# Embedded Machine Learning Lab

Max Schik, Darius Schefer | March 20, 2024

# Contents

Data

Adaptation and optimization
oooo

Integration and pipeline
oo

Live demo
o

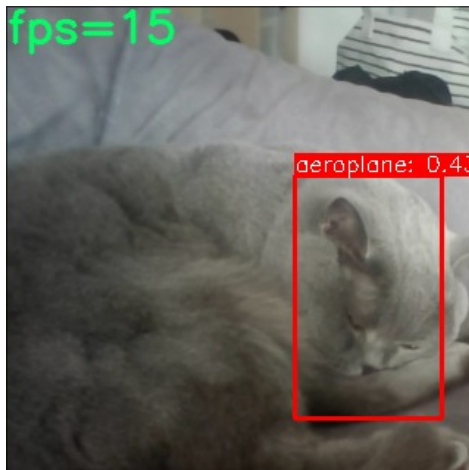**2**/11    03/20/2024    Max Schik, Darius Schefer: EML-lab                                                                                                    EML-lab

# Data

- More data
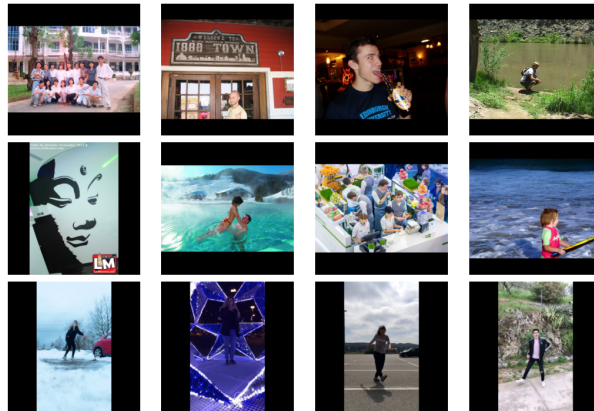  - Human Dataset [a] (17,300 images)
  - Tiktok Dancing [b] (2615 images)
- Data augmentation
  - Albumentations[c] library
  - Rotation, flipping, contrast

---

[a]https://www.kaggle.com/datasets/fareselmenshawii/human-dataset
[b]https://www.kaggle.com/datasets/tapakah68/segmentation-full-body-tiktok-dancing-dataset
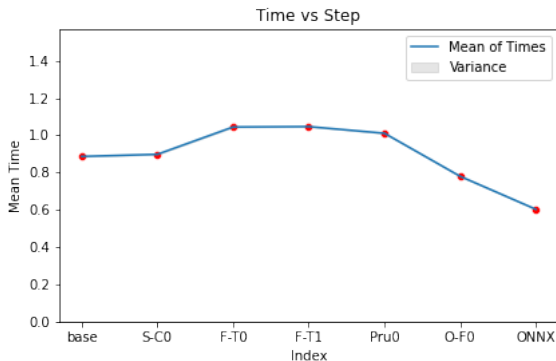[c]https://albumentations.ai/



Data
●

Adaptation and optimization
○○○○

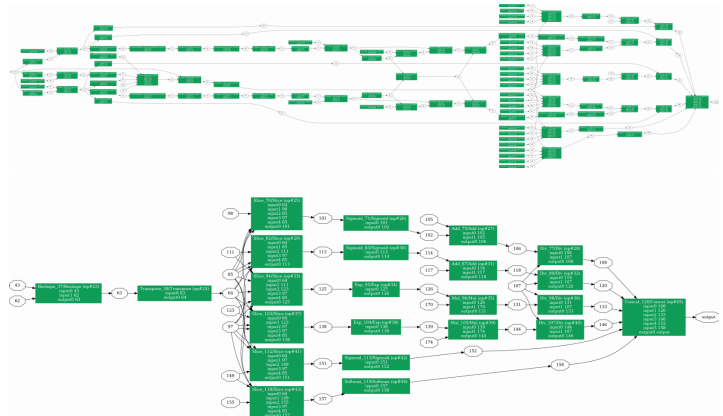Integration and pipeline
○○

Live demo
○

# Adagtation and optimization

- Person-only-detection, fine-tuning
- Iterative pruning
- Batch norm optimization
- Inference

Data
○

Adaptation and optimization
●○○○

Integration and pipeline
○○

Live demo
○

**5/11**   03/20/2024   Max Schik, Darius Schefer: EML-lab                                                    EML-lab
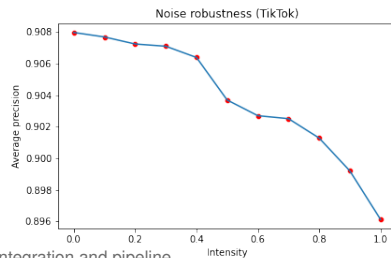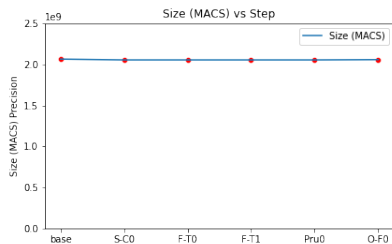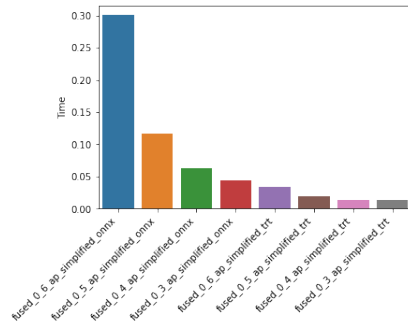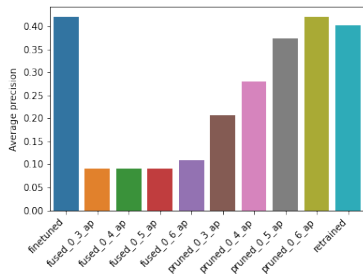
# Inference

- ONNX
- TensorRT
  - TensorRT supports only a subset of ONNX spec
  - `ReflectivePad` → `ConstantPad`
  - Simplification of graph with `onnx-simplifier`[a]

---

[a]https://github.com/daquexian/onnx-simplifier

Size (MACS) vs Step



Noise robustness (TikTok)



Data
○

Adaptation and optimization
○○●○

Integration and pipeline
○○

Live demo
○

# Detections

# Integration and pipeline

- Configuration
  - File containing steps and parameters
  - Reproducibility
  - Version control

```yaml
1  --- !experiment
2  start_weights_path: "./weights/voc_pretrained.pt"
3  augment: True
4
5  steps:
6  - !strip_classes
7      finetune: true
8      finetune_epochs: 15
9  - !pruning
10     target_acc: 0.3
11     prune_ratio: 0.05
12     batch_size: 64
13     num_train_epochs: 10
14     num_eval_batches: 10
15 - !operator_fusion {}
16
```

Data
○

Adaptation and optimization
○○○○

**Integration and pipeline**
●○

Live demo
○

# Pipeline

- jetcam
- inference using framework
    - separate thread for camera creates problem for tensorrt
- faster NMS by using torch implementation (30% faster than default implementation)

# Live Demo

Data

Adaptation and optimization
○○○○

Integration and pipeline
○○

Live demo
●

**11**/11   03/20/2024    Max Schik, Darius Schefer: EML-lab                                    EML-lab