



# Probe-Based Volume Estimation Using Machine Learning Techniques

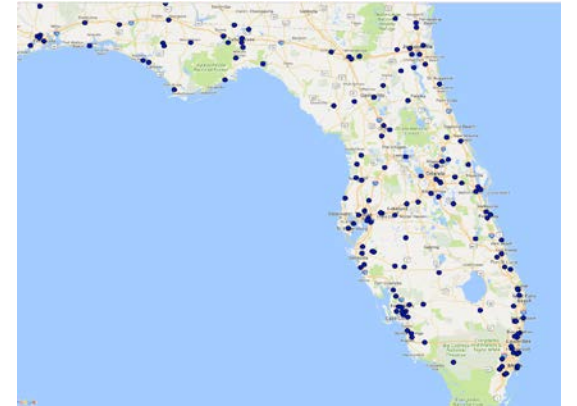
**Kaveh Farokhi Sadabadi, Przemyslaw Sekula, Zachary Vander Laan**  
**Center for Advanced Transportation Technology (CATT)**  
**University of Maryland**

Presented by:  
**Zachary Vander Laan**

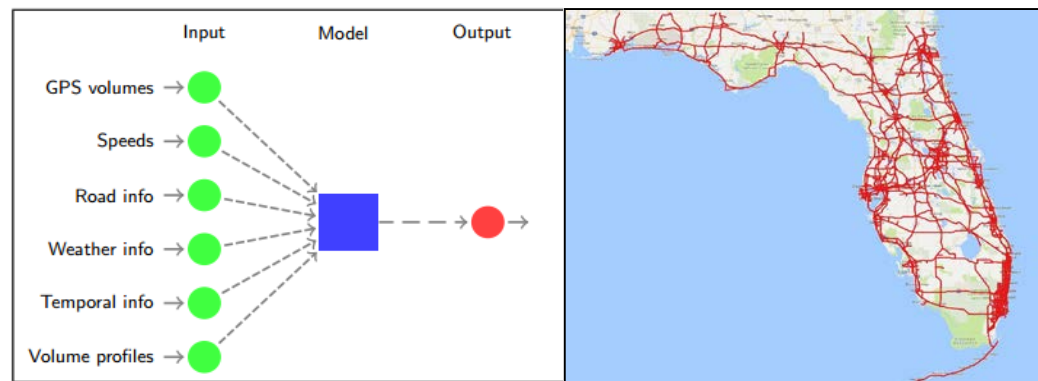
**2018 NaTMEC Conference**  
**Irvine, CA**

# Objectives

- Given the following:
  - Probe volumes (processed from GPS traces of a subset of vehicles),
  - Other archived data (speeds, road geometry, weather, etc.)
  - Continuous count data from select locations



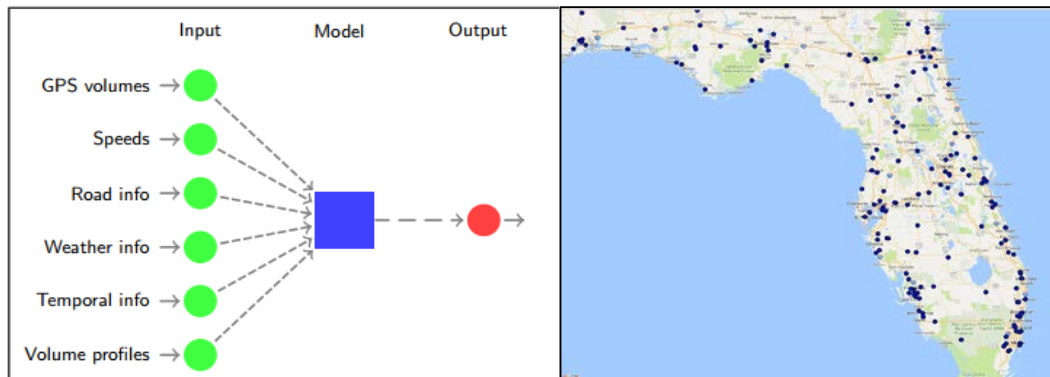
- Can we build a model to accurately estimate statewide volumes?



