# Real-time Processing of Noisy Data from Connected Vehicles
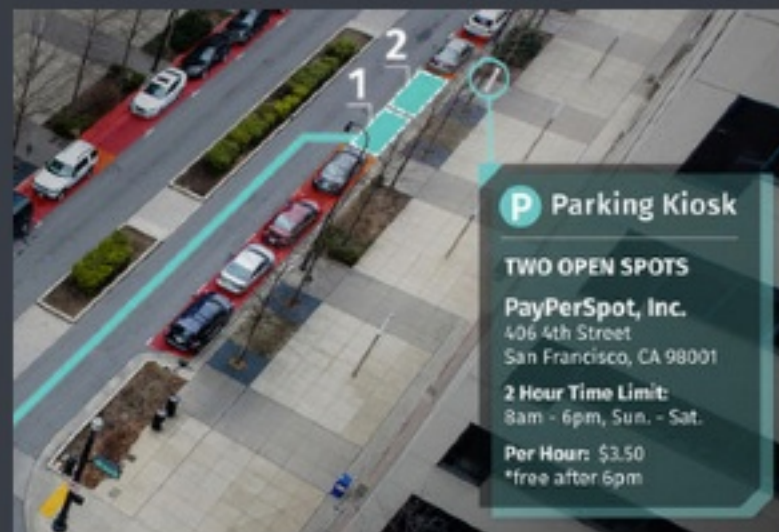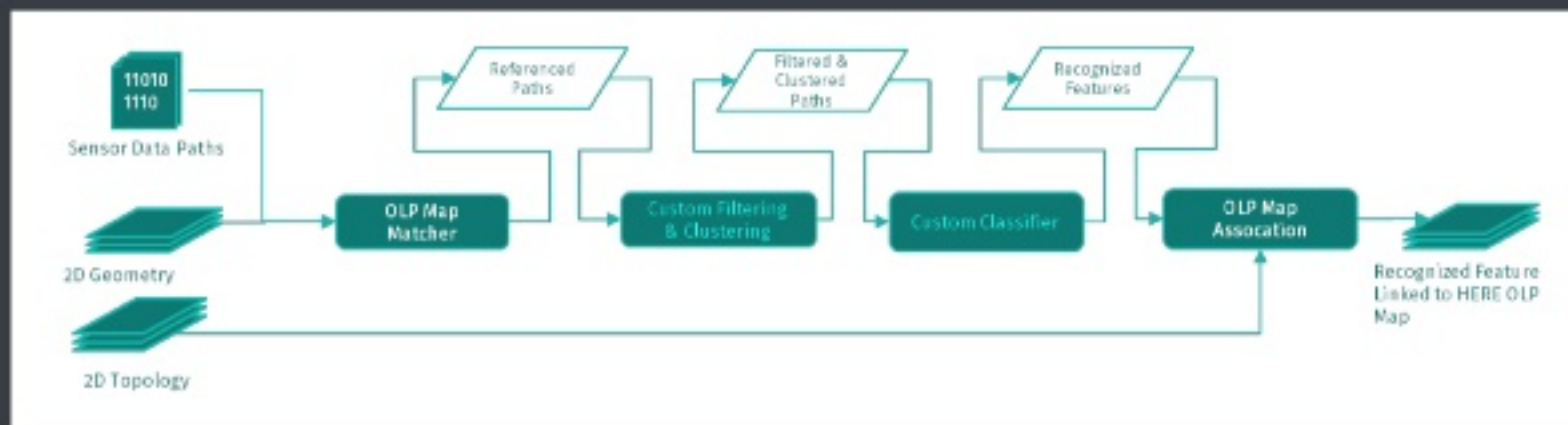
Robin * Slomkowski
Flink Forward 2018-09

