# MicroNets

Neural Network Architectures for Deploying TinyML Application on Commodity Microcontrollers

**Colby Banbury***[1,2], Chuteng Zhou*[1], Igor Fedorov*[1,] Ramon Matas Navarro[1], Urmish Thakker[3], Dibakar Gope[1], Vijay Janapa Reddi[1], Matthew Mattina[1], Paul N. Whatmough[1]

1 **arm**

2 HARVARD
John A. Paulson
School of Engineering
and Applied Sciences

3 SambaNova
SYSTEMS

# What is TinyML?

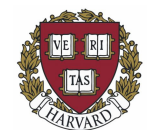## ML Inference at <1mWatt

# IoT Paradigm
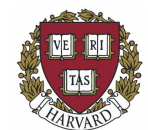
"Smart" Devices
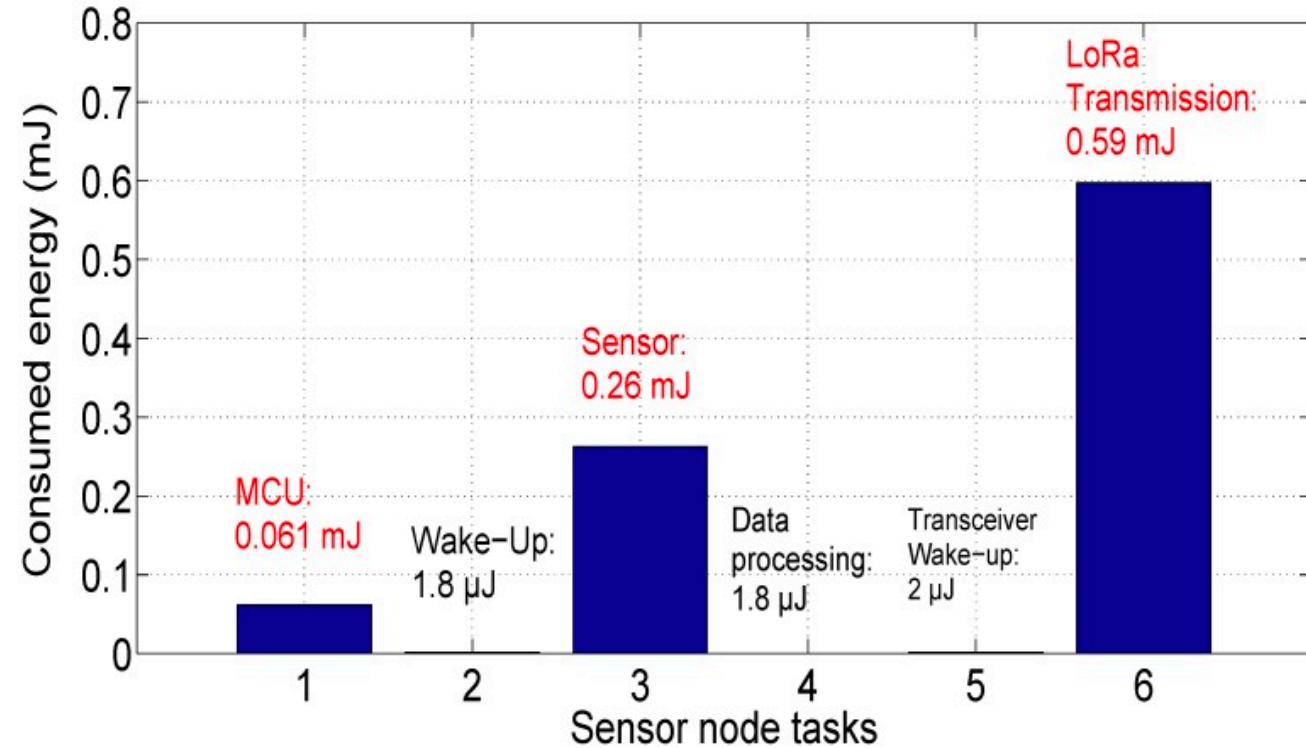=
Everything is collecting data

# IoT's Fatal Flaw

# <1%
## of IoT data is analyzed or used at all

Source: Mckinsey Global Institute. "The Internet of Things: Mapping the Value Beyond the Hype." mckinsey.com

# IoT's Fatal Flaw



Bouguera, Taoufik et al. "Energy Consumption Model for Sensor Nodes Based on LoRa and LoRaWAN." *Sensors (Basel, Switzerland)* vol. 18,7 2104. 30 Jun. 2018, doi:10.3390/s18072104

# Transmission is Energy Hungry

# The On-Device Advantage

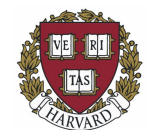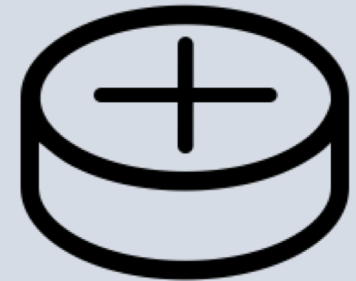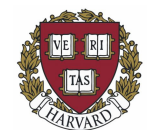+ Energy Efficiency

+ Responsiveness

+ Privacy

+ Mobility

# The On-Device Advantage

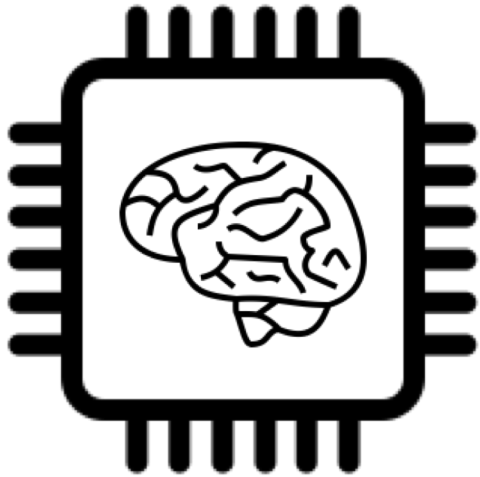**+ Energy Efficiency**

# Goal

## MicroNets

Create **efficient** and **deployable** model
architectures for **tinyML** applications

# Executive Summary

TinyML

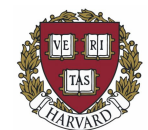**TinyML systems have severe constraints and require highly tuned model architectures**

SRAM

Flash

Latency
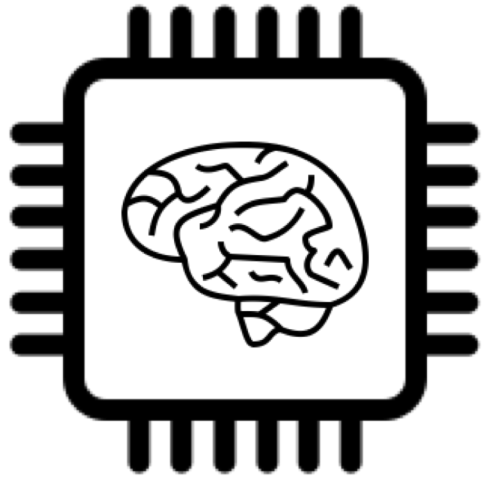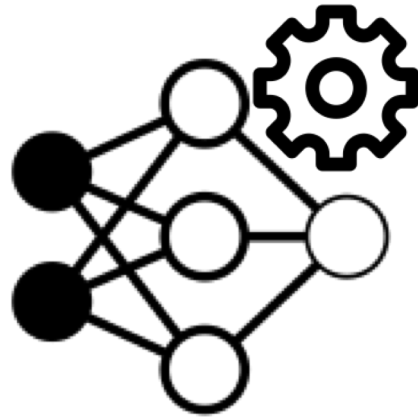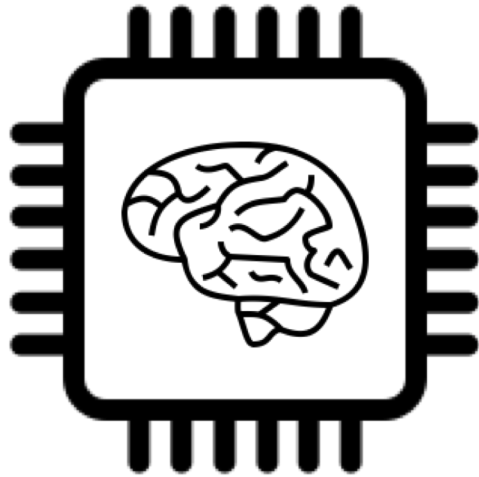
Energy

# Executive Summary

TinyML

Differentiable
Neural Architecture
Search

Differentiable Neural
Architecture Search (DNAS)
can **rapidly** find models that
**meet the constraints** given
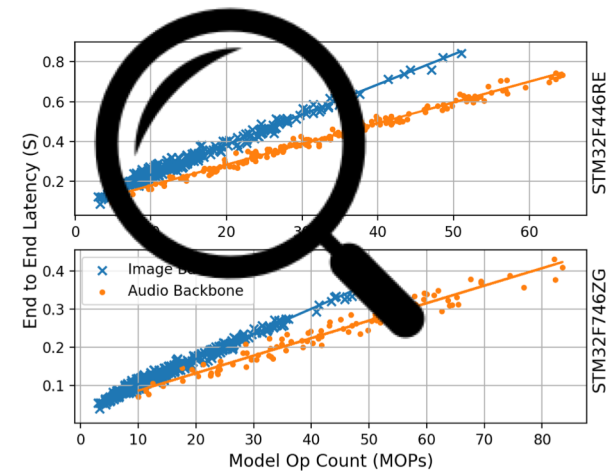**viable proxies**

# Executive Summary

TinyML

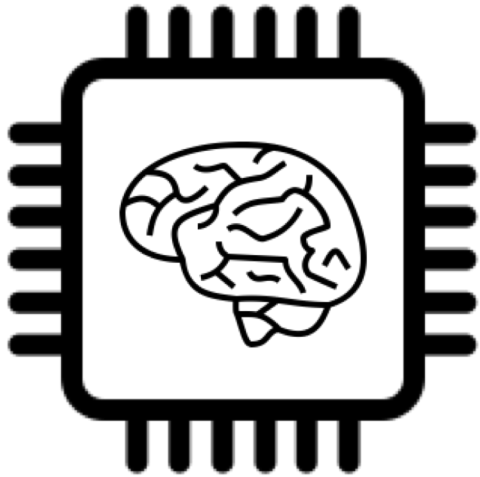Differentiable Neural Architecture Search
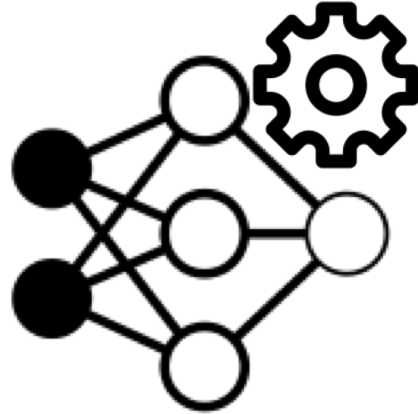
Hardware Characterization



SRAM and Flash are easily calculated while **Op count is a viable proxy** latency and energy
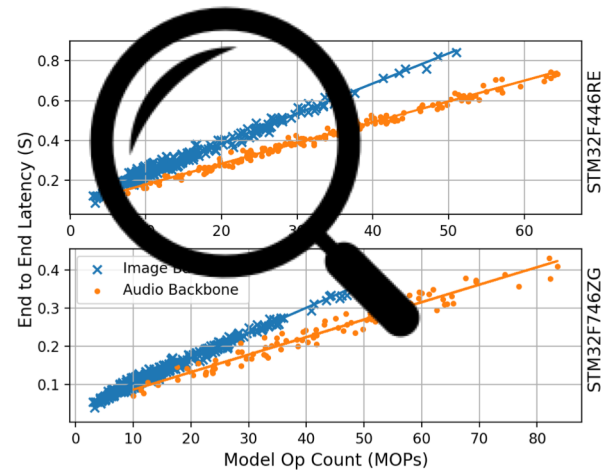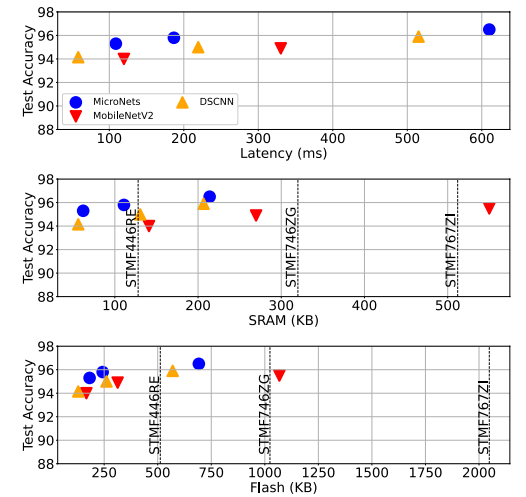
# Executive Summary

TinyML

Differentiable Neural Architecture Search
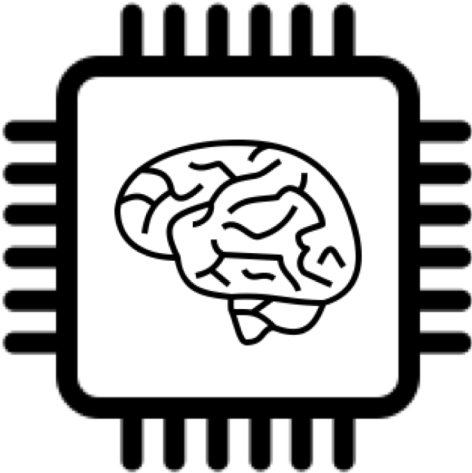
Hardware Characterization

