



# Ensuring Quality Data for Deep Learning in Varied Application Domains: Data Collection, Curation and Annotation

Gaurav Singh

System Architect and Perception Lead

Nemo @ Ridecell

# Talk outline

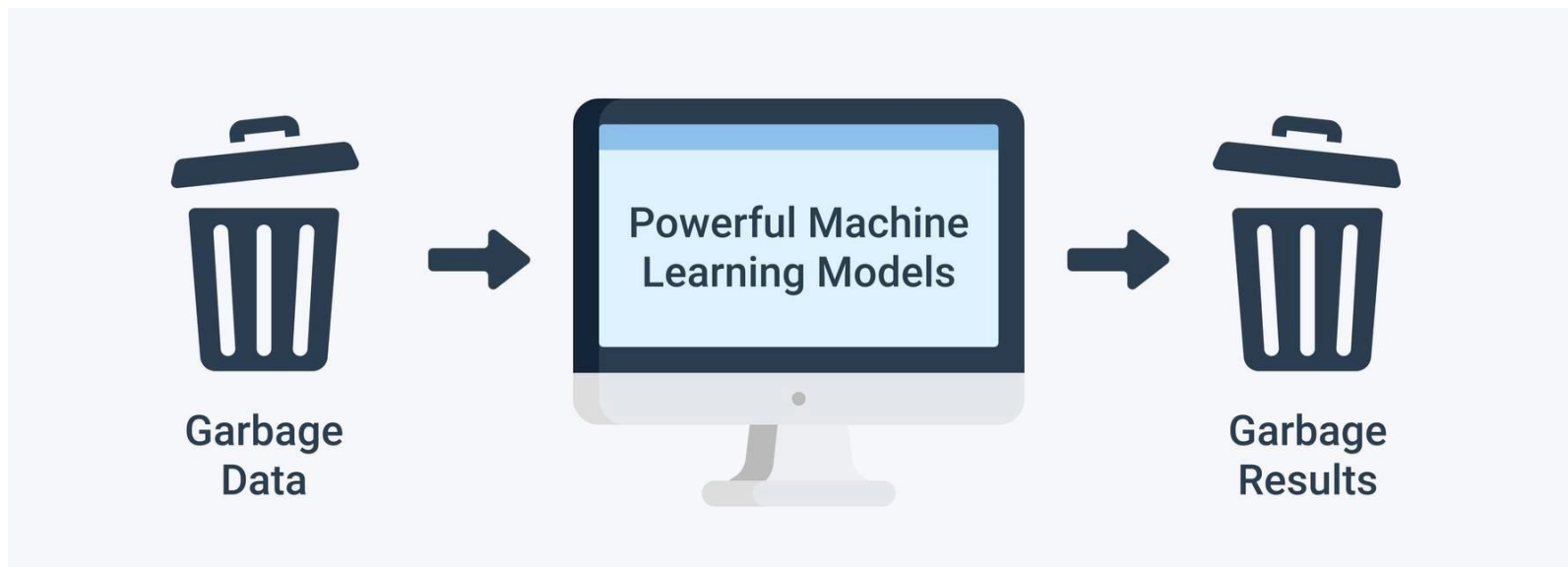


- Introduction and motivation
- Data collection
- Data curation
- Data annotation
- Conclusion

# Introduction and motivation



Collection, curation and annotation key to good deep learning











Stage	Questions to address
Collection	Collect/buy, post-production data collection, synthetic data?
Curation	Active learning, tagging, curation tooling
Annotation	Instructions, auditing techniques, in-house vs managed service

# Data collection

# Data collection - datasets for deep learning



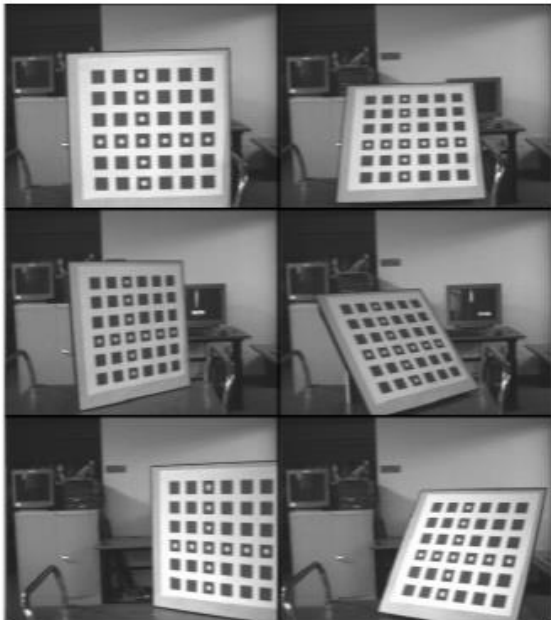
	<i>Open and free datasets</i>	<i>Paid Licensed datasets</i>	<i>Own datasets</i>	<i>Synthetic datasets</i>
<b>Description</b>	Free datasets under Creative Commons Attribution 4.0 licenses or similar	Paid licenses. Sold by companies for commercial gain	Collected on your own or with partners	Generated for specific use cases / corner cases by simulation companies.
<b>Examples</b>	  	  	Collected by your own company and annotated.	 
<b>Pros</b>	Lots of research interest - publications and networks built on top of these datasets. Labeling done usually. Government datasets usually here. Can be relabeled with required class.	Some research interest from universities etc. Generally higher quality labels. Some filtering to pick interesting sequences.	Can be heavily filtered for use case, full flexibility in labeling. Quality under company collecting and using data.	Generated for specific corner cases where algorithms are not working well. Easily annotated.
<b>Cons</b>	Can be limited size, labeling accuracy can be suspect.	Cost, Can't be re-labelled.	Costly to collect and annotate	Reality - simulation gap

# Data collection – key challenges for sensor data



Sensor data (Lidar, Camera, GPS etc.) collection challenges

- Sensor rig setup
- Sensor rig calibration
- Time syncing sensors
- Maintenance of rig – recalibration etc.



