively underdeveloped and counterintuitive direction in the literature.
1897 Second, sociodemographic inference from behavioral signals is intrinsically difficult. While
1898 some features of travel (e.g., trip chaining or mode diversity) can correlate with variables like
1899 gender or income, the relationships are often weak, noisy, and highly context dependent. In
1900 practice, models trained on one dataset may exhibit strong predictive performance yet fail
1901 to generalize when applied to another dataset, especially across geographies with different
1902 urban forms and social norms (Sheller and Urry, 2006). These challenges complicate efforts
1903 to construct transferable models and limit the extent to which empirical findings can be
1904 applied universally.

1905 This study aims to overcome these obstacles through the following contributions. First,
1906 to confront the limited theoretical grounding for inferring demographics from mobility, we
1907 introduce a family of higher-order descriptors based on directed mobility graphs in which
1908 vertices encode activity purposes and edges connect chronologically adjacent trips. By
1909 *higher-order*, we mean measures that go beyond first-order metrics (i.e., visit counts, mode
1910 frequencies) to capture relations across sequences of trips and destinations (e.g., how evenly
1911 travel is spread across activities, whether trips form loops or tours, whether they mix modes,
1912 etc.). We demonstrate that these descriptors are strongly and interpretablely associated with
1913 sociodemographic attributes and that they substantially increase the predictive power of
1914 imputation models beyond classical and spatiotemporal features (defined in Section 4.4.1).

1915 Second, to mitigate the inherent difficulty and uncertainty in generalizing sociodemographic
1916 inference models across diverse contexts, we operationalize uncertainty quantification and
1917 model calibration for multi-class classification. Specifically, we leverage metrics that encourage
1918 models to match their prediction confidences with their accuracies. We also use visual
1919 diagnostic tools like reliability diagrams to identify calibration gaps in our experiments.
1920 Applying this protocol, we find that the proposed higher-order descriptors consistently
1921 improve out-of-sample likelihoods, but their effect on calibration is mixed: they close gaps in
1922 several settings but can yield conservative (under-confident) probabilities, possibly due to the
1923 risk of overfitting.

1924 Third, we establish the value of a multitask (MT) learning strategy that predicts multiple
1925 sociodemographic attributes simultaneously. Under standard statistical learning assump-
1926 tions, parameter sharing across related tasks reduces estimator variance and improves out
1927 of sample performance (Baxter, 2000; Caruana, 1997). Thus, jointly modeling targets like
1928 age, gender, income, and household size allows the network to exploit shared structure in

1929 the mapping from mobility to sociodemographic attributes, yielding a more data-efficient
 1930 and robust estimator than training separate models. Moreover, MT learning can improve
 1931 transferability across contexts by reducing sensitivity to distributional shift (i.e., changes in
 1932 the input–output relationship between training and test data). We corroborate these claims
 1933 with cross-temporal generalization experiments in which the multi-task variant consistently
 1934 outperforms single-task baselines under a range of sample-size constraints.

1935

1936 **Our Contributions:**

- 1937 • We introduce a behaviorally grounded family of **higher-order descriptors based**
 1938 **on mobility graphs**, in which vertices represent travel purposes and edges represent
 1939 temporally ordered trips between them. When appended to classical mobility features,
 1940 our feature set consistently raises the out-of-sample accuracy of sociodemographic
 1941 inference (e.g., age, gender, income, etc.) across multiple experimental setups.
- 1942 • We operationalize **uncertainty quantification and calibration methods** for mobility-
 1943 based sociodemographic inference, using metrics that align predicted confidence with
 1944 observed accuracy. Visual tools like reliability diagrams reveal calibration gaps and
 1945 help diagnose model behavior.
- 1946 • We demonstrate the benefits of **multi-task learning** for transferability and data
 1947 efficiency. Training a unified model to predict multiple sociodemographic attributes
 1948 jointly improves sample efficiency and enhances generalization to new data compared
 1949 to single-task baselines. In both data sparse regimes and those in which the test
 1950 set systematically differs from the training set, the multi-task variant consistently
 1951 outperforms single-task counterparts, indicating enhanced robustness and practical
 1952 transfer across contexts.

1953 The rest of this chapter is organized as follows: Section 4.2 provides a review of relevant
 1954 literature, both in the context of classical travel behavior studies and those that seek to
 1955 answer the same question as ours. Section 4.3 outlines the datasets we leverage, while Section
 1956 4.4 details our methodological approach. Section 4.5 presents our experimental design and
 1957 the numerical results. We conclude with a discussion of our findings, limitations, and ideas
 1958 for future work in Section 4.6.

1959 4.2 Literature Review

1960 In this section, we review previous efforts to infer sociodemographic backgrounds of individuals
1961 based on their mobility behavior. We distinguish between these studies and a parallel body of
1962 work that predicts demographics from mobile phone communication metadata (call records,
1963 texting patterns, app usage, etc.). Both lines of research share the premise that digital
1964 behavioral traces contain sociodemographic signals, but they differ in the type of behavior
1965 analyzed. Mobility focused studies use *where people go* as the primary predictor, whereas
1966 communication focused studies use who people connect with and how they use their devices.

1967 Most mobility-inference studies begin by transforming raw coordinates into interpretable
1968 spatial or temporal descriptors. Spatial metrics quantify the extent and diversity of an
1969 individual’s movements: the radius of gyration (Ding et al., 2019), the heterogeneity of visited
1970 locations (Wu et al., 2019), the number of unique visited locations (Wu et al., 2019; Zhong
1971 et al., 2015), and those that relate to the distance traversed (Wu et al., 2019). Temporal
1972 features characterize regularity and rhythm, such as the day-to-day similarity of location
1973 sequences or the distribution of departure times for commutes (Ding et al., 2019), under
1974 the hypothesis that highly routinized commuters differ demographically from, for example,
1975 students or gig-economy workers. Semantic signals add further discriminatory power. Studies
1976 extract the frequency of visits to certain points-of-interest (POIs). For example, frequent
1977 stops at beauty salons or supermarkets can proxy for gender (Doi et al., 2021), while regular,
1978 short school drop-offs can signal parental status. From a classical travel behavior lens, stop
1979 durations, tour counts, and land-use contexts also prove to carry predictive power on whether
1980 or not a person drives, their work status, and education level (Auld et al., 2015). Despite the
1981 range of mobility descriptors, it remains unclear which signals reliably predict which attributes
1982 and under what conditions. Reported associations are often ad hoc and context-specific. Most
1983 studies focus on overall accuracy using bundled features, without isolating the marginal value
1984 of specific groups (spatial, temporal, semantic) or testing their consistency across settings.
1985 To fill this gap, we systematically assess the incremental contribution and robustness of each
1986 feature family.

1987 To interpret these descriptors, researchers have moved from statistical models to more
1988 data-driven approaches. Early work mapped those hand-crafted features into conventional
1989 statistical or rule-based models. Auld et al. (2015) combined fuzzy clustering of classical
1990 travel descriptors with decision trees and nested-logit models, while Zhong et al. (2015)
1991 used tensor factorization to decompose location check-in patterns. More recent research

1992 has gravitated toward representation-learning. [Solomon et al. \(2018\)](#) interpreted each day’s
1993 trajectory as a “sentence” and learned Word2Vec embeddings prior to classification, whereas
1994 [Ding et al. \(2019\)](#) employed long short-term memory networks on smart-card sequences. [Xu](#)
1995 [et al. \(2020\)](#) modeled the city as a heterogeneous mobility network and learned embeddings
1996 by preserving both physical co-visitation and semantic similarity between users. Generally,
1997 these approaches aim to capture the temporal structure and contextual regularities of daily
1998 movement, letting the model implicitly learn what aspects of movement are informative for
1999 demographics. However, most of these models train separately for each attribute, ignoring
2000 shared structure across related targets (i.e., age and the number of children). This makes
2001 them data-hungry and less transferable across settings with distribution shifts. We address
2002 this by using a multitask learning framework that jointly predicts all attributes from a shared
2003 representation, improving sample efficiency and out-of-sample robustness.

2004 When it comes to the output of prediction models, the uncertainty that is quantified
2005 needs to be interpreted with care. Many works simply present a confusion matrix or error
2006 rate, which is an aggregate uncertainty measure. These tell us, for example, that the model
2007 is wrong 20% of the time overall, but not *which* 20% of cases or how to flag an uncertain
2008 individual. The tendency of earlier models to overfit their training data (as in [Auld et al.](#)
2009 [\(2015\)](#)) highlights the risks of relying on point predictions without accounting for variance or
2010 confidence. A few studies have taken steps toward more meaningful uncertainty estimates.
2011 For example, some inference models incorporate cross-validation accuracy into probability
2012 estimates, essentially adjusting predicted class probabilities downward to reflect the model’s
2013 known error rate ([Zhang et al., 2024a](#)). This can prevent overconfidence in the predictions by
2014 “baking in” the chance of error. Moreover, the recent work by [Zhao et al. \(2022\)](#) introduces a
2015 theoretical framework to estimate beforehand how separable the classes might be, given the
2016 covariance structure of mobility data. Building on these advances, we calibrate predicted
2017 probabilities using metrics that encourage honesty between confidence and accuracy. We
2018 further show that our feature set and multitask learning framework yield more reliable and
2019 discriminative uncertainty estimates compared to baselines. These steps support a shift
2020 toward uncertainty-aware inference, where decisions can reflect the model’s confidence.

2021 **4.3 Datasets**

2022 We analyze three of the four most recent waves of the Puget Sound Regional Council (PSRC)
2023 Household Travel Survey (HTS), with the exception being the 2021 wave during which

Table 4.1: Descriptive statistics of the PSRC HTS in the three waves used in this study (post-processing; values in parentheses denote the strata percentage associated with wave)

Category	Variable	2017	2019	2023
Wave	Field Dates	04/10 – 06/15	03/11 – 05/30	04/24 – 05/29
	Households	3,160	2,902	3,504
	Persons	5,545	5,116	5,959
	Trips	51,029	71,913	56,028
Gender	Male	2,714 (48.9%)	2,475 (48.4%)	2,597 (43.6%)
	Female	2,724 (49.1%)	2,533 (49.5%)	2,854 (47.9%)
	Non-binary	17 (0.31%)	23 (0.45%)	121 (2.03%)
Age	0–11	495 (9.74%)	477 (10.1%)	549 (10.8%)
	12–17	151 (2.97%)	159 (3.36%)	197 (3.93%)
	18–34	1,739 (34.2%)	1,514 (31.9%)	1,335 (26.7%)
	35–54	1,562 (30.7%)	1,480 (31.3%)	1,517 (30.3%)
	55–74	970 (19.1%)	941 (19.9%)	1,152 (23.0%)
	75+	165 (3.25%)	163 (3.44%)	267 (5.33%)
HH Income	Under \$25,000	410 (8.07%)	314 (6.63%)	354 (7.07%)
	\$25k–49,999	642 (12.6%)	588 (12.4%)	533 (10.6%)
	\$50k–74,999	744 (14.6%)	715 (15.1%)	646 (12.9%)
	\$75k–99,999	728 (14.3%)	657 (13.9%)	531 (10.6%)
	\$100k+	2,558 (50.3%)	2,460 (51.9%)	2,943 (58.8%)
Children in HH	0	3,538 (69.6%)	3,297 (69.6%)	3,371 (67.3%)
	1	702 (13.8%)	586 (12.4%)	604 (12.1%)
	2	712 (14.0%)	580 (12.2%)	777 (15.1%)
	3+	130 (2.57%)	250 (5.28%)	277 (5.53%)

2024 travel behavior was heavily confounded by the COVID-19 pandemic. The 2017, 2019, and
2025 2023 surveys were fielded biennially using an address-based probability sample covering the
2026 four-county central Puget Sound region (King, Kitsap, Pierce, and Snohomish). The 2021
2027 wave uniquely included an additional opt-in online panel, which we omit here for consistency.
2028 Table 4.1 highlights relevant descriptive statistics from our post-processed version of this
2029 dataset.

2030 During data cleaning we applied a set of exclusion rules to ensure that only complete
2031 and internally consistent trip records entered the analysis file. First, any trip lacking a valid
2032 origin or destination purpose code was deleted, because purpose is central to several of our
2033 behavioral indicators. We likewise removed trips for which the travel-mode field was blank;

2034 mode choice is a key outcome variable and cannot be reliably imputed when entirely missing.
 2035 Third, observations without usable spatial information—specifically, trips whose destination
 2036 coordinate could not be geocoded or whose reported distance was zero—were discarded,
 2037 as they preclude computation of distance-based measures. Finally, we excluded the small
 2038 number of records reporting negative travel durations, which are symptomatic of data-entry
 2039 or time-stamp errors. These filters leave a sample of trips with complete purpose, mode,
 2040 spatial and temporal attributes suitable for subsequent modeling.

2041 Although our motivation is to enable sociodemographic inference from passively-collected
 2042 mobile data, we conduct this study using HTS data due to two key advantages. First, it
 2043 provides ground-truth sociodemographic labels, which are essential for supervised learning and
 2044 evaluation. Second, because HTS is fielded biannually, it enables testing under distribution
 2045 shift by evaluating models across different survey waves. While our features are derived from
 2046 structured trip diaries, recent advances in imputation algorithms now make it feasible to
 2047 extract similar semantic information from raw GPS traces (Gao et al., 2024; Merikhipour
 2048 et al., 2024). Thus, the methods developed here can largely be applied to passive data once
 2049 suitably enriched.

2050 4.4 Methodology

2051 Section 4.4 details our methodology. 4.4.1 formalizes the mobility-graph representation and
 2052 defines the activity- and trip-level covariates used in prediction. In 4.4.2, we outline definitions
 2053 and metrics that allow us to operationalize uncertainty quantification in our context. Finally,
 2054 4.4.3 describes the theory behind the MT approach as well as the specific neural architecture
 2055 we leverage.

2056 4.4.1 Characterizing travel behavior with mobility graphs

2057 Daily activity chains are encoded as directed graphs $G = (V, E)$, where vertices V represent
 2058 unique destination purposes (e.g., home, gym, school), and edges E represent temporally
 2059 ordered trips between them. From this representation, we extract two layers of information:
 2060 (1) *activity (node) features*, which describe how frequently, evenly, and in what sequence
 2061 specific activities are pursued; and (2) *trip (edge) features*, which summarize mode diversity
 2062 (i.e., tendency to use different travel modes), co-travel composition (i.e., share of trips taken
 2063 alone vs. with others), and daily travel motifs (i.e., minimal, recurring subgraphs that capture

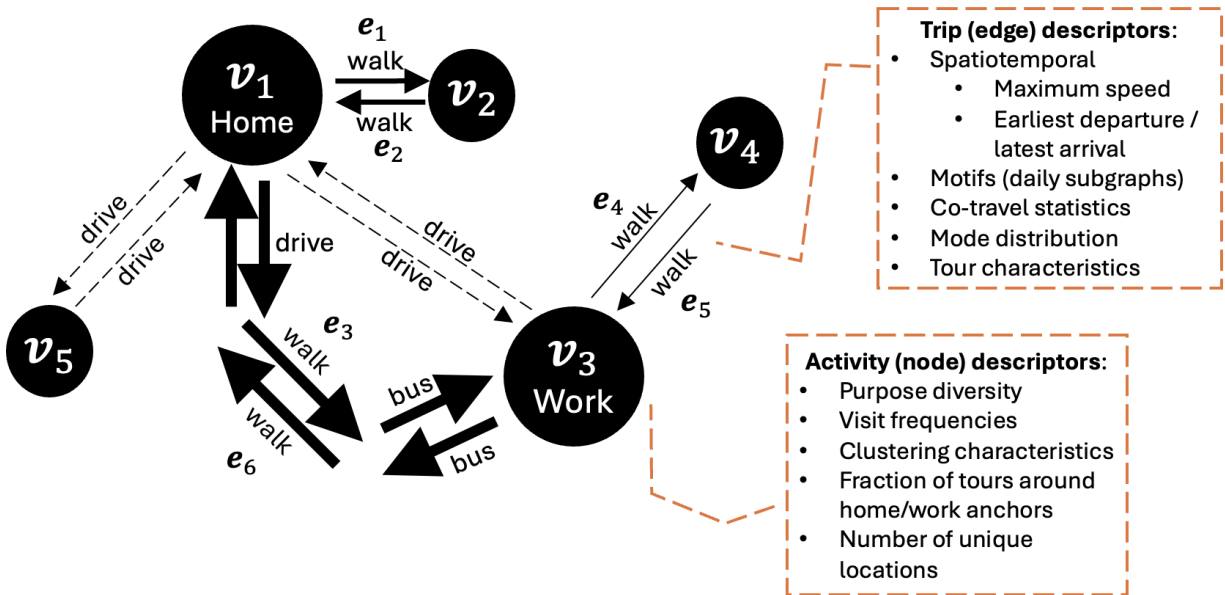


Figure 4.1: Example daily mobility graph where the edges are chronologically numbered. Width of trip arrows corresponds to frequency of visits over observation period. Dashed arrows denote known trips that are not observed on this day.

2064 the daily structure of trip sequences). We give precise definitions for metrics to quantify the
 2065 above in Section 4.4.1.2 (see Figures 4.1–4.3).

2066 At its simplest, the fraction of one’s trips dedicated to specific activities is a classical
 2067 indicator of who they may be (Lu and Pas, 1999). Persons with children spend more time
 2068 at schools, while older individuals spend more time shopping. Similarly, the fraction of
 2069 trips taken with modes can be indicative of sociodemographic background. In the Seattle
 2070 context, those with higher incomes tend to drive more, while men tend to bike more than
 2071 women (we detail more correlations in Section 4.5.1). Though this type of knowledge is not
 2072 readily available from raw GPS data, there are now imputation algorithms that can infer
 2073 multiple class trip purposes and modes with F1 scores up to 76% (Gao et al., 2024) and 92%
 2074 (Merikhipour et al., 2024), respectively.

2075 4.4.1.1 Activity (node) features

2076 Let $\mathbf{x} = (x_1, \dots, x_N)$ be the frequency vector of an individual’s N distinct origin-destination
 2077 (OD) purpose pairs, where each pair connects two activity purposes (e.g., going home, to

2078 the gym, to school). The total trip count is $T = \sum_i x_i$. To quantify how evenly trips are
 2079 distributed across these OD purpose pairs, we compute the Shannon entropy (Shannon, 1948):

$$H_t = - \sum_{i=1}^N \frac{x_i}{T} \log_2 \left(\frac{x_i}{T} \right), \quad (4.1)$$

2080 which is always non-negative and upper-bounded by $\log_2(N)$. High entropy indicates a diverse
 2081 and balanced use of OD purpose pairs, while low entropy reflects a concentration of trips on just
 2082 a few repeated purposes. While Shannon entropy measures the diversity of trip distribution,
 2083 the Gini coefficient captures the degree of inequality in how frequently OD purpose pairs are
 2084 used. Let the edge counts be sorted in non-decreasing order, $x_{(1)} \leq \dots \leq x_{(N)}$, and define
 2085 the cumulative trip count to rank k as $C_k = \sum_{i \leq k} x_{(i)}$. Then the Gini coefficient (Dorfman,
 2086 1979) is:

$$G = 1 + \frac{1}{N} - \frac{2}{NT} \sum_{k=1}^N C_k. \quad (4.2)$$

2087 A higher Gini indicates that a small number of purposes account for most trips, while a lower
 2088 Gini suggests more equitable use across destinations. This reflects the extent to which an
 2089 individual’s travel is exploratory versus habitual, connecting to classical notions of travel
 2090 regularity (Alessandretti et al., 2018; Kitamura and Van Der Hoorn, 1987).

2091 To quantify global cohesion in the mobility graph, we compute the *global clustering*
 2092 *coefficient* C_{glob} (Opsahl and Panzarasa, 2009). This metric captures the tendency for travel
 2093 purposes to be mutually connected through sequences of trips, forming tightly knit triangular
 2094 structures. Let τ_{Δ} denote the number of closed triplets (triangles) and τ_{\wedge} the number of
 2095 connected triplets (wedges). Then:

$$C_{\text{glob}} = \frac{3\tau_{\Delta}}{\tau_{\wedge}}. \quad (4.3)$$

2096 In our context, where nodes are activity purposes and edges are trips, high clustering
 2097 implies that individuals frequently travel between triplets of destinations in looping patterns
 2098 (e.g., home to gym to store to home), rather than taking out-and-back trips. This metric
 2099 overlaps with the concept of trip chaining, where multiple activities are linked into a single
 2100 tour rather than occurring as separate out-and-back trips from a primary anchor (e.g., home
 2101 or work) (Ellegard et al., 1977; Hanson, 1980). To complement the global clustering coefficient,
 2102 which reflects the overall cohesion of the network, we also compute the *mean local clustering*
 2103 *coefficient* \bar{c} (Kaiser, 2008). This measure averages each node’s local transitivity, placing

2104 greater emphasis on peripheral or low-degree destinations:

$$\bar{c} = \frac{1}{N} \sum_{v=1}^N \frac{2t_v}{k_v(k_v - 1)} \quad (4.4)$$

2105 where k_v is the degree of node $v \in V(G)$ (i.e., the number of other destinations directly
2106 connected to it), and t_v is the number of triangles that include v . In mobility terms, a high \bar{c}
2107 indicates that even less frequently visited destinations are embedded in tightly connected
2108 clusters. For example, this may suggest that auxiliary stops (e.g., a coffee shop, child’s
2109 school) are routinely integrated into a larger, cohesive tour structure rather than occurring in
2110 isolation.

2111 4.4.1.2 Trip (edge) features

2112 To characterize the edge layer of each mobility graph we compute three non-overlapping
2113 subfamilies of descriptors: spatiotemporal characteristics, daily travel motifs, and social
2114 co-travel mix. Let n_{trips} and n_{tour} denote, respectively, the number of observed trips and the
2115 number of closed tours (roundtrips anchored at either home or work) accumulated over the
2116 study horizon for a given individual. On the spatiotemporal side, we compute the fraction of
2117 trips taken during peak hours (f_{rush} ; defined as 7–9 AM or 4–6 PM) and the fraction taken
2118 on weekends (f_{weekend}). We also extract the earliest, average, and latest departure times
2119 across all observed days. In addition, we calculate the average and maximum values for trip
2120 duration, speed, and distance over the study period.

2121 The heterogeneity of daily mobility behavior tends to boil down to a handful of unique
2122 patterns, or “motifs” (Schneider et al., 2013). Inspired by the success attained by Wu et al.
2123 (2019), we derive motif counts after collapsing consecutive duplicate travel purposes (e.g.,
2124 *home* → *home* → *store* becomes *home* → *store*). The resulting sequence is parsed into
2125 one or more motifs drawn from the canonical set {single-no-return, out-and-back, chain,
2126 single-cycle, double-cycle, cycle-chain} (see Figure 4.2). Let m_j be the count of motifs of type
2127 j accumulated over all observed days and $M = \sum_j m_j$ the individual’s total motif count. We
2128 retain the motif fractions $f_k = m_j/M$ together with the motif entropy $H_m = -\sum_j \frac{m_j}{M} \log_2(\frac{m_j}{M})$
2129 which summarizes how evenly the six patterns are observed.

2130 Trips can also be composed of multiple modes, the use of which can correlate with income
2131 and age. Let $\mathcal{K}[mm]$ be the indicator that a given tour involves at least two distinct transport

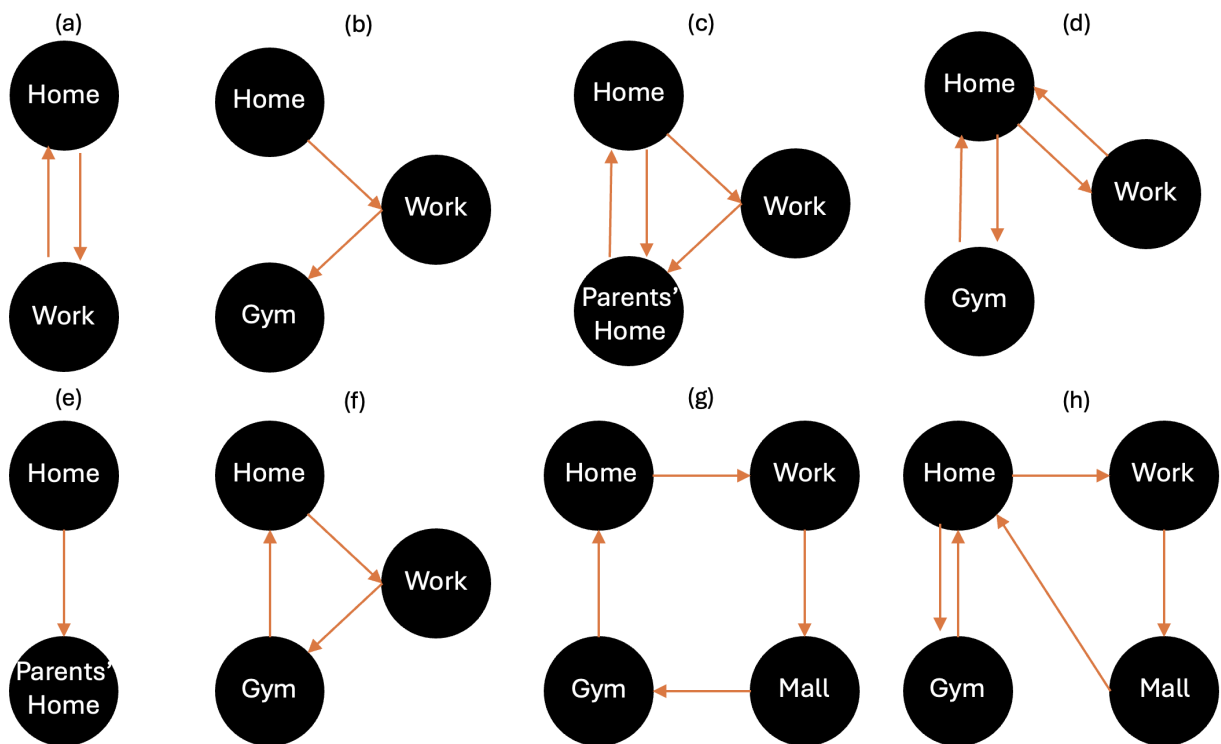


Figure 4.2: Illustration of daily travel motifs (Schneider et al., 2013; Wu et al., 2019); (a) Out-and-back; (b) Chain; (c) Cycle-chain; (d, h) Double-cycle; (e) Single-no-return; (f, g) Single-Cycle

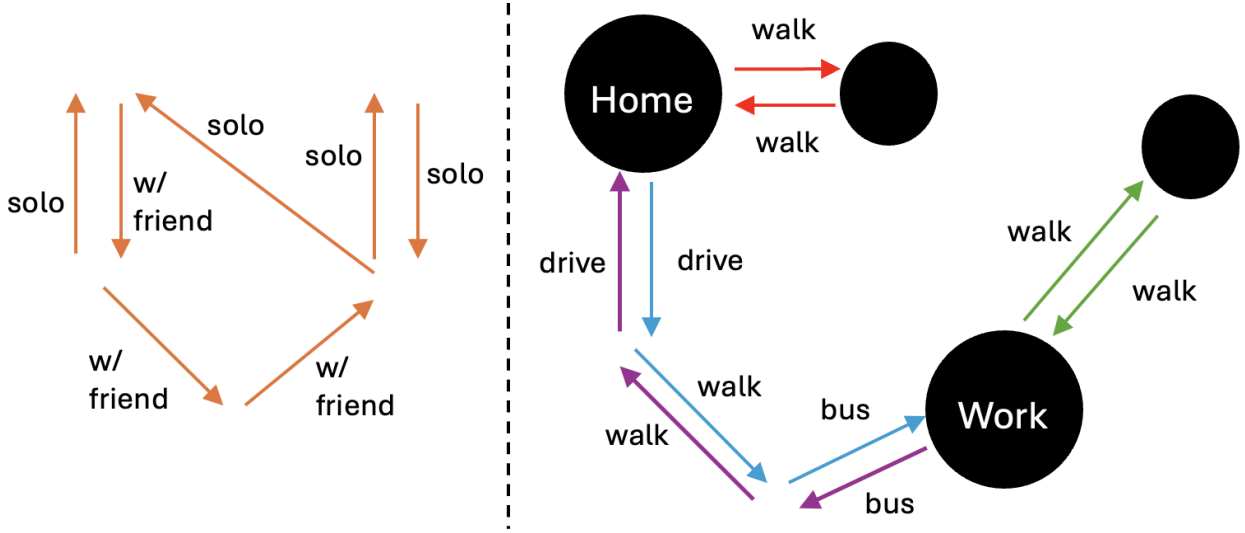


Figure 4.3: Illustration of selected metrics; (left) example travel day. In this case, $n_{\text{trips}} = 7$ and $f_{\text{comp}} = 3/7$; (right) each color denotes a tour to/from the anchors. In this case, $n_{\text{tour}} = 4$ and $f_{\text{mm}} = 2/4$.

2132 modes. The *multi-modal fraction* is then

$$f_{mm} = \frac{\mathcal{K}[mm]}{n_{\text{tour}}} \quad (4.5)$$

2133 A higher value reflects broader access to, or preference for, heterogeneous mode combinations
 2134 (e.g., walk-bus-walk), which we find is positively correlated with income in the Seattle context
 2135 (see Figure 4.5).

2136 Finally, whether one travels with companions correlates strongly with multiple sociode-
 2137 mographic labels. Thus, we tag each trip as *solo*, *with household companion(s)*, or *with*
 2138 *non-household companion(s)*. Let n_{solo} , n_{hh} , n_{nonhh} be the respective counts. We define

$$f_{\text{solo}} = \frac{n_{\text{solo}}}{n_{\text{trips}}}, \quad f_{\text{hh}} = \frac{n_{\text{hh}}}{n_{\text{trips}}}, \quad f_{\text{nonhh}} = \frac{n_{\text{nonhh}}}{n_{\text{trips}}}, \quad (4.6)$$

2139 and composite fraction with companions $f_{\text{comp}} = 1 - f_{\text{solo}}$. These metrics embed information
 2140 on household structure, caregiving roles, and wider social engagement. An illustrative
 2141 computation for one travel day appears in Figure 4.3.

