



How Do We Enable Edge ML Everywhere? Data, Reliability, and Silicon Flexibility

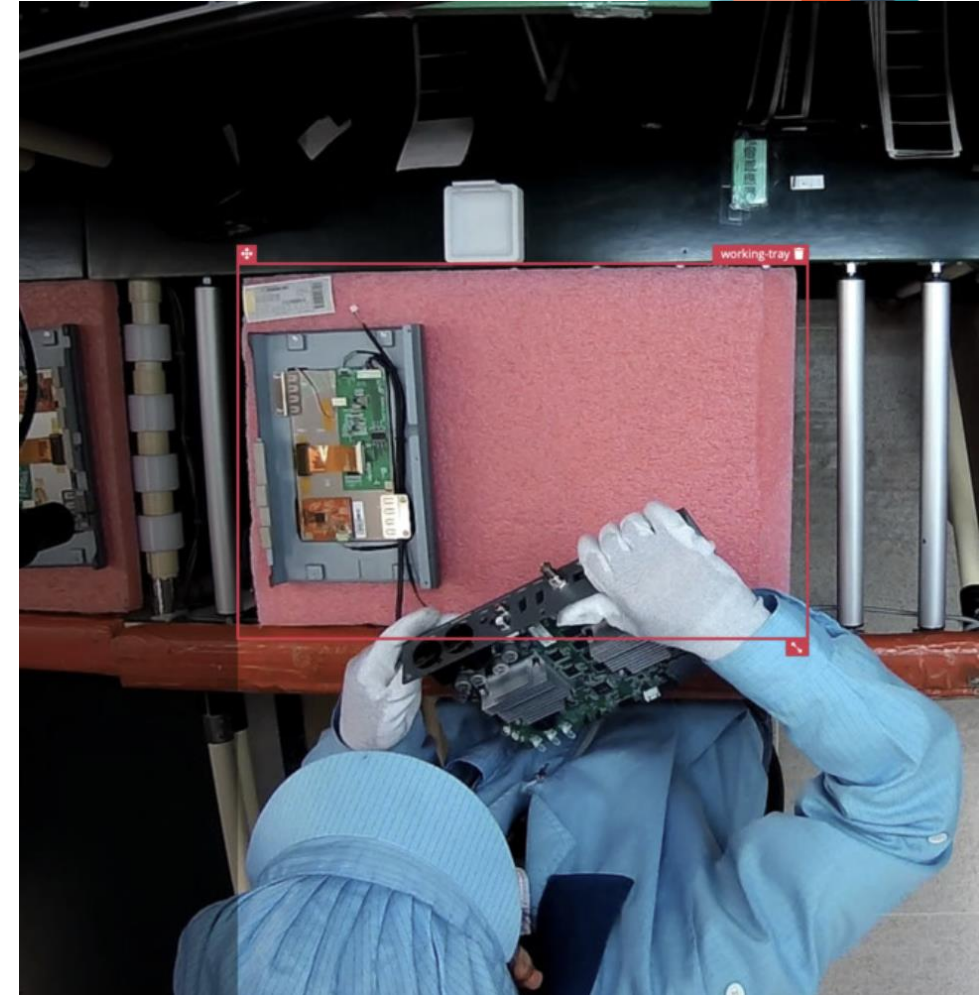
Zach Shelby
Co-founder and CEO
Edge Impulse

Advantech increases manufacturing productivity by 15%



Visual inspection system to monitor worker safety and flag delays on the production line in real-time.

- A reported 15% overall increase in production line efficiency
- Faster detection of idle time raises assembly-line productivity
- Managers free up time to focus on production planning and operations



State of the edge ML 2022

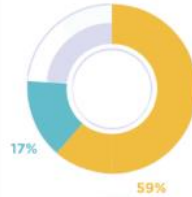


User population

1K+ People polled
91 Countries represented



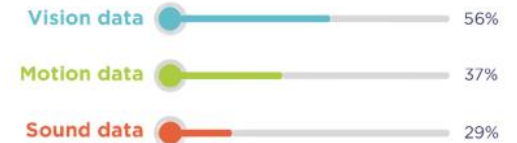
Frameworks



59% of participants with ML experience use **TensorFlow** and **17%** use **PyTorch** as their ML framework.

Top 3 Data Types

Participants with ML Experience



Professional Experience

37% of survey participants have 1+ years of professional ML experience.



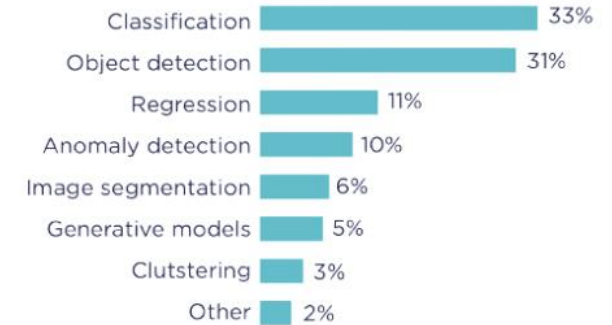
53% of survey participants have 1+ years of professional embedded systems experience.



Top 5 Favourite Boards for Edge ML Projects

- 1 Raspberry Pi 4
- 2 Arduino Nano 33 BLE Sense
- 3 ESP32
- 4 Raspberry Pi Pico
- 5 NVIDIA Jetson Nano

Common ML Tasks



Dataset Source

Participants with ML Experience



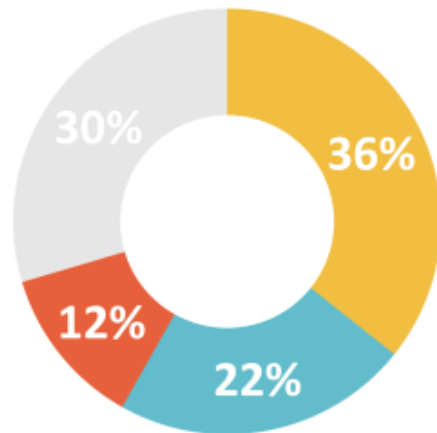
What are the barriers to edge ML scale?