# Data collection – key challenges contd.



More challenges –

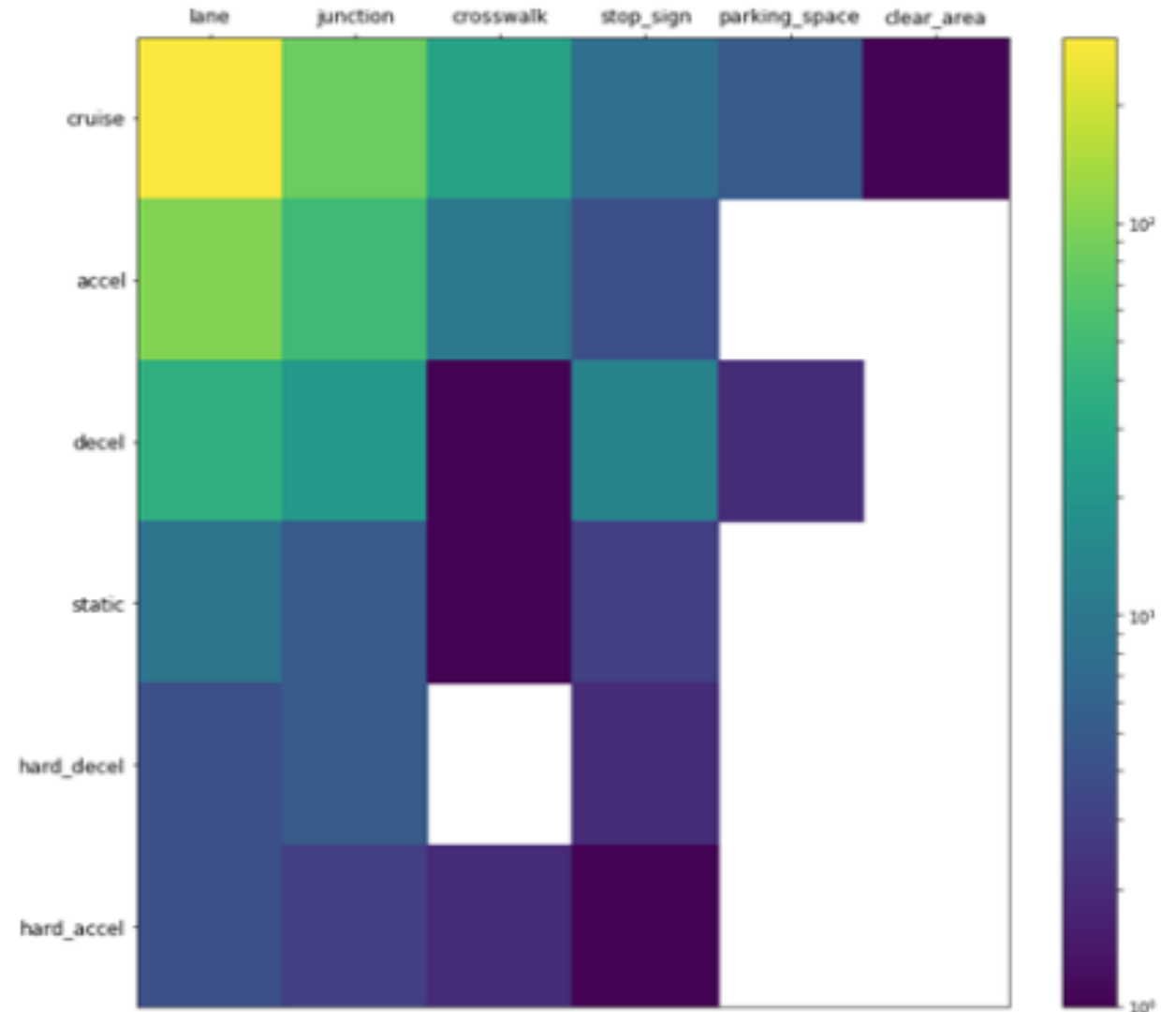
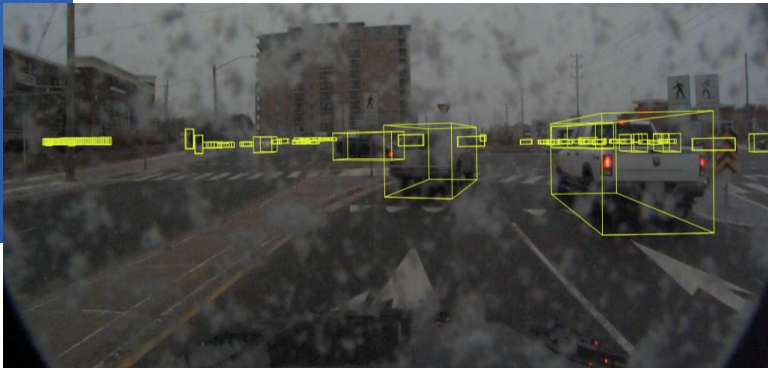
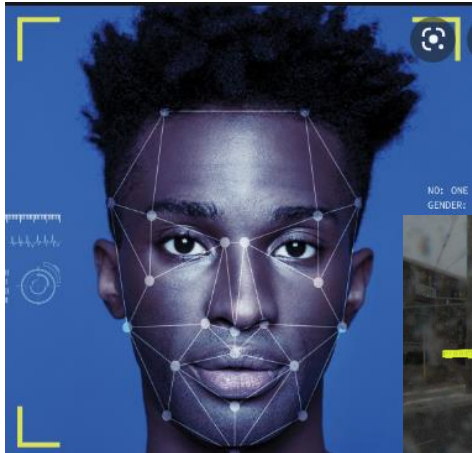
- Planning of collection routes
- Cost of operations
- Data pipelines and data infrastructure



# Data collection – pre-launch best practices



- Metadata important during data collection
  - Pre-trained networks used
  - Other sources - GPS or map data
- Metadata gives 3Ws (where, when and what) -
  - Help explain gaps in data collection
  - Target under-represented scenarios





# Data collection – production best practices



- Use ML confidence scores.
- Use customer feedback, e.g., customer spoke twice to voice recognition.
- Use ambiguous data, e.g., traffic sign detected close. Earlier frames selected.



# Synthetic data – for corner cases



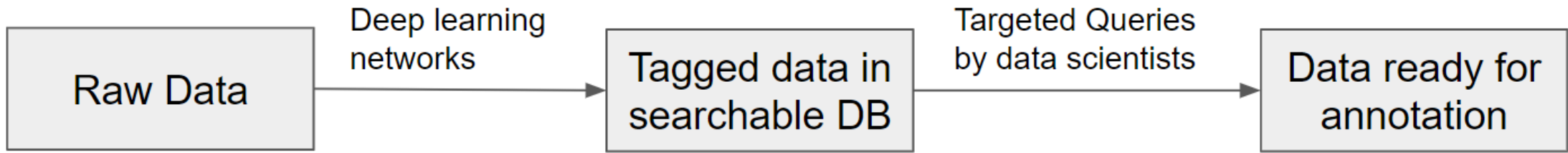
- Synthetic data for rarer conditions/difficult to capture real world scenarios.
- Synthetic data mixed with real data gives good performance.



