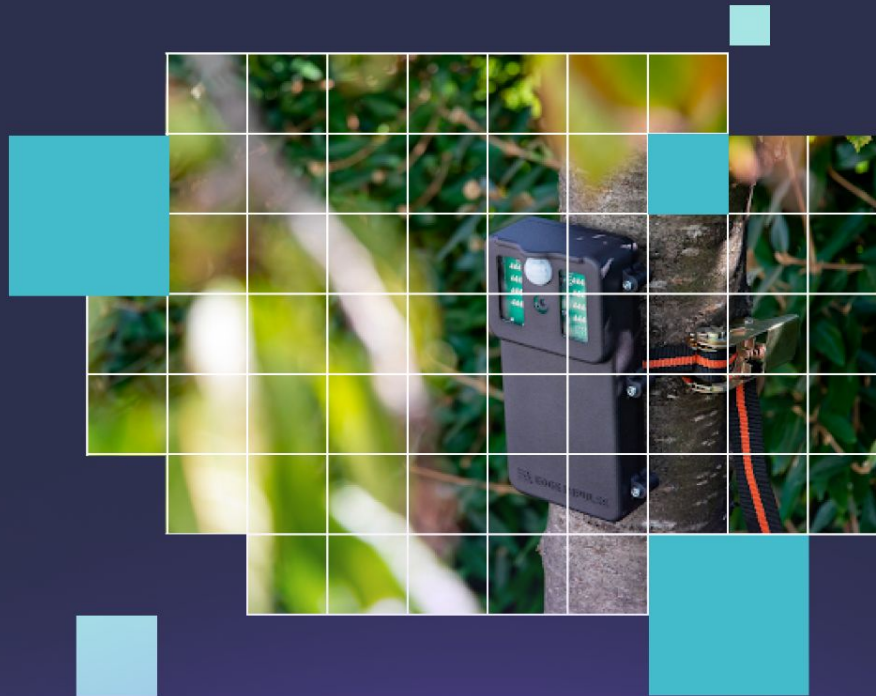EDGE IMPULSE

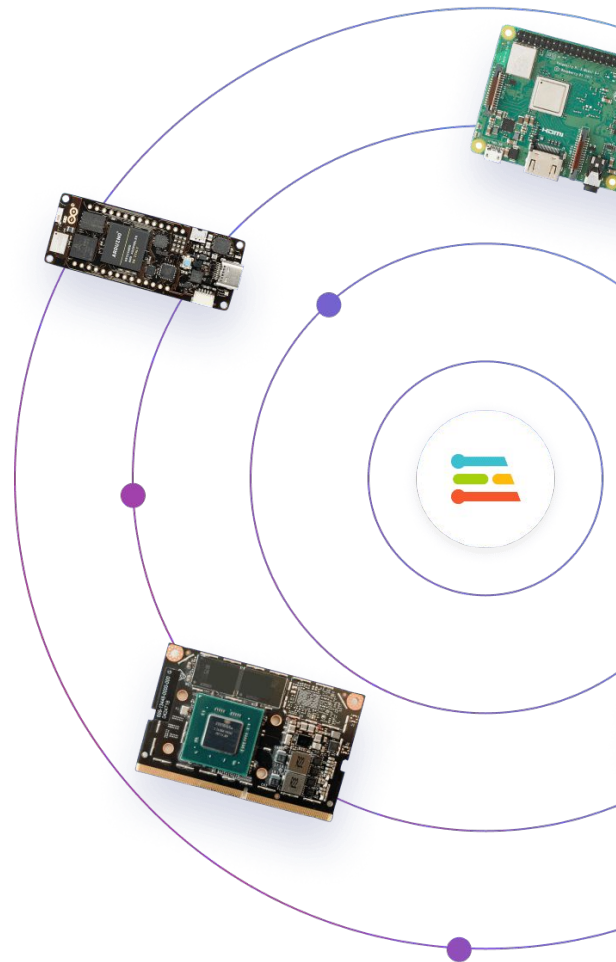# The Future of Embedded ML

Alessandro Grande
Head of Product

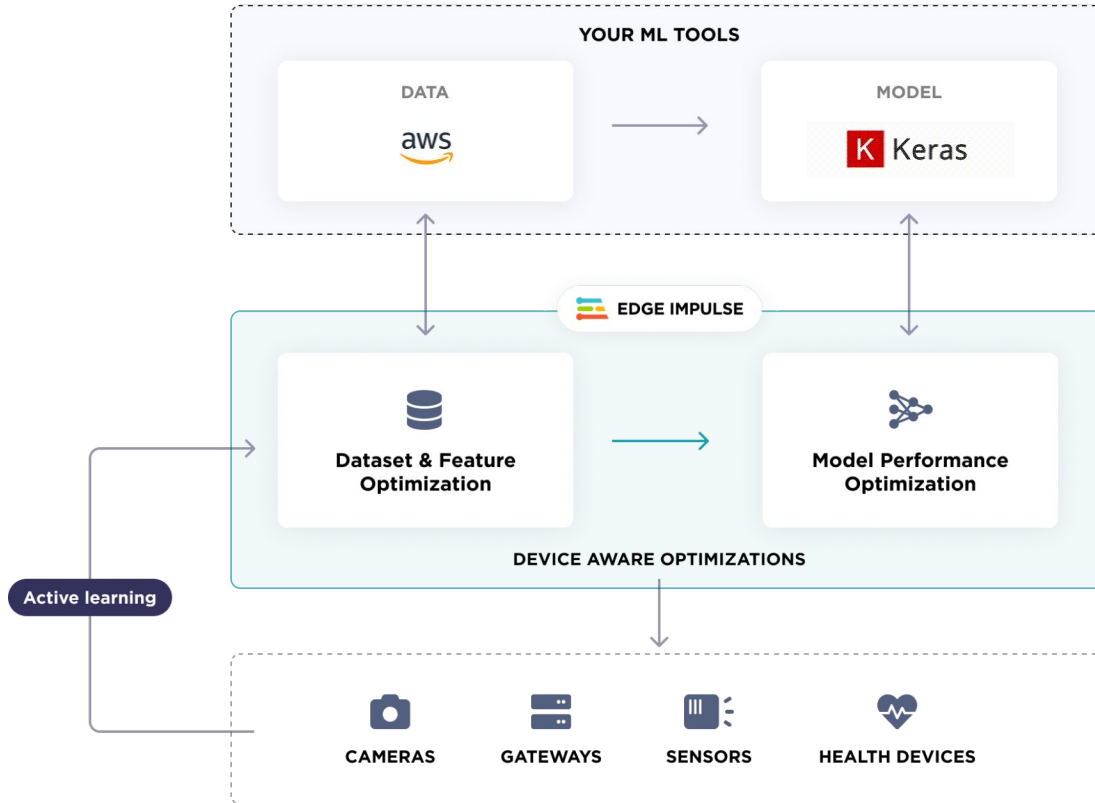*ICTP, Trieste - July 3, 2023*

# Agenda

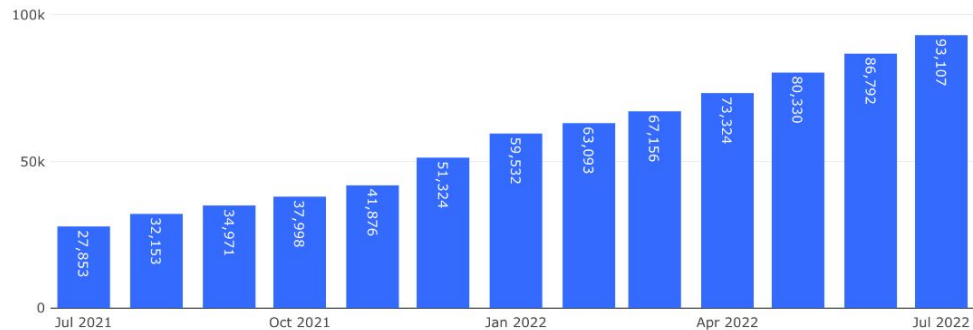1. Intro to Edge Impulse

2. Customer challenges

3. What's beneath the surface

4. Resources

5. Next steps

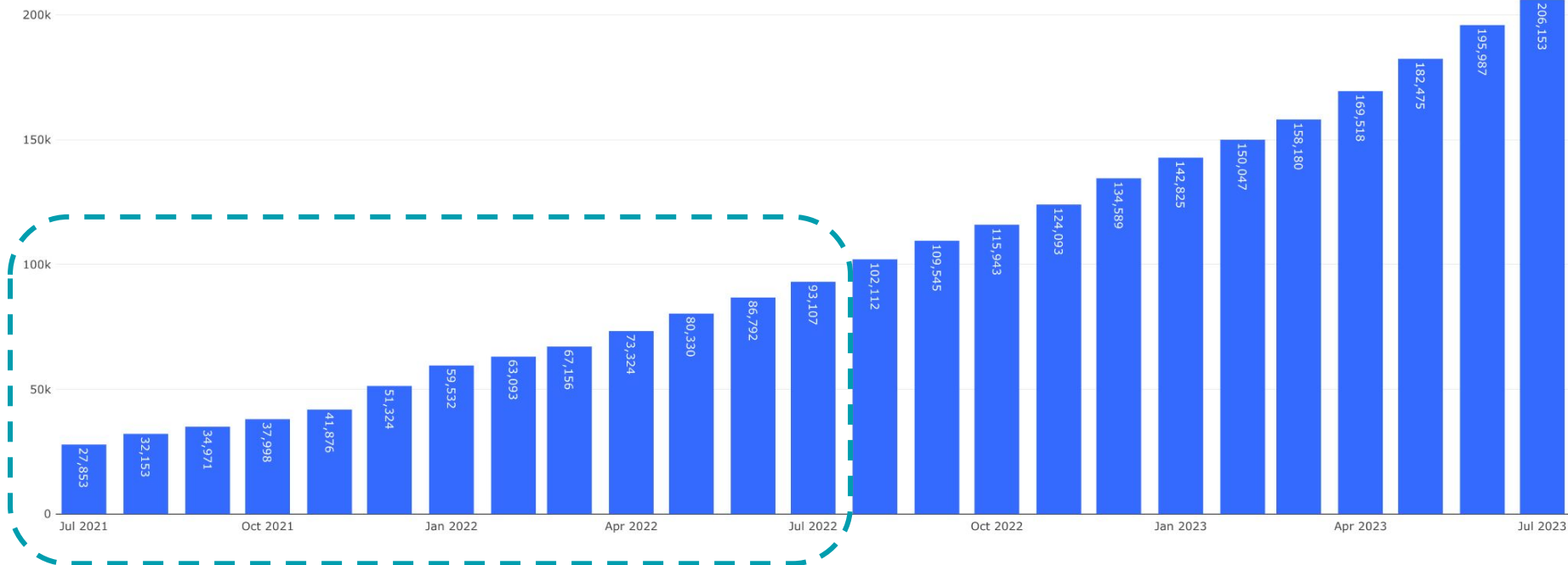# The edge AI platform



**YOUR ML TOOLS**

DATA

aws

MODEL

K Keras

EDGE IMPULSE

Dataset & Feature Optimization

Model Performance Optimization

**DEVICE AWARE OPTIMIZATIONS**

Active learning

CAMERAS    GATEWAYS    SENSORS    HEALTH DEVICES

# Number of Projects on Edge Impulse

100k

50k

27,853 | 32,153 | 34,971 | 37,998 | 41,876 | 51,324 | 59,532 | 63,093 | 67,156 | 73,324 | 80,330 | 86,792 | 93,107

0

