

January 16, 2025  
Daniel Bogdoll

# McCity Data Engine

## From Data Stream to AI Model



# About Me



*2011 - 2019* • M.Sc.  
**Computational Engineering Science**



*2016* • Research Intern  
**Autonomous Driving**



*2018 - 2021* • CEO & Co-Founder  
**Carpooling Startup**



*2020 - 2025* • PhD Candidate  
**Autonomous Driving**



*2024 - 2025* • Research Scholar  
**Data Engine**

# About Mcity

*2021* • Expansion  
Smart Intersections

*2015* • Foundation  
32-acre “mini-city”



# About Mcity

*2021* • Expansion  
**Smart Intersections**

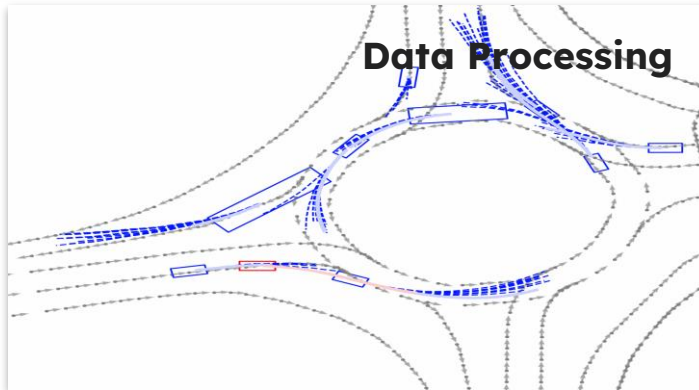
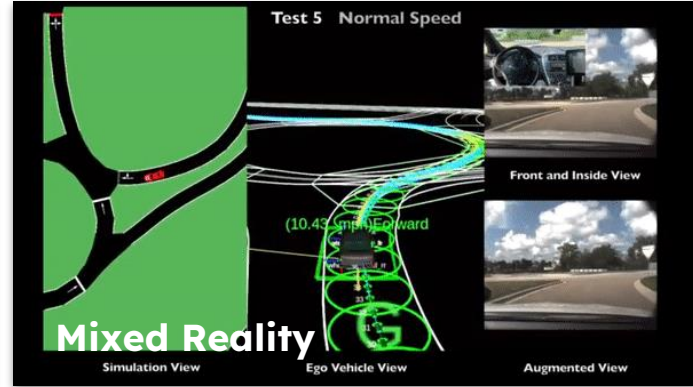
*2024* • Digital Twin  
**Mcity 2.0**

*2015* • Foundation  
**32-acre “mini-city”**





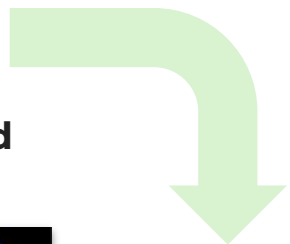
# About Mcity



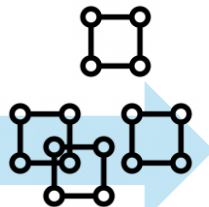
# From Data to Model



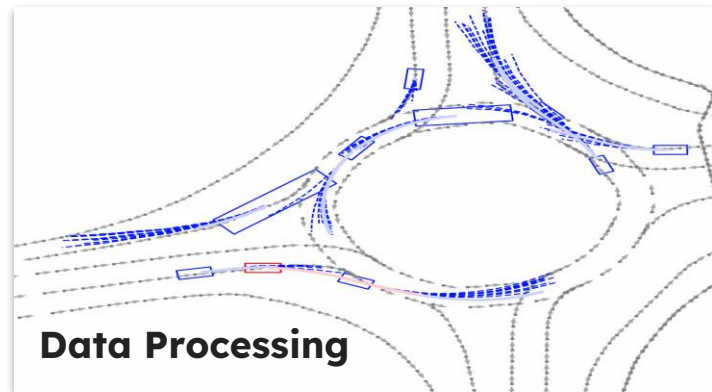
**Labeled  
Data**



**Neural  
Network**



**Abstract  
Data**

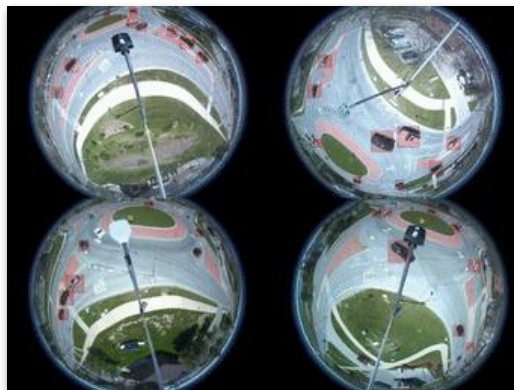
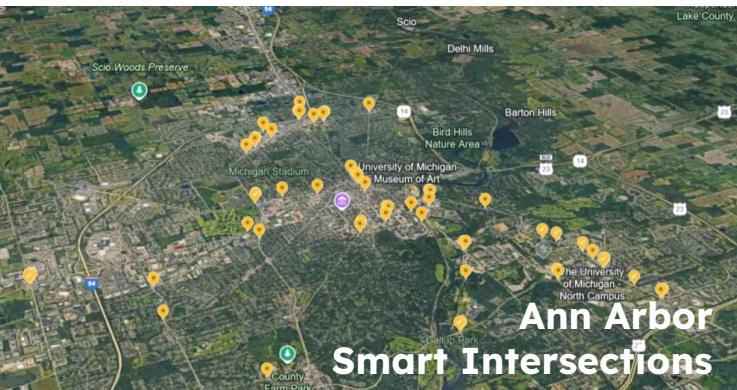
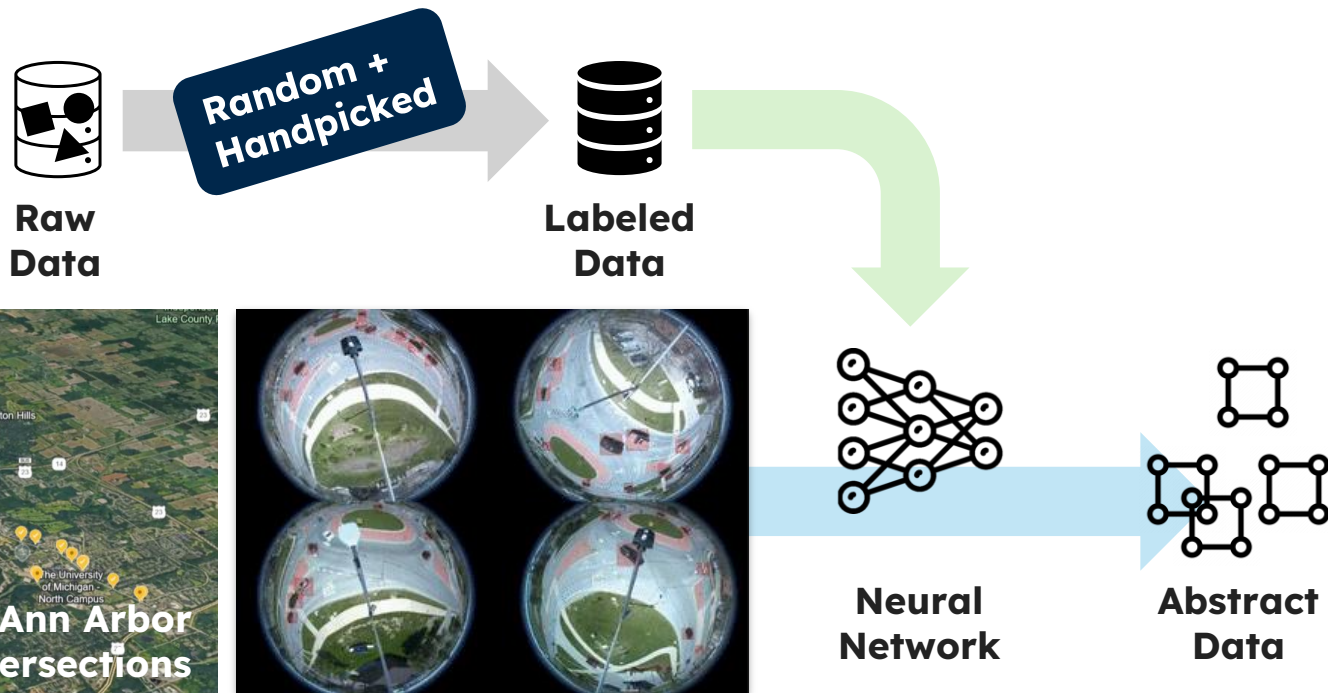


**Data Processing**

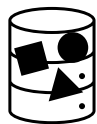
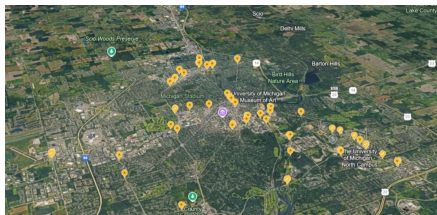


**Real World Data**

# From Data to Model



# From Data to Model

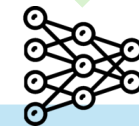
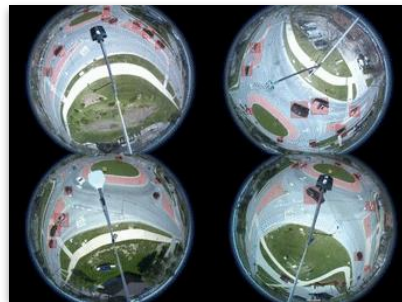
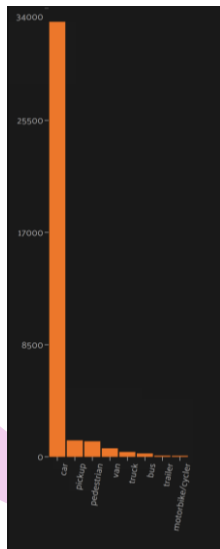