# Volume Estimation: General Approach

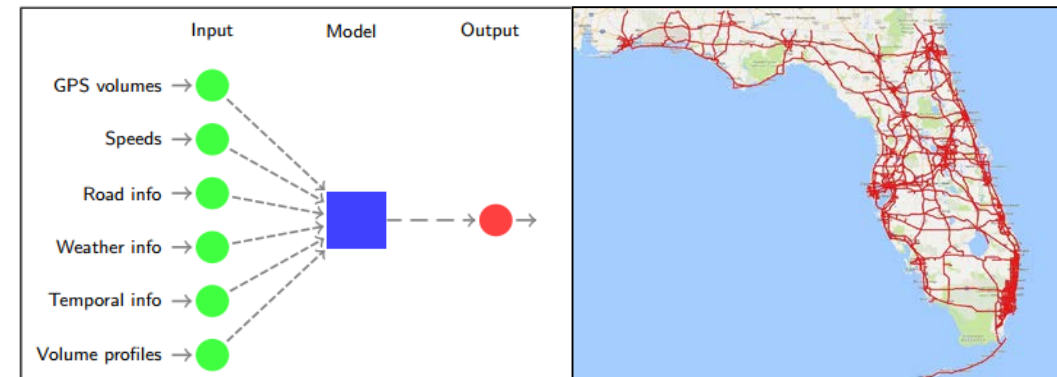
## Develop and Train Model

- Where? TMC segments associated with continuous count stations
- How? Construct machine learning model to learn relation between input variables and continuous count volumes



## Apply model to state road network

- Where? All TMCs on road network
- How? Apply trained model to input variables from any TMC segment on the network



# Data

## Data needed at all TMCs

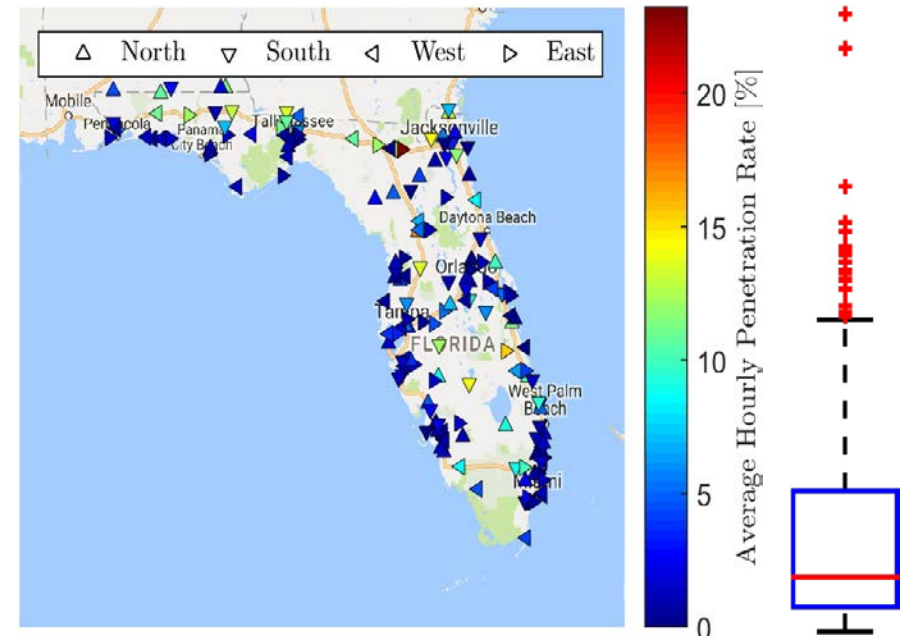
- **GPS probe data (INRIX)**

Dataset	Trips	Waypoints	Median Pen. Rate
Maryland	20 M	1.4 B	1.9%
Florida	75 M	3.4 B	2.1%
New Hampshire	7 M	595 M	2.3%

- **Probe speeds**
- **Road characteristics**
  - # lanes, speed limit, facility type, etc.
- **Weather**
- **TTI hourly volume estimates** (optional)

