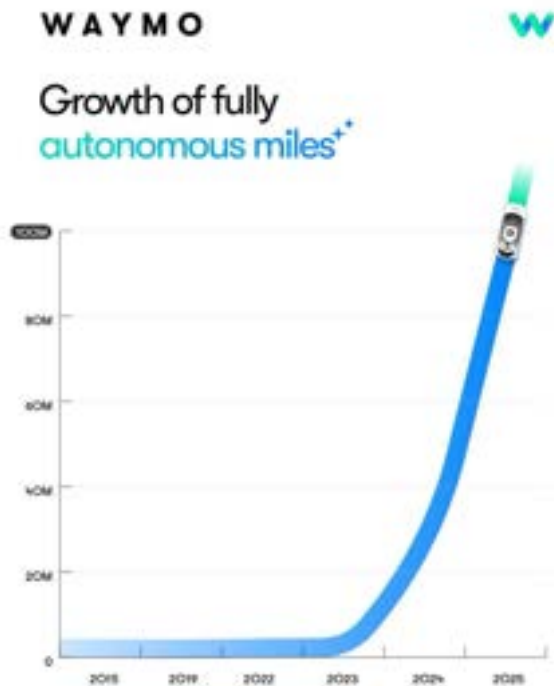




Time-Series Processing at Waymo

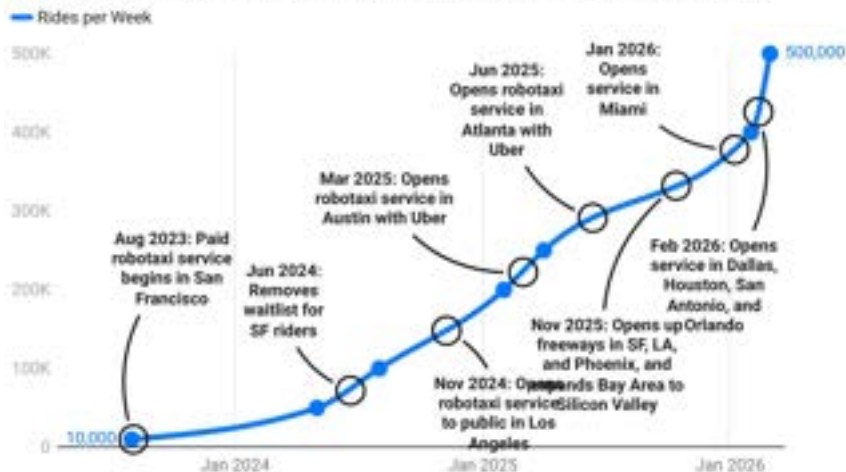
Joey Raso
Software Engineer @ Waymo

Waymo's Rapid Scaling



Waymo's Quick Ascent to Half a Million Weekly Robotaxi Rides

The company's paid autonomous vehicle rides per week accelerated quickly, as it opened up service to new cities over the last year. Waymo now operates in 10 cities across the U.S.



Hover over the dots for ridership milestones.

Source: TechCrunch reporting, Waymo public statements • Created with Datawrapper

Bringing Waymo to more places



- ★ Waymo operating cities
- ◆ Operation coming soon

Early Event Discovery

As Waymo scales, we move further and further into the long-tail in terms of unique road conditions. This means that it's more important than ever that we're collecting interesting road events for further analysis.

Further, we'd like this collection to happen as quickly as possible to establish confidence as we roll out new software to the fleet.