Jul 2021          Oct 2021          Jan 2022          Apr 2022          Jul 2022

# Number of Projects on Edge Impulse



Jul 2021: 27,853
32,153
34,971
Oct 2021: 37,998
41,876
51,324
Jan 2022: 59,532
63,093
67,156
Apr 2022: 73,324
80,330
86,792
Jul 2022: 93,107
102,112
109,545
Oct 2022: 115,943
124,093
134,589
Jan 2023: 142,825
150,047
158,180
Apr 2023: 169,518
182,475
195,987
Jul 2023: 206,153

# TinyML Use Cases



Health — Industrial — Wearables — Infrastructure — Buildings

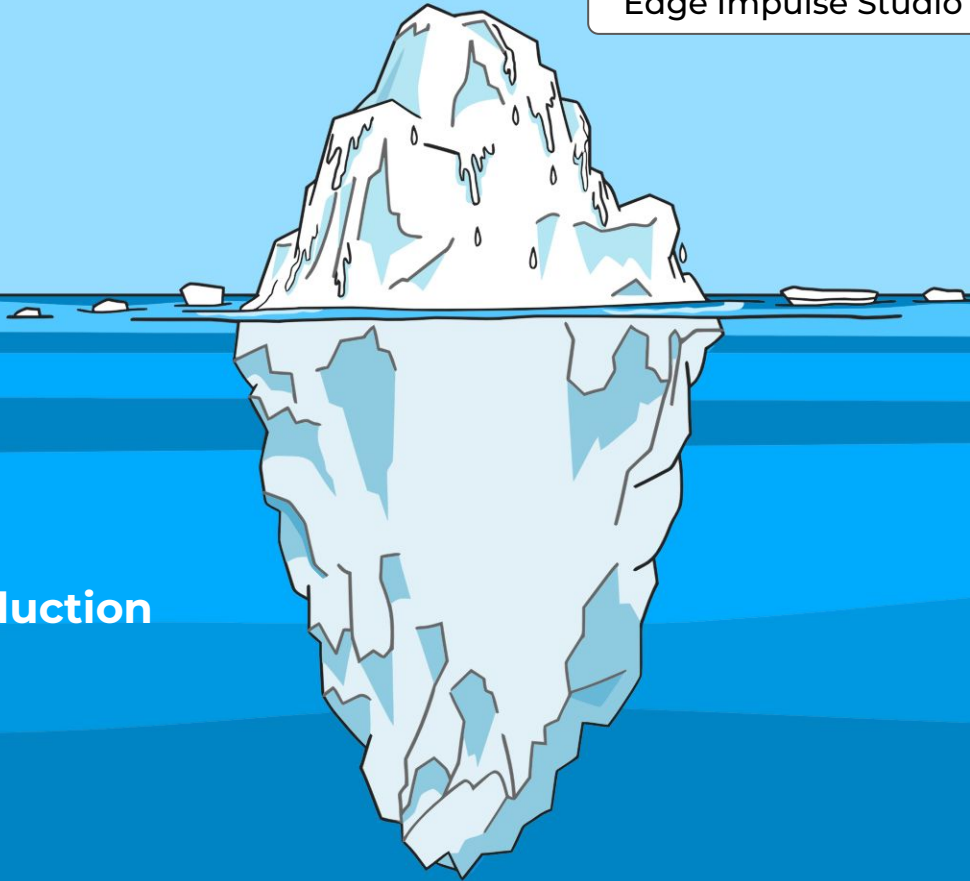ŌURA — KNOW LABS — NOWATCH — Brambles — NASA — Lexmark

Hyfe AI — ECOLAB — IZOELEKTRO

# Production Challenges
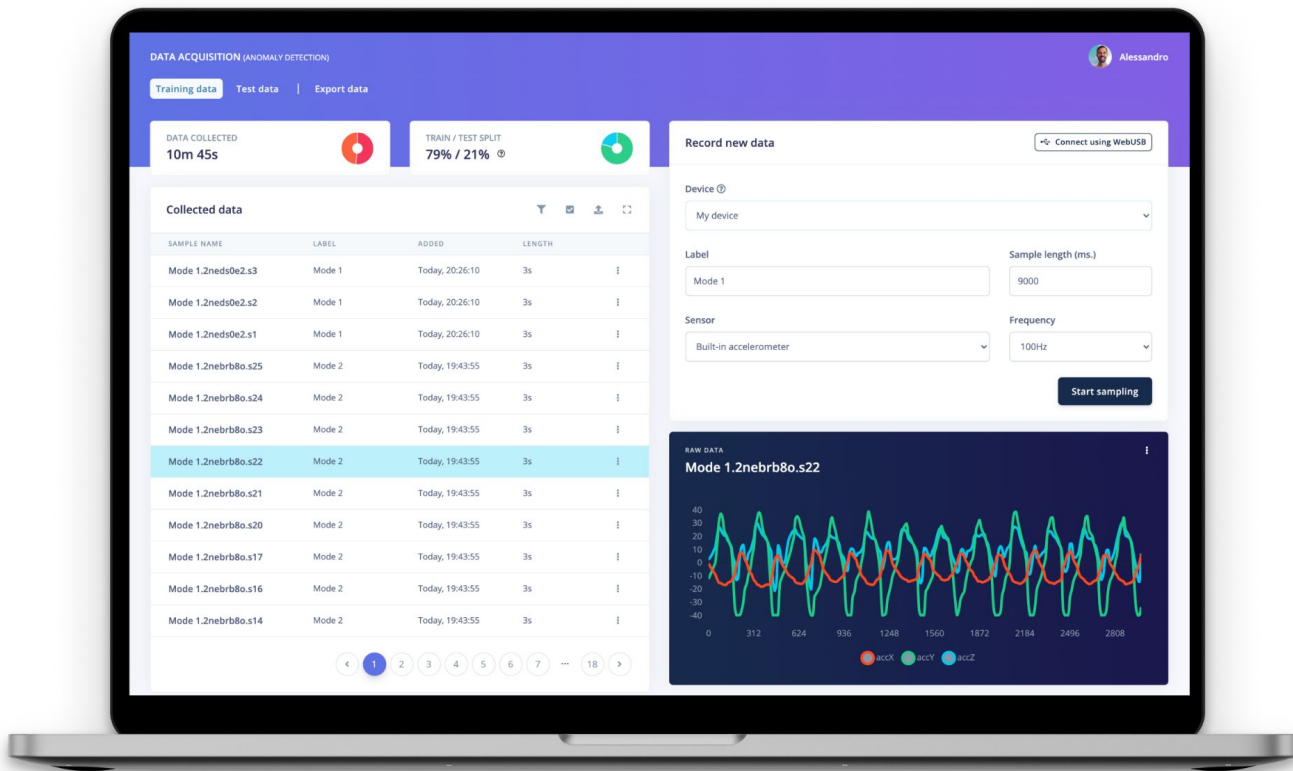
1. Data collection

2. Data quality analysis

3. Feature extraction and DSP

4. Deployment

5. Monitoring performance

studio.edgeimpulse.com/evaluate
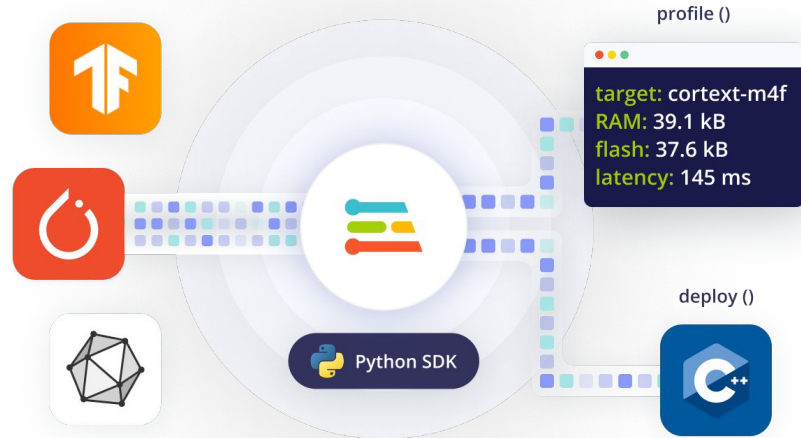
# BYOM & Python SDK