## Data needed only at continuous count stations

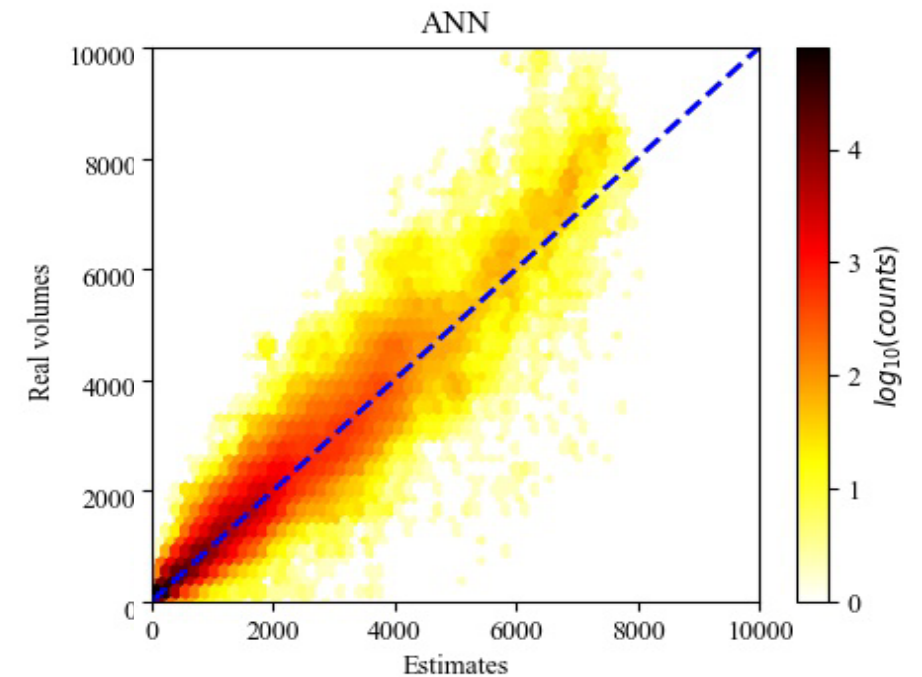
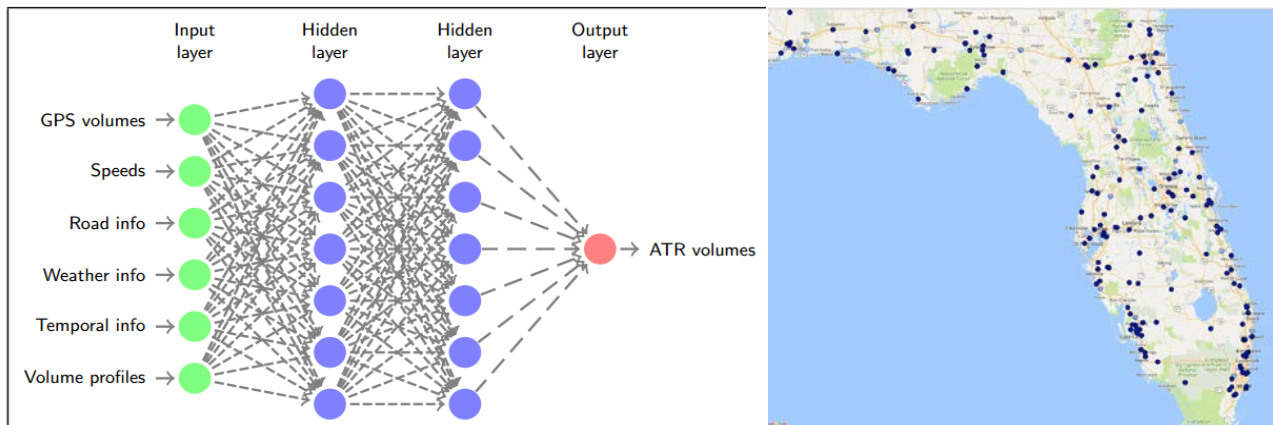
- **Ground truth count data**
  - Used for model training / evaluation
  - Used to estimate probe penetration rate



*Continuous count stations in Florida, with corresponding GPS probe penetration rate*

# Model Evaluation

- Model: “Dense” Artificial Neural Network (ANN)
- Cross validation : Repeat N times (N = number of stations)
  - Train model using data from all but one continuous count station
  - Generate model predictions using data from remaining station



- Evaluation: Compare estimates with actual volumes & generate metrics

# Quantifying Model Accuracy

$y_i$  = observed volume,  $\bar{y}_i$  = average observed volume,  $\hat{y}_i$  = model volume estimate,  $y_{\max}$  = max observed volume