# Data curation

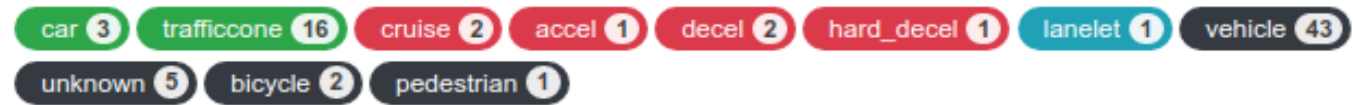
# Data curation – prelaunch

## Tagging data for targeted search



Deep learning networks/algos

### Tags



Metadata tags stored in databases with timestamps.  
Exposed for querying through a GUI

# Data curation – prelaunch

## Tools to select batch for annotation



Tool to search tags

E.g., Search stop sign, construction sign and other key scenarios

Data scientists use such tools

Criteria

Min Duration: 1  
Min Gap: 0.5

Terminal

camera2d\_objects in

trafficcone

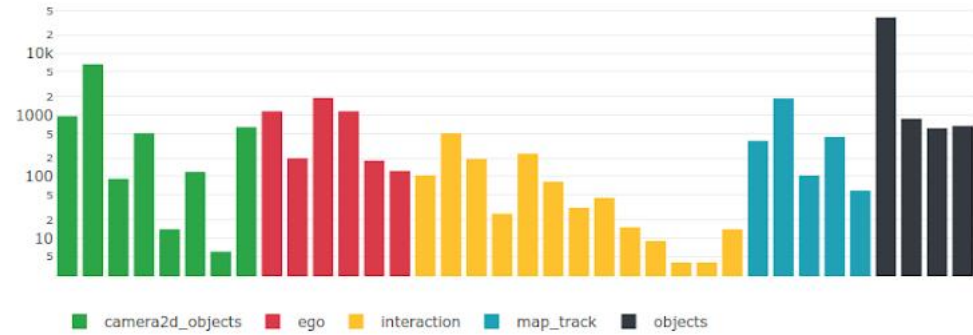
car

**stop\_sign**

truck