**Raw Data**

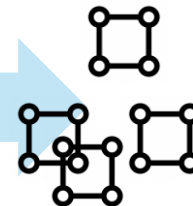
**Random + Handpicked**



**Labeled Data**



**Neural Network**



**Abstract Data**



**Model Issues**



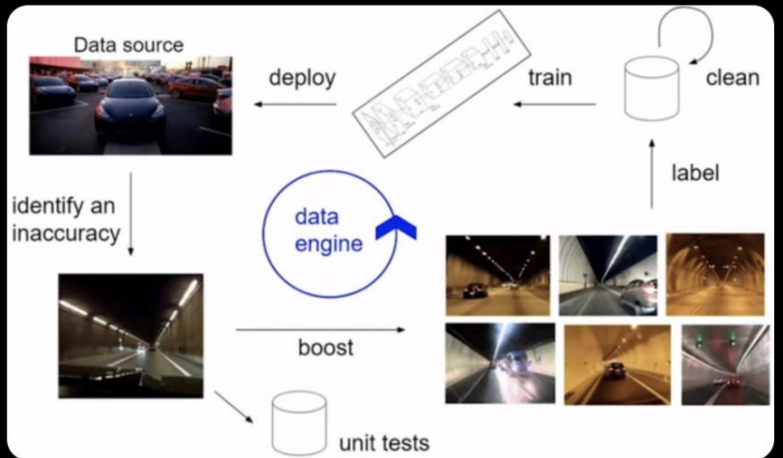
# Data Engine



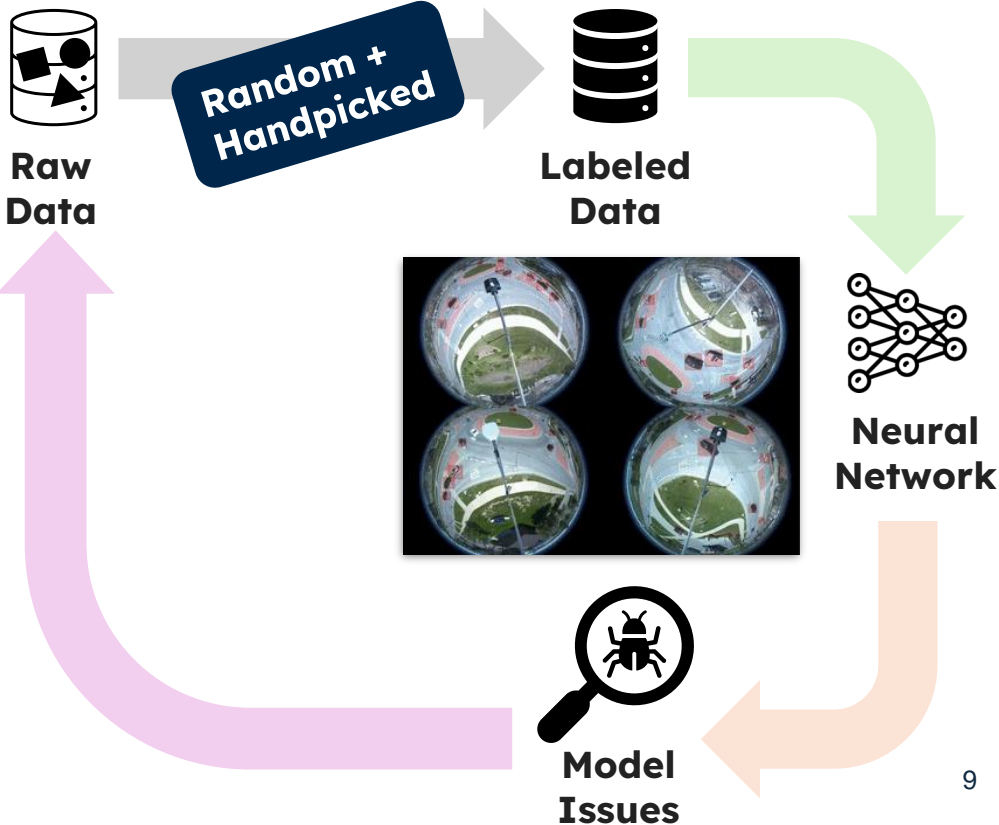
Andrej Karpathy  
@karpathy

Follow

Potentially nitpicky but competitive advantage in AI goes not so much to those with data but those with a data engine: iterated data acquisition, re-training, evaluation, deployment, telemetry. And whoever can spin it fastest. Slide from Tesla to ~illustrate but concept is general



