

# tinyML<sup>®</sup> Talks

*Enabling Ultra-low Power Machine Learning at the Edge*

“Minimizing resource usage in microcontrollers for cost effective solutions”

Ilya Gozman – Senior Fellow, Chief AI Architect, Grovety

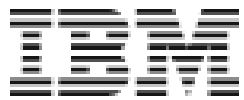
September 26, 2023



[www.tinyML.org](http://www.tinyML.org)



Thank you, **tinyML Strategic Partners**,  
for committing to take tinyML to the next Level, together



T I N Y



TALKS  
*webcast*

# Executive Strategic Partners

**Qualcomm**  
AI research

# Advancing AI research to make efficient AI ubiquitous

## Power efficiency

Model design, compression, quantization, algorithms, efficient hardware, software tool

## Personalization

Continuous learning, contextual, always-on, privacy-preserved, distributed learning

## Efficient learning

Robust learning through minimal data, unsupervised learning, on-device learning

## A platform to scale AI across the industry



### Perception

Object detection, speech recognition, contextual fusion



### Reasoning

Scene understanding, language understanding, behavior prediction



### Action

Reinforcement learning for decision making



Edge cloud



Cloud



IoT/IIoT



Automotive



Mobile



Accelerate Your Edge Compute

**SYNTIANT**

Making Edge AI A Reality

[www.syntiant.com](http://www.syntiant.com)

# Platinum Strategic Partners



**DEPLOY VISION AI  
AT THE EDGE AT SCALE**

**SONY**

# Gold Strategic Partners





AHEAD OF WHAT'S POSSIBLE™



AHEAD OF WHAT'S POSSIBLE™

Where what if  
becomes what is.

Witness potential made possible at [analog.com](http://analog.com).

Build the  
Future of tinyML

on **arm**



T I N Y



TALKS  
*webcast*



**EDGE IMPULSE**

# The Leading Development Platform for Edge ML

[edgeimpulse.com](https://edgeimpulse.com)

Decarbonization

Digitalization



Driving decarbonization and digitalization. Together.

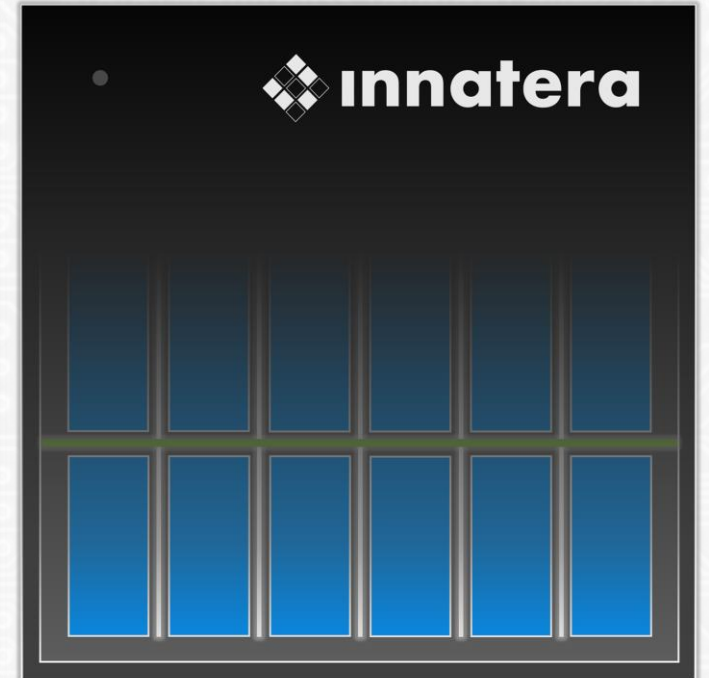
**Infineon serving all target markets as**  
**Leader in Power Systems and IoT**

[www.infineon.com](http://www.infineon.com)





