

# Deep Learning for detection on a phone:

*how to stay sane  
and build a pipeline you can trust*

Irina Vidal Migallón  
*Viorama*

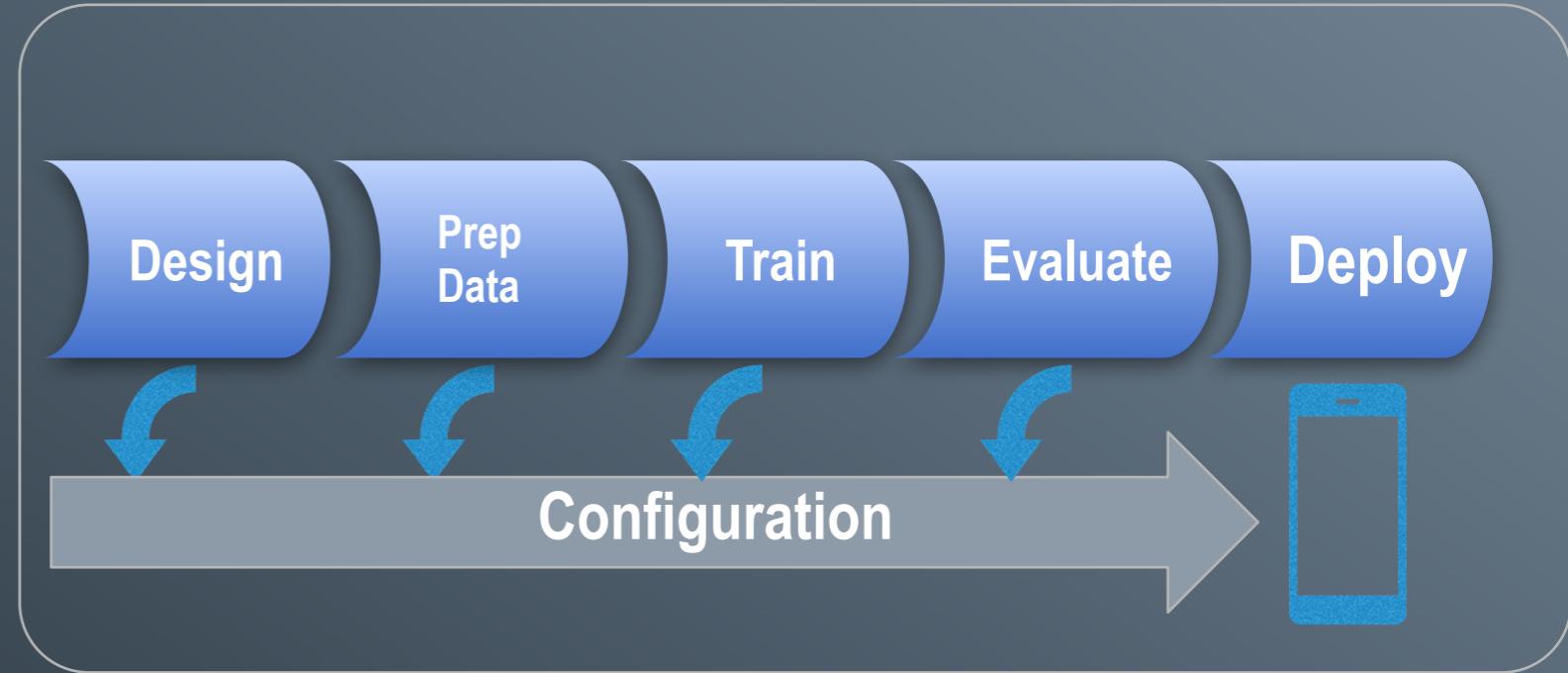




“Work to deploy a Neural Network model on a phone starts **long before** you write the first line of native code.”

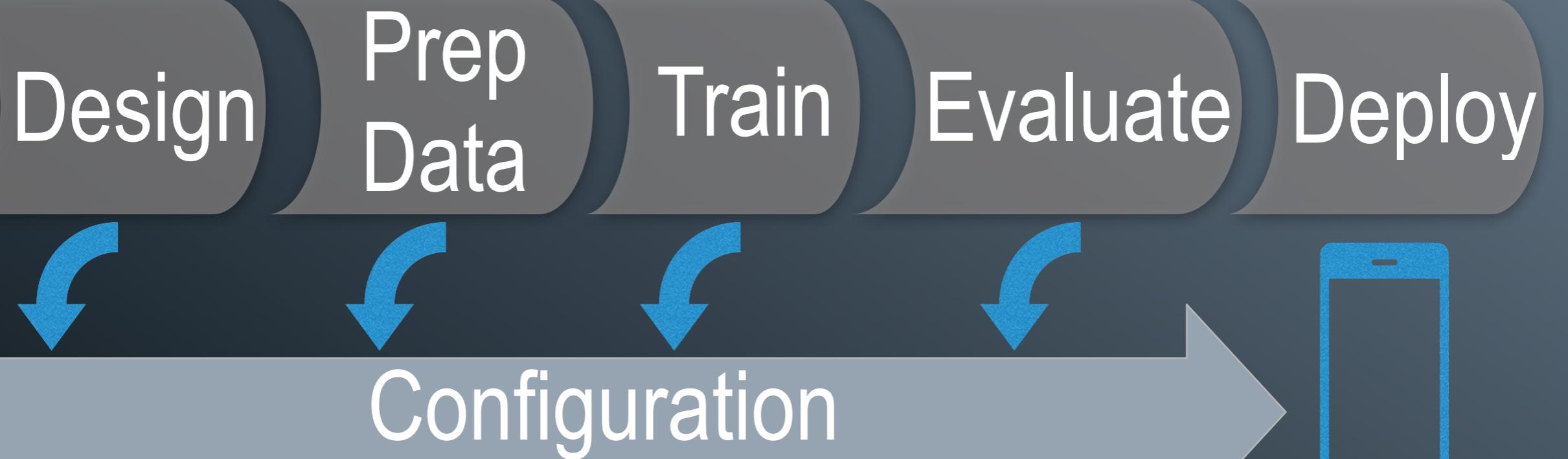
– *Sun Tzu (500 BC, apocryphal)*

Which architecture?  
What weights?  
What deployment parameters?



