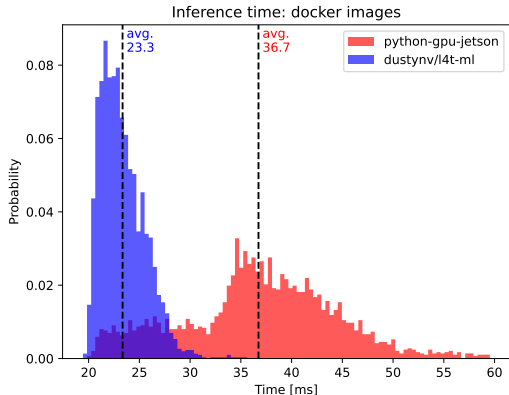# Live camera flow

**Computation steps**

1. capture image with camera (disabled JPEG encode for Jupyter)
2. run model on image
3. filter anchors, apply non-maximum suppression
4. draw bounding boxes and fps
5. compress image to JPEG (OpenCV should use libjpeg-turbo already)
6. show image to user (MJPEG stream)

**Jupyter notebook unreliable**: kernel crashes, unresponsive, . . .
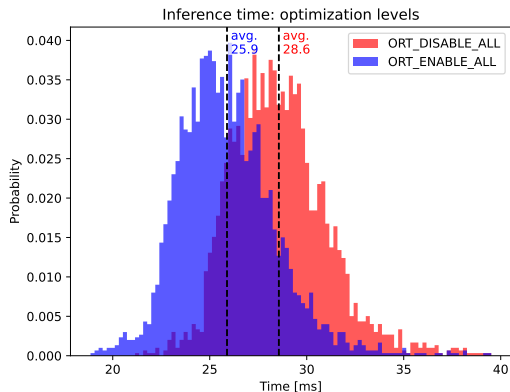$\rightarrow$ **run on a custom web server inside docker**

# Docker container

- provided image outdated (2021)
- use "l4t-ml" from dustynv (11/2023)
  `https://github.com/dusty-nv/jetson-containers`
- chose matching NVIDIA JetPack/L4T
- newer libraries, **faster inference**
- unchanged performance on other parts

- simple web server with Flask
- low resource usage
- has endpoint serving MJPEG stream
- UI to start/stop camera
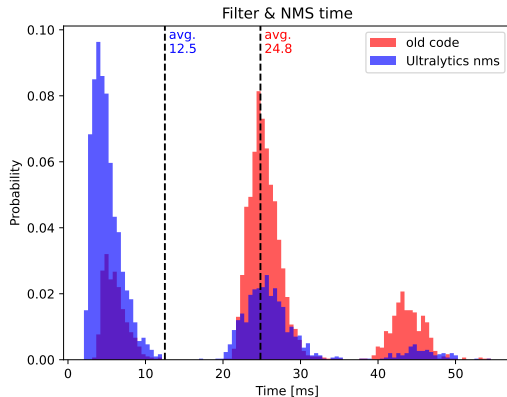


Inference time: docker images

# ONNX runtime

- use for inference on GPU
- supposedly faster than PyTorch
- represents model computations as graph
- use optimization `ORT_ENABLE_ALL`
  $\rightarrow$ fuses conv and batch-norm layers

- Quantization
  - int8 $\rightarrow$ not supported by GPU
  - float16 $\rightarrow$ no speed up

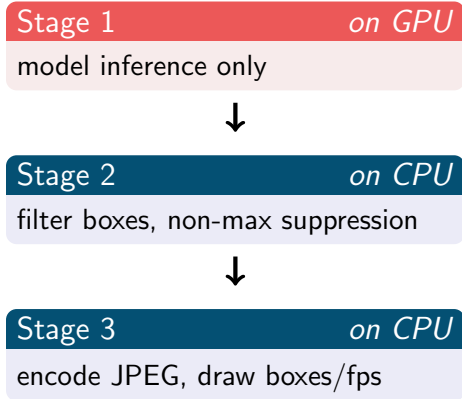

Inference time: optimization levels

# Filter boxes and NMS

- original code: slow, pure Python impl.
- use `non_max_suppression()` from Ultralytics library

- licensed as AGPL v3
- everything combined in one function
- based on native `torchvision.ops.nms()`

# Pipelining

- split work between GPU and CPU
- parallelize steps
- as Python `threading.Thread`
  $\rightarrow$ GIL mostly no issue
- connected by queues, length limited

- **Input**: image from camera callback
- **Output**: byte stream for HTTP response

| Stage 1 | *on GPU* |
|---|---|
| model inference only | |

$\downarrow$

| Stage 2 | *on CPU* |
|---|---|
| filter boxes, non-max suppression | |

$\downarrow$

| Stage 3 | *on CPU* |
|---|---|
| encode JPEG, draw boxes/fps | |

# Pipelining: results