Biggest Concerns

Users with ML experience

- Model accuracy
- Not enough data
- Model size
- Other



Challenges with Data

Users with ML experience

- 1 Extracting meaningful features from data
- 2 Cleaning and preparing data
- 3 Obtaining data
- 4 Understanding what data to focus on
- 5 Managing large datasets

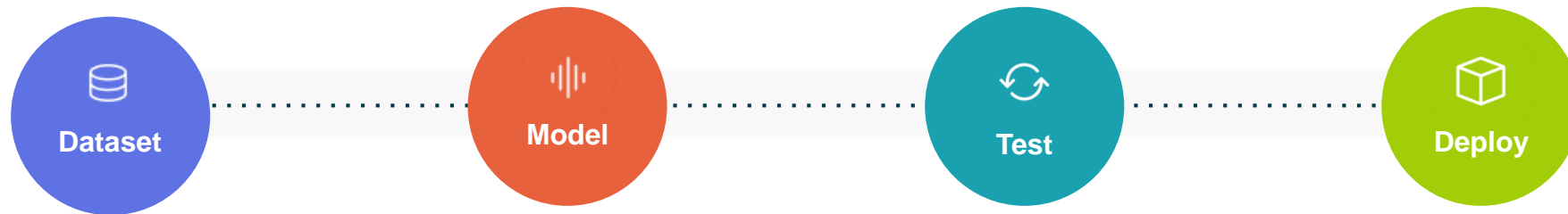
Top 5 Challenges

Preventing Users from Getting Involved in Edge ML

- 1 Lack of time
- 2 Lack of expertise in Data Science
- 3 Lack of resources
- 4 Difficulty in collecting the right data
- 5 Challenge with developing models optimized for the edge

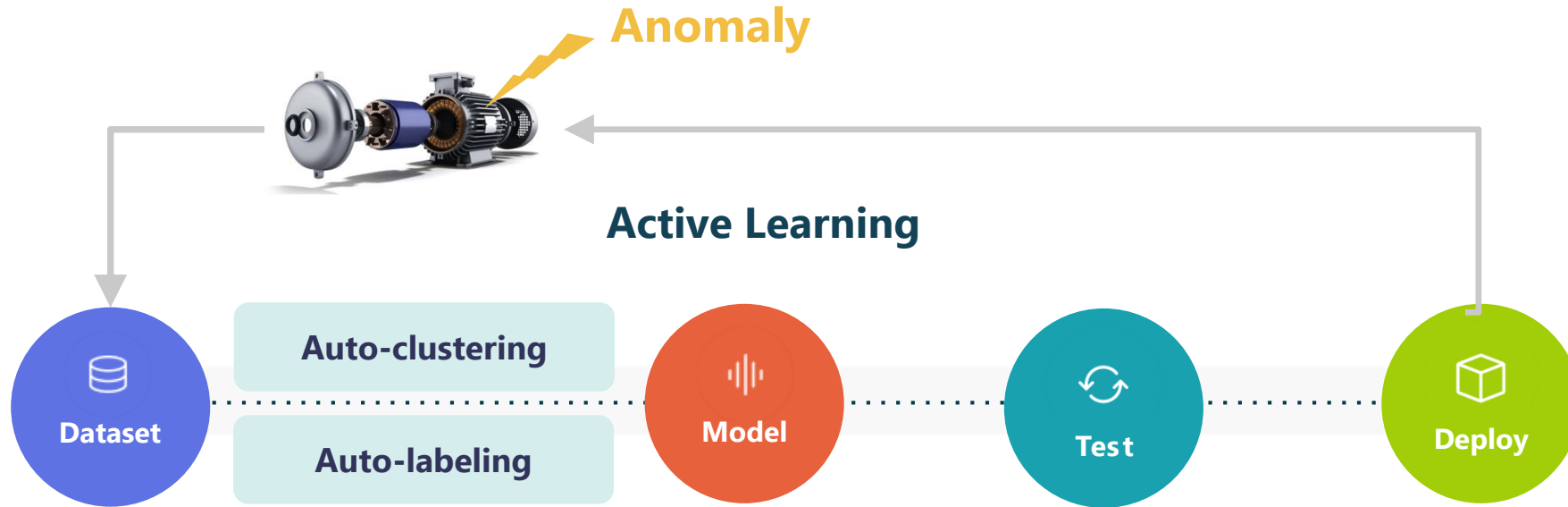
Unleashing the data

Big data, big model is a problem



- State of the art: Monolithic batch, from big data to big model
- Developed for ML research, unsuitable for most edge ML applications

Data-Centric ML – from zero to hero



- Auto-clustering, using feature analysis to help experts quickly label data
- Auto-labeling, using pre-trained expert models to suggest labels
- Active learning, using inference to drive the data collection process

