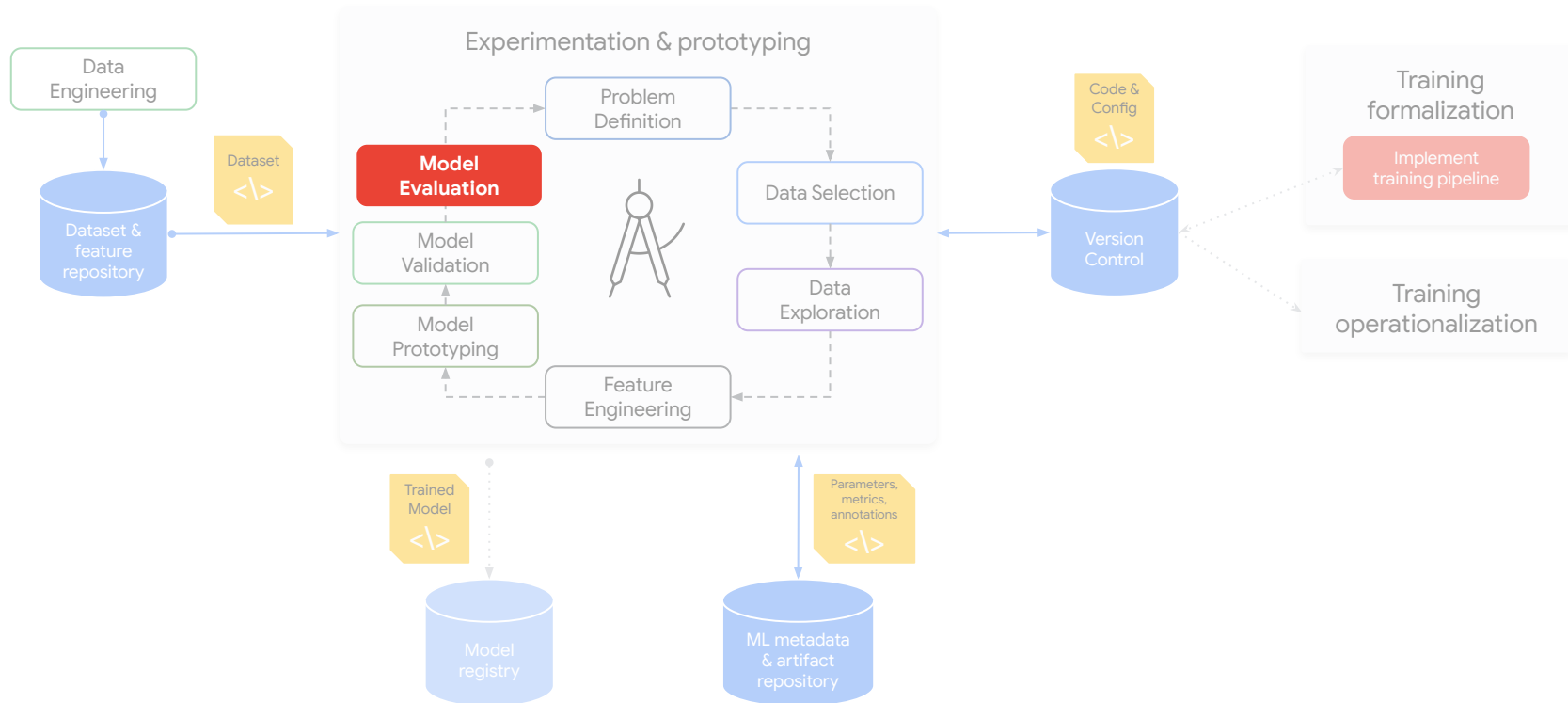# ML Development: Model Evaluation

—

# **MLOps:** ML Development

# The MLOps **Personas**

ML
Engineer

**ML
Researcher**

Data
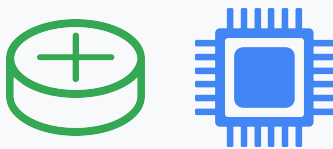Scientist

Data
Engineer

Software
Engineer

DevOps
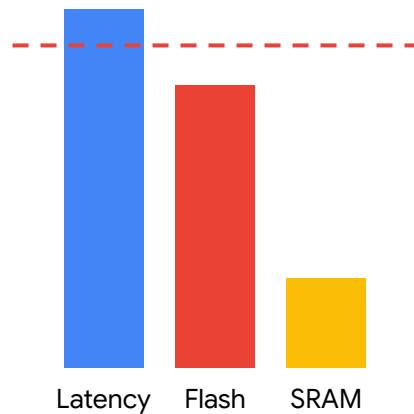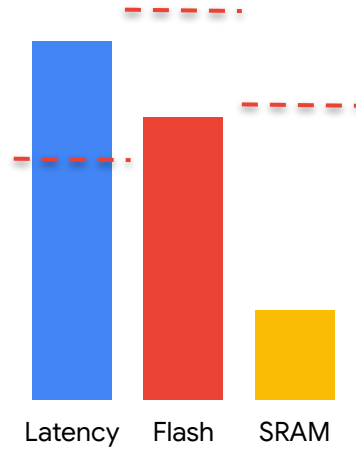
Business
Analyst

# Constraints for **on-device computing**
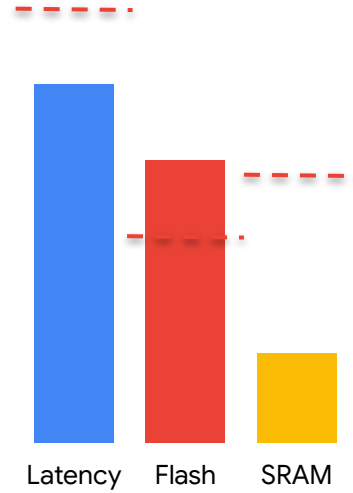
**Latency**

**Limited Devices**

**Battery**

# ML Workflow

| Collect Data | Preprocess Data | Design a Model | Train a Model | Evaluate Optimize | Convert Model | Deploy Model | Make Inferences |