2142 *4.4.2 Probabilistic scoring and calibration for categorical targets*

2143 We next tackle the retention of imputation uncertainty in supervised multi-class classification.
 2144 Here, we try to answer two methodological questions: (1) How can a model provide both
 2145 class probabilities and a principled measure of confidence? (2) Which evaluation criteria
 2146 reward not only accuracy but also well-calibrated uncertainty? We begin with definitions
 2147 and then describe metrics and visual tools to understand the output of various models.

2148 *4.4.2.1 Definitions*

2149 Let $Y \in \{1, \dots, K\}$ denote the sociodemographic label and X the mobility graph descriptors,
 2150 both drawn from the ground truth joint distribution $\pi(X, Y) = \pi(Y | X)\pi(X)$. A model m
 2151 outputs $m(X) = (\hat{Y}, \hat{P})$, where \hat{Y} is the predicted class and \hat{P} the associated confidence (e.g.,
 2152 probability of correctness). We would like calibrated confidence estimates, meaning that if m
 2153 outputs a confidence of 0.6 for 10 predictions, roughly 6 should be correct. Formally, *perfect*
 2154 *calibration* is defined as

$$\Pr(\hat{Y} = Y | \hat{P} = p) = p \quad \text{for all } p \in [0, 1], \quad (4.7)$$

2155 where the probability is over the joint distribution. A lack of calibration can lead to biased
 2156 share estimates and misleading uncertainty, whereas well-calibrated models ensure that
 2157 confidence values are useful and interpretable. However, the probability above cannot be
 2158 computed directly from finite samples, as \hat{P} is a continuous random variable. We therefore
 2159 rely on empirical approximations to estimate calibration.

2160 *4.4.2.2 Metrics and reliability diagrams*

2161 We evaluate predictions with complementary measures of separability, probability quality,
 2162 and calibration. Top 1 accuracy is the fraction of test cases for which the class with highest
 2163 predicted probability coincides with the true label, which is interpretable as “percentage
 2164 correct”. However, it evaluates predictions at a single decision rule (i.e., choose the argmax),
 2165 and is therefore sensitive to class imbalance and agnostic to the quality of the full probability
 2166 vector. To assess separability independent of any fixed threshold, we also report the area
 2167 under the ROC curve (AUROC), which measures how often the model ranks the true class
 2168 above competing classes (0.5 = random, 1 = perfect). For multiclass tasks, we compute the
 2169 one-against-rest macro-AUROC to ensure equal contribution from each class ([Hand and Till](#),

2170 2001). AUROC generalizes accuracy by integrating over all possible confidence thresholds,
 2171 but it still does not assess probability calibration. A model can have high AUROC while
 2172 producing poorly calibrated probabilities. Thus, we report AUROC alongside metrics that
 2173 reward well-calibrated probability estimates to evaluate uncertainty quality.

2174 One such metric is the negative log-likelihood (NLL), a standard measure of a probabilistic
 2175 model’s quality (Hastie et al., 2001). It is also commonly called the cross-entropy loss in the
 2176 context of deep learning (LeCun et al., 2015). Given $\{(x_i, y_i)\}_{i=1}^N$ as inputs and K classes, a
 2177 probabilistic classifier specifies a categorical distribution q , where each sample has predicted
 2178 class probabilities $\hat{\mathbf{p}}_i = (\hat{p}_{i1}, \dots, \hat{p}_{iK})$ satisfying $\sum_k \hat{p}_{ik} = 1$. The likelihood of the observed
 2179 label y_i under the model for input x_i is the scalar $q(y_i | x_i) = \hat{p}_{i, y_i}$. Then, the NLL averages
 2180 the negative log of those probabilities:

$$\text{NLL} = -\frac{1}{N} \sum_{i=1}^N \log q(y_i | x_i) = -\frac{1}{N} \sum_{i=1}^N \log \hat{p}_{i, y_i}. \quad (4.8)$$

2181 NLL is minimized when the predicted probabilities match the true conditional distribution
 2182 (Gneiting and Raftery, 2007). Intuitively, NLL rewards placing high probability on the correct
 2183 class and heavily penalizes overconfident mistakes.

2184 On the other hand, the Expected Calibration Error (ECE) summarizes how well confidences
 2185 match accuracies. For each sample i , let $\hat{y}_i = \arg \max_k \hat{p}_{ik}$ be the predicted label, and let the
 2186 confidence of the prediction be the highest predicted probability: $\hat{p}_i = \max_k \hat{p}_{ik}$. To derive
 2187 ECE, we partition the predictions into M confidence bins (each of size $1/M$). Let B_m be the
 2188 set of indices of samples whose prediction confidence falls into the interval $I_m = \left(\frac{m-1}{M}, \frac{m}{M}\right]$.
 2189 The accuracy and average confidence within bin B_m are:

$$\text{acc}(B_m) = \frac{1}{|B_m|} \sum_{i \in B_m} \mathbb{1}(\hat{y}_i = y_i), \quad (4.9)$$

$$\text{conf}(B_m) = \frac{1}{|B_m|} \sum_{i \in B_m} \hat{p}_i, \quad (4.10)$$

2190 where y_i is the true class label for sample i . Given these definitions, the expected calibration
 2191 error is the weighted average of the absolute accuracy-confidence gaps (Naeini et al., 2015):

$$\text{ECE} = \sum_{m=1}^M \frac{|B_m|}{N} |\text{acc}(B_m) - \text{conf}(B_m)|. \quad (4.11)$$

2192 The difference between accuracy and confidence for a given bin is called the *calibration gap*.
 2193 Thus, a perfectly calibrated model attains $\text{ECE} = 0$. In our experiments, we pre-specify
 2194 $M = 15$ and use equal-width bins, following common practice in the calibration literature
 2195 (Naeini et al., 2015). This number provides a practical balance between resolution and
 2196 statistical stability: too few bins obscure fine-grained miscalibration, while too many yield
 2197 noisy estimates due to small sample counts per bin. Equal-width binning ensures that
 2198 confidence intervals are evenly spaced and facilitates comparability across models and tasks.

2199 One related diagnostic tool we leverage is the *reliability diagram* (e.g., shown in Figure 4.7),
 2200 which visually represents model calibration. These diagrams plot expected sample accuracy
 2201 (Eq. 4.9) as a function of confidence (Eq. 4.10). If the model is perfectly calibrated, the
 2202 diagram should trace the identity line. Points below the diagonal indicate over-confidence,
 2203 while those above indicate under-confidence. Note that reliability diagrams do not display the
 2204 proportion of samples in each bin, and thus cannot be used to estimate how many samples
 2205 are calibrated.

2206 4.4.3 Multitask Learning

2207 Multitask learning improves generalization by *inductive transfer*, the idea that learning
 2208 several related tasks together leads to better performance than learning each one in isolation
 2209 (Caruana, 1997). In the classical “hard parameter sharing” formulation, tasks share part of
 2210 the model’s parameters, which acts as a data-dependent regularizer that limits the effective
 2211 hypothesis space. When tasks are related, this sharing reduces sample complexity: fewer
 2212 labeled examples per task are needed to reach a given level of accuracy (Baxter, 2000).

2213 Let T tasks be indexed by $t \in \{1, \dots, T\}$, where each task provides samples $(X, y_t) \sim \mathcal{D}_t$
 2214 and a task-specific loss ℓ_t . We learn a shared representation $h = g(X; \theta)$ and task-specific
 2215 predictors $f_t(h; \phi_t)$ by minimizing the weighted multi-task risk:

$$\min_{\theta, \{\phi_t\}} \sum_{t=1}^T w_t \mathbb{E}_{(X, y_t) \sim \mathcal{D}_t} [\ell_t(f_t(g(X; \theta); \phi_t), y_t)] + \Omega(\theta, \phi_t), \quad (4.12)$$

2216 where w_t are task weights, θ are the shared parameters of the representation h , ϕ_t are the
 2217 task-specific parameters of head f_t , and Ω is a standard parameter regularization term. In
 2218 practice, we optimize the empirical version using mini-batches sampled from the joint dataset
 2219 and backpropagate the sum of per-task losses; the shared parameters θ receive gradients from
 2220 all tasks and thereby encode common structure.

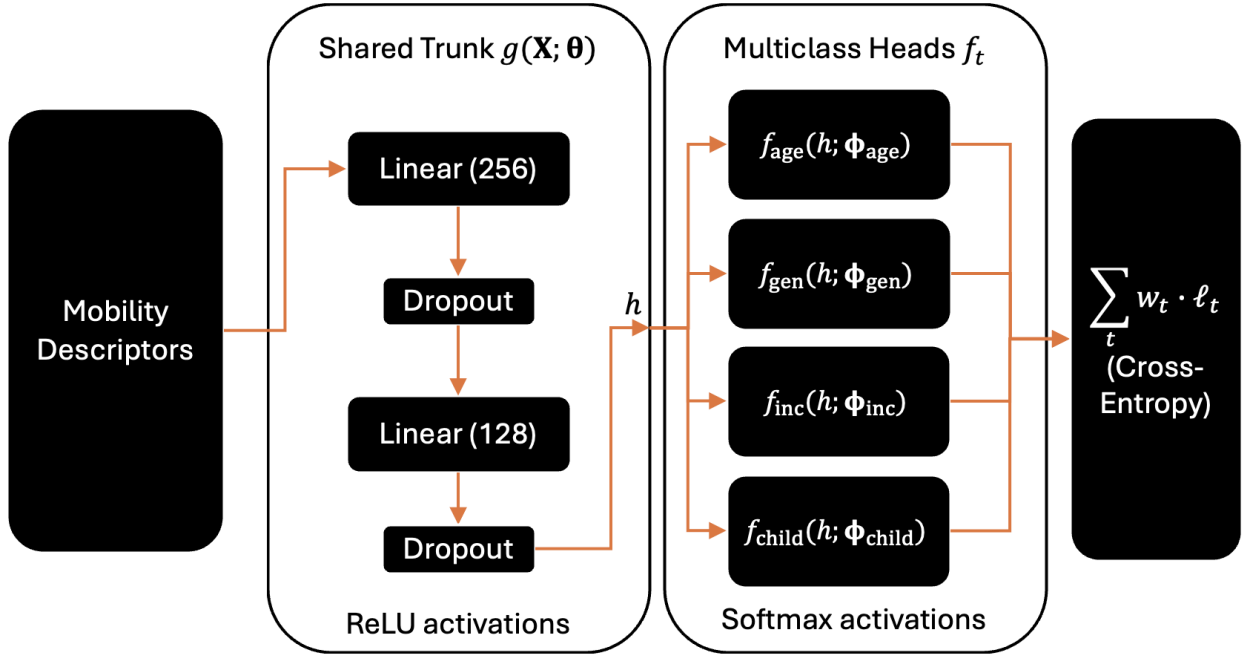


Figure 4.4: Shared-trunk multitask architecture. Mobility descriptors are mapped to a shared representation $h = g(X; \theta)$ by a two-layer feed-forward network with ReLU activations and dropout. Four task-specific heads f_t produce class probabilities via softmax. Training minimizes the cross-entropy loss $\sum_t w_t \cdot \ell_t$ where we set the weights to be equal.

2221 In this study, we adopt hard parameter sharing in a feed-forward network (shown in
 2222 Figure 4.4). A shared trunk g maps mobility features to a latent representation via three
 2223 connected layers with rectified linear units (ReLU) and dropout. ReLU enables the network
 2224 to learn complex nonlinear functions by zeroing out negative inputs, while dropout improves
 2225 generalization by randomly disabling activations during training and encouraging distributed
 2226 representations (LeCun et al., 2015). Four multiclass task heads f_t branch from this represen-
 2227 tation to predict age, gender, income, and number of children. This constitutes a genuine
 2228 MT setup because (i) all tasks update the shared trunk during training, and (ii) only a small,
 2229 task-specific set of parameters is unique to each head. The training objective is the weighted
 2230 sum of task losses, where we use equal weights for each task.

2231 The mobility descriptors given in Section 4.4.1 may capture latent routines that are jointly
 2232 informative for several sociodemographic attributes. Thus, sharing this representation across
 2233 the learning tasks pools statistical evidence, improving sample efficiency for sparsely labeled

2234 or imbalanced tasks, and regularizing confidence so that probabilities remain conservative
 2235 under modest distributional shifts. Consistent with the theory above, our experiments show
 2236 that the shared-trunk model often matches or exceeds single-task baselines in AUROC while
 2237 delivering lower NLL and ECE, particularly when training data are limited or the test set
 2238 differs from the training set.

2239 **4.5 Experiments**

2240 We showcase our results. In 4.5.1, we analyze correlations between our descriptors and
 2241 sociodemographic targets, which fit well-known tropes in the literature. 4.5.2 assesses the
 2242 extent to which our features improve the predictability of demographics and model calibration.
 2243 In 4.5.3, we demonstrate the value of multitask learning, which improves outcomes in data
 2244 scarce regimes or on test sets that differ from the training set.

2245 *4.5.1 Linkages between mobility descriptors and sociodemographics*

2246 Figure 4.5 shows the largest magnitude correlations between selected sociodemographics and
 2247 our feature set. Age is most strongly related to the share of school-purpose travel ($\rho \approx -0.50$)
 2248 and the drive-alone fraction ($\rho \approx +0.30$), with younger travelers exhibiting more bus use and
 2249 group car travel. Income is negatively associated with the share of shopping trips ($\rho < 0$)
 2250 but positively related to trip speed, distance, and car use, consistent with longer, faster
 2251 trips among higher-income travelers. Interestingly, higher income also correlates with greater
 2252 transit use, likely reflecting the prevalence of white-collar commuters in Seattle who rely on
 2253 transit for downtown access. Gender effects are modest ($|\rho| \leq 0.08$) but systematic: men tend
 2254 to bike more, travel solo, and have a higher fraction of their trips be work-related, while
 2255 women have more complex tours, more errand/appointment trips, and more shopping trips.
 2256 By contrast, the number of children shows much stronger signal (up to $|\rho| \approx 0.5$): households
 2257 with more children are characterized by more carpooling (3+ in vehicle) and higher shares of
 2258 school-bus, school, and escort travel, alongside lower prevalence of solo driving, transit, and
 2259 walking, fewer work-anchored tours, shorter trip durations, and fewer shopping trips. These
 2260 patterns reflect the child-serving, time-constrained nature of family travel.

2261 Table 4.2 shows the results of linear regression models with the sociodemographic attributes
 2262 as the independent variables and our features as the dependent variables. Due to high
 2263 negative correlation between age and household size (older respondents tend to live in smaller
 2264 households as children move out), we try two variants of each model. Multimodality rises with

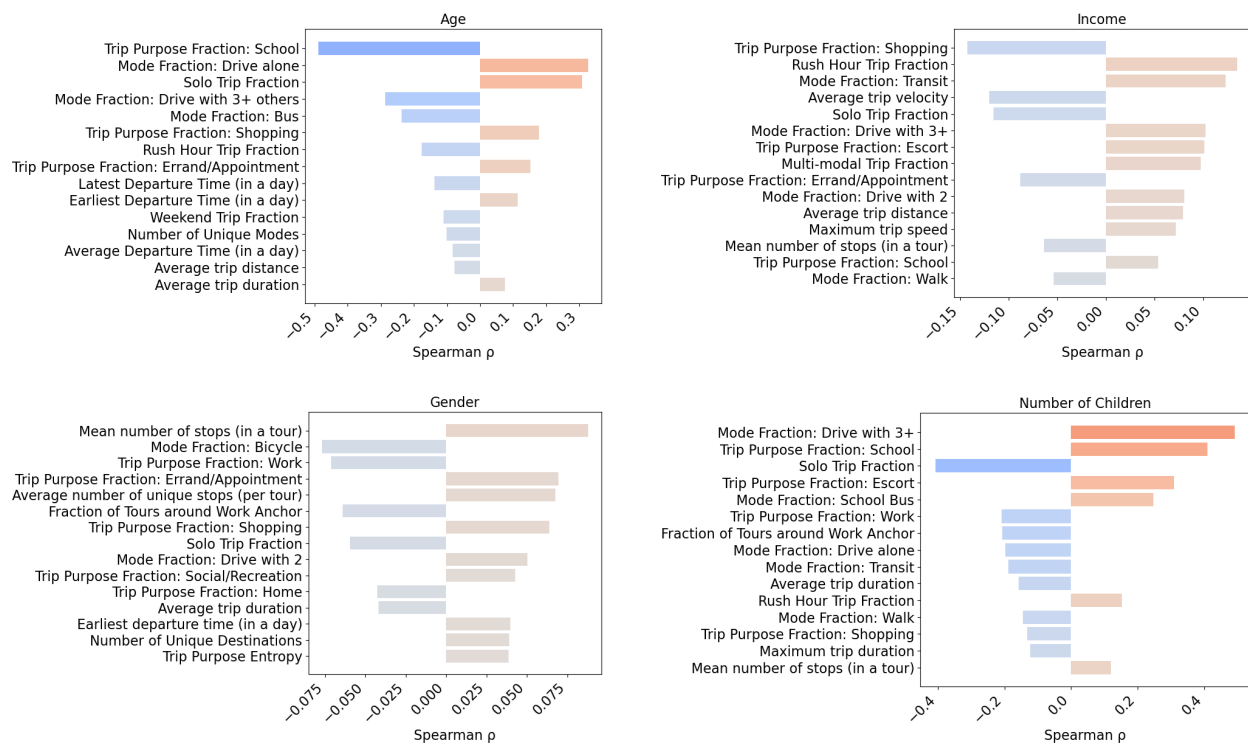


Figure 4.5: Largest magnitude Spearman rank correlations (all $p < 0.001$) between mobility descriptors and demographics. (top left): age; (top right): household income; (bottom left): gender; (bottom right): number of children. Bars to the left indicate negative associations; to the right, positive.

2265 age, income, and number of children, but the explained variance is small. The fraction with
 2266 companions shows the clearest sociodemographic signal: it decreases with age and increases
 2267 with income, being female, and number of children. By contrast, average local clustering
 2268 is only weakly related to demographics, tending to be lower at higher incomes and slightly
 2269 higher among women, with small effects. Out-and-back motif fractions decline with age and
 2270 among women. Overall, co-travel patterns carry the most sociodemographic information in
 2271 these linear models, while clustering and motif shares add subtle signals.

2272 4.5.2 Uplift of mobility graph features on predictability

2273 We evaluated the incremental predictive value of mobility features using nested feature sets
 2274 that comprise of classical variables (C), spatiotemporal attributes (ST), diversity-related
 2275 metrics (D), daily motifs (M), and co-travel statistics (CT), all discussed in Section 4.4.1.

Table 4.2: Selected linear regression models

Dependent Variable Model / Ind. Variable	Multi-modal Fraction		Fraction with Companions		Avg. Local Clustering Coeff.		Out-and-back Fraction (motif)	
	M1	M2	M1	M2	M1	M2	M1	M2
Intercept	0.50**	0.53**	0.46**	0.61**	0.35**	0.35**	0.42**	0.41**
Age	0.02**	0.01**	-0.08**	-0.12**	-0.00	0.00	-0.02**	-0.02**
Income	0.03**	0.03**	0.01**	0.02**	-0.00*	-0.01*	-0.00	-0.01
Gender	0.01	0.01	0.05**	0.06**	0.02**	0.02**	-0.02**	-0.02**
Household Size	0.03**	—	0.12**	—	-0.01	—	-0.01	—
R^2	0.01	0.01	0.22	0.16	0.002	0.001	0.004	0.003

[**, and * signify a p-value less than 0.01, and 0.05, respectively.]

2276 Classical (C) covariates include purpose shares (e.g., work, shopping, leisure), mode shares
 2277 (drive, transit, walk), and simple tour statistics (e.g., number of tours and the share anchored
 2278 at home/work). Spatiotemporal (ST) adds departure and arrival timing (e.g., first departure,
 2279 latest arrival), trip durations, shares by period (peak/off-peak, weekend), and basic speeds
 2280 and distances. Diversity (D) captures how spread-out travel is across activities, modes, and
 2281 origin-destination pairs (e.g., purpose entropy, the fraction of multi-modal tours, and triadic
 2282 closure among destinations). Motifs (M) comprise the fractions of canonical day-level patterns
 2283 (out-and-back, chains, cycles, etc.) together with a summary of their evenness. Co-travel (CT)
 2284 records with whom trips are taken: the shares of solo, with-household, and with-non-household
 2285 trips, and the overall accompanied fraction. In our setup, each set builds on the previous
 2286 one—for example, the +ST set includes all classical variables plus spatiotemporal attributes,
 2287 while the +CT set includes all preceding groups and adds co-travel statistics.

2288 All experiments followed the same data protocol, which had two evaluation splits: on the
 2289 *overall* split, we created a 70/10/20 train/validation/test partition on the combined waves of
 2290 the PSRC survey. On the *cross-temporal* split, we trained and validated on the 2017 and 2019
 2291 waves and tested on the 2023 wave, a period with documented shifts in mobility behavior due
 2292 to the aftermath of COVID-19 (see differences in wave composition in Table 4.1). For both
 2293 splits, we used five-fold multilabel cross-validation on the pooled training and validation sets,
 2294 with final performance reported on the fixed test set (either a random 20% holdout or the
 2295 2023 wave, respectively).

2296 We evaluated performance on top-1 accuracy, area under the receiver operating curve
2297 (AUROC), negative log-likelihood (NLL), and expected calibration error (ECE), as defined
2298 previously. Figure 4.6 presents the main results for the multi-task DNN, which is our primary
2299 model due to its compatibility with shared representation learning. To assess the robustness of
2300 the proposed feature sets across model types, we additionally evaluated three other classifiers:
2301 random forests (RF), gradient boosted machines (GBM), and support vector machines (SVM).
2302 Results for these models are reported in C.1 (Tables C.1-C.8), and show that the proposed
2303 descriptors broadly improve separability and likelihood across models. However, the most
2304 comprehensive feature set (+CT) does not always yield the best performance—approximately
2305 44% (14 out of 32) of the top metric values (shown in red) come from simpler sets. Performance
2306 trends are largely stable across model types, suggesting that the observed improvements
2307 stem primarily from the covariates rather than model-specific effects. That no single model
2308 dominates across all tasks and metrics further underscores the absence of a one-size-fits-all
2309 solution: practitioners must select models and feature sets based on their specific performance
2310 goals.

2311 For the multi-task DNN shown in Figure 4.6, we note the following observations of interest.
2312 First, extra covariates tend to help more under the *overall* split than the *cross-temporal*
2313 split. This suggests that some of the added features may capture time-specific patterns that
2314 do not generalize well when the test data come from a different distribution. In general,
2315 while richer feature sets can increase expressiveness, they may also introduce redundancy
2316 or collinearity, amplifying variance without adding meaningful new signal—especially when
2317 model capacity is not properly regularized. This effect is not uniform across models. As
2318 shown in Tables C.1-C.8, gradient-boosted trees (GBMs) appear less prone to overfitting
2319 under the cross-temporal split, possibly due to their implicit regularization and ability to
2320 ignore noisy or uninformative splits.

2321 A second trend relates to model calibration. The expected calibration error (ECE) does not
2322 consistently improve with feature richness. In two of the four tasks, the most comprehensive
2323 set (+CT) yields the best-calibrated predictions, but in the others, simpler feature sets
2324 perform better. GBMs again stand out, often producing the lowest ECE scores overall, which
2325 may reflect their ability to learn conservative margins or resist overconfidence in low-signal
2326 settings. These findings highlight that while added features can increase model expressiveness,
2327 their value depends on stability across contexts and their interaction with model calibration.

2328 Finally, the proposed descriptors generally improve both class separability (AUROC) and

2329 model fit (NLL). Under the overall split, gains tend to increase with each feature group, with
 2330 the +CT set frequently yielding the best performance and rarely underperforming simpler
 2331 sets. In the cross-temporal setting, improvements are more muted and task-dependent. For
 2332 Age, +CT performs best for RF and GBM and is competitive for DNNs. For Number of
 2333 Children, +CT achieves the highest AUROC, while +D yields the best NLL. Results for
 2334 Household Income and Gender are more variable, with no single feature set dominating across
 2335 metrics or models. These patterns suggest that feature utility is context- and task-dependent,
 2336 with diminishing returns or instability emerging under distribution shift.

2337 Reliability diagrams provide hints towards clarifying the variability in ECE performance.
 2338 Figure 4.7 highlights that, under several settings (e.g., Age in the overall split and Number
 2339 of Children under the cross-temporal split), the empirical accuracy within certain confidence
 2340 bins lies well above the diagonal, indicating under-confidence: the model predicts correctly
 2341 more often than its stated probabilities would suggest. This conservatism leaves AUROC
 2342 and accuracy unchanged or improved but increases ECE, which penalizes the magnitude of
 2343 the accuracy–confidence gap irrespective of sign. Surprisingly, this seems to happen both
 2344 when more features are added (as in the top row) *and* when features are removed (as in the
 2345 bottom row).

2346 4.5.3 Impact of Multitask Learning

2347 We compared a shared trunk multitask (MT) network with matched single task variants (STV)
 2348 across age, gender, income, and number of children, while holding architecture, optimization,
 2349 and regularization fixed (see C.1 for full hyperparameter tuning details). To stabilize learning
 2350 and reduce task interference, we optionally applied a *per-task layer normalization* module
 2351 before each task head, normalizing activations separately for each prediction branch rather
 2352 than sharing normalization statistics across tasks. To probe sample efficiency, we randomly
 2353 subsampled the training split to fractions $\{1.0, 0.1, 0.01, 0.001\}$ and evaluated on the untouched
 2354 validation and test sets. To avoid overfitting as training data size got small, we used only the
 2355 classical and spatiotemporal features.

2356 Figure 4.8 displays our findings. On the overall split (top four rows of Figure 4.8), the MT
 2357 network and its STVs achieve comparable top-1 accuracies and AUROC at full data. As
 2358 the training fraction shrinks, the relative patterns diverge by task. Age, gender, and number
 2359 of children exhibit a modest MT advantage in accuracy and AUROC at intermediate and low
 2360 fractions, indicating effective transfer of shared structure when labels are scarce. By contrast,

2361 HH income shows mixed patterns, particularly at the smallest fractions, suggesting that
2362 shared representation can sometimes blur task-specific distinctions (i.e., negative transfer).
2363 In these cases, the single-task networks appear better able to specialize when the shared
2364 structure among targets is weak.

2365 The uncertainty metrics favor the multitask learning-based model. NLL remains lower
2366 for MT than STVs as data become scarce, and ECE is lower for MT for three of the four
2367 fractions over each sociodemographic target. This indicates that even when discrimination
2368 gains are inconsistent, the shared representation can still regularize learning by preventing
2369 overconfidence and smoothing probability estimates.

2370 Under the cross-temporal generalization setting (bottom four rows of Figure 4.8), the
2371 MT network is again competitive with STVs in accuracy and AUROC across most fractions,
2372 though the advantages are uneven. Gains are most apparent for HH income, number of
2373 children, and to a lesser extent, age, while gender shows little benefit or mild negative
2374 transfer. The improvements in NLL are more pronounced than in the overall split, suggesting
2375 that multitask learning helps regularize against overconfidence under distribution shift. By
2376 pooling representational strength across targets, the model can better withstand changes in
2377 the marginal distribution of mobility features. This effect is especially helpful when some
2378 individual targets have relatively few informative examples in the new distribution (i.e., with
2379 rare subgroups or behaviors that shift over time). This stabilizing effect is especially useful
2380 in real-world deployment, where shifts in behavior over time (e.g., due to pandemics, fuel
2381 prices, or policy changes) can erode model reliability.

2382 **4.6 Conclusion**

2383 This study advances sociodemographic inference from mobility traces along three fronts:
2384 (i) it introduces a behaviorally grounded family of higher-order mobility descriptors that
2385 move beyond first-order counts to encode trip sequencing, cohesion, and social co-travel,
2386 among other characteristics; (ii) it operationalizes uncertainty-aware evaluation for multi-class
2387 prediction in this context, a dimension largely absent in prior studies, which tend to focus on
2388 point estimates without assessing the reliability of model confidence, and; (iii) it examines how
2389 a shared-trunk multitask (MT) architecture compares to matched single-task variants in data
2390 efficiency and calibration quality. Empirically, the proposed descriptors generally raise out-of-
2391 sample accuracy and lower NLL across attributes. The benefits of MT learning, however, are
2392 nuanced. While it often improves calibration and robustness under data scarcity and temporal

2393 distribution shift, these gains are not uniform across targets. In some cases—particularly
2394 for attributes with weaker behavioral overlap or higher label noise—task interference leads
2395 to mild negative transfer, causing MT to underperform specialized single-task networks.
2396 Overall, shared representations can regularize predictions and enhance reliability, but their
2397 value depends on the degree of cross-task relatedness and the balance between shared and
2398 task-specific learning.

2399 This work has limitations that motivate future research. We do not derive features
2400 directly from raw GPS/LBS signals; instead, we rely on processed trip diaries with purpose,
2401 mode, and co-travel labels. Many descriptors could in principle be computed from passive
2402 traces (e.g., reverse-geocoded activity locations, dwell-time anchors, and mode inference
2403 from speed/acceleration and network context), but their fidelity will depend on upstream
2404 imputation methods (Gao et al., 2024; Merikhipour et al., 2024). Furthermore, our cross-
2405 temporal evaluation is confined to a single region; a natural extension is a cross-city study to
2406 assess geographic portability.

2407 More broadly, this paper does not take a model- or architecture-centric view of predictive
2408 performance. Our emphasis is on deriving behaviorally grounded features, quantifying their
2409 contributions, and examining model uncertainty rather than pursuing architectural novelty.
2410 The only modeling variation we explore is multitask learning, chosen to test whether shared
2411 representations help under data sparsity. This is a well-established line of inquiry rather than
2412 a new algorithmic proposal. That said, there is considerable potential in moving beyond
2413 hand-crafted descriptors toward models that can automatically discover structure in large-
2414 scale mobility data. Recent advances in Transformer-based and self-supervised architectures
2415 point to promising directions, particularly for unlabeled PCM with long temporal depth,
2416 where rich behavioral regularities could be captured through pretraining (Wu et al., 2024b)
2417 and later adapted to sociodemographic inference.