# Making Safer Roads
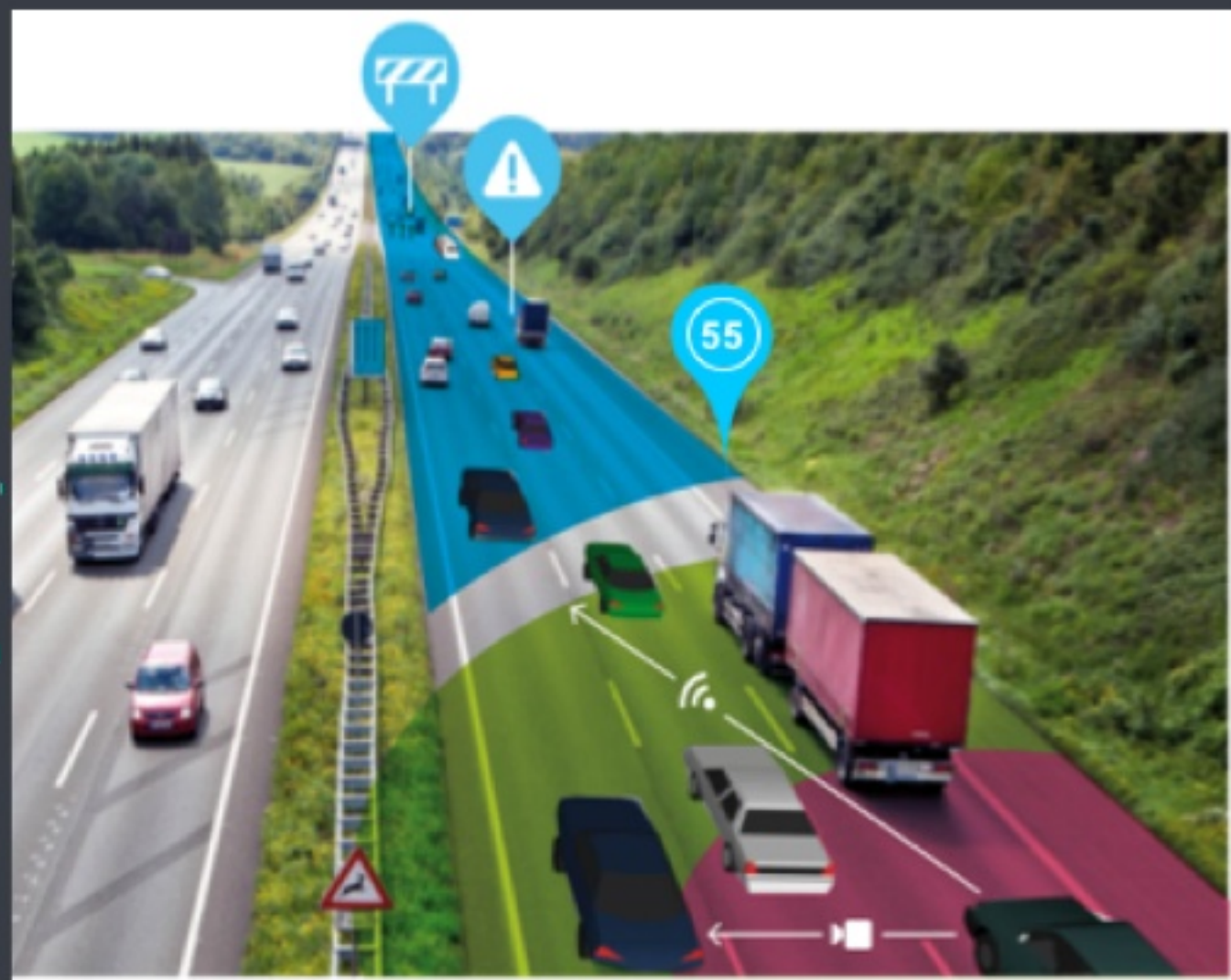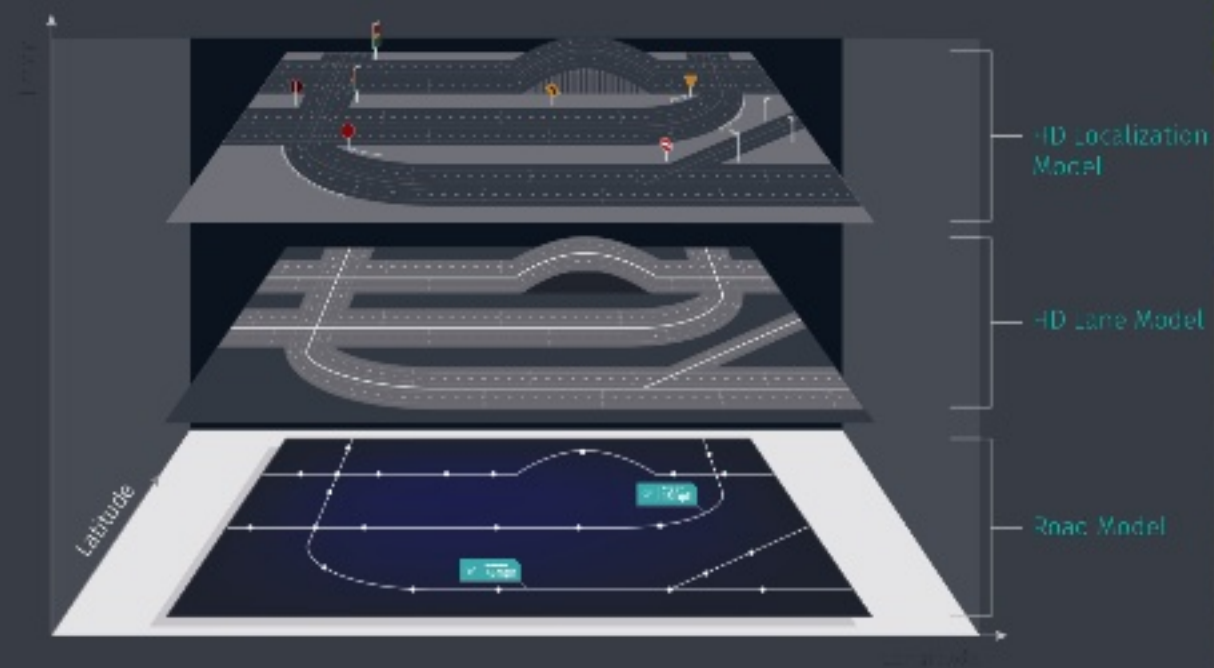
Slippery roads

Variable Signs

Accidents and other Events
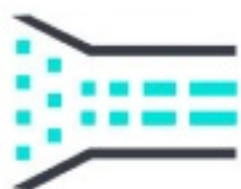
Live Parking Data

# Building Living Maps

Annotating Sensors against an authoritative reference.