- Profile on-device performance of any trained model

- Analyze the impact of architectural decisions

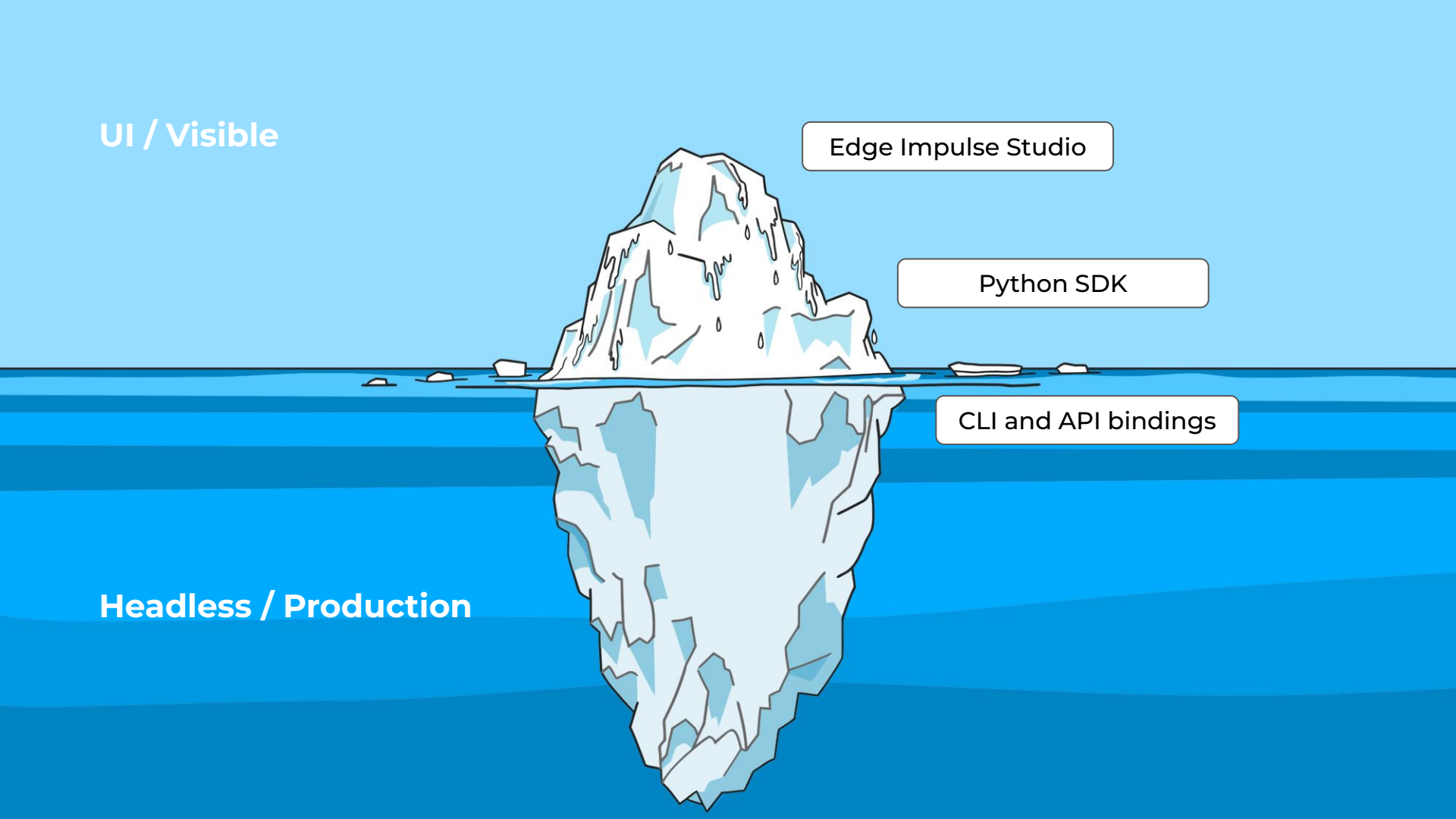- Generate optimized C++ libraries

- Deploy to any edge device

profile ()

target: cortext-m4f
RAM: 39.1 kB
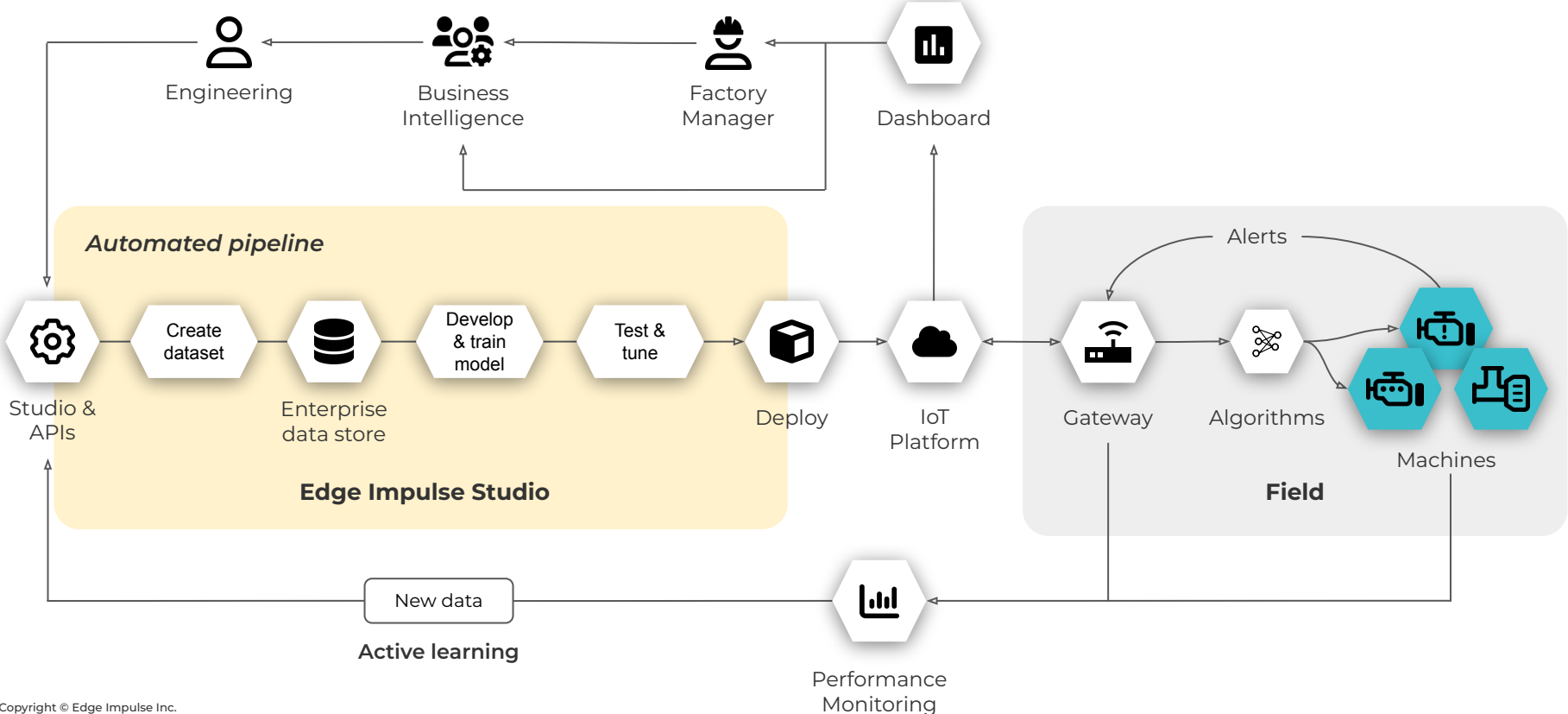flash: 37.6 kB
latency: 145 ms

Python SDK

deploy ()

# CLI

## Edge Impulse CLI tools

Command-line interface tools for Edge Impulse. We make things smarter by enabling developers to create the next generation of intelligent device solutions with embedded Machine Learning.

This package consists of four tools (click to see their respective documentation):

- edge-impulse-daemon - configures devices over serial, and acts as a proxy for devices that do not have an IP connection.
- edge-impulse-uploader - allows uploading and signing local files.
- edge-impulse-data-forwarder - a very easy way to collect data from any device over a serial connection, and forward the data to Edge Impulse.
- edge-impulse-run-impulse - show the impulse running on your device.
- edge-impulse-blocks - create organizational transformation blocks.
- eta-flash-tool - to flash the Eta Compute ECM3532 AI Sensor.
- himax-flash-tool - to flash the Himax WE-I Plus development board.