# NEUROMORPHIC INTELLIGENCE FOR THE SENSOR-EDGE

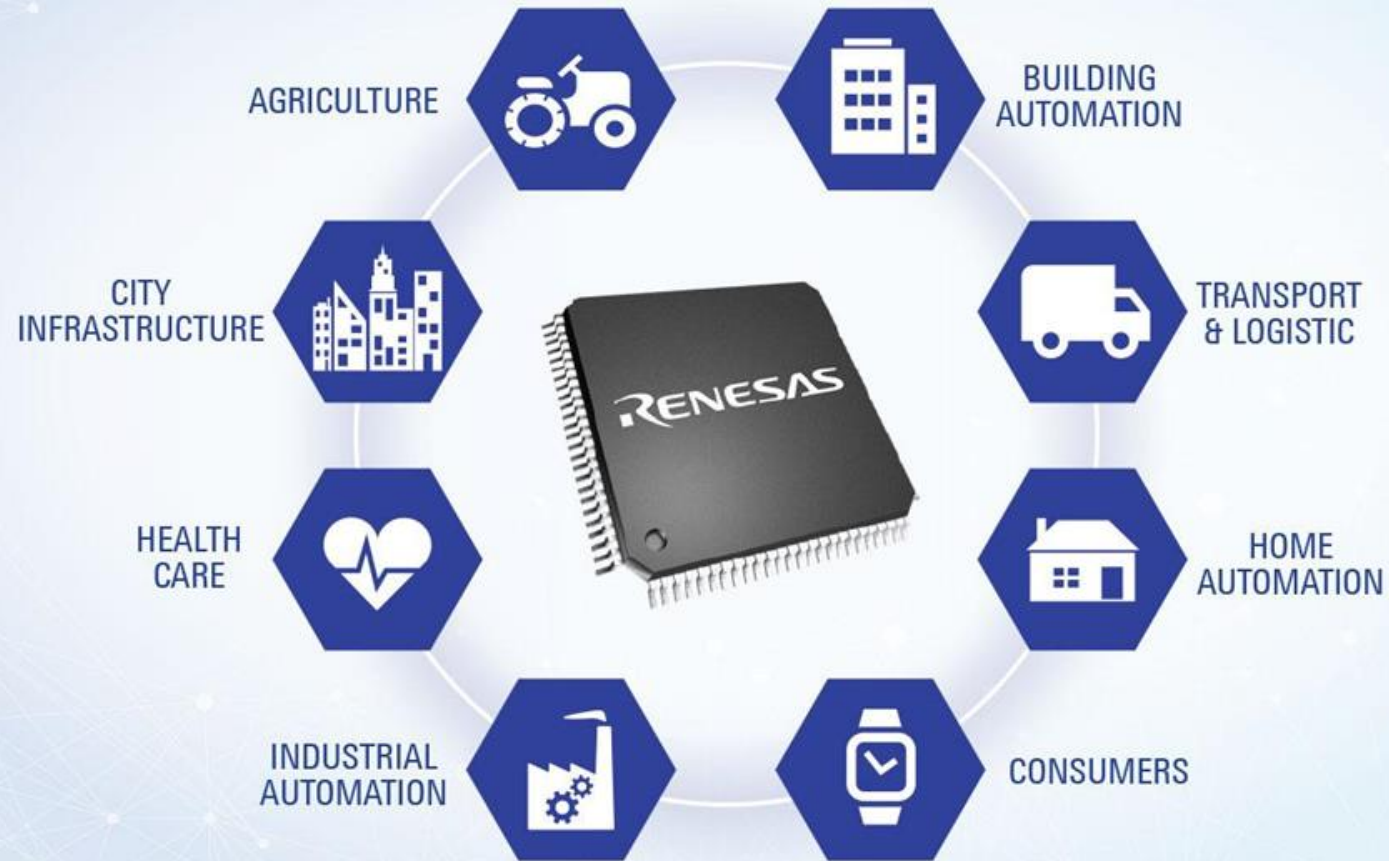


[www.innatera.com](http://www.innatera.com)



Microsoft

**Renesas is enabling the next generation of AI-powered solutions that will revolutionize every industry sector.**



[renesas.com](https://www.renesas.com)



life.augmented

**STMicroelectronics provides extensive solutions to make tiny Machine Learning easy**





# ENGINEERING EXCEPTIONAL EXPERIENCES

We engineer exceptional experiences for consumers in the home, at work, in the car, or on the go.

[www.synaptics.com](http://www.synaptics.com)



T I N Y



# Silver Strategic Partners



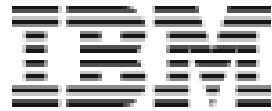
brainchip



GREENWAVES  
TECHNOLOGIES



⚡ Grovety Inc.