# HERE Open Location Platform

## Rich capabilities to create intelligent location-centric products and services from your data

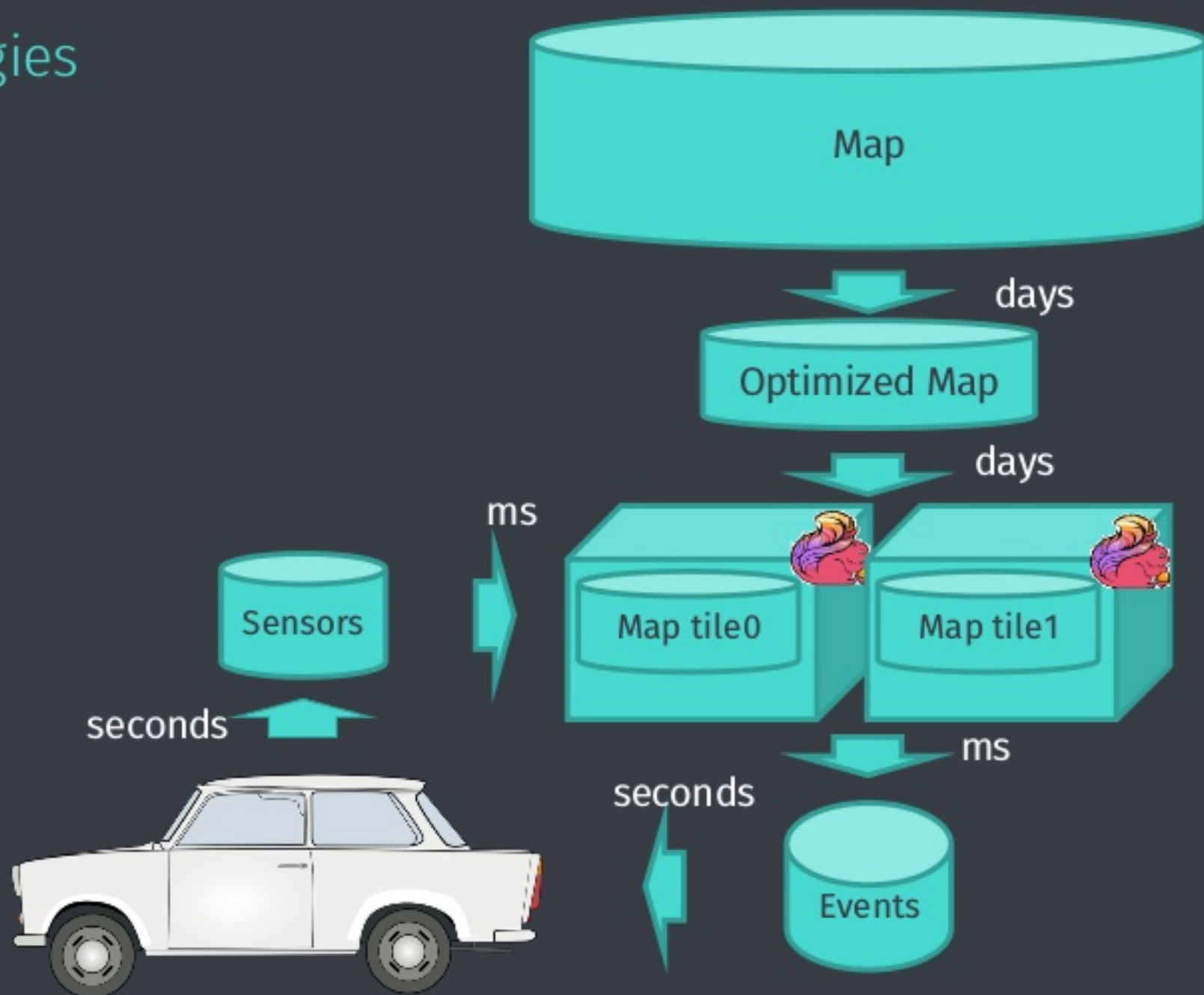| **Explore and ingest** | **Transform and enrich** | **Analyze and model** | **Scale and distribute** |
|---|---|---|---|
| Powerful tools to access HERE data and bring your own data into OLP. | Enrich your data with location context for additional value. | Analyze data to derive insights and create machine learning models. | Publish your enriched datasets or intelligence into production at scale. |

*Focused on*

http://openlocation.here.com

# Stream Processing Strategies

**Goals**

- Lower Latency
- Less Engineering

**Techniques**

- Pre-Computing
- Deterministic Data Access
- Extensible Infrastructure Optimized Functions