# Embedded ML in the Real World



Engineering

Business Intelligence

Factory Manager

Dashboard

**Automated pipeline**

Studio & APIs

Create dataset

Enterprise data store

Develop & train model

Test & tune

Deploy

IoT Platform

Gateway

Algorithms

Alerts

Machines

**Edge Impulse Studio**

**Field**

New data

**Active learning**

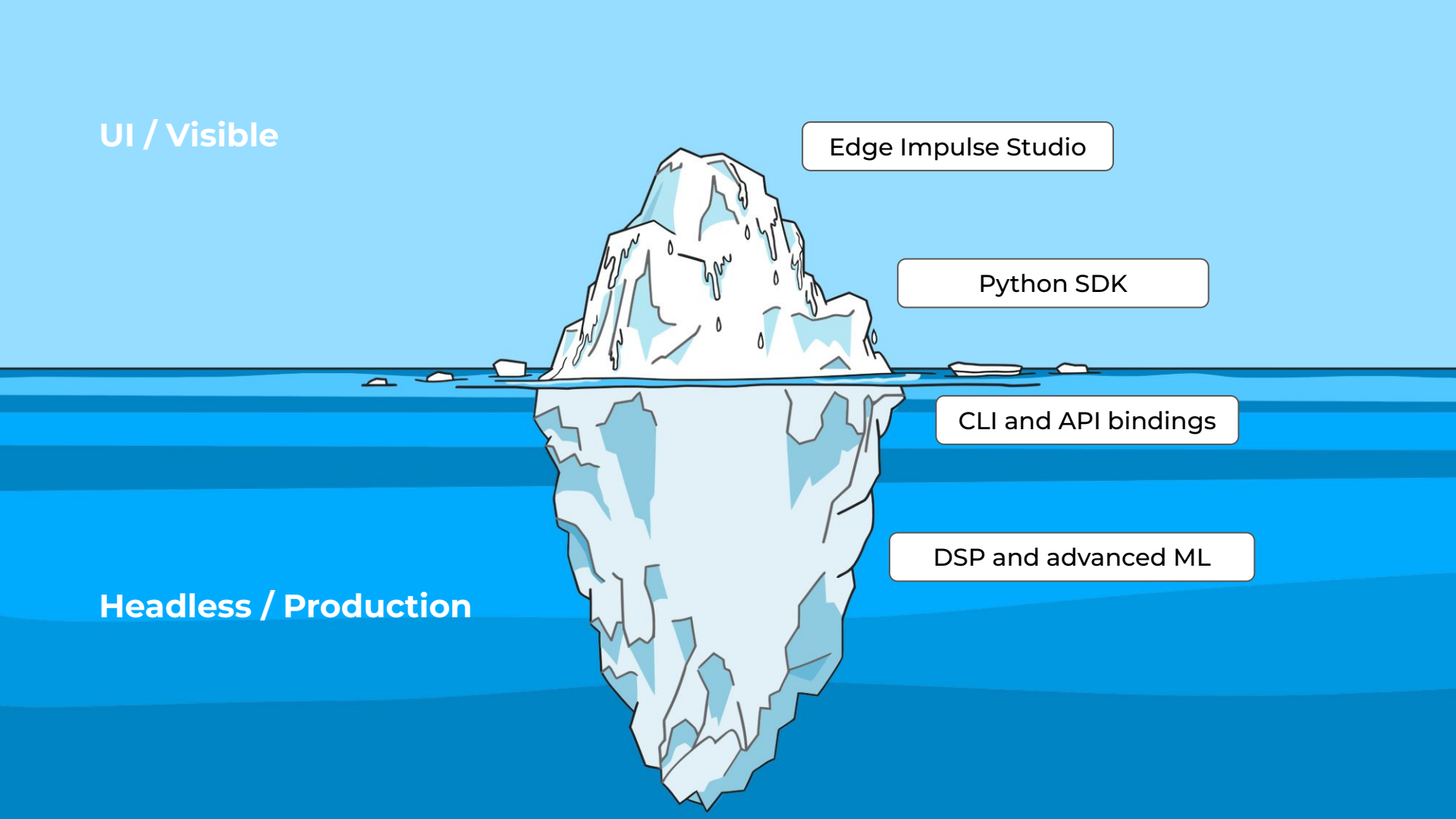Performance Monitoring

UI / Visible

Edge Impulse Studio

Python SDK
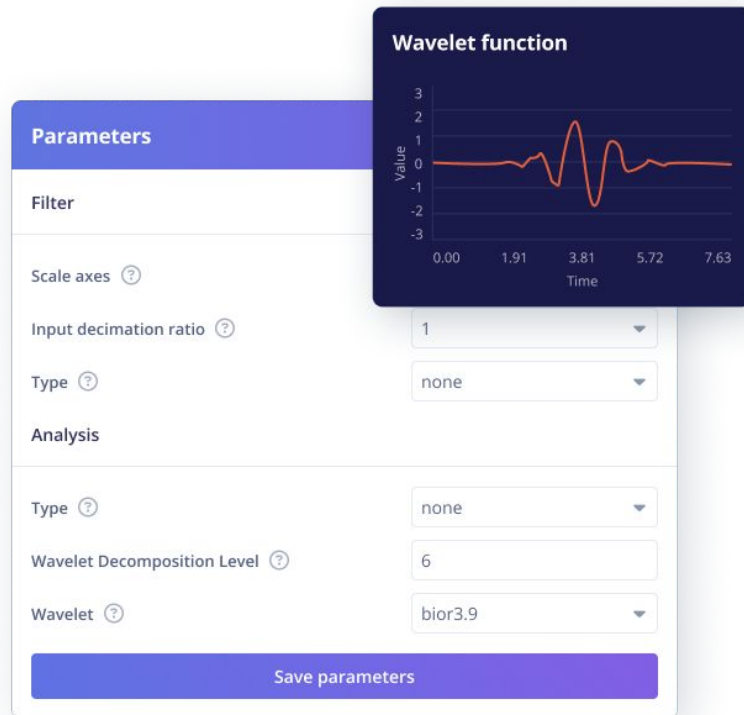
CLI and API bindings

DSP and advanced ML

Headless / Production

# Interactive Feature Engineering
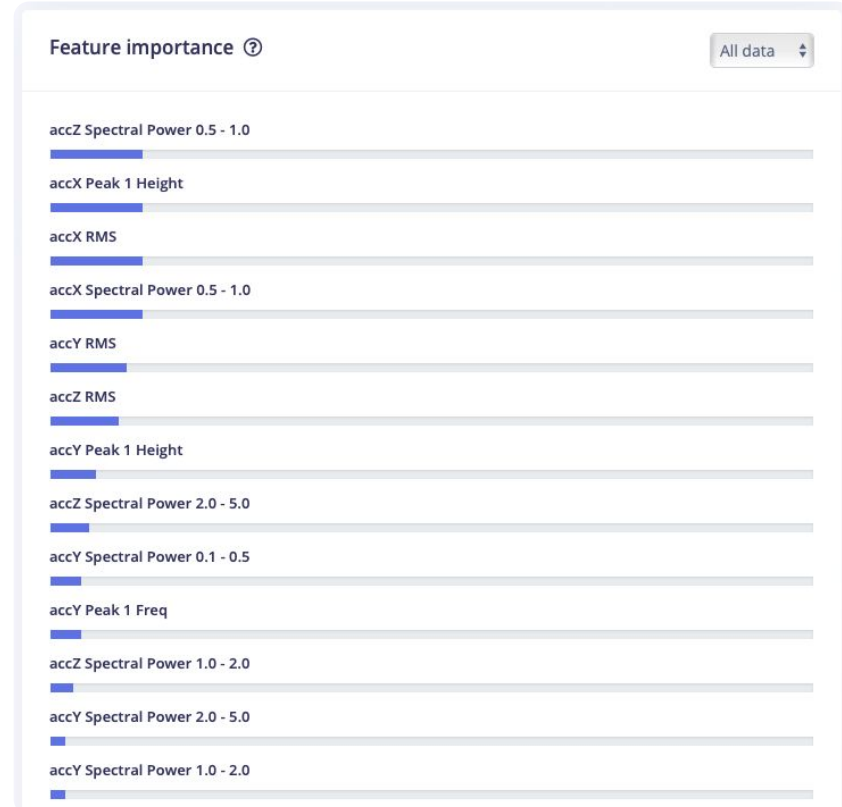
**Real-time visualization of DSP**

- Immediate feedback loop enabling tactile exploration by domain expert

- Service-based architecture for real time DSP on individual samples (separate from job-based system for batched data)
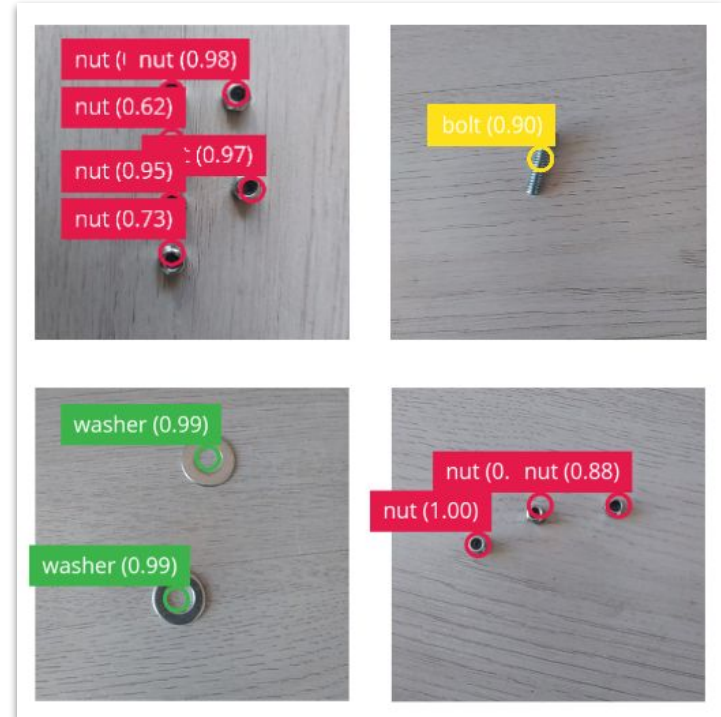
# Feature Importance