NotaAI





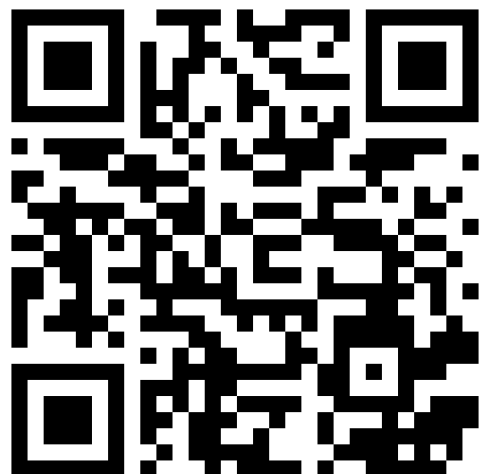
# Join Growing tinyML Communities:



16.9k members in  
49 Groups in 41 Countries

**tinyML - Enabling ultra-low Power ML at the Edge**

<https://www.meetup.com/tinyML-Enabling-ultra-low-Power-ML-at-the-Edge/>



4k members  
&  
13k followers

**The tinyML Community**

<https://www.linkedin.com/groups/13694488/>





Subscribe to  
**tinyML YouTube Channel**  
for updates and notifications  
*(including this video)*

[www.youtube.com/tinyML](http://www.youtube.com/tinyML)



**tinyML**  
4.33K subscribers

**10.5k subscribers, 628 videos with 380k views**

HOME VIDEOS PLAYLISTS COMMUNITY CHANNELS ABOUT

 13:24 On Device Learning Forum - Professors... 106 views · 4 days ago	 33:27 On Device Learning - Manuel Roveri: Is on-... 138 views · 4 days ago	 32:39 On Device Learning Forum - Warren Gros... 54 views · 4 days ago	 36:41 On Device Learning Forum - Yiran Chen... 47 views · 4 days ago	 34:03 On Device Learning Forum - Hiroku... 132 views · 4 days ago	 34:58 On Device Learning Forum - Song Han: O... 137 views · 4 days ago
 1:13 tinyML Smart Weather Station Challenge - ... 122 views · 4 days ago	 1:07:43 tinyML Talks Singapore... 262 views · 2 weeks ago	 53:41 tinyML Talks Shenzhen: Data... 511 views · 3 weeks ago	 45:46 tinyML Talks Singapore... 229 views · 3 weeks ago	 51:01 tinyML Smart Weather Station with Syntiant... 265 views · 3 weeks ago	 1:03:24 tinyML Trailblazers August with Vijay... 286 views · 1 month ago
 58:50 tinyML Auto ML Tutorial with SensiML 351 views · 1 month ago	 34:36 tinyML Auto ML Tutorial with Qeexo 462 views · 2 months ago	 55:01 tinyML Talks Germany: Neural network... 374 views · 2 months ago	 59:51 tinyML Trailblazers with Yoram Zylberberg 133 views · 2 months ago	 59:48 tinyML Auto ML Tutorial with Nota AI 287 views · 2 months ago	 58:09 tinyML Auto ML Tutorial with Neuton 336 views · 2 months ago
 1:02:30 tinyML Challenge 2022: Smart weather... 378 views · 2 months ago	 34:31 tinyML Talks South Africa - What is... 214 views · 2 months ago	 1:00:30 tinyML Talks: The new Neuromorphic Anal... 448 views · 2 months ago	 1:06:44 tinyML Talks Shenzhen: 分享主题... 159 views · 2 months ago	 1:53:07 tinyML Auto ML Forum - Paneldiscussion 190 views · 2 months ago	 42:13 tinyML Auto ML Forum - Demos 545 views · 2 months ago



# tinyML Asia Technical Forum

**November 16, 2023  
Seoul, South Korea**



**Register now**  
**<https://www.tinyml.org/event/asia-2023/>**

# 2023 Edge AI Technology Report

The guide to understanding the state of the art in hardware & software in Edge AI.



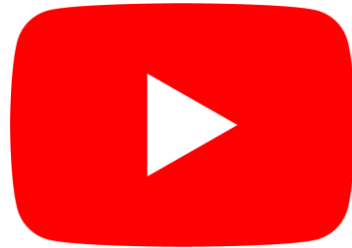


# Reminders

Slides & Videos will be posted tomorrow



[tinyml.org/forums](https://tinyml.org/forums)



[youtube.com/tinyml](https://youtube.com/tinyml)



Please use the Q&A window for your questions





## Ilya Gozman



Ilya is a Senior Fellow and a Chief AI Architect at Grovety, where he worked out his way from a rising talent developer to a veteran expert in AI, a frontline and prospective trend in IT-industry in recent years. He acquired extensive experience in developing general and AI compilers, and chip architectures both in LLVM and TVM backend optimizations; he also led teams working on compiling-related projects, video processing, and protocols support for IP cameras (C/C++). Ilya received Master degree in Applied Mathematics and Computer Science in 2007. Wide range of projects and profound research activity makes Ilya's experience valuable and demanded.