Map

days

Optimized Map

days

ms

Sensors

Map tile0

Map tile1

seconds

seconds

ms

Events

# Managing Pipelines on the HERE's Platform
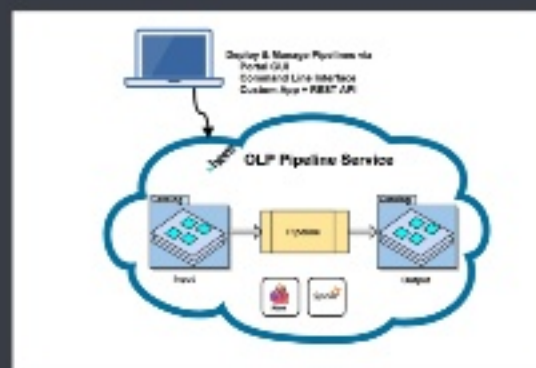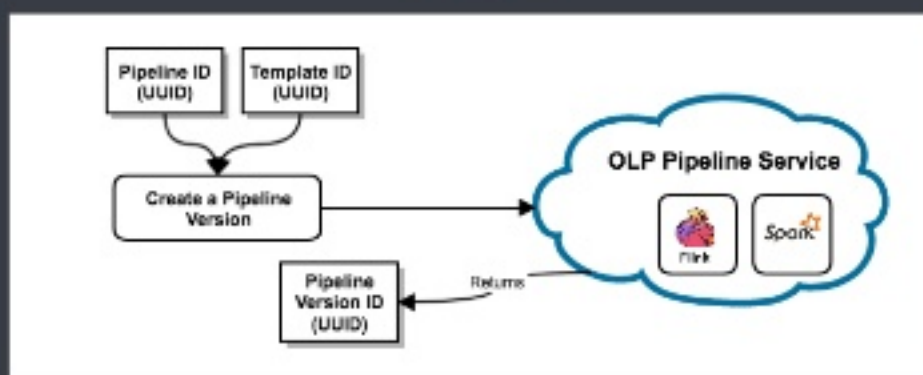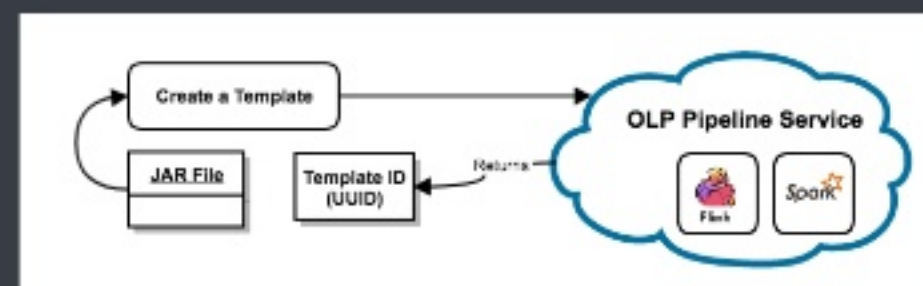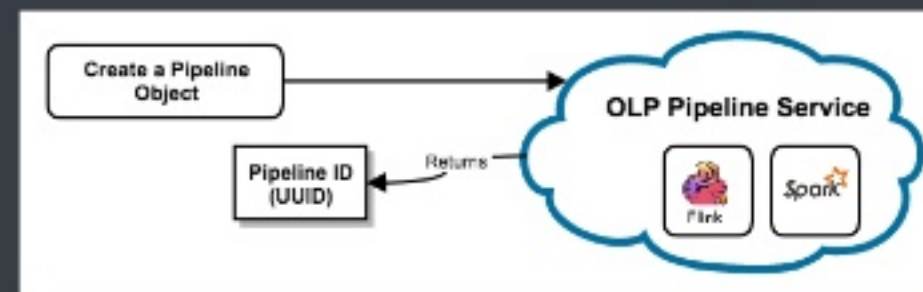
1. **Prepare**
   1. Create a Pipeline (Framework)
   2. Create a Pipeline Template (Jar)
   3. Create a Pipeline Version (Config)

2. **Deploy**
   1. Activate (Pipeline Version)

3. **Run**
   1. Cancel, Pause, Resume, Upgrade

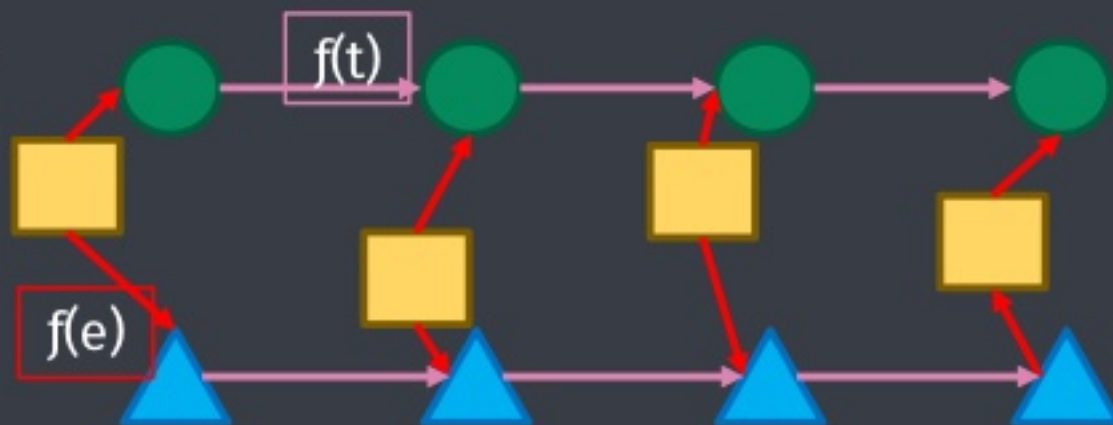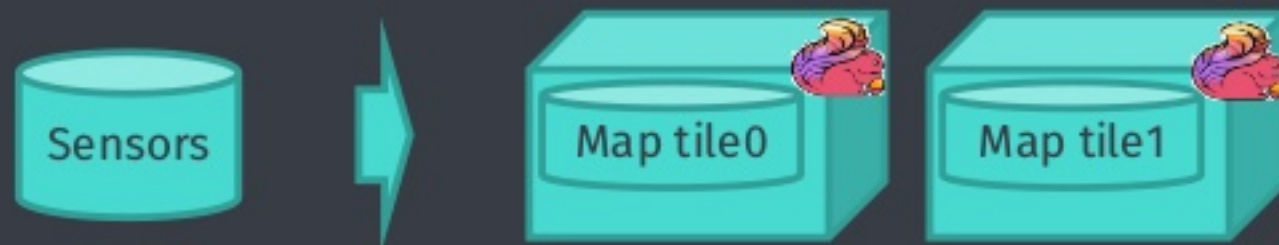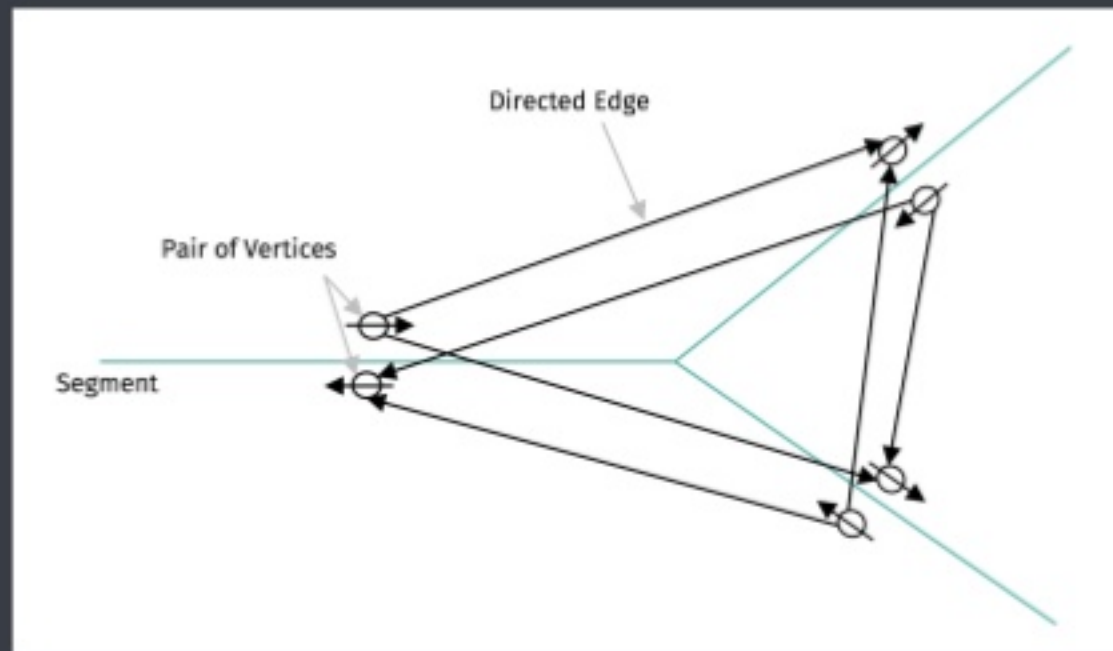# Implementation of Strategy