We'll see:

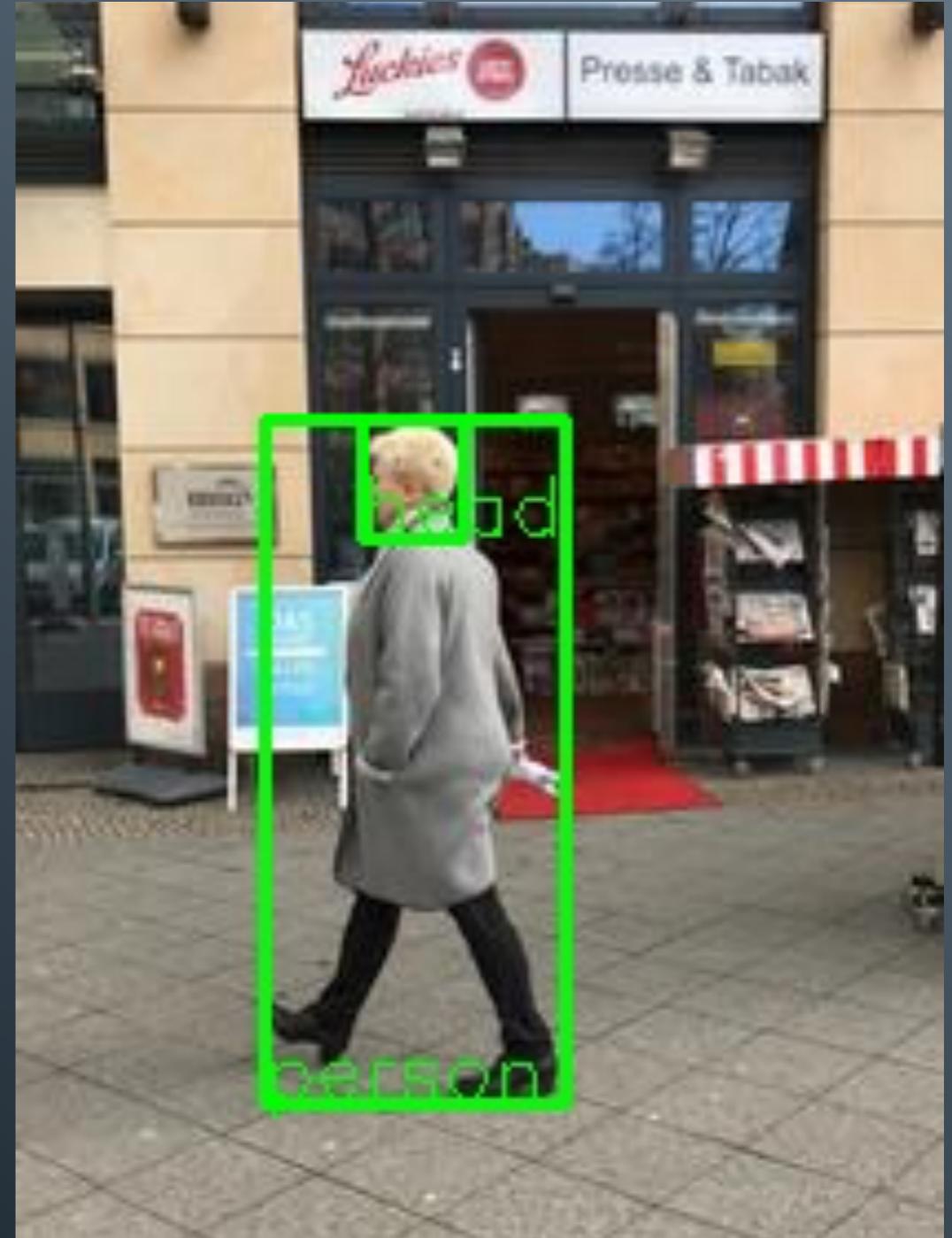
- Pipeline phases
- How to avoid pitfalls