Auto-clustering in action

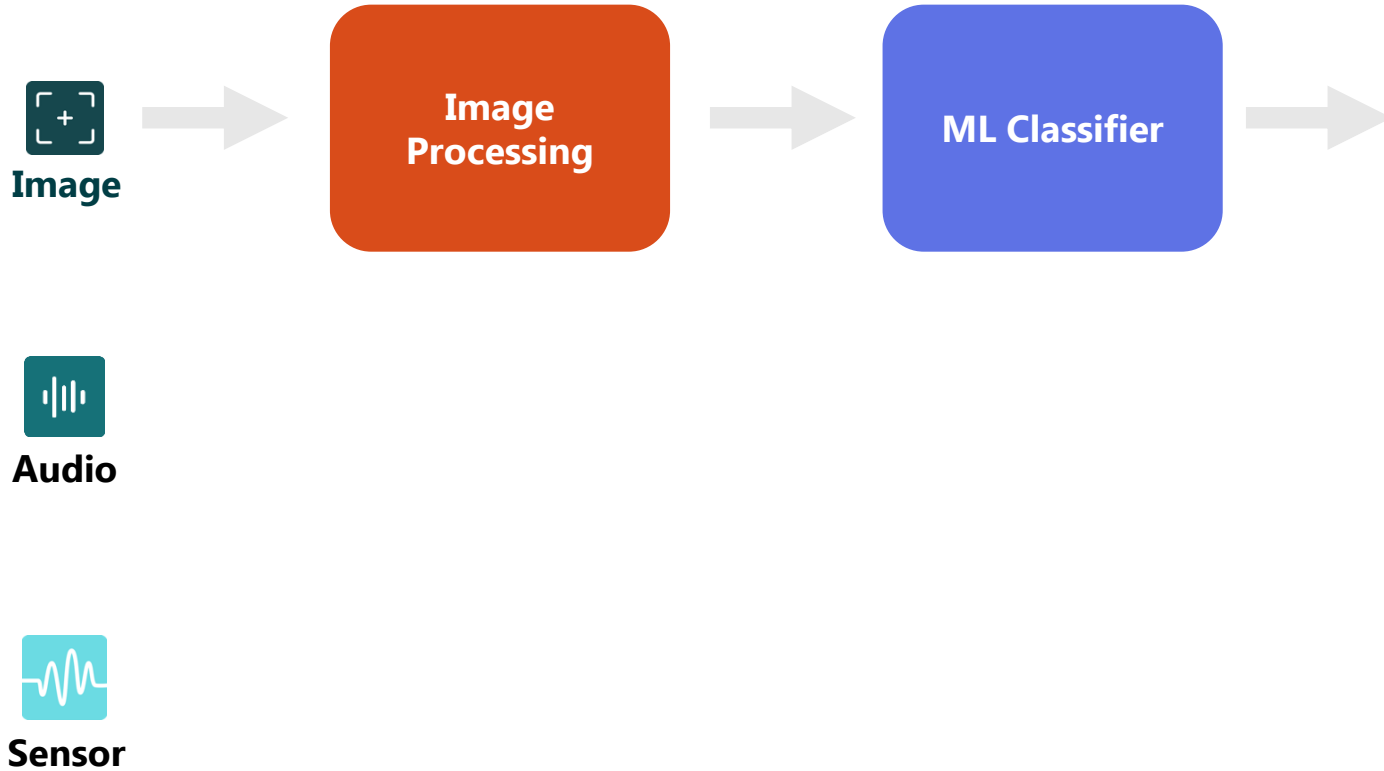


- Dashboard
 - Devices
 - Data acquisition
 - Create impulse
 - Spectral features
 - NN Classifier
 - Anomaly detection
 - EON Tuner
 - Retrain model
 - Live classification
 - Model testing
 - Versioning
 - Deployment
- GETTING STARTED
- Documentation
 - Forums