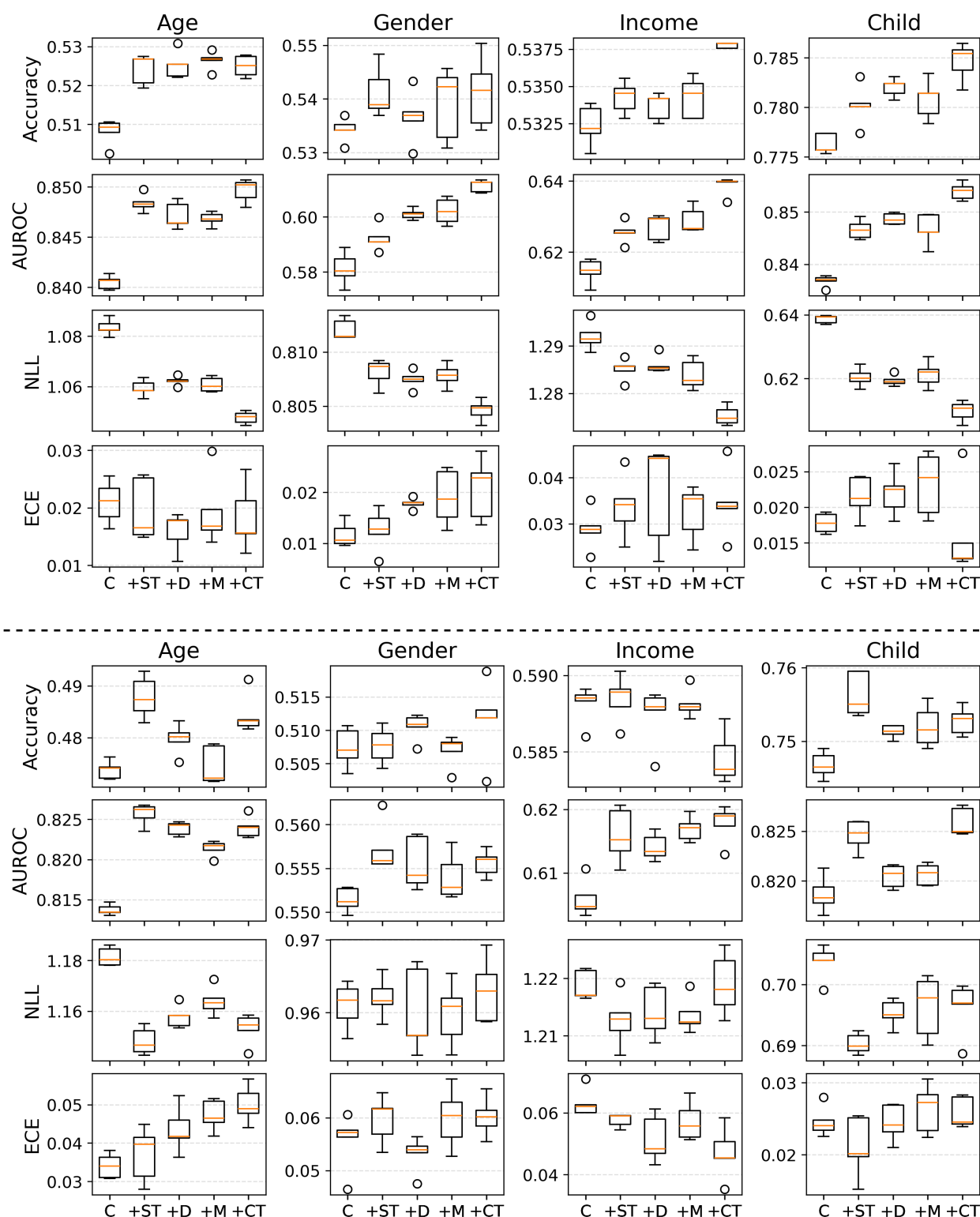


Figure 4.6: Marginal changes in top-1 accuracy (higher = better), AUROC (higher = better), NLL (lower = better), ECE (lower = better) as more features are added. (top four rows) Overall split; (bottom four rows) training with the 2017/2019 data and testing on the 2023 data.

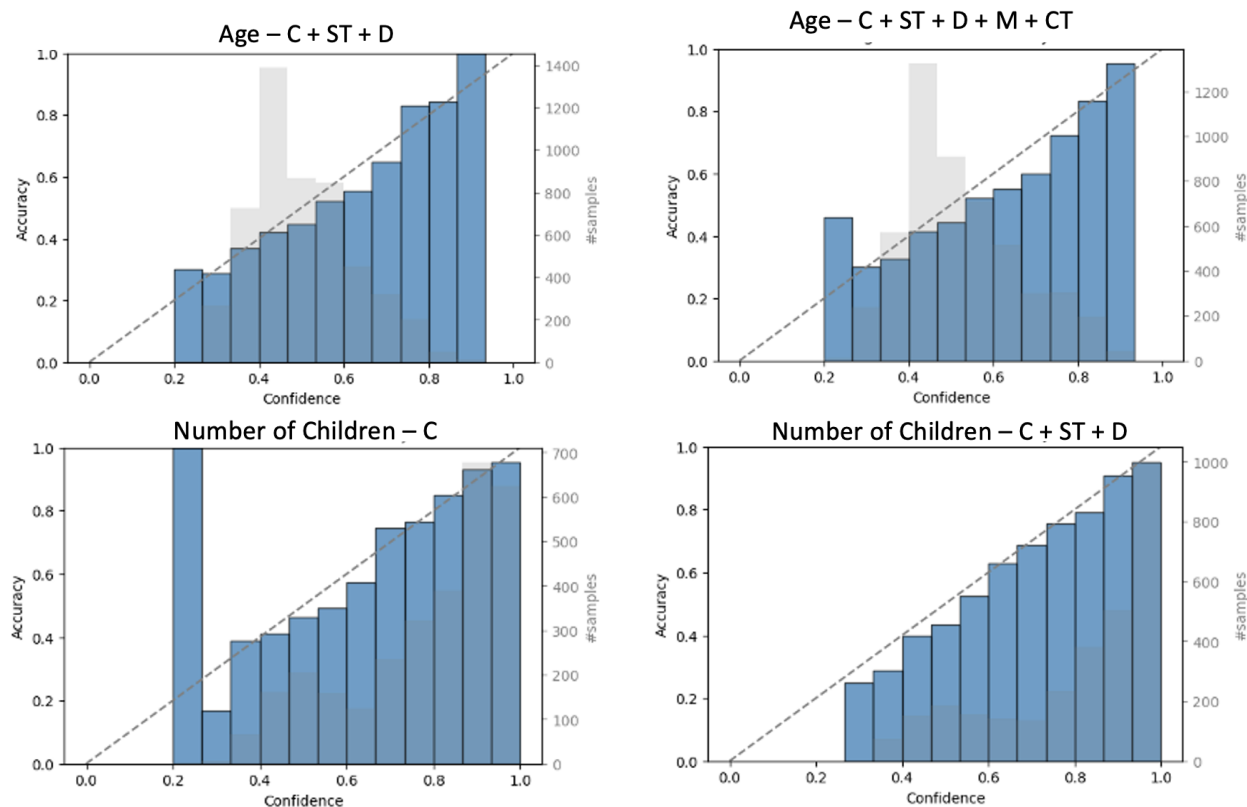


Figure 4.7: Reliability diagrams for representative settings. (TOP) Overall split, Age task with C+ST+D features (left) and All features (right). (BOTTOM) Cross-temporal split, Number of Children task with C (left) and C+ST+D (right). Bars show empirical accuracy within 15 equal-width confidence bins; the dashed line is the identity (perfect calibration). Grey histograms (right axes) give the number of samples per bin. Points above (below) the diagonal indicate under- (over-) confidence.

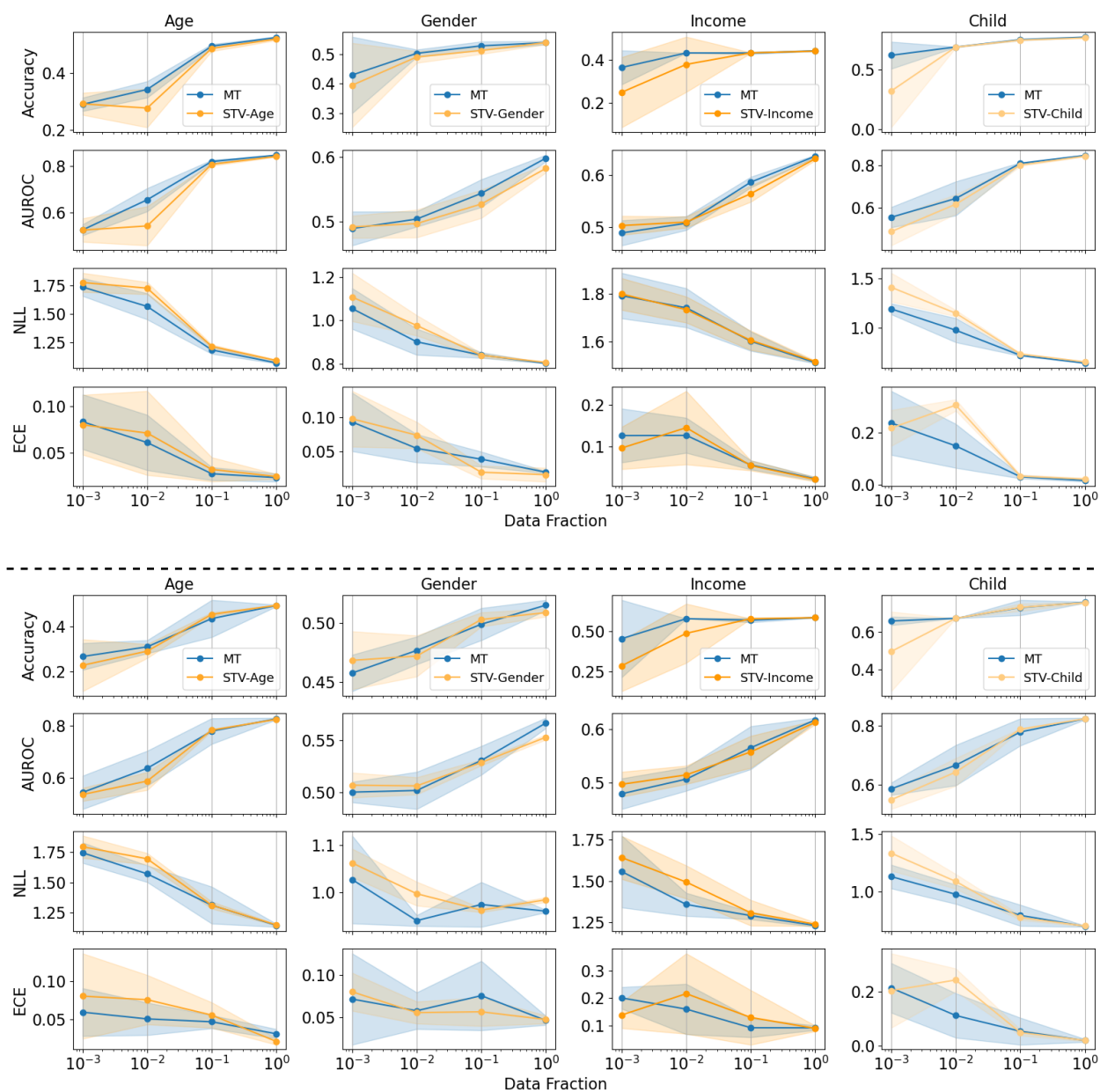


Figure 4.8: Performance of the MT variant (in blue) compared to ST variants (in orange) at different fractions of training data. (top four rows) Overall split; (bottom four rows) training with the 2017/2019 data and testing on the 2023 data.

Chapter 5

MPO'S USES OF AND NEEDS FOR BIG DATA

Big data products offer a new paradigm to understand and analyze human mobility patterns, a primary interest of long-range transportation planners. However, questions remain about how widely these datasets are used in practice and what factors limit their broader adoption in planning workflows and decision-making. In this chapter, we first provide a comprehensive review of the literature on the use of big data in metropolitan planning, with a focus on practical applications, perceptions, and institutional barriers to adoption. Then, we present the perspectives of more than 50 planners from MPOs across the United States, collected at a workshop in May 2024. While we found a range of use cases, there was also a tendency to focus on a narrow set of applications. *Transparency, regulation, and legitimacy* emerged as the primary factors influencing adoption decisions.

5.1 Introduction

Long-range transportation planners (LRTPs) are tasked with supporting decision-makers in shaping multimillion-dollar investments in future mobility. In recent years, the increased availability and commercialization of big data have resulted in the rise of mobility-analytics vendors who offer a suite of products, hereafter referred to as big data products (BDPs), designed to help planning agencies make more informed decisions. These companies aggregate passively collected data from multiple sources (such as smartphones, in-vehicle GPS trackers, connected-vehicle data, etc.), to provide insights into fine-grained human mobility behavior. Their data products typically include key information on individuals' movements over time, capturing patterns such as travel routes, trip frequencies, and timing, which can be invaluable for understanding mobility trends and planning future infrastructure investments.

While BDPs offer a new way to understand and analyze mobility patterns, we know little about their actual use in transportation planning. It is unclear how widely these datasets or products are being utilized by planners or to what extent they are influencing decision-making processes. Moreover, the full potential of these datasets, especially in areas such as forecasting, equity analysis, and resilience planning, remains largely unexplored. Understanding the gaps

2446 in usage and application is key to unlocking their value for long-term transportation planning.

2447 This chapter is motivated by two recurring themes identified through ad-hoc conversations
2448 with various practitioners. These themes are: (A) the use of BDPs has been limited to a
2449 narrow range of applications and they are not the preferred solution to many problems faced
2450 by LRTPs; and (B) only a small number of MPOs (primarily ones with sufficient resources)
2451 use BDPs. Thus, this chapter addresses the following questions:

- 2452 1. How do planners perceive big data products offered by mobility analytics companies?
- 2453 2. What factors influence an MPO's decision to adopt big data products and services?
- 2454 3. What are the various use cases for big data in regional transportation planning?

2455 **5.2 Literature Review**

2456 The interest MPOs have in leveraging big data for transportation planning has generated
2457 considerable discussion (in the form of studies, reports, and memorandums). Here, we first
2458 outline common use cases noted in the literature. We then examine the key barriers to
2459 adoption, focusing on how planners perceive big data products and the institutional, technical,
2460 and organizational factors that shape their uptake.

2461 *5.2.1 Use Cases*

2462 A primary use of big data is to calibrate and validate travel demand models. Many MPOs
2463 now incorporate archived traffic speed/travel-time data and volume estimates from private
2464 vendors to improve their regional models (Bauer et al., 2014). For example, the Maricopa
2465 Association of Governments (Phoenix area) has purchased probe-based speed data (e.g.,
2466 products offered by INRIX) for several years and uses the National Management Research
2467 Dataset, a national traffic speed dataset provided by the FHWA, to validate model travel
2468 times and congestion patterns (FHWA, 2020). These passive data allow agencies to check
2469 that model outputs (e.g. travel speeds, delays) reflect observed conditions, leading to model
2470 improvements. The Capital District Transportation Committee (Albany, NY MPO) likewise
2471 compared model-predicted speeds and delays with an archived traffic sensor dataset and
2472 discovered the model had been underestimating off-peak and non-recurring delay (FHWA,
2473 2020). Thus, insights from big data help refine model parameters and ensure forecasts are
2474 grounded in real-world conditions

2475 Another key application is using big data to fill gaps in traffic counts and OD data for
2476 planning studies. Some MPOs face declining coverage or frequency of traditional traffic
2477 counts (i.e., due to limited budgets or disruptions like the pandemic) and have turned to
2478 private “Big Data” vendors to estimate measures like Annual Average Daily Traffic (AADT)
2479 on road segments. In one case, the Tulsa, OK MPO obtained crowdsourced traffic volume
2480 estimates from a vendor to populate their roadway network map when local agencies stopped
2481 collecting counts, and validated these estimates against the MPO’s own data (StreetLight,
2482 2025). Big data derived from mobile devices are also used to generate OD matrices and travel
2483 flows by time of day and mode. Several MPOs have experimented with big data products
2484 to derive regional trip distributions or corridor-specific OD patterns, which can supplement
2485 or reduce reliance on infrequent household travel surveys (KimleyHorn, 2021). These OD
2486 datasets have been applied for corridor studies, sub-area analyses, and to understand regional
2487 travel trends (e.g. inter-county commuting flows) that would be expensive to capture via
2488 surveys. However, as discussed later, agencies often use these new data cautiously, usually as
2489 a “second source” to compare against traditional survey-based OD data or model outputs.

2490 LRTPs also see potential for big data in analyzing emerging transportation modes and
2491 trends that are not well covered by traditional data. Big data from smartphones and app-based
2492 services could theoretically reveal the magnitude and patterns of rideshare demand. Indeed,
2493 many MPOs aspire to use big data to understand the impacts of transportation network
2494 companies (TNCs), yet implementation has lagged. Interviews with large MPOs found that
2495 few had obtained or analyzed TNC trip data despite interest (KimleyHorn, 2021). Only a
2496 couple of agencies (e.g. New York Metro and Chicago regions) were exploring data-sharing
2497 agreements to acquire TNC trip records through local governments, and even those had no
2498 concrete analysis completed due to lack of internal expertise to manage the data (KimleyHorn,
2499 2021).

2500 Big data is also emerging as a tool to evaluate transportation accessibility, although efforts
2501 are relatively nascent. Because passively collected datasets can be disaggregated spatially
2502 and temporally, they hold promise for analyzing differences in travel times, mobility patterns,
2503 and service levels across different communities. For instance, location-based data can be
2504 used to do distributional comparisons across various criteria, including those based on spatial
2505 factors (e.g., the location of a proposed project in relation to the community it affects) as
2506 well as user-based factors (e.g., the number of users of the project that will benefit) (Boyd
2507 et al., 2024). Nonetheless, researchers caution that these data pose limitations, particularly

2508 concerns over representation and bias, as not all population groups are equally captured by
2509 smartphones or connected vehicles (Richardson, 2021; Wang et al., 2025).

2510 5.2.2 Factors Influencing Adoption

2511 Though use cases are plenty, adoption among MPOs remains limited. Many agencies cite a
2512 combination of financial, organizational, and technical barriers. The following subsections
2513 examine key factors shaping adoption decisions, including concerns about data quality,
2514 resource constraints, staffing capacity, and leadership support. As one NCHRP workshop
2515 attendee put it bluntly, "Our big data issues are straightforward: we don't have the technology,
2516 money, or the skills." (Pecheux et al., 2020).

2517 A central barrier to big data adoption among MPOs is concern over data quality, trans-
2518 parency, and methodological validity. Most vendors rely on proprietary algorithms to infer
2519 trips, scale samples, or impute missing information, yet rarely disclose these processes in
2520 detail. Elements like sampling timeframes, sample bias corrections, or definitions of trip
2521 purpose often diverge from the standards used in regional models (KimleyHorn, 2021). As a
2522 result, many MPOs avoid relying on these data for core planning tasks outside of limited
2523 domains like congestion monitoring.

2524 The situation is exacerbated by proprietary database structures and software platforms,
2525 which limit interoperability and reinforce dependence on specific vendors (Pecheux et al.,
2526 2020). These restrictions make it difficult for agencies to conduct independent validation or
2527 integrate the data into existing workflows. Because the raw inputs are transformed through
2528 opaque processes, MPOs cannot easily adjust for known biases in the way they might with
2529 traditional data sources like travel surveys (Singh et al., 2022). Furthermore, differences in
2530 methodology across vendors or changes over time can undermine the comparability of key
2531 performance measures, particularly in long-range planning contexts (Bauer et al., 2014).

2532 Cost is one of the most frequently cited barriers to the adoption of big data among MPOs.
2533 High upfront and ongoing subscription fees can be prohibitive, especially for smaller agencies
2534 operating under tight budgets (KimleyHorn, 2021). MPOs must assess whether the perceived
2535 benefits of big data justify diverting resources from other essential planning tasks. In many
2536 cases, the cost-benefit calculus is unfavorable, especially when planners are uncertain how
2537 much of the data will ultimately be useful for investment decision-making (Bauer et al., 2014).

2538 Some agencies have attempted creative workarounds. For example, one MPO used federal
2539 Transportation Management Area (TMA) funds to subsidize vendor subscriptions. Yet even

2540 when funding is available, the lack of transparent pricing from vendors can deter adoption.
2541 Opaque service-based fee structures make it difficult to anticipate long-term costs or compare
2542 alternatives (KimleyHorn, 2021).

2543 Perceptions of cost are further shaped by technical and staffing needs. Cloud-based
2544 platforms may be viewed as more expensive than local storage options, partly because
2545 they often require hiring or contracting with specialized personnel (i.e., data analysts and
2546 engineers) to securely manage data acquisition and processing (Pecheux et al., 2020; Singh
2547 et al., 2022). Moreover, big data platforms often provide far more information than agencies
2548 can immediately use. Filtering that data down into actionable insights requires significant
2549 staff time and effort, adding another layer of hidden cost (Bauer et al., 2014).

2550 Furthermore, the successful use of big data within MPOs depends not only on acquiring
2551 the data itself, but also on having the staff capacity to manage, interpret, and apply it
2552 effectively. Implementing big data typically requires new information technology systems,
2553 staff training, and sometimes even revised hiring strategies (KimleyHorn, 2021). Generally,
2554 however, smaller transportation agencies lack the dedicated personnel needed to carry out
2555 these tasks. Specialized roles such as data engineers, database administrators, or GIS analysts
2556 may be absent, or distributed across partner organizations without centralized oversight
2557 (Richardson, 2021).

2558 Agencies also report difficulty recruiting and retaining staff with experience in big data
2559 management. Personnel trained in traditional data systems are more widely available and
2560 better understood by leadership, while staff familiar with modern analytics pipelines, cloud-
2561 based platforms, or large-scale data integration are scarcer and harder to retain. Furthermore,
2562 siloed data cultures and internal concerns over data security or potential misuse often inhibit
2563 both intra-agency collaboration and external data sharing, limiting the agency's ability to
2564 adopt modern data practices (Pecheux et al., 2020).

2565 Training and knowledge sharing also matter. Some MPOs have learned from peer agencies
2566 or attended workshops (e.g., FHWA's training on using archived operations data FHWA 2020)
2567 which build confidence to adopt big data. Nonetheless, in many agencies, the lack of staff
2568 time and resources remains a core constraint. Without sustained investment in staffing and
2569 training, and without leadership recognition of the importance of data governance, agencies
2570 are unlikely to move beyond pilot projects or basic uses of big data.

2571 Finally, leadership support plays a pivotal role in whether big data tools are adopted
2572 within MPOs. If executives or board members do not understand or trust the value of big

2573 data, planning staff often struggle to move initiatives forward. A national review noted that
2574 “leadership often does not fully understand the value of big data,” which creates institutional
2575 hesitation (Pecheux et al., 2020). In many cases, skepticism stems from doubts about
2576 data quality, limited familiarity with new technologies, or uncertainty about how big data
2577 contributes to planning outcomes (Richardson, 2021).

2578 Without a visible champion for data sharing at the executive level, it can be difficult for
2579 staff to secure the resources and flexibility needed to experiment with or integrate new data
2580 products. Even when analysts make the case for cloud-based platforms or advanced analytics
2581 tools, they often fail to demonstrate an immediate return on investment to decision-makers,
2582 especially in agencies that prioritize cost-efficiency or regulatory compliance (Pecheux et al.,
2583 2020).

2584 Indeed, literature on organizational psychology may provide helpful guidance in this
2585 context. Teams adopt new practices more readily when leaders create psychological safety
2586 (i.e., a shared belief that it is safe to take interpersonal risks) because staff can experiment,
2587 surface errors, and learn without fear of blame. This enables careful piloting of unfamiliar
2588 data products (Edmondson, 1999). Agencies also need absorptive capacity, the ability to
2589 recognize the value of external knowledge, assimilate it, and apply it. Leadership investment
2590 in related skills and routines strengthens this capacity (Cohen and Levinthal, 1990). Finally,
2591 technology-acceptance research shows that executive messaging and support affect user beliefs
2592 about performance gains, required effort, social norms, and facilitating conditions, which
2593 together shape adoption intentions and use (Venkatesh et al., 2003).

2594 **5.3 Methods**

2595 We conducted an interactive workshop at the Association of Metropolitan Planning Organiza-
2596 tions (AMPO) Planning Tools & Training Symposium in Albuquerque, New Mexico, in May
2597 2024. The workshop was attended by over 50 participants, all of whom held roles related to
2598 LRTPs in some capacity.

2599 We used Mentimeter to facilitate real-time engagement and gather background information.
2600 Then, based on the answers to the question, “What best describes your role in your MPO?”
2601 (Table 5.1), we divided the room into three roughly even-sized groups and conducted focus
2602 group discussions tailored to each group’s role and expertise.



Figure 5.1: Bar plot of answers mentioned in response to “In which area does your MPO operate?”

2603 **5.4 Findings**

2604 *5.4.1 Descriptive Statistics of Participants*

2605 The workshop saw participation from MPO staff from across every geographical region in
 2606 the United States, with the Northeast and the Pacific Northwest (PNW) having the most
 2607 participants (see Figure 5.1). When asked about their MPOs’ experience with BDPs, most of
 2608 the participants hinted at some proficiency through in-house usage or via consultants—only
 2609 17 percent of the respondents had no prior hands-on experience with BDPs (see Table 5.1).
 2610 The prior exposure of workshop participants suggests their insights would be more informed
 2611 and valuable to the discussion.

2612 We also asked the respondents about their role within the MPO as shown in Table 5.1.
 2613 They were almost evenly divided between technical staff who primarily work with data (33%),
 2614 managers who primarily supervise people, projects, or consultants (23%), and staff who do
 2615 both (37%). Hearing from both technical and managerial staff is crucial, as their perspectives
 2616 on the use of BDPs may differ. Technical staff are typically more engaged with the practical,
 2617 day-to-day data challenges, while managerial staff are more focused on broader strategic
 2618 concerns such as resource allocation and decision-making processes. Both viewpoints offer
 2619 complementary insights into how BDPs can be integrated into MPO operations.

2620 We note that the goal of the workshop was not necessarily to attain a representative

Table 5.1: Distribution of answers to exploratory questions about the audience

Answer	Percentage*
”What experience does your MPO have with big data?”	
None.	7%
We’re learning about it but don’t have hands-on experience.	10%
We’ve experimented with it through trial subscriptions or purchases.	17%
We’ve used it in-house or through consultants to support work programs or planning/policy decisions.	66%
”What best describes your role in your MPO?”	
I primarily manage people, projects, or consultants; focus on the big picture; and do NOT handle data.	23%
I primarily work with data to support my MPO’s work programs and projects.	33%
I perform some management, consider the big picture, AND directly process/analyze data.	37%
I perform other roles that little to do with data.	7%

**rounded to nearest integer*

2621 sample of MPOs, but rather to present a wide range of perspectives that exist.

2622 5.4.2 Perceptions of Big Data Products

2623 Most planners view big data as complementary to traditional data collection methods—neither
 2624 a complete replacement nor without value. During the focus group discussions, we found
 2625 some differences in how managers perceived BDPs compared to that of technical MPO
 2626 staff. Managers perceived big data as a tool for updating functional classification of various
 2627 roadways and doing related analyses, such as deriving mode splits, normalizing crash rates, and

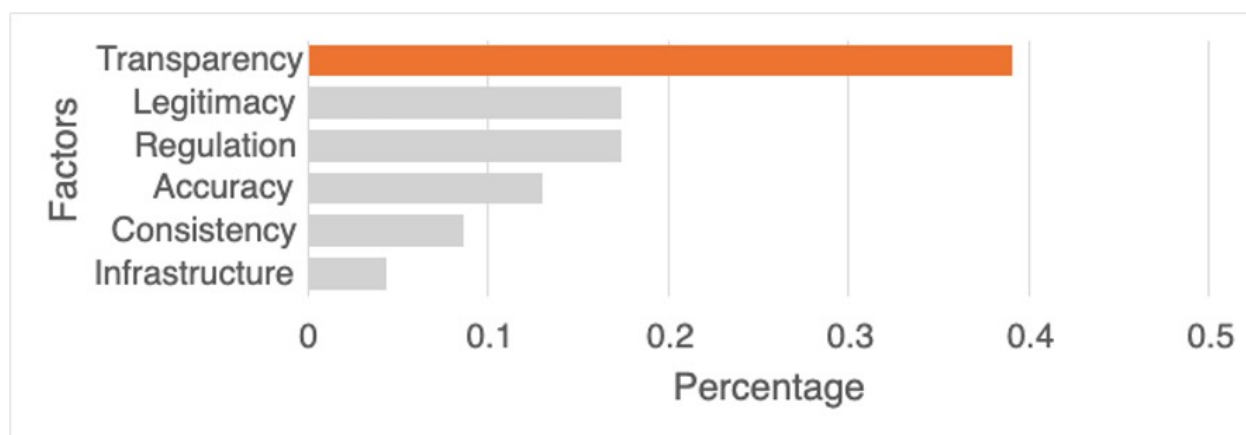


Figure 5.3: Bar plot of factors mentioned in response to “What would help you gain confidence or have trust in using big data?”

2640 academic studies (thereby reducing risk). Regulation ranked third, referring to data quality
 2641 guidance from transportation authorities and decision makers with vested interest.

2642 As a follow-up, we asked participants what they would like to know about BDPs that
 2643 they do not already know. The responses included questions regarding the quality of the
 2644 data like “How accurate is it?” and “Is the accuracy consistent over time?”, questions that
 2645 expand on transparency from the supplier side like “What steps were taken between the raw
 2646 and processed datasets?”, as well as operational-level questions about the process of adopting
 2647 BDPs like “Will costs change over time?” and “How do we ‘keep’ the data we buy after our
 2648 contract ends?”.

2649 5.4.4 Use Cases

2650 We found a wide range of BDP use cases among MPOs, as detailed in Table 5.2. However,
 2651 the majority of LRTPs focus on a small number of BDP applications. The most frequently
 2652 mentioned benefits of partnering with mobility-analytics companies were the ability to easily
 2653 derive peak traffic volumes for low functional classification roads, understand turn counts
 2654 at intersections, and measure pedestrian and cyclist activity. This seemed to agree with
 2655 our hypothesis that many MPOs who adopt BDPs only utilize them for a narrow scope of
 2656 problems.

2657 Among the use cases shown in Table 5.2, we found little-to-no prior literature on using big
 2658 data to study work-housing imbalances and job accessibility. Other equity-related applications,

2659 such as identifying locations of high noise and vehicle emissions on highway networks and
2660 assessing community resiliency, were similarly under-explored. More research in these areas
2661 could provide valuable insights to address transportation inequities and improve public health
2662 and community resilience.