Latency  Flash  SRAM
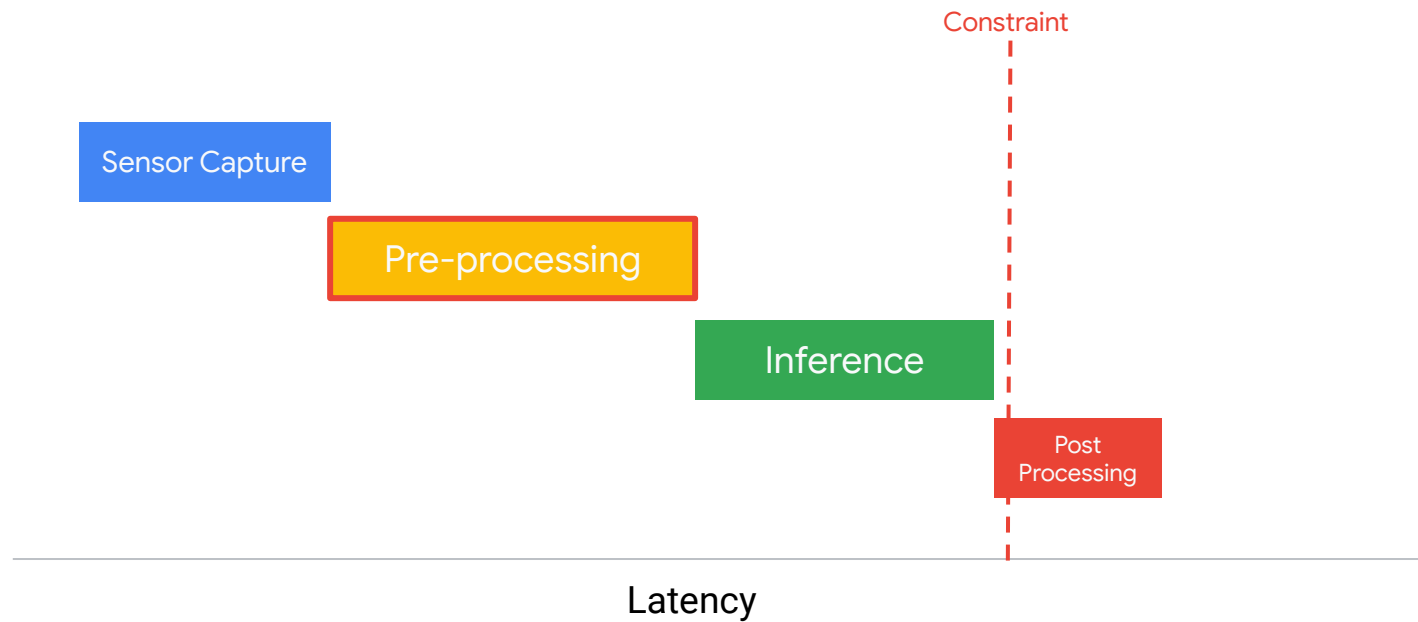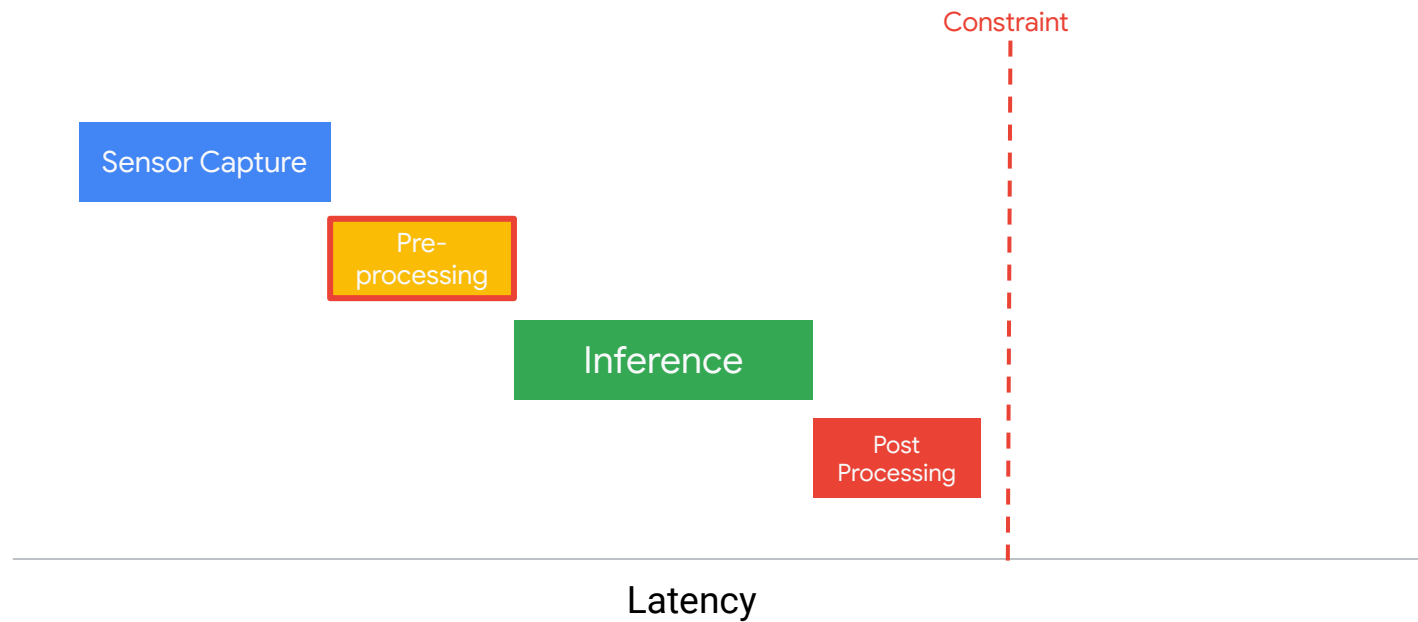
MCU

Latency  Flash  SRAM

DSP

Constraint

Sensor Capture

Pre-processing

Inference

Post Processing

Latency

# Spectrograms v. MFCCs



DSP result

Spectrogram

**Accuracy 80.87%**

Processed features

0.3688, 0.4828, 0.2290, 0.3349, 0.3003, 0.2865, 0.2555,…

On-device performance ⓘ

PROCESSING TIME
123 ms.