MicroNets



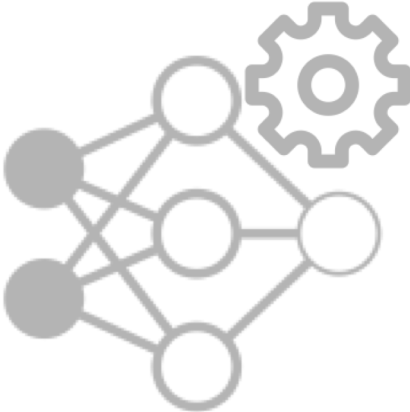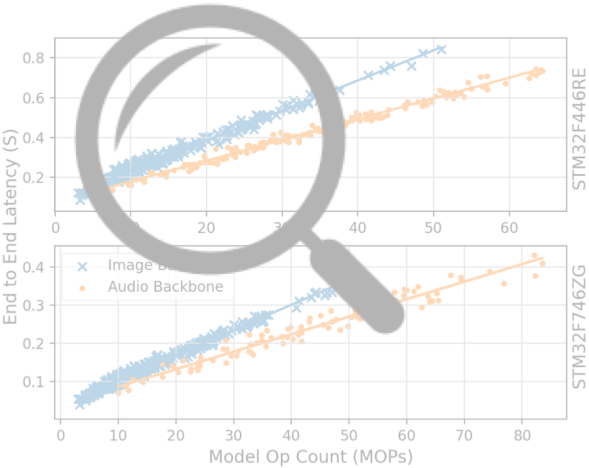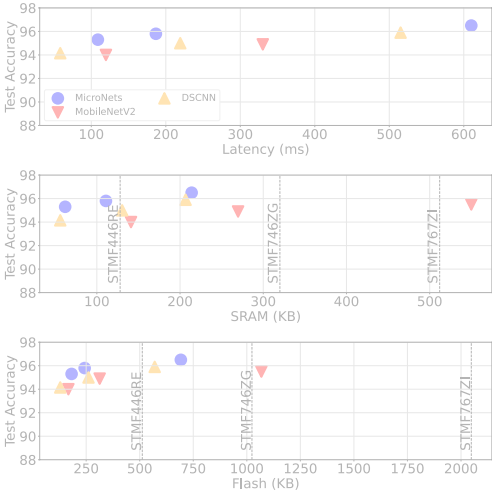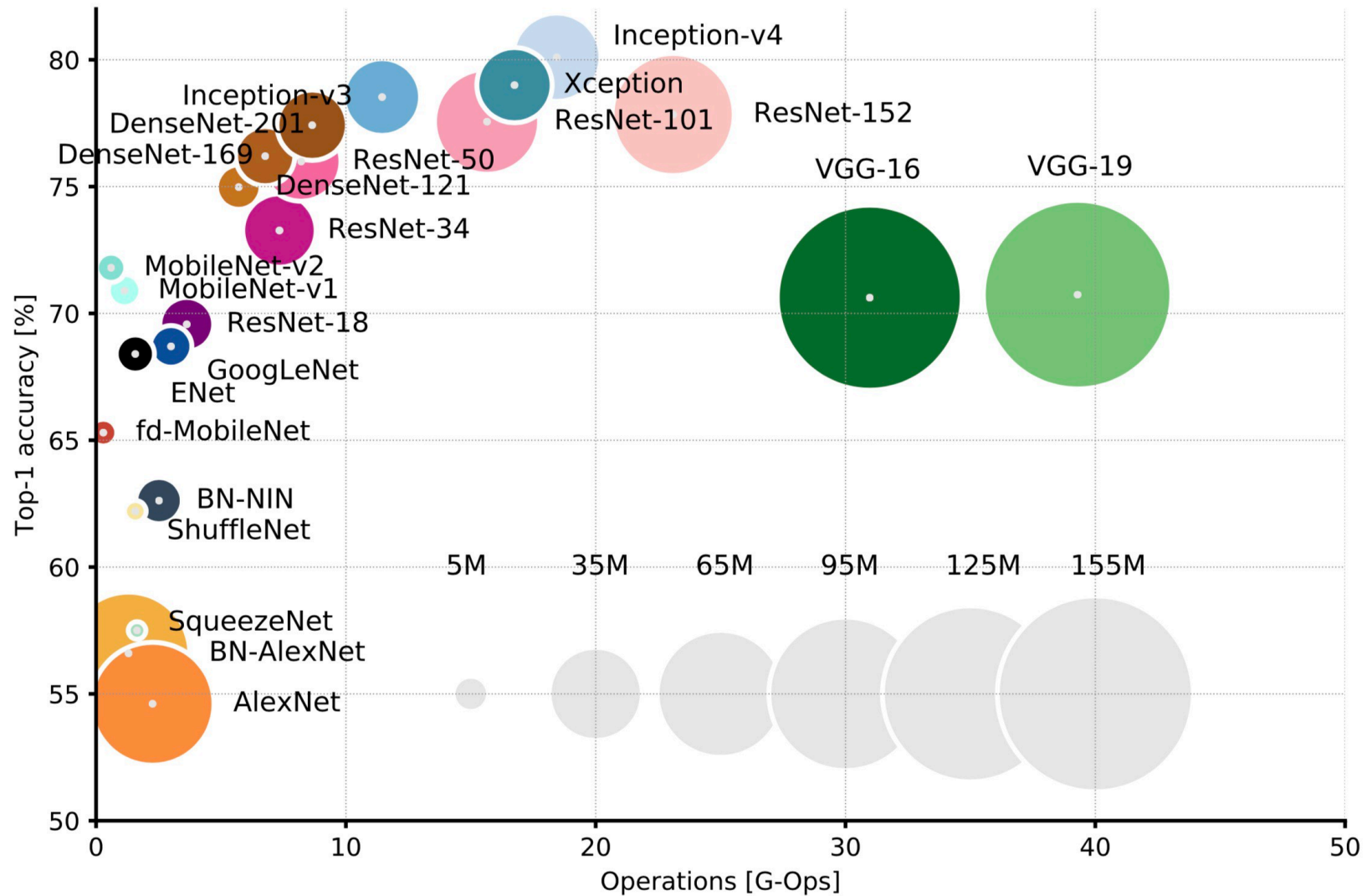We achieve **state of the art performance** on three TinyML tasks

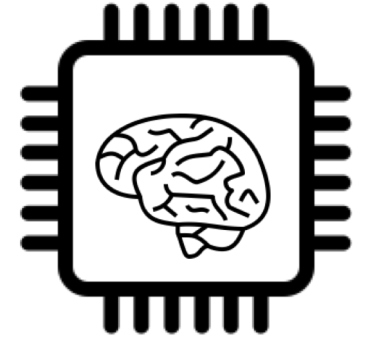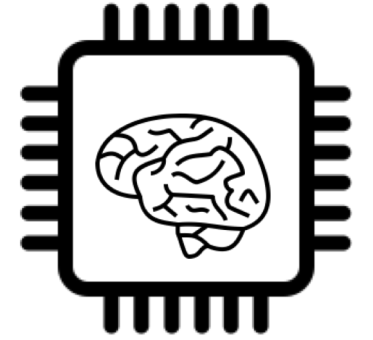# Executive Summary

TinyML

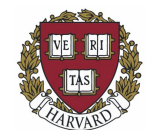Differentiable Neural Architecture Search

Hardware Characterization

MicroNets

# The Challenge



Source: https://culurciello.medium.com/analysis-of-deep-neural-networks-dcf398e71aae

# The Challenge



Source: https://culurciello.medium.com/analysis-of-deep-neural-networks-dcf398e71aae

# The Constraints

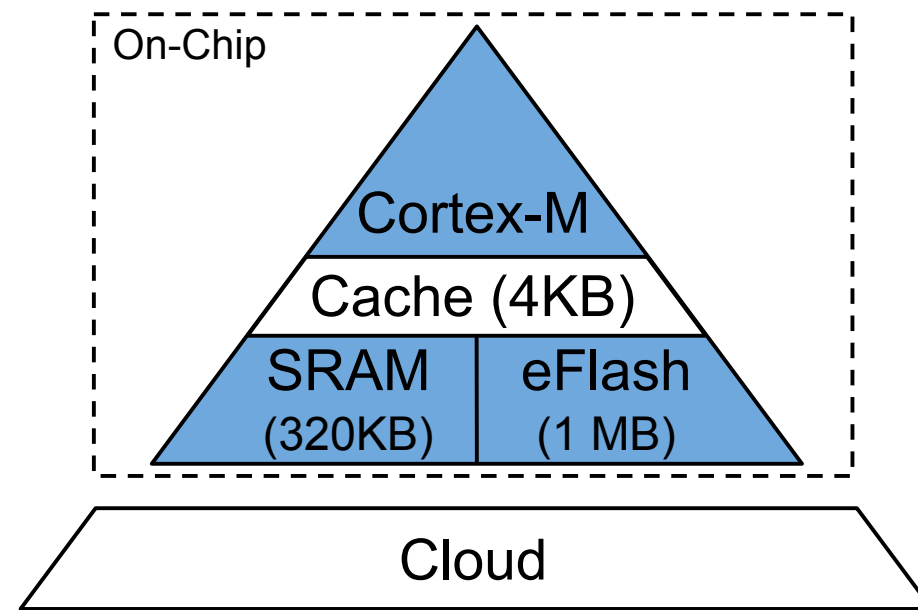| Platform | Architecture | Memory | Storage | Power | Price |
|---|---|---|---|---|---|
| **CloudML**<br>Nvidia V100 | GPU<br>Nvidia Volta | HBM<br>16GB | SSD/Disk<br>TB~PB | 250W | $9K |
| **MobileML**<br>Cell Phone | CPU<br>Mobile CPU | DRAM<br>4GB | Flash<br>64GB | ~8W | ~$750 |
| **TinyML**<br>F446RE<br>F746ZG<br>F767ZI | MCU<br>Arm M4<br>Arm M7<br>Arm M7 | SRAM<br>128KB<br>320KB<br>512KB | eFlash<br>0.5MB<br>1MB<br>2MB | 0.1W<br>0.3W<br>0.3W | $3<br>$5<br>$8 |

# The Constraints

Mobile

Tiny

On-Chip

Cortex-A
L1 (48KB)
L2 (4 MB)
DRAM (4GB)
Flash (64GB)

On-Chip

Cortex-M
Cache (4KB)
SRAM (320KB)  eFlash (1 MB)

Cloud

arm
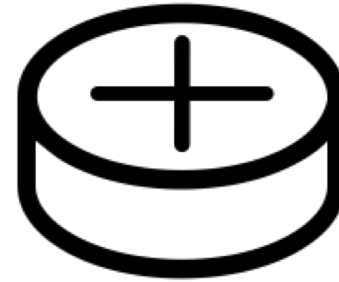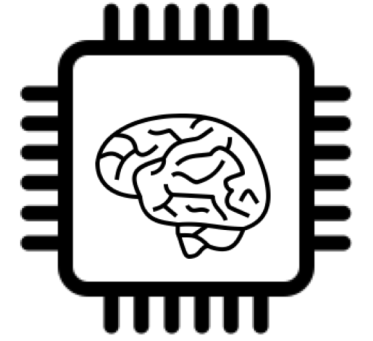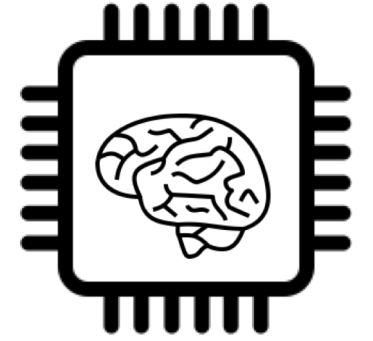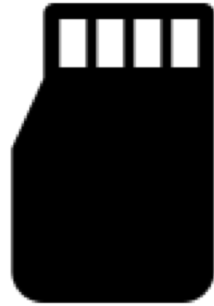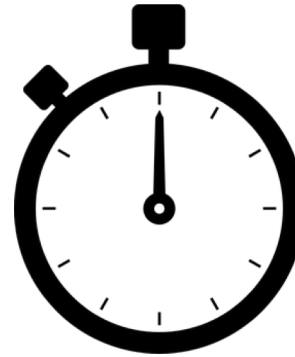
# The Constraints

# The Constraints

# TinyML Constraints
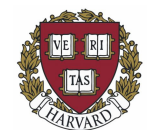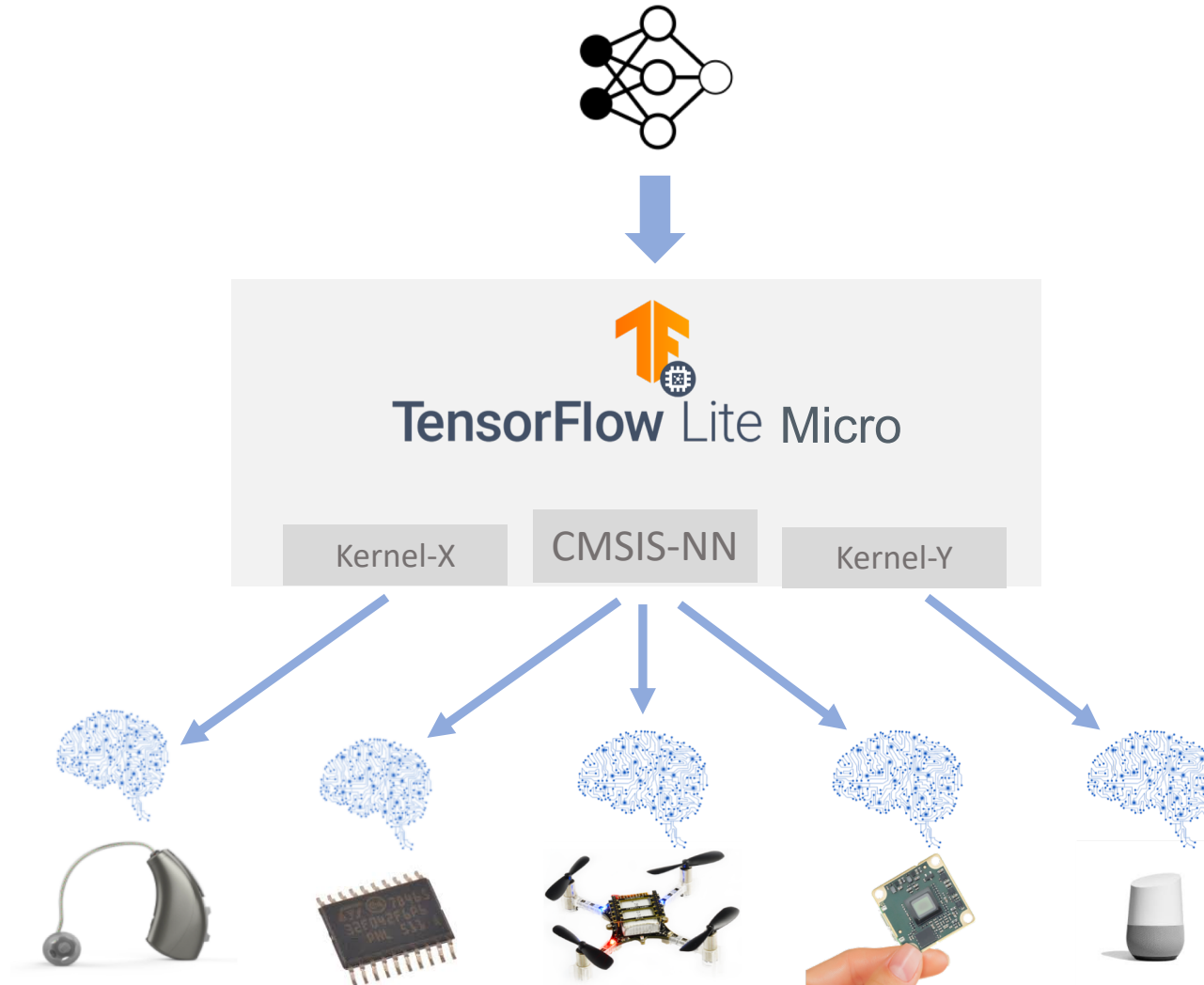
SRAM             Flash             Latency             Energy
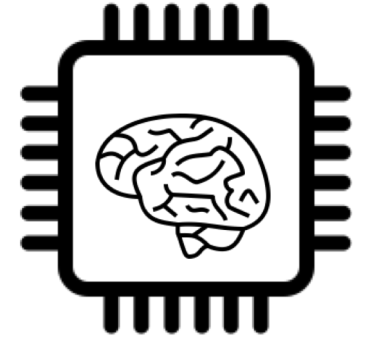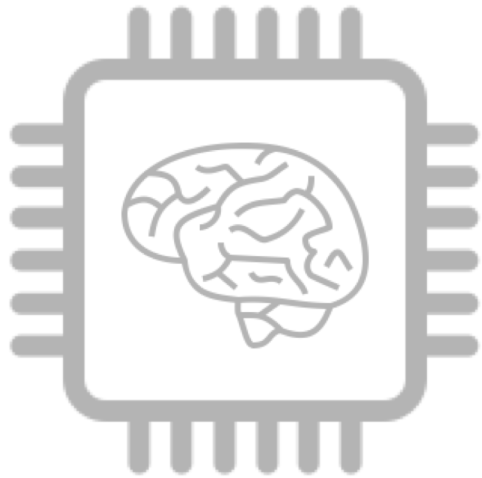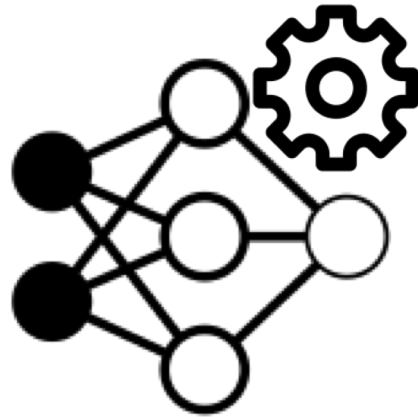
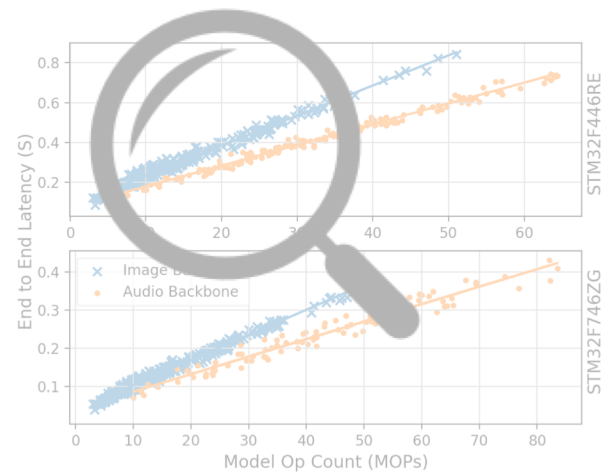On-Chip