2:47 PM · Dec 5, 2022



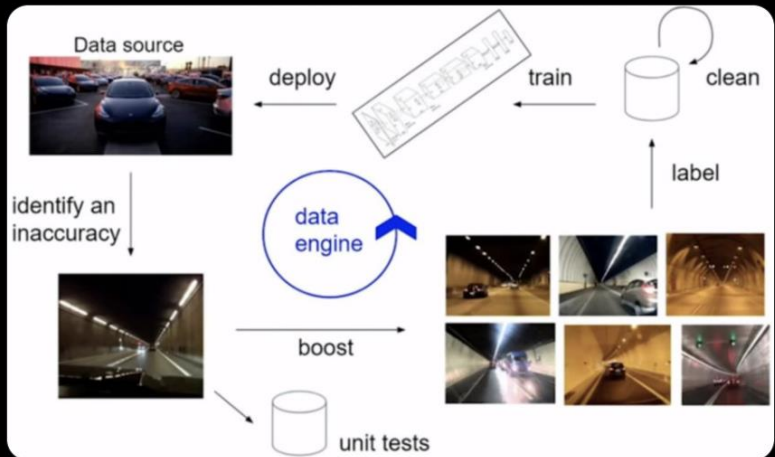
# Data Engine



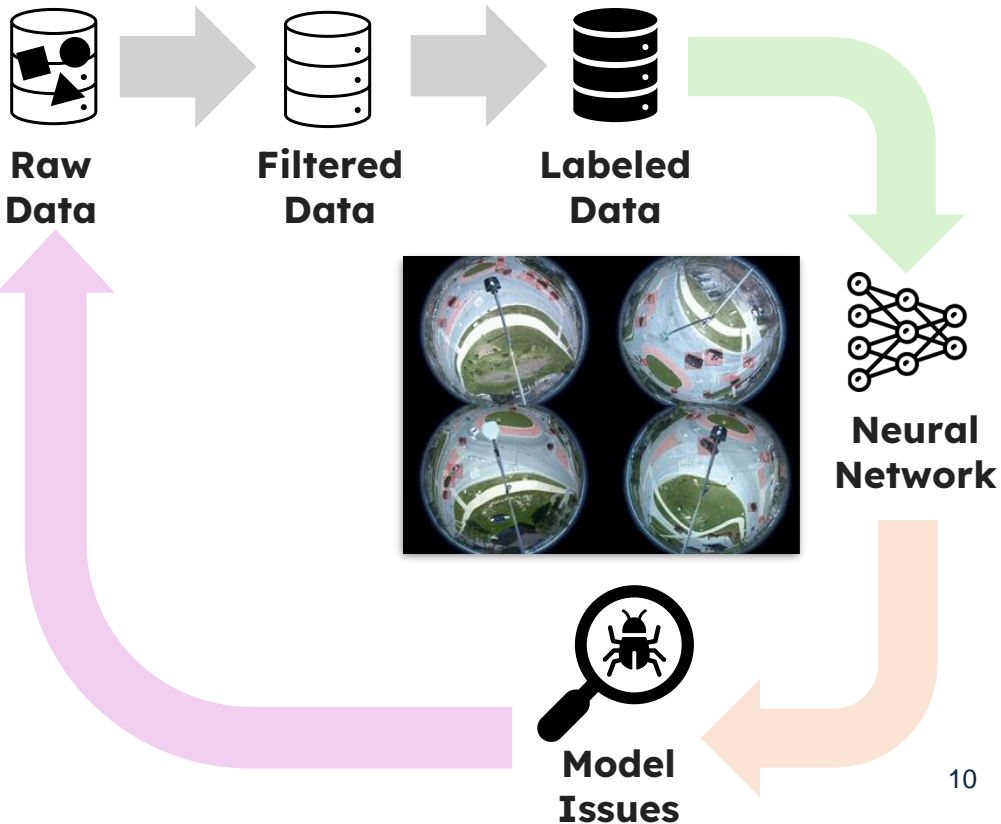
Andrej Karpathy ✓  
@karpathy

Follow

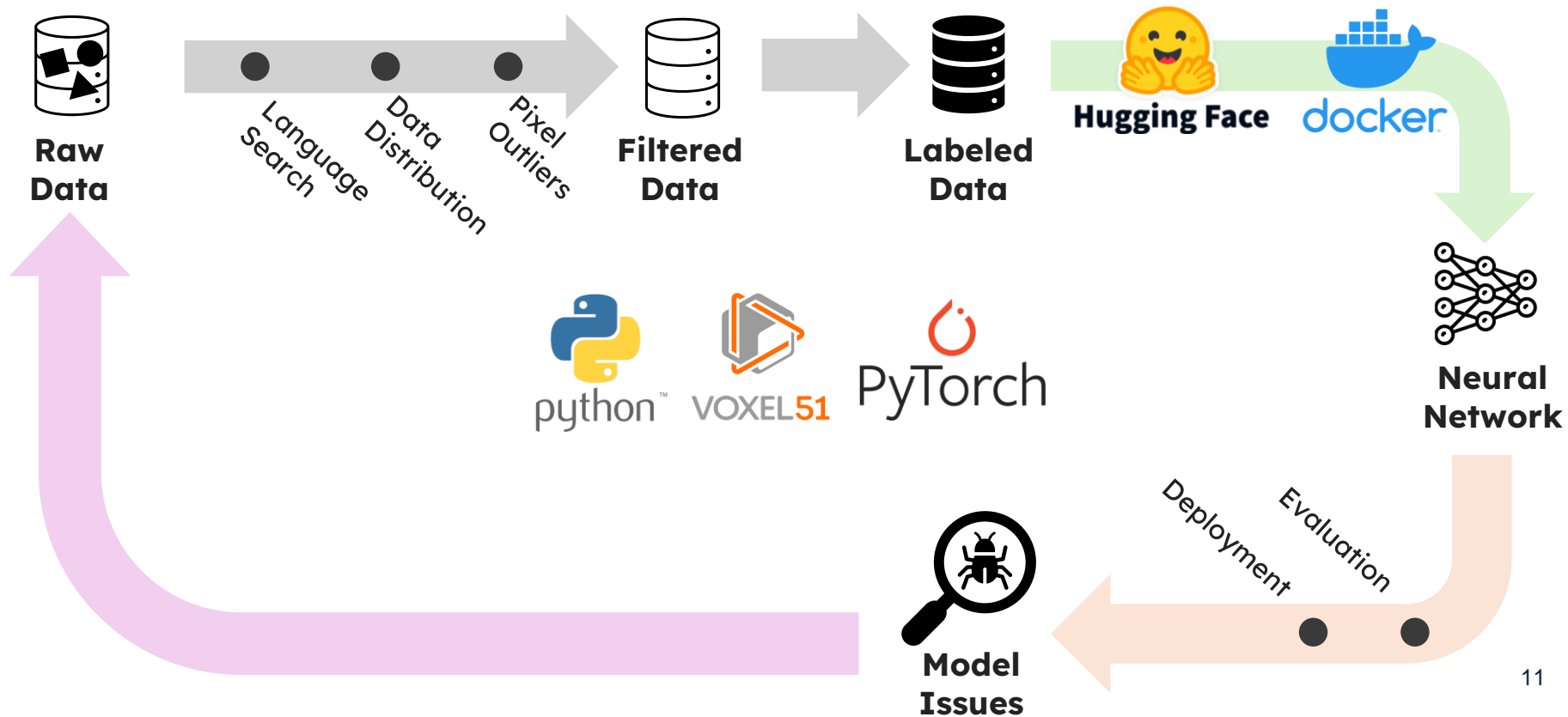
Potentially nitpicky but competitive advantage in AI goes not so much to those with data but those with a data engine: iterated data acquisition, re-training, evaluation, deployment, telemetry. And whoever can spin it fastest. Slide from Tesla to ~illustrate but concept is general