PEAK RAM USAGE
28 KB



DSP result

Cepstral Coefficients

**Accuracy 89.59%**

Processed features

−0.2280, −0.5125, −0.5620, −0.2370, −0.4003, 1.0088, 0.…

On-device performance ⓘ

PROCESSING TIME
229 ms.

PEAK RAM USAGE
17 KB

Constraint

Sensor Capture

Pre-processing

Inference

Post Processing

Latency

# Profiling **Metrics**

**Latency**

**Memory**

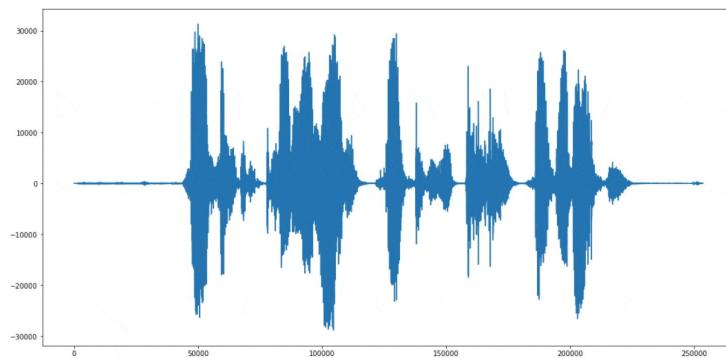**Energy**

# Profiling **Metrics**

**Latency**

Memory

Energy

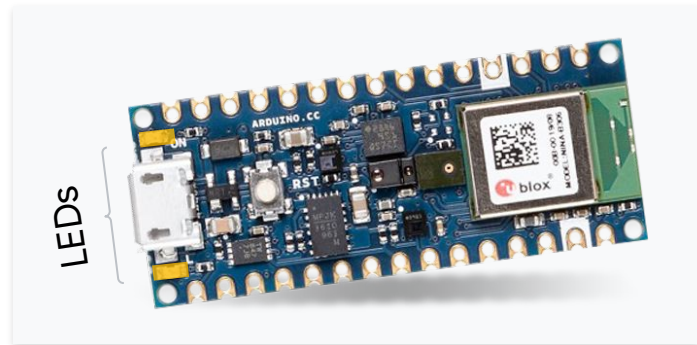**Desired**
Latency



**Deployed**
Latency