FALLING TREES



FLOOD



FIRST RESPONDERS



PEOPLE EXITING MOVING CARS



BOUNCE HOUSE



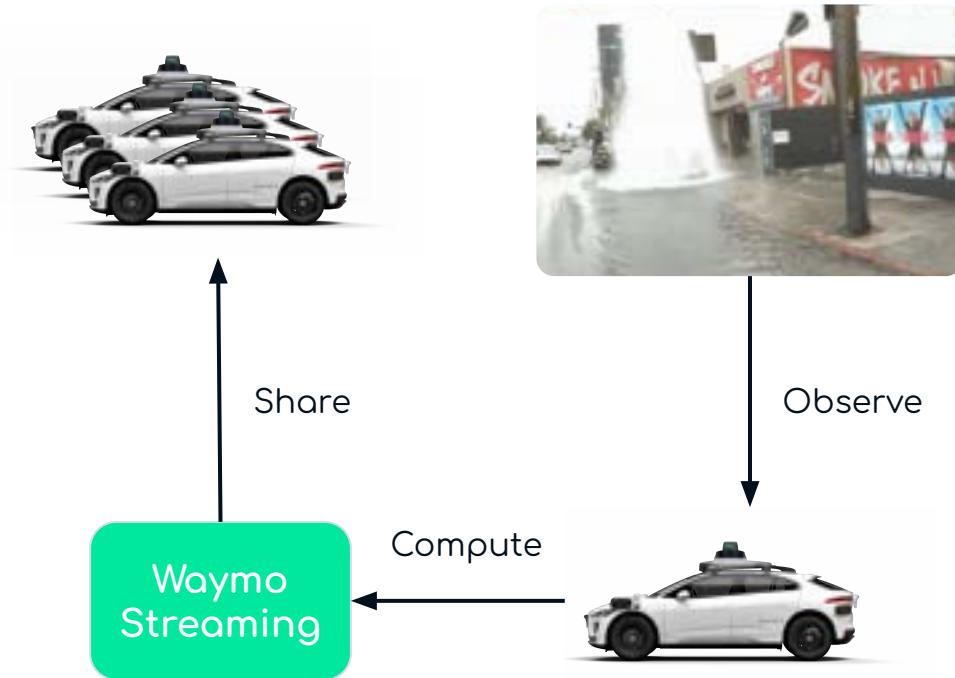
ROAD DEBRIS



Fleet Information Sharing

A higher concentration of vehicles on the road means that there are likely to be road features that we are traversing multiple times per day across the fleet.

By propagating and sharing information amongst the fleet, we can approach tricky situations with a stronger prior and enhance the Waymo Driver's ability to respond.



Importance Sampling

While Waymo has access to a massive amount of data storage, as we've solved the challenges associated with autonomous driving, we try to be efficient with our data storage usage and are seeing opportunities for further efficiencies as we scale.

As we scale and drive millions more miles in more places, we need to build systems to refine and grow our datasets to reflect the world we're driving today and tackle the remaining long-tail challenges while ensuring an efficient approach to data storage.

 Mashable

Thanks a lot, AI: Hard drives are already sold out for the entire year, says Western Digital

Western Digital says its all sold out of hard drives for 2026, less than two months into the year.



 Tom's Hardware

Hard drives on backorder for two years as AI data centers trigger HDD shortage — delays forcing rapid transition to QLC SSDs

The AI boom might help QLC overtake TLC in the next two years.



 Data Center Dynamics

Storage wars: Is this the end for hard drives in the data center?

Could the shortages facing the HDD market ultimately lead to the technology being replaced?



To enable these Early Detection, Observability, and Importance Sampling use cases, we need a capable Streaming Data Processing Framework to **process the Waymo driving data, make that data available to our engineers, and, when applicable, propagate signal to the fleet.**

System Requirements

Fleet Scale: Our Streaming systems should be able to handle the data load of Waymo's rapidly expanding fleet.

Late-Arriving / Out-Of-Order Data: The fleet data is subject to network latencies and other delays. Whatever system we build needs to be able to handle these delays gracefully.

Low Latency: To get the most out of our data, we need to get it in front of our engineers as quickly as possible. Therefore, the system must be able to process at low latency.

Stateful Processing: Much of the signal that we need to generate relies on timelines / windows of data rather than point observations. The framework should have primitives for stateful processing.