+ Add condition

Run



Download 0 Selected | Export for Annotation | Total duration 5m 20.3s

Start	Duration	Tags	Trip ID
	0:36 / 12:35		

Replay

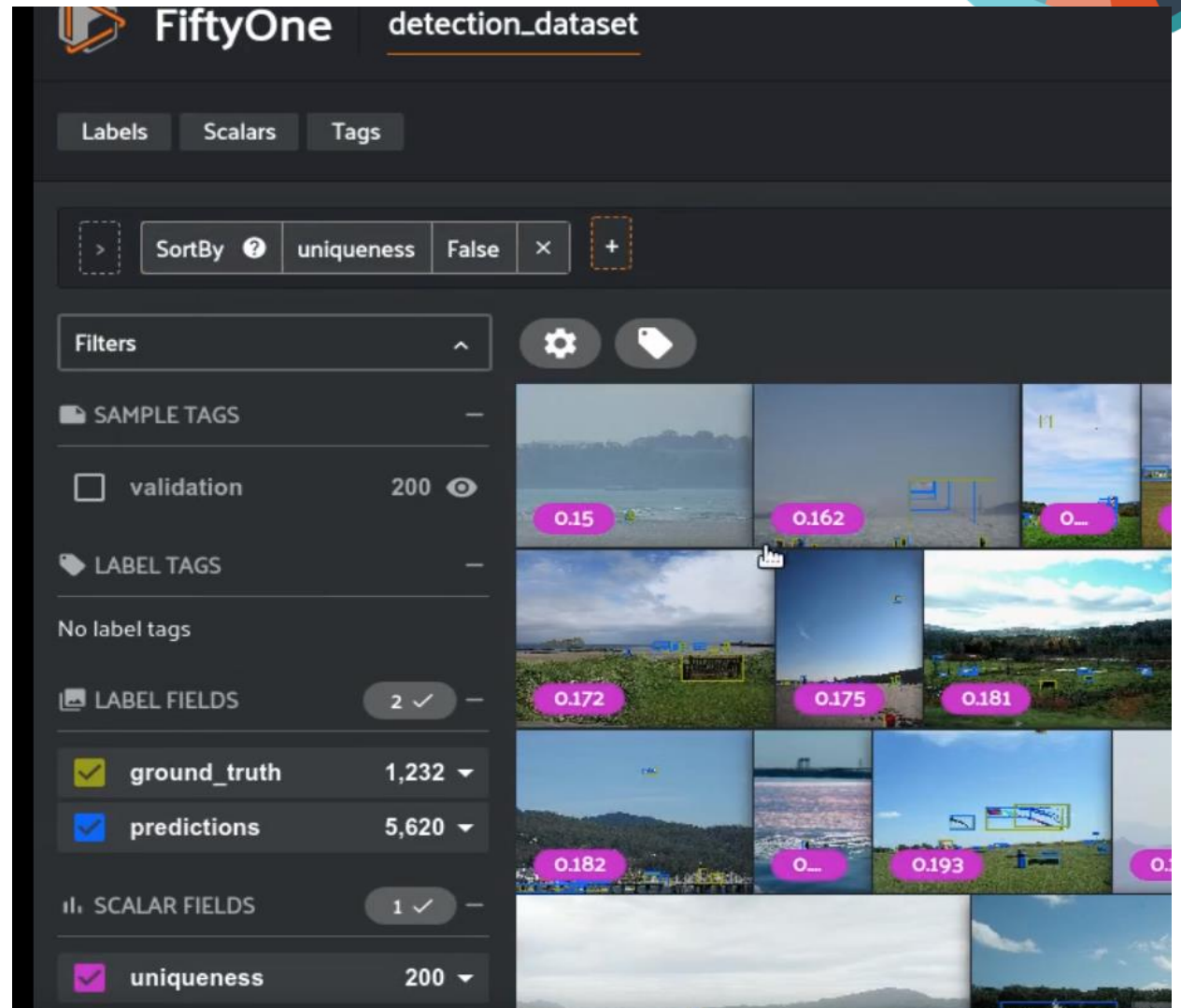
# Data curation – prelaunch

## Tools to select batch for annotation

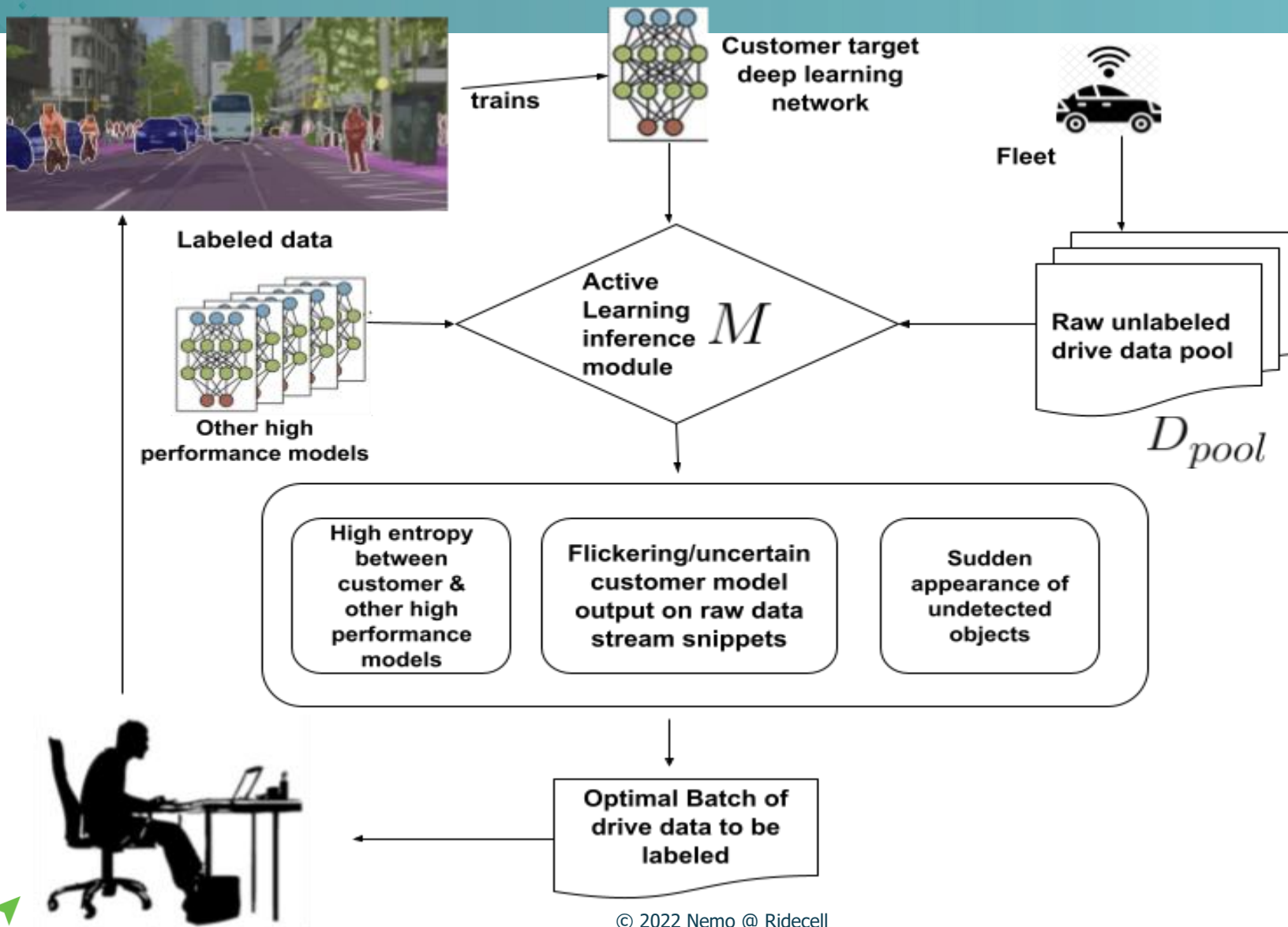


### Key features:

- Fast metadata queries – varied tags overlapping
- Re-run tagging with new models
- Support for text/image similarity
- Compare ground truth and ML inference visually



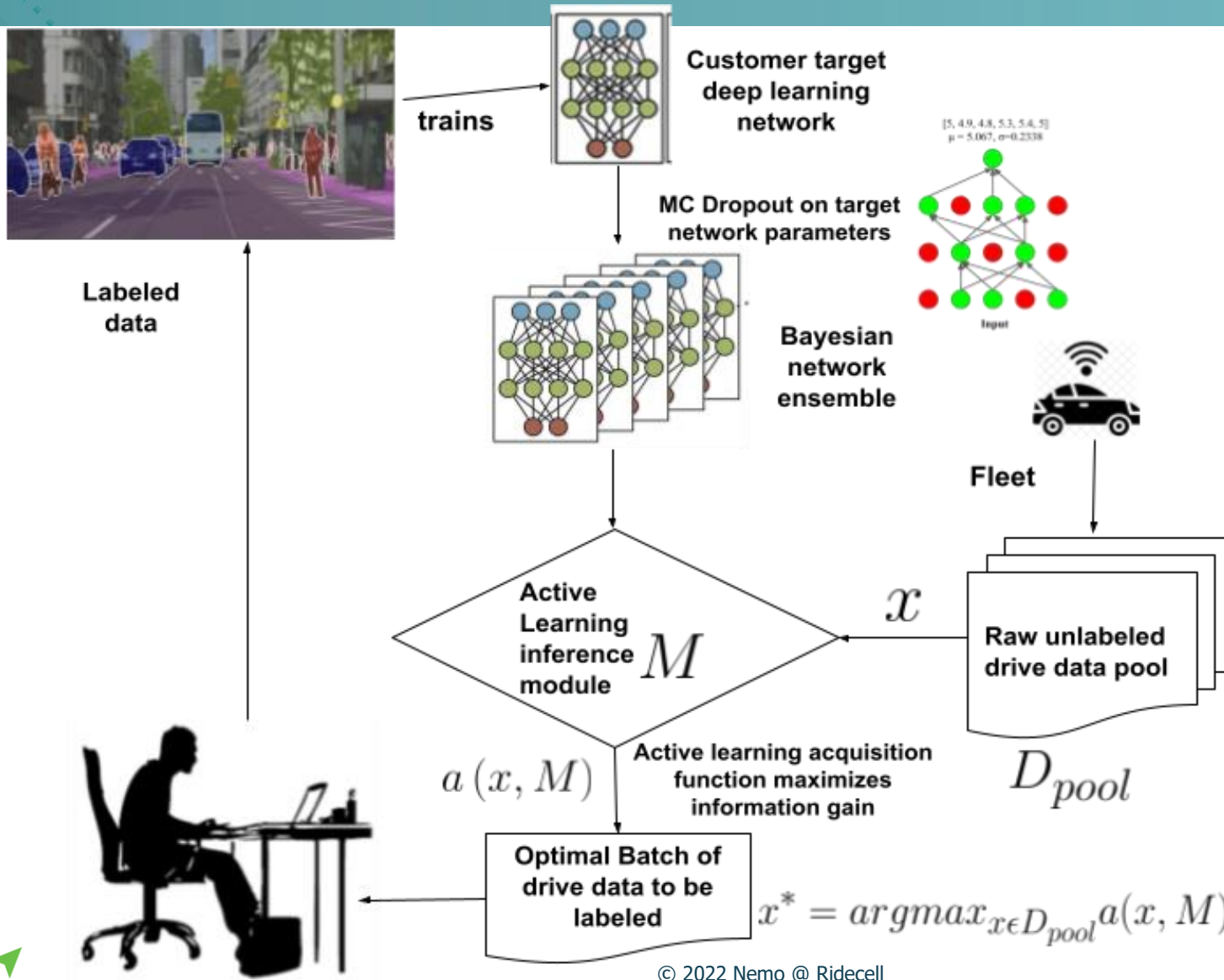
# Data curation – active learning in production systems