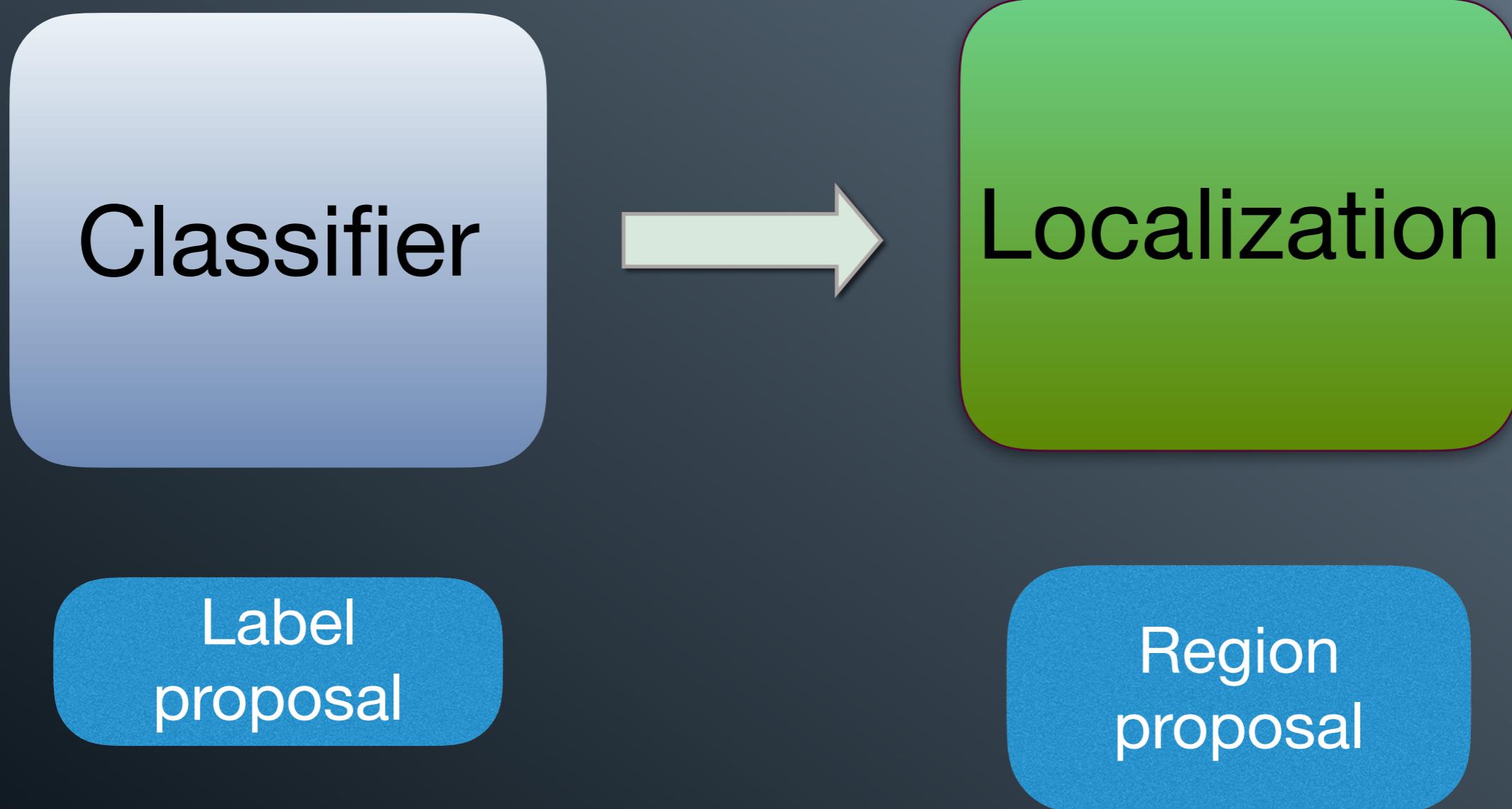
# Classifier



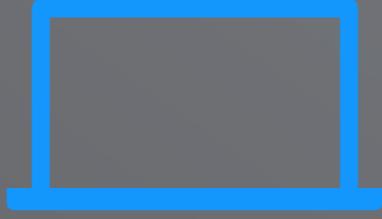
# Detector



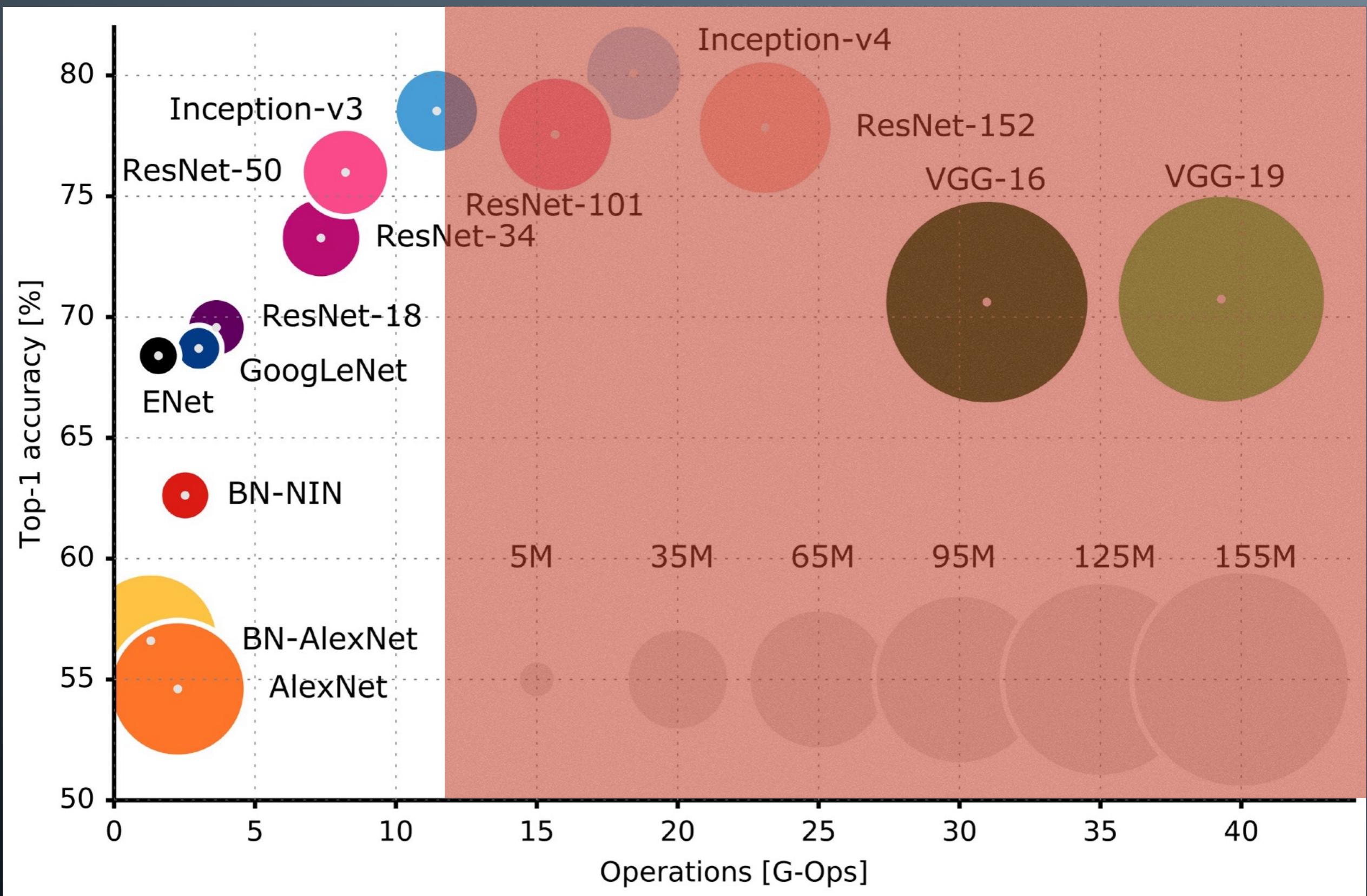
# Detector



# Constraints

		
<b>RAM</b>	<b>1 GB</b>	<b>32 GB</b>
<b>GFLOPS (GPU)</b>	<b>150</b>	<b>4000</b>
<b>Clock (GPU)</b>	<b>600 MHz</b>	<b>1500 MHz</b>





(A. Canziani, A. Paszke, E. Culurciello, 2016)

# Slimmer classifier



- ▶ Memory
  - ▶ Convolutional vs Fully Connected
- ▶ Inference time
  - ▶ Input & layer sizes