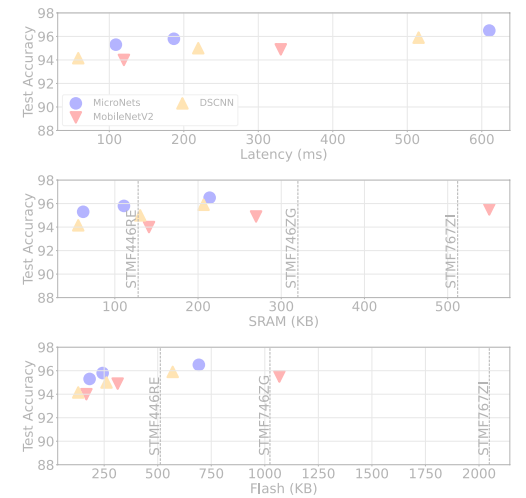# TensorFlow Lite for Microcontrollers

# Executive Summary
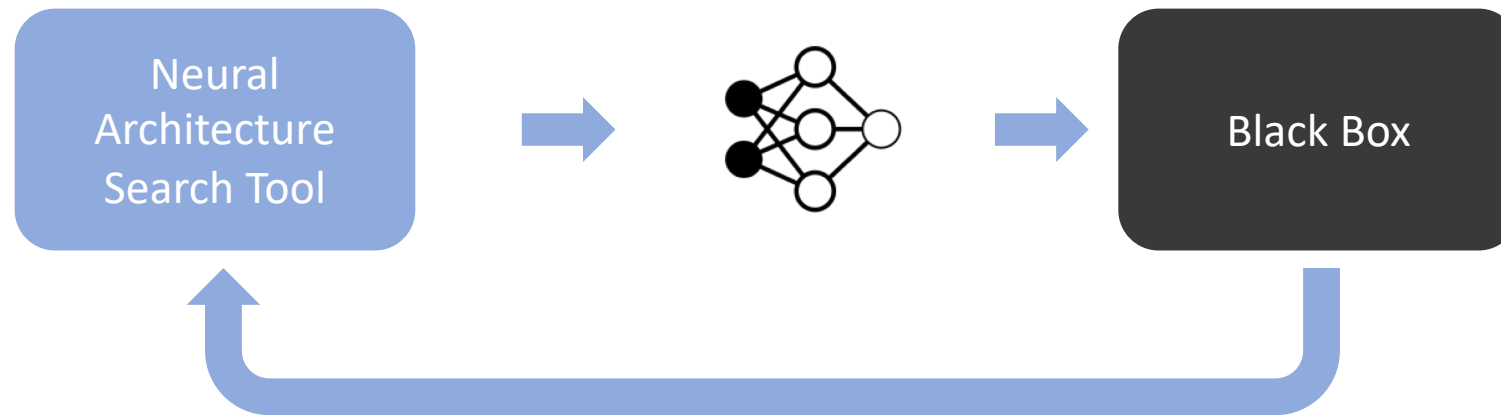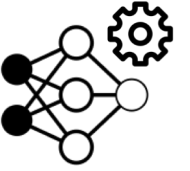
TinyML

Differentiable Neural Architecture Search
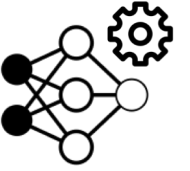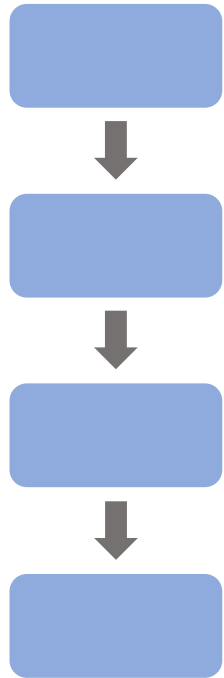
Hardware Characterization

MicroNets

# Neural Architecture Search (NAS)
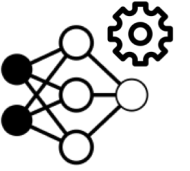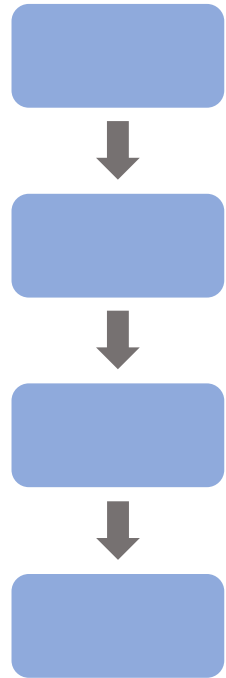
# Differentiable Neural Architecture Search (DNAS)

Existing Model

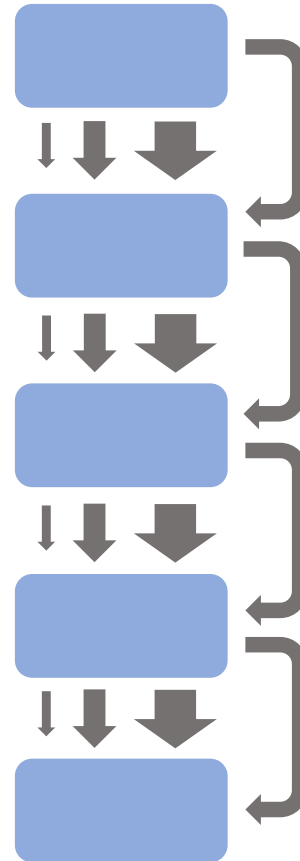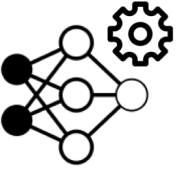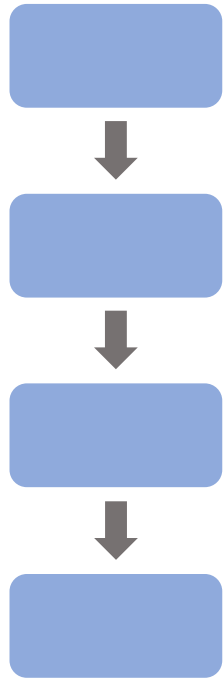# Differentiable Neural Architecture Search (DNAS)

Existing Model

Super Net Backbone

Relaxation

# Differentiable Neural Architecture Search (DNAS)

Existing Model
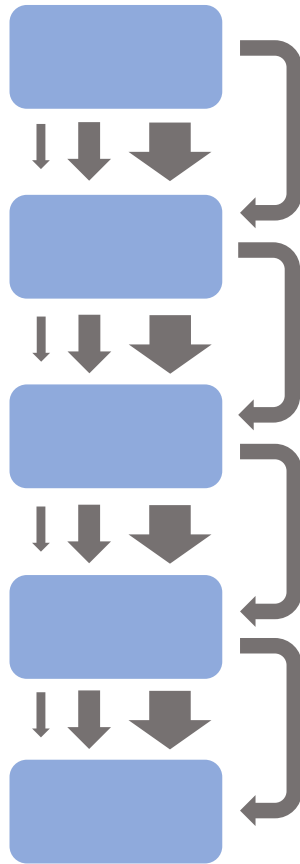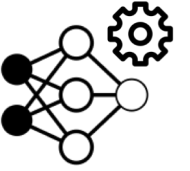
Super Net Backbone

Final Architecture

Relaxation

Gradient
Descent
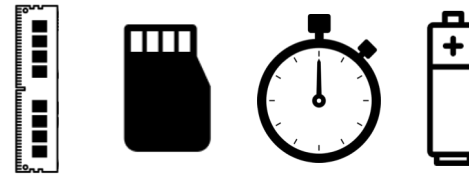
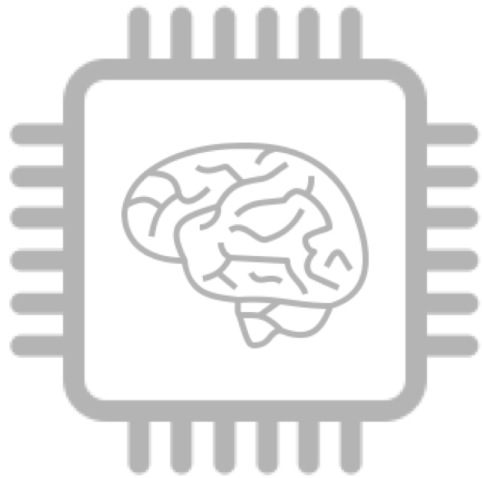# Differentiable Neural Architecture Search (DNAS)

DNAS is **fast** but needs
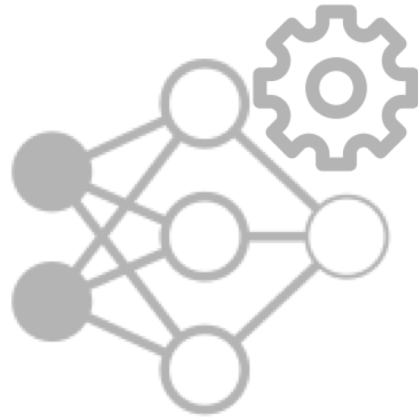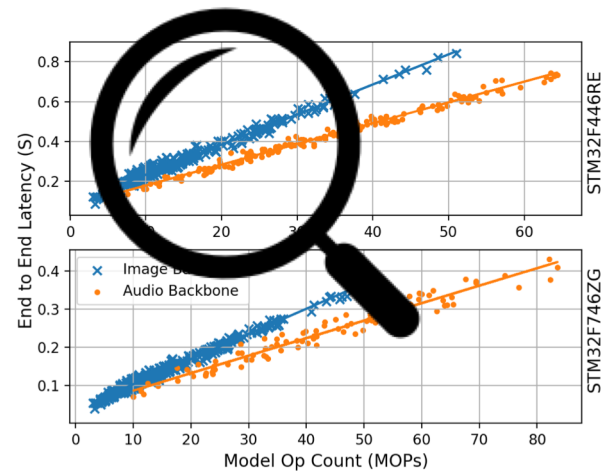**continuous functions** for the
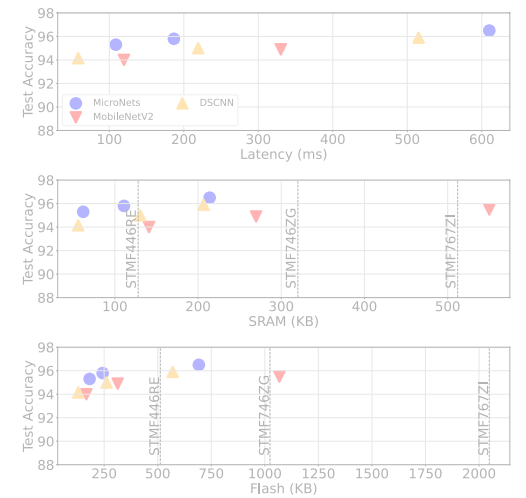hardware objectives

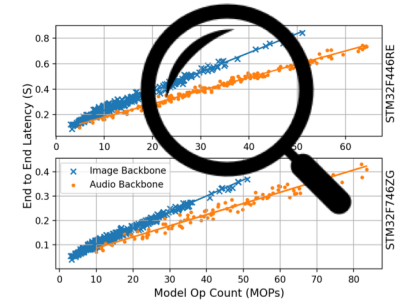# Executive Summary
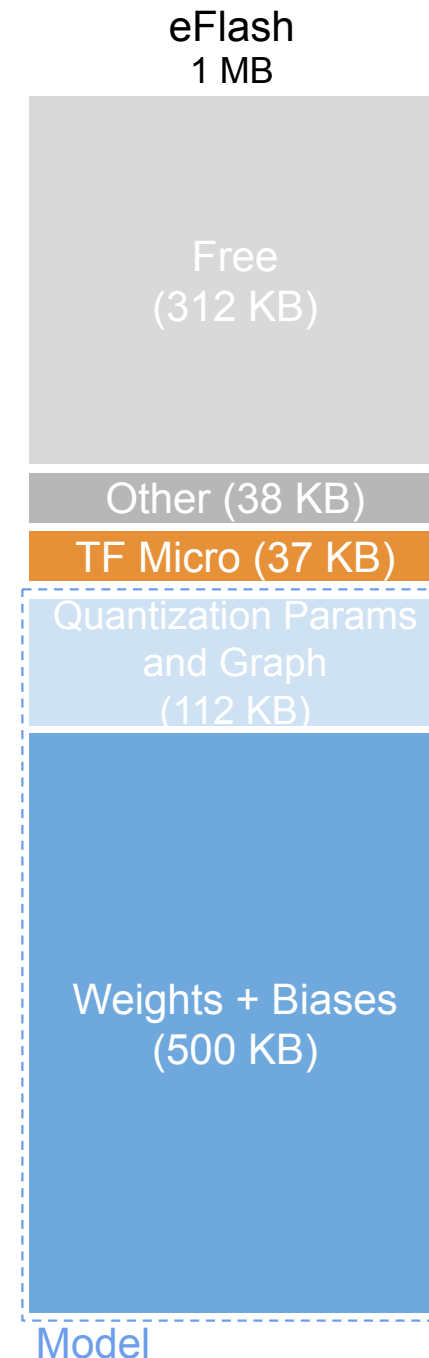