- Mean Absolute Percentage Error (MAPE)

- Reflects absolute volume accuracy
- *Good: 10-15% (high volume),  
15-25% (mid volume)  
25-??% (low volume)*

$$MAPE = \left( \frac{1}{n} \sum_{i=1}^n \left| \frac{\hat{y}_i - y_i}{y_i} \right| \right) \times 100$$

- Error to Max Flow Ratio (EMFR)

- Captures accuracy relative to capacity (max observed flow)
- *< 10% becomes useful, < 5% target*

$$EMFR = \left( \frac{1}{n} \sum_{i=1}^n \left| \frac{\hat{y}_i - y_i}{y_{\max}} \right| \right) \times 100$$

- Coefficient of Determination (R<sup>2</sup>)

- Shows explanatory power of model
- *> 0.70 good, > 0.80 better, > 0.90 best*

$$R^2 = 1 - \frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{\sum_{i=1}^n (y_i - \bar{y})^2}$$

# Results: Overview

## Summary

- Promising model performance, even across multiple scenarios
- Stable accuracy levels across multiple datasets

## Observations

- ↑ Road class = ↑ Accuracy
- ↑ Avg. hourly volume = ↑ Accuracy
- ↑ Avg. hourly GPS counts = ↑ Accuracy

**Median Error Metrics by Dataset**

Dataset	R2	MAPE (%)	EMFR (%)
Maryland	0.85	22.6	6.6
Florida	0.83	24.8	6.6
New Hampshire	0.82	27.6	7.3

# Flagging Unusual Behavior

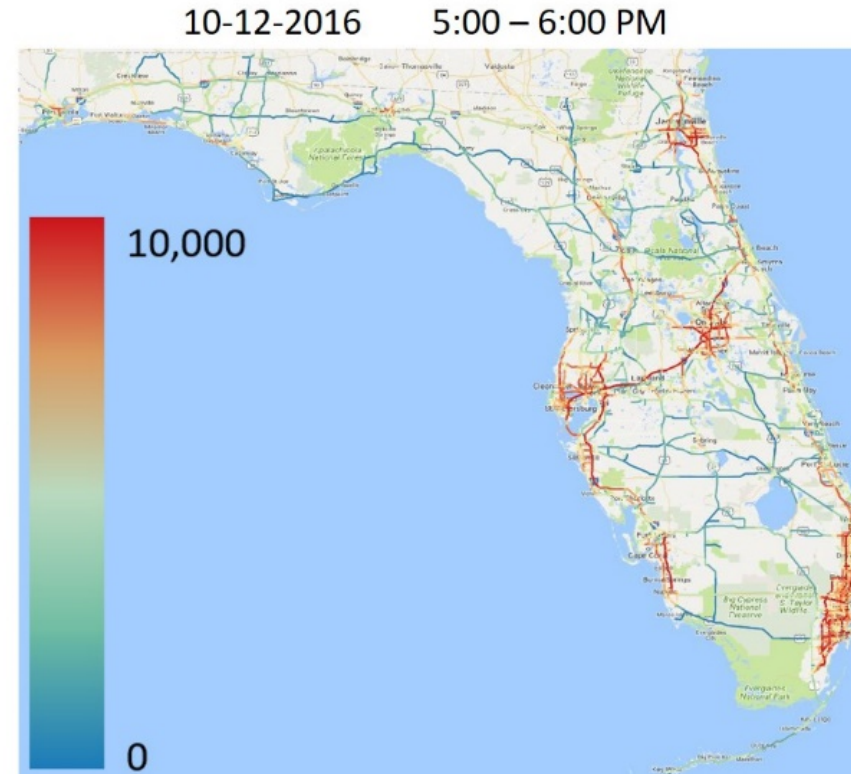
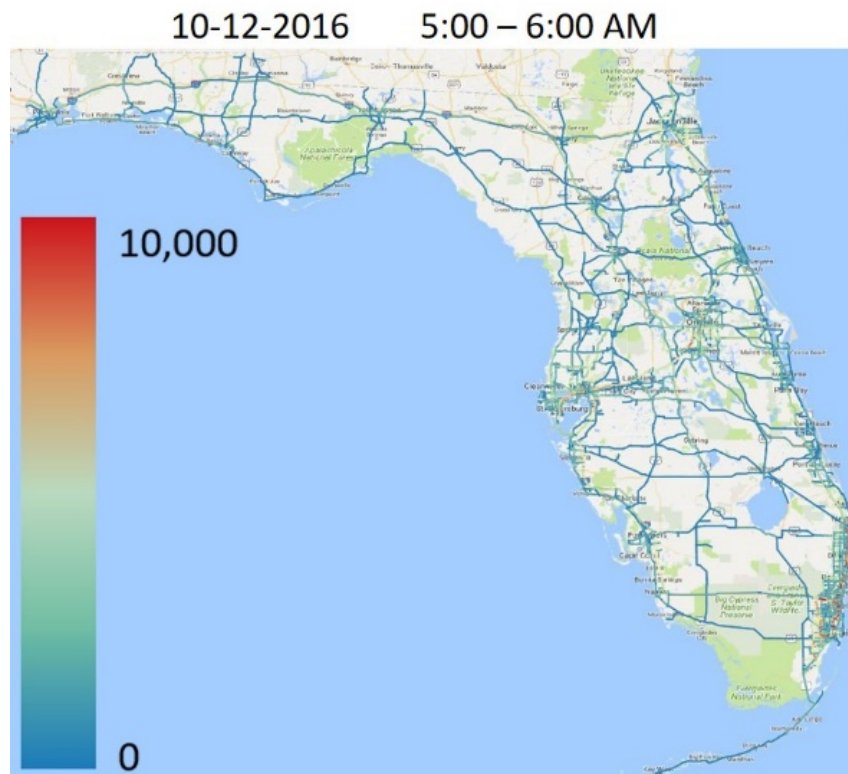
**Goal:** Develop flags to highlight unusual input data and output model estimates