# Operation Count



**Math OPs** ~ **Latency**

# Latency Profiling Methods

LEDs

**Internal Timer**

**other tools**

# Latency Profiling Methods



LEDs

Internal Timer

other tools

**Blink**

**Blink**

Sensor Capture

Preprocessing

Inference

Post Processing

Estimated Total Latency

**Blink**

**Blink**

**Blink**

Preprocessing

Inference

Post Processing

Estimated Total Latency

Estimated Total Latency
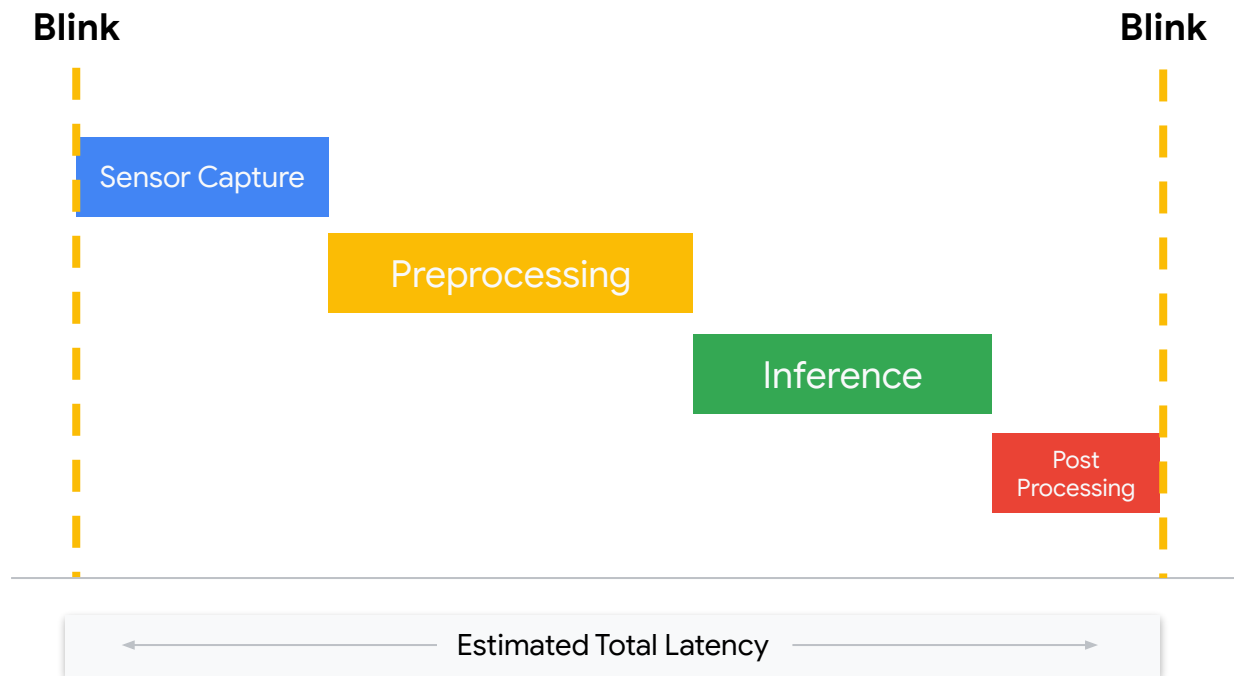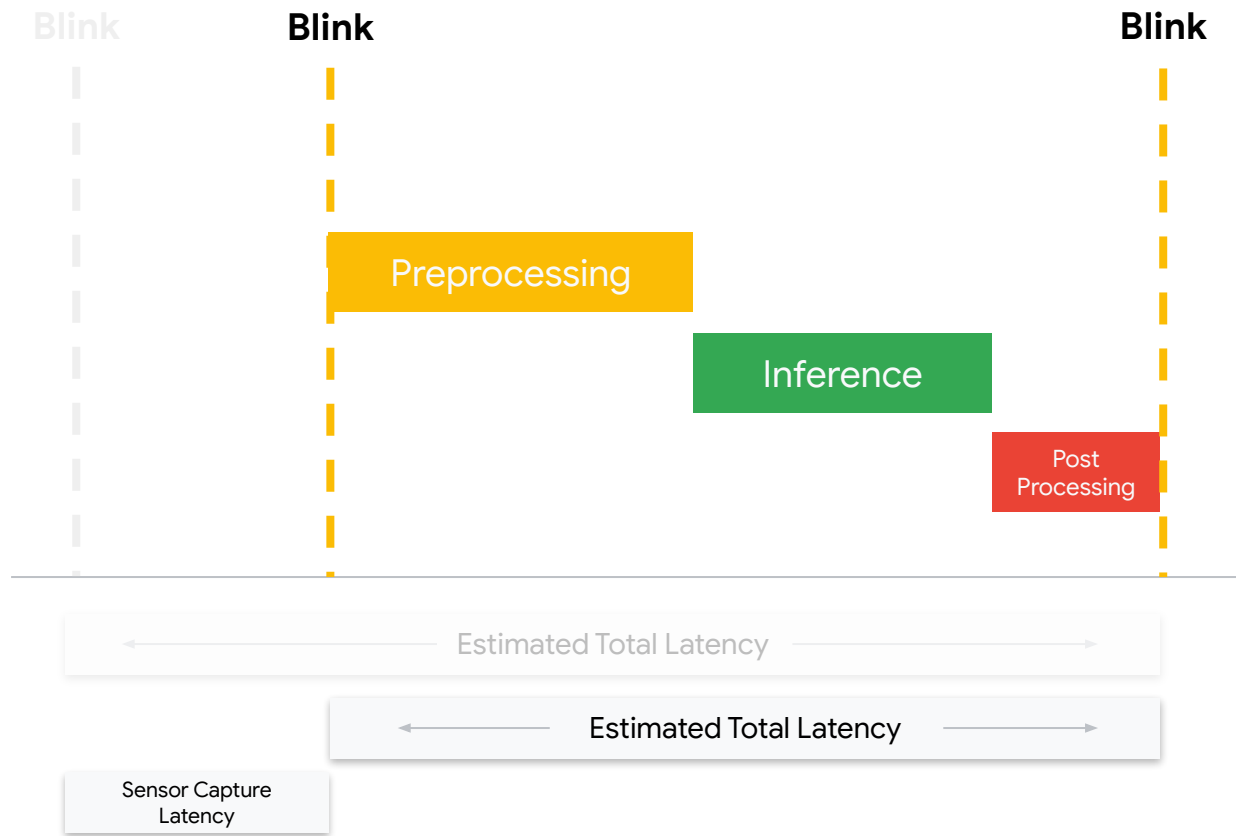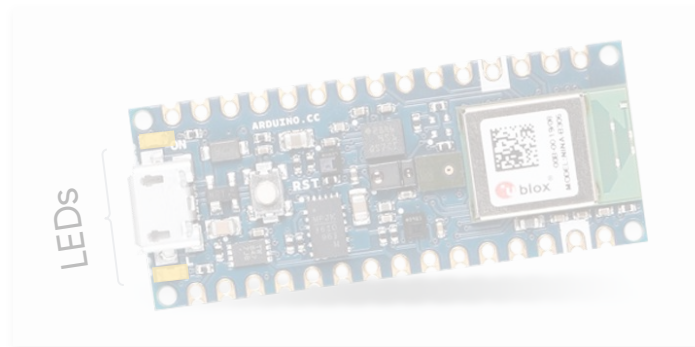
Sensor Capture Latency