# TinyML - growing interest

Edge AI allows business to improve the AI applications' overall cost-effectiveness by optimal use of NNs, computing resources and power consumption reduction.

At the same time, the numerous potential benefits of Edge AI face several challenges associated with its implementation and usability. [1]

# Fine-Tuning Strategies

## **Model modification:**

- Compress off-the-shelf networks by pruning and quantization
- Simplify unsupported operations to primitive blocks
- Transform and merge network layers
- Optimize resource-intensive layers

## **Inference time optimization:**

- Use hardware-specific acceleration instructions

## **Memory requirements optimization:**

- Optimize schedule of the operation flow
- Store weights on external storage

## **Energy Efficiency**

- Throttling MCU/NPU operating frequencies
- Use advantages of heterogenous systems
- Intelligent power management

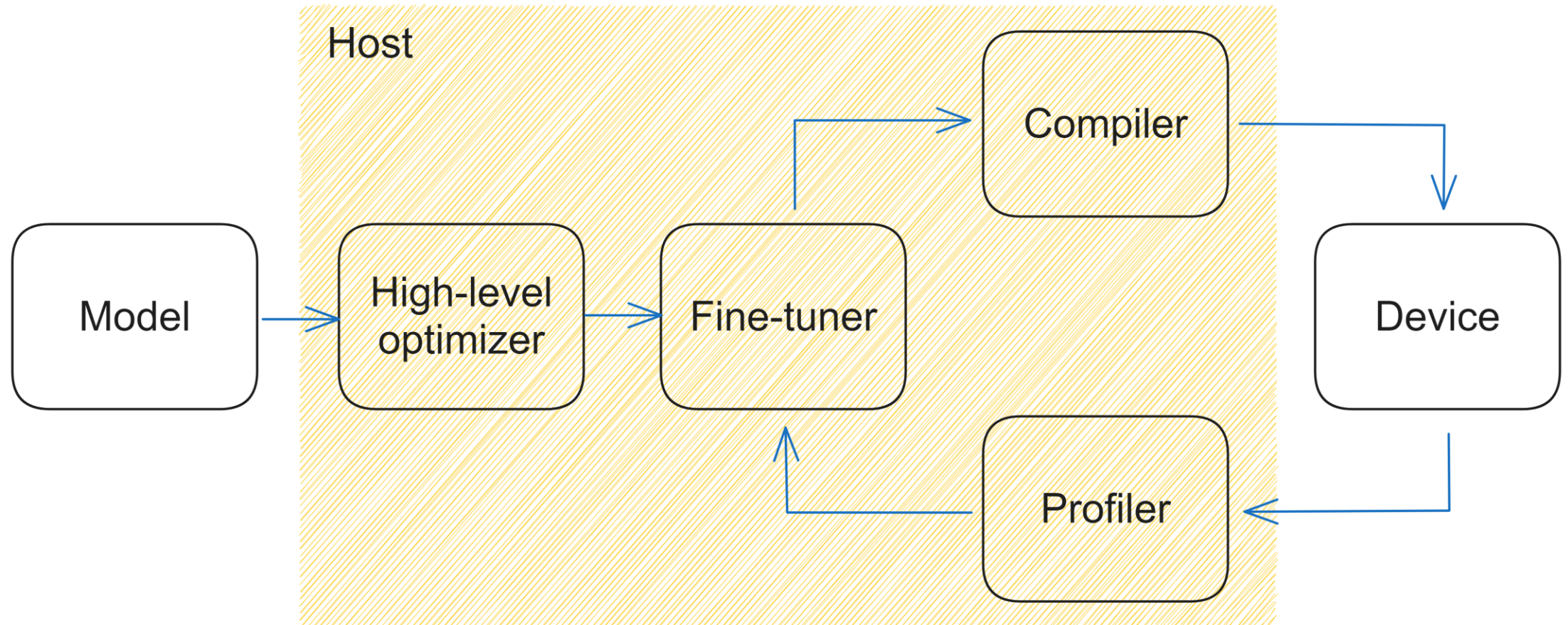
# Approaches to handle Cost Challenges in TinyML