## Don't use everything

- Uses recursive feature elimination with cross-validation (RFECV)

- Only computed for relatively low-dimensionality data



Feature importance ⓘ                    All data ⌄

accZ Spectral Power 0.5 - 1.0

accX Peak 1 Height

accX RMS

accX Spectral Power 0.5 - 1.0

accY RMS

accZ RMS

accY Peak 1 Height

accZ Spectral Power 2.0 - 5.0

accY Spectral Power 0.1 - 0.5

accY Peak 1 Freq

accZ Spectral Power 1.0 - 2.0

accY Spectral Power 2.0 - 5.0

accY Spectral Power 1.0 - 2.0

# FOMO: Faster Objects, More Objects

- **20x average performance improvement**
- Object detection on MCUs
- Ultra fast on embedded Linux
- Better at detecting smaller and more numerous objects
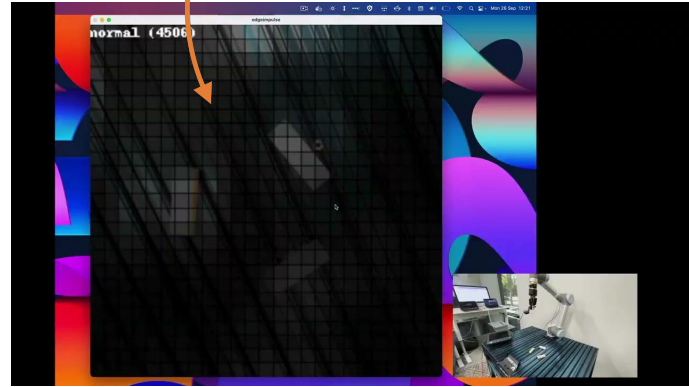- Capable of segmentation and counting objects

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

# FOMO: Faster Objects, More Objects

- Remove classification head, replace with GMMs

- Only requires training on normal data

- Each cell tells you the chance that it's an anomaly

- Same performance as FOMO:
  Up to 30fps. on Cortex-M7, <200K RAM

Each cell is an anomaly detector

# Keras Expert Mode

- For advanced users

- Use Keras standard API

- Customize NN architecture and take full control over training procedure

### Neural Network settings

⋮

#### Training settings

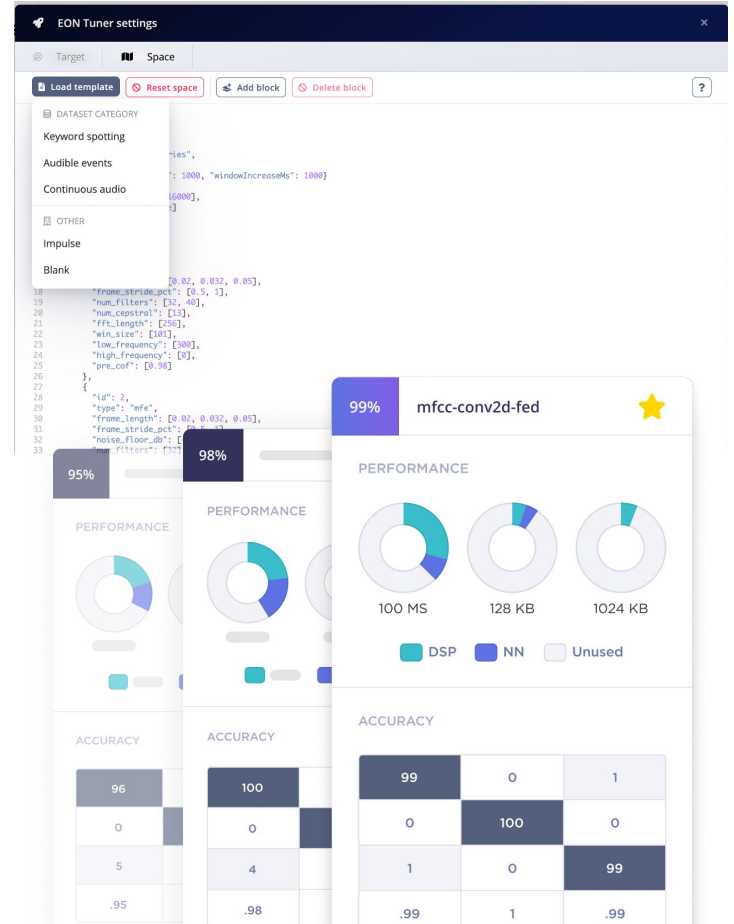Validation set size ?                 `20`   %

#### Neural network architecture

```python
import tensorflow as tf
from tensorflow.keras.models import Sequential
from tensorflow.keras.layers import Dense, InputLayer, Dropout, Conv1D, Conv2D, Flatten, Reshape, MaxPooling1D,
    MaxPooling2D, BatchNormalization, TimeDistributed
from tensorflow.keras.optimizers import Adam

# model architecture
model = Sequential()
model.add(Dense(20, activation='relu',
    activity_regularizer=tf.keras.regularizers.l1(0.00001)))
model.add(Dense(10, activation='relu',
    activity_regularizer=tf.keras.regularizers.l1(0.00001)))
model.add(Dense(classes, activation='softmax', name='y_pred'))

# this controls the learning rate
opt = Adam(learning_rate=0.0005, beta_1=0.9, beta_2=0.999)
# this controls the batch size, or you can manipulate the tf.data.Dataset objects yourself
BATCH_SIZE = 32
train_dataset = train_dataset.batch(BATCH_SIZE, drop_remainder=False)
validation_dataset = validation_dataset.batch(BATCH_SIZE, drop_remainder=False)
callbacks.append(BatchLoggerCallback(BATCH_SIZE, train_sample_count))

# train the neural network
model.compile(loss='categorical_crossentropy', optimizer=opt, metrics=['accuracy'])
model.fit(train_dataset, epochs=30, validation_data=validation_dataset, verbose=2, callbacks=callbacks)

# Use this flag to disable per-channel quantization for a model.
# This can reduce RAM usage for convolutional models, but may have
# an impact on accuracy.
disable_per_channel_quantization = False
```