- **Flag 1** - based on GPS input data (key model “ingredient”)
  - *Typical:* Observed GPS counts within X std. dev of mean GPS counts during same day of week and hour
  - *Low:* Less than *Typical* range
  - *High:* Greater than *Typical* range
- **Flag 2** - based on output model estimates
  - *Typical:* Observed hourly estimates within X std. dev of mean estimates during same day of week and hour
  - *Low:* Less than *Typical* range
  - *High:* Greater than *Typical* range



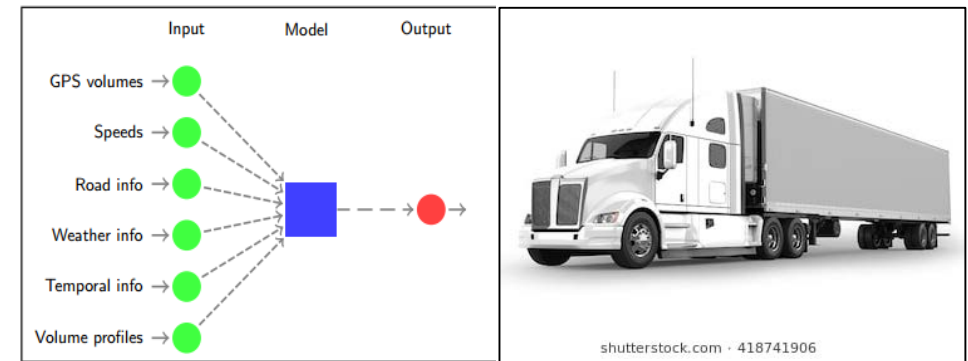
# Statewide Model

- Apply trained model to entire road network in Florida
  - Requires 3 months of hourly input data at ~20k TMCs
  - Generate hourly volume estimates at each input time/location



# Hourly Freight Volume Estimation

- Use same approach to estimate hourly freight volumes
  - Train model using only ground truth truck counts
  - Need continuous count data by FHWA weight class!
- Initial Florida freight results look promising on high FRC
  - FRC 1 results comparable to all-vehicle model accuracy
  - Not enough data to reliably estimate low FRC



**Median Error Metrics: Florida Truck Volume Estimation**

FHWA Class 5-13	R <sup>2</sup>	MAPE (%)	EMFR (%)
<b>Overall</b>	0.77	37.9	7.5
FRC 1	0.83	23.5	6.3
FRC 2	0.76	42.2	7.9
FRC 3 & 4	0.65	48.9	9.2