# Latency Profiling Methods

LEDs

**Internal Timer**

other tools
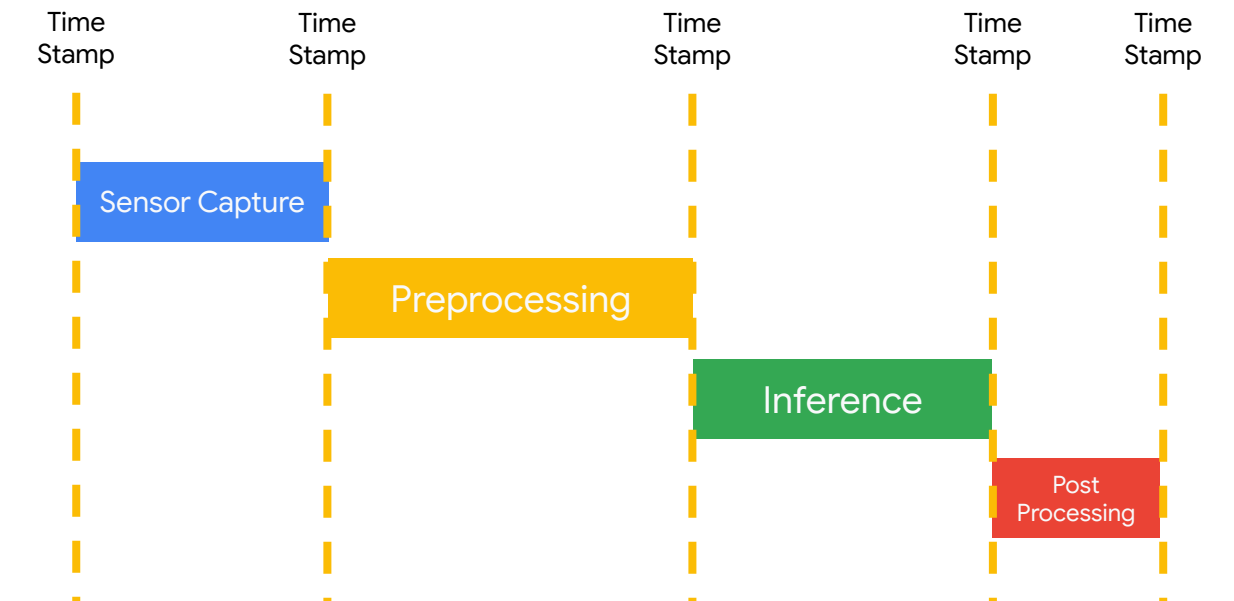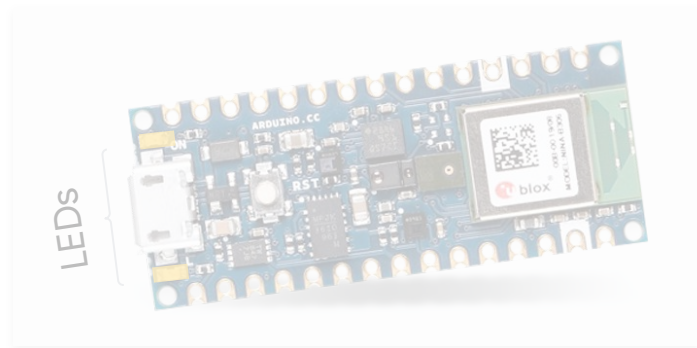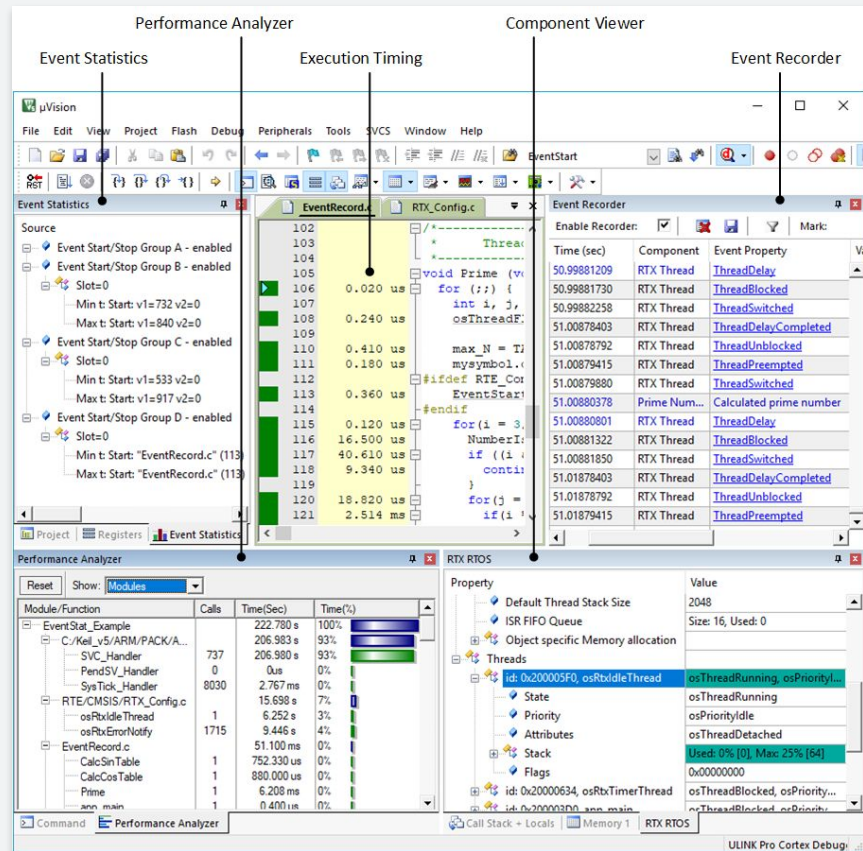
# Latency Profiling Methods