Table 5.2: Types of questions MPO staffers have solved or plan on solving, categorized

Category	Short Description	Long Description
Equity-related	Quantifying work-housing imbalance and job accessibility (Zhang et al., 2017; Zhou et al., 2018)	Analyzing the distribution of jobs relative to housing to ensure equitable access to employment opportunities.
	Identifying locations of high noise and vehicle emissions on highway network (Lan et al., 2020)	Locating areas with high levels of noise and emissions to address environmental justice concerns.
	Assessing community resilience (Sarker et al., 2020)	Assessing the transportation network's ability to withstand and recover from adverse events such as natural disasters.
	Origin-destination studies (Alexander et al., 2015b; Iqbal et al., 2014b; Liu et al., 2022b)	Studying travel patterns between different origins and destinations to understand mobility flows.
	Mode choice analysis (Gong et al., 2012; Phithakitnukoon et al., 2017; Yang et al., 2022)	Analyzing the factors influencing the selection of different modes of transport (e.g., car, public transit).
Safety-related	Safety planning and project identification (Ambros et al., 2024)	Identifying and prioritizing safety projects to reduce accidents and enhance road safety.
	Crash hotspot analysis (Xie et al., 2017)	Pinpointing locations with high accident rates to implement targeted safety measures.
	Mobile device usage while driving (Ahlström et al., 2020; Khurana and Goel, 2020)	Monitoring mobile device usage to identify distracted driving behaviors and develop interventions.
	Lane change analysis (Park et al., 2018)	Analyzing lane change behaviors to improve road design and traffic flow management.

Planning-related	Complete Street studies (Bian et al., 2023)	Planning and designing streets that accommodate all users, including pedestrians, cyclists, and motorists.
	Visitor analysis (Miah et al., 2017 ; Reif and Schmücker, 2020)	Understanding the travel behaviors and patterns of visitors to optimize transportation services for tourism.
	Validation of observed travel behavior data (Chen et al., 2016b ; Liu et al., 2014 ; Uğurel et al., 2024)	Comparing big data insights with observed travel behaviors to ensure accuracy and reliability.
	EV charging demand estimation (Yang et al., 2017b)	Forecasting the demand for electric vehicle charging stations to support the growing number of EVs.
	Transit planning and routing (Hadjidimitriou et al., 2021 ; Lu et al., 2021)	Designing and optimizing public transit routes to improve efficiency and coverage.
	Bicycle and pedestrian trip planning (Lee and Sener, 2020 ; Yu et al., 2020)	Tracking bicycle and pedestrian trips to support infrastructure improvements and promote active transportation.
	Freight movement tracking (Akter et al., 2023)	Monitoring the movement of goods to enhance freight logistics and reduce congestion.
	Land use changes (Yin et al., 2021)	Assessing the impact of land use changes on transportation networks and planning accordingly.
Operational-related	Land border crossing monitoring (Sakhare et al., 2024)	Providing real-time information on wait times and traffic conditions at land border crossings.
	Through trip volume estimation (Huntsinger and Ward, 2015)	Measuring the volume of trips passing through a region without stopping to better understand transit patterns.
	Turn count estimation (Kan et al., 2019)	Counting vehicle turns at intersections to optimize signal timing and improve traffic flow.
	Demand estimation during special events (Pereira et al., 2015)	Analyzing transportation demand and patterns during special events to enhance planning and resource allocation.

Chapter 6

CONCLUSION, DISCUSSION, AND FUTURE WORK

2663

2664

2665 This dissertation examines how behaviorally informed machine learning methods can
2666 improve the reliability of mobility data and, in turn, support more transparent and defensible
2667 transportation planning decisions. Chapter 2 demonstrates that both short- and long-duration
2668 gaps in GPS traces can be addressed using a Bayesian learning framework that explicitly
2669 captures the temporal periodicities underlying human travel patterns. Chapter 3 extends
2670 this idea by constructing individualized function spaces, enabling the generation of synthetic
2671 mobility data that more faithfully replicates observed behavior. Chapter 4 then investigates
2672 the challenges of predicting sociodemographic attributes from mobility signals, showing
2673 that limited sample sizes and distribution shifts across time and geography impose inherent
2674 performance ceilings and underscore the need for uncertainty quantification and alternative
2675 learning strategies such as multi-task learning. Finally, Chapter 5 draws on a qualitative
2676 workshop with long-range transportation planners to document practical use cases, points of
2677 friction, and the institutional and organizational barriers that shape the adoption of big data
2678 products in real planning workflows.

2679

2680 This chapter situates these contributions within the broader landscape of transportation
2681 planning and engineering. It examines data governance and transparency challenges that
2682 arise when relying on proprietary pipelines, the growing mismatch between vendor practices
2683 and the scientific evidence surrounding imputation, uncertainty, and representativeness, and
2684 the implications of PCM based products for travel demand modeling and forecasting. It
2685 also highlights the increasing importance of ML literacy for planners and engineers as big
2686 data tools become embedded in everyday workflows. The chapter concludes by outlining
2687 several avenues for future research, including the development of foundation models for
2688 human mobility, methods to expand the effective amount of training data, privacy preserving
2689 computational frameworks, and the institutional infrastructure needed to support responsible
and scalable adoption of big data in public sector planning.

2690 **6.1 Discussion**

2691 A central implication of this work is the need for stronger data governance frameworks that
2692 support transparency without compromising privacy. The individualized models proposed in
2693 Chapter 3 can be shared through privacy-preserving protocols that avoid exposing raw traces
2694 while still enabling reproducibility and third-party auditing. Establishing such protocols
2695 would give agencies more confidence in relying on these tools and would facilitate comparative
2696 evaluations across regions. In the longer term, the field may benefit from a national mobility
2697 data baseline: a shared, standardized, and continually updated data pool that public agencies
2698 and researchers can build upon. This idea, analogous to the National Weather Service’s
2699 mandate, would lower entry barriers, reduce dependence on proprietary pipelines, and create
2700 an open benchmark for evaluating vendor products.

2701 The findings in this dissertation also point to a growing mismatch between vendor practices
2702 and the scientific evidence on how mobility data should be handled. Many commercial
2703 providers treat imputation as a purely computational task. In contrast, Chapters 2 and
2704 3 show that transforming raw GPS traces into behaviorally coherent trajectories is far
2705 from straightforward. It requires explicit assumptions about trip-making rhythms, anchor
2706 point structures, daily cycles, and route feasibility. When these assumptions are hidden
2707 behind proprietary code, agencies effectively lose control over the behavioral content of the
2708 datasets they rely on. Similarly, vendors often distribute aggregate traffic volumes with no
2709 associated uncertainty and no demographic auditing, masking representativeness biases that
2710 can distort downstream analyses. The methods developed here demonstrate that uncertainty
2711 and demographic skew are not nuisances to suppress but properties that must be quantified,
2712 communicated, and carried forward through planning workflows. Ignoring uncertainty can
2713 encourage a false sense of precision and can skew policy decisions toward overconfident
2714 interpretations of noisy signals.

2715 Chapter 4 speaks directly to the issue of representativeness. It shows that accurate
2716 prediction of sociodemographic attributes such as age, gender, income, and household structure
2717 is constrained by small survey samples and distribution shifts across time and geography.
2718 There is a ceiling on achievable accuracy even under careful model design, which means that
2719 inferred labels should be interpreted probabilistically rather than as ground truth. At the
2720 same time, these models provide a practical tool for diagnosing who is under represented in
2721 PCM datasets. By comparing the predicted sociodemographic composition of PCM users
2722 to survey or census benchmarks, agencies can identify systematic blind spots and design

2723 weighting schemes, targeted outreach, or complementary data collections. In this sense,
2724 sociodemographic inference is less a silver bullet fix for bias than a diagnostic layer that
2725 makes representativeness issues visible and actionable.

2726 These methodological advances also have direct implications for travel demand modeling
2727 and forecasting. More complete and behaviorally coherent trajectories, as developed in
2728 Chapters 2 and 3, can strengthen activity based models by improving how daily rhythms,
2729 tour structures, and multi stop chains are represented and calibrated. PCM derived traces
2730 and synthetic data can be aggregated into OD flows that capture temporal and spatial
2731 heterogeneity at resolutions that traditional surveys cannot support. This enables a wider
2732 range of applications, including anomaly detection, special event demand forecasting, and
2733 quasi experimental evaluations of network changes or policy interventions using before and
2734 after PCM signals. Incorporating uncertainty from imputation and sociodemographic inference
2735 into these models would allow long range forecasts and project appraisals to reflect the true
2736 confidence bounds of the underlying data, rather than treating datasets as fixed and error
2737 free.

2738 Finally, these discussions and some of the results we show in Chapter 5 suggest that
2739 ML literacy must become a core competency within transportation agencies. The goal is
2740 not for planners or engineers to train new neural networks, but to critically interrogate the
2741 assumptions, data transformations, confidence scores, and representativeness characteristics
2742 of the products they procure. As PCM-based tools become more influential in shaping
2743 planning analyses and downstream investments, the profession needs a workforce capable
2744 of evaluating models not only on accuracy but on transparency, behavioral validity, and
2745 uncertainty quantification.

2746 6.2 Future Work

2747 This dissertation opens several promising directions for advancing behaviorally informed
2748 mobility modeling and its integration into planning practice. In the era of large language
2749 models, a natural next step is the development of foundation models for human mobility.
2750 Emerging large multimodal transformers can jointly encode mobility traces, points of interest,
2751 text-rich contextual data, and sociodemographic attributes, offering a flexible backbone
2752 for a wide range of downstream tasks (Wu et al., 2024b; Zhang et al., 2024b). Training
2753 such models on continental-scale datasets and fine-tuning them for specific metropolitan
2754 regions would provide a shared representational layer that captures both universal and region-
2755 specific mobility structure. Embedding privacy-preserving mechanisms, such as differential
2756 privacy during training or federated learning across data providers, will be essential to ensure
2757 that these models are deployable in environments where data sensitivity and institutional
2758 constraints remain significant.

2759 At the modeling level, a key opportunity lies in the integration of PCM-enhanced datasets
2760 with public-sector ABMs. Many existing ABMs rely on synthetic populations built from
2761 infrequent surveys or censuses (Brinckerhoff, 2010), which age quickly and struggle to reflect
2762 nuanced telecommuting patterns or emerging modes of travel. Calibration is often limited by
2763 sparse validation data and by aggregate count measures that hide within-day heterogeneity.
2764 PCM-derived trip data can help address these gaps by capturing temporal patterns, tour
2765 structures, and neighborhood-level variability that surveys may miss (Tozluoğlu et al., 2025).
2766 Future work should develop methods for fusing PCM with HTS-based populations, updating
2767 or reweighting agents as conditions evolve, and incorporating PCM-based OD estimates into
2768 scenario analysis and network design. An important piece of this agenda will be explicit
2769 uncertainty propagation, so that ABM forecasts reflect both structural model uncertainty
2770 and the noise inherent in high-frequency mobility data.

2771 The practical deployment of these models raises important questions around privacy-
2772 preserving computation. Federated learning offers a pathway for training individualized
2773 imputation or inference models without centralizing raw traces (Ezequiel et al., 2022), and
2774 on-device imputation and feature extraction can reduce privacy risks by ensuring that only
2775 aggregated or derived quantities leave the device (Qi et al., 2017). These approaches, however,
2776 come with tradeoffs (Charles and Konečný, 2021). Centralized training is often simpler to
2777 implement, can be more data efficient, and may achieve higher accuracy when data are highly
2778 unbalanced across participants. Distributed frameworks introduce communication overhead

2779 and coordination challenges when data are not independent across clients. They also face
2780 cybersecurity risks such as model inversion and poisoning attacks, which can undermine
2781 both privacy and data integrity (Wang et al., 2024). Future work will need to clarify when
2782 federated or on-device approaches offer meaningful privacy gains over well-governed centralized
2783 systems, and how to design architectures, legal agreements, and monitoring tools that balance
2784 competing interests (i.e., utility vs. privacy) in public-sector planning contexts.

2785 Another question that warrants deeper investigation is how to overcome the fundamental
2786 sample size limitations of rich behavioral datasets such as household travel surveys. As shown
2787 in Chapter 4, predictive performance tends to have a moderate upper bound when trained
2788 solely on HTS data, reflecting both the small size of these surveys and their limited coverage
2789 of the joint mobility-sociodemographic space. Semi-supervised learning (SSL) methods,
2790 specifically *self-training*, offer a principled path forward by expanding the effective training
2791 set through the incorporation of pseudo-labeled observations (Lee, 2013; Van Rooyen and
2792 Williamson, 2018). Building calibrated, uncertainty-aware SSL pipelines could meaningfully
2793 enhance our ability to enrich PCM datasets with respondent-like attributes, thereby improving
2794 downstream planning analyses without requiring substantially larger surveys.

2795 Finally, an obvious research direction involves the institutional infrastructure surrounding
2796 big data products. As PCM datasets become embedded in planning workflows, indepen-
2797 dent audits of vendor products, much like financial or algorithmic audits in other sectors
2798 (Costanza-Chock et al., 2022; Liu et al., 2022a), will be critical. These audits should eval-
2799 uate data completeness, imputation logic, uncertainty communication, and demographic
2800 representativeness, providing agencies with an evidence-based basis for procurement and
2801 use. The methodological tools introduced in this dissertation offer building blocks for such a
2802 framework.

BIBLIOGRAPHY

- Abkowitz, M. D. (1981). An analysis of the commuter departure time decision. *Transportation* 10, 283–297
- Ahlström, C., Wachtmeister, J., Nyman, M., Nordenström, A., and Kircher, K. (2020). Using smartphone logging to gain insight about phone use in traffic. *Cognition, Technology & Work* 22, 181–191
- Ak, C., Ergönül, Ö., and Gönen, M. (2018). Structured gaussian processes with twin multiple kernel learning. In *Asian Conference on Machine Learning (ACML)*. 65–80
- Akter, T., Hernandez, S., and Camargo, P. V. (2023). Freight Operational Characteristics Mined from Anonymous Mobile Sensor Data. *Transportation Research Record* 2677, 236–247. doi:10.1177/03611981231158639. Publisher: SAGE Publications Inc
- Alberini, A., Burra, L. T., Cirillo, C., and Shen, C. (2021). Counting vehicle miles traveled: What can we learn from the nhts? *Transportation research part D: transport and environment* 98, 102984
- Alessandretti, L., Aslak, U., and Lehmann, S. (2020). The scales of human mobility. *Nature* 587, 402–407
- Alessandretti, L., Sapiezynski, P., Sekara, V., Lehmann, S., and Baronchelli, A. (2018). Evidence for a conserved quantity in human mobility. *Nature Human Behaviour* 2, 485–491. doi:10.1038/s41562-018-0364-x. Publisher: Nature Publishing Group
- Alexander, L., Jiang, S., Murga, M., and González, M. C. (2015a). Origin–destination trips by purpose and time of day inferred from mobile phone data. *Transportation Research Part C: Emerging Technologies* 58, 240–250. doi:10.1016/j.trc.2015.02.018

- Alexander, L., Jiang, S., Murga, M., and González, M. C. (2015b). Origin–destination trips by purpose and time of day inferred from mobile phone data. *Transportation Research Part C: Emerging Technologies* 58, 240–250. doi:10.1016/j.trc.2015.02.018
- Alvarez, M. A., Luengo, D., and Lawrence, N. D. (2013). Linear latent force models using gaussian processes. *IEEE transactions on pattern analysis and machine intelligence* 35, 2693–2705
- Ambros, J., Elgner, J., Valentova, V., Bak, R., and Kiec, M. (2024). Proactive safety assessment of urban through-roads based on GPS data. *Archives of Transport* 69, 113–125. doi:10.61089/aot2024.gpa7v104. Number: 1
- Apple (2021). *Apple*
- Auld, J., Mohammadian, A. K., Oliveira, M. S., Wolf, J., and Bachman, W. (2015). Demographic Characterization of Anonymous Trace Travel Data. *Transportation Research Record* 2526, 19–28. doi:10.3141/2526-03. Publisher: SAGE Publications Inc
- Ban, X. J., Chen, C., Wang, F., Wang, J., Zhang, Y., et al. (2018). *Promises of data from emerging technologies for transportation applications: Puget sound region case study*. Tech. rep., United States. Federal Highway Administration
- Barbosa, H., Barthelemy, M., Ghoshal, G., James, C. R., Lenormand, M., Louail, T., et al. (2018). Human mobility: Models and applications. *Physics Reports* 734, 1–74
- Barla, A., Odone, F., and Verri, A. (2003). Histogram intersection kernel for image classification. In *Proceedings 2003 international conference on image processing (Cat. No. 03CH37429)* (IEEE), vol. 3, III–513
- Batista, S. F., Cantelmo, G., Menéndez, M., and Antoniou, C. (2022). A gaussian sampling heuristic estimation model for developing synthetic trip sets. *Computer-Aided Civil and Infrastructure Engineering* 37, 93–109
- Bauer, J., Pack, M., Giragosian, A., Kehoe, N., Evans, J., Voorhies, K., et al. (2014). *Use of Archived Operations Data in Transportation Planning*. Tech. rep., Saxton Transportation Operations Laboratory

- Baxter, J. (2000). A Model of Inductive Bias Learning. *Journal of Artificial Intelligence Research* 12, 149–198. doi:10.1613/jair.731. ArXiv:1106.0245 [cs]
- Bayarma, A., Kitamura, R., and Susilo, Y. O. (2007). Recurrence of daily travel patterns: stochastic process approach to multiday travel behavior. *Transportation Research Record* 2021, 55–63
- Ben-Moshe, B., Elkin, E., Levi, H., and Weissman, A. (2011). Improving accuracy of gns devices in urban canyons. In *CCCG*. 511–515
- Berke, A., Doorley, R., Larson, K., and Moro, E. (2022). Generating synthetic mobility data for a realistic population with RNNs to improve utility and privacy. In *Proceedings of the 37th ACM/SIGAPP Symposium on Applied Computing* (New York, NY, USA: Association for Computing Machinery), SAC '22, 964–967. doi:10.1145/3477314.3507230
- Bhat, C. R. and Koppelman, F. S. (1994). A structural and empirical model of subsistence activity behavior and income. *Transportation* 21, 71–89. doi:10.1007/BF01119635
- Bhat, C. R. and Koppelman, F. S. (1999). *Activity-Based Modeling of Travel Demand* (Boston, MA: Springer US). 35–61. doi:10.1007/978-1-4615-5203-1_3
- Bian, R., Tolford, T., Liu, S., and Gangireddy, S. (2023). Lessons learned from evaluating complete streets project outcomes with emerging data sources. *Transportation Planning and Technology* 46, 754–772. doi:10.1080/03081060.2023.2214136. Publisher: Routledge .eprint: <https://doi.org/10.1080/03081060.2023.2214136>
- Bonilla, E. V., Chai, K., and Williams, C. (2007). Multi-task gaussian process prediction. *Advances in neural information processing systems* 20
- Boyd, T., Dinehart, T., and Williams, K. (2024). *Uses and Limitations of Big Data for Evaluating Transportation Equity*. Tech. rep., National Institute for Congestion Reduction, Tampa, FL. doi:10.5038/cutr-nicr-y3-1-11
- [Dataset] Brinckerhoff, P. (2010). Design and Development Plan for the MAG CT-RAMP Activity-Based Model (ABM)

- Cao, W., Wu, Z., Wang, D., Li, J., and Wu, H. (2016). Automatic user identification method across heterogeneous mobility data sources. In *2016 IEEE 32nd International Conference on Data Engineering (ICDE)*. 978–989. doi:10.1109/ICDE.2016.7498306
- Caruana, R. (1997). Multitask Learning. *Machine Learning*, 41–75doi:10.1023/A:1007379606734
- Charles, Z. and Konečný, J. (2021). Convergence and Accuracy Trade-Offs in Federated Learning and Meta-Learning. In *Proceedings of The 24th International Conference on Artificial Intelligence and Statistics* (PMLR), 2575–2583. ISSN: 2640-3498
- Chen, C., Ma, J., Susilo, Y., Liu, Y., and Wang, M. (2016a). The promises of big data and small data for travel behavior (aka human mobility) analysis. *Transportation Research Part C: Emerging Technologies* 68, 285–299. doi:10.1016/j.trc.2016.04.005
- Chen, C., Ma, J., Susilo, Y., Liu, Y., and Wang, M. (2016b). The promises of big data and small data for travel behavior (aka human mobility) analysis. *Transportation Research Part C: Emerging Technologies* 68, 285–299. doi:10.1016/j.trc.2016.04.005
- Chen, C., Wang, R., Bansal, P., Chen, L., Ugurel, E., Zhang, Y., et al. (2026). From biases to opportunities: leveraging location-based-service (lbs) data for next-generation transportation planning. *Transportation Research Part C: Emerging Technologies* 182, 105416
- Chen, Y., Hosseini, B., Owhadi, H., and Stuart, A. M. (2021). Solving and learning nonlinear pdes with gaussian processes. *Journal of Computational Physics* 447, 110668
- Cipra, T., Trujillo, J., and Robio, A. (1995). Holt-winters method with missing observations. *Management Science* 41, 174–178
- Cohen, W. M. and Levinthal, D. A. (1990). Absorptive Capacity: A New Perspective on Learning and Innovation. *Administrative Science Quarterly* 35, 128–152. doi:10.2307/2393553. Publisher: [Sage Publications, Inc., Johnson Graduate School of Management, Cornell University]

- Costanza-Chock, S., Raji, I. D., and Buolamwini, J. (2022). Who audits the auditors? recommendations from a field scan of the algorithmic auditing ecosystem. In *Proceedings of the 2022 ACM Conference on Fairness, Accountability, and Transparency*. 1571–1583
- Cuomo, S., Di Cola, V. S., Giampaolo, F., Rozza, G., Raissi, M., and Piccialli, F. (2022). Scientific machine learning through physics-informed neural networks: Where we are and what’s next. *Journal of Scientific Computing* 92, 88
- Daisy, N. S., Millward, H., and Liu, L. (2018). Trip chaining and tour mode choice of non-workers grouped by daily activity patterns. *Journal of Transport Geography* 69, 150–162
- Deng, T., Zhang, K., and Shen, Z.-J. M. (2021). A systematic review of a digital twin city: A new pattern of urban governance toward smart cities. *Journal of Management Science and Engineering* 6, 125–134
- Ding, S., Huang, H., Zhao, T., and Fu, X. (2019). Estimating Socioeconomic Status via Temporal-Spatial Mobility Analysis - A Case Study of Smart Card Data. In *2019 28th International Conference on Computer Communication and Networks (ICCCN)*. 1–9. doi: 10.1109/ICCCN.2019.8847051. ISSN: 2637-9430
- Doi, S., Mizuno, T., and Fujiwara, N. (2021). Estimation of socioeconomic attributes from location information. *Journal of Computational Social Science* 4, 187–205. doi: 10.1007/s42001-020-00073-w. PMID: 32838050 PMID: PMC7271143
- Dorfman, R. (1979). A Formula for the Gini Coefficient. *The Review of Economics and Statistics* 61, 146–149. doi:10.2307/1924845. Publisher: The MIT Press
- Douglas, D. H. and Peucker, T. K. (1973). Algorithms for the reduction of the number of points required to represent a digitized line or its caricature. *Cartographica: the international journal for geographic information and geovisualization* 10, 112–122
- Duvenaud, D. (2014). *Automatic model construction with Gaussian processes*. Ph.D. thesis, University of Cambridge

- Duvenaud, D., Lloyd, J., Grosse, R., Tenenbaum, J., and Zoubin, G. (2013). Structure discovery in nonparametric regression through compositional kernel search. In *International Conference on Machine Learning* (PMLR), 1166–1174
- Eagle, N. and Pentland, A. S. (2009). Eigenbehaviors: Identifying structure in routine. *Behavioral ecology and sociobiology* 63, 1057–1066
- Edmondson, A. (1999). Psychological Safety and Learning Behavior in Work Teams. *Administrative Science Quarterly* 44, 350–383. doi:10.2307/2666999. Publisher: SAGE Publications Inc
- Ellegard, K., Hagerstrand, T., and Lenntorp, B. (1977). Activity Organization and the Generation of Daily Travel: Two Future Alternatives. *Economic Geography* 53, 126. doi:10.2307/142721
- Ezequiel, C. E. J., Gjoreski, M., and Langheinrich, M. (2022). Federated Learning for Privacy-Aware Human Mobility Modeling. *Frontiers in Artificial Intelligence* 5. doi:10.3389/frai.2022.867046. Publisher: Frontiers
- Feng, J., Yang, Z., Xu, F., Yu, H., Wang, M., and Li, Y. (2020). Learning to Simulate Human Mobility. In *Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining* (New York, NY, USA: Association for Computing Machinery), KDD '20, 3426–3433. doi:10.1145/3394486.3412862
- [Dataset] FHWA (2020). Applying Archived Operations Data in Transportation Planning
- [Dataset] Fiedler, D., Čáp, M., Nykl, J., Žilecký, P., and Schaefer, M. (2019). Map Matching Algorithm for Large-scale Datasets. ArXiv:1910.05312 [cs, eess] version: 1
- Frias-Martinez, V., Soguero, C., and Frias-Martinez, E. (2012). Estimation of urban commuting patterns using cellphone network data. In *Proceedings of the ACM SIGKDD international workshop on urban computing*. 9–16
- Gammelli, D., Peled, I., Rodrigues, F., Pacino, D., Kurtaran, H. A., and Pereira, F. C. (2020). Estimating latent demand of shared mobility through censored gaussian processes. *Transportation Research Part C: Emerging Technologies* 120, 102775

- Gao, L., Huang, H., Ye, J., and Wang, D. (2024). Activity type detection of mobile phone data based on self-training: Application of the teacher–student cycling model. *Transportation Research Part C: Emerging Technologies* 161, 104550. doi:10.1016/j.trc.2024.104550
- Gardner, J., Pleiss, G., Weinberger, K. Q., Bindel, D., and Wilson, A. G. (2018). Gpytorch: Blackbox matrix-matrix gaussian process inference with gpu acceleration. *Advances in neural information processing systems* 31
- Gibbs, M. N. (1998). *Bayesian Gaussian processes for regression and classification*. Ph.D. thesis, Citeseer
- Gneiting, T. and Raftery, A. E. (2007). Strictly Proper Scoring Rules, Prediction, and Estimation. *Journal of the American Statistical Association* 102, 359–378. doi:10.1198/016214506000001437. Publisher: Informa UK Limited
- Gönen, M. and Alpaydın, E. (2011). Multiple kernel learning algorithms. *The Journal of Machine Learning Research* 12, 2211–2268
- Gong, H., Chen, C., Bialostozky, E., and Lawson, C. T. (2012). A GPS/GIS method for travel mode detection in New York City. *Computers, Environment and Urban Systems* 36, 131–139. doi:10.1016/j.compenvurbsys.2011.05.003
- Gong, X., Huang, Z., Wang, Y., Wu, L., and Liu, Y. (2020). High-performance spatiotemporal trajectory matching across heterogeneous data sources. *Future Generation Computer Systems* 105, 148–161
- Gonzalez, M. C., Hidalgo, C. A., and Barabasi, A.-L. (2008). Understanding individual human mobility patterns. *nature* 453, 779–782
- Hadjidimitriou, N. S., Lippi, M., and Mamei, M. (2021). A Data Driven Approach to Match Demand and Supply for Public Transport Planning. *IEEE Transactions on Intelligent Transportation Systems* 22, 6384–6394. doi:10.1109/TITS.2020.2991834. Conference Name: IEEE Transactions on Intelligent Transportation Systems

- Hand, D. J. and Till, R. J. (2001). A Simple Generalisation of the Area Under the ROC Curve for Multiple Class Classification Problems. *Machine Learning* 45, 171–186. doi:10.1023/A:1010920819831
- Hanson, S. (1980). The importance of the multi-purpose journey to work in urban travel behavior. *Transportation* 9, 229–248. doi:10.1007/BF00153866
- Hanson, S. and Huff, J. O. (1988). Repetition and day-to-day variability in individual travel patterns: Implications for classification. *Behavioural Modelling in Geography and Planning, Croom Helm, London* , 368–398
- Hao, Q., Chen, L., Xu, F., and Li, Y. (2020). Understanding the urban pandemic spreading of covid-19 with real world mobility data. In *Proceedings of the 26th ACM SIGKDD international conference on knowledge discovery & data mining*. 3485–3492
- Harris, C. R., Millman, K. J., van der Walt, S. J., Gommers, R., Virtanen, P., Cournapeau, D., et al. (2020). Array programming with NumPy. *Nature* 585, 357–362. doi:10.1038/s41586-020-2649-2
- Hastie, T., Friedman, J., and Tibshirani, R. (2001). *The Elements of Statistical Learning*. Springer Series in Statistics (New York, NY: Springer New York). DOI: 10.1007/978-0-387-21606-5 ISSN: 0172-7397
- Hills, T. T., Todd, P. M., Lazer, D., Redish, A. D., and Couzin, I. D. (2015). Exploration versus exploitation in space, mind, and society. *Trends in cognitive sciences* 19, 46–54
- Hochreiter, S. and Schmidhuber, J. (1997). Long short-term memory. *Neural computation* 9, 1735–1780
- Hunter, J. D. (2007). Matplotlib: A 2d graphics environment. *Computing in Science & Engineering* 9, 90–95. doi:10.1109/MCSE.2007.55
- Huntsinger, L. F. and Ward, K. (2015). Using Mobile Phone Location Data to Develop External Trip Models. *Transportation Research Record* 2499, 25–32. doi:10.3141/2499-04. Publisher: SAGE Publications Inc