TinyML

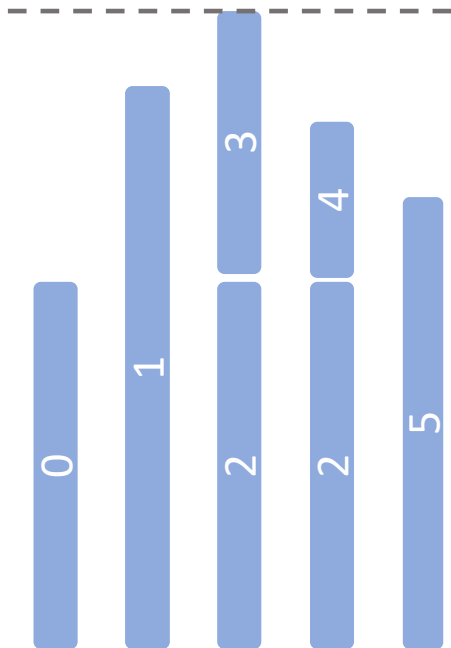Differentiable
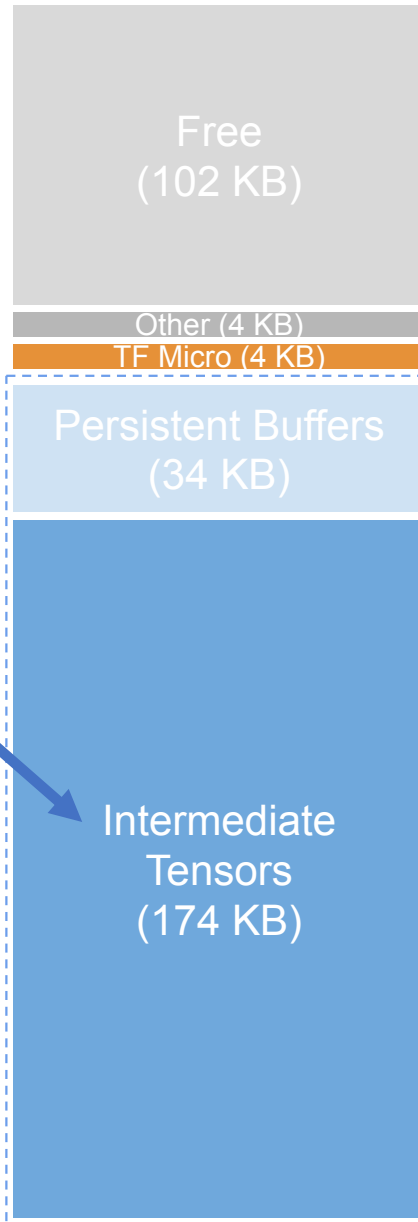Neural Architecture
Search

Hardware
Characterization

MicroNets

# SRAM and Flash

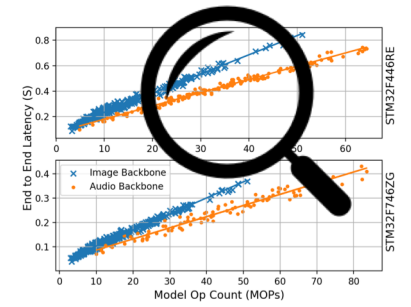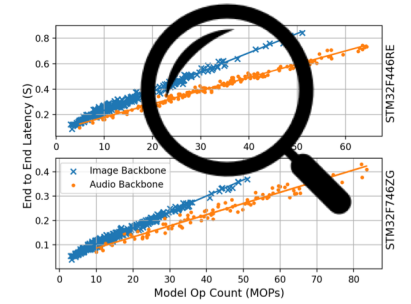**SRAM**
320 KB

- Free (102 KB)
- Other (4 KB)
- TF Micro (4 KB)
- Persistent Buffers (34 KB)
- Intermediate Tensors (174 KB)

Model

**eFlash**
1 MB

- Free (312 KB)
- Other (38 KB)
- TF Micro (37 KB)
- Quantization Params and Graph (112 KB)
- Weights + Biases (500 KB)

Model

# SRAM



SRAM
320 KB

- Free (102 KB)
- Other (4 KB)
- TF Micro (4 KB)
- Persistent Buffers (34 KB)
- Intermediate Tensors (174 KB)

Model

eFlash
1 MB

- Free (312 KB)
- Other (38 KB)
- TF Micro (37 KB)
- Quantization Params and Graph (112 KB)
- Weights + Biases (500 KB)

Model

# Flash

SRAM
320 KB

Free
(102 KB)

Other (4 KB)

TF Micro (4 KB)

Persistent Buffers
(34 KB)

Intermediate
Tensors
(174 KB)

Model

eFlash
1 MB

Free
(312 KB)

Other (38 KB)

TF Micro (37 KB)

Quantization Params
and Graph
(112 KB)

Weights + Biases
(500 KB)

Model

# SRAM and Flash

**SRAM**
320 KB

Free
(102 KB)

Other (4 KB)

TF Micro (4 KB)

Persistent Buffers
(34 KB)

Intermediate
Tensors
(174 KB)

Model

**eFlash**
1 MB

Free
(312 KB)

Other (38 KB)

TF Micro (37 KB)

Quantization Params
and Graph
(112 KB)

Weights + Biases
(500 KB)

Model

Overhead
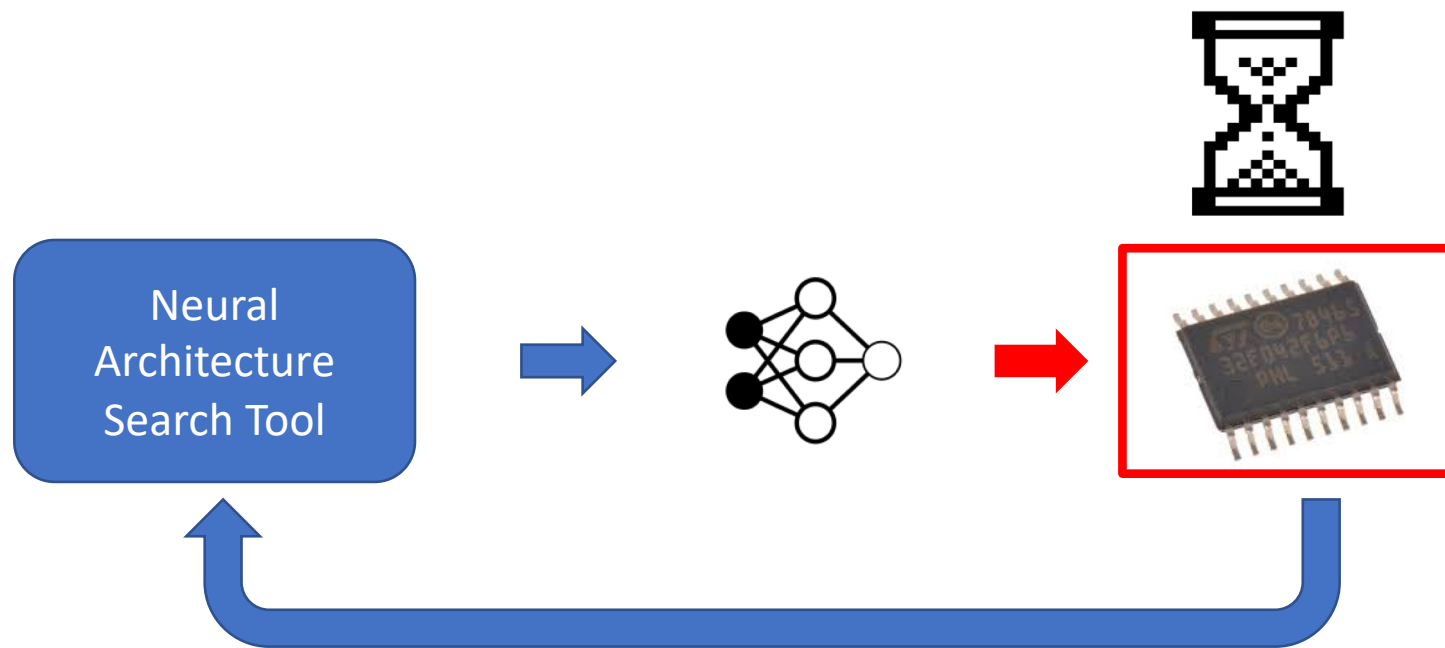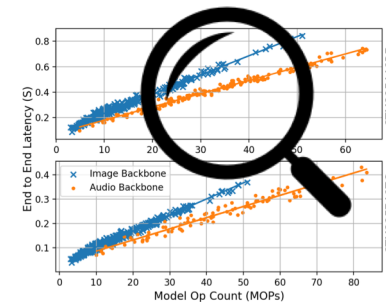
Determined by
the Model
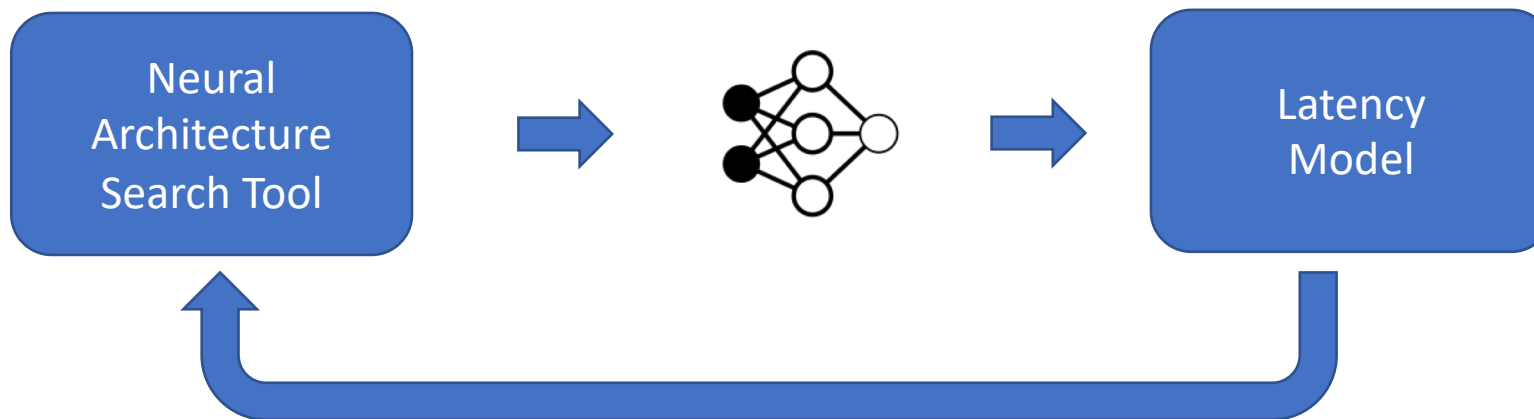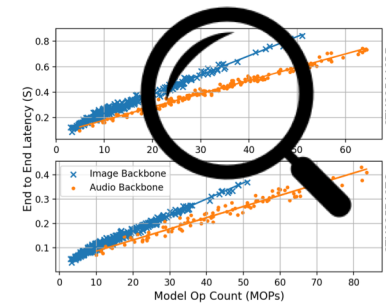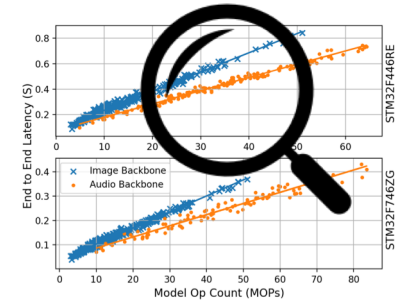Architecture

# TinyML Constraints

✓ SRAM     ✓ Flash     Latency     Energy
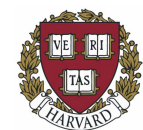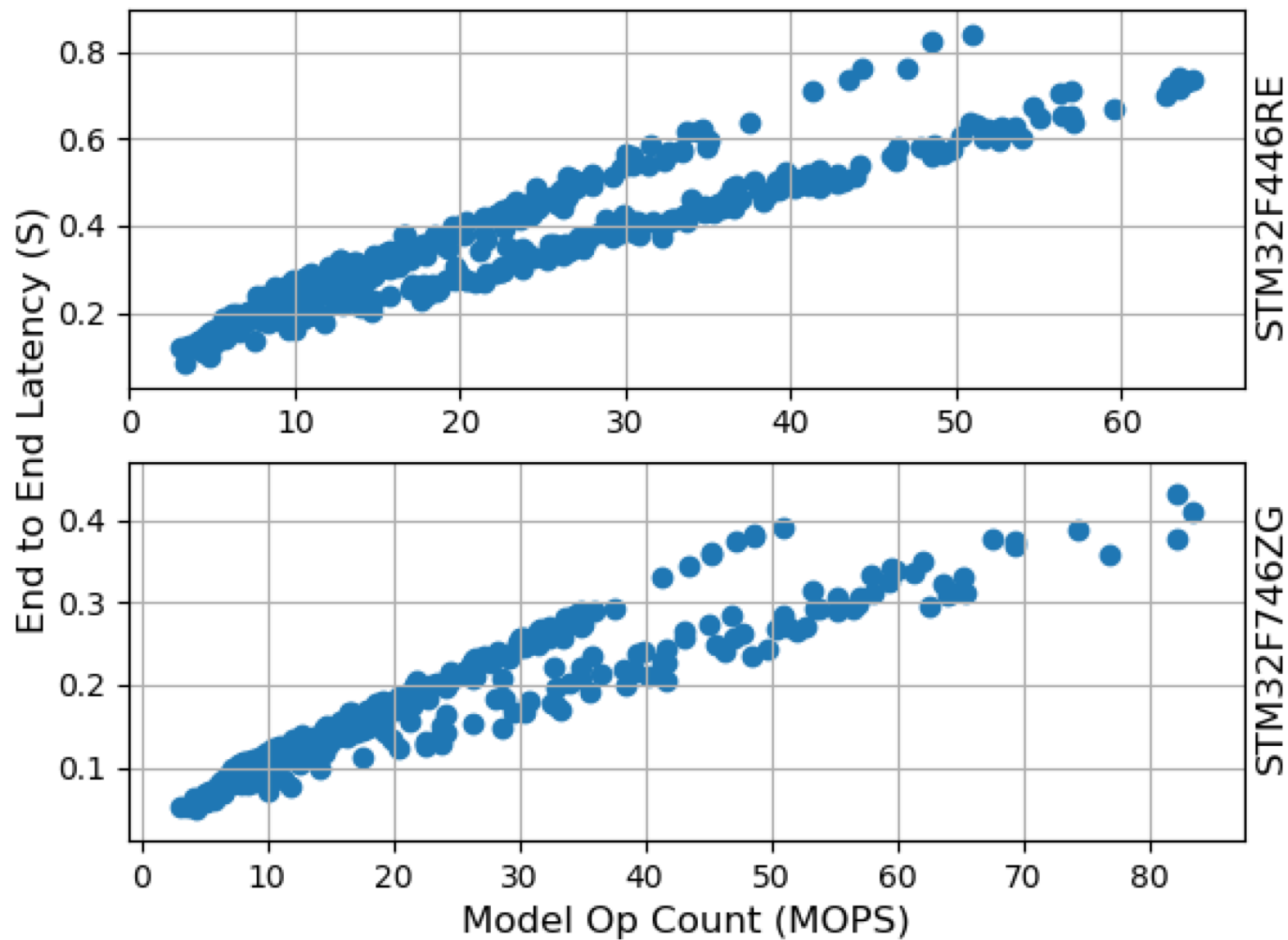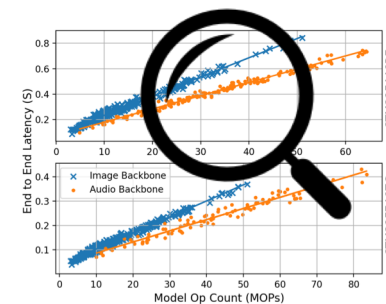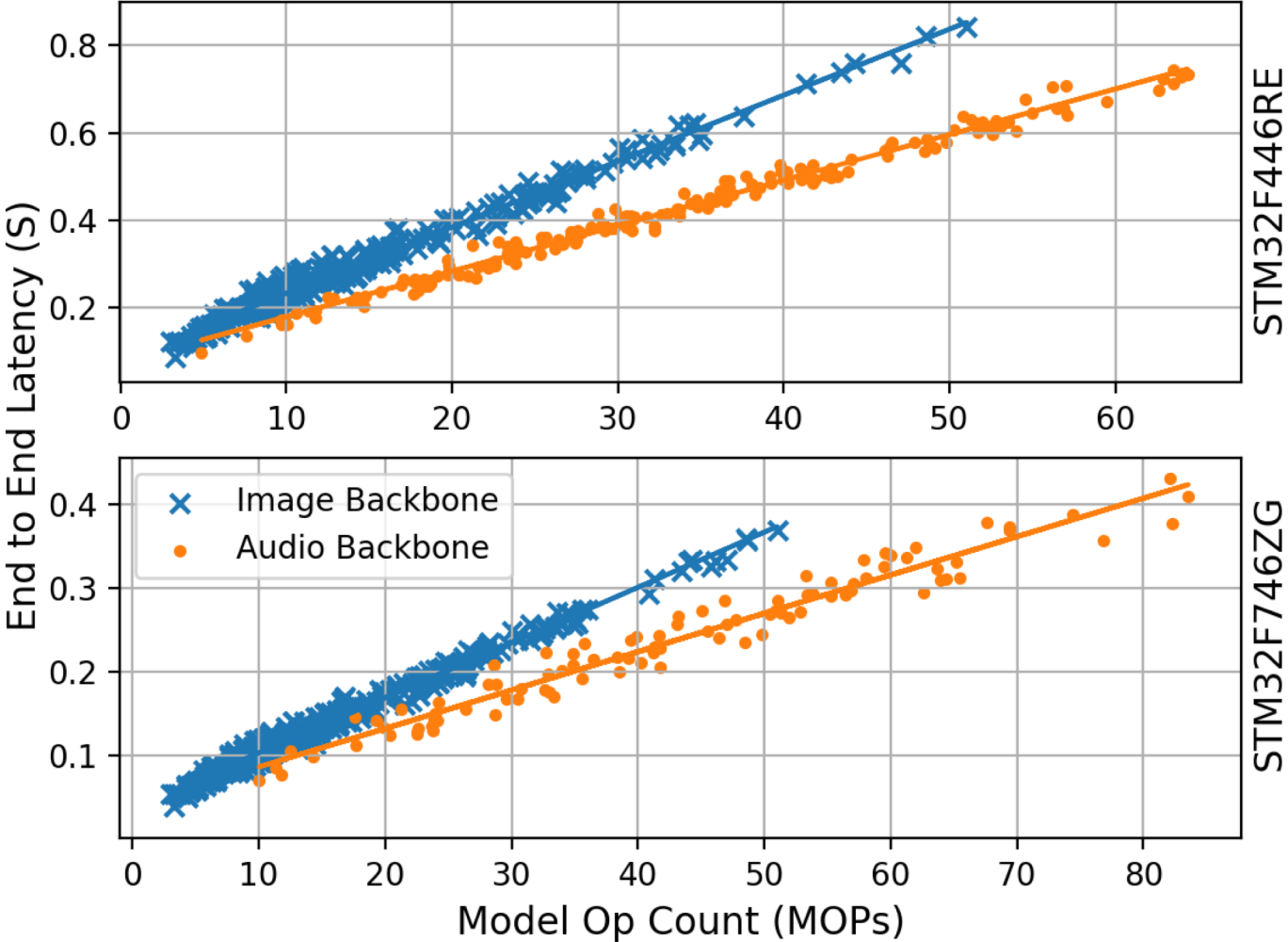
arm
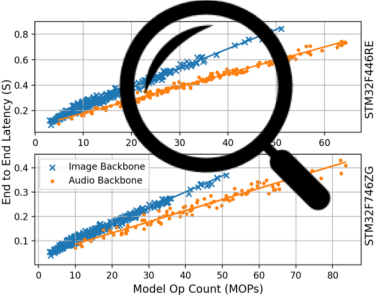
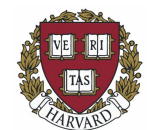# Direct Latency Benchmarking