LEDs

Internal
Timer

**other
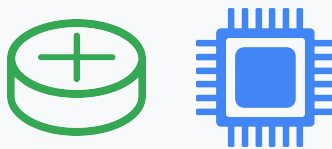tools**

# Profiling **Tools**

- More advanced tools can be used to understand the latency impact at a very **Fine grained level**.
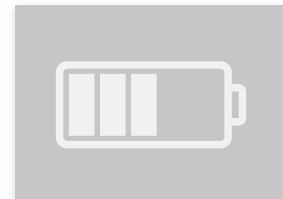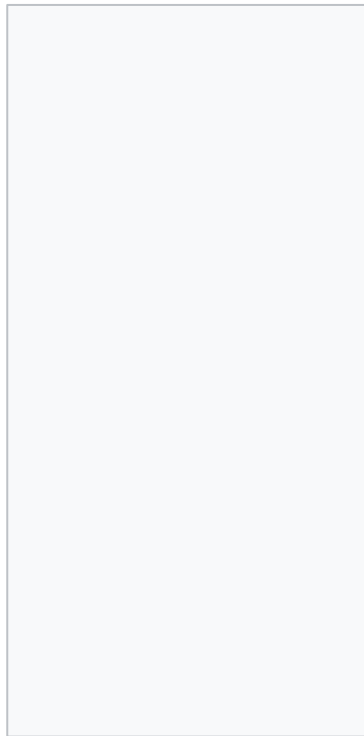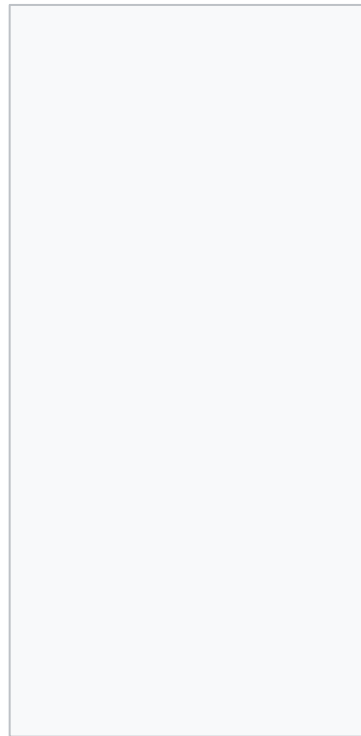
# Profiling **Metrics**
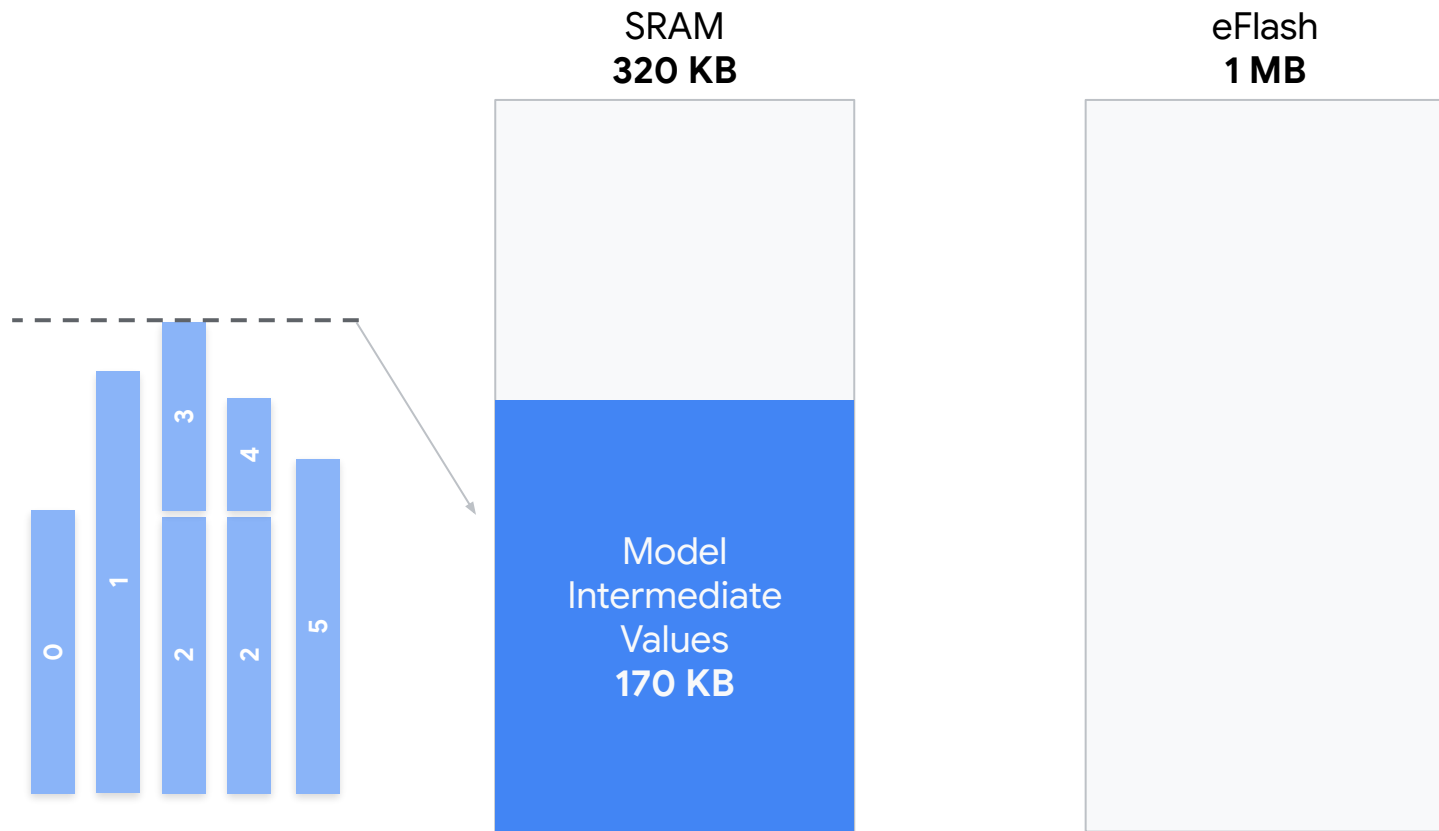
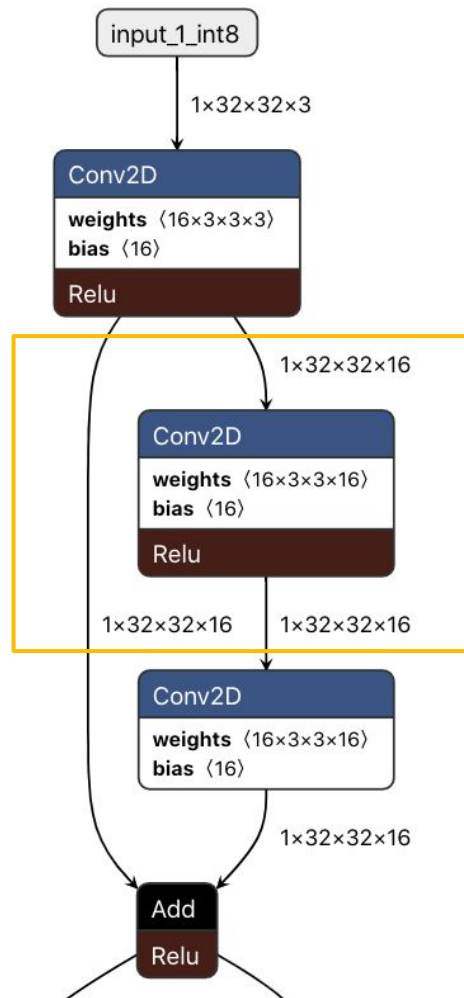Latency