Example is image data from production car fleet.

This workflow can be modified for any type of data.

# Data curation – active Bayesian learning



Deep learning confidence scores not reliable.

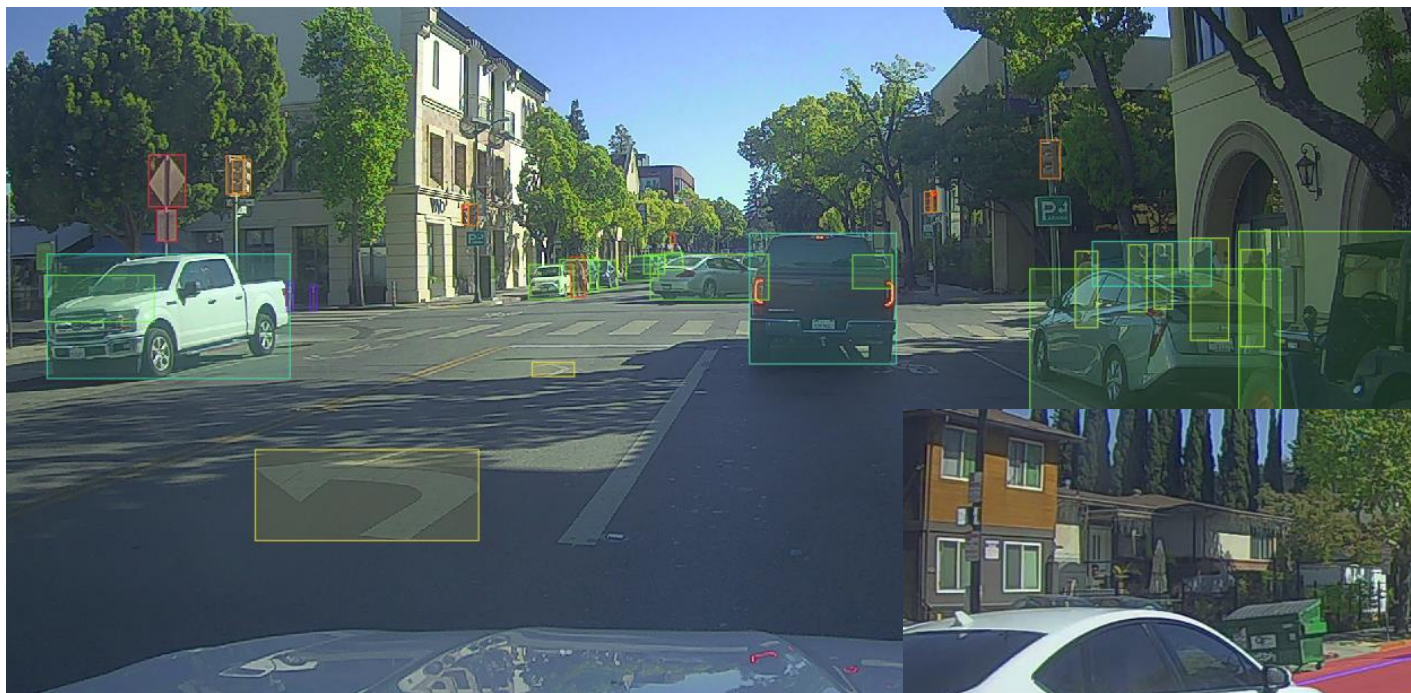
Convert a network to Bayesian network using MC Dropout.