Making ML industrial grade

Sensor fusion for industrial-grade ML



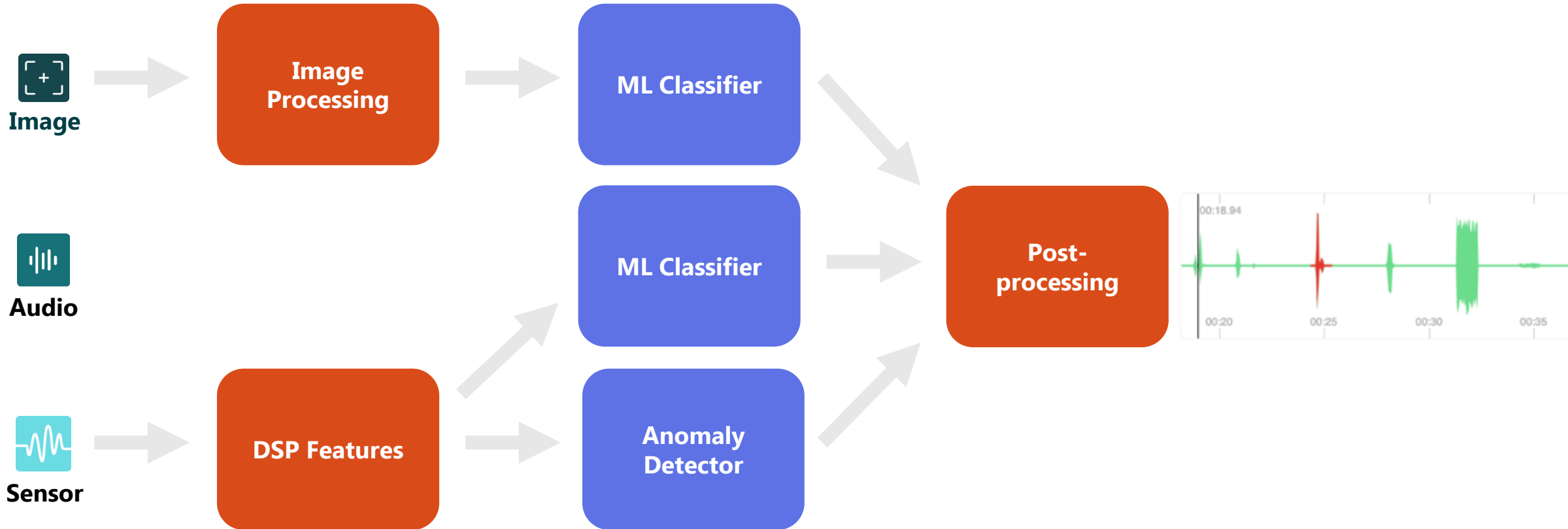
ACCURACY
97.6%

LOSS
0.08

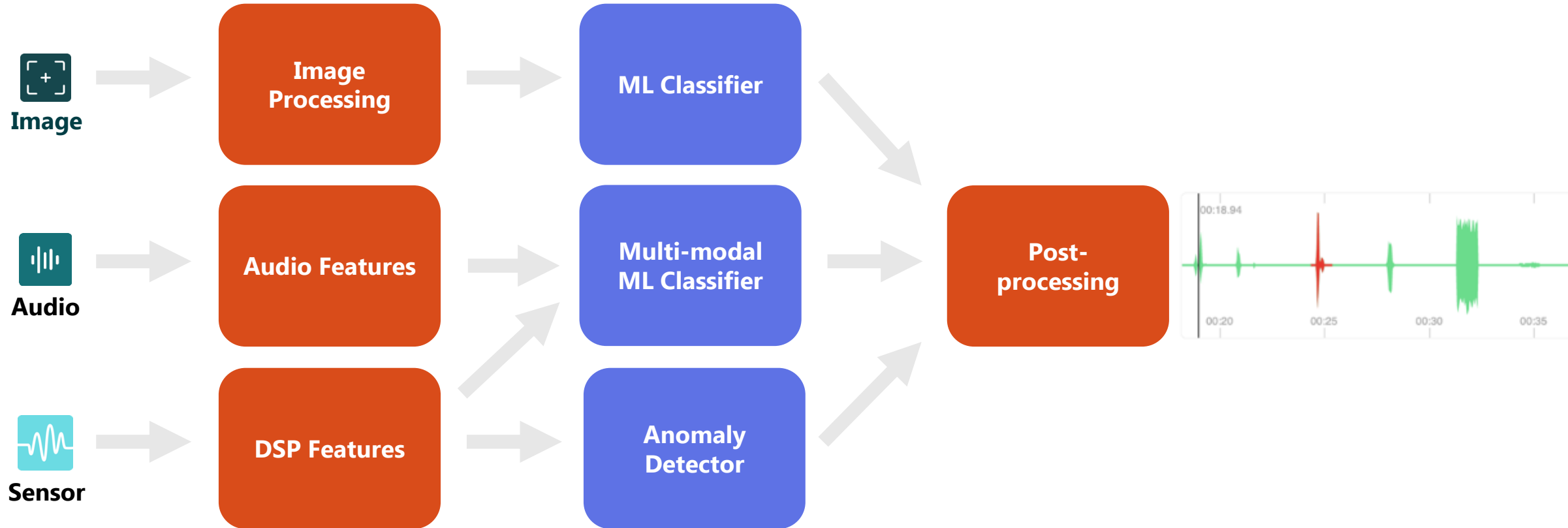
Confusion matrix (validation set)

	_BACKGROUND	_UNKNOWN	
_BACKGROUND	97.7%	2.3%	
_UNKNOWN	2.3%	96.9%	
IMPULSE	0%	0%	
F1 SCORE	0.98	0.97	

Sensor fusion for industrial-grade ML



Sensor fusion for industrial-grade ML



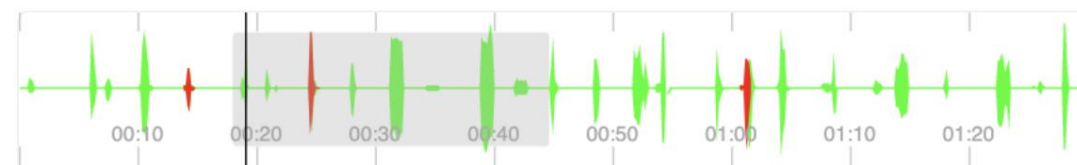
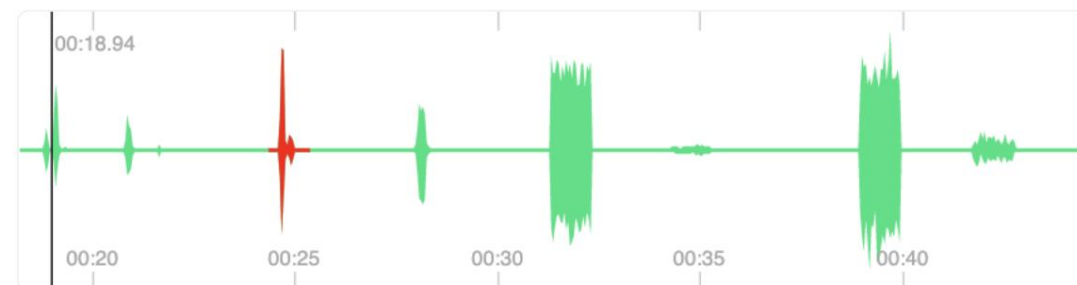
Calibrating performance at scale



- Model training validation \neq performance
- Requires testing on the entire algorithm with real-world data for realistic performance
- Understand the impact of post-processing while accounting for device constraints and latency
- Choose the ideal balance between false activations and false rejections
- Leverage genetic algorithms to design optimal post-processing configuration

Generated audio

Below you can see and play with the generated audio file, which also shows where false positives and negatives appear in the audio.