# Direct Latency Benchmarking

# Latency Model

# Per Layer Latency
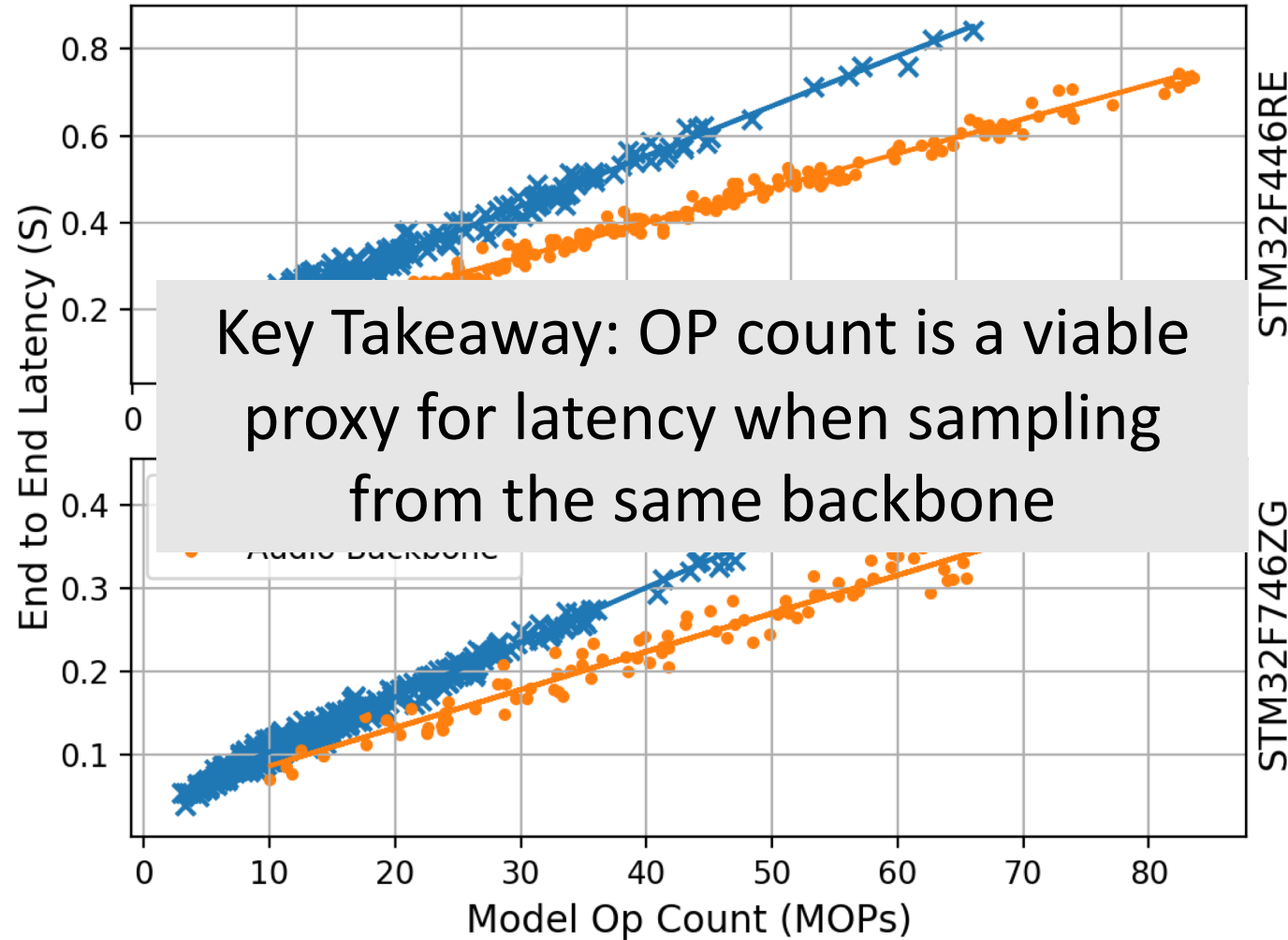
# Model Latency

# Model Latency

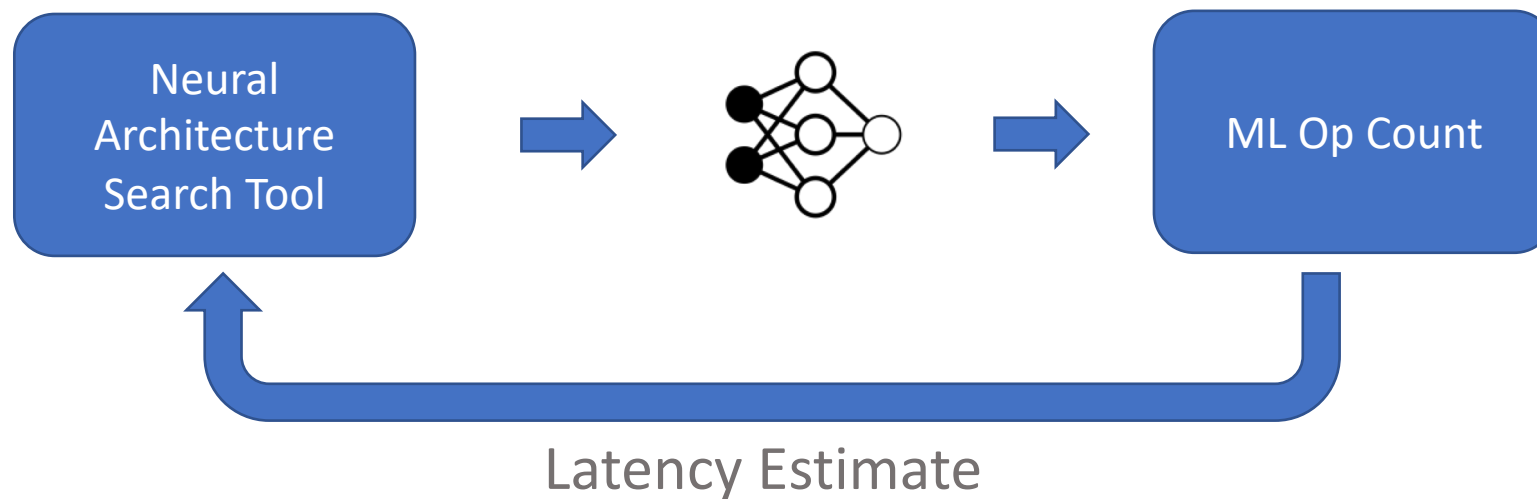# Model Latency



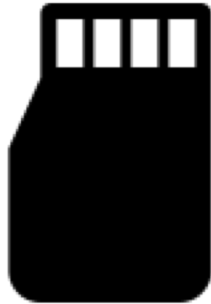Key Takeaway: OP count is a viable proxy for latency when sampling from the same backbone