Memory

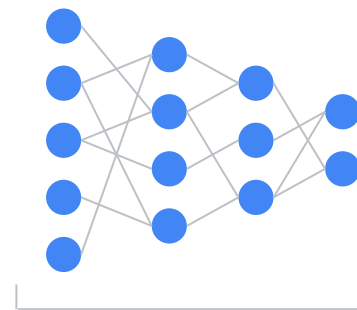Energy

# Memory and Storage

SRAM
**320 KB**

eFlash
**1 MB**

# Memory and Storage

SRAM
**320 KB**

eFlash
**1 MB**

Model
Intermediate
Values
**170 KB**

# Memory and Storage

SRAM
**320 KB**

eFlash
**1 MB**

Model
Intermediate
Values
**170 KB**

Model
Weights
**500 KB**

# Memory and Storage



SRAM
**320 KB**

Bare Metal OS

TFLite Micro

Other Buffers

Model
Intermediate
Values
**170 KB**

eFlash
**1 MB**

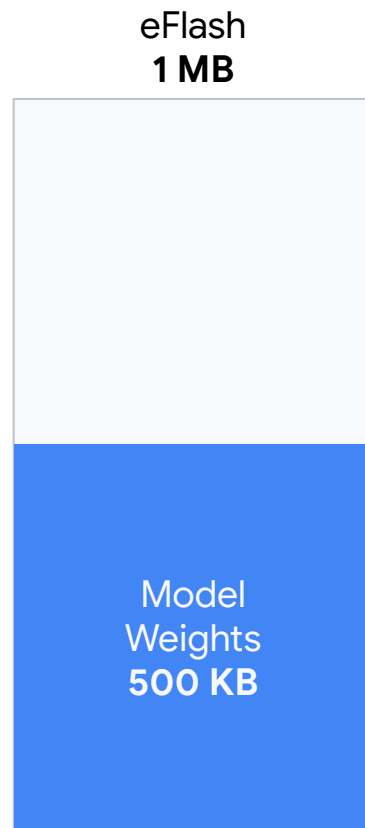Bare Metal OS

TFLite Micro

Quantization Parameters

Model
Weights
**500 KB**

# Memory and Storage

# Memory and Storage

**SRAM**
**320 KB**

Bare Metal OS

TFLite Micro

Other Buffers

Model
Intermediate
Values
**170 KB**

Measure

**eFlash**
**1 MB**

Bare Metal OS

TFLite Micro

Quantization Parameters

Model
Weights
**500 KB**

Measure

# Memory and Storage

eFlash
**1 MB**

Bare Metal OS

TFLite Micro

Quantization Parameters

Model
Weights
**500 KB**

**Size**

.tflite

# Memory and Storage



SRAM
**320 KB**

Bare Metal OS

TFLite Micro

Other Buffers

Model
Intermediate
Values
**170 KB**

# Memory and Storage

SRAM
**320 KB**

Bare Metal OS

TFLite Micro

Other Buffers

Model
Intermediate
Values
**170 KB**

Recording
Micro
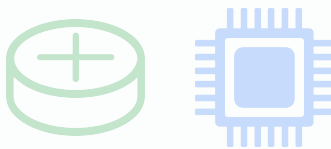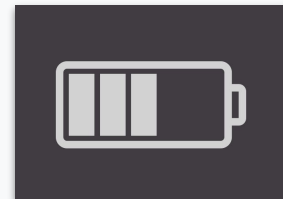Interpreter