# Data annotation

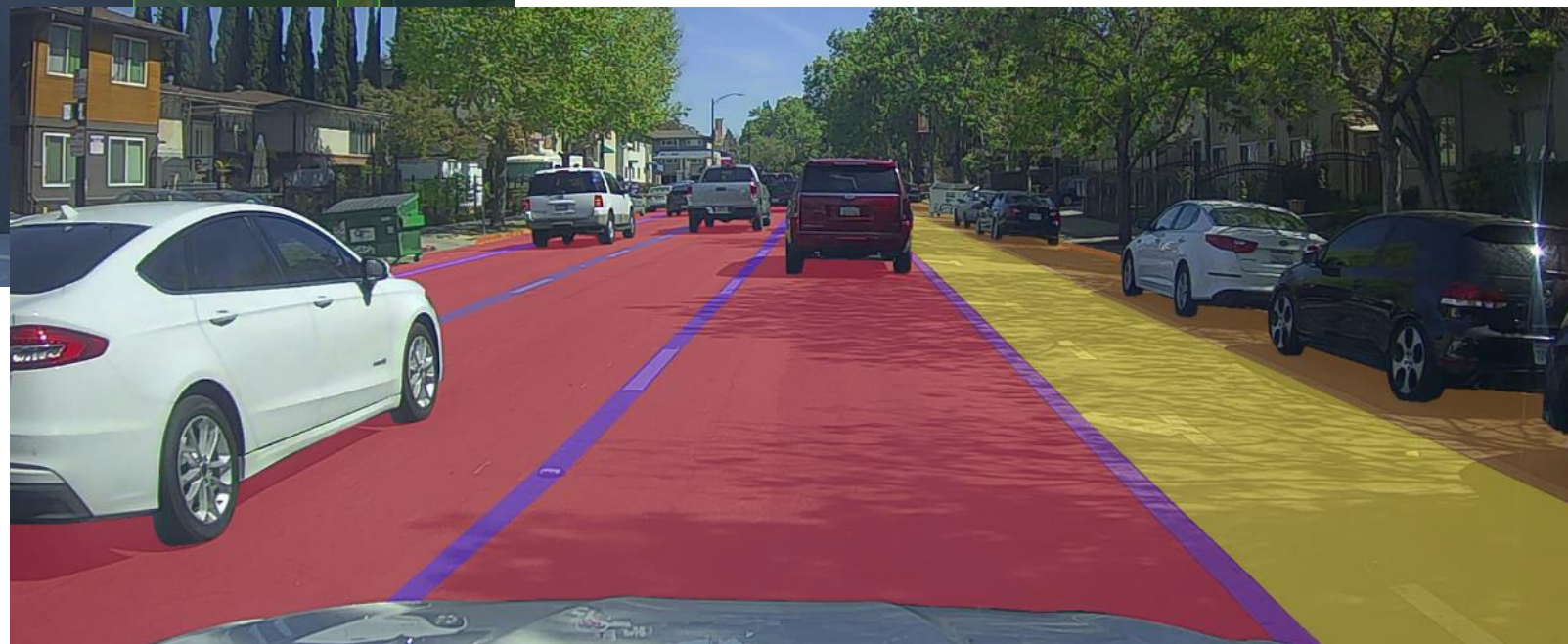
# Data annotation – types of annotation

## Image annotation



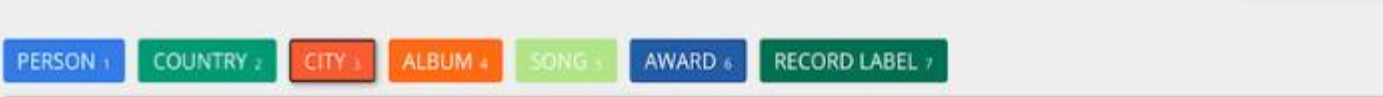
Bounding box

Semantic segmentation annotation



# Data annotation – types of annotation

## Other types

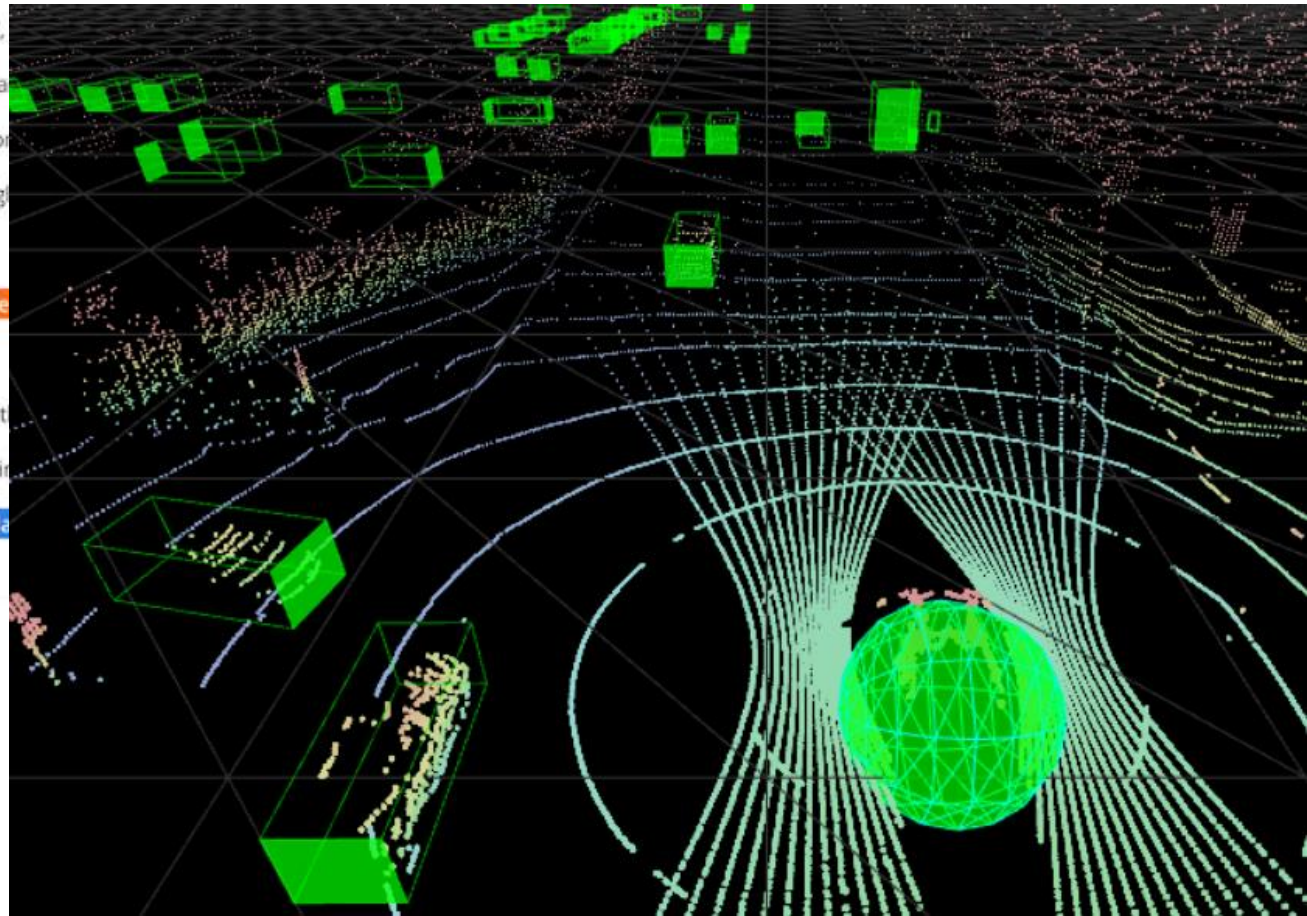


**Sia Kate Isobelle Furler** (/ˈsiːtə/ SEE-ə; born 18 December 1975) is an Australian singer, songwriter, record producer and music video director.[1] She started her career as a singer in the acid jazz band **Crisp** in the mid-1990s in Adelaide. In 1997, **Crisp** disbanded, she released her debut studio album titled **OnlySee** in **Australia**. She moved to **London, England**, and provided lead vocals for the British duo **Zero 7**. In 2000, **Sia** released her second studio album, **Healing Is Difficult**, on **Columbia** label the following year, and her third studio album, **Colour the Small One**, in 2004, but all of these struggled to connect with a mainstream audience.