- Huo, J., Cox, C. D., Seaver, W. L., Robinson, R. B., and Jiang, Y. (2010). Application of two-directional time series models to replace missing data. *Journal of Environmental Engineering* 136, 435–443
- Idé, T. and Kato, S. (2009). Travel-time prediction using gaussian process regression: A trajectory-based approach. In *Proceedings of the 2009 SIAM International Conference on Data Mining* (SIAM), 1185–1196
- Iqbal, M. S., Choudhury, C. F., Wang, P., and González, M. C. (2014a). Development of origin–destination matrices using mobile phone call data. *Transportation Research Part C: Emerging Technologies* 40, 63–74. doi:10.1016/j.trc.2014.01.002
- Iqbal, M. S., Choudhury, C. F., Wang, P., and González, M. C. (2014b). Development of origin–destination matrices using mobile phone call data. *Transportation Research Part C: Emerging Technologies* 40, 63–74. doi:10.1016/j.trc.2014.01.002
- Jahani, E., Sundsøy, P., Bjelland, J., Bengtsson, L., Pentland, A. t., and De Montjoye, Y.-A. (2017). Improving official statistics in emerging markets using machine learning and mobile phone data. *EPJ Data Science* 6, 3. doi:10.1140/epjds/s13688-017-0099-3
- Jiang, L., Chen, C.-X., and Chen, C. (2023). L2mm: learning to map matching with deep models for low-quality gps trajectory data. *ACM Transactions on Knowledge Discovery from Data* 17, 1–25
- Kaiser, M. (2008). Mean clustering coefficients: the role of isolated nodes and leafs on clustering measures for small-world networks. *New Journal of Physics* 10, 083042. doi:10.1088/1367-2630/10/8/083042
- Kan, Z., Tang, L., Kwan, M.-P., Ren, C., Liu, D., and Li, Q. (2019). Traffic congestion analysis at the turn level using Taxis’ GPS trajectory data. *Computers, Environment and Urban Systems* 74, 229–243. doi:10.1016/j.compenvurbsys.2018.11.007
- Kapp, A., Hansmeyer, J., and Mihaljević, H. (2023). Generative Models for Synthetic Urban Mobility Data: A Systematic Literature Review. *ACM Comput. Surv.* 56, 93:1–93:37. doi:10.1145/3610224

- Khurana, R. and Goel, M. (2020). Eyes on the Road: Detecting Phone Usage by Drivers Using On-Device Cameras. In *Proceedings of the 2020 CHI Conference on Human Factors in Computing Systems* (New York, NY, USA: Association for Computing Machinery), CHI '20, 1–11. doi:10.1145/3313831.3376822
- Kim, D., Park, K., Park, Y., and Ahn, J.-H. (2019). Willingness to provide personal information: Perspective of privacy calculus in iot services. *Computers in Human Behavior* 92, 273–281
- Kim, E.-K. and MacEachren, A. (2014). An index for characterizing spatial bursts of movements: A case study with geo-located twitter data. In *GIScience 2014 Workshop on Analysis of Movement Data* (Citeseer)
- KimleyHorn (2021). *Framework for Evaluating Big Data in Regional Travel and Mobility Analyses*. Tech. rep.
- Kingma, D. P. and Ba, J. (2014). Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*
- Kitamura, R. (1988). An evaluation of activity-based travel analysis. *Transportation* 15, 9–34
- Kitamura, R. and Kermanshah, M. (1983). Identifying time and history dependencies of activity choice. *Transportation Research Record* 944, 22–30
- Kitamura, R. and Van Der Hoorn, T. (1987). Regularity and irreversibility of weekly travel behavior. *Transportation* 14, 227–251
- Kohn, R. and Ansley, C. F. (1986). Estimation, prediction, and interpolation for arima models with missing data. *Journal of the American statistical Association* 81, 751–761
- Kondor, D., Hashemian, B., de Montjoye, Y.-A., and Ratti, C. (2020). Towards Matching User Mobility Traces in Large-Scale Datasets. *IEEE Transactions on Big Data* 6, 714–726. doi:10.1109/TBDDATA.2018.2871693

- Kroesen, M. (2014). Modeling the behavioral determinants of travel behavior: An application of latent transition analysis. *Transportation Research Part A: Policy and Practice* 65, 56–67. doi:10.1016/j.tra.2014.04.010
- Kulkarni, V. and Garbinato, B. (2017). Generating synthetic mobility traffic using RNNs. In *Proceedings of the 1st Workshop on Artificial Intelligence and Deep Learning for Geographic Knowledge Discovery* (New York, NY, USA: Association for Computing Machinery), GeoAI '17, 1–4. doi:10.1145/3149808.3149809
- Lan, Z., He, C., and Cai, M. (2020). Urban road traffic noise spatiotemporal distribution mapping using multisource data. *Transportation Research Part D: Transport and Environment* 82, 102323. doi:10.1016/j.trd.2020.102323. ADS Bibcode: 2020TRPD...8202323L
- Lapidus, L. and Pinder, G. F. (1999). *Numerical solution of partial differential equations in science and engineering* (John Wiley & Sons)
- Lasserre, J. A., Bishop, C. M., and Minka, T. P. (2006). Principled hybrids of generative and discriminative models. In *2006 IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR'06)* (IEEE), vol. 1, 87–94
- Le, T. V., Oentaryo, R., Liu, S., and Lau, H. C. (2016). Local gaussian processes for efficient fine-grained traffic speed prediction. *IEEE Transactions on Big Data* 3, 194–207
- LeCun, Y., Bengio, Y., and Hinton, G. (2015). Deep learning. *Nature* 521, 436–444. doi: 10.1038/nature14539. Publisher: Nature Publishing Group
- Lee, D.-H. (2013). Pseudo-Label : The Simple and Efficient Semi-Supervised Learning Method for Deep Neural Networks
- Lee, J., Bahri, Y., Novak, R., Schoenholz, S. S., Pennington, J., and Sohl-Dickstein, J. (2017). Deep neural networks as gaussian processes. *arXiv preprint arXiv:1711.00165*
- Lee, K. and Sener, I. N. (2020). Emerging data for pedestrian and bicycle monitoring: Sources and applications. *Transportation Research Interdisciplinary Perspectives* 4, 100095. doi:10.1016/j.trip.2020.100095

- Lee, M. and McNally, M. G. (2006). An empirical investigation on the dynamic processes of activity scheduling and trip chaining. *Transportation* 33, 553–565
- Lee, M. S. and McNally, M. G. (2003). On the structure of weekly activity/travel patterns. *Transportation Research Part A: Policy and Practice* 37, 823–839
- Lee, Y., Hickman, M., and Washington, S. (2007). Household type and structure, time-use pattern, and trip-chaining behavior. *Transportation Research Part A: Policy and Practice* 41, 1004–1020. doi:10.1016/j.tra.2007.06.007
- Li, M., Wang, K., and Nurul Habib, K. (2025). Investigating recall & proxy bias in household travel surveys under the core-satellite fusion paradigm. *Transportation* doi:10.1007/s11116-025-10654-1
- Li, Z., Ning, H., Jing, F., and Lessani, M. N. (2024). Understanding the bias of mobile location data across spatial scales and over time: A comprehensive analysis of SafeGraph data in the United States. *PLOS ONE* 19, e0294430. doi:10.1371/journal.pone.0294430. Publisher: Public Library of Science
- Liu, F., Janssens, D., Cui, J., Wang, Y., Wets, G., and Cools, M. (2014). Building a validation measure for activity-based transportation models based on mobile phone data. *Expert Systems with Applications* 41, 6174–6189. doi:10.1016/j.eswa.2014.03.054
- Liu, G. and Onnela, J.-P. (2021). Bidirectional imputation of spatial gps trajectories with missingness using sparse online gaussian process. *Journal of the American Medical Informatics Association* 28, 1777–1784
- Liu, X., Glocker, B., McCradden, M. M., Ghassemi, M., Denniston, A. K., and Oakden-Rayner, L. (2022a). The medical algorithmic audit. *The Lancet Digital Health* 4, e384–e397
- Liu, Z., Liu, Z., and Fu, X. (2022b). Dynamic Origin-Destination Flow Prediction Using Spatial-Temporal Graph Convolution Network With Mobile Phone Data. *IEEE Intelligent Transportation Systems Magazine* 14, 147–161. doi:10.1109/MITS.2021.3082397. Conference Name: IEEE Intelligent Transportation Systems Magazine

- Liu, Z., Lyu, C., Huo, J., Wang, S., and Chen, J. (2022c). Gaussian process regression for transportation system estimation and prediction problems: The deformation and a hat kernel. *IEEE Transactions on Intelligent Transportation Systems* 23, 22331–22342
- Lu, K., Liu, J., Zhou, X., and Han, B. (2021). A Review of Big Data Applications in Urban Transit Systems. *IEEE Transactions on Intelligent Transportation Systems* 22, 2535–2552. doi:10.1109/TITS.2020.2973365. Conference Name: IEEE Transactions on Intelligent Transportation Systems
- Lu, X. and Pas, E. I. (1999). Socio-demographics, activity participation and travel behavior. *Transportation Research Part A: Policy and Practice* 33, 1–18. doi:10.1016/S0965-8564(98)00020-2
- McCool, D., Lugtig, P., and Schouten, B. (2022). Maximum interpolable gap length in missing smartphone-based gps mobility data. *Transportation* , 1–31
- McGuckin, N. and Murakami, E. (1999a). Examining trip-chaining behavior: Comparison of travel by men and women. *Transportation Research Record* 1693, 79–85. doi:10.3141/1693-12
- McGuckin, N. and Murakami, E. (1999b). Examining Trip-Chaining Behavior: Comparison of Travel by Men and Women. *Transportation Research Record: Journal of the Transportation Research Board* 1693, 79–85. doi:10.3141/1693-12
- Merikhipour, M., Khanmohammadidoustani, S., and Abbasi, M. (2024). Transportation mode detection through spatial attention-based transductive long short-term memory and off-policy feature selection. *Expert Systems With Applications* 267. doi:10.1016/j.eswa.2024.126196. [Online; accessed 2024-12-25]
- Merrill, N. H., Atkinson, S. F., Mulvaney, K. K., Mazzotta, M. J., and Bousquin, J. (2020). Using data derived from cellular phone locations to estimate visitation to natural areas: An application to water recreation in new england, usa. *PloS one* 15, e0231863
- Miah, S. J., Vu, H. Q., Gammack, J., and McGrath, M. (2017). A Big Data Analytics Method for Tourist Behaviour Analysis. *Information & Management* 54, 771–785. doi: 10.1016/j.im.2016.11.011

- Mohamed, A., Zhu, D., Vu, W., Elhoseiny, M., and Claudel, C. (2022). Social-implicit: Rethinking trajectory prediction evaluation and the effectiveness of implicit maximum likelihood estimation. In *European Conference on Computer Vision* (Springer), 463–479
- Mokhtarian, P. L. and Chen, C. (2004). Ttb or not TTB, that is the question: a review and analysis of the empirical literature on travel time (and money) budgets. *Transportation Research Part A: Policy and Practice* 38, 643–675. doi:10.1016/j.tra.2003.12.004
- Morris, D. H., Rossine, F. W., Plotkin, J. B., and Levin, S. A. (2021). Optimal, near-optimal, and robust epidemic control. *Communications Physics* 4, 78
- Müller, M. (2007). Dynamic time warping. *Information retrieval for music and motion* , 69–84
- Naeini, M. P., Cooper, G., and Hauskrecht, M. (2015). Obtaining Well Calibrated Probabilities Using Bayesian Binning. *Proceedings of the AAAI Conference on Artificial Intelligence* 29. doi:10.1609/aaai.v29i1.9602. Number: 1
- Nasernejad, P., Sayed, T., and Alsaleh, R. (2021). Modeling pedestrian behavior in pedestrian-vehicle near misses: A continuous gaussian process inverse reinforcement learning (gp-irl) approach. *Accident Analysis & Prevention* 161, 106355
- Nevin, J. W., Vaquero-Caballero, F., Ives, D. J., and Savory, S. J. (2021). Physics-informed gaussian process regression for optical fiber communication systems. *Journal of Lightwave Technology* 39, 6833–6844
- Nishii, K., Kondo, K., and Kitamura, R. (1988). An empirical analysis of trip chaining behavior
- Ong, C., Williamson, R. C., and Smola, A. (2002). Hyperkernels. *Advances in neural information processing systems* 15
- Opsahl, T. and Panzarasa, P. (2009). Clustering in weighted networks. *Social Networks* 31, 155–163. doi:10.1016/j.socnet.2009.02.002

- Pappalardo, L. and Simini, F. (2018). Data-driven generation of spatio-temporal routines in human mobility. *Data Mining and Knowledge Discovery* 32, 787–829. doi:10.1007/s10618-017-0548-4
- Pappalardo, L., Simini, F., Barlacchi, G., and Pellungrini, R. (2019). scikit-mobility: A python library for the analysis, generation and risk assessment of mobility data. *arXiv preprint arXiv:1907.07062*
- Park, H., Oh, C., Moon, J., and Kim, S. (2018). Development of a lane change risk index using vehicle trajectory data. *Accident Analysis & Prevention* 110, 1–8. doi:10.1016/j.aap.2017.10.015
- Pecheux, K. K., Pecheux, B. B., and Ledbetter, G. (2020). *Framework for Managing Data from Emerging Transportation Technologies to Support Decision-Making* (Washington, D.C.: Transportation Research Board). doi:10.17226/25965
- Pendyala, R. M., Kitamura, R., Chen, C., and Pas, E. I. (1997). An activity-based microsimulation analysis of transportation control measures. *Transport Policy* 4, 183–192. doi:https://doi.org/10.1016/S0967-070X(97)00005-X
- Pereira, F. C., Rodrigues, F., and Ben-Akiva, M. (2015). Using Data From the Web to Predict Public Transport Arrivals Under Special Events Scenarios. *Journal of Intelligent Transportation Systems* 19, 273–288. doi:10.1080/15472450.2013.868284. Publisher: Taylor & Francis _eprint: https://doi.org/10.1080/15472450.2013.868284
- Phithakkitnukoon, S., Sukhvibul, T., Demissie, M., Smoreda, Z., Natwichai, J., and Bento, C. (2017). Inferring social influence in transport mode choice using mobile phone data. *EPJ Data Science* 6, 11. doi:10.1140/epjds/s13688-017-0108-6
- Qi, B., Kang, L., and Banerjee, S. (2017). A vehicle-based edge computing platform for transit and human mobility analytics. In *Proceedings of the Second ACM/IEEE Symposium on Edge Computing* (New York, NY, USA: Association for Computing Machinery), SEC '17, 1–14. doi:10.1145/3132211.3134446

- Raissi, M., Perdikaris, P., and Karniadakis, G. E. (2019). Physics-informed neural networks: A deep learning framework for solving forward and inverse problems involving nonlinear partial differential equations. *Journal of Computational physics* 378, 686–707
- Rasmussen, C. and Ghahramani, Z. (2000). Occam’ s Razor. In *Advances in Neural Information Processing Systems* (MIT Press), vol. 13
- Rasmussen, C. E. and Williams, C. K. I. (2006). *Gaussian processes for machine learning*. Adaptive computation and machine learning (Cambridge, Mass: MIT Press). OCLC: ocm61285753
- Razavi, R., Xue, G., and Akpan, I. J. (2024). Predicting Sociodemographic Attributes from Mobile Usage Patterns: Applications and Privacy Implications. *Big Data* 12, 213–228. doi:10.1089/big.2022.0182. Publisher: Mary Ann Liebert, Inc., publishers
- Reif, J. and Schmücker, D. (2020). Exploring new ways of visitor tracking using big data sources: Opportunities and limits of passive mobile data for tourism. *Journal of Destination Marketing & Management* 18, 100481. doi:10.1016/j.jdmm.2020.100481
- Ren, H., Ruan, S., Li, Y., Bao, J., Meng, C., Li, R., et al. (2021). Mtrajrec: Map-constrained trajectory recovery via seq2seq multi-task learning. In *Proceedings of the 27th ACM SIGKDD Conference on Knowledge Discovery & Data Mining*. 1410–1419
- Richardson, A. J., Ampt, E. S., and Meyburg, A. H. (1996). Nonresponse issues in household travel surveys. In *Transportation Research Board* (National Research Council), vol. 10, 79–114
- Richardson, H. (2021). Finding a Win-Win: Planning and Data-Sharing Partnerships between Governments and Public Land Management Agencies
- Rittel, H. W. and Webber, M. M. (1973). Dilemmas in a general theory of planning. *Policy sciences* 4, 155–169
- Rodrigues, F., Henrickson, K., and Pereira, F. C. (2018). Multi-output gaussian processes for crowdsourced traffic data imputation. *IEEE Transactions on Intelligent Transportation Systems* 20, 594–603

- Sakhare, R. S., Desai, J., Saldivar-Carranza, E. D., and Bullock, D. M. (2024). Methodology for Monitoring Border Crossing Delays with Connected Vehicle Data: United States and Mexico Land Crossings Case Study. *Future Transportation* 4, 107–129. doi:10.3390/futuretransp4010007. Number: 1 Publisher: Multidisciplinary Digital Publishing Institute
- Sarker, M. N. I., Peng, Y., Yiran, C., and Shouse, R. C. (2020). Disaster resilience through big data: Way to environmental sustainability. *International Journal of Disaster Risk Reduction* 51, 101769. doi:10.1016/j.ijdr.2020.101769
- Schneider, C. M., Belik, V., Couronné, T., Smoreda, Z., and González, M. C. (2013). Unravelling daily human mobility motifs. *Journal of The Royal Society Interface* 10, 20130246. doi:10.1098/rsif.2013.0246. Publisher: Royal Society
- Schölkopf, B. and Smola, A. J. (2002). *Learning with kernels: support vector machines, regularization, optimization, and beyond* (MIT press)
- Schwarz, G. (1978). Estimating the dimension of a model. *The annals of statistics* , 461–464
- Shannon, C. E. (1948). A mathematical theory of communication. *The Bell System Technical Journal* 27, 379–423. doi:10.1002/j.1538-7305.1948.tb01338.x
- Sheller, M. and Urry, J. (2006). The New Mobilities Paradigm. *Environment and Planning A: Economy and Space* 38, 207–226. doi:10.1068/a37268. Publisher: SAGE Publications Ltd
- Shin, Y., Darbon, J., and Karniadakis, G. E. (2020). On the convergence of physics informed neural networks for linear second-order elliptic and parabolic type pdes. *arXiv preprint arXiv:2004.01806*
- Singh, G., Sivaraman, V., and Hard, E. (2022). *State of Emerging Mobility Big Data Sources and its Applications — Task 1: Evaluate Mobility Datasets*. Tech. rep., Texas A&M Transportation Institute
- Snelson, E. and Ghahramani, Z. (2005). Sparse gaussian processes using pseudo-inputs. *Advances in neural information processing systems* 18

- Solomon, A., Bar, A., Yanai, C., Shapira, B., and Rokach, L. (2018). Predict Demographic Information Using Word2vec on Spatial Trajectories. In *Proceedings of the 26th Conference on User Modeling, Adaptation and Personalization* (New York, NY, USA: Association for Computing Machinery), UMAP '18, 331–339. doi:10.1145/3209219.3209224. [Online; accessed 2025-06-20]
- Song, C., Koren, T., Wang, P., and Barabási, A.-L. (2010a). Modelling the scaling properties of human mobility. *Nature physics* 6, 818–823
- Song, C., Qu, Z., Blumm, N., and Barabási, A.-L. (2010b). Limits of predictability in human mobility. *Science* 327, 1018–1021
- [Dataset] Spectus (2022). Sensitive Points of Interest Policy
- Stentoft, A., Lee, B.-S., and Schläpfer, M. (2024). Quantifying the uncertainty of mobility flow predictions using gaussian processes. *Transportation* 51, 2301–2322
- Stopher, P. R., Kockelman, K., Greaves, S. P., and Clifford, E. (2008). Reducing Burden and Sample Sizes in Multiday Household Travel Surveys. *Transportation Research Record* 2064, 12–18. doi:10.3141/2064-03. Publisher: SAGE Publications Inc
- Storm, P. J., Mandjes, M., and van Arem, B. (2022). Efficient evaluation of stochastic traffic flow models using gaussian process approximation. *Transportation research part B: methodological* 164, 126–144
- [Dataset] StreetLight (2025). Tulsa MPO Fills Traffic Count Gaps - AADT.pdf
- Sulis, P., Manley, E., Zhong, C., and Batty, M. (2018). Using mobility data as proxy for measuring urban vitality. *Journal of Spatial Information Science* 2018, 137–162. doi: 10.5311/JOSIS.2018.16.384
- Sun, H., Yang, C., Deng, L., Zhou, F., Huang, F., and Zheng, K. (2021). Periodicmove: Shift-aware human mobility recovery with graph neural network. In *Proceedings of the 30th ACM International Conference on Information & Knowledge Management*. 1734–1743

- Teixeira, D. d. C., Almeida, J. M., and Viana, A. C. (2021). On estimating the predictability of human mobility: the role of routine. *EPJ Data Science* 10, 49
- Timmermans, H. and Arentze, T. A. (2011). Transport models and urban planning practice: Experiences with albatross. *Transport Reviews* 31, 199–207. doi:10.1080/01441647.2010.518292
- Titsias, M. (2009). Variational learning of inducing variables in sparse gaussian processes. In *Artificial intelligence and statistics* (PMLR), 567–574
- Tobler, W. R. (1970). A computer movie simulating urban growth in the detroit region. *Economic geography* 46, 234–240
- Tozluoğlu, Ç., Liao, Y., and Sprei, F. (2025). Mobile phone application data for activity plan generation. *Transportation* , 1–33
- Ugurel, E. (2023). Time-varying transition matrices with multi-task gaussian processes. *arXiv preprint arXiv:2306.11772*
- Ugurel, E., Guan, X., Wang, Y., Huang, S., Wang, Q. R., and Chen, C. (2024a). Correcting missingness in passively-generated mobile data with multi-task gaussian processes. *Transportation Research Part C* 161
- Ugurel, E., Wu, X., Wang, R., Lee, B. H. Y., and Chen, C. (2024b). Metropolitan Planning Organizations’ Uses of and Needs for Big Data. *Findings* doi:10.32866/001c.127143. Publisher: Findings Press
- Uğurel, E., Huang, S., and Chen, C. (2024). Learning to generate synthetic human mobility data: A physics-regularized Gaussian process approach based on multiple kernel learning. *Transportation Research Part B: Methodological* 189, 103064. doi:10.1016/j.trb.2024.103064
- Van Rooyen, B. and Williamson, R. C. (2018). A theory of learning with corrupted labels. *Journal of Machine Learning Research* 18, 1–50
- van Rossum, G. (1995). *Python tutorial*. Tech. Rep. CS-R9526, Centrum voor Wiskunde en Informatica (CWI), Amsterdam

- Venkatesh, Morris, Davis, and Davis (2003). User Acceptance of Information Technology: Toward a Unified View. *MIS Quarterly* 27, 425. doi:10.2307/30036540. Publisher: JSTOR
- Wallace, B., Barnes, J., and Rutherford, G. S. (2000). Evaluating the effects of traveler and trip characteristics on trip chaining, with implications for transportation demand management strategies. *Transportation Research Record* 1718, 97–106. doi:10.3141/1718-13
- Wang, F. and Chen, C. (2018). On data processing required to derive mobility patterns from passively-generated mobile phone data. *Transportation Research Part C: Emerging Technologies* 87, 58–74
- Wang, F., Wang, J., Zhang, Y., Chen, C., and Ban, X. J. (2021a). *Travelers' Adaptive Behaviors in Response to Seattle's Alaskan Way Viaduct Replacement*. Tech. rep.
- Wang, F., Wang, X., and Ban, X. J. (2024). Data poisoning attacks in intelligent transportation systems: A survey. *Transportation Research Part C: Emerging Technologies* 165, 104750. doi:10.1016/j.trc.2024.104750
- Wang, S., Wang, H., and Perdikaris, P. (2021b). On the eigenvector bias of fourier feature networks: From regression to solving multi-scale pdes with physics-informed neural networks. *Computer Methods in Applied Mechanics and Engineering* 384, 113938
- Wang, S., Yu, X., and Perdikaris, P. (2022a). When and why pinns fail to train: A neural tangent kernel perspective. *Journal of Computational Physics* 449, 110768
- Wang, Y., Guan, X., Ugurel, E., Chen, C., Huang, S., and Wang, Q. R. (2025). Exploring biases in travel behavior patterns in big passively generated mobile data from 11 U.S. cities. *Journal of Transport Geography* 123, 104108. doi:10.1016/j.jtrangeo.2024.104108
- Wang, Z., Xing, W., Kirby, R., and Zhe, S. (2022b). Physics informed deep kernel learning. In *International Conference on Artificial Intelligence and Statistics* (PMLR), 1206–1218
- Wang, Z., Zhang, S., and James, J. (2020). Reconstruction of missing trajectory data: a deep learning approach. In *2020 IEEE 23rd International Conference on Intelligent Transportation Systems (ITSC)* (IEEE), 1–6

- Wes McKinney (2010). Data Structures for Statistical Computing in Python. In *Proceedings of the 9th Python in Science Conference*, eds. Stéfan van der Walt and Jarrod Millman. 56 – 61. doi:10.25080/Majora-92bf1922-00a
- Wesolowski, A., Eagle, N., Noor, A. M., Snow, R. W., and Buckee, C. O. (2013). The impact of biases in mobile phone ownership on estimates of human mobility. *Journal of The Royal Society Interface* 10, 20120986. doi:10.1098/rsif.2012.0986
- Williams, N. E., Thomas, T. A., Dunbar, M., Eagle, N., and Dobra, A. (2015). Measures of human mobility using mobile phone records enhanced with gis data. *PloS one* 10, e0133630
- Wilson, A. and Adams, R. (2013). Gaussian process kernels for pattern discovery and extrapolation. In *International conference on machine learning* (PMLR), 1067–1075
- Wilson, A. G. (2014). *Covariance kernels for fast automatic pattern discovery and extrapolation with Gaussian processes*. Ph.D. thesis, University of Cambridge Cambridge, UK
- Wu, F., Cheng, Z., Chen, H., Qiu, Z., and Sun, L. (2024a). Traffic state estimation from vehicle trajectories with anisotropic gaussian processes. *Transportation Research Part C: Emerging Technologies* 163, 104646
- Wu, L., Yang, L., Huang, Z., Wang, Y., Chai, Y., Peng, X., et al. (2019). Inferring demographics from human trajectories and geographical context. *Computers, Environment and Urban Systems* 77, 101368. doi:10.1016/j.compenvurbsys.2019.101368
- Wu, X., He, H., Wang, Y., and Wang, Q. (2024b). Pretrained mobility transformer: A foundation model for human mobility. *arXiv preprint arXiv:2406.02578*
- Wu, X., Wang, Y., Ugurel, E., Chen, C., Huang, S., and Wang, Q. R. (2024c). Location-based service (lbs) data quality metrics and effects on mobility inference. *arXiv preprint arXiv:2411.16595*
- Xie, K., Ozbay, K., Kurkcu, A., and Yang, H. (2017). Analysis of Traffic Crashes Involving Pedestrians Using Big Data: Investigation of Contributing Factors and Identification of Hotspots. *Risk Analysis* 37, 1459–1476. doi:10.1111/risa.12785. eprint: <https://onlinelibrary.wiley.com/doi/pdf/10.1111/risa.12785>

- Xie, Y., Zhao, K., Sun, Y., and Chen, D. (2010). Gaussian processes for short-term traffic volume forecasting. *Transportation Research Record* 2165, 69–78
- Xu, F., Lin, Z., Xia, T., Guo, D., and Li, Y. (2020). Sume: Semantic-enhanced Urban Mobility Network Embedding for User Demographic Inference. *Proc. ACM Interact. Mob. Wearable Ubiquitous Technol.* 4, 98:1–98:25. doi:10.1145/3411807
- Yang, M., Pan, Y., Darzi, A., Ghader, S., Xiong, C., and Zhang, L. (2022). A data-driven travel mode share estimation framework based on mobile device location data. *Transportation* 49, 1339–1383. doi:10.1007/s11116-021-10214-3
- Yang, T., Xu, X., Guo, Q., Zhang, L., and Sun, H. (2017a). Ev charging behaviour analysis and modelling based on mobile crowdsensing data. *IET Generation, Transmission & Distribution* 11, 1683–1691. doi:10.1049/iet-gtd.2016.1200. eprint: <https://onlinelibrary.wiley.com/doi/pdf/10.1049/iet-gtd.2016.1200>
- Yang, T., Xu, X., Guo, Q., Zhang, L., and Sun, H. (2017b). EV charging behaviour analysis and modelling based on mobile crowdsensing data. *IET Generation, Transmission & Distribution* 11, 1683–1691. doi:10.1049/iet-gtd.2016.1200. eprint: <https://onlinelibrary.wiley.com/doi/pdf/10.1049/iet-gtd.2016.1200>
- Yang, X., Barajas-Solano, D., Tartakovsky, G., and Tartakovsky, A. M. (2019). Physics-informed cokriging: A gaussian-process-regression-based multifidelity method for data-model convergence. *Journal of Computational Physics* 395, 410–431
- Yin, J., Dong, J., Hamm, N. A. S., Li, Z., Wang, J., Xing, H., et al. (2021). Integrating remote sensing and geospatial big data for urban land use mapping: A review. *International Journal of Applied Earth Observation and Geoinformation* 103, 102514. doi:10.1016/j.jag.2021.102514
- Yu, Q., Zhang, H., Li, W., Sui, Y., Song, X., Yang, D., et al. (2020). Mobile phone data in urban bicycle-sharing: Market-oriented sub-area division and spatial analysis on emission reduction potentials. *Journal of Cleaner Production* 254, 119974. doi:10.1016/j.jclepro.2020.119974

- Yuan, Y. and Raubal, M. (2012). Extracting dynamic urban mobility patterns from mobile phone data. In *Geographic Information Science: 7th International Conference, GIScience 2012, Columbus, OH, USA, September 18-21, 2012. Proceedings 7* (Springer), 354–367
- Yuan, Y., Zhang, Z., Yang, X. T., and Zhe, S. (2021). Macroscopic traffic flow modeling with physics regularized gaussian process: A new insight into machine learning applications in transportation. *Transportation Research Part B: Methodological* 146, 88–110
- Zeidan, A., Lagerspetz, E., Zhao, K., Nurmi, P., Tarkoma, S., and Vo, H. T. (2020). GeoMatch: Efficient Large-scale Map Matching on Apache Spark. *ACM/IMS Transactions on Data Science* 1, 1–30. doi:10.1145/3402904
- Zhang, B., Rasouli, S., and Feng, T. (2024a). Social demographics imputation based on similarity in multi-dimensional activity-travel pattern: A two-step approach. *Travel Behaviour and Society* 37, 100843. doi:10.1016/j.tbs.2024.100843
- Zhang, K., Pang, Y., Zhang, Y., and Sekimoto, Y. (2024b). MobGLM: A Large Language Model for Synthetic Human Mobility Generation. In *Proceedings of the 32nd ACM International Conference on Advances in Geographic Information Systems* (New York, NY, USA: Association for Computing Machinery), SIGSPATIAL '24, 629–632. doi:10.1145/3678717.3691311
- Zhang, P., Zhou, J., and Zhang, T. (2017). Quantifying and visualizing jobs-housing balance with big data: A case study of Shanghai. *Cities* 66, 10–22. doi:10.1016/j.cities.2017.03.004
- Zhao, Y., Pawlak, J., and Sivakumar, A. (2022). Theory for socio-demographic enrichment performance using the inverse discrete choice modelling approach. *Transportation Research Part B: Methodological* 155, 101–134. doi:10.1016/j.trb.2021.11.004
- Zheng, Y., Li, Q., Chen, Y., Xie, X., and Ma, W.-Y. (2008). Understanding mobility based on gps data. In *Proceedings of the 10th international conference on Ubiquitous computing*. 312–321
- Zheng, Y., Xie, X., Ma, W.-Y., et al. (2010). Geolife: A collaborative social networking service among user, location and trajectory. *IEEE Data Eng. Bull.* 33, 32–39