Hypothetical: Red Light Runners



Hypothetical: Red Light Runners

```
perceptions = (  
  p  
  | "ReadStream" >> beam.io.ReadFromPubSub(  
    subscription=input_subscription)  
)
```

Consume Data Streams from
Waymo Fleet

```
red_light_runners = (  
  perceptions  
  | "WindowIntoFixed" >> beam.WindowInto(window.FixedWindows(5.0))  
  | "KeyByObjectId" >> beam.Map(lambda obj: (obj.object_id, obj))  
  | "GroupTrajectories" >> beam.GroupByKey()  
  | "CalculateKinematics" >> beam.ParDo(DetectRedLightRunnersFn())  
)
```

Aggregate and perform
Kinematics Calculation

```
red_light_runners  
  | "EmitSignals" >> beam.Map(print)  
)
```

Emit to your Sink of choice

Hypothetical: Red Light Runners

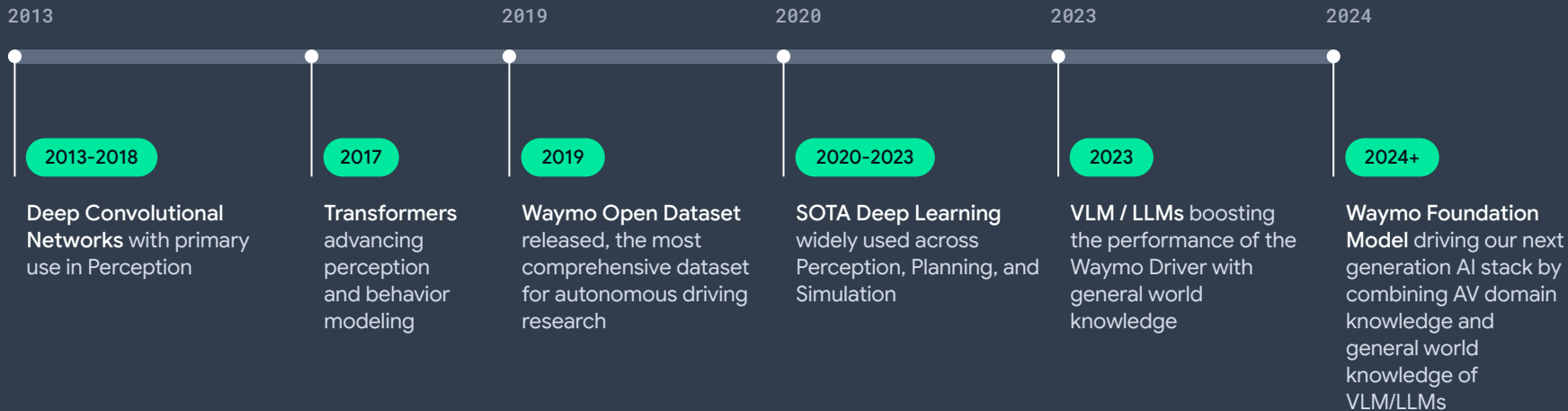
With our new Dataflow-powered Red light runner detector, we could:

- **Discover red light runner events** in the logs to generate test sets.
- Mark log segments containing red light runners for **extended retention**.
- **Share these events event** back to the fleet to inform the Waymo Driver of potentially dangerous intersections.