# Latency Model



Neural Architecture Search Tool → → ML Op Count

Latency Estimate

# TinyML Constraints

| SRAM | Flash | Latency | Energy |
|:---:|:---:|:---:|:---:|
| ✔ | ✔ | ✔ | |

# Direct Energy Benchmarking

# Model Energy

# Model Energy

# Model Energy



Key Takeaway: Energy consumption is only a function of latency for a given MCU
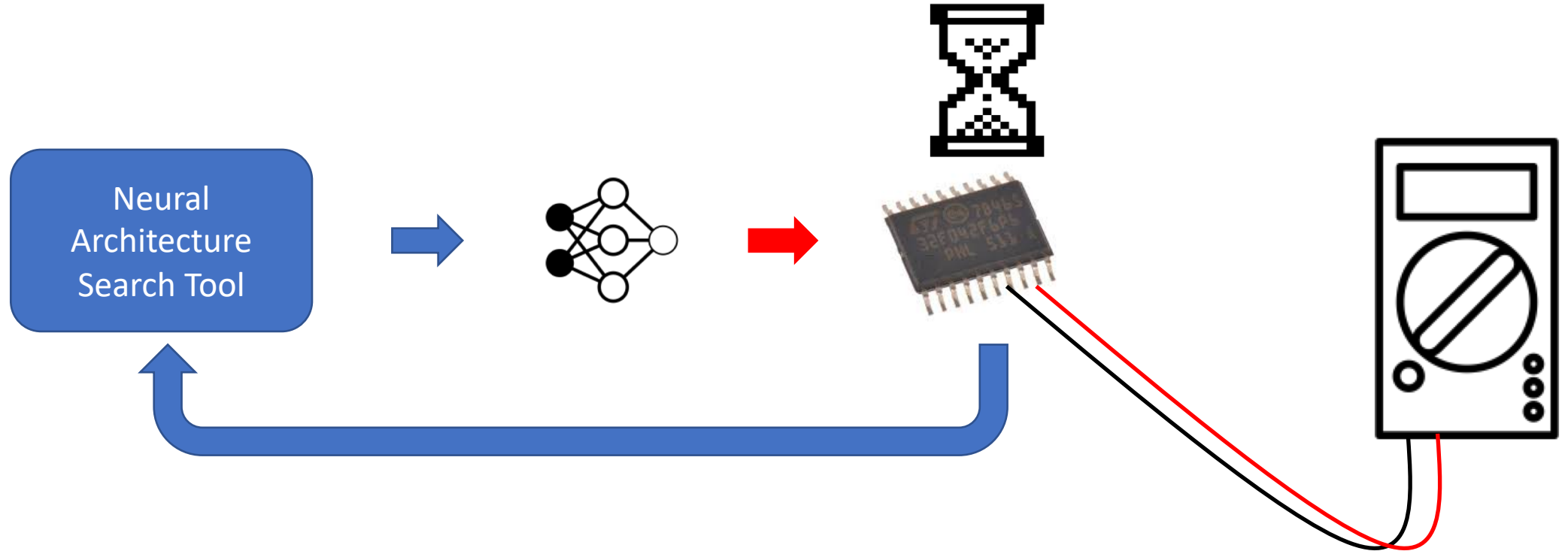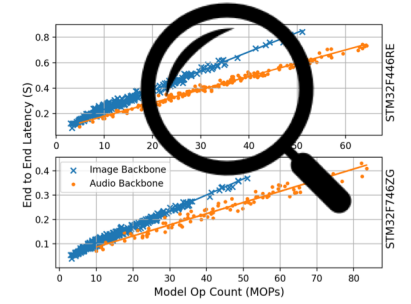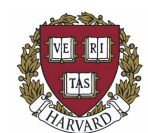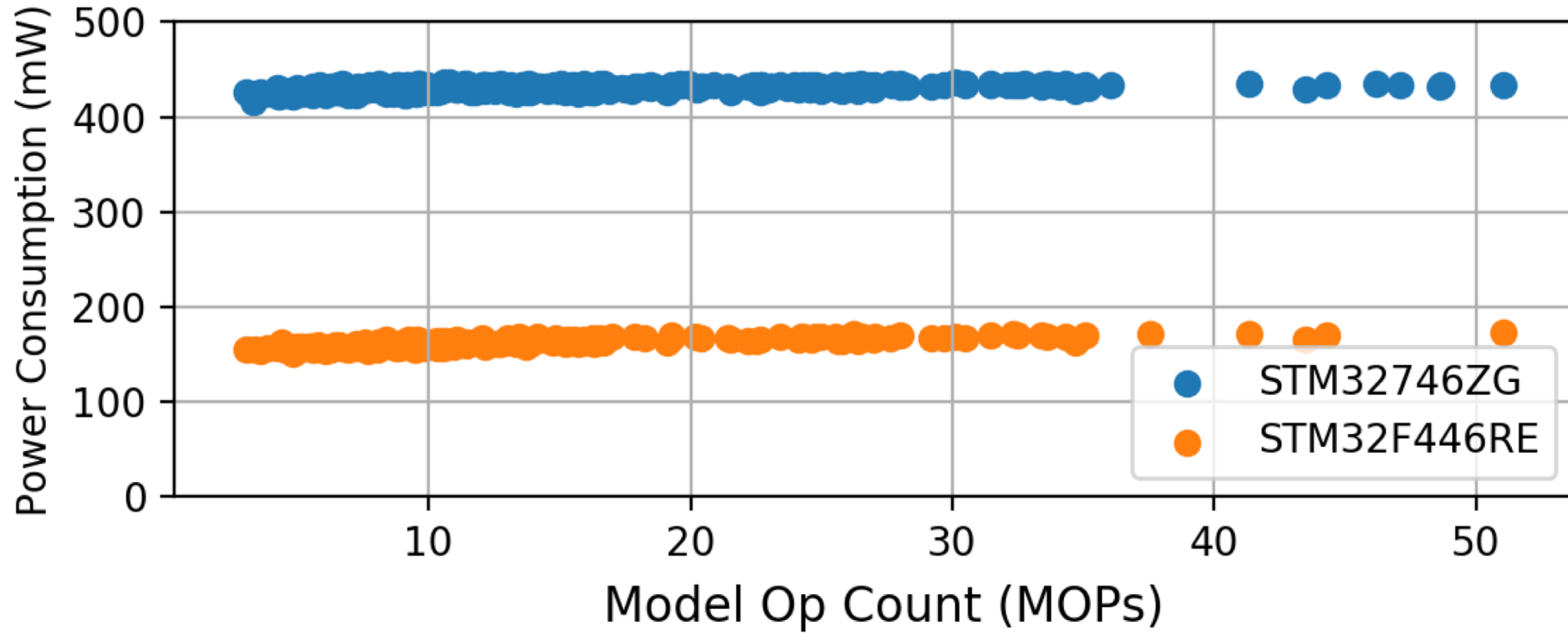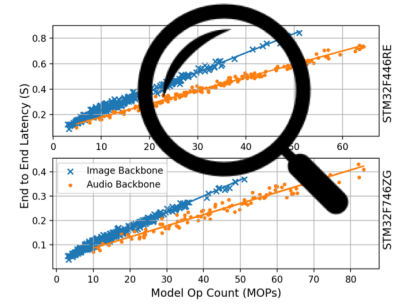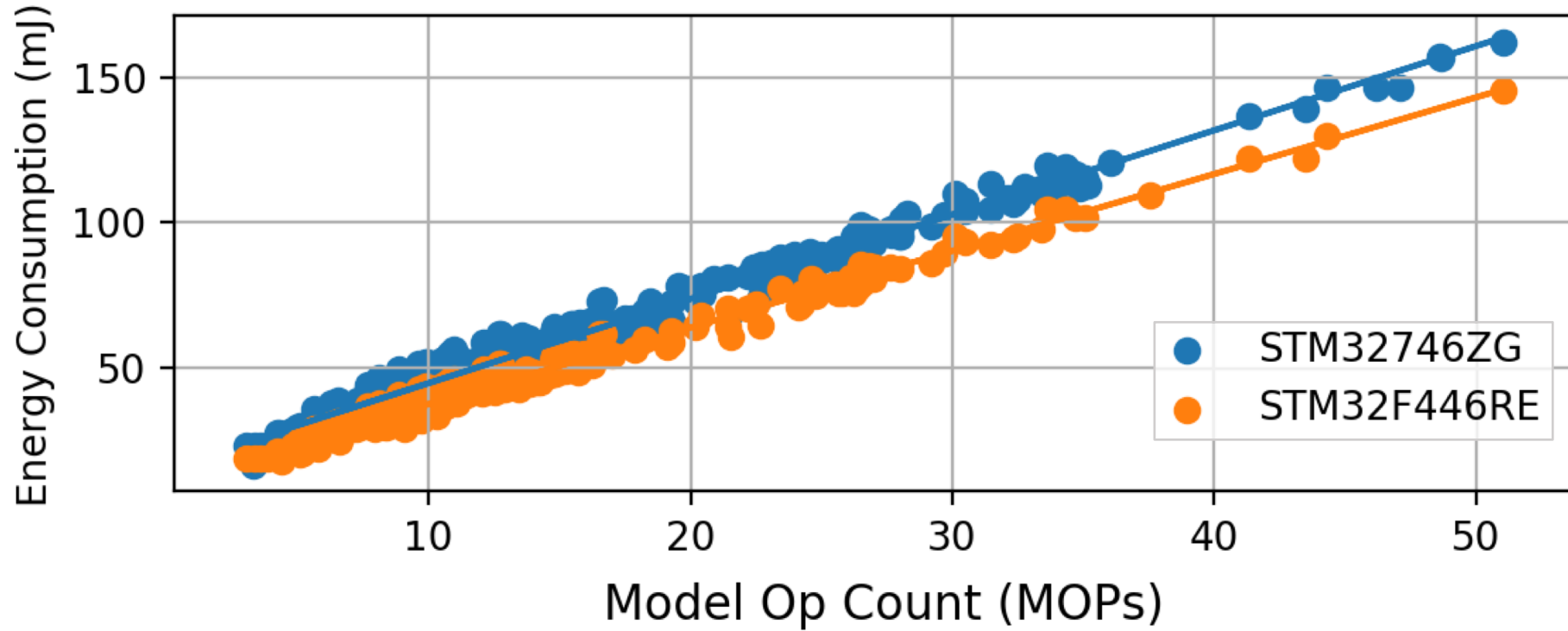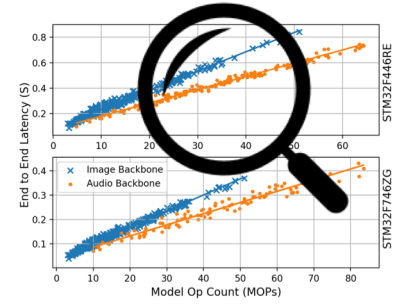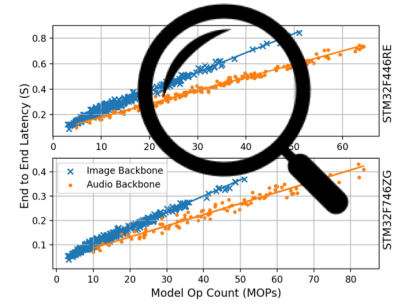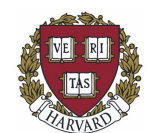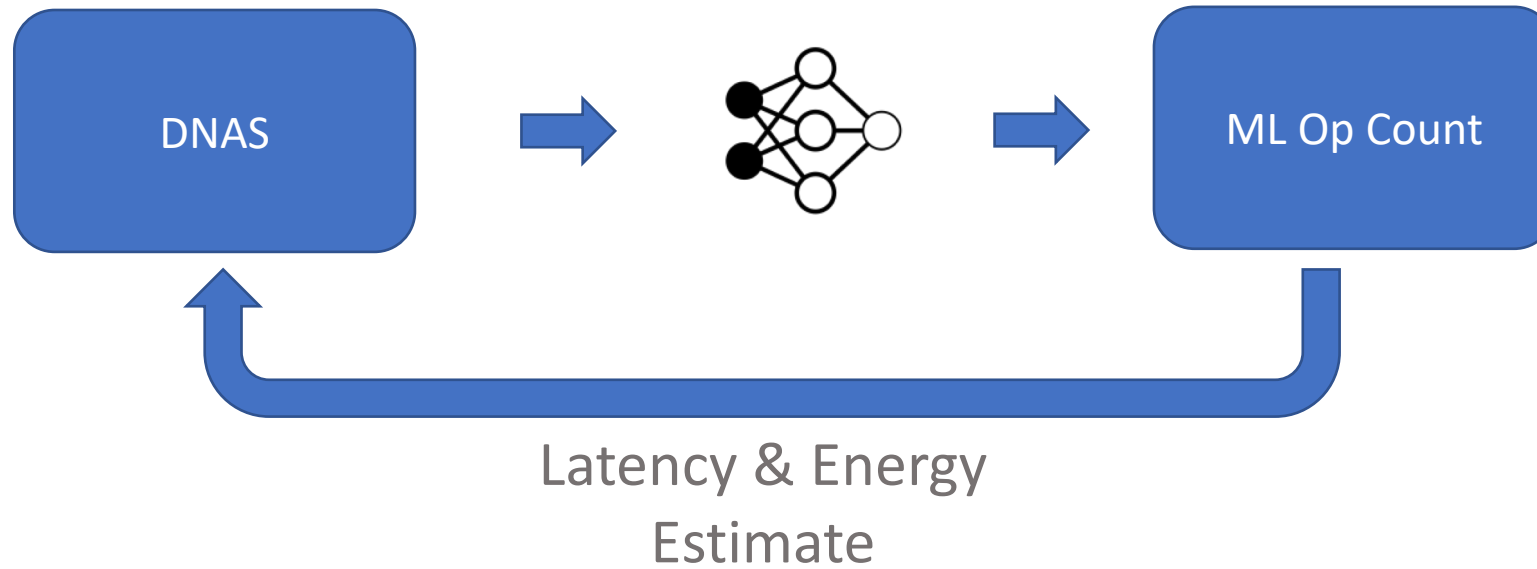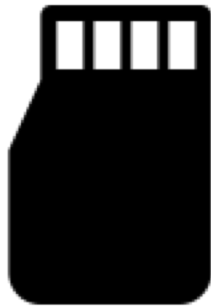
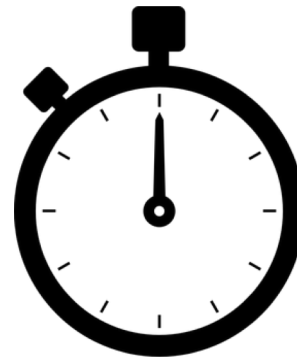# Latency & Energy Model
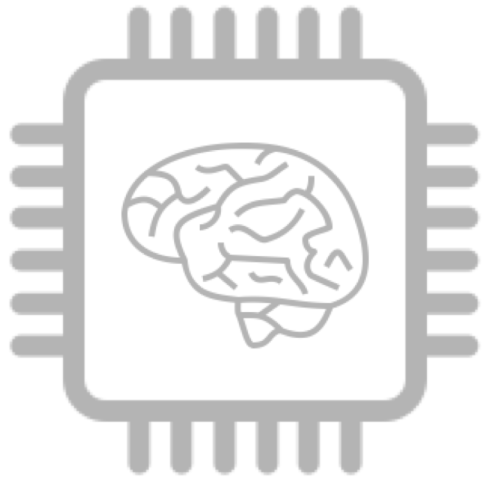
# TinyML Constraints



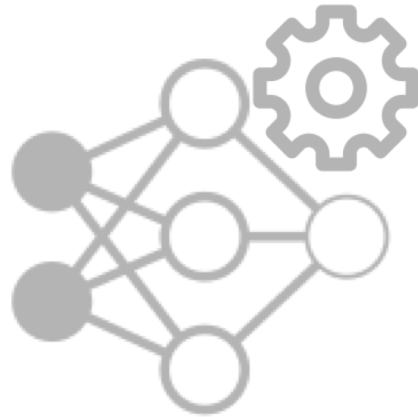SRAM       Flash       Latency       Energy
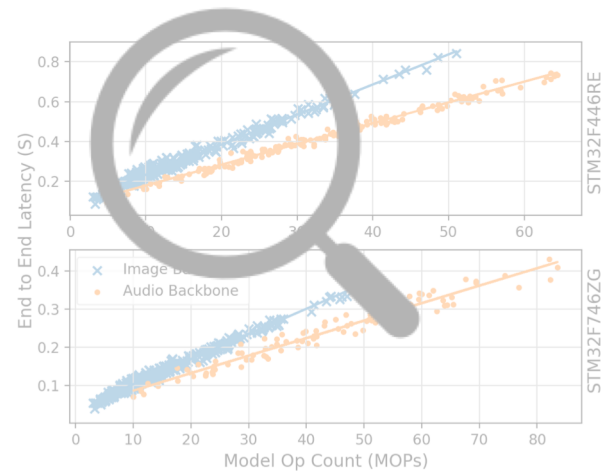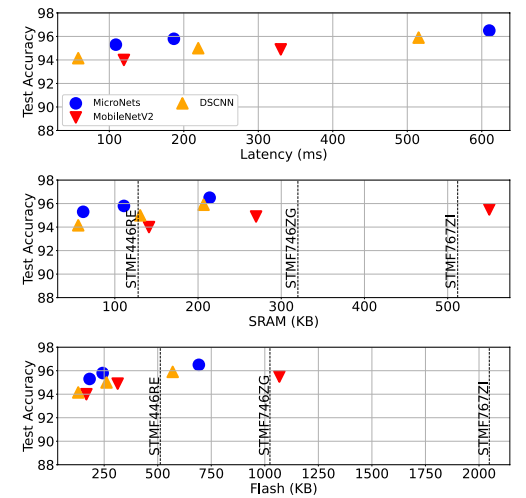
arm

# Executive Summary

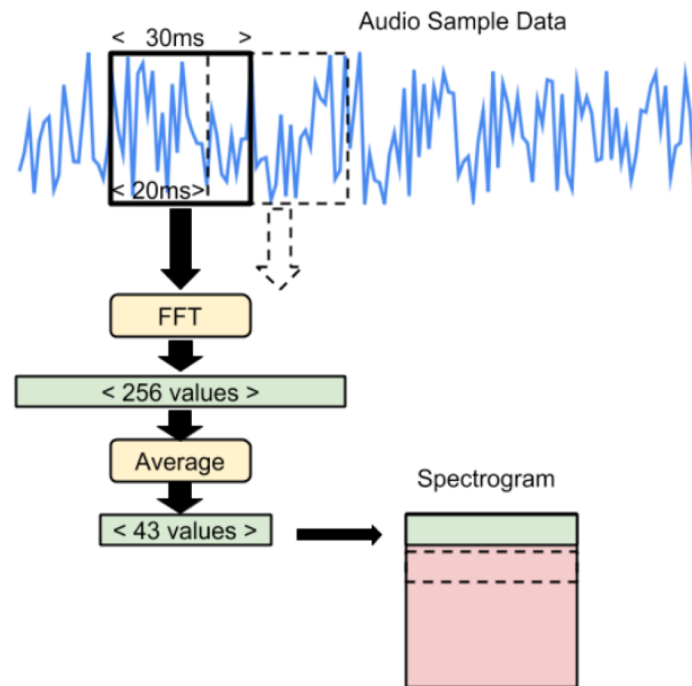TinyML

Differentiable
Neural Architecture
Search

Hardware
Characterization

MicroNets

# TinyMLPerf Use Cases

Keyword Spotting



Warden, Pete. "Speech commands: A dataset for limited-vocabulary speech recognition." *arXiv preprint arXiv:1804.03209* (2018).

# TinyMLPerf Use Cases

### Keyword Spotting

### Anomaly Detection



Warden, Pete. "Speech commands: A dataset for limited-vocabulary speech recognition." *arXiv preprint arXiv:1804.03209* (2018).