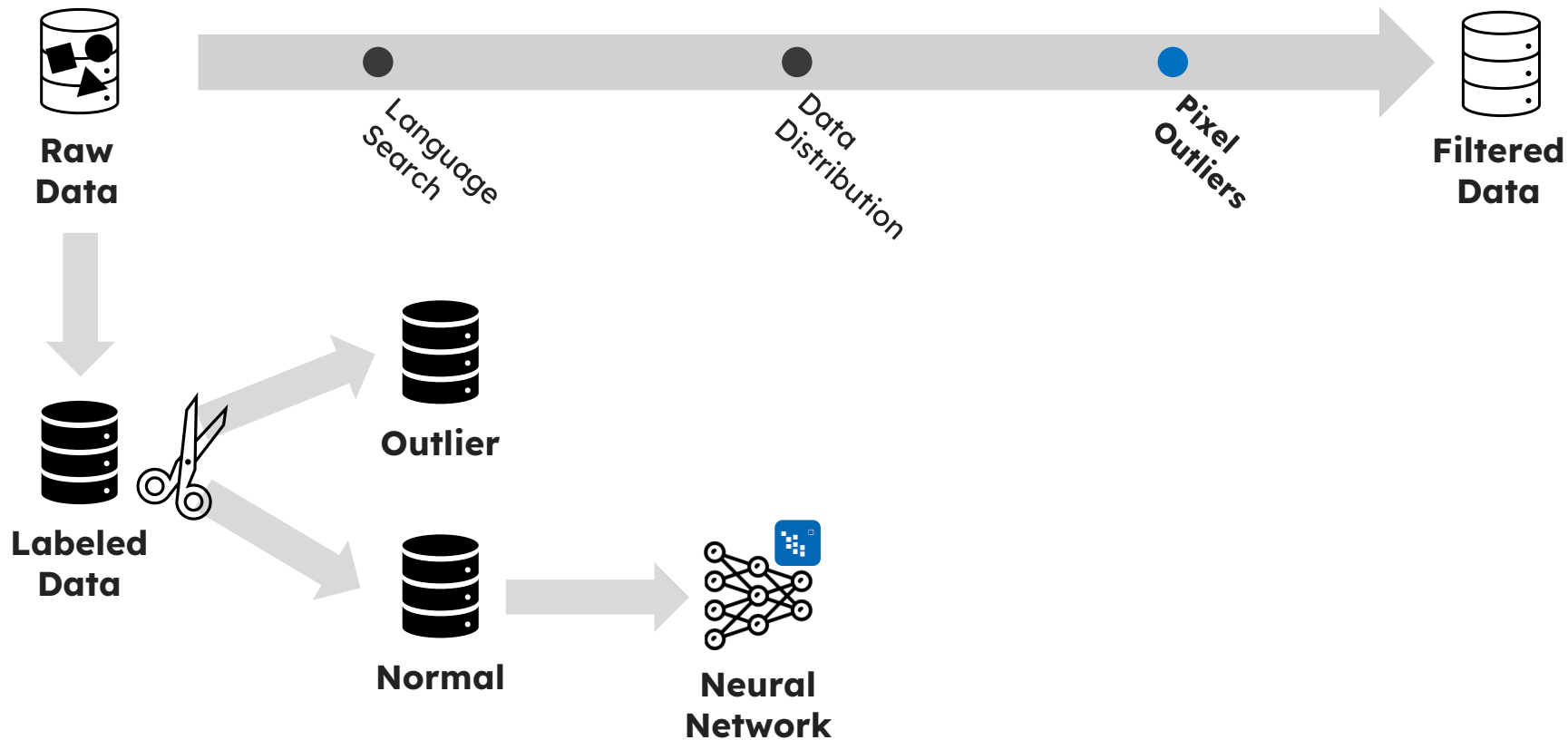
2:47 PM · Dec 5, 2022



# McCity Data Engine



# Outlier Detection

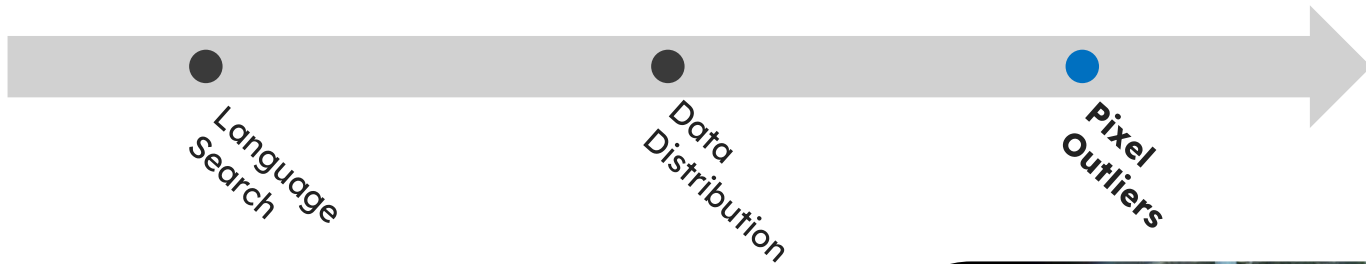




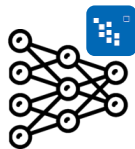
# Outlier Detection



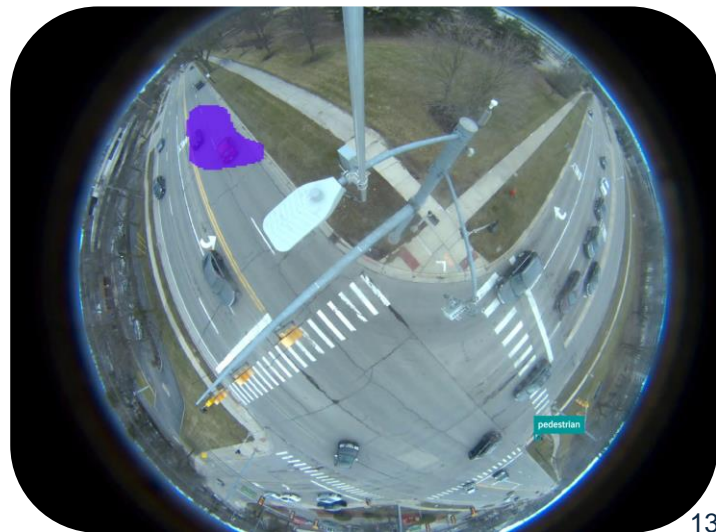
Raw  
Data



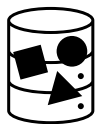
Filtered  
Data



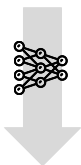
Neural  
Network



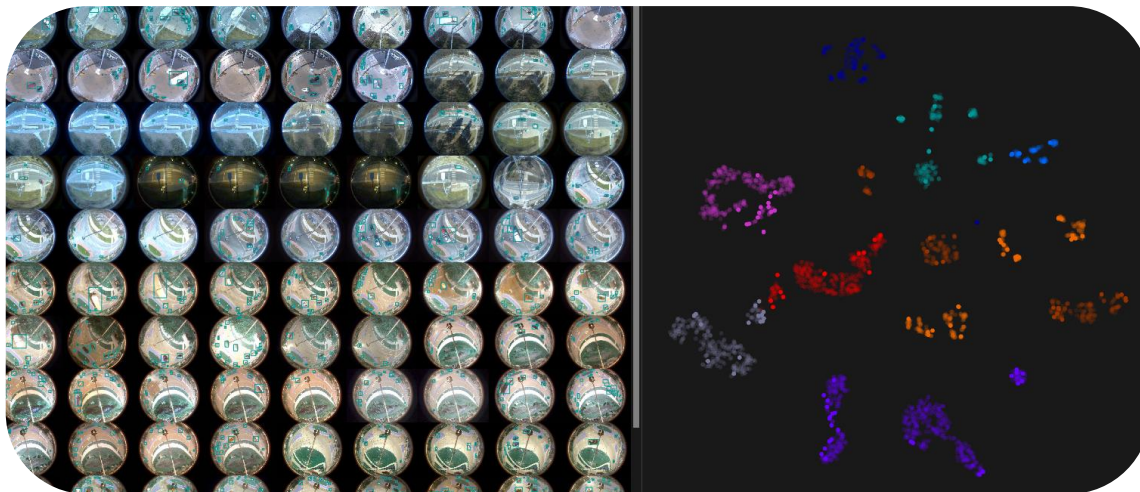
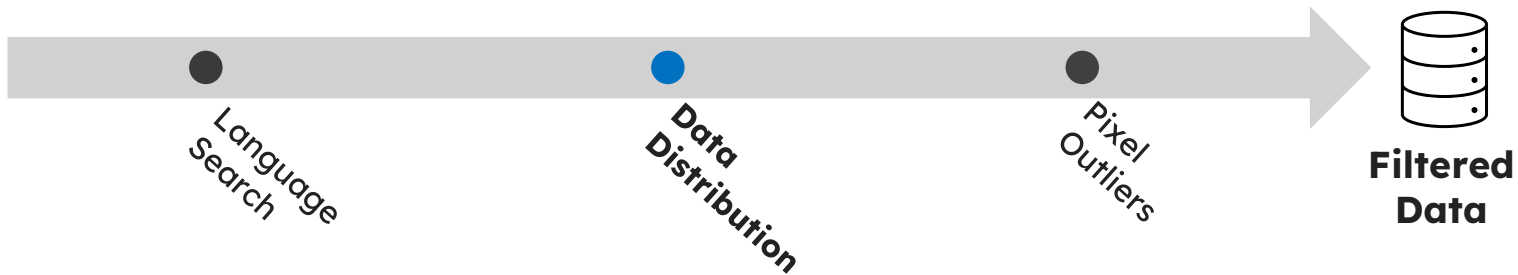
# Data Distribution



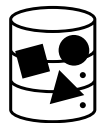
Raw  
Data



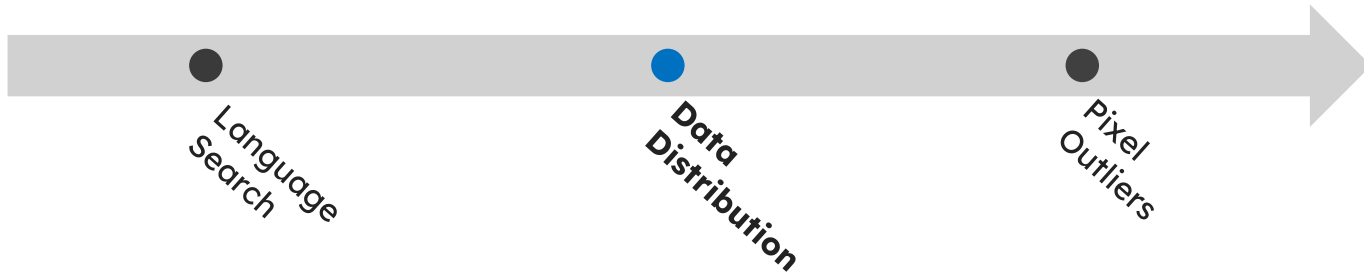
Embeddings



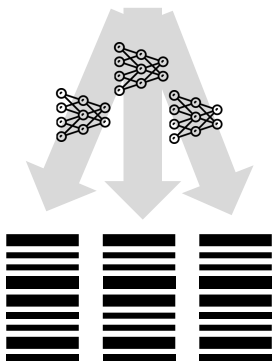
# Data Distribution



Raw Data



Filtered Data



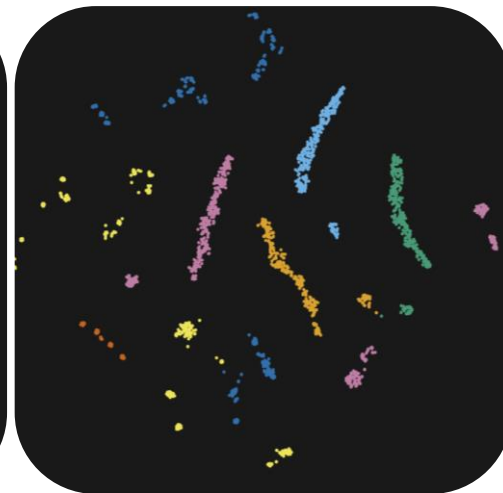
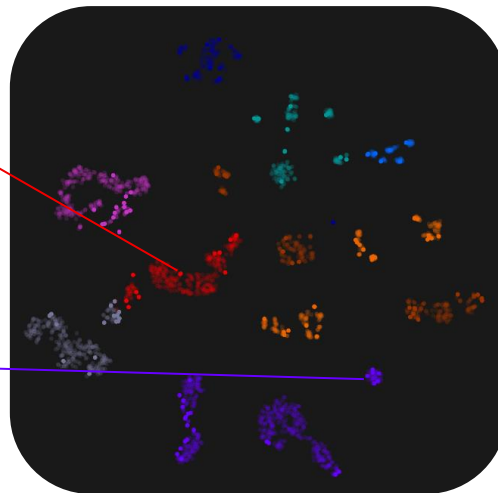
Embedding Ensemble