# Slimmer detector

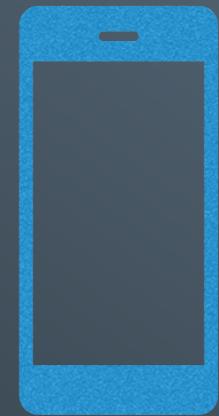


- ▶ What determines the size of your detector?
- ▶ Keep only what you need

Design      Prep Data      Train      Evaluate      Deploy

Arch

Configuration



Data



- ✓ Intended use & format
- ✓ Annotation format
- ✓ Consistency

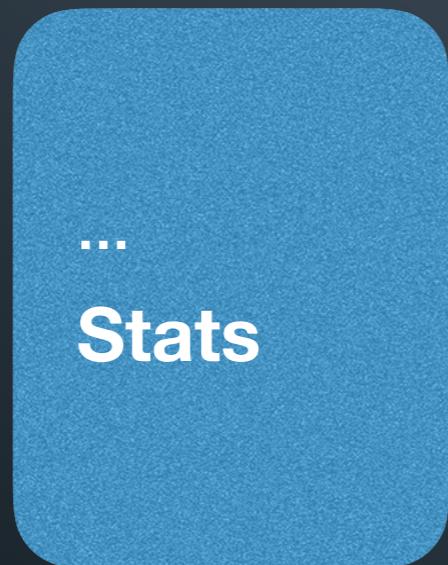
Free\*

- MS COCO
- PASCAL VOC
- KITTI

\* Pay with your time instead.

Explore your data

Analyse your data:

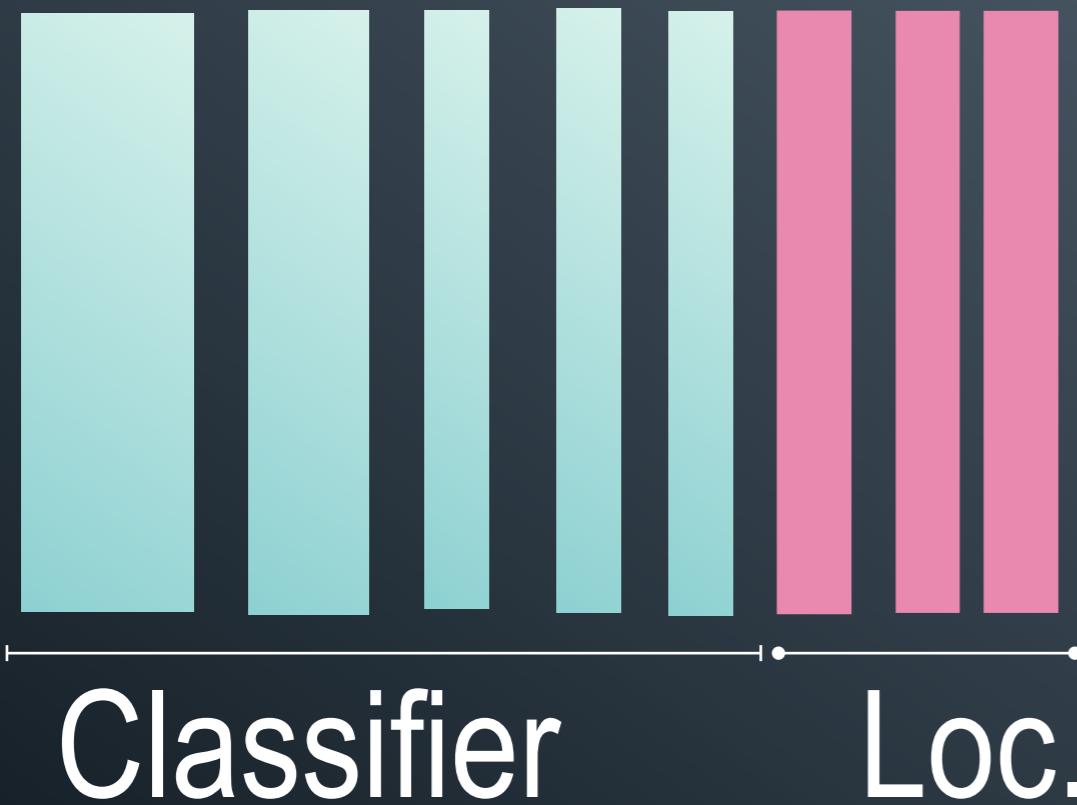


Track:

- Subsets
- Augmentations



# Training

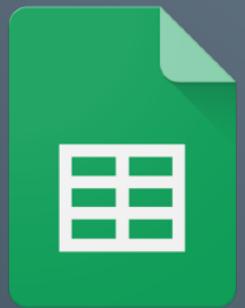


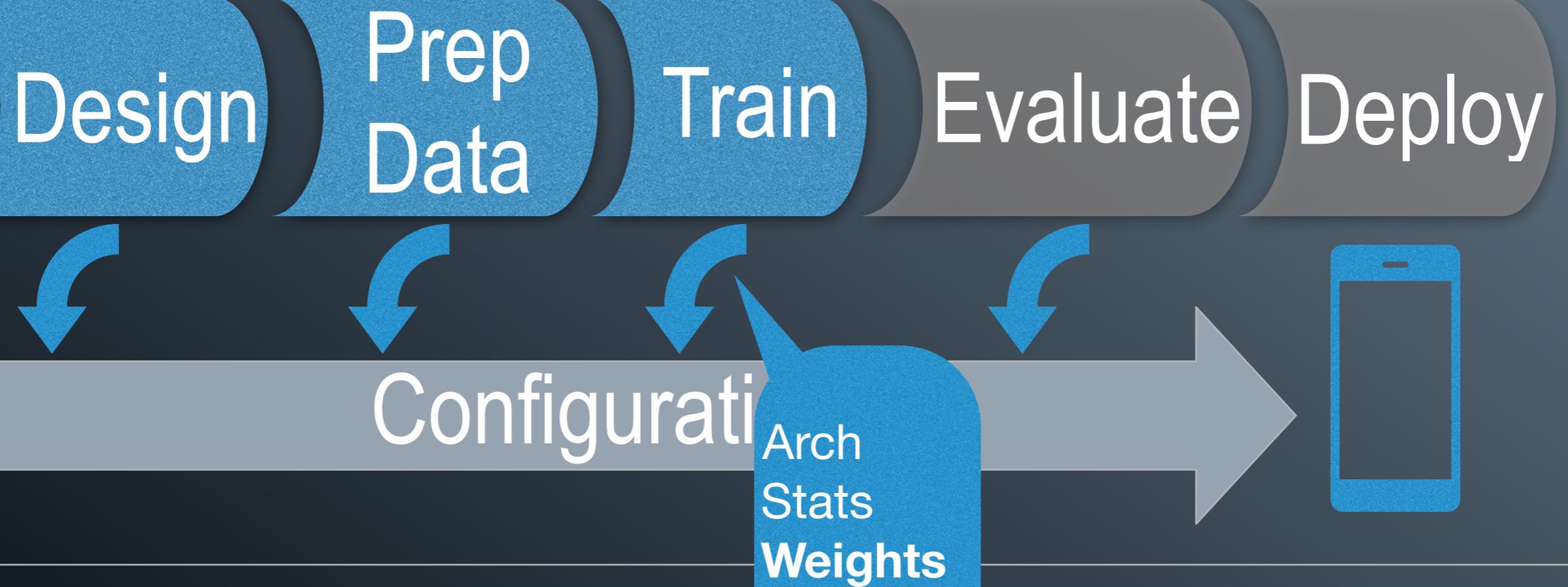
- ▶ Pre-trained models
- ▶ Transfer learning