# EON Tuner

**Establish a baseline quickly**

- Search space based on prior knowledge of data modalities

- Reusable workers to minimize startup cost

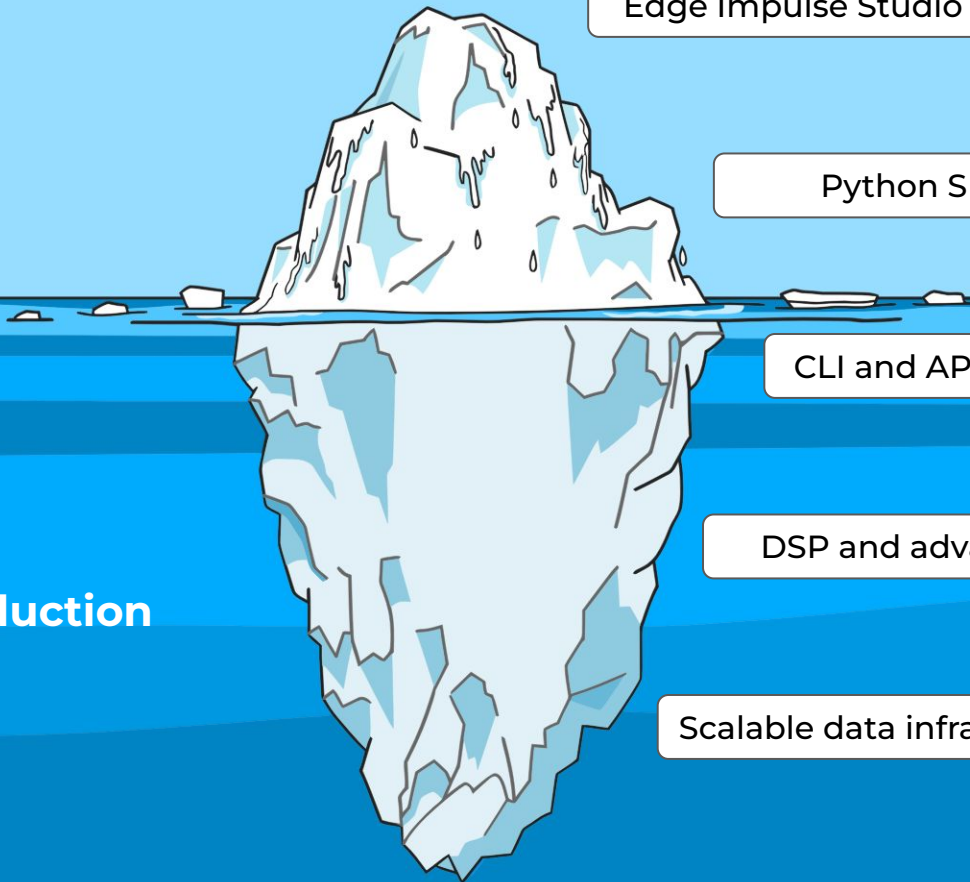- Customize search space

UI / Visible
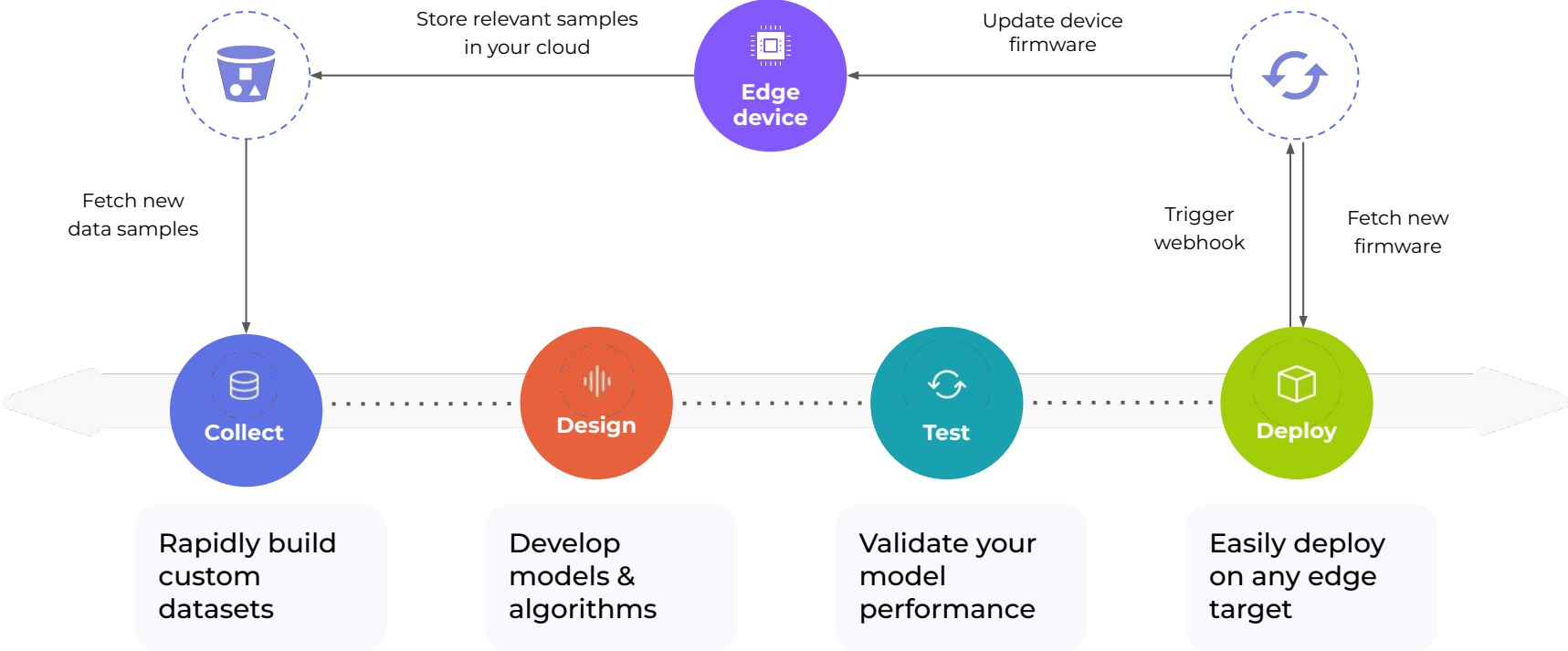
Edge Impulse Studio

Python SDK

CLI and API bindings

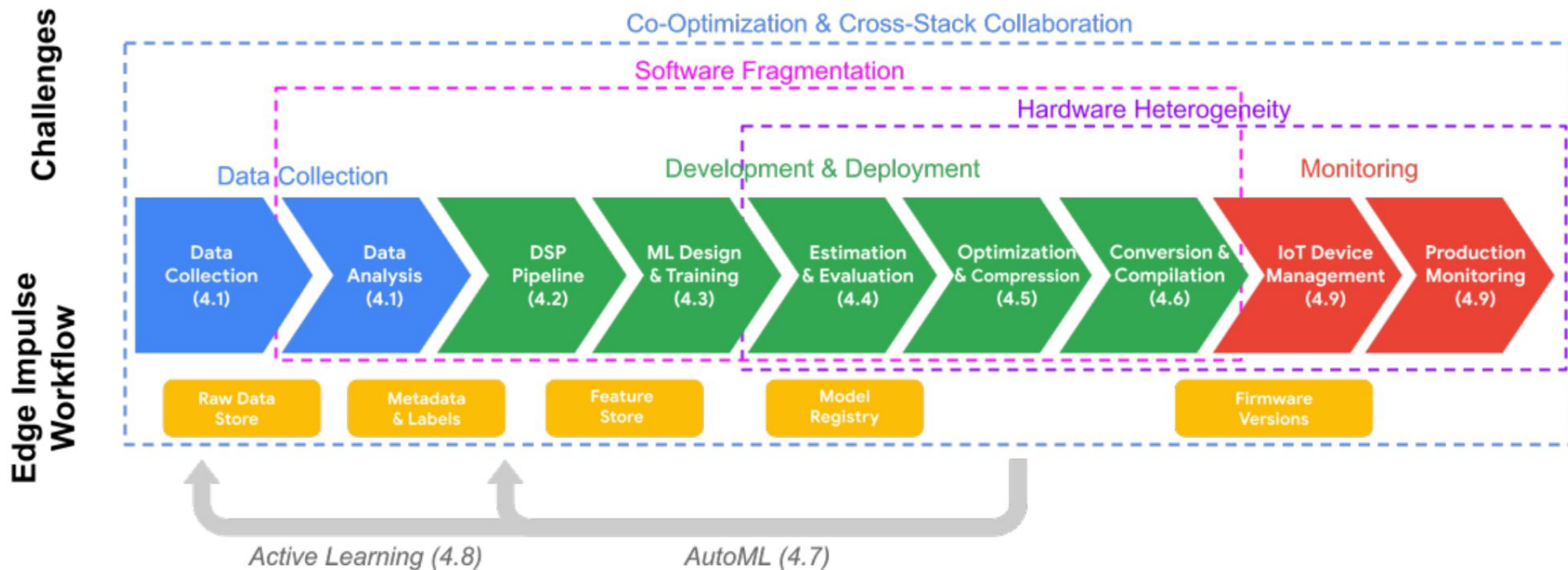DSP and advanced ML

Headless / Production
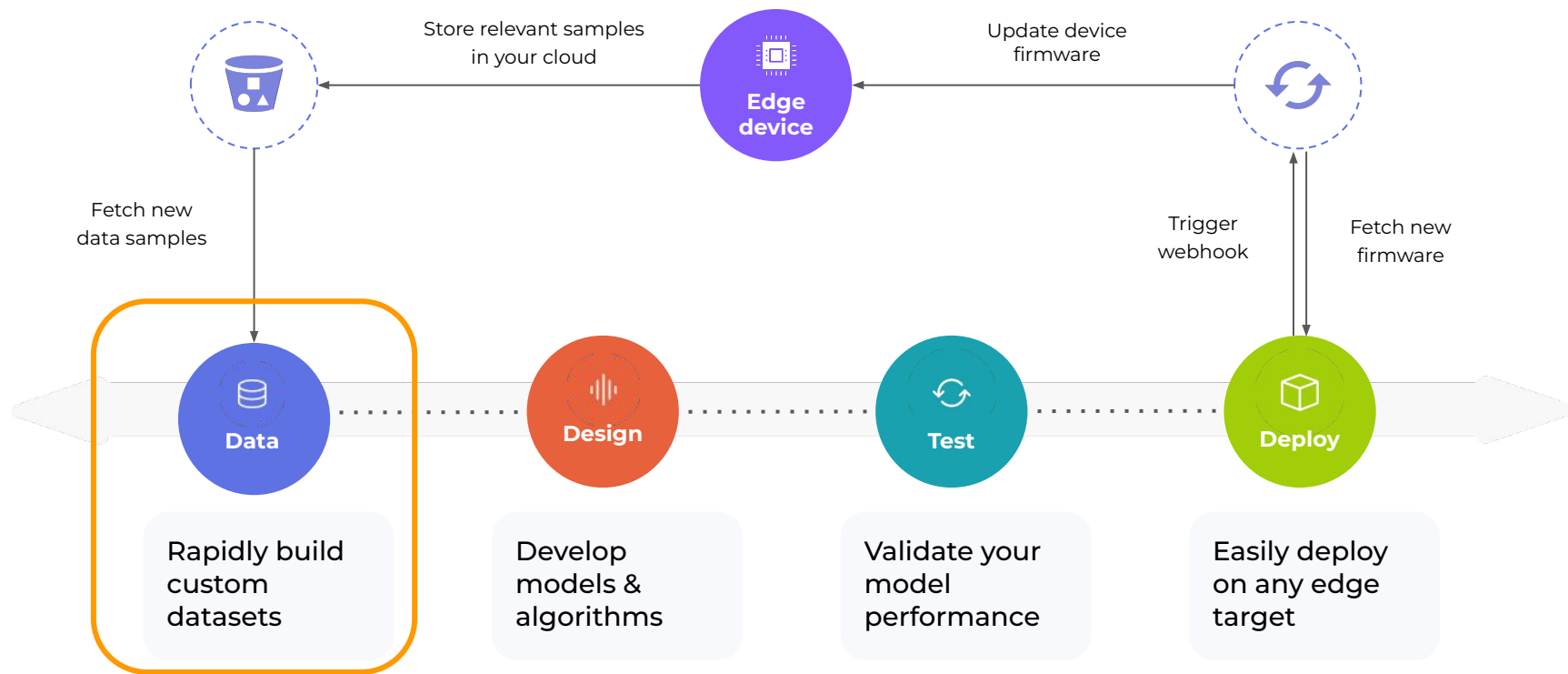
Scalable data infrastructure