- ## Optimize Data for Use Case
  - Don't need all data for sensor attribution
  - We know how the data will be used

- ## Key-By Map Tiling
  - Location Data does not arrive randomly, make caches work

- ## Extensible Map-Matching Implementation
  - We can use stateless evaluation to choose the right result



Directed Edge

Pair of Vertices

Segment

Sensors

Map tile0    Map tile1

f(t)

f(e)

# Map Matching Code

Data Client gets the Optimized Map

Data Client gets the Stream for Flink

Parse SDII proto

Default Path Map Matcher

# What Kind of Data is SDII?

## Series of Timestamped KV Pairs

- A series of GPS values with time allows identifying a point better
- A collection of other sensors, slope, wheels, signs, etc... with time

UTC Time:    2014-03-25 15:37:46.500Z
Type:        MAP_MATCHED_HD_MAP
Lon / Lat / Alt: ...

UTC Time:    2014-03-25 15:37:47.200Z
Type:        FILTERED
Lon / Lat / Alt: ...

UTC Time:    2014-03-25 15:37:46.334Z
Type:        RAW_GPS
Lon / Lat / Alt: 9.1234567/48.1234567/123.456
Speed/Heading/Accuracies: ...

## GPS Point to Sensor Estimation is the User's Problem

- DBSCAN clustering is supported by HERE, but custom logic may be added.