Representative Samples



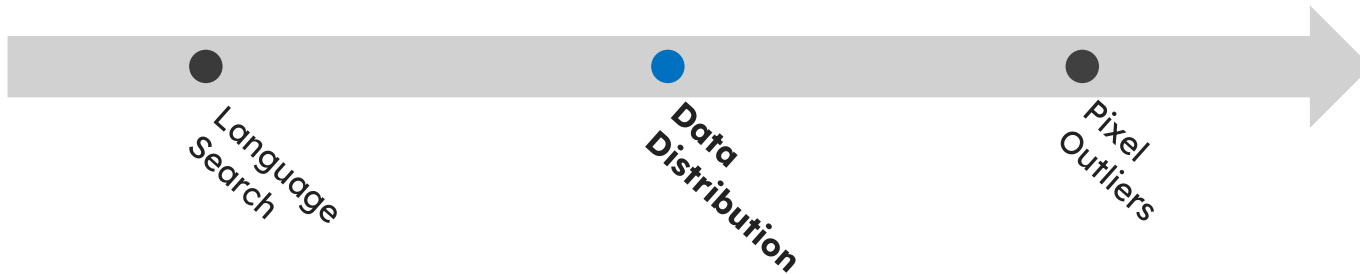
Rare Samples



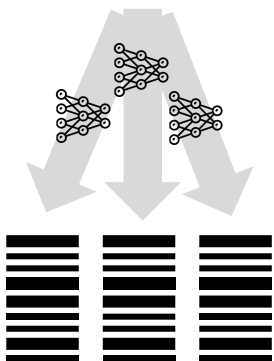
# Data Distribution



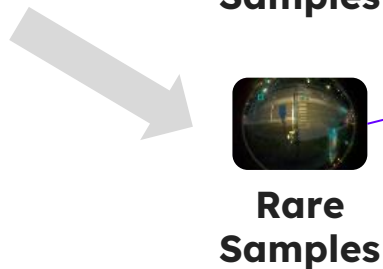
Raw Data



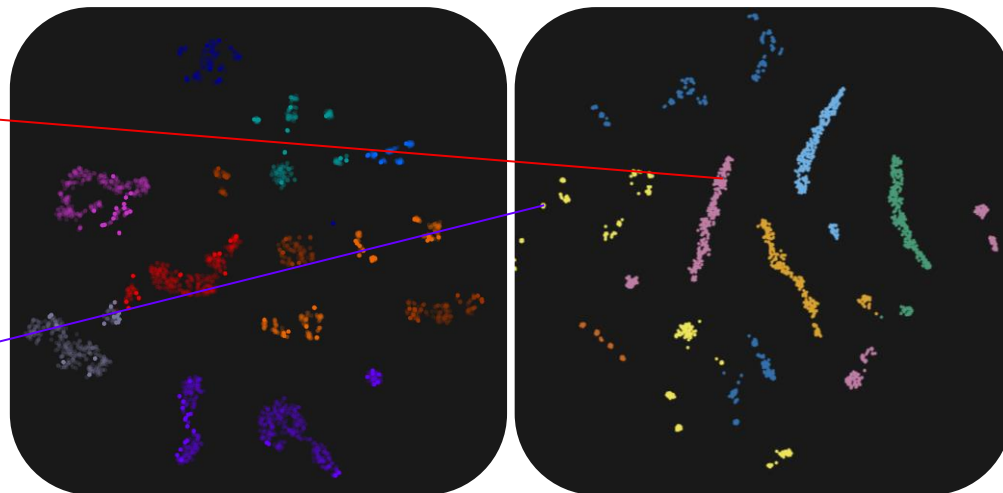
Filtered Data



Embedding Ensemble

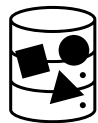


Rare Samples

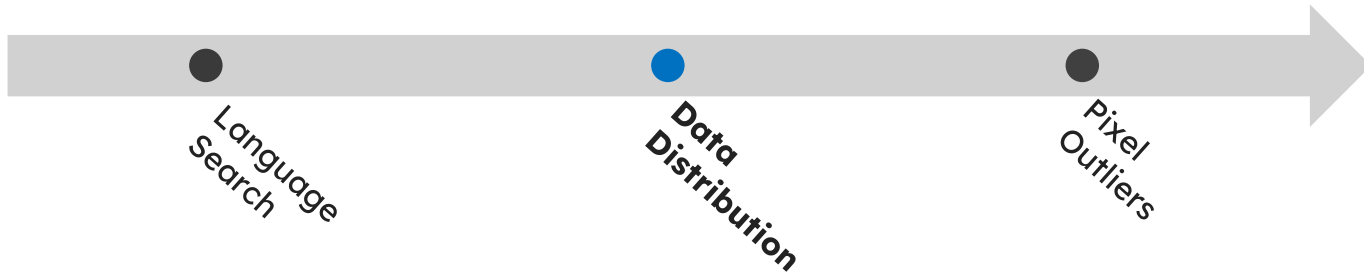




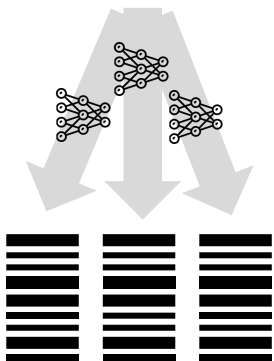
# Data Distribution



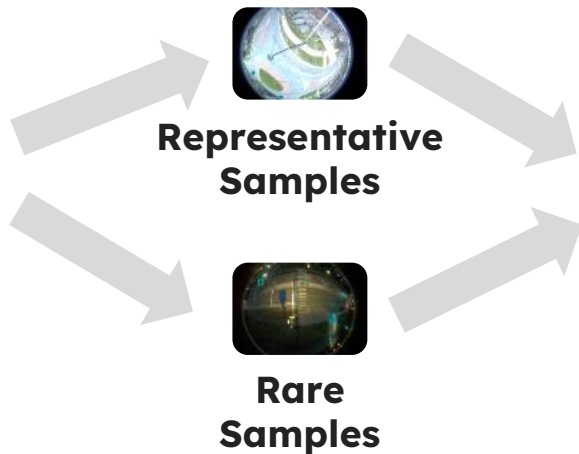
**Raw Data**



**Filtered Data**



**Embedding Ensemble**

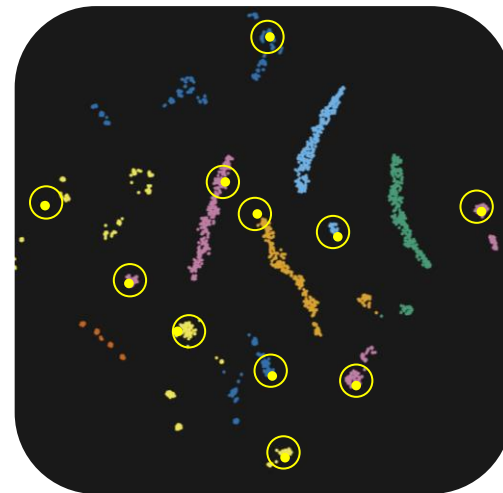


**Representative Samples**

**Rare Samples**



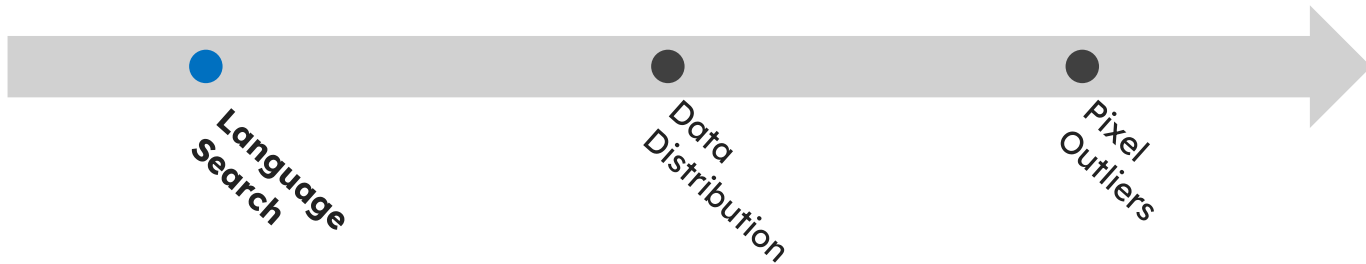
**Greedy Neighbor Selection**



# Language Search

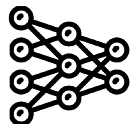


Raw Data



Filtered Data

“cyclist”

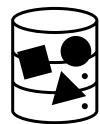


Zero Shot Detection



Many  
false positives  
and  
false negatives

# Language Search



Raw Data

“cyclist”

Language Search



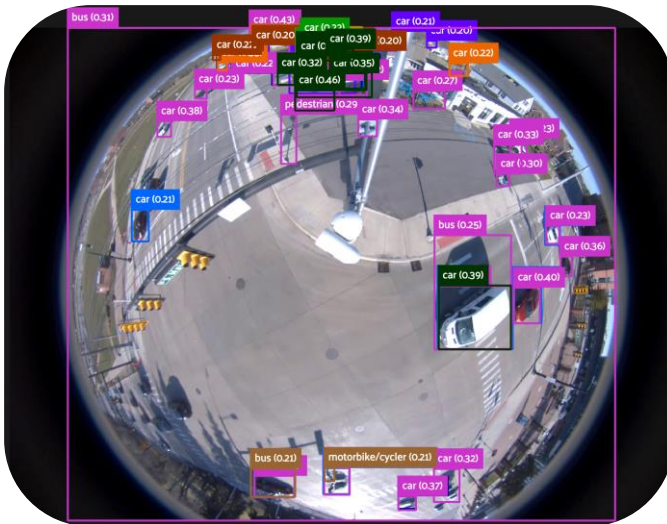
Zero Shot Ensemble

Data Distribution

Pixel Outliers



Filtered Data



Instances Filter



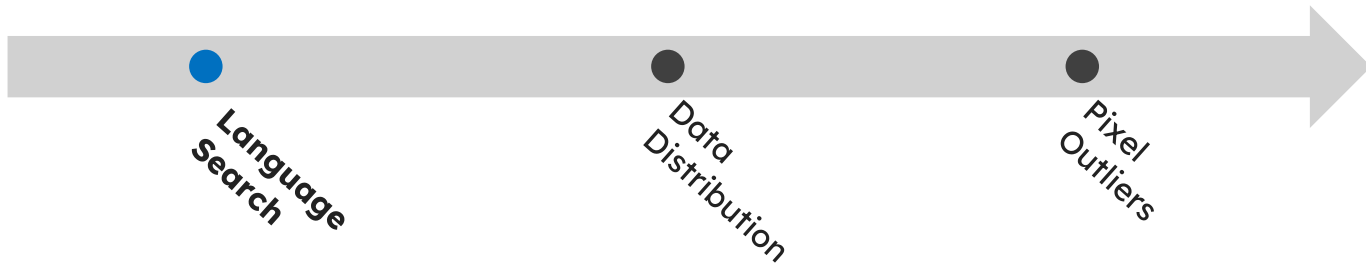
3/5

Agreement Filter

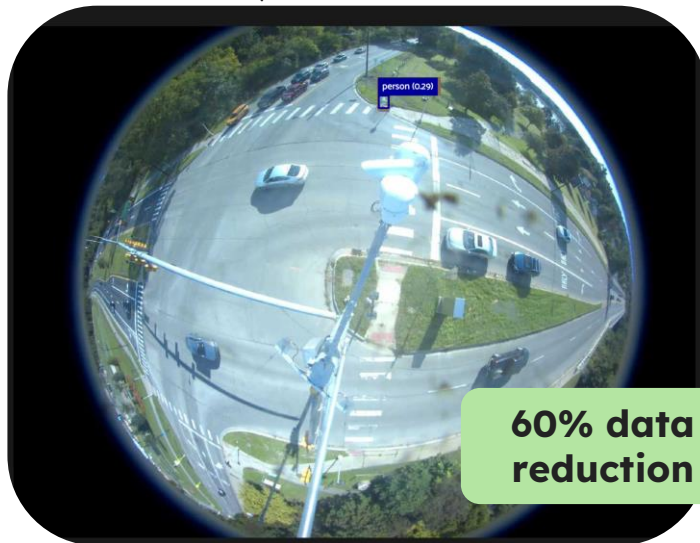
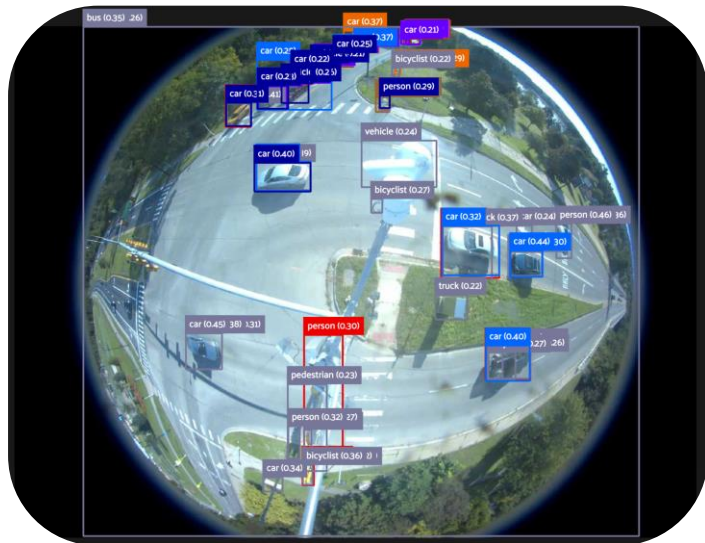
# Language Search



Raw Data



Filtered Data



Instances Filter

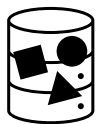


**3/5**

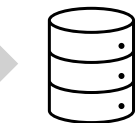
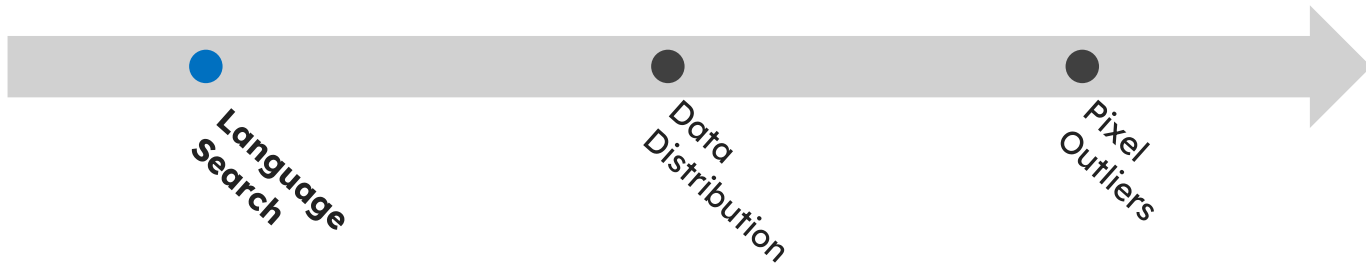
Agreement Filter