- Zheng, Y., Zhang, L., Xie, X., and Ma, W.-Y. (2009). Mining interesting locations and travel sequences from gps trajectories. In *Proceedings of the 18th international conference on World wide web*. 791–800
- Zhong, Y., Yuan, N. J., Zhong, W., Zhang, F., and Xie, X. (2015). You Are Where You Go: Inferring Demographic Attributes from Location Check-ins. In *Proceedings of the Eighth ACM International Conference on Web Search and Data Mining* (New York, NY, USA: Association for Computing Machinery), WSDM '15, 295–304. doi:10.1145/2684822.2685287. [Online; accessed 2025-06-20]
- Zhou, H., Wang, H., Zhou, Y., Luo, X., Tang, Y., Xue, L., et al. (2020). Demystifying diehard android apps. In *Proceedings of the 35th IEEE/ACM International Conference on Automated Software Engineering*. 187–198
- Zhou, X., Yeh, A. G. O., and Yue, Y. (2018). Spatial variation of self-containment and jobs-housing balance in Shenzhen using cellphone big data. *Journal of Transport Geography* 68, 102–108. doi:10.1016/j.jtrangeo.2017.12.006
- Zhu, Y., Ye, Y., Wu, Y., Zhao, X., and Yu, J. (2023). SynMob: Creating High-Fidelity Synthetic GPS Trajectory Dataset for Urban Mobility Analysis. *Advances in Neural Information Processing Systems* 36, 22961–22977
- Zhu, Z., Xu, M., Di, Y., and Yang, H. (2022). Fitting spatial-temporal data via a physics regularized multi-output grid gaussian process: case studies of a bike-sharing system. *IEEE Transactions on Intelligent Transportation Systems* 23, 21090–21101
- Álvarez, M., Luengo, D., and Lawrence, N. D. (2009). Latent force models. In *Proceedings of the Twelfth International Conference on Artificial Intelligence and Statistics*, eds. D. van Dyk and M. Welling (Hilton Clearwater Beach Resort, Clearwater Beach, Florida USA: PMLR), vol. 5 of *Proceedings of Machine Learning Research*, 9–16

Appendices

Appendix 1

SUPPLEMENTARY MATERIAL FOR CHAPTER 2

A.1 Initialization Strategy & Parameter Optimization

The optimization problem of a GP is non-convex. Therefore, kernel parameter initialization can help avoid bad local optima, leading to better model estimation and more accurate prediction results. We achieve this in two ways. First, if the gap is short enough, we take advantage of training points near the gap to infer trip characteristics. These include metrics like average velocity, total distance traveled, trip duration, and stop rate, among others. Our first set of experiments suggest that trips with differing characteristics tend to land at different optimal parameters. Therefore, we suggest initializing the lengthscale parameter of continuous temporal dimensions (Unix time in our case) according to the characteristics of the trip with missing data. If using a monotonically-increasing temporal dimension as an input, one alternative we have identified is to initialize the lengthscale parameter with the average length of gap between observations in the training data. Ak et al. (2018) employed this strategy without any optimization. We however do further optimize the parameter beyond this initial heuristic using the Adaptive Moment Estimation algorithm (Kingma and Ba, 2014), a preferred choice among various machine learning frameworks due to its computational efficiency and ability to manage sparse gradients on noisy data. We initialize binary variables (i.e., one-of-k encoded categorical variables) with a lengthscale of 1. This is largely done to avoid model misspecification during the optimization stage as lengthscale parameters are constrained to be nonnegative. In our implementation, a small amount of noise (called jitter) is added to the diagonal of the covariance matrix. This is done to ensure numerical stability when performing matrix operations like inversion and decomposition. Simultaneously, however, jitter can lead to parameter values dipping below 0 at the very first iteration of the optimization algorithm¹.

¹If the learning rate is small, the first iteration of the optimization algorithm may result in a value that is very close to 0, and the associated jitter may cause the parameter to flip signs.

A.2 Algorithm to Determine Training Data for Experiments in 2.6.2

To choose the training dataset in 2.6.2, we used Algorithm 2, which simulates gaps in trajectory data based on the current temporal occupancy with respect to the bin length being tested. The following is a step-by-step breakdown of that algorithm:

1. **Binning the Trajectory Data:** The algorithm starts by creating an array of integers called bins that span from $\mathbf{t}(u, 1)$ to $\mathbf{t}(u, N)$ with a specified step size τ . This step is effectively dividing the entire temporal range of the trajectory data into bins of equal length.
2. **Mapping Data Points to Bins:** The *bins_dict* is then created using dictionary comprehension to map each data point to the relevant bin. This means that each data point in the trajectory is associated with a specific bin based on its timestamp.
3. **Selecting Non-Empty Bins:** Another dictionary, *non_empty_bins_dict*, is created to store only the non-empty bins. This step involves filtering out bins that have data points associated with them.
4. **Setting Target Occupancy:** The algorithm then sets a target occupancy level (*target_occup*) by randomly selecting a value between 0 and the current temporal occupancy (*q_curr*). The selected value is used to determine the desired temporal occupancy.
5. **Gap Simulation Loop:** The algorithm enters a loop to simulate gaps in the trajectory data until the temporal occupancy falls below the target level (*target_occup*). Within this loop:
 - It randomly chooses a bin from the non-empty bins using *np.random.choice*.
 - All data points associated with the chosen bin are removed from the original data, and the *bins_dict* and *non_empty_bins_dict* are updated accordingly.
6. **Calculating New Temporal Occupancy:** After the loop, the algorithm calculates the new temporal occupancy (*new_occup*) with the gapped data.
7. **Return Output:** Finally, the algorithm returns the gapped trajectory data, as well as the new temporal occupancy (q_τ).

A.3 Detailed Imputation Accuracy Metrics for Robustness Experiments

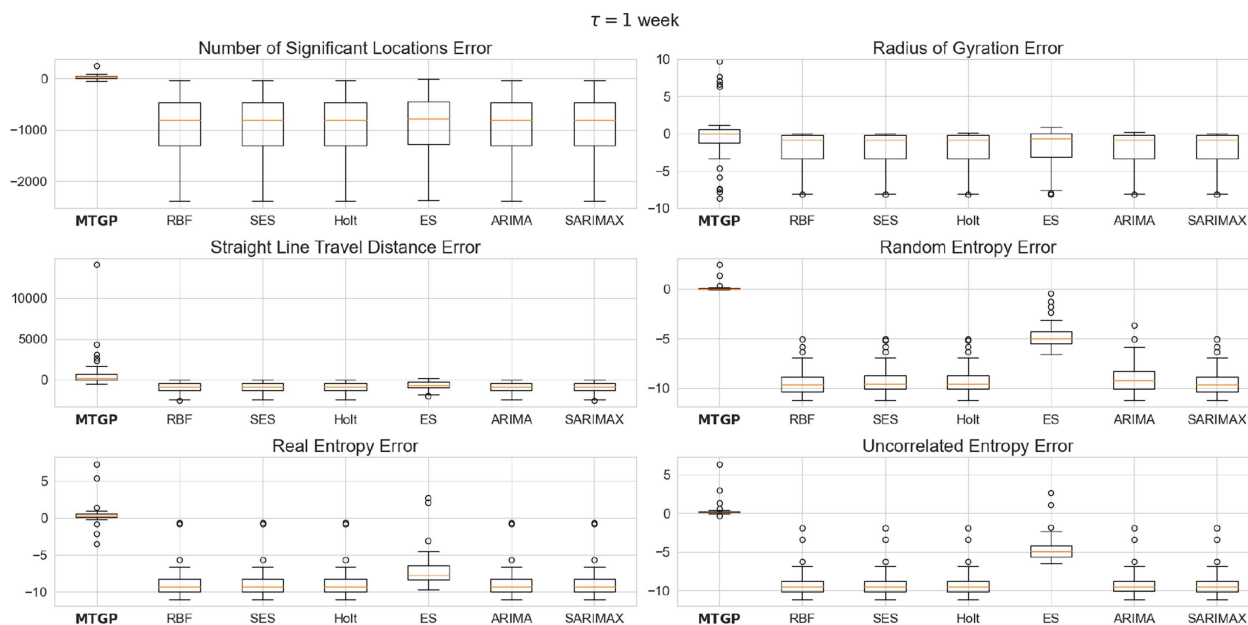


Figure A.1: Mobility error metrics for $\tau = 1$ week

A.4 Benchmark Methods and Model Parameters

A.5 Max Speed Threshold Sensitivity Analysis

To determine the optimal choice for a max speed threshold during the data preprocessing stage, we created a kernel density estimation (KDE) plot to visualize the number of observations eliminated by different speed limits. The KDE plot shows that the 200 km/h speed limit captures all of the points that the 300 km/h speed limit captures and some more. Because the speeds we calculate are based on “as the crow flies” distances, it is better to be more conservative in this limit—therefore, the additional points captured by the 200 km/h limit are most likely erroneous and have actual speeds much higher than 200 km/h (which is already high). This suggests that the 200 km/h speed limit is a better limit for eliminating outliers and reducing noise in the dataset. The 100 km/h speed limit, on the other hand, eliminates too many observations and reduces the data quality. Therefore, we chose the 200 km/h speed limit as the optimal threshold for our model (Figure A.15).

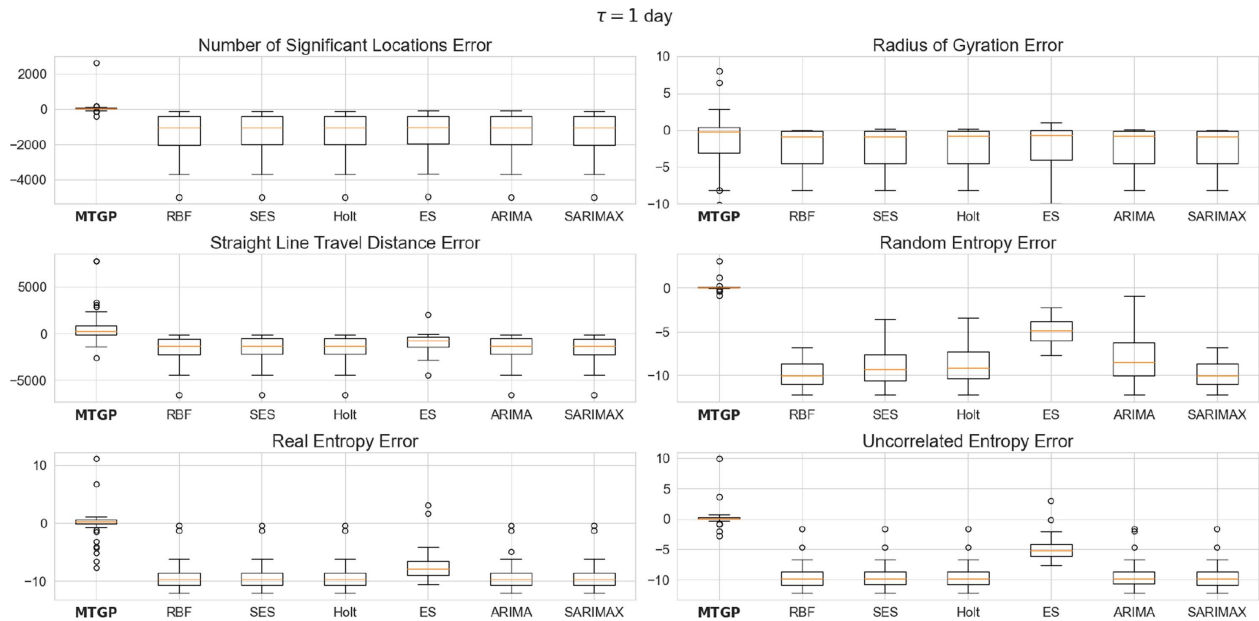


Figure A.2: Mobility error metrics for $\tau = 1$ day

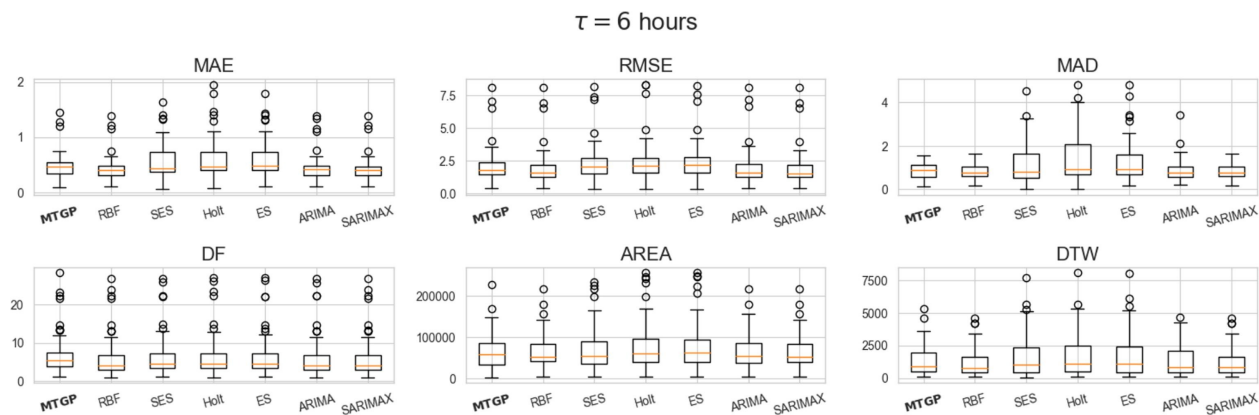
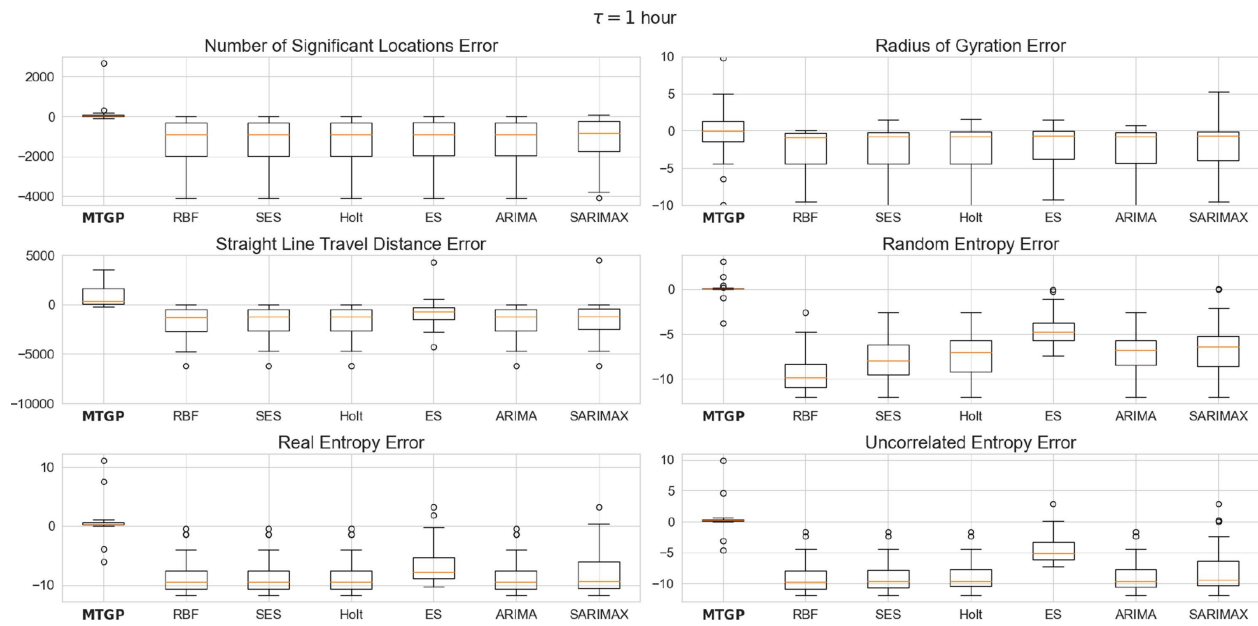
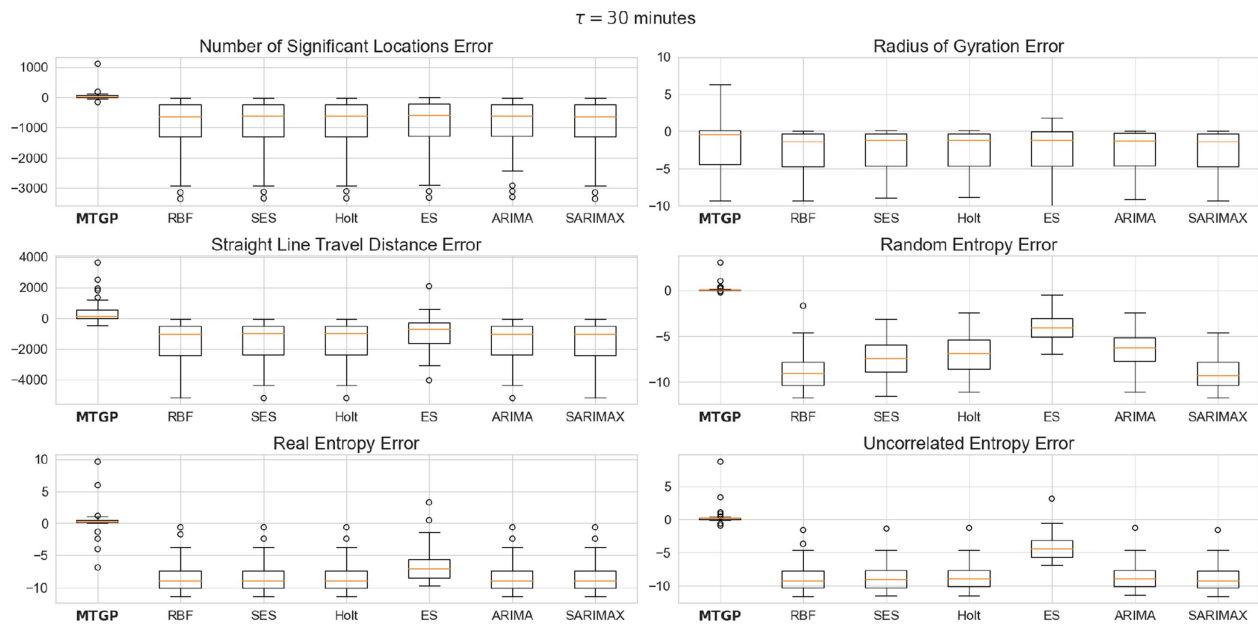
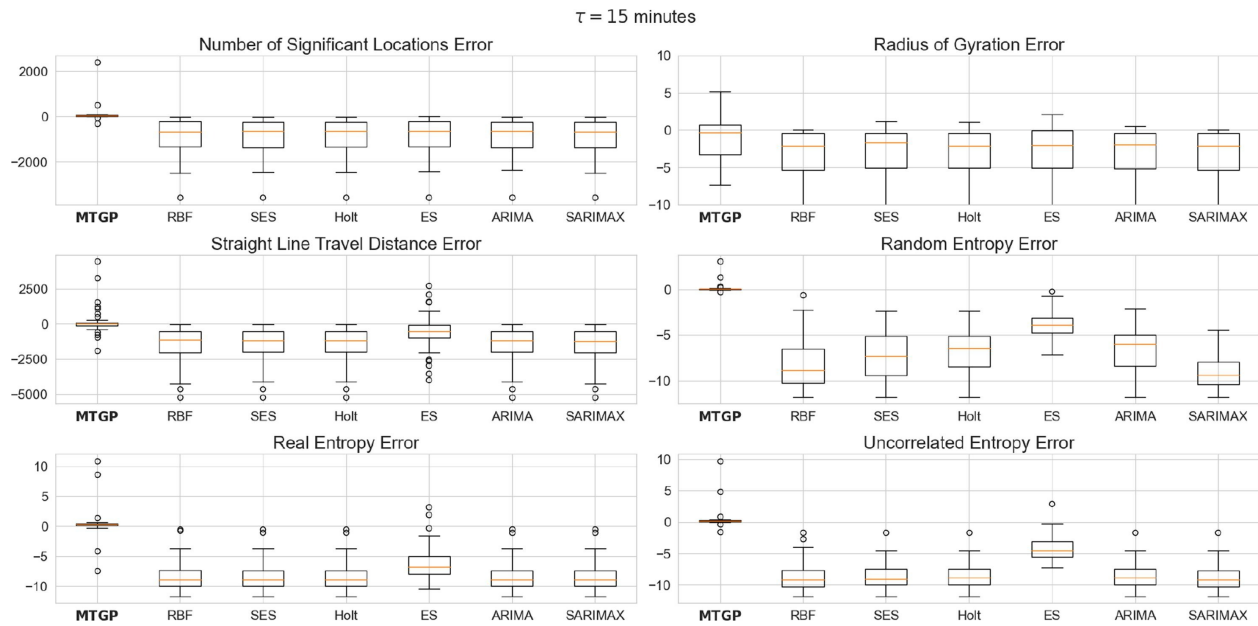
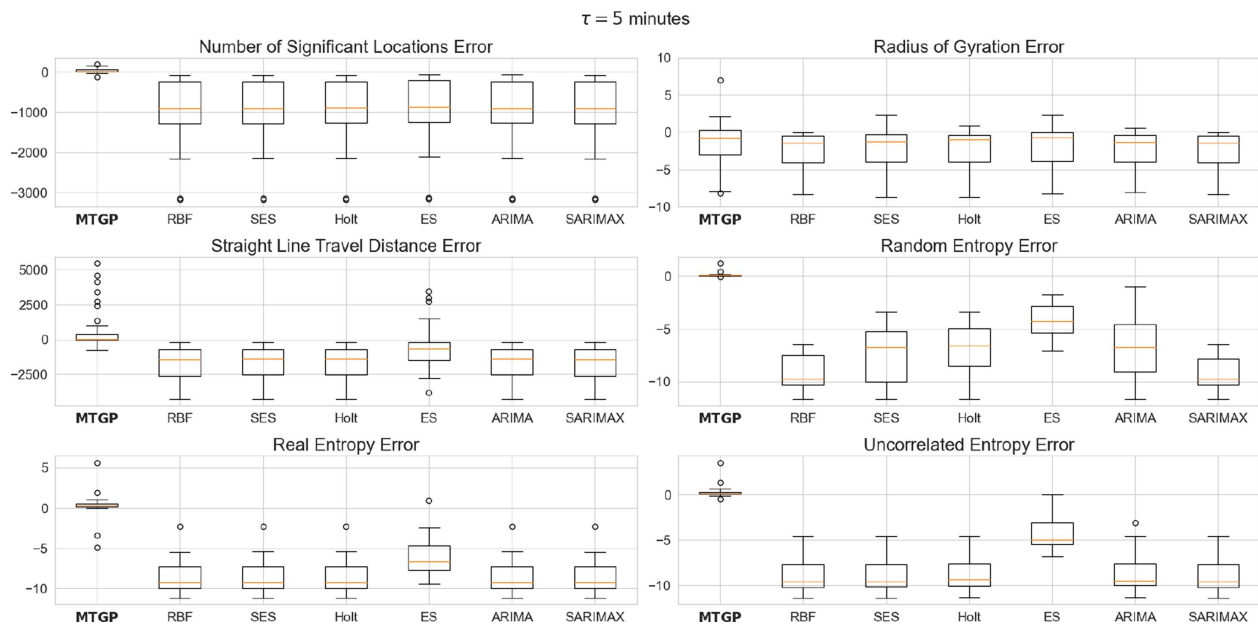
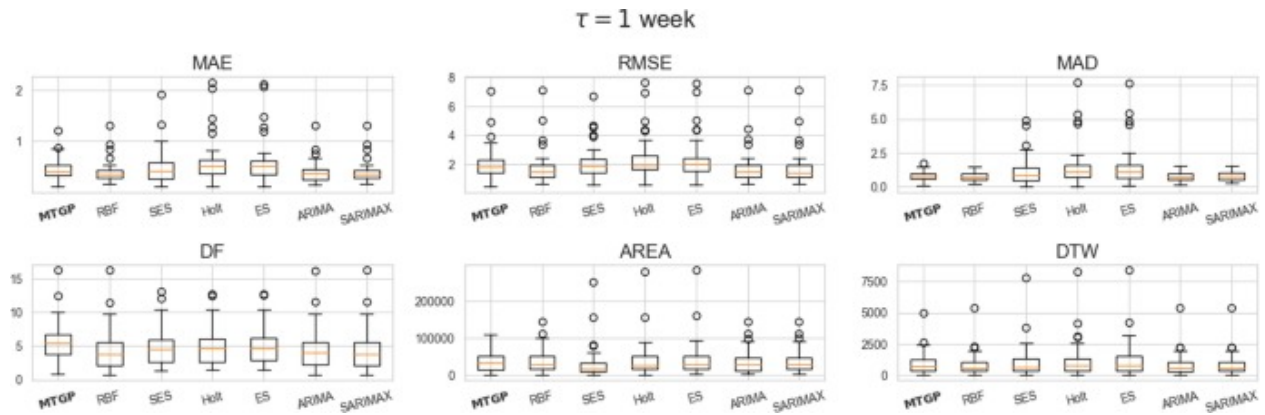
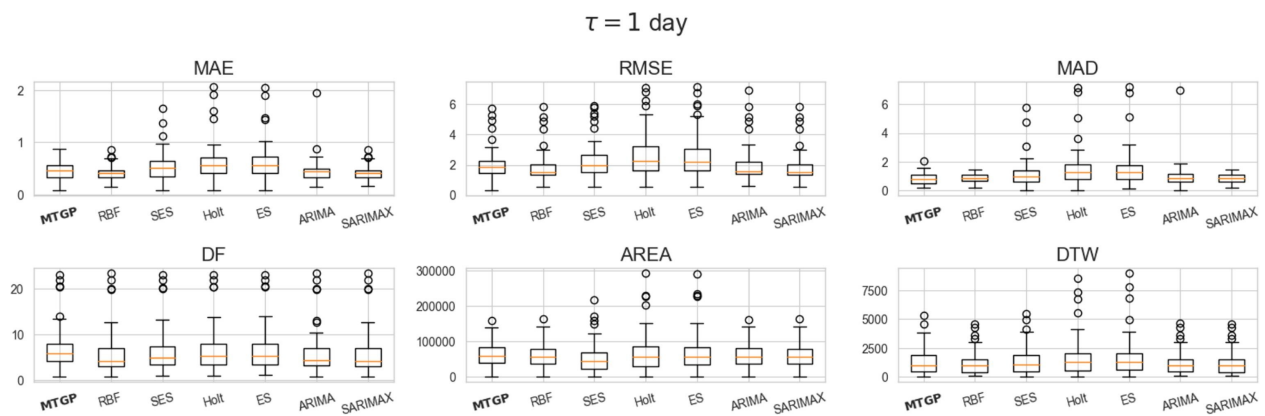
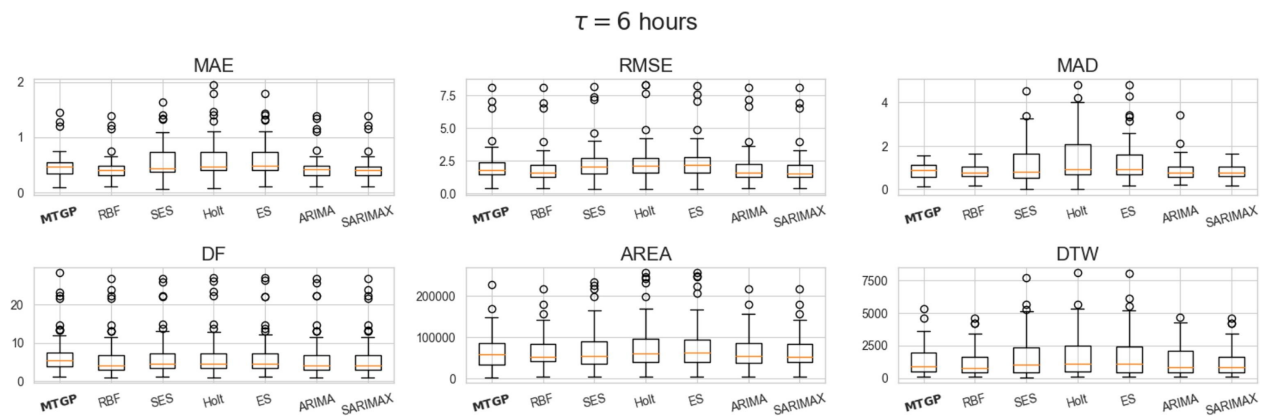


Figure A.3: Mobility error metrics for $\tau = 6$ hours

Figure A.4: Mobility error metrics for $\tau = 1$ hourFigure A.5: Mobility error metrics for $\tau = 30$ minutes

Figure A.6: Mobility error metrics for $\tau = 15$ minutesFigure A.7: Mobility error metrics for $\tau = 5$ minutes

Figure A.8: Classical error metrics for $\tau = 1 \text{ week}$ Figure A.9: Classical error metrics for $\tau = 1 \text{ day}$ Figure A.10: Classical error metrics for $\tau = 6 \text{ hours}$