Zoom in

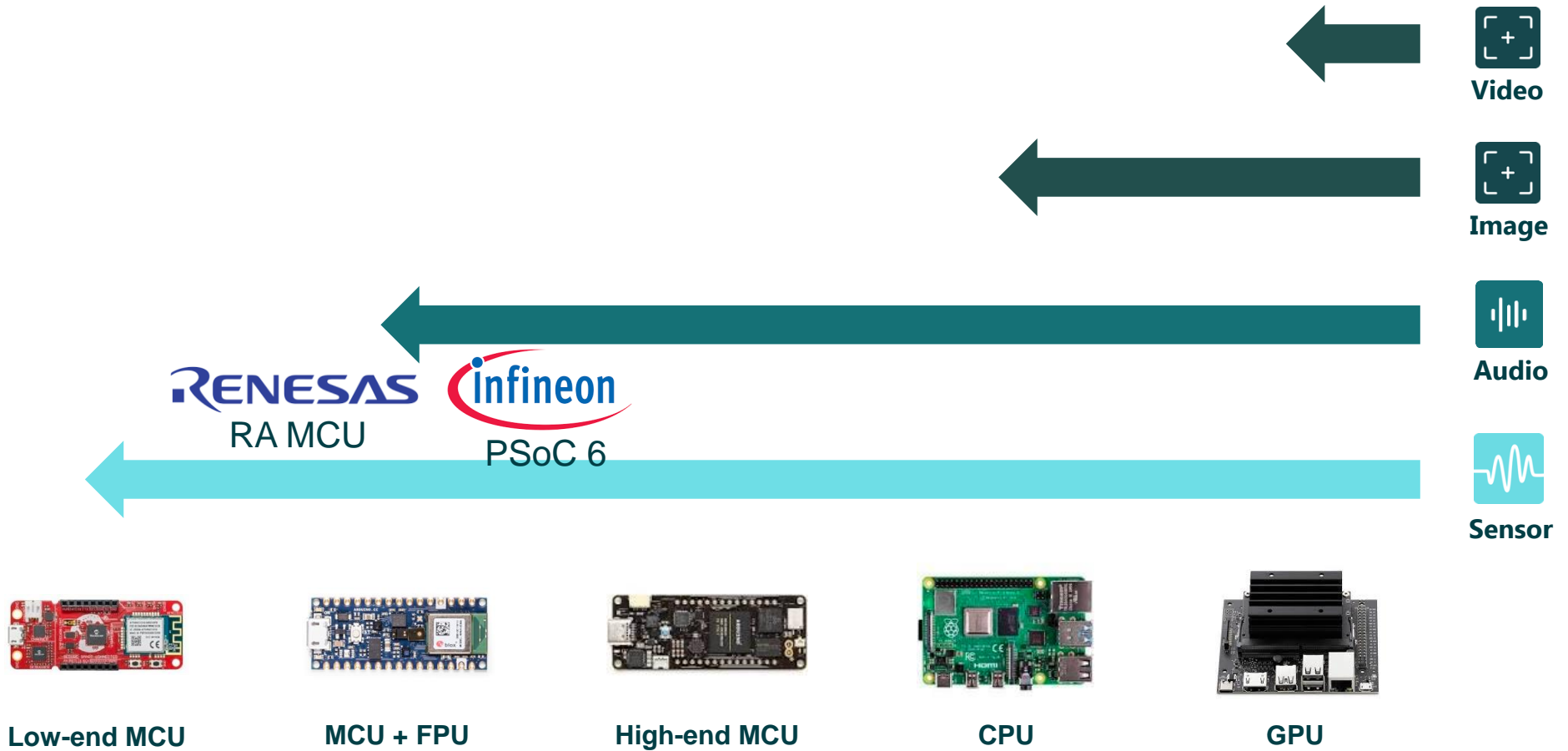
Zoom out



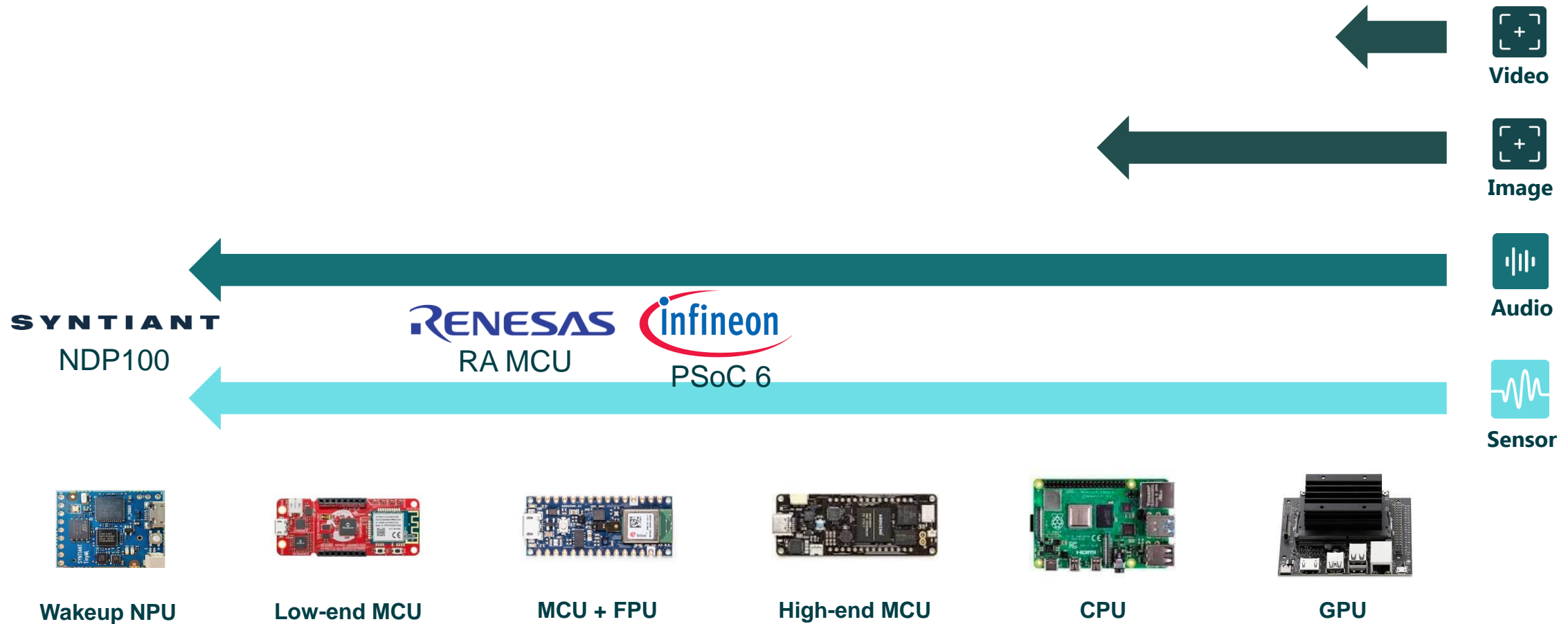
ID	LABEL	START TIME	END TIME	
9	yes => no	13.914	14.914	Play
27	yes => no	60.807	61.807	Play