Position Estimates

UTC Time: 2014-03-25 15:37:46.500Z

UTC Time: 2014-03-25 15:37:47.000Z

Slope: 1.1%

Slope: 1.5%

Slope: 1.4%

Slope: 1.3%

Slope: 1.3%

Slope: 1.2%

Wheels Slipping

Sign Recognized

here

# Optimized Map

Java Integer Arrays to avoid scala boxing/unboxing

Add attributed Index for constant time access to in memory data structures

## Why so much effort?

- These data must be scanned thousands of times to build emission and transition probabilities.
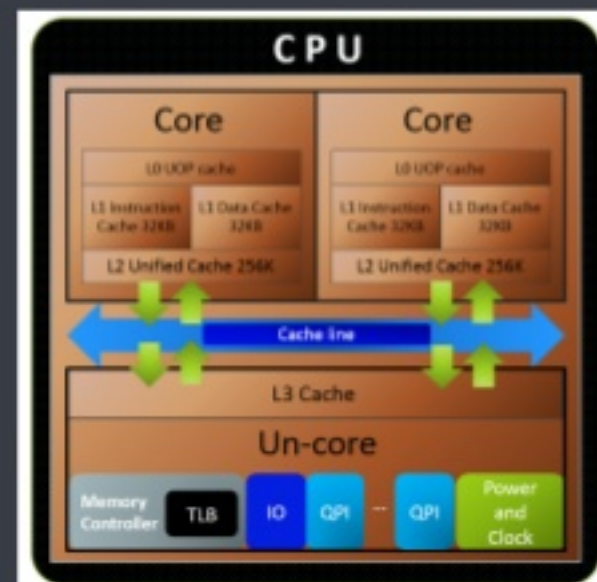- Series of points are near each other so cache hits work well
- We want stream processing light as possible, do this work 1x
- Multiple Orders of magnitude faster than the naïve approach
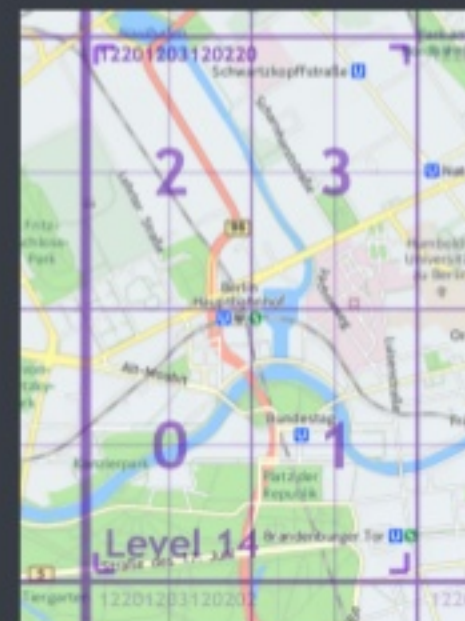
9/4/18
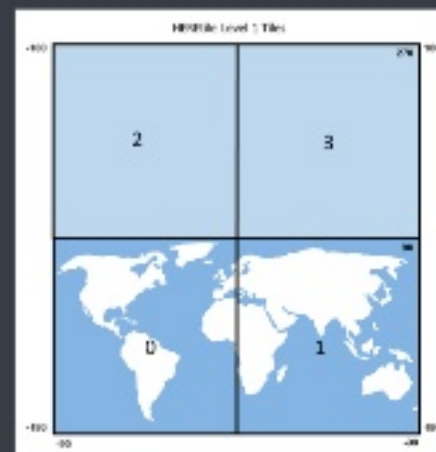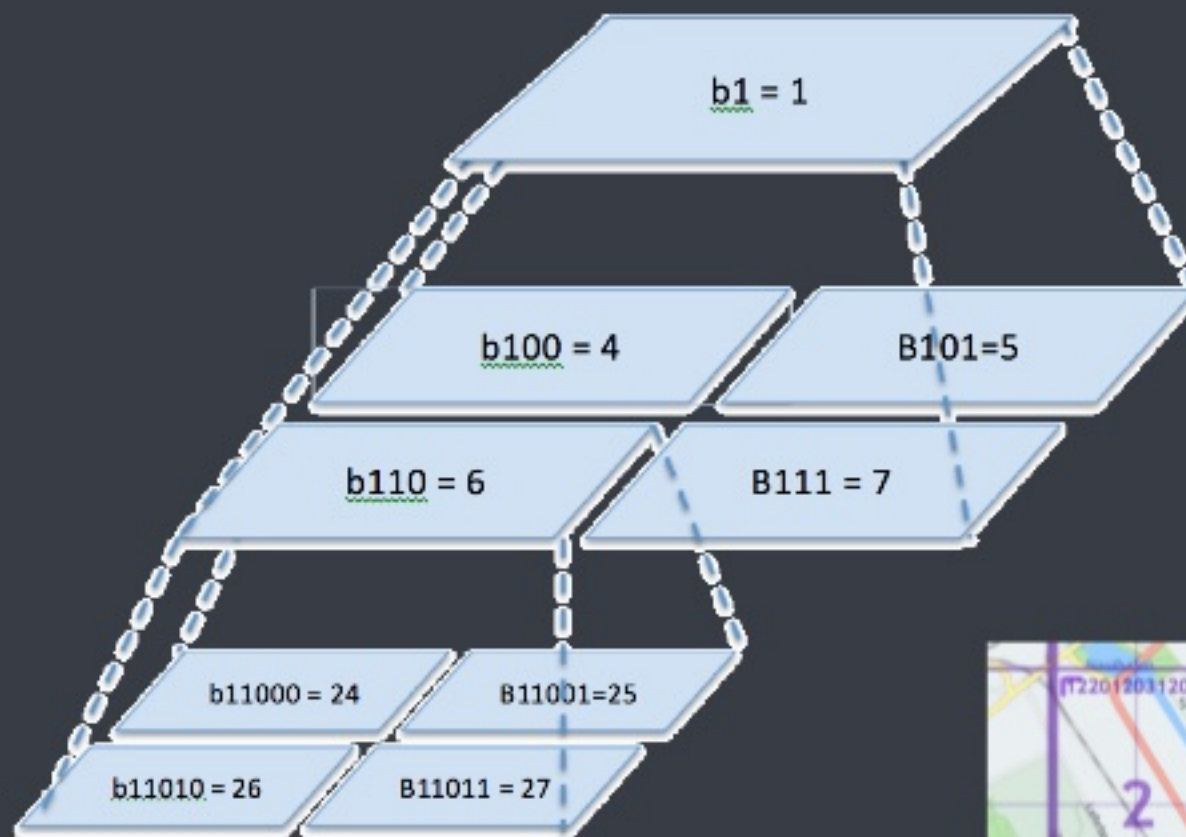
# What is Map Tiling?