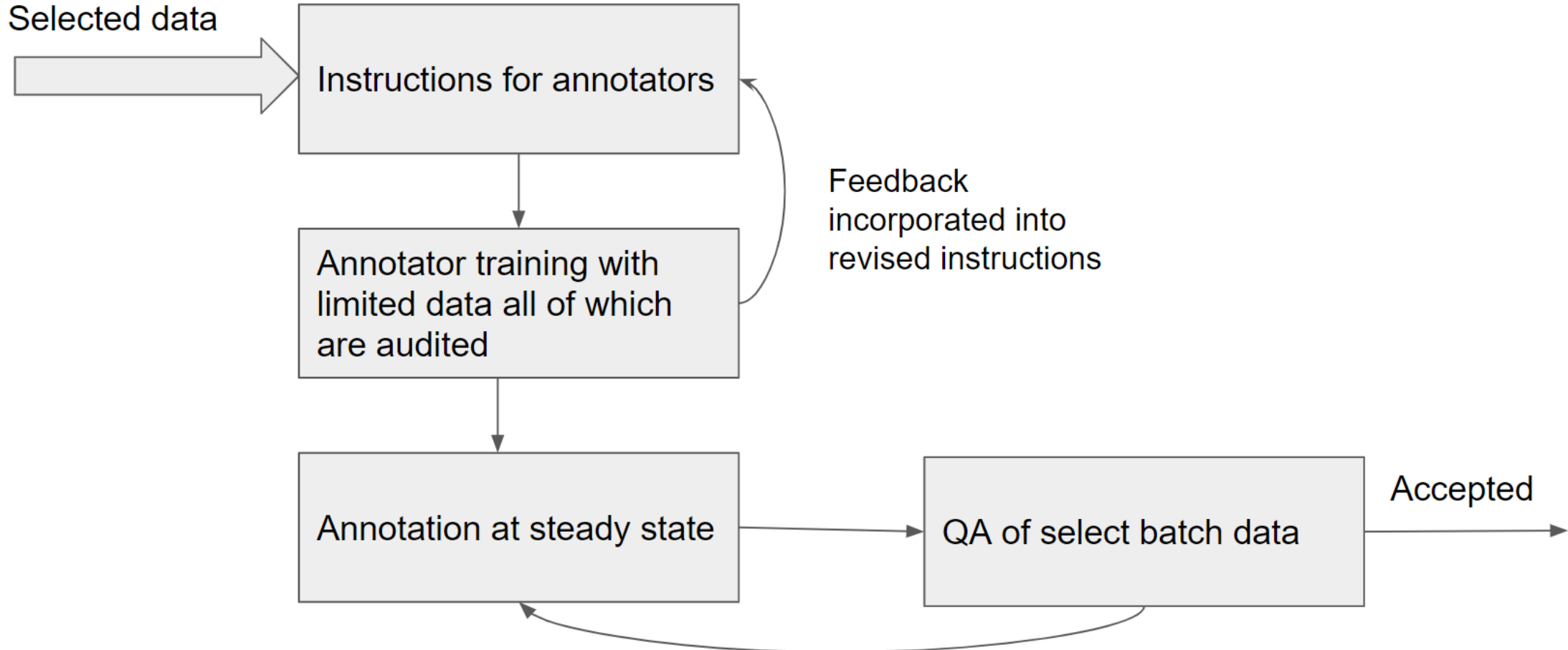
**Sia** relocated to **New York City** in 2005 and toured in the **United States**. Her fourth and fifth studio albums, **Some People Have Real Problems** and **We Are Born**, were released in 2008 and 2010, respectively. Each was certified gold by the Australian Recording Industry Association and attracted wider notice than her earlier albums. Uncomfortable with her growing fame, **Sia** took a hiatus from performing, during which she focused on songwriting for other artists, producing successful collaborations "Titanium" (with **David Guetta**), "Diamonds" (with **Rihanna**) and "Wild Ones" (with **Flo Rida**).

Text annotation image from Analytics

Lidar 3D bounding box annotation



# Data annotation – annotation process



Batch Rejected based on Service-level Agreement (SLA)

# Data annotation – annotation instructions

Ambiguous situations explained



## Correct sun glare

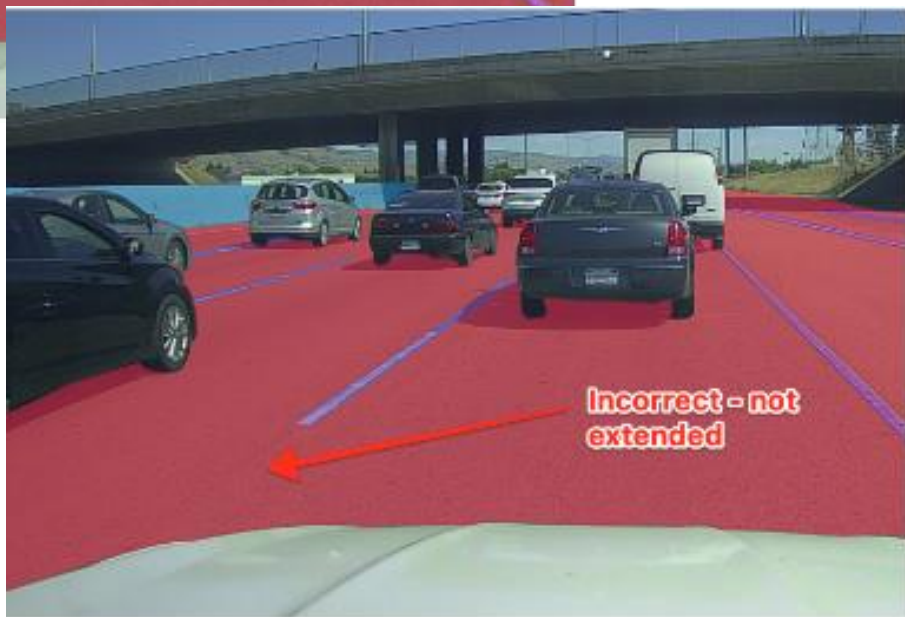
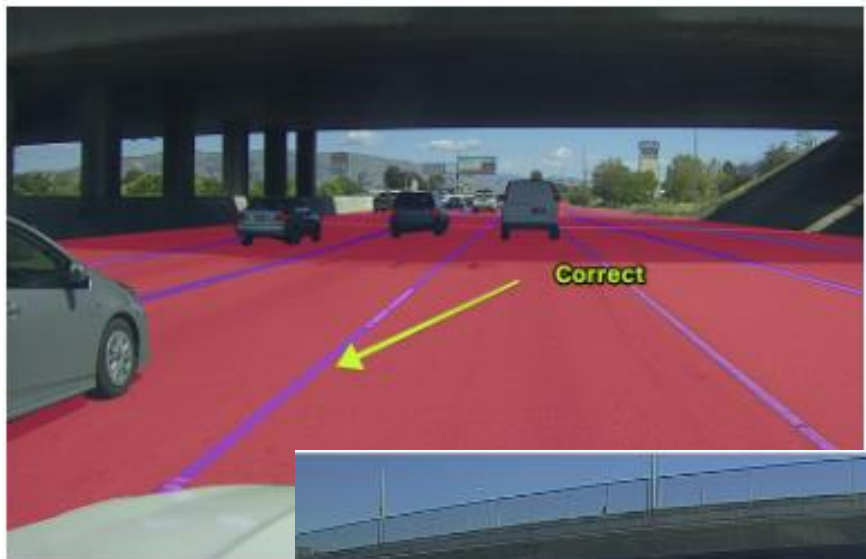