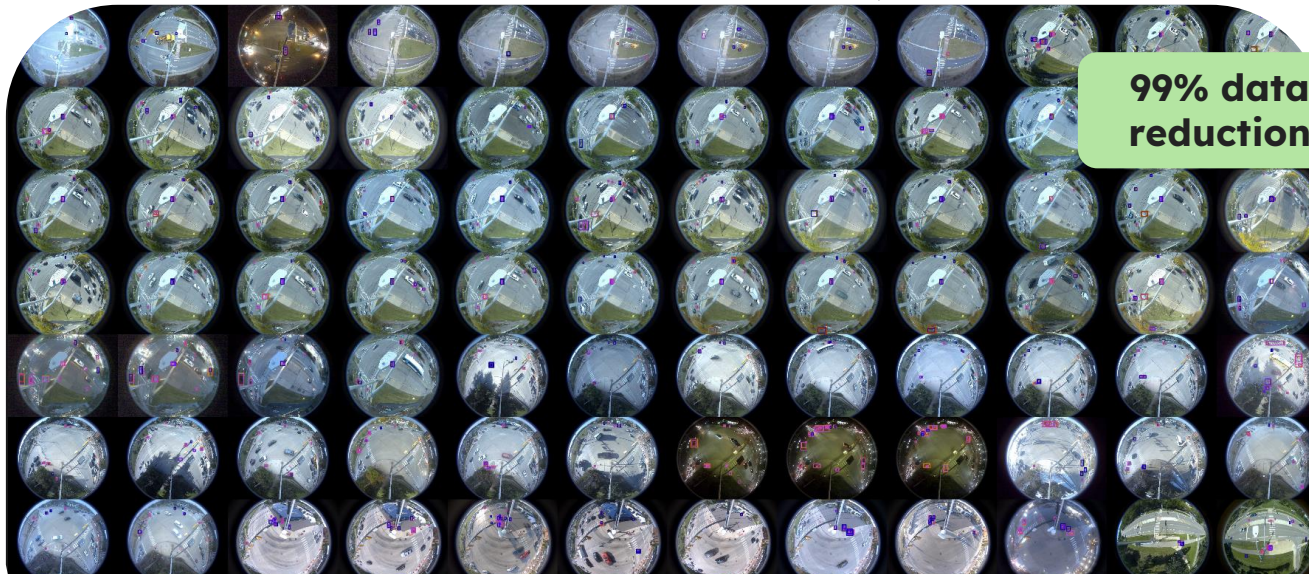
# Language Search



Raw  
Data



Filtered  
Data



99% data  
reduction



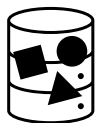
Instances  
Filter



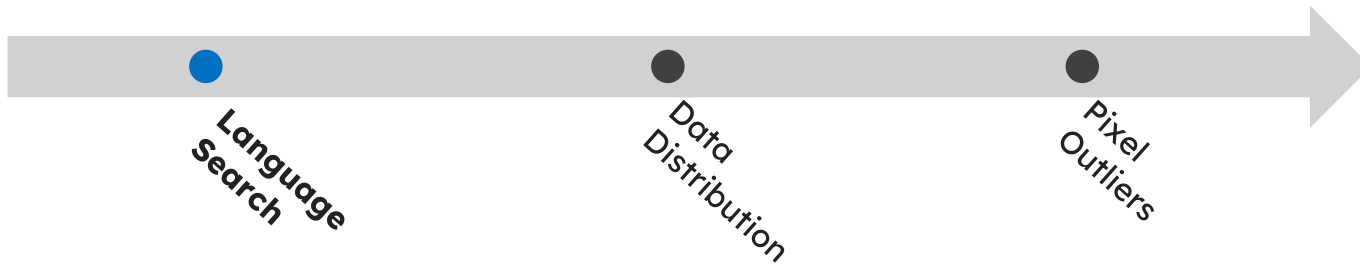
**3/5**

Agreement  
Filter

# Language Search



Raw  
Data



Filtered  
Data

Voxel51  
Data-Centric AI competition

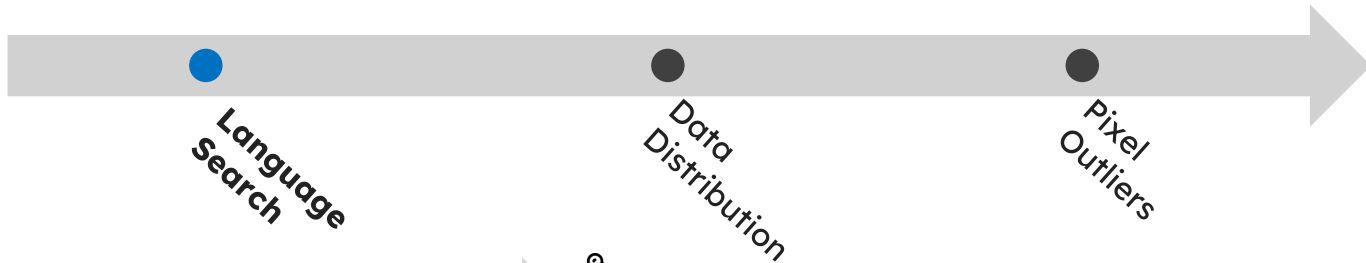


**3/5**  
Agreement  
Filter

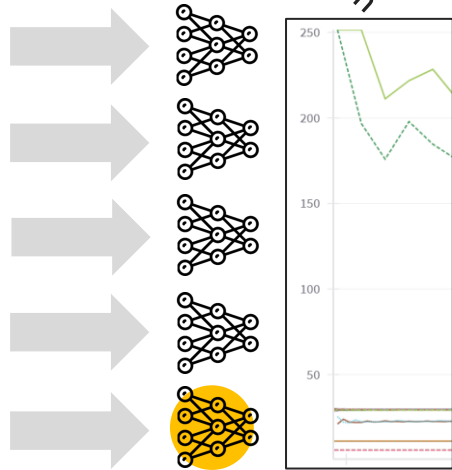
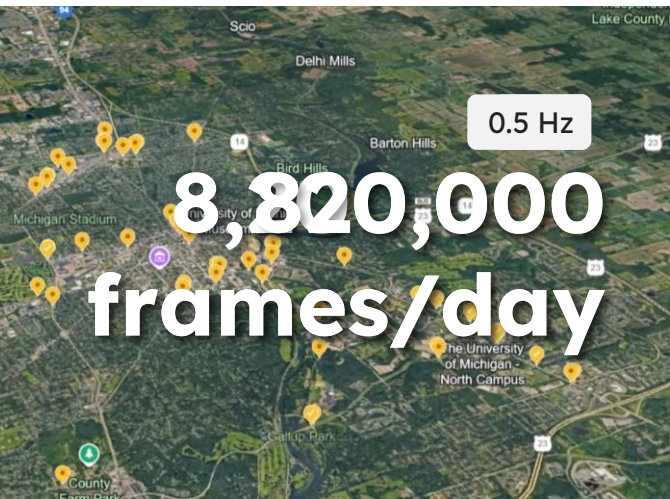
# Language Search