## HERE Tile

## Strategies for categorizing data
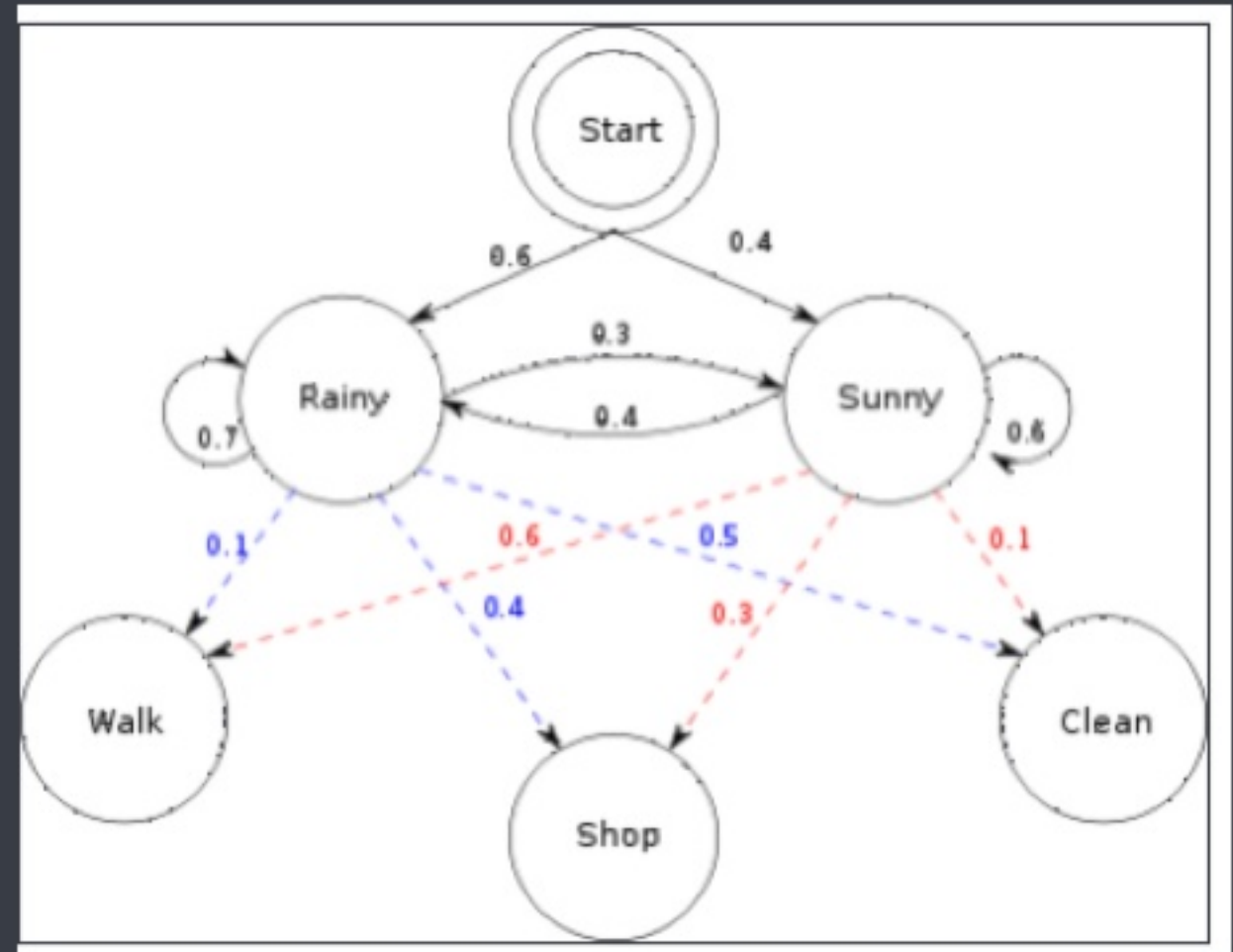- Partitioning at arbitrary depth

## Why Key-By makes sense
- People do not jump randomly between streets or around the world, so you tend to get multiple results near by each other.
- Enhances cache hits immensely
- Yes we talk about origin so given envelope may cross tiles, but the caching works better this way anyway.
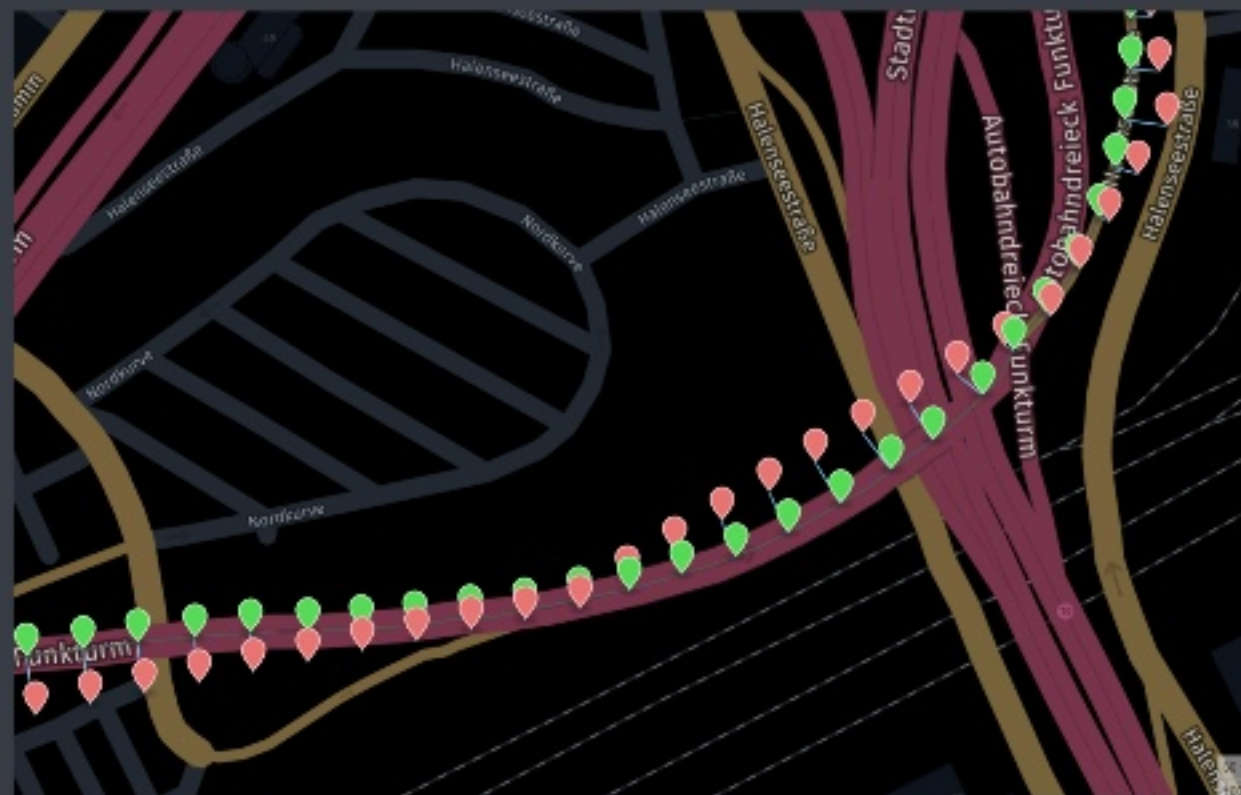
# Map Matching: Maintaining Algorithms