ML on today and tomorrow's silicon

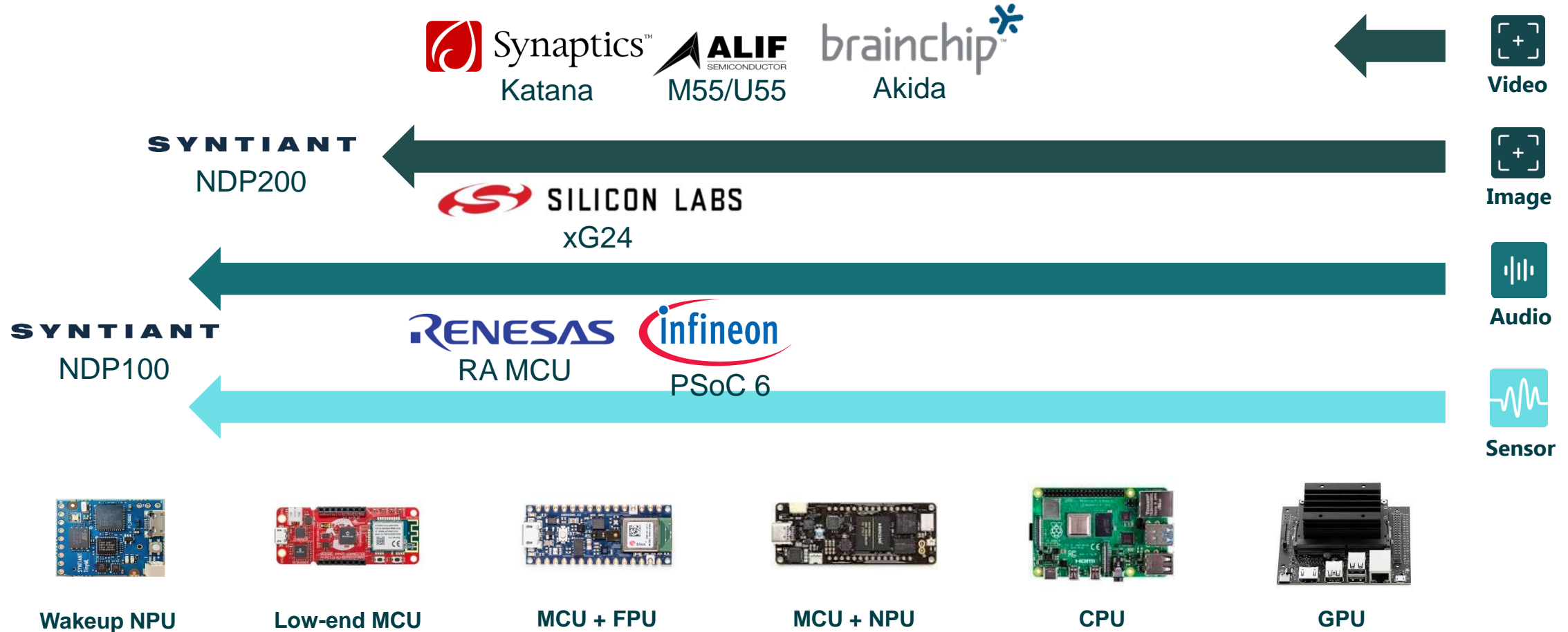
Hardware is sexy again!