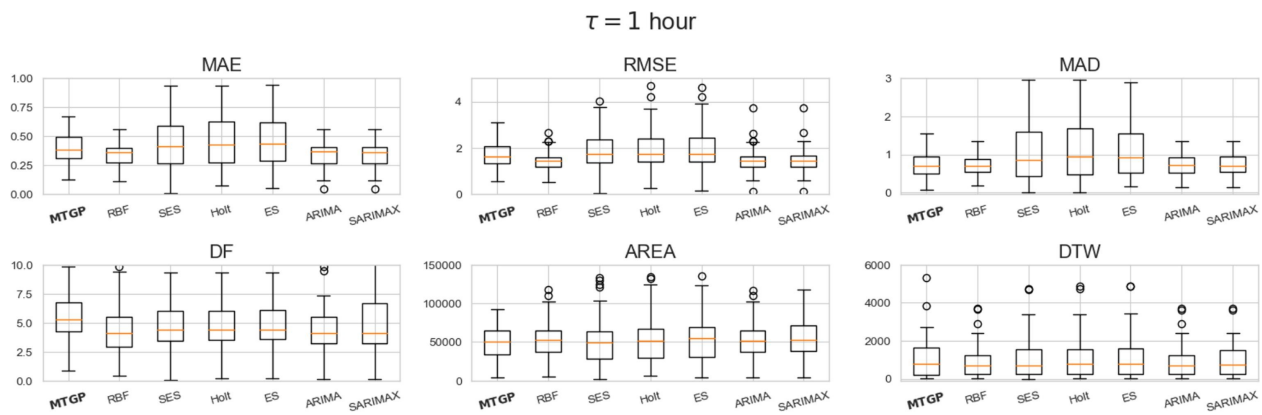


Figure A.11: Classical error metrics for $\tau = 1$ hour

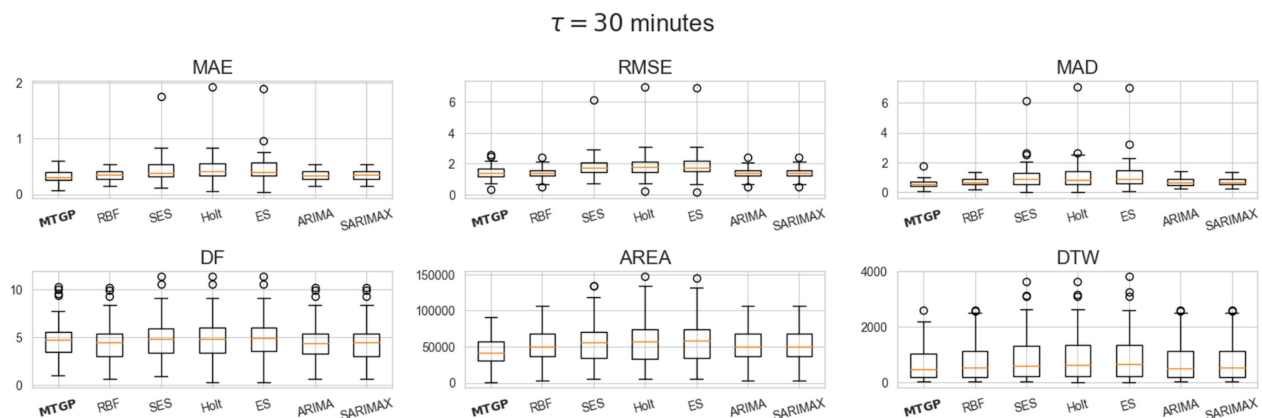


Figure A.12: Classical error metrics for $\tau = 30$ minutes

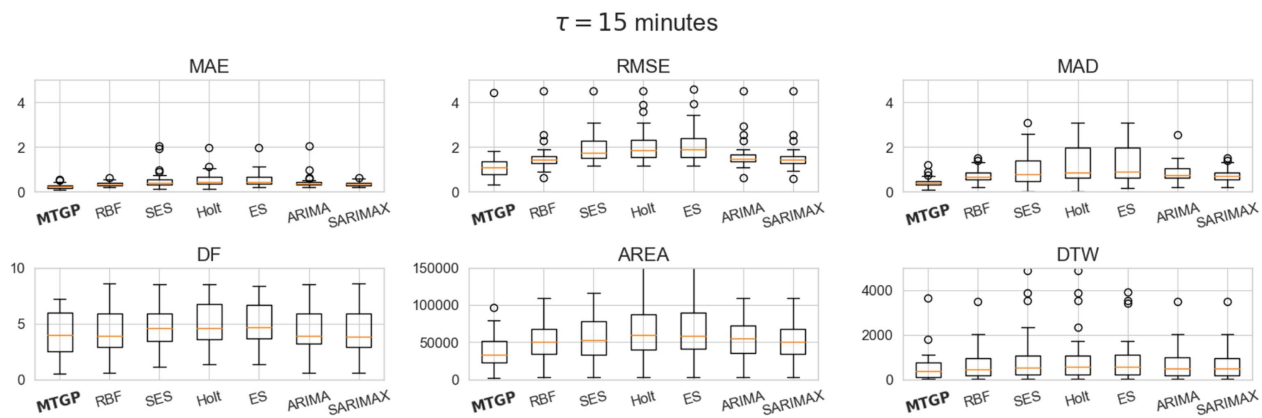


Figure A.13: Classical error metrics for $\tau = 15$ minutes

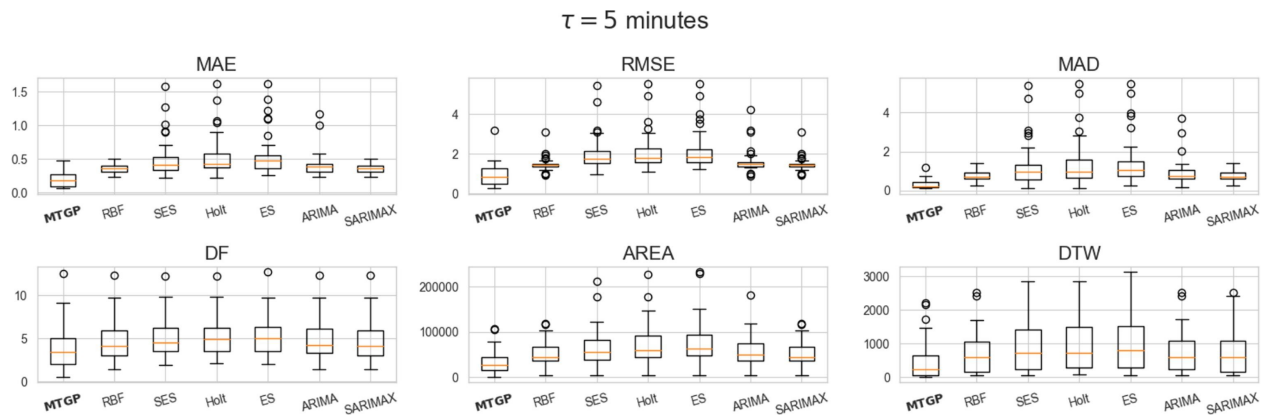


Figure A.14: Classical error metrics for $\tau = 5$ minutes

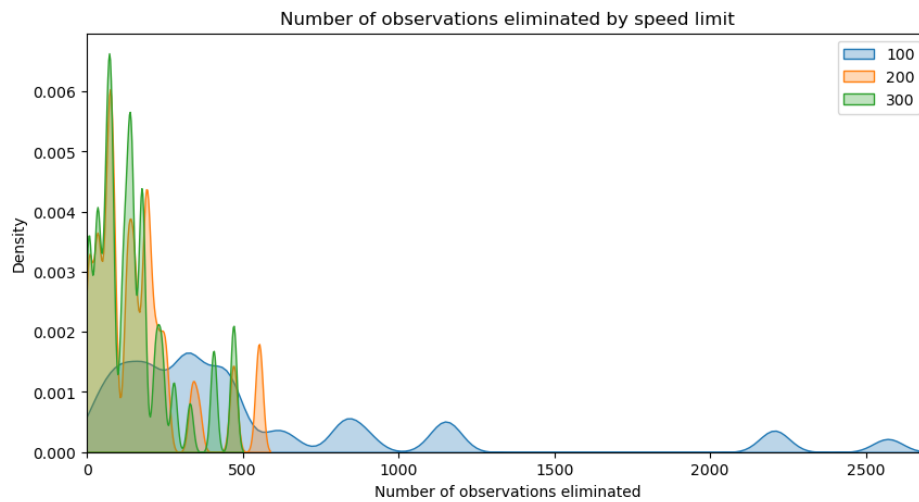


Figure A.15: The number of observations eliminated by varying speed limits across the users' data analyzed in Section 2.6

Table A.1: Parameters of benchmark methods.

Method	Parameter Count	Parameters	Notes
SES	1	α	Controls the weight given to the most recent observation when estimating the baseline of the time series.
Holt–Winters	3	α	Determines the weight given to the recent change in the level of the data.
		β_s	Determines the weight given to the recent change in the level of the data.
		β_d	Determines the dampening of the trend.
ES	4	α	Same as above.
		β	Determines the weight given to the recent change in the level of the data.
		γ	Controls the weight assigned to the seasonal component.
		m	Represents the length of the seasonal cycle (number of seasonal periods).
ARIMA	5–10	p	Number of lagged observations used to model the current value.
		d	Number of differencing operations to make the series stationary.
		q	Number of lagged forecast errors used to model the time series.
SARIMAX	8–20	p, d, q	Same as above.
		P	Number of seasonal autoregressive terms.
		D	Number of seasonal differencing operations.
		Q	Number of seasonal moving average terms.
		s	Determines periodicity (number of periods in season).
MTGP–RBF	16	ℓ	Controls extrapolation distance; one lengthscale per input dimension (14 total).
		η	Weight of a kernel component (output variance).
		ε	Homoscedastic Gaussian noise.
MTGP (ours)	21	ℓ, η, ε	Same as above (two per dimension).
		p	Determines the distance between repetitions of the function (two total).
		α	Controls the weighting of large-scale and small-scale variations (two total).

Table A.2: Parameters used during robustness experiments.

Method	Parameter	Value
Preprocessing	Max Speed (km/h)	200
	Spatial Radius for Compression (km)	0.3
	Uncertainty Filter (m)	100
Adam	Learning Rate for Adam	0.3
	Max Number of Training Iterations	150
Periodic Kernels	Initial Period Length for $K_{\text{PER},1}$	1440
	Initial Period Length for $K_{\text{PER},2}$	10080
Rational Quadratic Kernels	Initial Lengthscale for t_u	$\frac{1}{2N} \sum_{i=1}^N t_{u,i} - t_{u,i-1}$
	Initial Lengthscale for $[t_{s_1}, \dots, t_{s_m}]$	1
Search Range for ARIMA	Order of the AR, p	[0, 3]
	Order of the differencing, d	[0, 2]
	Order of the MA, q	[0, 3]
Search Range for SARIMAX	Order of the AR, p	[0, 3]
	Order of the differencing, d	[0, 2]
	Order of the MA, q	[0, 3]
	Order of the seasonal AR, P	[0, 3]
	Order of the seasonal differencing, D	[0, 2]
	Order of the seasonal MA, Q	[0, 3]
	Order of the periodicity, s	24

Appendix 2

SUPPLEMENTARY MATERIAL FOR CHAPTER 3

B.1 List of Notation Used

Table [B.1](#) lists every variable and its definition used in Chapter 3.

B.2 Dataset and Software Implementation

For Sections [3.4.1](#) through [3.4.3](#), we employ privacy-protected, passively-generated GPS data from Spectus, a U.S.-based data solution provider specializing in geospatial analytics. The dataset contains discrete GPS points of 2,000 anonymous, opted-in individuals in the Greater Seattle Area between December 2019 and July 2020. The location data recorded includes timestamps, unique device identifiers, latitudes and longitudes, and a measure of data precision (i.e., a spatial radius for which the provider has 95% confidence in the reported coordinates). For the experiment in Section [3.4.3](#), we sequentially increased the number of training points in each split while keeping the number of testing points constant, as shown in Figure [B.1](#).

For the experiment in Section [3.4.4](#), we use Microsoft Asia’s GeoLife dataset, a comprehensive collection of location-based data gathered from GPS-enabled devices. This dataset spans from April 2007 to April 2012, amassed through the voluntary participation of users who carried GPS-enabled devices, such as smartphones or dedicated GPS trackers. These devices continuously recorded the users’ location coordinates, capturing their movements and activities throughout the day. Besides being publically available, a huge benefit of this data source is that data points are more evenly distributed in time (as GPS-deviced were pinged regularly) giving a more accurate estimation of individuals’ route choices. Our experiments specifically used the subset of all GeoLife users who had mode labels available. We discarded training sets with less than 100 points (due to numerical instabilities in matrix operations) and more than 2,000 points (due to time constraints). We only considered trips made via buses, bikes, and by walking, as the other modes had too few trips to attain meaningful results.

Table B.1: Variable dictionary

General Notation	
λ	Longitude
ϕ	Latitude
v	Velocity
β	Bearing
Subscripts i and g	Indices for observations
Subscript j	Index for task
Subscript u	Index for input dimension
GP Notation	
Y	Matrix of outputs
T	Matrix of temporal variables
P	Matrix of physical variables
X	Matrix of inputs including T and P
ϵ	Gaussian white noise
d	Number of input dimensions
n	Number of training data points
m	Number of inferred data points
c	Number of inducing points
k	Covariance function of a GP
K	Covariance matrix
K^f	Inter-task covariance matrix
k	Vector of covariances
\otimes	Kronecker product
μ_*	Posterior mean function of the multi-task GP
σ_*^2	Posterior variance of the multi-task GP
μ_{**}	Posterior mean function of the physics-regularized multi-task GP
σ_{**}^2	Posterior variance of the physics-regularized multi-task GP
Θ	Set of GP parameters
l_k	Lengthscale parameter of kernel k
α	Scale mixture parameter of an RQ kernel
p	Period length parameter for a PER kernel
MKL Notation	
B	Set of base kernels
M	Number of MKL branches
A	Set of algebraic operations
η_h	Weight of kernel component h

bold denotes vectors.

Capital Bold denotes matrices.

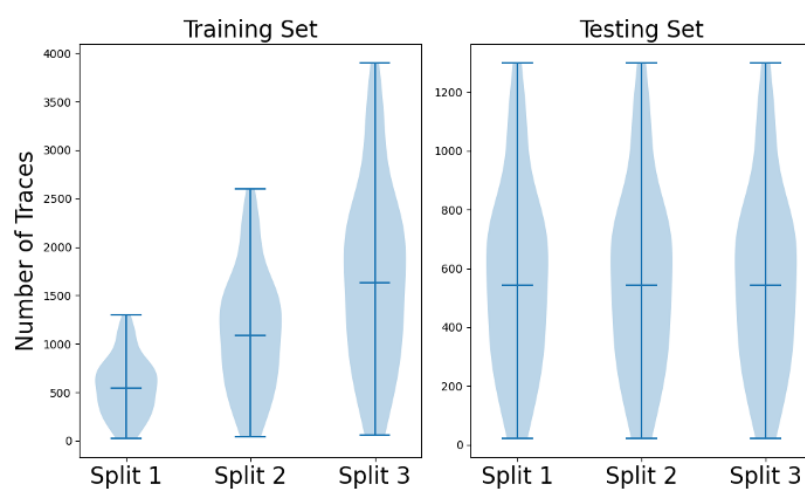


Figure B.1: Violin plots of training and testing set sizes across our passively-collected mobile data experiments

We implemented our methodology in the PYTHON programming language (van Rossum, 1995). We used GPYTORCH (Gardner et al., 2018) to increase efficiency and speed up matrix inversions within our GP framework. Furthermore, we used PANDAS (Wes McKinney, 2010) to read and wrangle the mobile dataset from which we derived the GP models, NUMPY (Harris et al., 2020) for a variety of array operations, and SCIKIT-MOBILITY (Pappalardo et al., 2019) to conduct data cleaning and preprocessing. Last but not least, we used MATPLOTLIB (Hunter, 2007) to produce all visualizations of our results.

B.3 Descriptive Statistics

Most data points exhibit a high level of precision, with approximately 95% of all traces located within a 65-meter radius of their true position (Figure B.2a). This represents a significant improvement compared to a previous study conducted in 2018 using a similar dataset, where only around 7% of observations had a precision worse than 1,000 meters (Ban et al., 2018). Regarding the temporal distribution of observations throughout the day (Figure B.2b), our findings align with that of Ban et al. (2018). We observe two distinct peaks on weekdays: one in the morning, around 7 am, and another in the evening, around 6 pm. On weekends, a single mid-day peak is evident. Moreover, the number of observations exhibits a positive trend as the hours progress, which holds true for every day.

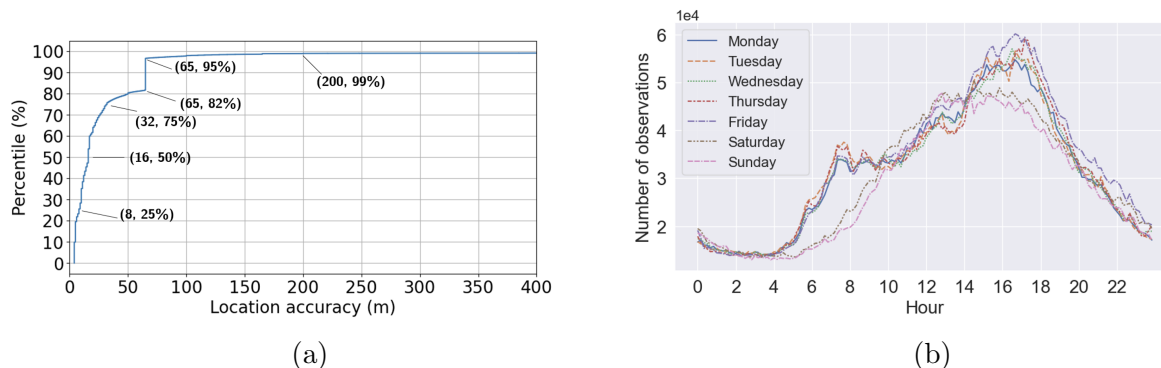


Figure B.2: Descriptive analysis of Spectus data for all 2,000 users. (a) Cumulative distribution of location accuracy. (b) Distribution of observations within each day of the week.

We find trivial differences between the entire dataset and the 33 users sampled for the experiments in Section 3.4.3, confirmed by looking at various metrics. Figure B.3 highlights the distribution of location accuracies between the two sets, showing that the experiments sample does not contain any outliers. Figure B.4 shows the distribution of observation counts by time of day (normalized by the total count of observations) between the two datasets, where we observe slight differences in the weekday distributions. Finally, Figure B.5 shows the distribution of sampling rates over time in the dataset. We observe slightly higher intervals between consecutive observations in certain weeks, but overall the two distribution look fairly similar.

B.4 Data Preprocessing and Modeling Considerations

We preprocess raw GPS traces to remove noisy data in two ways. First, we filter out rows using an upper bound on velocity (200 km/h, as the crow flies), which removes oscillations and physically-impossible jumps that may be present in raw mobile data. We further remove observations with a precision radius greater than 300 meters. The potential of these noisy points to provide valuable location information is shadowed by the problems they cause for model calibration, which disrupts the continuity of smooth trajectories.

We also reduce the size of trajectories through a compression algorithm (Algorithm 4) that aims to generate an approximated trajectory that largely retains the shape of the original trajectory. If two points are within a certain Haversine distance (we use 100 meters), we

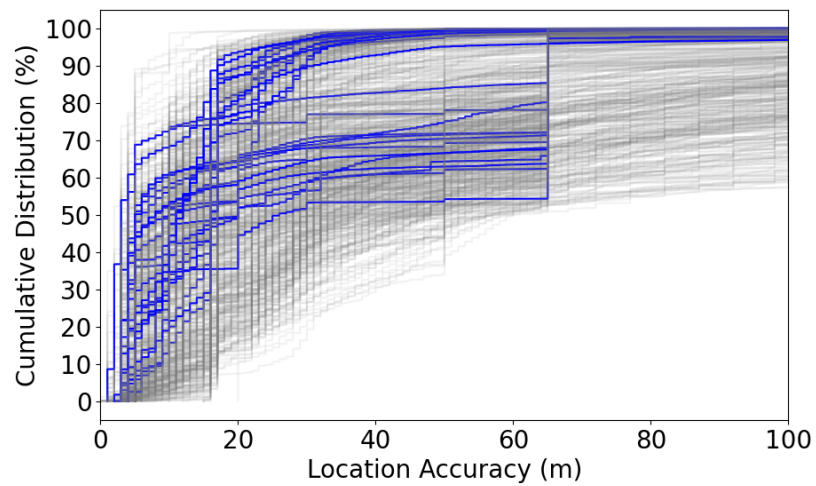


Figure B.3: Cumulative distribution of location accuracy for Spectus users in the Experiments sample (blue) and the whole dataset (grey)

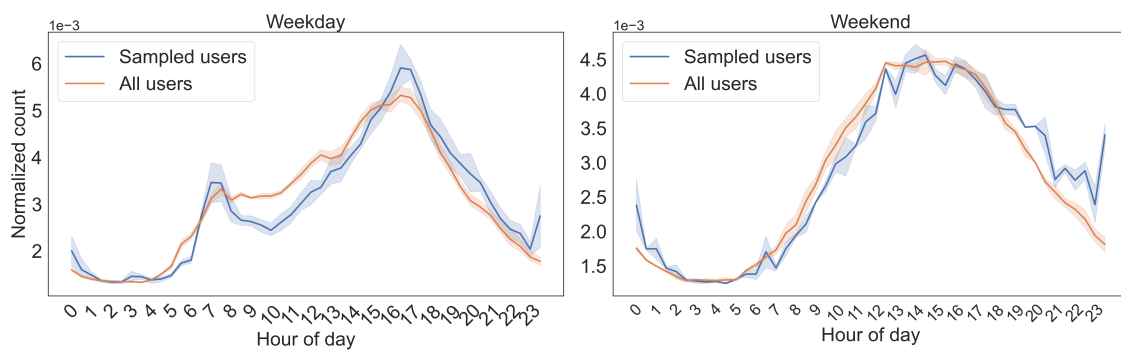


Figure B.4: Normalized observation counts by time of day for Spectus users in the Experiments sample (blue) and the whole dataset (orange)

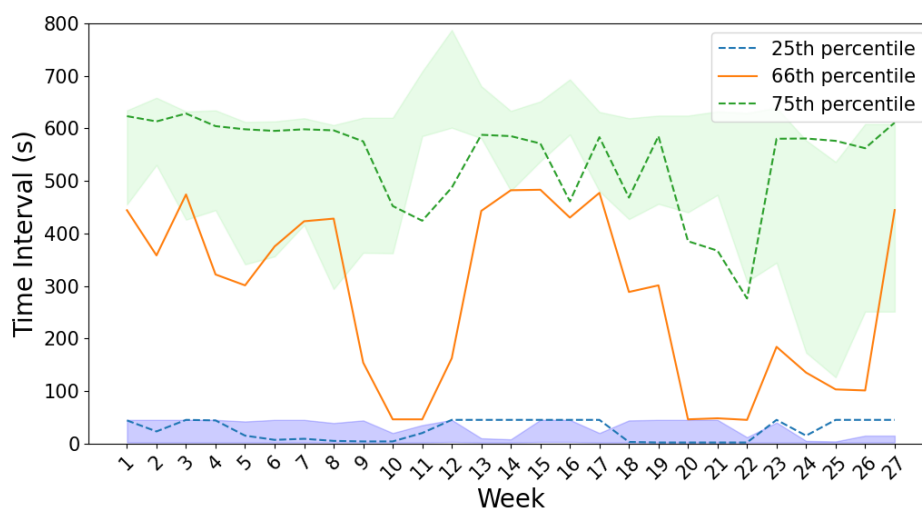


Figure B.5: Evolution of time interval between two consecutive points Week 1 to 27. Lines show the percentiles for the experiment sample, while the green and blue fillings shows the 70-80 and 15-35 percentile range for the entire dataset, respectively.

replace the two entries with their median—this is a variation of the Douglas-Peucker algorithm (Ugurel et al., 2024a). We do this to ease the computational burdens of training our GP framework, which can be burdensome for large trajectory sets. Because many GPS points are close in time and space, this allows us to significantly reduce the size of our training sample without compromising accuracy.

A practical consideration with using a greedy strategy to learn the form of the best-fitting kernel is that of computational complexity. For this, we suggest leveraging “the median trick” (Schölkopf and Smola, 2002)—instead of minimizing the negative marginal log likelihood to learn the parameters of the kernel, they can be manually set using apriori knowledge of the domain and the data. For example, the lengthscale parameter can be set as the median gap between observations in the input space. This idea gets tricky with parameters that do not relate to the smoothness of the function (i.e., the outputscale or the period length), but these too can be experimented with in a systematic way to reduce the computational burden.

B.5 Predictor Variables

Table B.2 shows the temporal and physical variables we use as predictors in our experiments. As discussed in Section 3.3.4, we leverage various calendar-based structures to encode more

Table B.2: Temporal and physical dimensions used in our experiments. Asterisk* denotes a dimension that was not employed in the GeoLife experiments.

Variable	Notation	Type	Model Inputs
Unix time (normalized)	\mathbf{t}_u	Continuous	$[0, 1, \dots, n]$
Seconds (after midnight) - Sine	\mathbf{t}_{ss}	Continuous	$[0, \dots, 1]$
Seconds (after midnight) - Cosine	\mathbf{t}_{sc}	Continuous	$[0, \dots, 1]$
Day of week	\mathbf{t}_d	Categorical (Binary)	$[0, 1]$
AM peak*	\mathbf{t}_{am}	Binary	$[0, 1]$
PM peak*	\mathbf{t}_{pm}	Binary	$[0, 1]$
Average segment velocity	\mathbf{v}	Continuous	$[0, \infty]$
Bearing	β	Continuous	$[0, 360]$

regularity into the temporal GP. We normalize Unix time by equating a user’s first observation to 0 and increasing it linearly thereafter. The second of the day is converted to cyclical encoding by calculating the sine and cosine values, wrapping smoothly around the boundaries (e.g., 23:59:59 and 00:00:00). Monday is the first day of the week and it is denoted with a 0. The AM peak variable is active in any observation between 6:00 AM and 10:00 AM, while the PM peak is between 3:00 PM and 7:00 PM (both in local time). The average segment velocity and bearing are calculated according to Equations 3.7 and 3.8, respectively.

B.6 Algorithms

Algorithms 3-6 are used inside Algorithm 7 (which is for the GeoLife experiments), while Algorithm 1 is also employed in the rest of the experiments.

Algorithm 3: Filter Trips

Input: Points

Output: Filtered points

Function *FilterTrips(points)*:

Filter points outside Beijing based on longitude and latitude constraints;
 Remove trips going outside the boundary;
 Filter points if segment velocity between two consecutive points is greater than 300 km/h;
return *Filtered points*;

Algorithm 4: Trajectory Data Compression

Input: Individual's trajectory data; spatial radius parameter r
Output: Compressed trajectory data
 LenTraj \leftarrow Count the total number of observations in the raw trajectory;
 Sorting: Sort by 'User ID' and 'Datetime';
 Initialization: Create lists $[\phi]$ and $[\lambda]$, initializing them with ϕ_i and λ_i . Initialize $i = 0$ and $j = 1$;
for $i < LenTraj$ **do**
 | **for** $j < LenTraj + 1$ **do**
 | | $d_{ij} \leftarrow$ Measure Haversine distance between (ϕ_i, λ_i) and (ϕ_j, λ_j) ;
 | | **if** $d_{ij} > r$ **then**
 | | | break // End current for loop;
 | | **end**
 | | Add ϕ_j and λ_j to lists $[\phi]$ and $[\lambda]$;
 | | $j \leftarrow j + 1$;
 | **end**
 | $(\phi, \lambda) \leftarrow$ np.median($[\phi]$), np.median($[\lambda]$) // Replace each point in $[\phi]$ and $[\lambda]$ with the median point of each list;
 | $i \leftarrow i + j$;
 | $j \leftarrow i + 1$ // Update indices so that the compressed points are skipped in the next iteration;
end
return *Output* // The compressed trajectory data;

Algorithm 5: Elbow Method

Input: Data, max_clusters, plot
Output: Optimal number of clusters
Function *FindOptimalTripClusters(data, max_clusters)*:
 | Initialize empty lists to store the inertia values and number of clusters;
 | inertia_values $\leftarrow []$;
 | num_clusters $\leftarrow []$;
 | **for** k **to** 2 **to** $max_clusters + 1$ **do**
 | | Apply K-means clustering;
 | | Compute inertia and number of clusters;
 | | Append inertia value and number of clusters to the lists;
 | **end**
 | Calculate differences in inertia values;
 | Calculate differences ratio;
 | Find the index of the elbow point;
 | **return** *Optimal number of clusters*;

Algorithm 6: Train-Test Split

Input: DataFrame, k , random_state**Output:** List of K folds**Function** *TripLabelBasedKFoldSplit*(df , k , $random_state$):

```

  Get unique trip IDs from the dataframe;
  Initialize KFold;
  List to hold the training and testing dataframes for each fold;
  Perform the k-fold split on the trip IDs;
  return List of  $K$  folds;

```

Algorithm 7: Trip Sampling

Input: Set of trips $Z_{w,e}$ for user w and mode e **Output:** Three folds for training and testing**Function** *TripSampling*($Z_{w,e}$):

```

  FilteredTrips  $\leftarrow$  FilterTrips( $Z_{w,e}$ , 300 km/h);
  SimplifiedTrips  $\leftarrow$  TrajectoryDataCompression( $FilteredTrips$ , 0.2 km radius);
  if  $|SimplifiedTrips| < 10$  then
    | Continue to next set of trips;
  end
  else if  $10 \leq |SimplifiedTrips| \leq 20$  then
    | End;
  end
  else
    |  $k \leftarrow$  ElbowMethod(KMeansClustering( $SimplifiedTrips$ ));
    | Clusters  $\leftarrow$  KMeansClustering( $k$ ,  $SimplifiedTrips$ );
    | Sort Clusters by size in descending order;
    | TopTwoClusters  $\leftarrow$  First two clusters in Clusters;
    | if  $|SimplifiedTrips| > 20$  then
      | | Shuffle SimplifiedTrips;
      | | SampledTrips  $\leftarrow$  First 20 trips in SimplifiedTrips;
    | end
    | else if  $|SimplifiedTrips| < 10$  then
      | | MoreClusters  $\leftarrow$  KMeansClustering( $SimplifiedTrips$ ,  $k + 1$ );
      | | SimplifiedTrips;
    | end
    | else
      | | End;
    | end
  end
  TrainTestSplit( $SimplifiedTrips$ );

```

B.7 Handling Nonstationarity with GPs

Human mobility has various spatial and temporal scales (Alessandretti et al., 2020). Specifically, as a person moves, time determines their location at the level of minutes and seconds, while when that person is stationary at an activity location, time impacts their travel behavior at the level of hours. For modeling purposes, this implies that in a traditional time-series GP with a single temporal input, the lengthscale parameter of the GP will not be able to capture the impact of time evenly between the two states (moving and staying).