# Profiling Metrics

Latency

Memory

**Energy**

# **Estimating** Energy

Estimate by **latency**
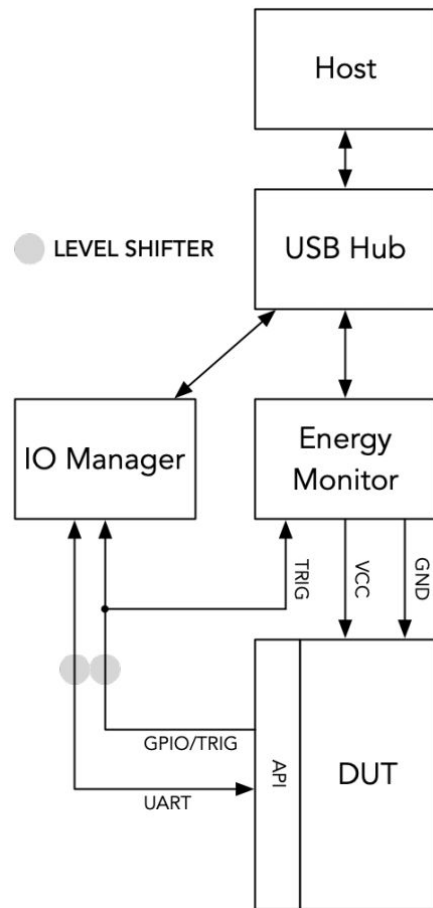- Rough estimate
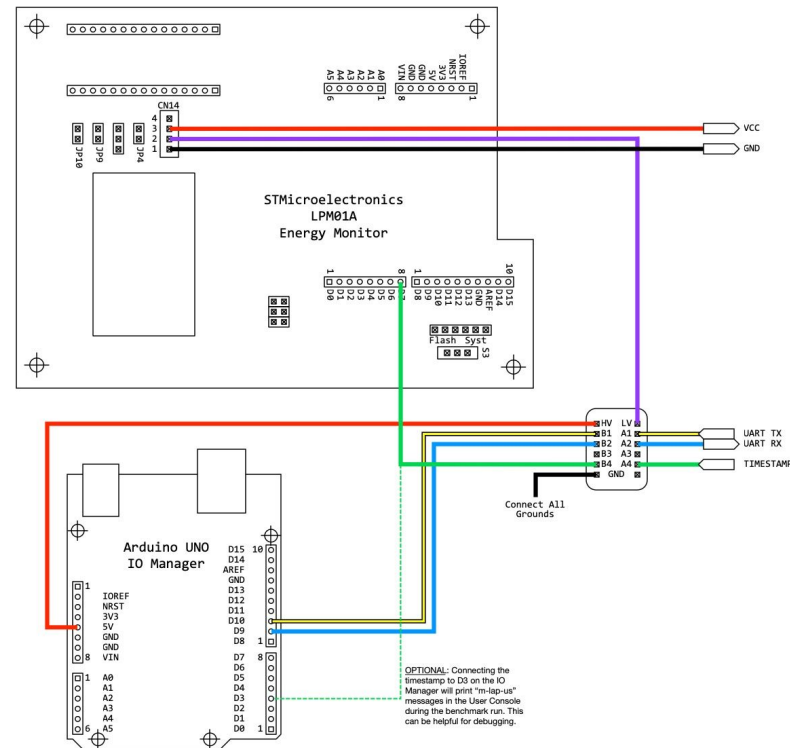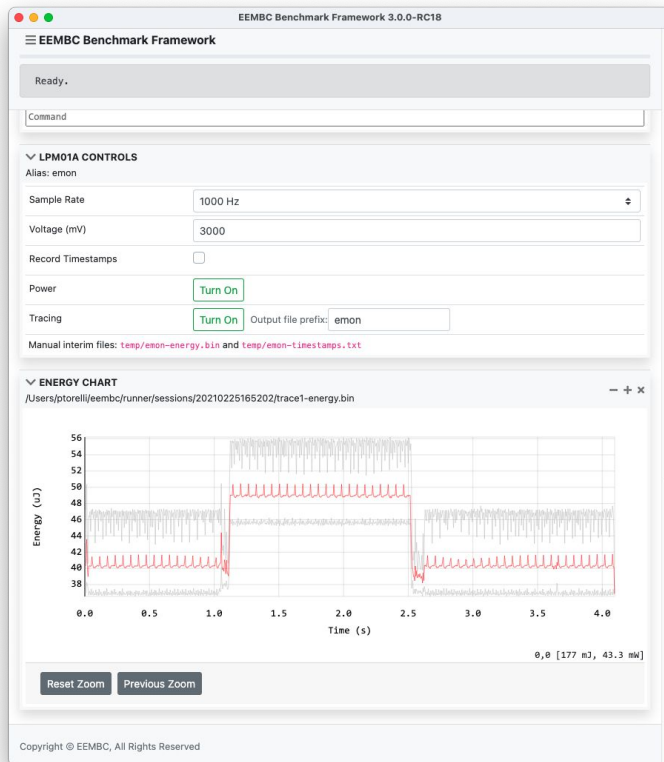- Only relative to other models



Latency

# **Measuring** Energy

Accurately measuring energy
is **complex**:

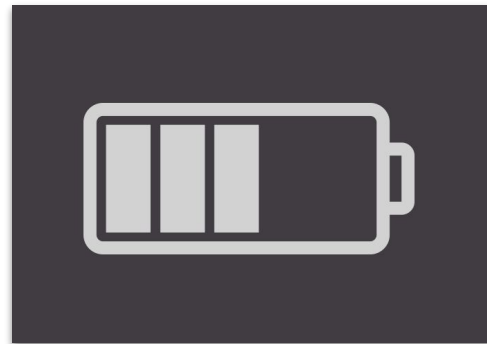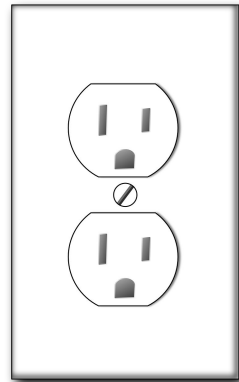- Isolating out I/O and power
  planes can be complicated
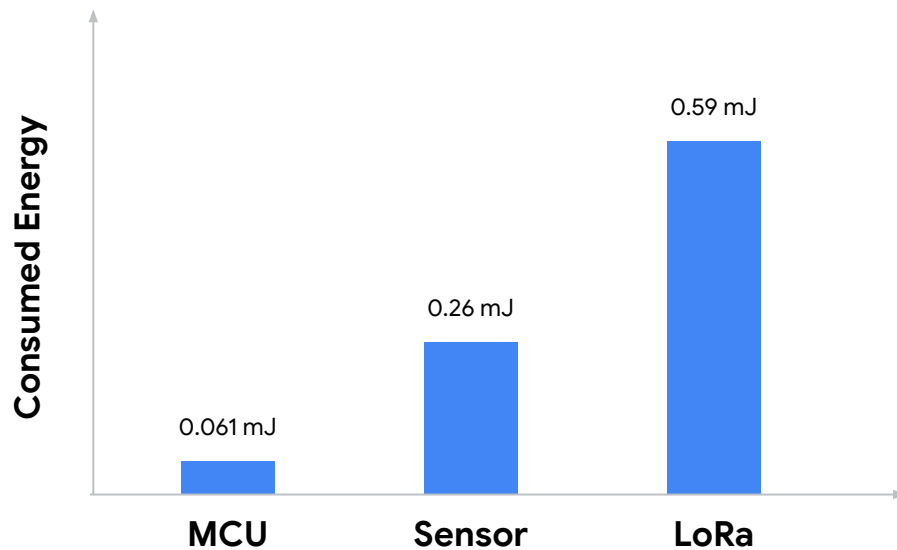
# EEMBC's EnergyRunner™

# **Deployment** Scenario

**Battery** Powered:

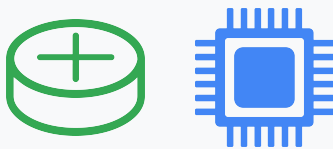- Size of battery?
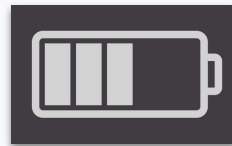- How often is it charged?
- Energy Harvesting?

# Other Factors



Bar chart titled "Consumed Energy" showing: MCU 0.061 mJ, Sensor 0.26 mJ, LoRa 0.59 mJ

# Constraints for **on-device computing**



**Latency**

**Limited Devices**

**Battery**