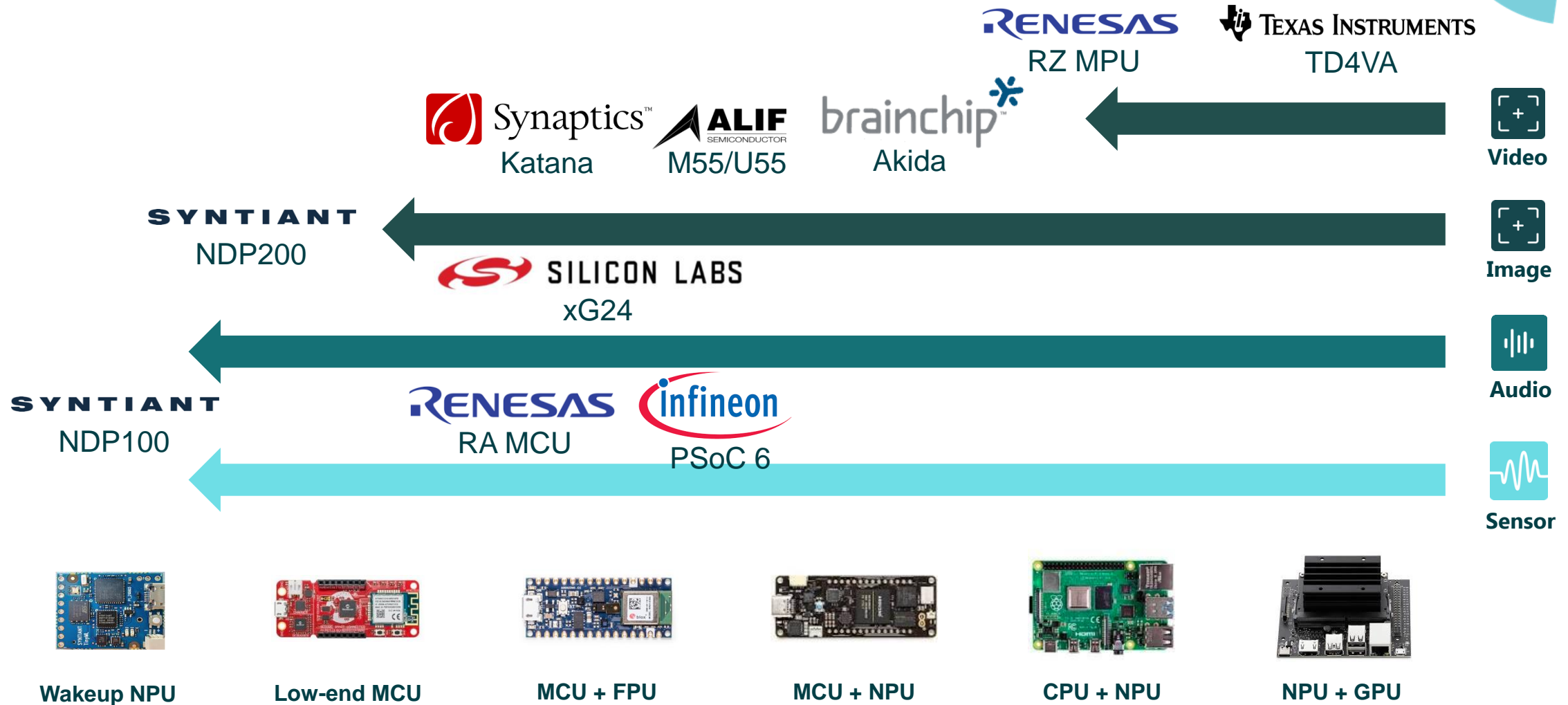
Hardware is sexy again!



Hardware is sexy again!



Hardware is sexy again!



Hardware profiling & tuning



- The power of hardware profiling
- Digital twin of ML on hardware
- We are combining hardware profiling with hyperparameter search – **EON Tuner**
- Hardware-aware AutoML across data, pre-processing and ML blocks

On-device performance ?



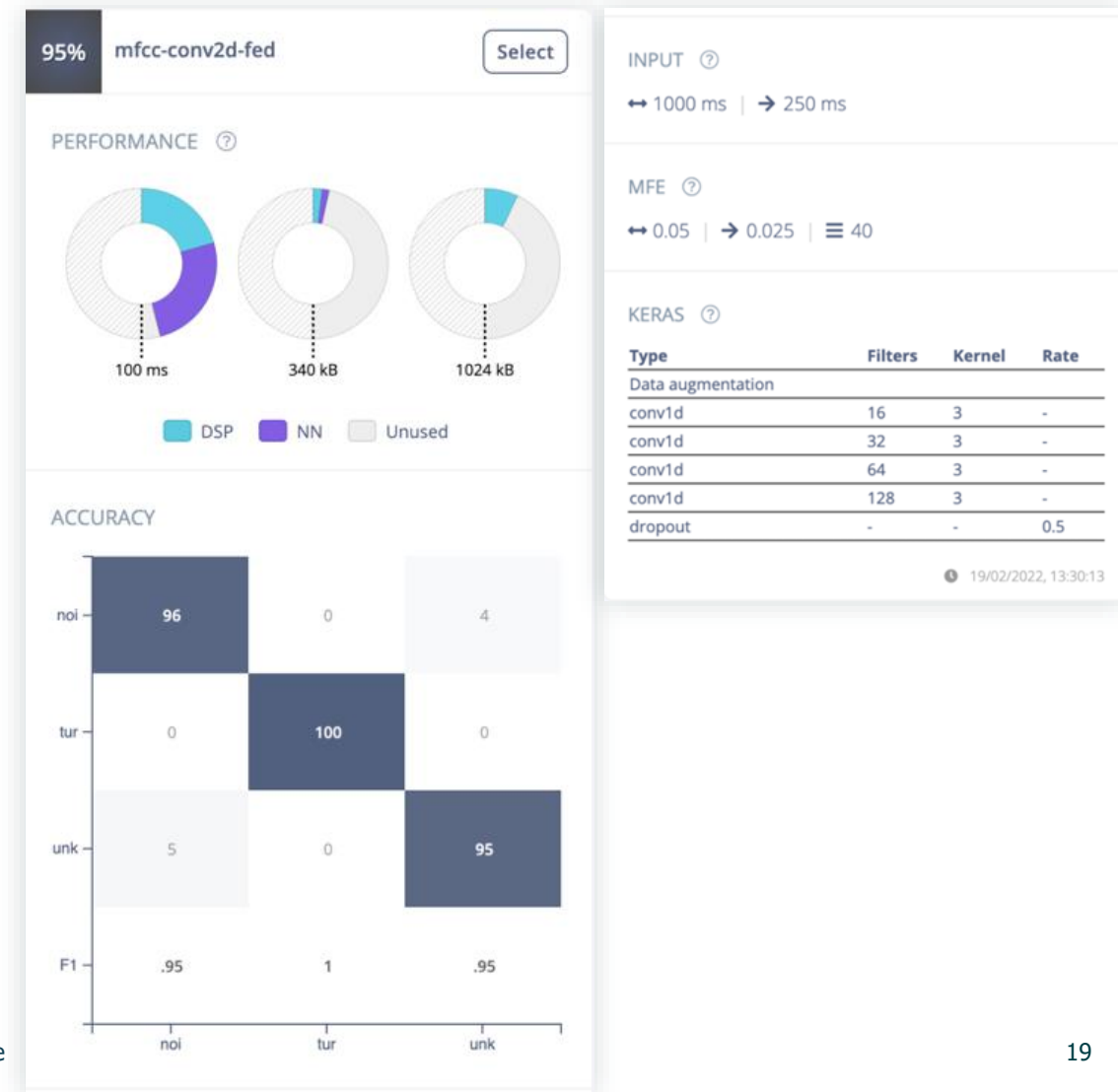
INFERENCE TIME
1 ms.