Raw Data



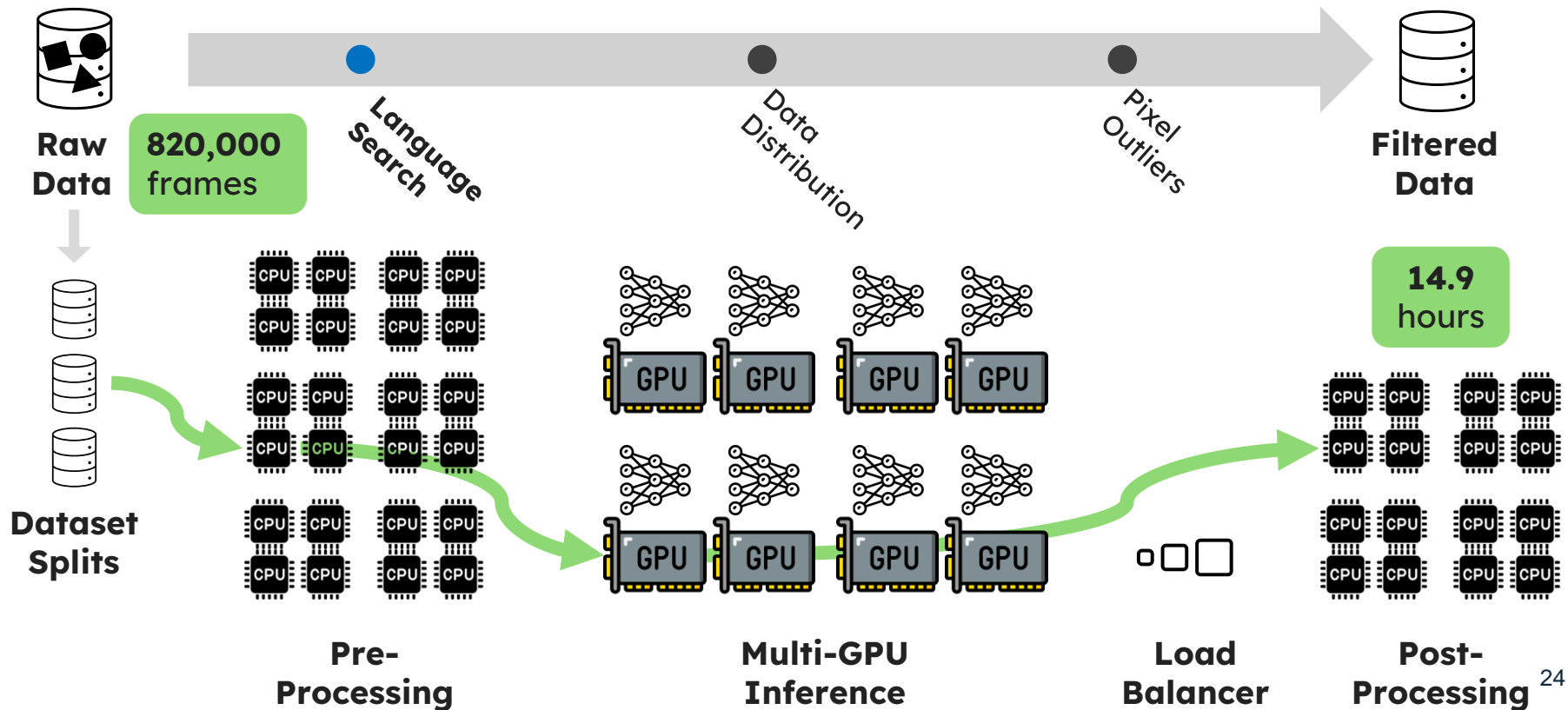
Filtered Data



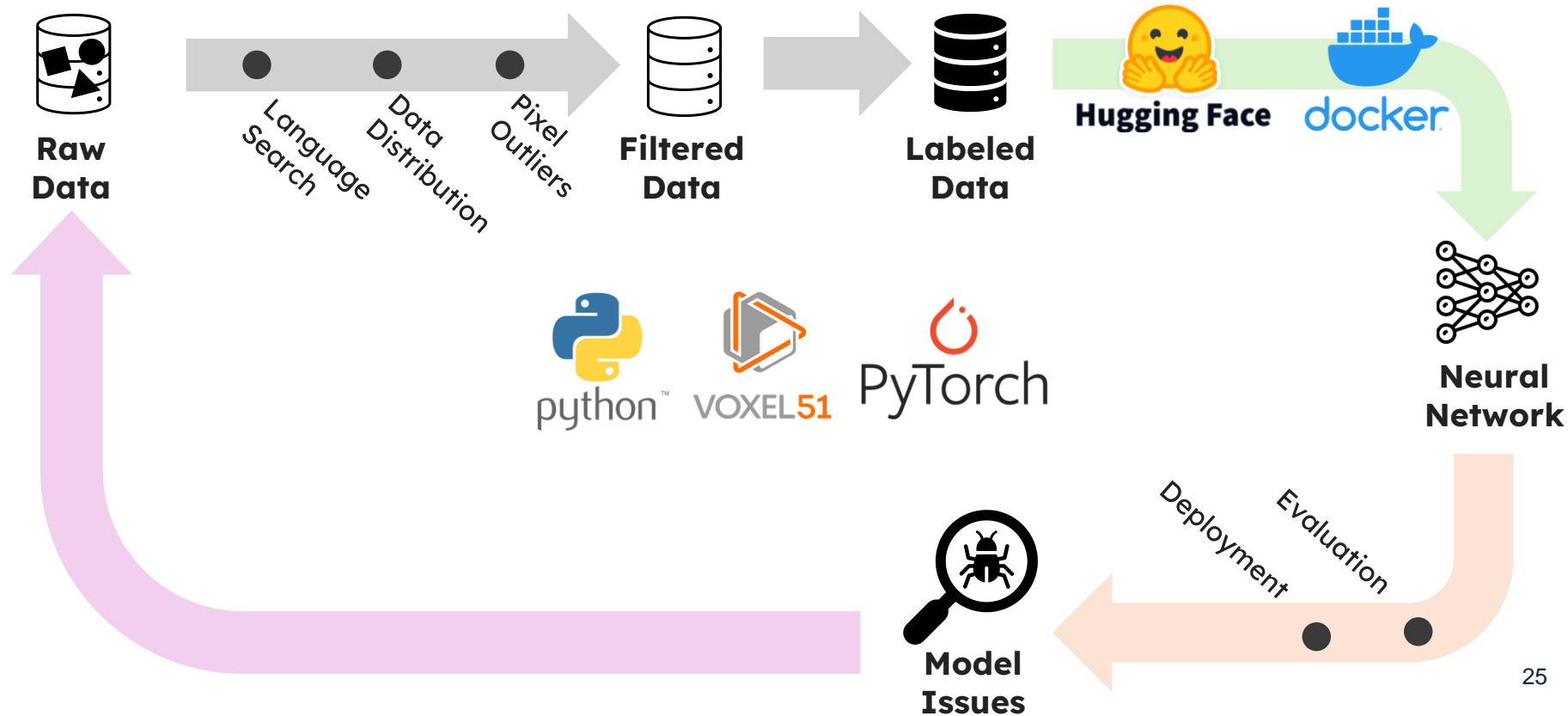
Zero Shot Ensemble

Google OWLv2  
6 frames/second

# Language Search

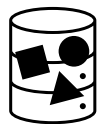


# McCity Data Engine

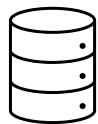
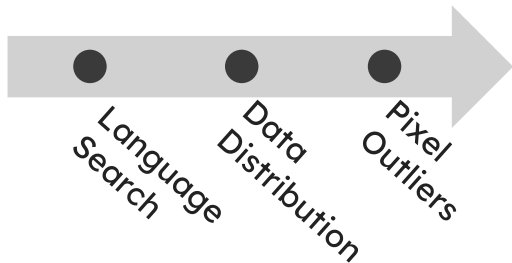




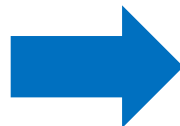
# Data Labeling



Raw Data



Filtered Data



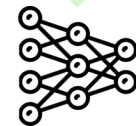
Labeled Data



Hugging Face



docker



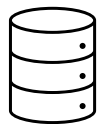
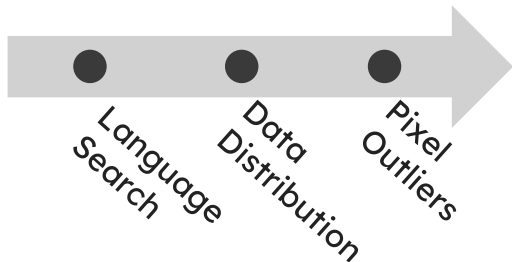
Neural Network



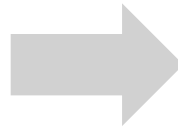
# Model Training



Raw  
Data



Filtered  
Data



Labeled  
Data



Hugging Face

docker

RT  
DETR (24)

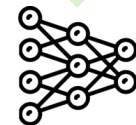
Conditional  
DETR (23)

DETA (22)

Deformable  
DETR (21)

DETR (20)

Yolos (21)

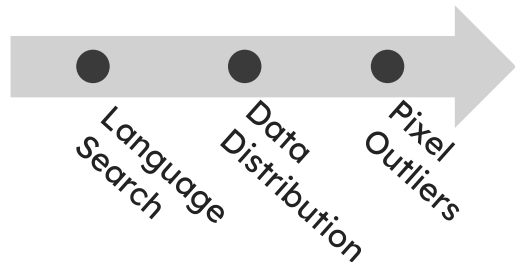


Neural  
Network

# Model Training



Raw Data



Filtered Data



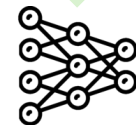
Labeled Data



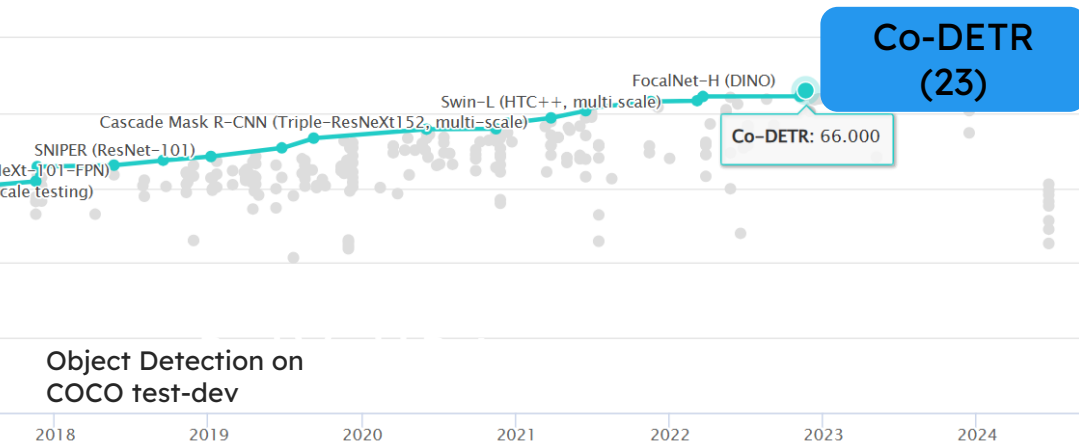
Hugging Face



docker



Neural Network



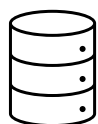
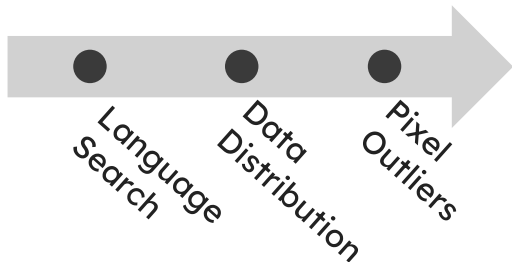
● Other models    ● Models with highest box mAP

- RT DETR (24)
- Conditional DETR (23)
- DETA (22)
- Deformable DETR (21)
- DETR (20)
- Yolos (21)