# Data-Centric ML



Store relevant samples
in your cloud

Edge
device

Update device
firmware

Fetch new
data samples

Trigger
webhook

Fetch new
firmware

Collect

Design

Test

Deploy

Rapidly build
custom
datasets

Develop
models &
algorithms

Validate your
model
performance

Easily deploy
on any edge
target

# Active Learning with Edge Impulse



*Source: Shawn Hymel, et al. Edge Impulse: An MLOps Platform for Tiny Machine Learning, November 2022. arXiv:2212.03332*

# Active Learning **with Edge Impulse**



Store relevant samples
in your cloud

Update device
firmware

**Edge
device**

Fetch new
data samples

Trigger
webhook

Fetch new
firmware

**Data**

**Design**

**Test**

**Deploy**

Rapidly build
custom
datasets

Develop
models &
algorithms

Validate your
model
performance

Easily deploy
on any edge
target

# Working with Data

- Data pipelines and transformation - enabling data preparation at scale

- Data campaign dashboards - optimize performance and share learnings



Data pipelines and transformations



Data dashboards

# Data preparation

- Fetch data

- **Basic checks**: Are all files present? Do all files start / end around the same time? All expected labels for the study present?

- **Advanced checks**: Correlation between different devices (e.g. HR from PPG, and HR from Polar)?

- Runs automatically at a set interval (or on-demand, or triggered from code)

- Sends email on new data

# Visualize data and uncover critical insights

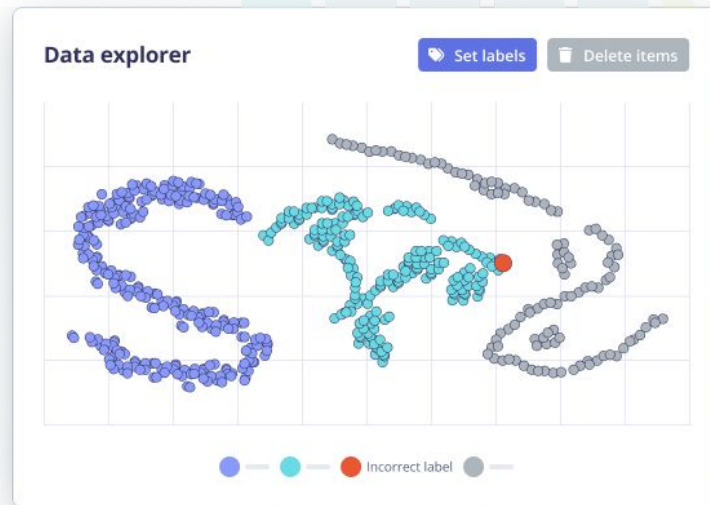## Data explorer

The data explorer shows a complete view of all data in your project. Use it to quickly label your data, or spot outliers. Learn more.
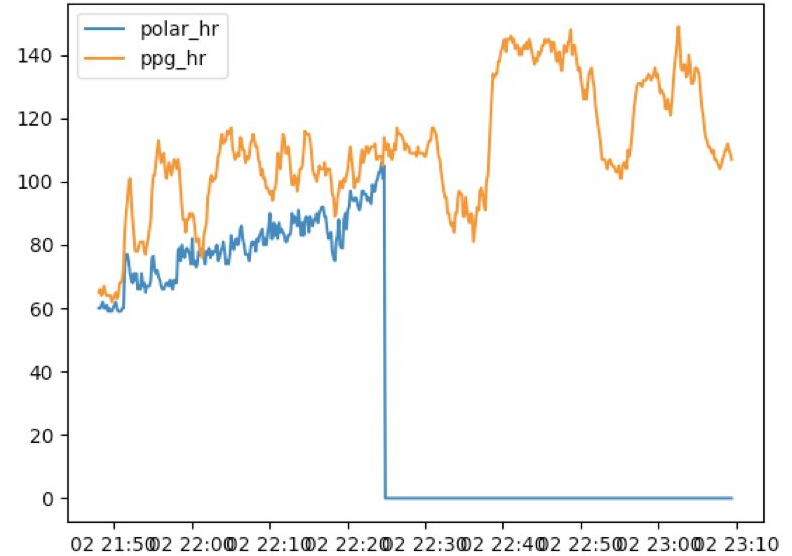
**How should we generate the data explorer?**

○ 🎤 **Using a pretrained keywords model**
Great for keywords that fit in a 1 second window.

○ 〰️ **Using your trained impulse**
Works great if you have collected some labeled data already and have a trained model.

○ 📶 **Using the preprocessing blocks in your impulse**
Use this if you don't have any labels for your data yet, and thus can't train a full model.

**Dimensionality reduction technique**

◉ 📊 **t-SNE**
Recommended for your dataset. Separates best, but takes a significant amount of time on large datasets.

○ 📊 **PCA**
Separates less well, but works on any dataset size.

## Data explorer

🏷️ Set labels    🗑️ Delete items

🔵 ━   🔵 ━   🔴 Incorrect label   ⚪ ━

# Validating data: correlation



Fixes many issues: data uploaded for wrong participant, device failure and
can be used to collect clock drift. Here implemented for PPG (derive HR) + Polar H10.

# Resources
## (Courses)

[tinyml.seas.harvard.edu](tinyml.seas.harvard.edu)

# Curriculum and Content

tinyml.seas.harvard.edu

## Full Courses

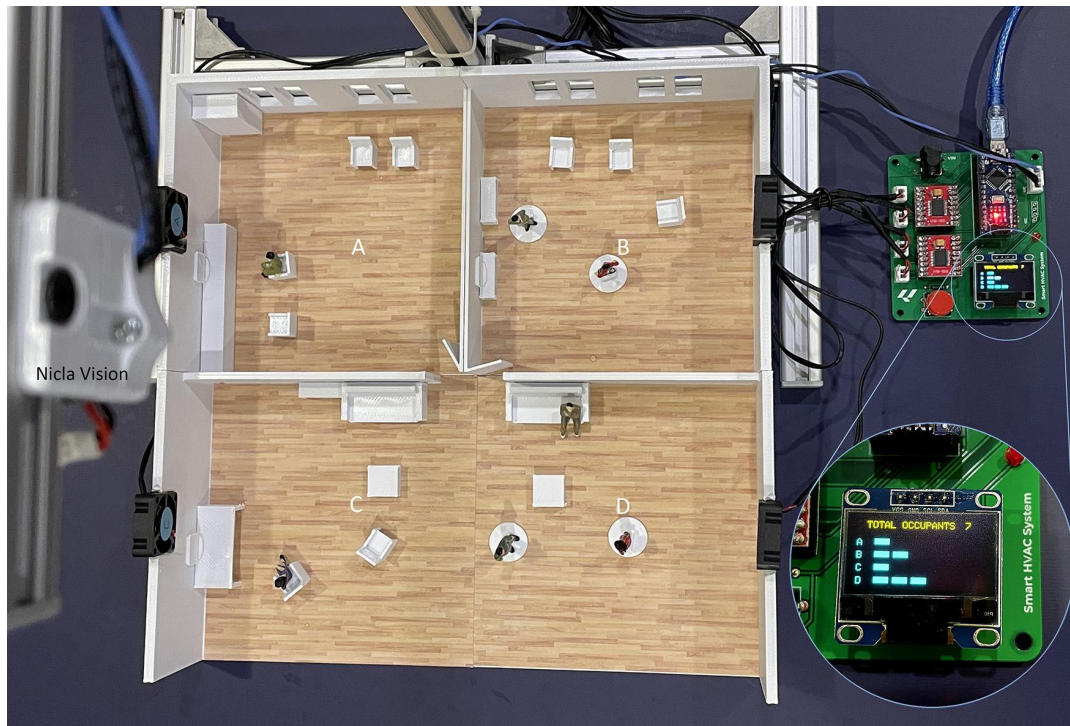| Organization | Course Name | Date of Course | Target Audience | Language of Instruction | Language of Materials | Links |
|---|---|---|---|---|---|---|
| edX | edX tinyML Specialization<br>by Harvard University | Launched 2020-2022 | Everyone | English | English | Course 1-3 Website<br>Course 4 Website<br>All Materials<br>All Colabs<br>Arduino Library |
| C | Embedded Machine Learning on Coursera<br>by Edge Impulse | Launched 2021-2022 | Everyone | English | English | Course 1<br>Course 2<br>All Materials |
| | ESE3600: Tiny Machine Learning<br>by the University of Pennsylvania | Fall 2022 | Undergraduate and Graduate Students | English | English | Website and Materials |
| MIT | MIT 6.S965<br>TinyML and Efficient Deep Learning | Fall 2022 | Graduate Students | English | English | Website<br>Materials |
| | UNIFEI IESTI01<br>TinyML - Machine Learning for Embedding Devices | Jan 2021 - Present | Undergraduate Students | Portuguese | English | 2022.1 Website and Materials<br>2021.2 Website and Materials<br>2021.1 Website and Materials |
| | Harvard CS249r<br>Tiny Machine Learning | Sept 2020 - Present | Graduate Students | English | English | 2022 Website and Assignments<br>2020 Website<br>2020 Assignments |