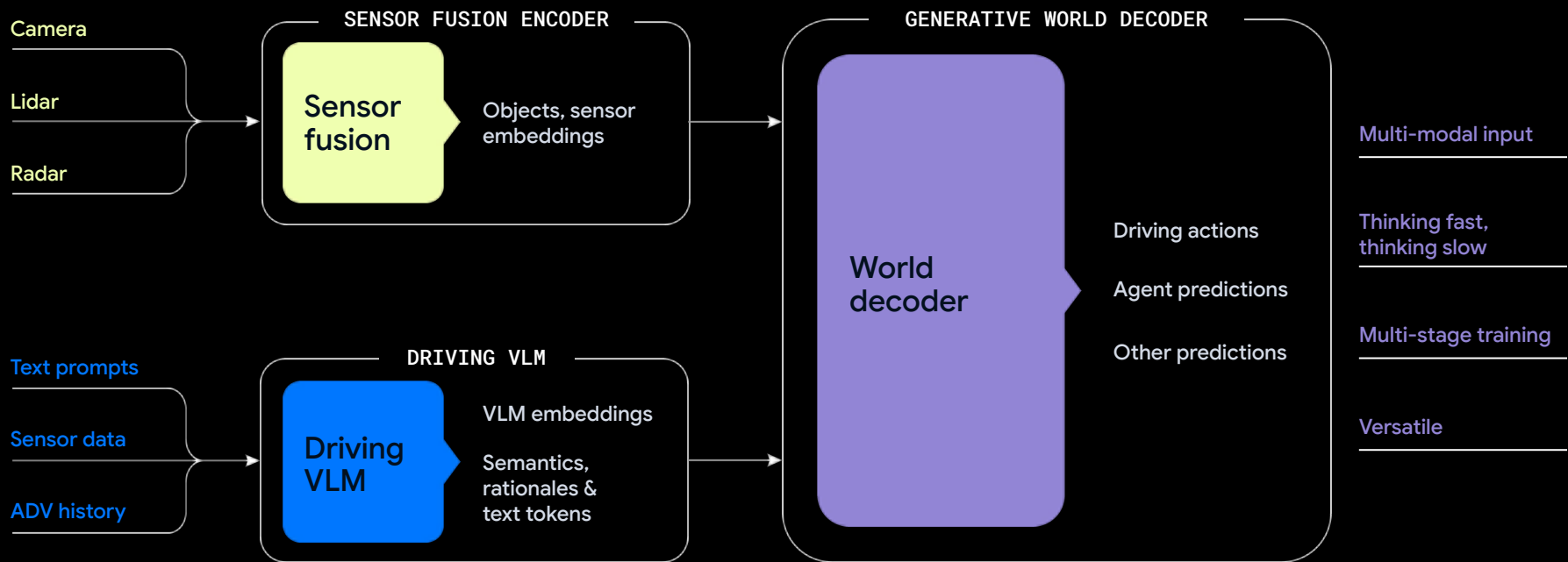
LLM-powered Signal Generation

We've seen how we can produce low-latency signal to power our Fleet Information Sharing, Event Discovery, and Importance Sampling use cases. As is the case with our Red Light Runner detector, our signals have traditionally been heuristics-based. Recent industry-wide AI developments, combined with Waymo's AI expertise, have opened up the door for AI / LLM-Powered Signal Generation.

Waymo is a pioneer in AI for autonomous driving



The Waymo Foundation Model



LLM-powered Signal Generation



Forward sensor + camera data over the cell network to Waymo's compute fleet.



The Waymo Streaming System performs windowed aggregations to generate Waymo FM features.

Waymo FM generates World Embeddings to as features in our importance classification.

Hypothetical: Red Light Runners

```
perceptions = (  
  p  
  | "ReadStream" >> beam.io.ReadFromPubSub(  
    subscription=input_subscription)  
)
```

Consume Data Streams from
Waymo Fleet

```
red_light_runners = (  
  perceptions  
  | "WindowIntoFixed" >> beam.WindowInto(window.FixedWindows(5.0))  
  | "KeyByObjectId" >> beam.Map(lambda obj: (obj.object_id, obj))  
  | "GroupTrajectories" >> beam.GroupByKey()  
  | "RedLightRunnerDetector" >> beam.ParDo(RunnerDetector())  
)
```

Aggregate and perform
Inference

```
red_light_runners  
  | "EmitSignals" >> beam.Map(print)  
)
```

Emit to your Sink of choice

Hypothetical: LLM Red Light Runners

- Composable APIs allow us to drop in an LLM Inference with minimal disruption to the rest of the pipeline.
- ML-first APIs enhance our heuristics-based approach with LLM Classification. With this API we can easily hook into Waymo's fleet of GPUs and TPUs.
- A combination of Batch and Streaming Execution modes allow us to easily backtest our model improvements.

QUESTIONS?