Purohit, Harsh, et al. "MIMII dataset: Sound dataset for malfunctioning industrial machine investigation and inspection." *arXiv preprint arXiv:1909.09347* (2019).
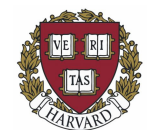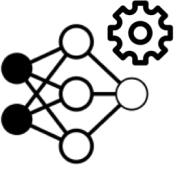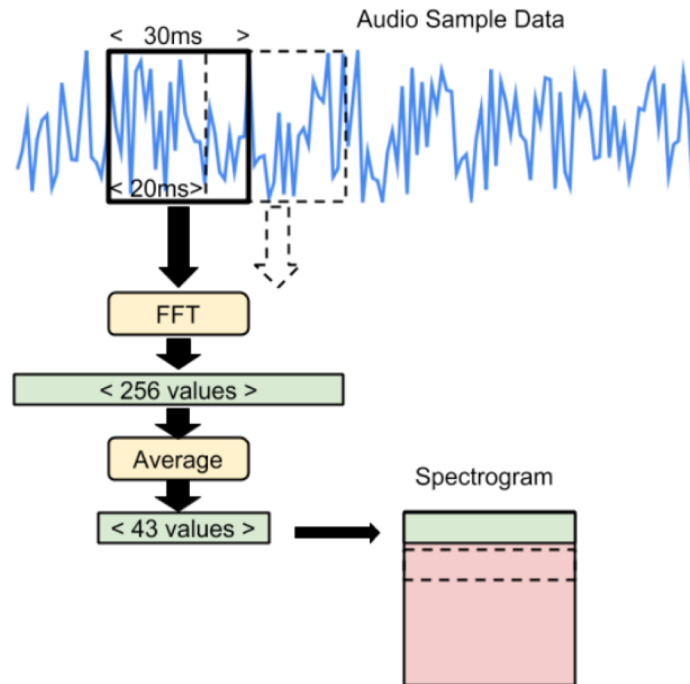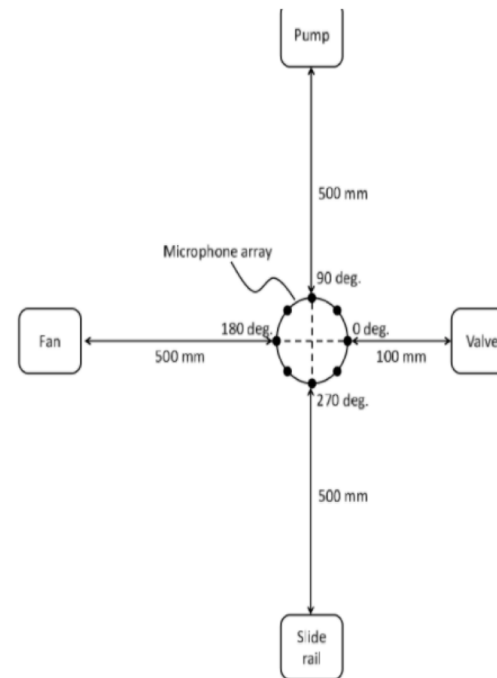
# TinyMLPerf Use Cases
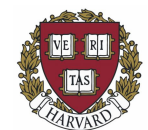
### Keyword Spotting



Warden, Pete. "Speech commands: A dataset for limited-vocabulary speech recognition." *arXiv preprint arXiv:1804.03209* (2018).
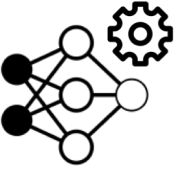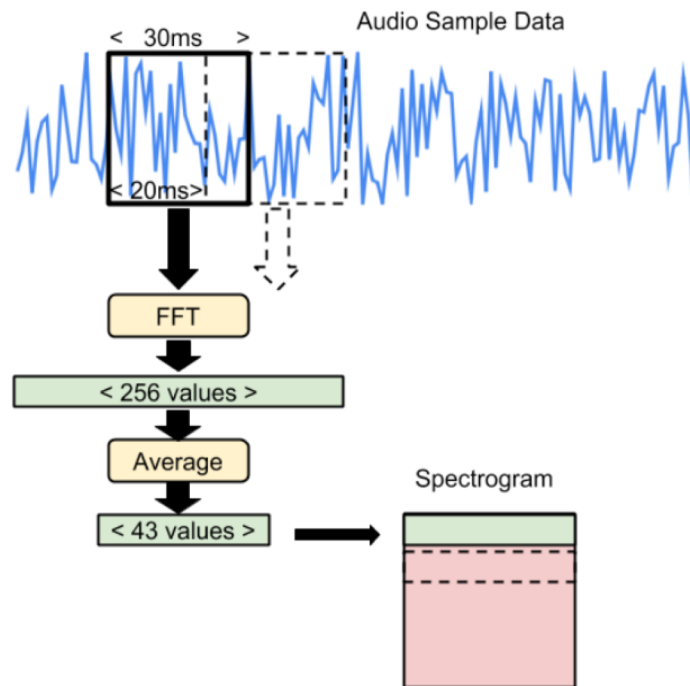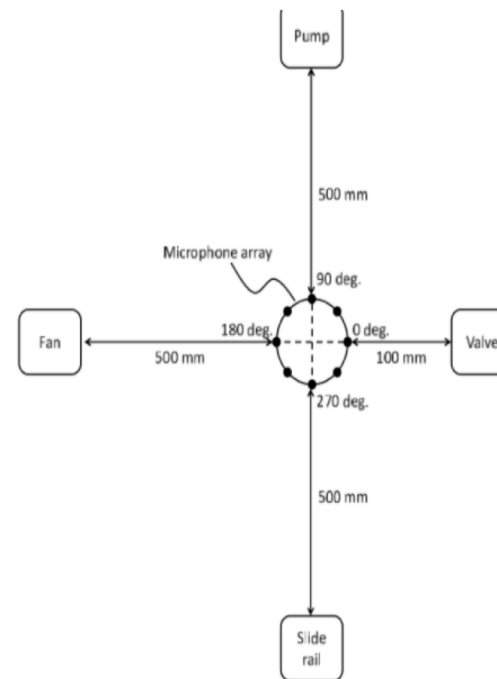
### Anomaly Detection



Purohit, Harsh, et al. "MIMII dataset: Sound dataset for malfunctioning industrial machine investigation and inspection." *arXiv preprint arXiv:1909.09347* (2019).
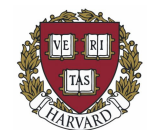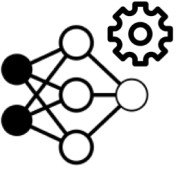
### Visual Wake Words



(a) 'Person'

(b) 'Not-person'

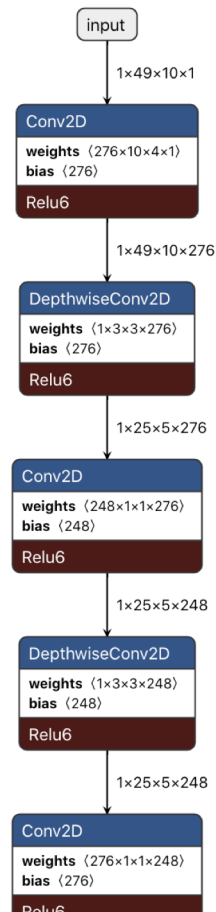Chowdhery, Aakanksha, et al. "Visual wake words dataset." *arXiv preprint arXiv:1906.05721* (2019).
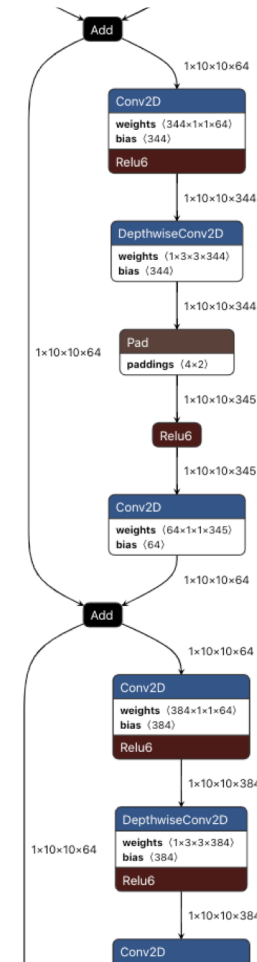
# Backbone Design
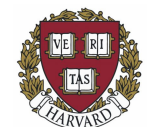


Keyword Spotting & Anomaly Detection
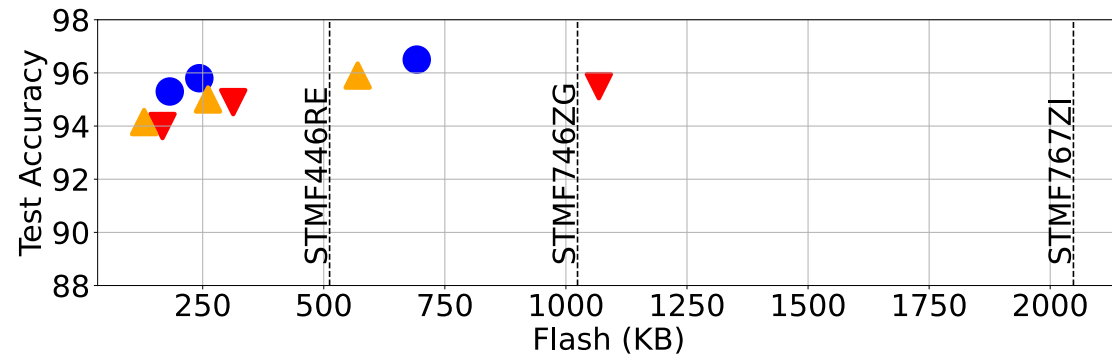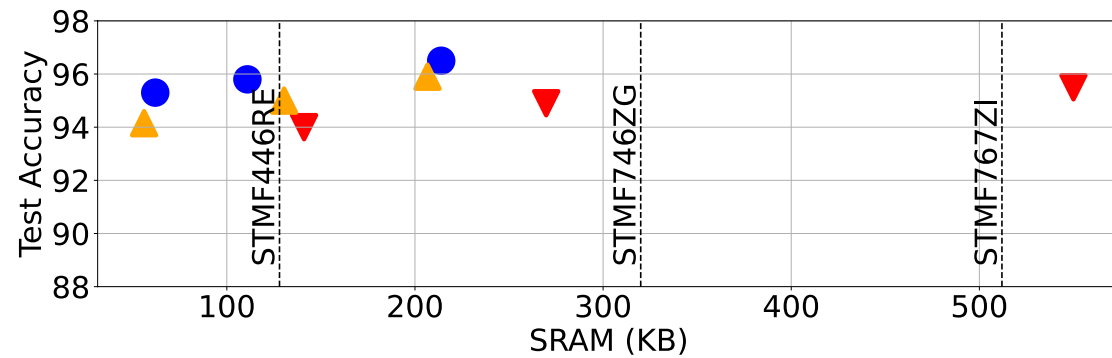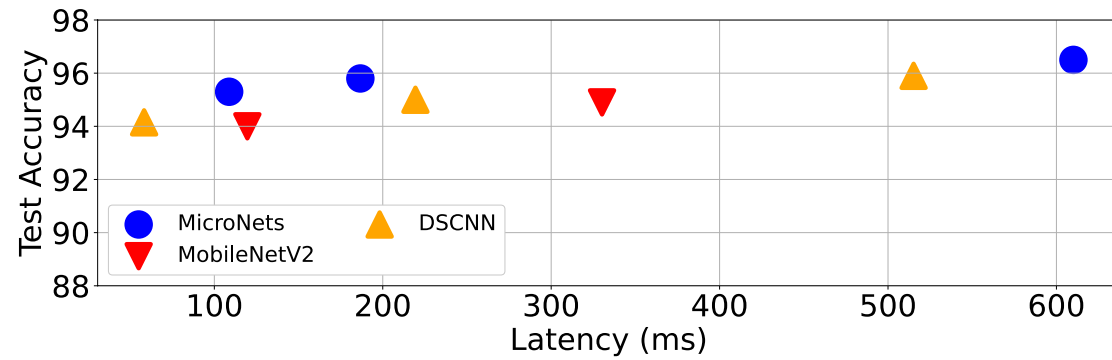
DSCNN-L[1]

Visual Wake Words

MobileNetV2[2]

[1] Zhang, Yundong, et al. "Hello edge: Keyword spotting on microcontrollers." *arXiv preprint arXiv:1711.07128* (2017).

[2] Sandler, Mark, et al. "Mobilenetv2: Inverted residuals and linear bottlenecks." *Proceedings of the IEEE conference on computer vision and pattern recognition*. 2018.
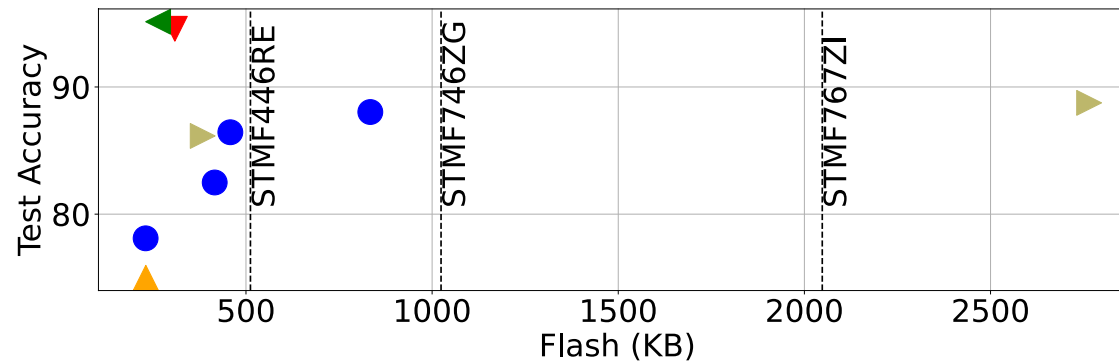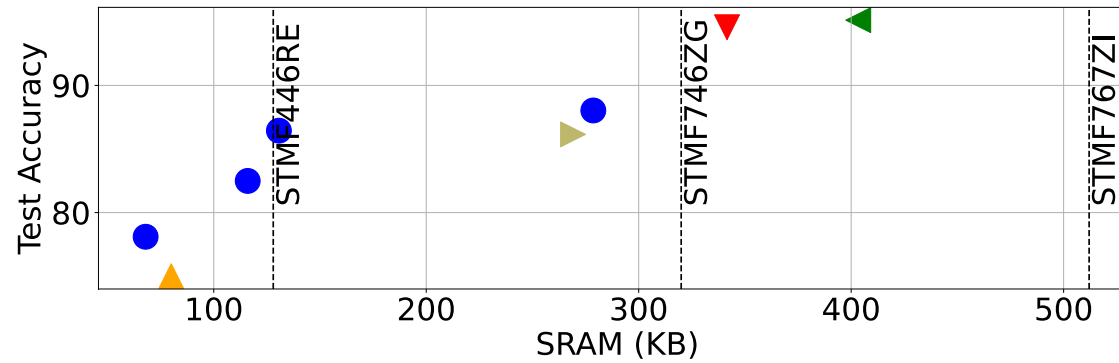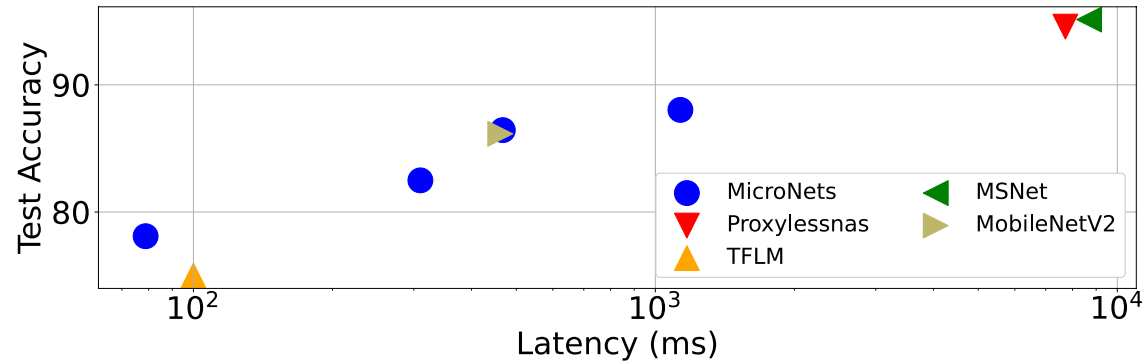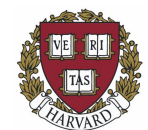
# Keyword Spotting
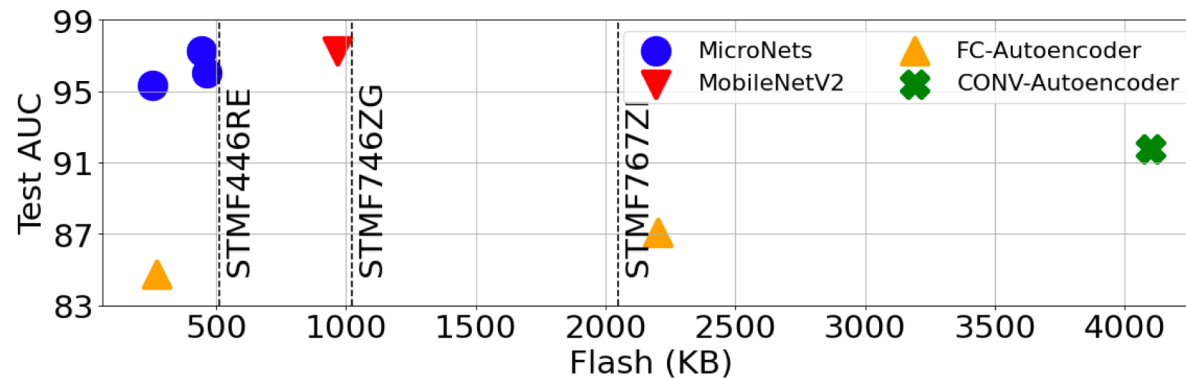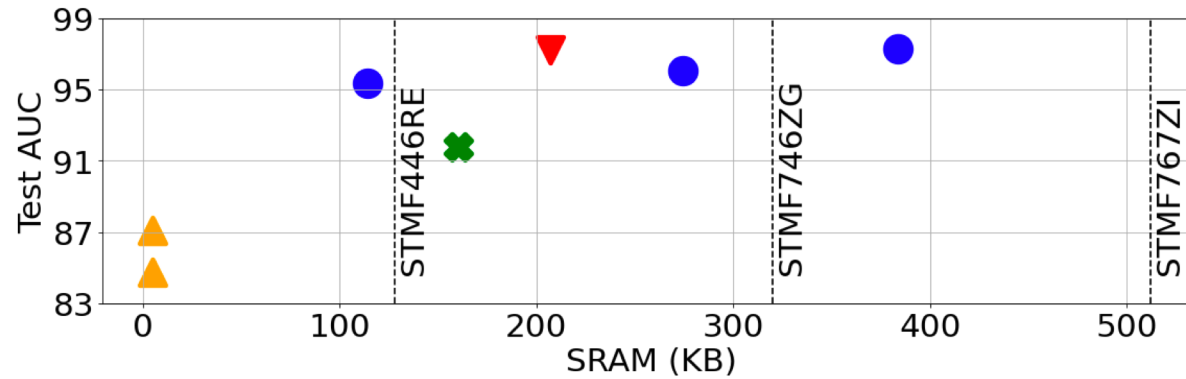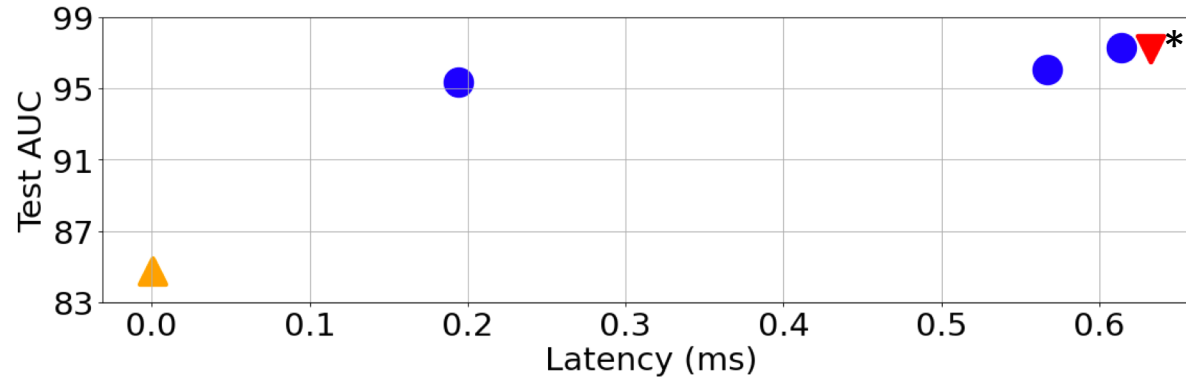
# Visual Wake Words