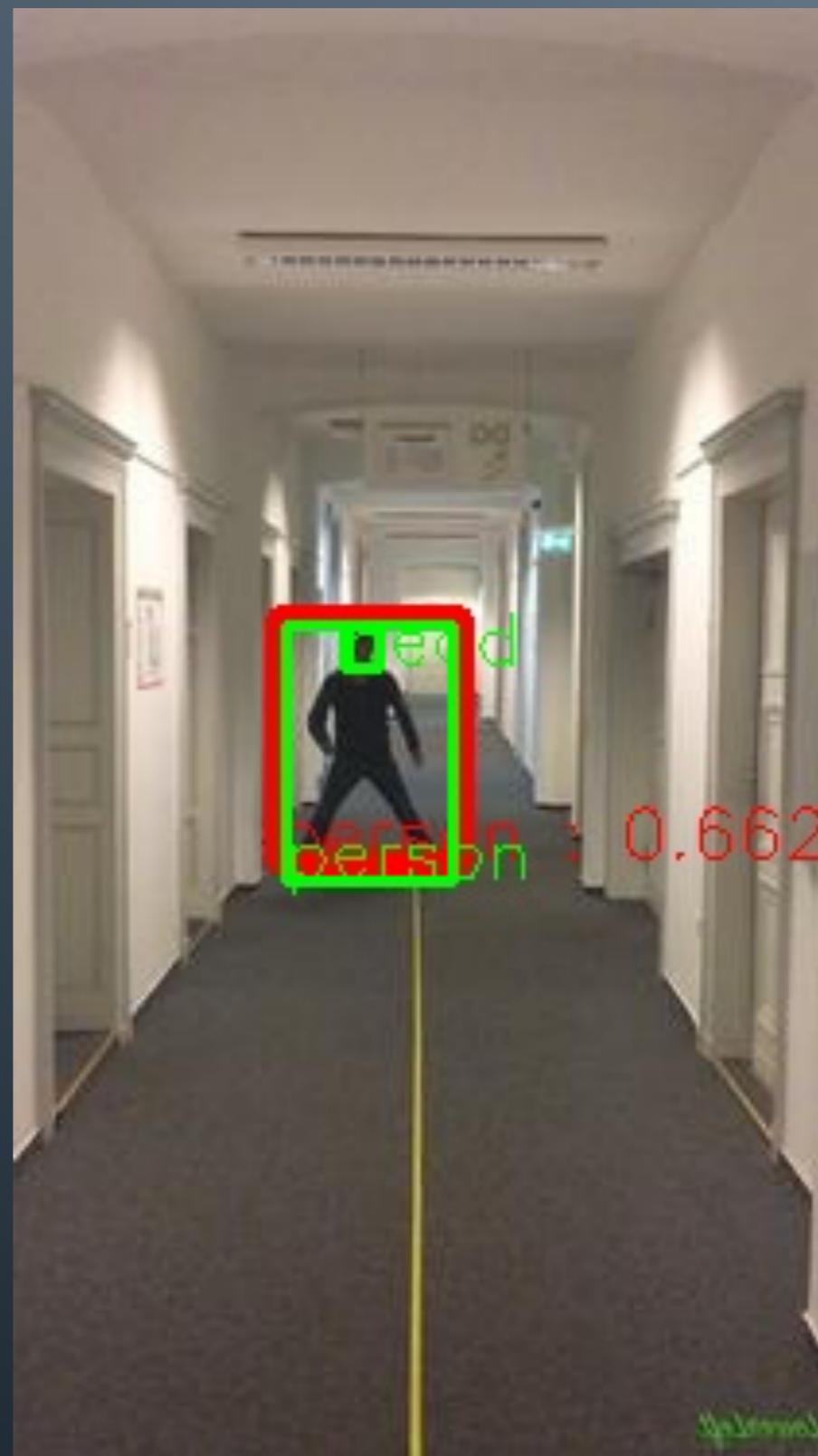
- Learning
- Frozen

# Training

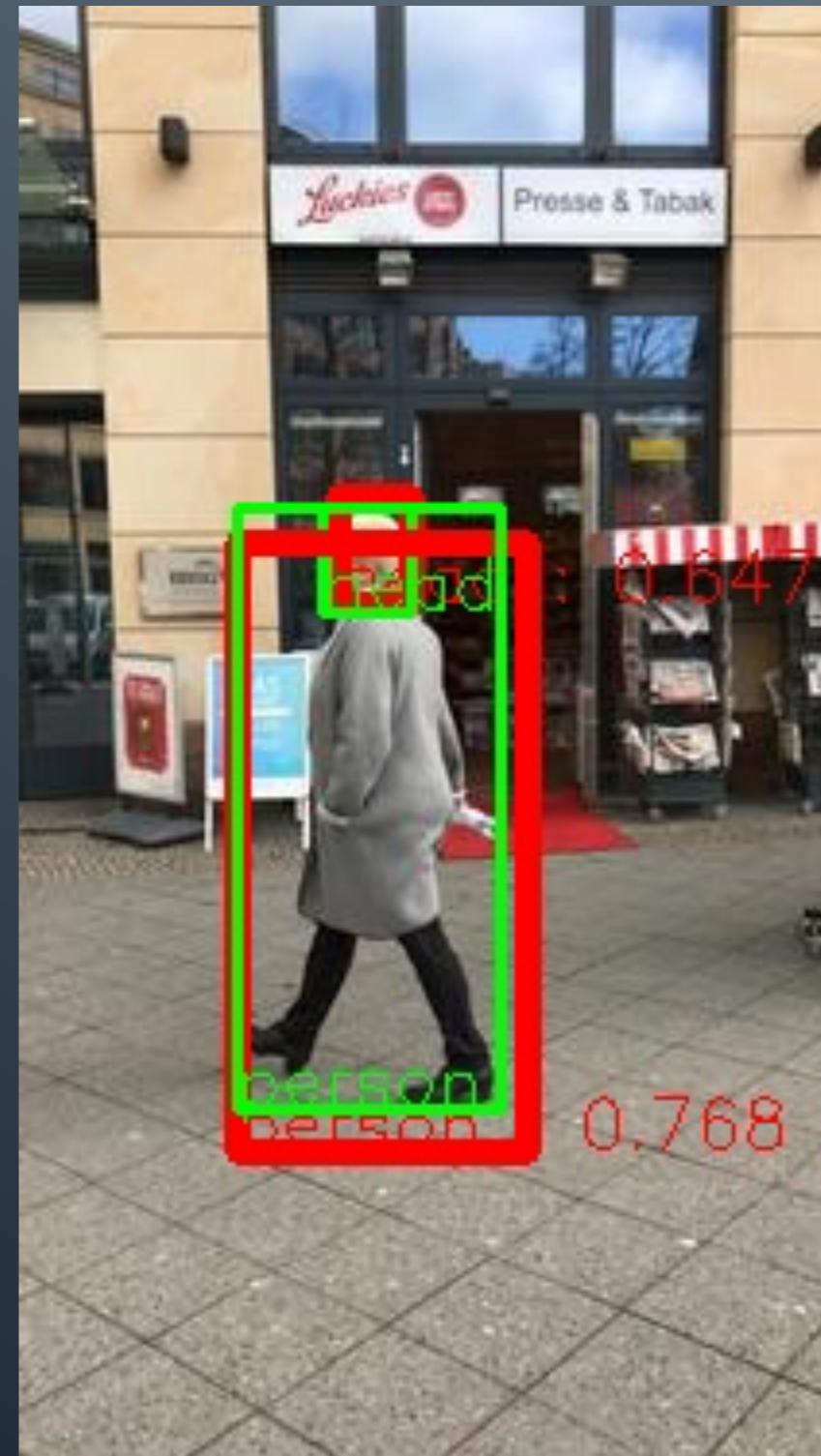