PEAK RAM USAGE
1.7K



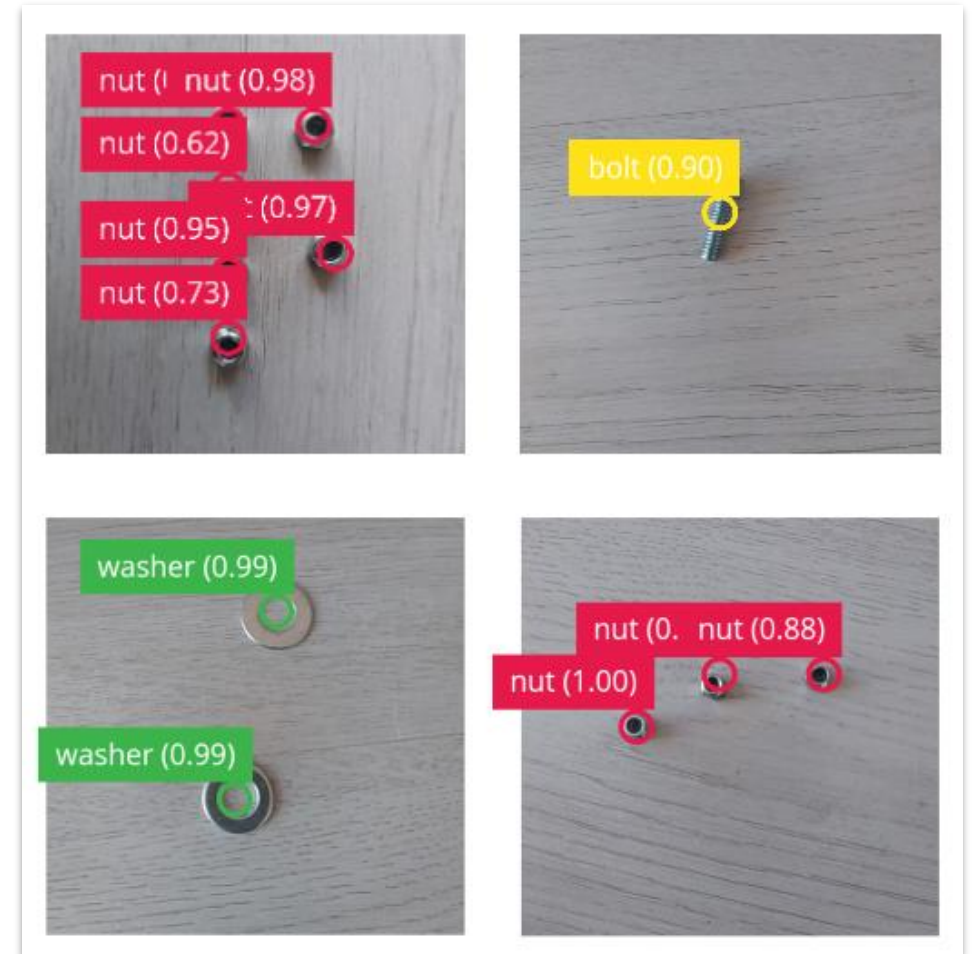
FLASH USAGE
19.4K



FOMO: Faster objects, more objects

- Object detection on constrained compute
- **20x** performance improvement
- Scales down to 100k RAM to enable MCUs
- Better at detecting smaller and more numerous objects
- Capable of segmentation and counting objects

	Cortex-M4	Cortex-M7	Cortex-A	Nvidia
FOMO	2 fps	15-30 fps	60+ fps	150+ fps
SSD	NA	NA	3 fps	20 fps



Increasing Manufacturing Efficiency

<https://casestudies.edgeimpulse.com/industrial-iot-with-advantech>

Constrained object detection: FOMO Blog

docs.edgeimpulse.com/docs/tutorials/fomo-object-detection-for-constrained-devices

Request a demo

edgeimpulse.com/schedule-a-demo

Don't miss these talks!



FOMO: Real-time Object Detection on Microcontrollers

Jan Jongboom, Co-founder and CTO, Edge Impulse

- **Date:** Wednesday, May 18
- **Start Time:** 10:50 am

Deep Dive: Develop and Deploy Advanced Edge Computer Vision—Fast!

Jenny Plunkett and Shawn Hymel, Snr. DevRel engineers

- **Date:** Thursday, May 19
- **Time:** 9 am – 12:00 pm