# Anomaly Detection
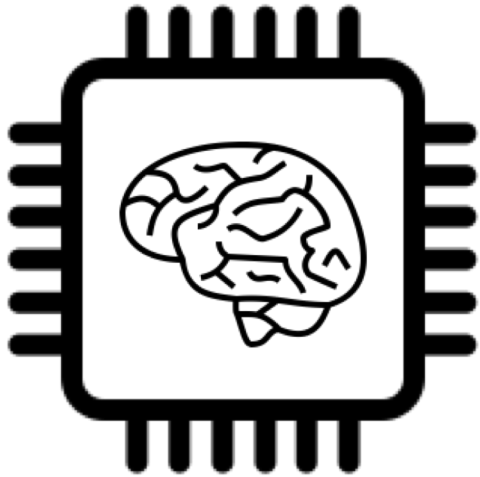
# Conclusion

TinyML

SRAM

Flash

TinyML systems have **severe constraints** and require **highly tuned** model architectures
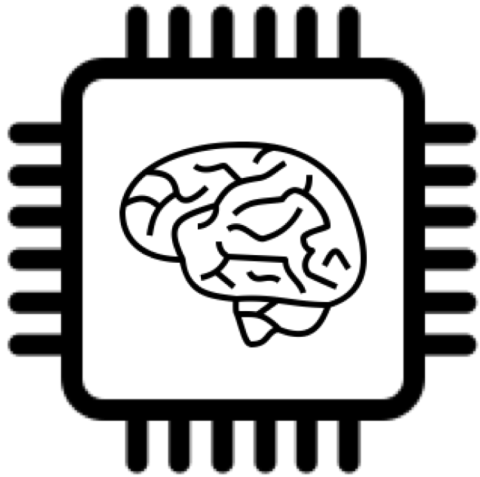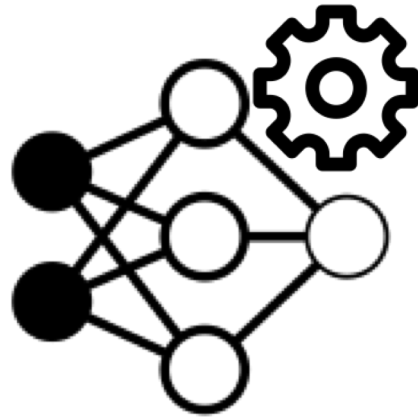
Latency
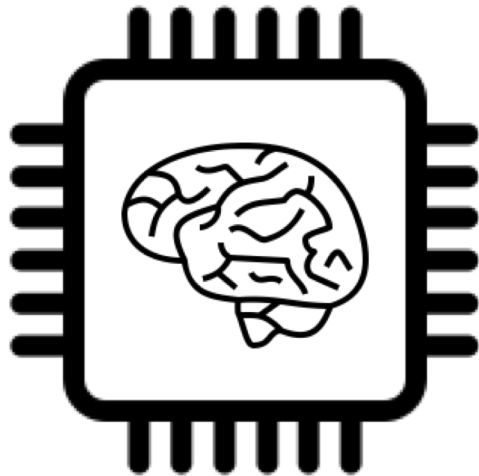
Energy

# Conclusion

TinyML
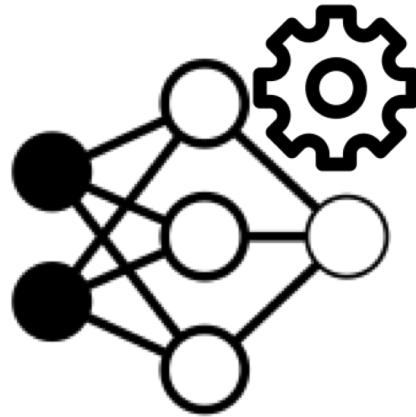
Differentiable
Neural Architecture
Search

Differentiable Neural
Architecture Search (DNAS)
can **rapidly** find models that
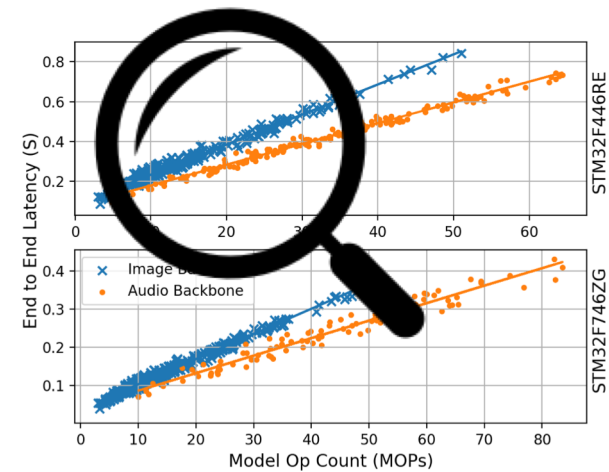**meet the constraints** given
**viable proxies**

# Conclusion

TinyML
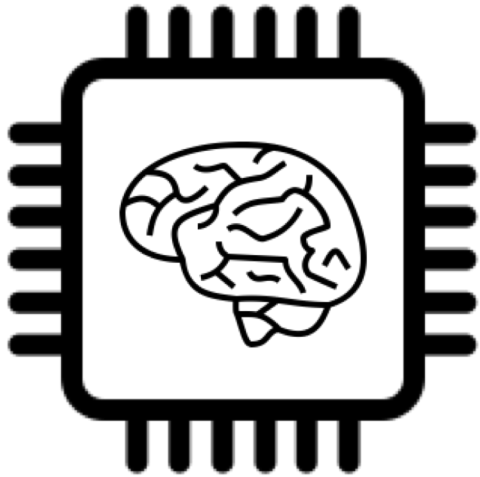
Differentiable Neural Architecture Search
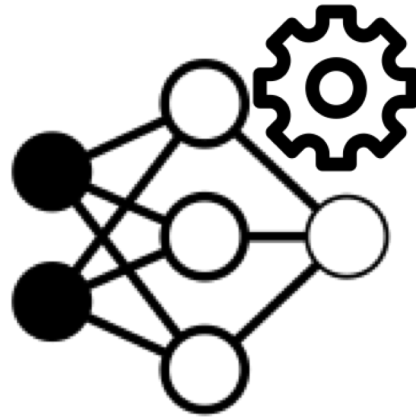
Hardware Characterization



SRAM and Flash are easily calculated while **Op count is a viable proxy** latency and energy
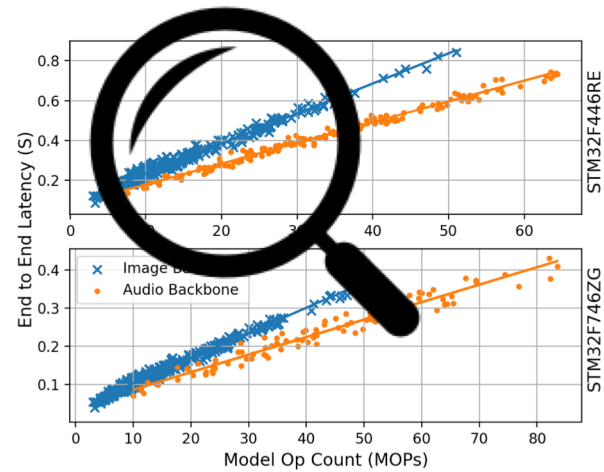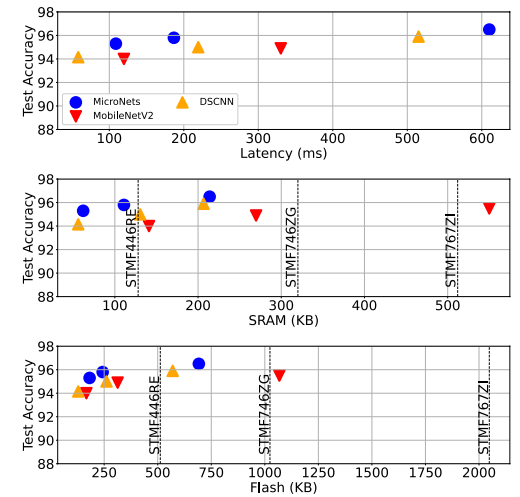
# Conclusion

TinyML

Differentiable
Neural Architecture
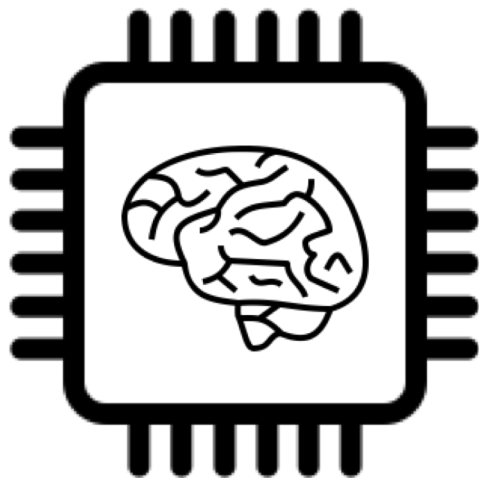Search

Hardware
Characterization

MicroNets



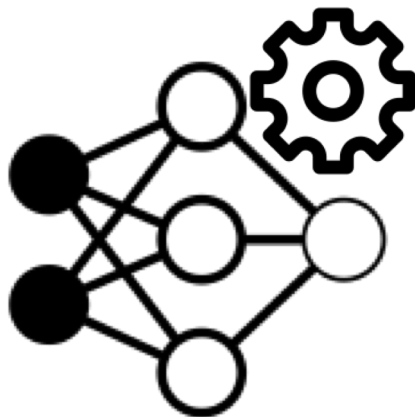We achieve **state of the art performance** on
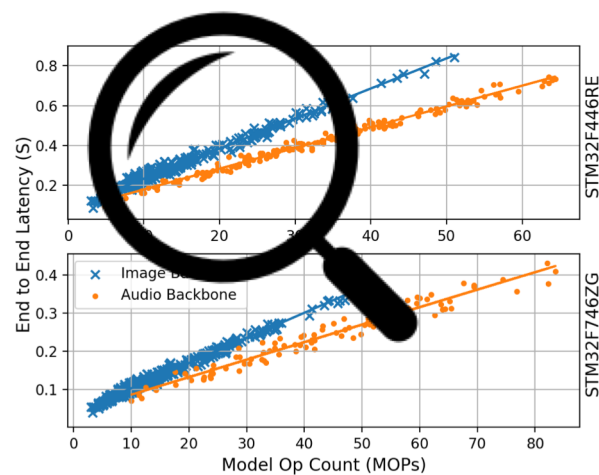three TinyML tasks

# Conclusion

TinyML
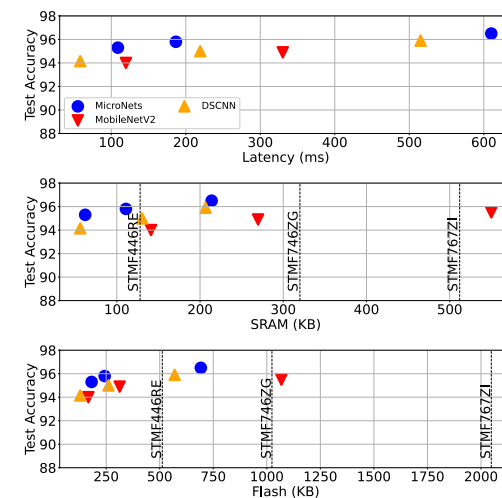
Differentiable
Neural Architecture
Search

Hardware
Characterization

MicroNets



Models and Training Scripts are available:
github.com/ARM-software/ML-zoo