## No sun glare – even though there is slight glare and flaring



# Data annotation – annotation instructions






Ambiguous situations explained



# Data annotation – annotation instructions

## Scenario instructions example



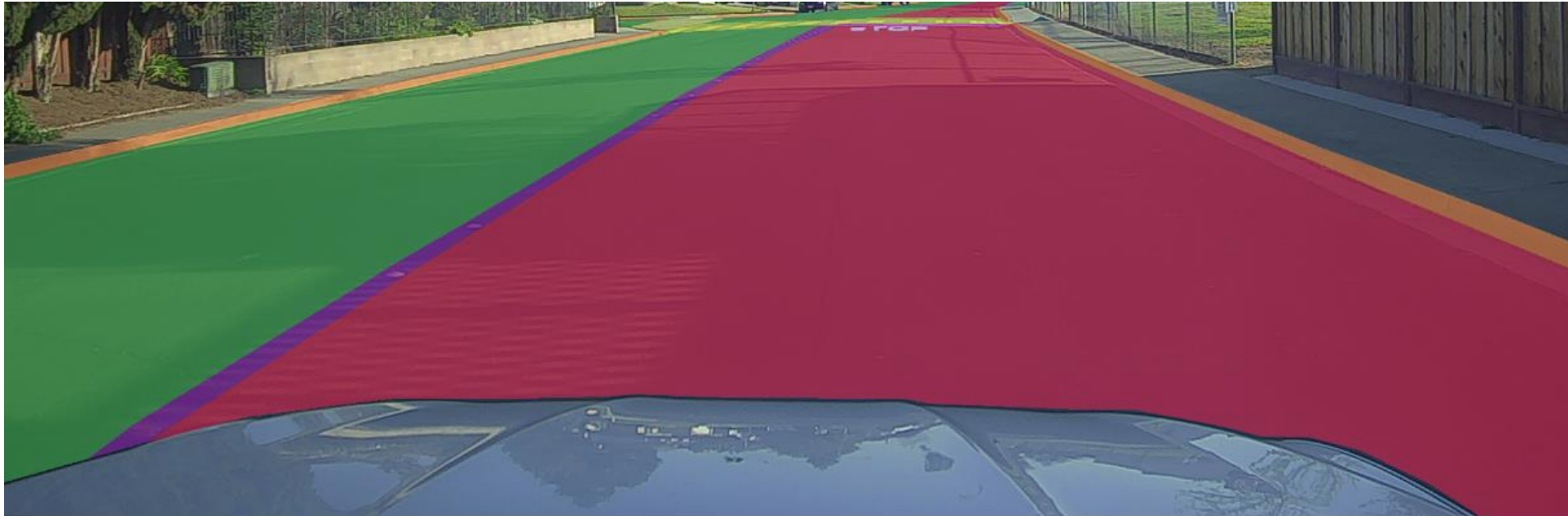
Interaction event	Description	Ego green/Vehicle blue
interaction:anti_parallel:cut_in:from_left	Object in adjacent lane, moving in opposite direction to ego and cutting into ego lane	
interaction:anti_parallel:turn_across:to_right	Object in adjacent lane, moving in opposite direction to ego and turning across ego lane to right	
interaction:anti_parallel:head_on	Object in adjacent lane, moving in opposite direction to ego and head on moving towards ego vehicle	
interaction:anti_parallel:cut_out:to_left	Object in adjacent lane, moving in opposite direction to ego and cutting out to left of ego	
interaction:anti_parallel:turn_out:to_left	Object in adjacent lane, moving in opposite direction to ego and turning out to perpendicular lane to left of ego	

Scenario annotation instructions example.

Clear visualizations and text description for annotators.

# Data annotation – auditing annotations

## Clear and precise feedback



Status: **REJECTED**      Audited: **7 MONTHS AGO**  
Comments: **STOP LINE 2 IS WRONGLY MARKED.**



# Data annotation – auditing annotations

## Clear and precise feedback



is: **REJECTED**      Audited: **A YEAR AGO**  
ments: **BOUNDING BOX FOR CAR B771 IS WRONG.. OTHER BOUNDING BOXES HAVE MINOR ERRORS**

# Data annotation – annotation in-house vs managed



	<b>In-house annotation</b>	<b>Managed service</b>
<b>Quality</b>	Generally, can be maintained higher	Depends on provider. If done with freelancers can be low
<b>Time to setup team</b>	Longer time to setup team.	Can be faster with managed service provider.
<b>Tooling</b>	Needs to purchased / developed in-house. Good options are now available	Good tooling through annotation companies.
<b>Scaling up</b>	Difficult	Easier as providers have bigger workforce.



## *Quality annotation is HARD!*

- Tooling issues
- ML assistance to annotation – can reduce quality!
- Bad annotation instructions
- QA in managed services is automated
- Annotator burn-out

# Conclusion

## Key Takeaways!

- Invest specifically in tooling for collection, tagging, annotation.
- Active learning is key during the curation stage.
  - Continuous improvement using data from production systems.
- Synthetic data for those hard-to-find corner cases.
- Pay close attention to annotation instructions

# Resources to learn more



## Resources

Active learning for deep learning

<https://jacobgil.github.io/deeplearning/activelearning>

Tools for ML data management

<https://scale.com/nucleus>

<https://nemosearch.ai/>

<https://voxel51.com/>