- Hidden Markov Model
  - Yes this is the best practice solution.
  - Optimizing the data needed to make the decisions is where you get the performance from

- Data Model Agnostic
  - Can be applied to any Vector Data

- Customizable
  - Emission Probability
  - Transition Probability

# Map Matching: Paths and Points

- Why are points together?

- Optimizing matching based on logical paths?

- What about anomalies?

- What about emission and transition logic?

# Map Matching: Emission Probability

- Points don't line up

- Shortest path doesn't match vector

- One algorithm to select the emission

- A second to select the route, to define the point.

Route 202
0.51 * 0.51 * 0.51 * 0.3 * 0.1  => **0.40 %**

West Swedesford Rd.
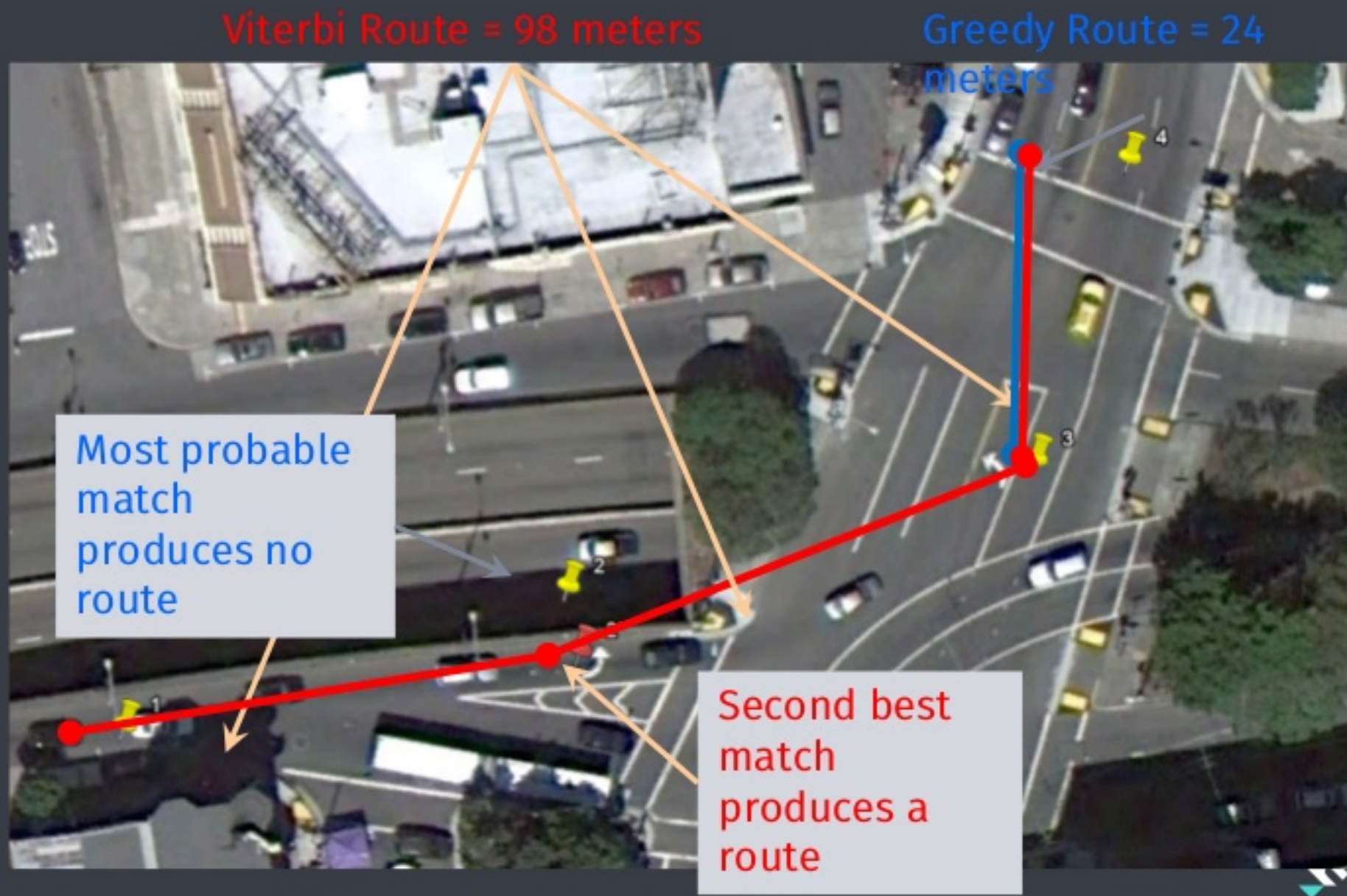0.49 * 0.49 * 0.49 * 0.7 *0.9 => **7.40 %**

# Map Matching: Transition Probability

- Viterbi as opposed to Greedy Match

- No Right Answer Every Time Bring your own algorithm

Viterbi Route = 98 meters

Greedy Route = 24 meters

Most probable match produces no route

Second best match produces a route

Thank you
http://openlocation.here.com

Contact