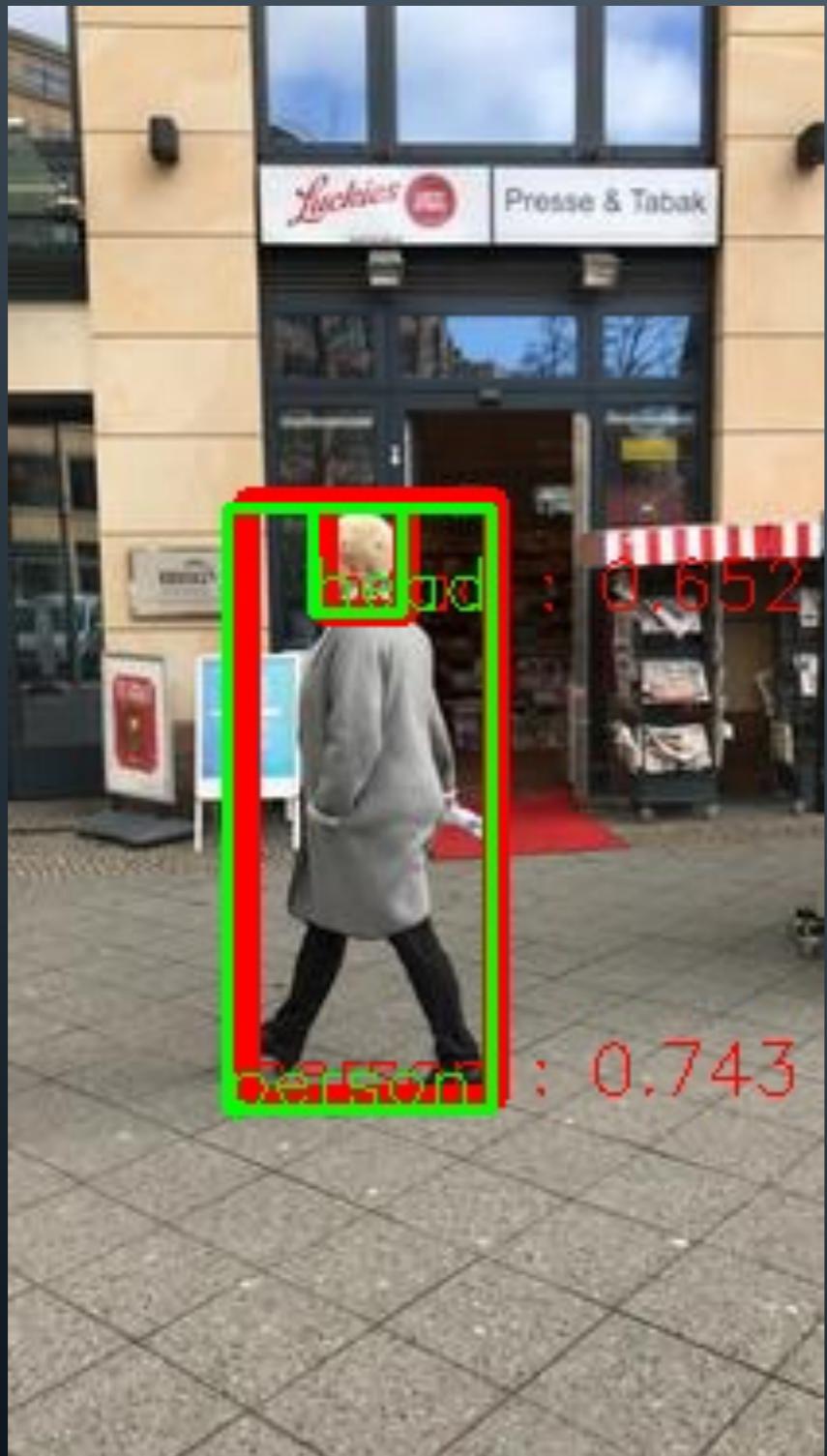
- ▶ Track, track, track
- ▶ Any tool - just do it.
  - ▶💡 you can feed CSVs to any other tool!





