



A USDOT NATIONAL  
UNIVERSITY TRANSPORTATION CENTER

**Carnegie Mellon University**



---

# Labeling Roads with Different Types of Automated Driving Functional Requirements using Machine Learning

Ding Zhao  
Carnegie Mellon University

<https://orcid.org/0000-0002-9400-8446>

FINAL RESEARCH REPORT

Contract # 69A3551747111

## DISCLAIMER

The contents of this report reflect the views of the authors, who are responsible for the facts and the accuracy of the information presented herein. This document is disseminated under the sponsorship of the U.S. Department of Transportation's University Transportation Centers Program, in the interest of information exchange. The U.S. Government assumes no liability for the contents or use thereof.

## **Project Investigator and Contributors:**

Project Investigator: Ding Zhao (Carnegie Mellon University, ORCID: Update)

Contributors:

1. Mengdi Xu (Carnegie Mellon University, ORCID: 0000-0001-9332-4175)
2. Rui Chen (Carnegie Mellon University, ORCID: 0000-0002-8671-8771)
3. Mansur Arief (Carnegie Mellon University, ORCID: pending)
4. Weiyang Zhang (Carnegie Mellon University, ORCID: 0000-0002-4347-1756)
5. Wenshuo Wang (Carnegie Mellon University, ORCID: 0000-0002-1860-8351)

## **Problem Statement:**

Automated vehicles (AVs) should be deployed gradually and geometrically selectively to ensure safety. Frequent collisions of AVs in certain driving scenarios, such as in dark streets or crowded areas, have posed wide concerns of AV safety. People want to know what kinds of driving circumstances or areas are easy for AVs and what are relatively hard and how to quantify the degree. In the new Federal guidance - *Automated Vehicles 3.0: Preparing for the Future of Transportation 3.0*, released in 2018 [1], the Department of Transportation proposed the concept of Operational Design Domain (ODD) to describe the driving complexity considering roadway types, geographic area, and speed range. This concept sheds a light on the evaluation of the difficulty of driving, but it is still not clear how to apply it in practice as neither the automation level nor the ODD provide a numerical solution, hence likely resulting in subjective, incomplete, and inherently somewhat ambiguous analysis to fully describe the complex nature of real-world traffic, and thus causing biased confidence and disqualification of AVs for public deployment.

In order to reduce the risks and lay the foundation of autonomous vehicles deployment, this project aims to label the roads of the city with different risk levels based on large scale real-world datasets of Pittsburgh city including multi-dimensional and multi-fidelity data [2]. We develop a framework to identify typical driving scenarios based on Nonparametric Bayesian (NPBayes) methods [3,4]. We further visualize the scenarios using the velocity fields which help provide a reference for characterizing the complexity of the possible situations and the risk level.

## **Methods:**

Our analysis is based on the Argoverse tracking Dataset [2] which contains Pittsburgh traffic information. The dataset is collected from onboard sensors such as lidars and cameras. In each trial, it contains the location, type, and bounding box of surrounding objects. In the data preprocessing process, we first remap the ego-centric trajectory data of all the detected vehicles into the Pittsburgh City map and further implement the Kalman filter to denoise.

We proposed to use a Gaussian Process (GP) [5] for modeling each type of multi-vehicle interaction scenario. GP is a nonparametric model that outperforms deep neural nets in terms of

uncertainty modeling and data efficiency. GP has been widely used in representing continuous-time trajectories and simple object dynamics. Note that the number of the possible scenarios corresponding to each map layout is hard to pre-define. Instead of heuristically select a number, we hope to use a data-driven method that can learn the number directly from data. Therefore, we use a Dirichlet Process (DP) [6] as the prior of the GP mixture model. The DP prior is able to model an infinite number of clusters and has a nice rich-gets-richer clustering property [6]. Although Dirichlet Process Gaussian Process mixture model (DPGP) has been implemented in analyzing traffic flow [7], we use a velocity field representation of the traffic flow instead of clustering the trajectory for a single-vehicle. In other words, the velocity field representation can capture the multi-vehicle interaction information which is suitable to describe scenarios. The number of clustered scenarios corresponding to the map layout is a sign of the risk level.

The algorithmic efficiency and DPGP model accuracy are improved. We implemented the DPGP algorithm in well-known probabilistic programming platforms, GPyTorch [8]. Instead of using Exact GP to do the inference, we use a sparse GP [9] to increase the scalability which is crucial when dealing with big data. We further developed a webpage based on Google Maps to visualize the risk heat map in Pittsburgh city. When clicking at a specific area, it visualizes the velocity field of each clustered scenario and the number of the clusters as the numerical risk level.

### **Findings:**

We analyzed the risk levels of different regions of Pittsburgh City based on the Argoverse Dataset. The risk level is identified as the clustered scenario number. Each scenario is a cluster learned by the DPGP model. Our results mainly include the traffic scenario clustering method, the learned scenarios and the visualization website.

### **Conclusion:**

To the best of our knowledge, we are the first group to use data-driven methods to define driving scenario risks. Instead of manually designing complex features comprising the overall risk, we build a framework leveraging the Dirichlet Process Gaussian Process as an end-to-end way that takes large amounts of naturalistic driving data and outputs the risk level directly. Matching different levels of automated vehicles (AVs) with map areas in different risk levels can help improve traffic efficiency, such as reducing the probability that catastrophes happen and avoiding the traffic congestion caused by incapable autonomous vehicles.

The visualization tool, the risk heat map, could help the public to get a numerical sense of the complexity of the local roads. It also serves as a reference for the autonomous companies to develop strategies to test and deploy their autonomous vehicles with safety considered. If the related corporations will take the risk levels identified into consideration, especially the autonomous companies, we believe it helps to increase the human trust level to AVs.

## Reference:

- [1] Vehicles, Automated. "3.0: Preparing for the Future of Transportation." Federal Policy Framework. National Highway Transportation Safety Administration, US Department of Transportation (2018).
- [2] Chang, Ming-Fang, et al. "Argoverse: 3d tracking and forecasting with rich maps." Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. 2019.
- [3] Müller, Peter, and Fernando A. Quintana. "Nonparametric Bayesian data analysis." *Statistical science* (2004): 95-110.
- [4] W. Wang and D. Zhao, "Extracting Traffic Primitives Directly From Naturalistically Logged Data for Self-Driving Applications," in *IEEE Robotics and Automation Letters*, vol. 3, no. 2, pp. 1223-1229, April 2018.
- [5] Rasmussen, Carl Edward. "Gaussian processes in machine learning." *Summer School on Machine Learning*. Springer, Berlin, Heidelberg, 2003.
- [6] Teh, Yee Whye. "Dirichlet Process." (2010): 280-287.
- [7] Sun, Shiliang, and Xin Xu. "Variational inference for infinite mixtures of Gaussian processes with applications to traffic flow prediction." *IEEE Transactions on Intelligent Transportation Systems* 12.2 (2010): 466-475.
- [8] Gardner, Jacob, et al. "Gpytorch: Blackbox matrix-matrix gaussian process inference with gpu acceleration." *Advances in Neural Information Processing Systems*. 2018.
- [9] Titsias, Michalis. "Variational learning of inducing variables in sparse Gaussian processes." *Artificial Intelligence and Statistics*. 2009.