# Model Training



Raw Data



Filtered Data



Labeled Data



Hugging Face



docker

Co-DETR (23)

RT DETR (24)

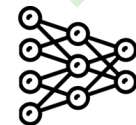
Conditional DETR (23)

DETA (22)

Deformable DETR (21)

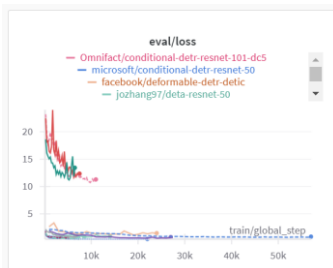
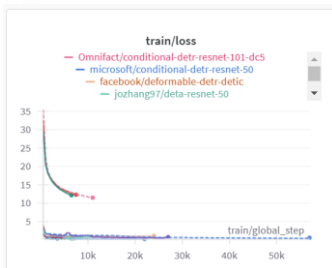
DETR (20)

Yolos (21)

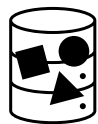


Neural Network

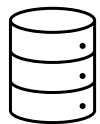
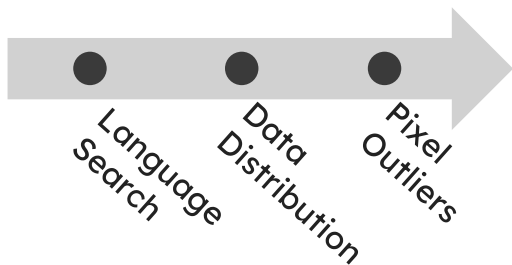
Weights & Biases



# McCity Data Engine



Raw Data



Filtered Data



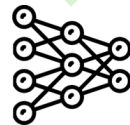
Labeled Data



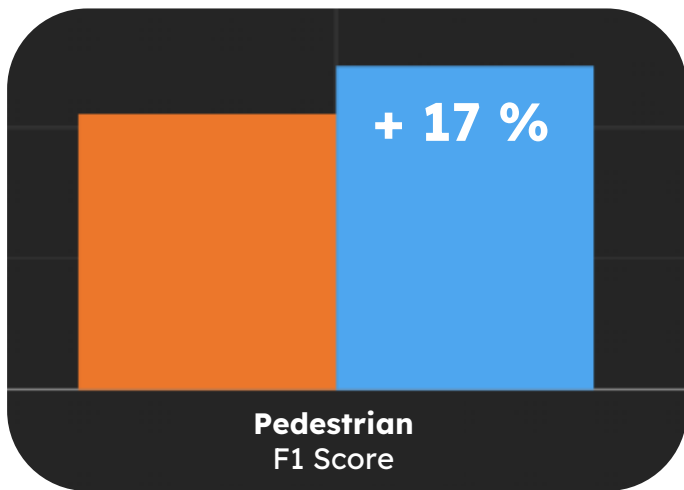
Hugging Face



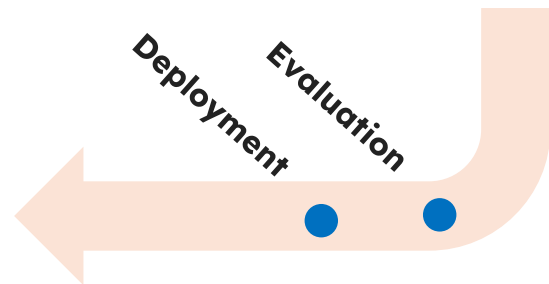
docker



Neural Network

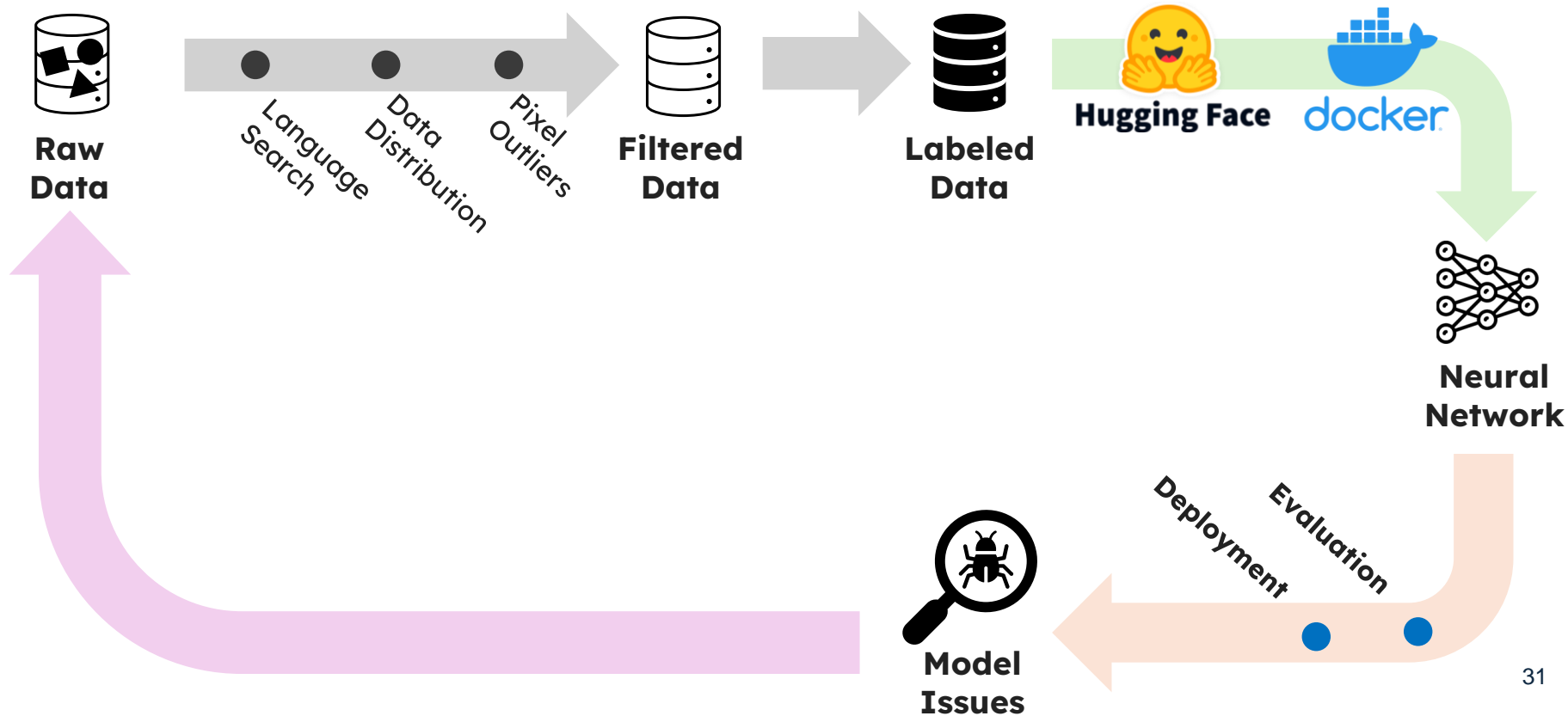


Model Issues

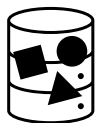




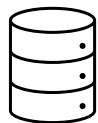
# McCity Data Engine



# McCity Data Engine



Raw  
Data



Filtered  
Data



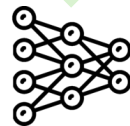
Labeled  
Data



Hugging Face



docker



Neural  
Network

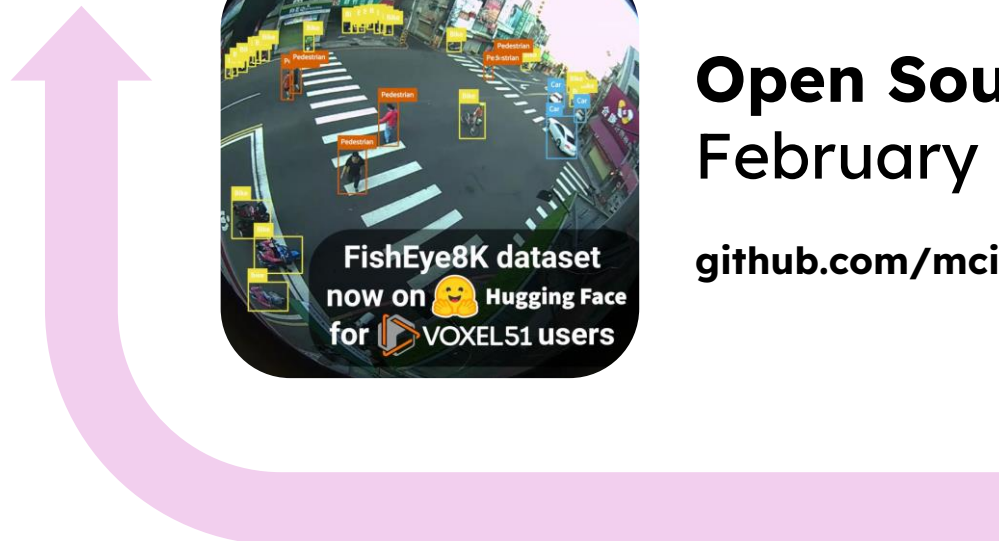
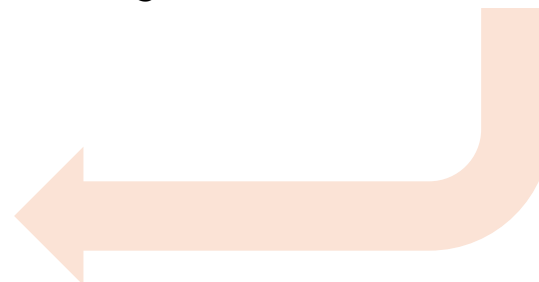


**Open Source Release**  
February 2025

[github.com/mcity/mcity\\_data\\_engine](https://github.com/mcity/mcity_data_engine)



Model  
Issues



# McCity Data Engine

Linked **in**

# Daniel Bogdoll

**Thank you!**  
**Questions?**