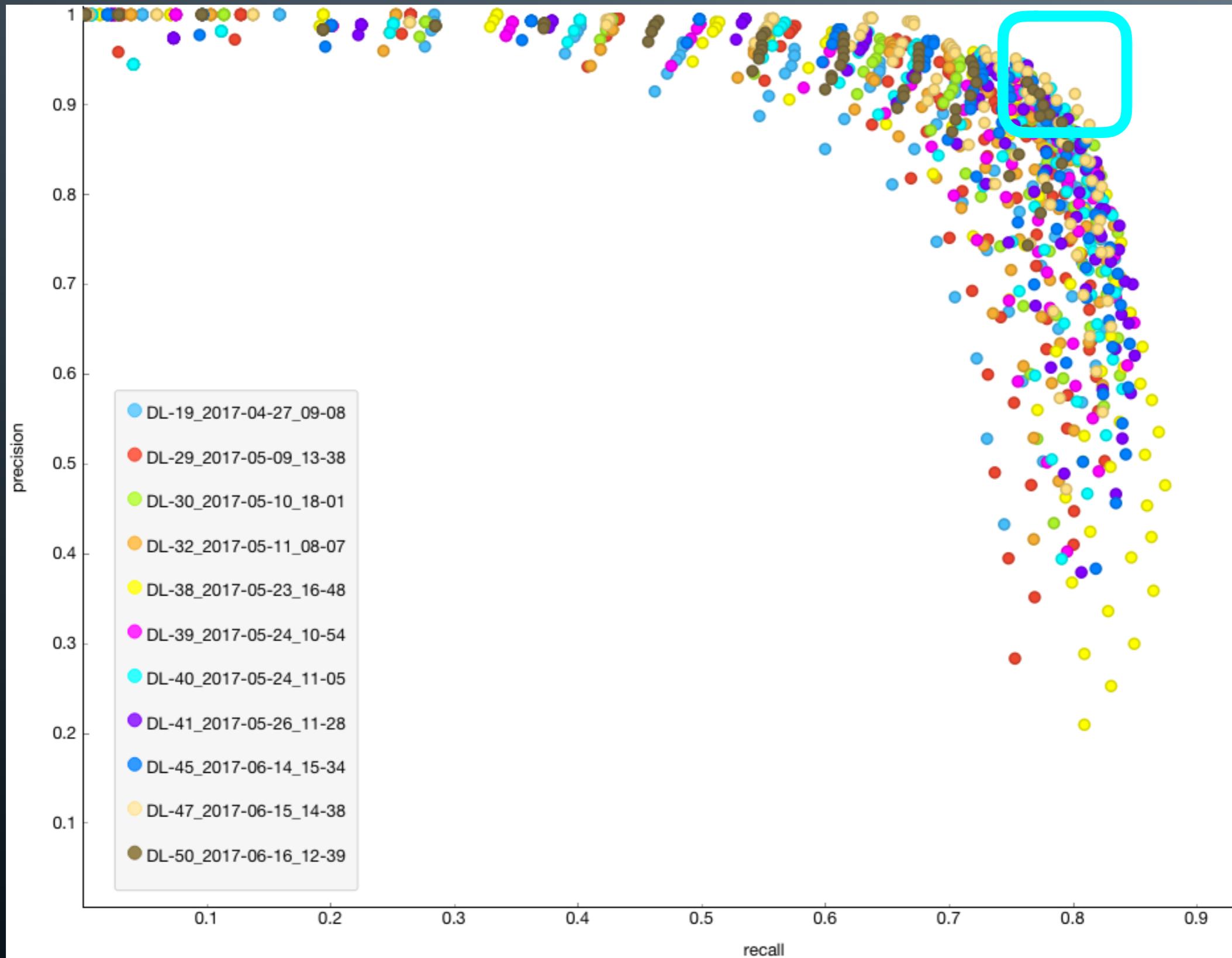


Detection  
Ground truth



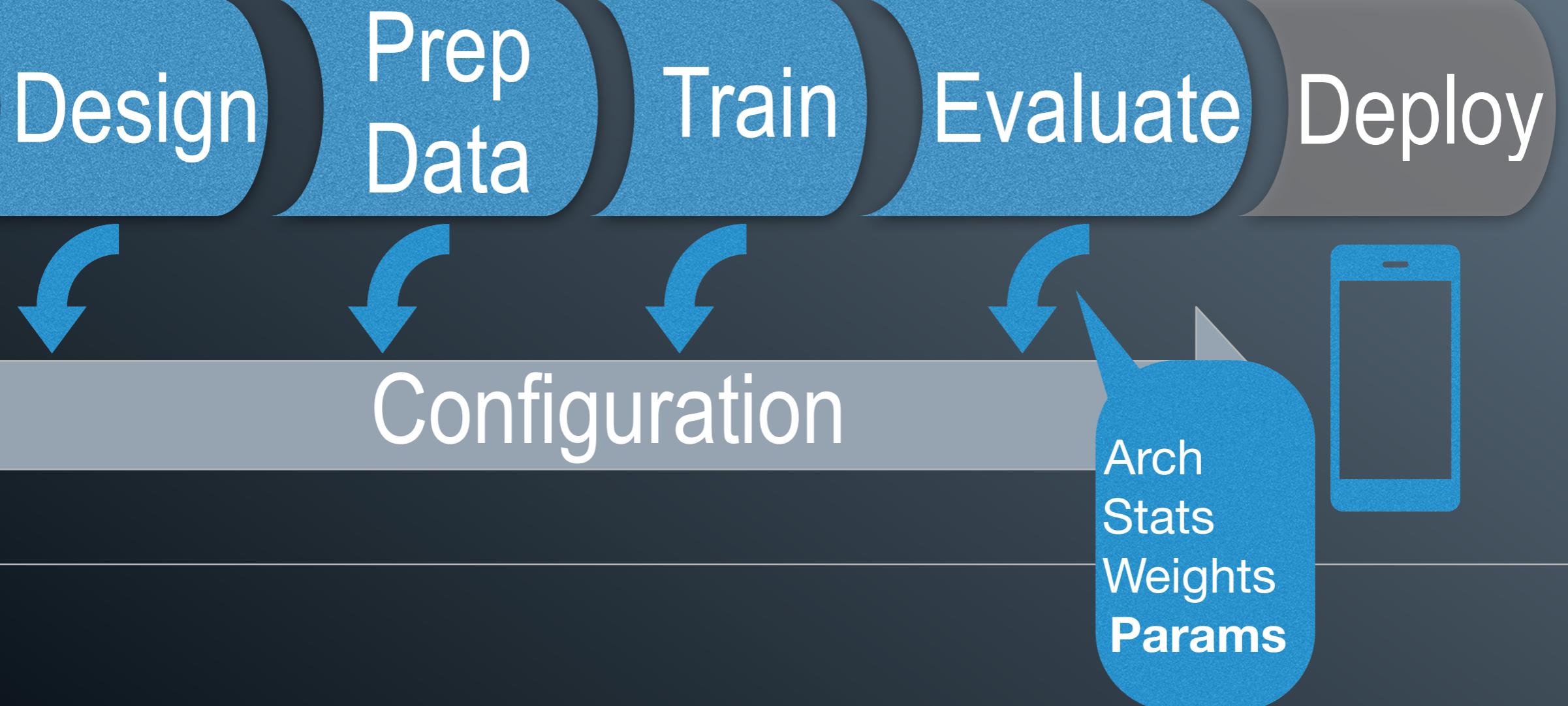
■ Detection  
■ Ground truth

# Evaluate



# Evaluate

- ▶ Architecture
- ▶ Datasets used
- ▶ Augmentations
- ▶ Hyperparameters
- ▶ Deployment parameters



Design      Prep Data      Train      Evaluate      Deploy



Configuration



Arch  
Stats  
Weights  
**Params**

# Deploy



TensorFlow



Caffe2

- ▶ Frameworks announced left & right
- ▶ Check for ops available