# Resources
## (Projects)

www.edgeimpulse.com/projects
docs.edgeimpulse.com/experts
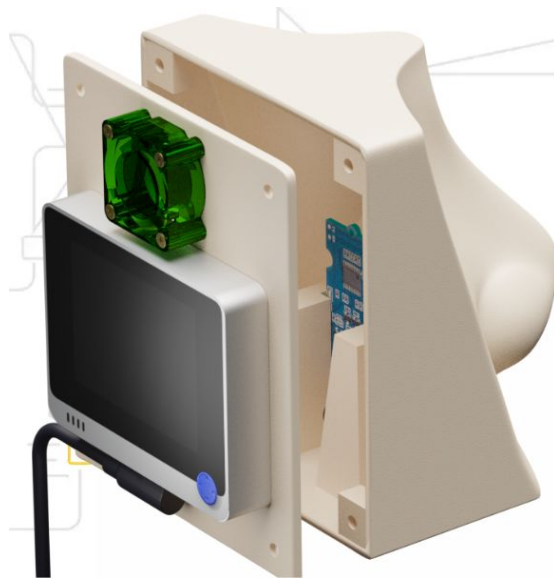
# Project: Smart HVAC

- **Creator**:
  Jallson Suryo
- **Description**:
  Set heating/cooling
  based on number of
  people in each room
- **Hardware**:
  Arduino Nicla Vision
- **Model**:
  FOMO



docs.edgeimpulse.com/experts/featured-machine-learning-projects/arduino-nicla-vision-smart-hvac

# Project: Artificial Nose

- **Creator**:
  Benjamin Cabé
- **Description**:
  Classify different odors
  based on gas data
- **Hardware**:
  Seeed Studio Wio
  Terminal
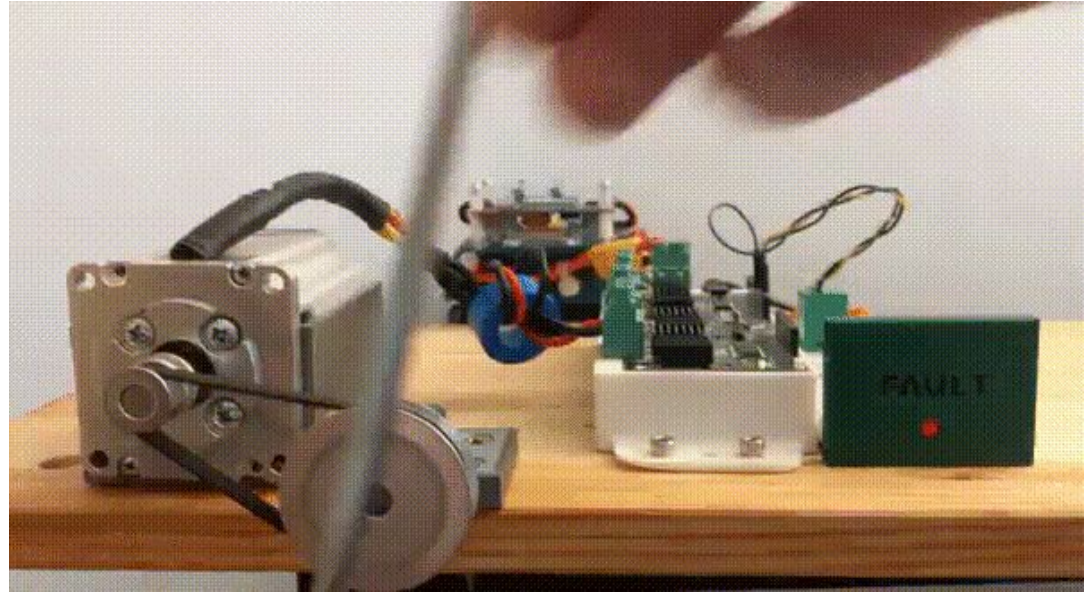- **Model**:
  DNN



TinyML-powered
artificial nose

🦷

kartben/artificial-nose

github.com/kartben/artificial-nose
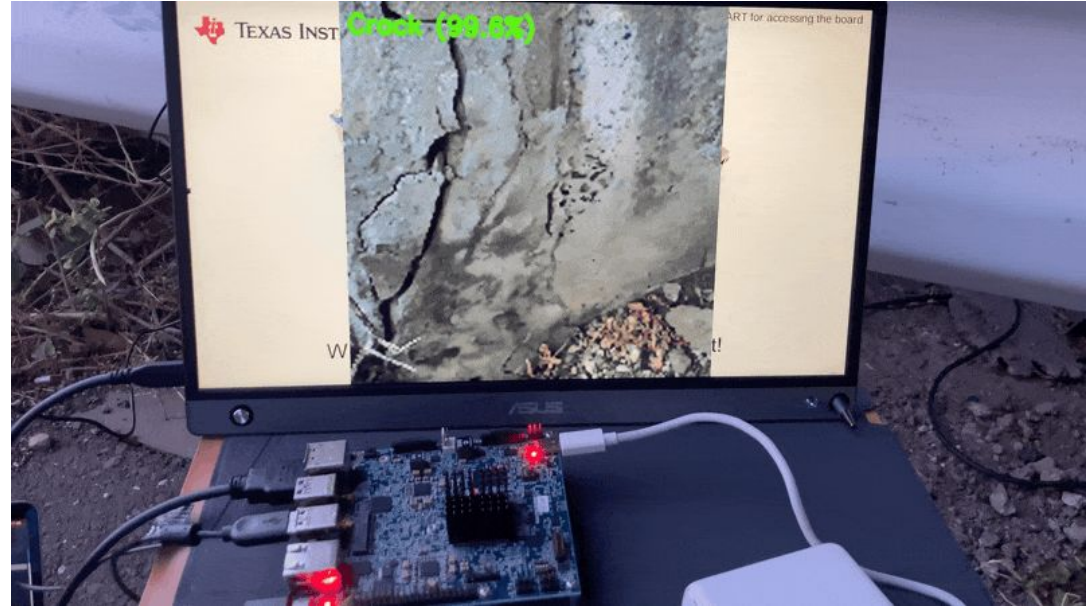
# Project: Motor Anomaly Detection

- **Creator**:
  Avi Brown
- **Description**:
  Identify anomalies
  based on motor current
  and voltage
- **Hardware**:
  Raspberry Pi Pico
- **Model**:
  K-means clustering



docs.edgeimpulse.com/experts/prototype-and-concept-projects/brushless-dc-motor-anomaly-detection

# Project: Concrete Surface Crack Detection

- **Creator**:
  Naveen Kumar
- **Description**:
  Identify surface cracks
  in concrete structures
- **Hardware**:
  TI TDA4VM
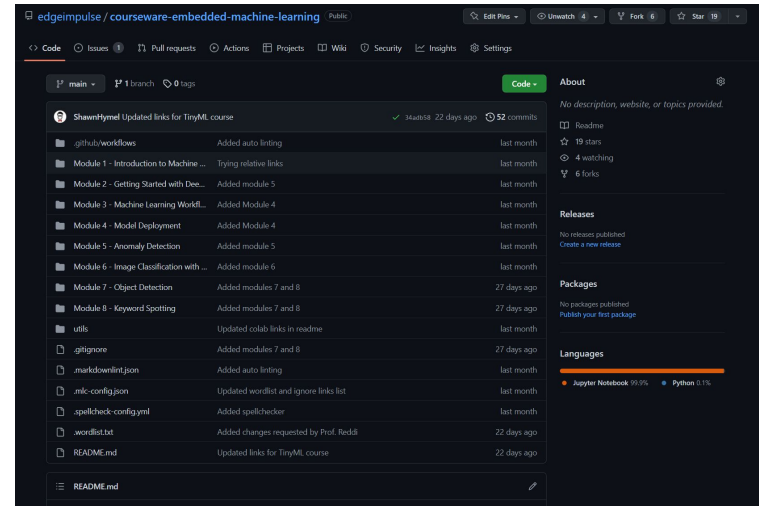- **Model**:
  MobileNetV2 with CAM



docs.edgeimpulse.com/experts/prototype-and-concept-projects/surface-crack-detection-ti-tda4vm

# University Program

**edgeimpulse.com/university**

1. Free hardware kits

2. Content to build curriculum

3. Access to expert network

4. Discount to enterprise edition

**Deadline July 16**

Let's simplify embedded ML for the next generation

of engineers together

Thanks!

**EDGE IMPULSE**

hello@edgeimpulse.com

3031 Tisch Way
110 Plaza West
San Jose, CA 95128
USA