*Similar accuracy to all-vehicle Florida model*

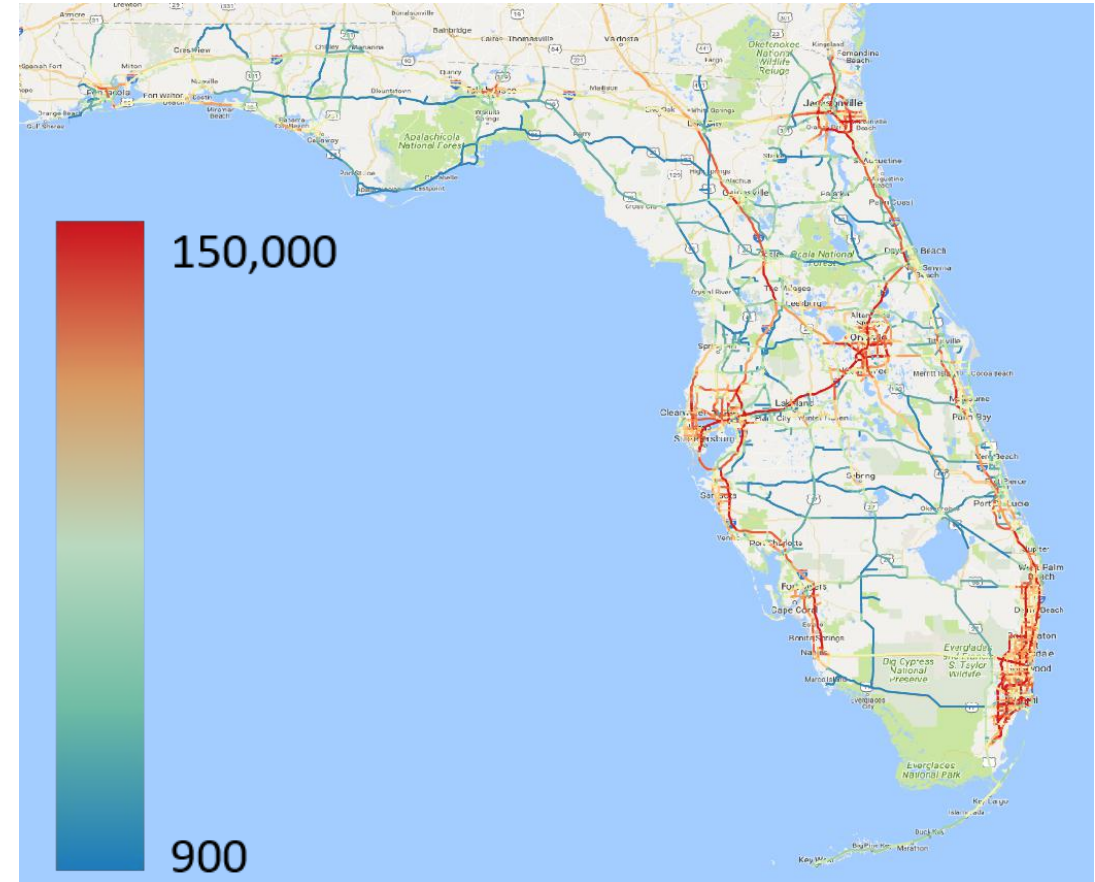


# AADT Estimation

- Possible Approaches:
  - Method 1: Aggregate hourly volume estimates
  - Method 2: Develop separate AADT estimation model
- Generated initial results via Method 1
  - Promising model performance!
  - Consistent with expectations along major highways and urban areas

R <sup>2</sup>	MAPE (%)
0.86	15

- Future work will compare approaches



# Next Steps

## Technical

- Refine hourly volume estimation models
- Scale approach to statewide networks in MD and NH
- Further investigate hourly truck volume and AADT estimation
- Explore transferability of models between different states

## Program Level

→ Get out of the lab and operationalize!

# Questions

## Contact Information

Kaveh Farokhi Sadabadi (PI)

[kfarokhi@umd.edu](mailto:kfarokhi@umd.edu)

Przemyslaw Sekula

[psekula@umd.edu](mailto:psekula@umd.edu)

Zachary Vander Laan

[zvanderl@umd.edu](mailto:zvanderl@umd.edu)

## Acknowledgments

- I-95 Corridor Coalition
- MDOT, FDOT, NHDOT