# *Avoiding Pitfalls*



*Facepalm Inception by [zantaff.deviantart.com](https://zantaff.deviantart.com)*

# Avoiding Pitfalls

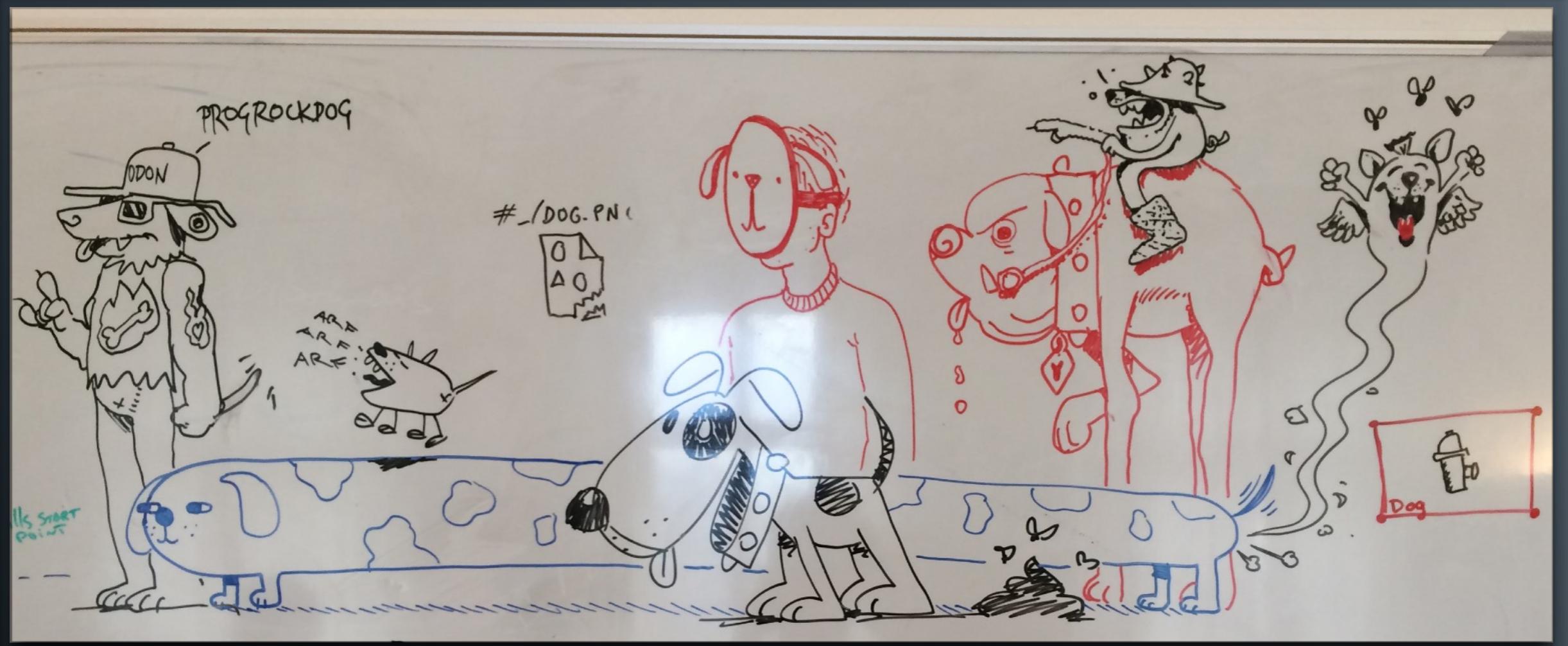
- ✓ Datasets:
- ✓ Consistent annotations
- ✓ Define detection
- ✓ Define box
- ✓ Define person

# Avoiding Pitfalls

- ✓ Evaluation: visual and metric
- ✓ Automatic configuration
- ✓ Tracking experiments
- ✓ Early profiling on device
- Feedback loop to design

# Summary

- ▶ Tips for models targeting mobile devices
- ▶ Configuration-driven pipeline
- ▶ Avoiding pitfalls!



Thank you!  
Any questions?