There are multiple ways to overcome this problem. The most common, which is the one we rely on in this paper, is to represent time using various calendar-based structures (e.g., seconds of the day, days of the week). This allows the GP to learn unique lengthscales across more degrees of freedom, providing the adaptive local smoothing necessary for non-stationary data. For example, while Unix time may have a small lengthscale (capturing second-level variations), other variables capturing the time of the day may have a lengthscale on the order of hours.

GPs are naturally well-suited to handle statistical nonstationarity as well. Duvenaud (2014) details various approaches to combining different kernels to attain nonstationary processes without having to detrend data. This tends to require usage of nonstationary kernels, like the Linear kernel (Equation B.1), which produces different predictions if the data were moved while the kernel parameters were kept fixed.

$$k_{LIN} = \mathbf{x}^\top \mathbf{x}' \tag{B.1}$$

To showcase GP’s flexibility in handling nonstationarity, we have included an example of modeling the Mauna Loa CO2 concentration dataset using a GP with feature engineering (and without detrending). This dataset exhibits an increasing mean in the observed values over years. The usual approach to modeling this type of data using linear models is to first detrend the data (i.e., as shown here). However, we show that by creating cyclical features like months and years and leveraging the right kernel structure, this type of data can be modeled accurately using GPs (as seen in Figure B.6).

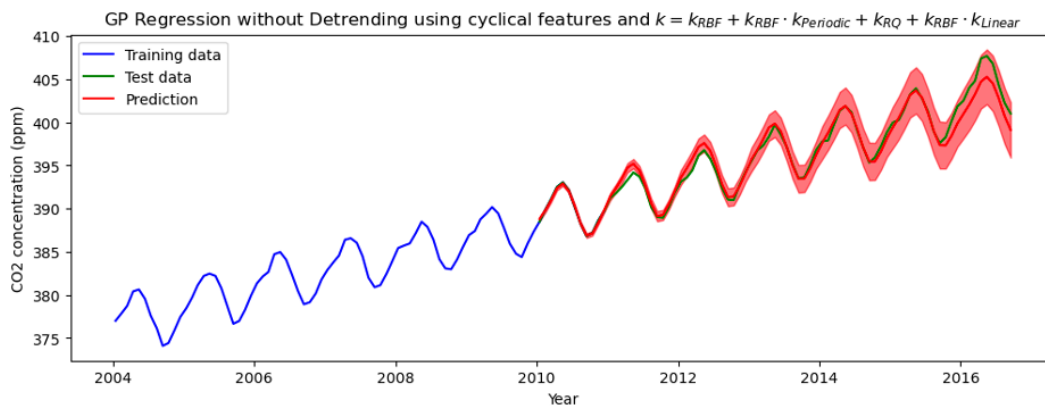


Figure B.6: GP regression using cyclical features and a composite kernel accurately captures nonstationary behavior in CO2 concentration of the Mauna Loa volcano

Appendix 3

SUPPLEMENTARY MATERIAL FOR CHAPTER 4

C.1 *Experimental Details and Full Model Results*

We detail the experimental setup for the rest of the classifiers used in the two experiments discussed in the main text. Both the multi-task DNNs and their single-task counterparts were built with two hidden layers, each separated by ReLU activations, and a uniform dropout rate of 0.3 between layers. For optimization, we performed a grid search over learning rates $\{10^{-3}, 10^{-4}, 5 \times 10^{-5}, 10^{-5}\}$, batch sizes $\{16, 32, 64, 128\}$, and weight decay values $\{10^{-3}, 10^{-4}, 10^{-5}\}$. The best configuration, determined via validation loss, used a learning rate of 5×10^{-5} , batch size of 64, and weight decay of 10^{-4} . Models were trained for up to 200 epochs with early stopping if the validation loss did not improve for 20 consecutive epochs.

To isolate the effect of parameter sharing (MT) from raw capacity, we matched aggregate hidden width across conditions. The unified MT model used a shared trunk of $256 \rightarrow 128$ units, while each single-task variant (STV) used $64 \rightarrow 32$ units. Since there are four STVs, their combined width ($4 \times 64, 4 \times 32$) is comparable to the MT trunk (256, 128), making representational capacity roughly parity-matched while differing in whether features are learned jointly or separately.

Despite the smaller per-model widths, STVs require training four independent networks (four forward/backward passes, four optimizers, four early-stopping loops) and thus incur higher wall-clock time than the single unified MT model, which amortizes the trunk compute across tasks. Empirically, STVs were consistently slower than MT across data fractions (see Figure C.1).

For the benchmark models leveraged in the *uplift* experiments, we used SKLEARN’s built-in **random forest (RF) classifier**, **gradient boosting (GB) classifier**, and **C-Support Vector (SVC) classifier**. For RFs and SVCs, we performed a small grid search over the number of estimators $\{100, 300, 500\}$ and regularization parameter $C \in \{0.1, 1, 10\}$, respectively. We selected 500 estimators for RFs and $C = 1$ for SVCs based on validation performance. Due to computational load, we kept the number of GB estimators at 100. We also enabled probability

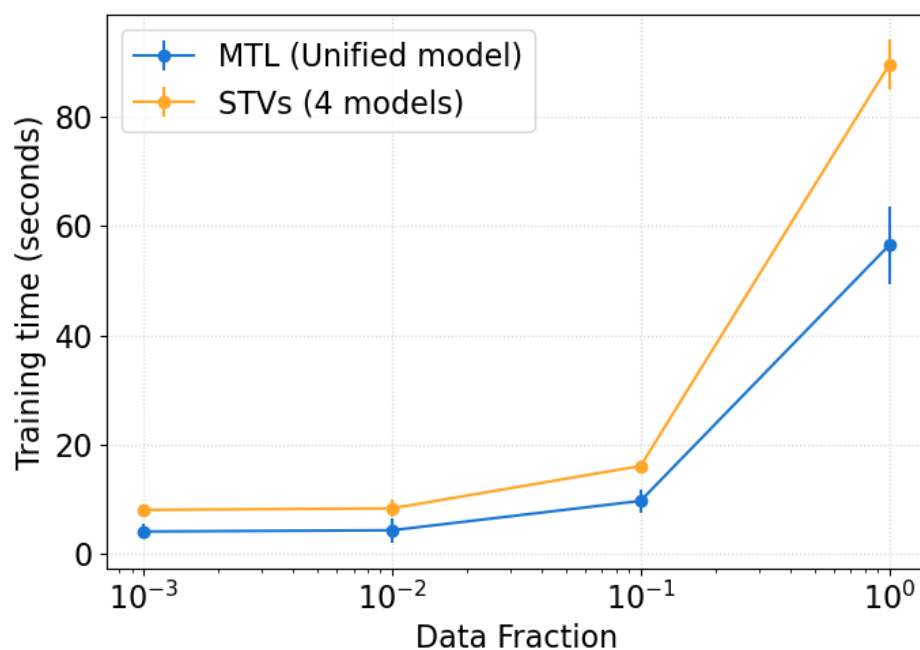


Figure C.1: Comparison of training times between the unified multi-task learning (MT) model and four separate single-task variants (STVs) across varying data fractions (log-scaled on x-axis). Despite its larger architecture, the MT model trains faster overall because its shared trunk amortizes computation across tasks, whereas STVs require four independent forward-backward passes. Error bars show the standard deviation across cross-validation folds.

Table C.1: Performance across feature sets and models for Age (overall split). Values are mean \pm sd across folds; best per model in **bold**; best in metric in **red**.

Metric	Feature set	DNN	RF	GBM	SVM
Accuracy	C	0.508 \pm 0.003	0.483 \pm 0.003	0.503 \pm 0.000	0.504 \pm 0.000
	+ST	0.524 \pm 0.004	0.516 \pm 0.004	0.513 \pm 0.000	0.525 \pm 0.001
	+D	0.525 \pm 0.004	0.517 \pm 0.003	0.515 \pm 0.000	0.523 \pm 0.000
	+M	0.526 \pm 0.002	0.518 \pm 0.003	0.514 \pm 0.000	0.520 \pm 0.000
	+CT	0.525 \pm 0.003	0.520 \pm 0.003	0.518 \pm 0.000	0.527 \pm 0.001
AUROC	C	0.840 \pm 0.001	0.815 \pm 0.000	0.835 \pm 0.000	0.824 \pm 0.000
	+ST	0.848 \pm 0.001	0.834 \pm 0.001	0.842 \pm 0.000	0.838 \pm 0.000
	+D	0.847 \pm 0.001	0.834 \pm 0.001	0.842 \pm 0.000	0.838 \pm 0.000
	+M	0.847 \pm 0.001	0.834 \pm 0.001	0.842 \pm 0.000	0.836 \pm 0.000
	+CT	0.850 \pm 0.001	0.837 \pm 0.001	0.843 \pm 0.000	0.838 \pm 0.000
NLL	C	1.084 \pm 0.003	1.277 \pm 0.007	1.099 \pm 0.000	1.134 \pm 0.000
	+ST	1.059 \pm 0.003	1.117 \pm 0.006	1.080 \pm 0.000	1.096 \pm 0.000
	+D	1.062 \pm 0.002	1.118 \pm 0.003	1.081 \pm 0.000	1.096 \pm 0.000
	+M	1.061 \pm 0.003	1.120 \pm 0.006	1.081 \pm 0.000	1.099 \pm 0.000
	+CT	1.048 \pm 0.002	1.098 \pm 0.003	1.071 \pm 0.000	1.086 \pm 0.000
ECE	C	0.021 \pm 0.004	0.060 \pm 0.004	0.028 \pm 0.001	0.034 \pm 0.001
	+ST	0.020 \pm 0.005	0.029 \pm 0.003	0.016 \pm 0.001	0.027 \pm 0.002
	+D	0.016 \pm 0.003	0.032 \pm 0.006	0.017 \pm 0.000	0.022 \pm 0.003
	+M	0.019 \pm 0.006	0.034 \pm 0.008	0.017 \pm 0.000	0.021 \pm 0.001
	+CT	0.018 \pm 0.006	0.034 \pm 0.003	0.016 \pm 0.000	0.023 \pm 0.002

estimates for the SVC classifier, which slowed down convergence. The rest of this section contains the full results tables in C.1 until C.8.

C.1.1 Overall Split

Tables C.1-C.4 show the full model results in the overall split.

C.1.2 Cross-temporal Split

Tables C.5-C.8 show the full model results in the overall split.

Table C.2: Performance across feature sets and models for Gender (overall split). Values are mean \pm sd across folds; best per model in **bold**; best in metric in **red**.

Metric	Feature set	DNN	RF	GBM	SVM
Accuracy	C	0.534 \pm 0.002	0.519 \pm 0.002	0.528 \pm 0.000	0.530 \pm 0.001
	+ST	0.541 \pm 0.005	0.541 \pm 0.002	0.547 \pm 0.000	0.538 \pm 0.000
	+D	0.537 \pm 0.005	0.541 \pm 0.006	0.543 \pm 0.000	0.541 \pm 0.000
	+M	0.539 \pm 0.007	0.541 \pm 0.004	0.545 \pm 0.000	0.537 \pm 0.000
	+CT	0.541 \pm 0.007	0.540 \pm 0.006	0.549 \pm 0.000	0.546 \pm 0.001
AUROC	C	0.581 \pm 0.006	0.558 \pm 0.002	0.588 \pm 0.000	0.588 \pm 0.001
	+ST	0.592 \pm 0.005	0.587 \pm 0.004	0.606 \pm 0.000	0.588 \pm 0.000
	+D	0.601 \pm 0.002	0.591 \pm 0.006	0.628 \pm 0.000	0.596 \pm 0.000
	+M	0.602 \pm 0.005	0.594 \pm 0.003	0.631 \pm 0.000	0.608 \pm 0.000
	+CT	0.611 \pm 0.002	0.593 \pm 0.007	0.627 \pm 0.000	0.612 \pm 0.000
NLL	C	0.812 \pm 0.001	0.929 \pm 0.002	0.813 \pm 0.000	0.814 \pm 0.000
	+ST	0.808 \pm 0.001	0.832 \pm 0.002	0.806 \pm 0.000	0.810 \pm 0.000
	+D	0.807 \pm 0.001	0.831 \pm 0.003	0.802 \pm 0.000	0.809 \pm 0.000
	+M	0.808 \pm 0.001	0.829 \pm 0.000	0.802 \pm 0.000	0.808 \pm 0.000
	+CT	0.805 \pm 0.001	0.829 \pm 0.002	0.801 \pm 0.000	0.806 \pm 0.000
ECE	C	0.012 \pm 0.002	0.077 \pm 0.002	0.019 \pm 0.000	0.008 \pm 0.001
	+ST	0.013 \pm 0.004	0.040 \pm 0.006	0.013 \pm 0.000	0.005 \pm 0.001
	+D	0.018 \pm 0.001	0.041 \pm 0.007	0.015 \pm 0.000	0.010 \pm 0.001
	+M	0.019 \pm 0.005	0.035 \pm 0.003	0.014 \pm 0.001	0.006 \pm 0.001
	+CT	0.021 \pm 0.006	0.039 \pm 0.004	0.013 \pm 0.000	0.013 \pm 0.002

Table C.3: Performance across feature sets and models for HH Income (overall split). Values are mean \pm sd across folds; best per model in **bold**; best in metric in **red**.

Metric	Feature set	DNN	RF	GBM	SVM
Accuracy	C	0.532 \pm 0.001	0.528 \pm 0.002	0.535 \pm 0.000	0.532 \pm 0.000
	+ST	0.534 \pm 0.001	0.561 \pm 0.001	0.531 \pm 0.000	0.535 \pm 0.001
	+D	0.534 \pm 0.001	0.563 \pm 0.001	0.534 \pm 0.000	0.536 \pm 0.001
	+M	0.534 \pm 0.001	0.561 \pm 0.001	0.532 \pm 0.000	0.536 \pm 0.001
	+CT	0.538 \pm 0.000	0.561 \pm 0.001	0.536 \pm 0.000	0.540 \pm 0.001
AUROC	C	0.615 \pm 0.003	0.617 \pm 0.001	0.612 \pm 0.000	0.599 \pm 0.001
	+ST	0.626 \pm 0.003	0.677 \pm 0.001	0.638 \pm 0.000	0.626 \pm 0.000
	+D	0.627 \pm 0.004	0.678 \pm 0.002	0.639 \pm 0.000	0.630 \pm 0.000
	+M	0.629 \pm 0.004	0.678 \pm 0.002	0.640 \pm 0.000	0.631 \pm 0.000
	+CT	0.639 \pm 0.003	0.691 \pm 0.001	0.650 \pm 0.000	0.644 \pm 0.000
NLL	C	1.292 \pm 0.003	1.556 \pm 0.013	1.289 \pm 0.000	1.305 \pm 0.000
	+ST	1.285 \pm 0.002	1.216 \pm 0.001	1.275 \pm 0.000	1.291 \pm 0.000
	+D	1.286 \pm 0.002	1.216 \pm 0.002	1.273 \pm 0.000	1.288 \pm 0.000
	+M	1.284 \pm 0.003	1.217 \pm 0.003	1.272 \pm 0.000	1.287 \pm 0.000
	+CT	1.275 \pm 0.002	1.202 \pm 0.002	1.263 \pm 0.000	1.275 \pm 0.000
ECE	C	0.029 \pm 0.004	0.048 \pm 0.002	0.029 \pm 0.000	0.025 \pm 0.001
	+ST	0.034 \pm 0.007	0.020 \pm 0.003	0.022 \pm 0.000	0.041 \pm 0.004
	+D	0.037 \pm 0.011	0.022 \pm 0.003	0.031 \pm 0.001	0.047 \pm 0.001
	+M	0.033 \pm 0.006	0.020 \pm 0.002	0.031 \pm 0.000	0.047 \pm 0.001
	+CT	0.035 \pm 0.007	0.027 \pm 0.004	0.024 \pm 0.000	0.064 \pm 0.001

Table C.4: Performance across feature sets and models for Number of Children (overall split). Values are mean \pm sd across folds; best per model in **bold**; best in metric in **red**.

Metric	Feature set	DNN	RF	GBM	SVM
Accuracy	C	0.776 \pm 0.001	0.783 \pm 0.001	0.779 \pm 0.000	0.774 \pm 0.000
	+ST	0.780 \pm 0.002	0.807 \pm 0.001	0.784 \pm 0.000	0.782 \pm 0.000
	+D	0.782 \pm 0.001	0.805 \pm 0.000	0.783 \pm 0.000	0.786 \pm 0.000
	+M	0.781 \pm 0.002	0.804 \pm 0.001	0.781 \pm 0.000	0.786 \pm 0.001
	+CT	0.785 \pm 0.002	0.808 \pm 0.001	0.787 \pm 0.000	0.791 \pm 0.000
AUROC	C	0.837 \pm 0.001	0.823 \pm 0.001	0.842 \pm 0.000	0.831 \pm 0.000
	+ST	0.847 \pm 0.002	0.861 \pm 0.001	0.852 \pm 0.000	0.844 \pm 0.000
	+D	0.849 \pm 0.001	0.862 \pm 0.002	0.853 \pm 0.000	0.845 \pm 0.000
	+M	0.847 \pm 0.003	0.861 \pm 0.001	0.854 \pm 0.000	0.844 \pm 0.000
	+CT	0.854 \pm 0.002	0.862 \pm 0.001	0.859 \pm 0.000	0.851 \pm 0.000
NLL	C	0.639 \pm 0.001	0.801 \pm 0.008	0.634 \pm 0.000	0.650 \pm 0.000
	+ST	0.620 \pm 0.003	0.595 \pm 0.003	0.615 \pm 0.000	0.625 \pm 0.000
	+D	0.619 \pm 0.002	0.597 \pm 0.008	0.613 \pm 0.000	0.623 \pm 0.000
	+M	0.621 \pm 0.004	0.596 \pm 0.005	0.613 \pm 0.000	0.623 \pm 0.000
	+CT	0.610 \pm 0.003	0.595 \pm 0.010	0.604 \pm 0.000	0.613 \pm 0.000
ECE	C	0.018 \pm 0.001	0.030 \pm 0.003	0.018 \pm 0.000	0.059 \pm 0.001
	+ST	0.021 \pm 0.003	0.055 \pm 0.001	0.019 \pm 0.000	0.045 \pm 0.002
	+D	0.022 \pm 0.003	0.057 \pm 0.002	0.023 \pm 0.000	0.035 \pm 0.001
	+M	0.023 \pm 0.004	0.058 \pm 0.001	0.020 \pm 0.000	0.038 \pm 0.001
	+CT	0.016 \pm 0.006	0.048 \pm 0.002	0.017 \pm 0.000	0.030 \pm 0.001

Table C.5: Performance across feature sets and models for Age (2017–2019 train, 2023 test split). Values are mean \pm sd across folds; best per model in **bold**; best in metric in **red**.

Metric	Feature set	DNN	RF	GBM	SVM
Accuracy	C	0.474 \pm 0.002	0.443 \pm 0.001	0.471 \pm 0.000	0.469 \pm 0.000
	+ST	0.488 \pm 0.004	0.480 \pm 0.002	0.489 \pm 0.000	0.484 \pm 0.001
	+D	0.480 \pm 0.003	0.479 \pm 0.004	0.485 \pm 0.000	0.484 \pm 0.001
	+M	0.475 \pm 0.004	0.478 \pm 0.003	0.483 \pm 0.000	0.480 \pm 0.001
	+CT	0.484 \pm 0.004	0.482 \pm 0.002	0.490 \pm 0.001	0.485 \pm 0.000
AUROC	C	0.814 \pm 0.001	0.781 \pm 0.001	0.813 \pm 0.000	0.802 \pm 0.000
	+ST	0.826 \pm 0.001	0.813 \pm 0.000	0.824 \pm 0.000	0.815 \pm 0.000
	+D	0.824 \pm 0.001	0.812 \pm 0.001	0.824 \pm 0.000	0.814 \pm 0.000
	+M	0.821 \pm 0.001	0.812 \pm 0.001	0.823 \pm 0.000	0.813 \pm 0.000
	+CT	0.824 \pm 0.001	0.815 \pm 0.001	0.826 \pm 0.000	0.815 \pm 0.000
NLL	C	1.181 \pm 0.004	1.505 \pm 0.006	1.176 \pm 0.000	1.216 \pm 0.001
	+ST	1.148 \pm 0.005	1.186 \pm 0.001	1.147 \pm 0.000	1.184 \pm 0.001
	+D	1.158 \pm 0.004	1.190 \pm 0.005	1.149 \pm 0.000	1.188 \pm 0.001
	+M	1.164 \pm 0.006	1.189 \pm 0.002	1.151 \pm 0.000	1.193 \pm 0.001
	+CT	1.153 \pm 0.006	1.173 \pm 0.004	1.135 \pm 0.000	1.180 \pm 0.000
ECE	C	0.034 \pm 0.003	0.087 \pm 0.001	0.039 \pm 0.000	0.055 \pm 0.001
	+ST	0.037 \pm 0.007	0.015 \pm 0.003	0.034 \pm 0.000	0.051 \pm 0.002
	+D	0.044 \pm 0.006	0.016 \pm 0.002	0.034 \pm 0.000	0.051 \pm 0.001
	+M	0.047 \pm 0.004	0.017 \pm 0.006	0.039 \pm 0.000	0.055 \pm 0.001
	+CT	0.050 \pm 0.005	0.019 \pm 0.003	0.036 \pm 0.001	0.057 \pm 0.001

Table C.6: Performance across feature sets and models for Gender (2017–2019 train, 2023 test split). Values are mean \pm sd across folds; best per model in **bold**; best in metric in **red**.

Metric	Feature set	DNN	RF	GBM	SVM
Accuracy	C	0.507 \pm 0.003	0.500 \pm 0.003	0.506 \pm 0.000	0.508 \pm 0.000
	+ST	0.508 \pm 0.003	0.517 \pm 0.002	0.524 \pm 0.000	0.511 \pm 0.001
	+D	0.511 \pm 0.002	0.513 \pm 0.002	0.523 \pm 0.000	0.513 \pm 0.000
	+M	0.507 \pm 0.002	0.509 \pm 0.003	0.523 \pm 0.000	0.509 \pm 0.001
	+CT	0.512 \pm 0.006	0.515 \pm 0.002	0.523 \pm 0.000	0.514 \pm 0.000
AUROC	C	0.551 \pm 0.001	0.514 \pm 0.001	0.546 \pm 0.000	0.542 \pm 0.000
	+ST	0.557 \pm 0.003	0.544 \pm 0.003	0.567 \pm 0.000	0.547 \pm 0.000
	+D	0.556 \pm 0.003	0.540 \pm 0.003	0.565 \pm 0.000	0.540 \pm 0.000
	+M	0.554 \pm 0.003	0.540 \pm 0.001	0.568 \pm 0.000	0.540 \pm 0.000
	+CT	0.556 \pm 0.001	0.544 \pm 0.002	0.570 \pm 0.000	0.545 \pm 0.000
NLL	C	0.961 \pm 0.003	1.396 \pm 0.008	0.976 \pm 0.000	0.953 \pm 0.001
	+ST	0.962 \pm 0.003	1.002 \pm 0.005	0.969 \pm 0.000	0.954 \pm 0.000
	+D	0.960 \pm 0.006	1.002 \pm 0.007	0.964 \pm 0.000	0.954 \pm 0.000
	+M	0.960 \pm 0.004	1.001 \pm 0.006	0.961 \pm 0.000	0.954 \pm 0.000
	+CT	0.963 \pm 0.004	1.007 \pm 0.011	0.962 \pm 0.000	0.954 \pm 0.001
ECE	C	0.056 \pm 0.005	0.108 \pm 0.003	0.058 \pm 0.000	0.039 \pm 0.002
	+ST	0.060 \pm 0.004	0.057 \pm 0.003	0.046 \pm 0.000	0.041 \pm 0.002
	+D	0.053 \pm 0.003	0.059 \pm 0.002	0.045 \pm 0.000	0.034 \pm 0.001
	+M	0.060 \pm 0.006	0.062 \pm 0.003	0.044 \pm 0.000	0.039 \pm 0.001
	+CT	0.060 \pm 0.004	0.058 \pm 0.002	0.044 \pm 0.000	0.037 \pm 0.002

Table C.7: Performance across feature sets and models for HH Income (2017–2019 train, 2023 test split). Values are mean \pm sd across folds; best per model in **bold**; best in metric in **red**.

Metric	Feature set	DNN	RF	GBM	SVM
Accuracy	C	0.588 \pm 0.001	0.542 \pm 0.002	0.582 \pm 0.000	0.586 \pm 0.001
	+ST	0.588 \pm 0.002	0.569 \pm 0.002	0.579 \pm 0.000	0.586 \pm 0.001
	+D	0.587 \pm 0.002	0.566 \pm 0.002	0.579 \pm 0.000	0.586 \pm 0.001
	+M	0.588 \pm 0.001	0.566 \pm 0.001	0.581 \pm 0.000	0.585 \pm 0.001
	+CT	0.585 \pm 0.002	0.562 \pm 0.002	0.577 \pm 0.000	0.581 \pm 0.001
AUROC	C	0.606 \pm 0.003	0.581 \pm 0.000	0.605 \pm 0.000	0.583 \pm 0.000
	+ST	0.616 \pm 0.004	0.610 \pm 0.002	0.617 \pm 0.000	0.599 \pm 0.000
	+D	0.614 \pm 0.002	0.610 \pm 0.001	0.620 \pm 0.000	0.598 \pm 0.000
	+M	0.617 \pm 0.002	0.610 \pm 0.002	0.615 \pm 0.000	0.604 \pm 0.000
	+CT	0.618 \pm 0.003	0.611 \pm 0.002	0.613 \pm 0.000	0.604 \pm 0.000
NLL	C	1.219 \pm 0.003	1.557 \pm 0.022	1.228 \pm 0.000	1.230 \pm 0.000
	+ST	1.213 \pm 0.005	1.244 \pm 0.002	1.225 \pm 0.000	1.224 \pm 0.000
	+D	1.214 \pm 0.005	1.244 \pm 0.002	1.221 \pm 0.000	1.225 \pm 0.000
	+M	1.214 \pm 0.003	1.244 \pm 0.002	1.226 \pm 0.000	1.222 \pm 0.000
	+CT	1.219 \pm 0.005	1.244 \pm 0.002	1.225 \pm 0.000	1.220 \pm 0.000
ECE	C	0.063 \pm 0.004	0.065 \pm 0.004	0.067 \pm 0.000	0.069 \pm 0.001
	+ST	0.058 \pm 0.002	0.080 \pm 0.001	0.068 \pm 0.000	0.070 \pm 0.001
	+D	0.052 \pm 0.008	0.079 \pm 0.001	0.064 \pm 0.000	0.072 \pm 0.001
	+M	0.057 \pm 0.006	0.080 \pm 0.002	0.070 \pm 0.000	0.072 \pm 0.001
	+CT	0.047 \pm 0.008	0.054 \pm 0.002	0.043 \pm 0.000	0.058 \pm 0.001

Table C.8: Performance across feature sets and models for Number of Children (2017–2019 train, 2023 test split). Values are mean \pm sd across folds; best per model in **bold**; best in metric in **red**.

Metric	Feature set	DNN	RF	GBM	SVM
Accuracy	C	0.747 \pm 0.002	0.741 \pm 0.003	0.745 \pm 0.000	0.754 \pm 0.000
	+ST	0.756 \pm 0.003	0.759 \pm 0.001	0.749 \pm 0.000	0.758 \pm 0.001
	+D	0.751 \pm 0.001	0.758 \pm 0.002	0.747 \pm 0.000	0.760 \pm 0.000
	+M	0.752 \pm 0.003	0.757 \pm 0.001	0.750 \pm 0.000	0.756 \pm 0.001
	+CT	0.753 \pm 0.002	0.758 \pm 0.001	0.747 \pm 0.000	0.757 \pm 0.000
AUROC	C	0.819 \pm 0.002	0.791 \pm 0.001	0.819 \pm 0.000	0.808 \pm 0.000
	+ST	0.825 \pm 0.002	0.809 \pm 0.001	0.825 \pm 0.000	0.814 \pm 0.000
	+D	0.820 \pm 0.001	0.811 \pm 0.002	0.826 \pm 0.000	0.814 \pm 0.000
	+M	0.821 \pm 0.001	0.810 \pm 0.001	0.826 \pm 0.000	0.814 \pm 0.000
	+CT	0.826 \pm 0.001	0.816 \pm 0.001	0.828 \pm 0.000	0.817 \pm 0.000
NLL	C	0.704 \pm 0.003	1.031 \pm 0.016	0.712 \pm 0.000	0.717 \pm 0.000
	+ST	0.690 \pm 0.002	0.729 \pm 0.005	0.695 \pm 0.000	0.705 \pm 0.000
	+D	0.695 \pm 0.002	0.722 \pm 0.006	0.695 \pm 0.000	0.704 \pm 0.000
	+M	0.696 \pm 0.005	0.720 \pm 0.005	0.695 \pm 0.000	0.706 \pm 0.001
	+CT	0.696 \pm 0.004	0.726 \pm 0.004	0.692 \pm 0.000	0.707 \pm 0.001
ECE	C	0.025 \pm 0.002	0.038 \pm 0.004	0.020 \pm 0.000	0.034 \pm 0.001
	+ST	0.021 \pm 0.004	0.040 \pm 0.001	0.016 \pm 0.000	0.023 \pm 0.001
	+D	0.024 \pm 0.003	0.040 \pm 0.002	0.016 \pm 0.000	0.022 \pm 0.001
	+M	0.026 \pm 0.003	0.041 \pm 0.001	0.013 \pm 0.000	0.017 \pm 0.001
	+CT	0.026 \pm 0.002	0.039 \pm 0.001	0.015 \pm 0.001	0.019 \pm 0.001