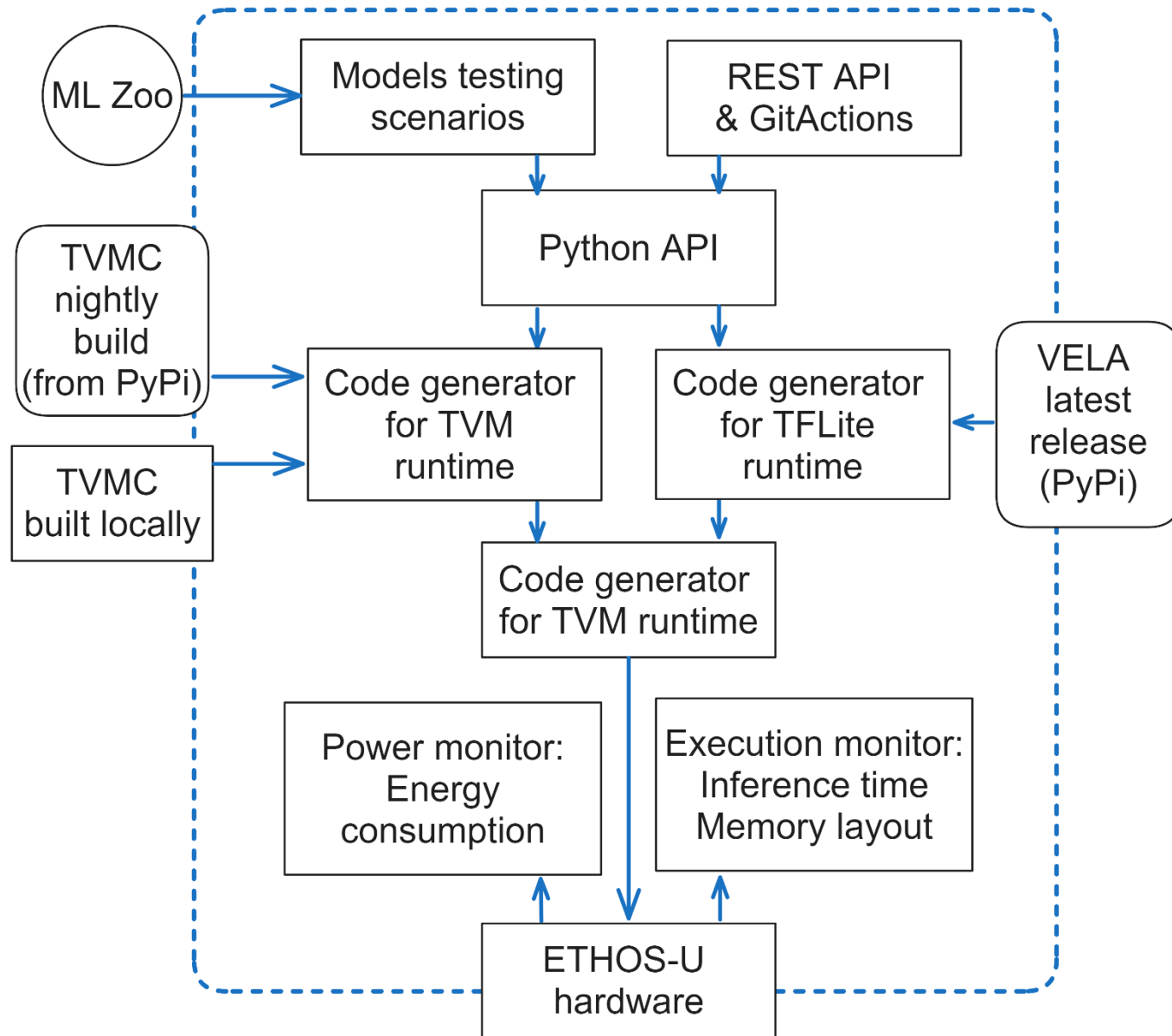


Minimizing Development costs and time of device

Reducing Device cost and its power consumption

# Fine-Tuning Strategy





## Fine tuning platform for CI tests and experiments on target HW

- NN inference on Alif hardware and FVP simulator
- Run on TFLiteMicro and TVM runtimes
- Support of any TVM commit
- Unified API for running NN inferences, various architectures and runtimes
- Actual inference time and power consumption measurements
- Model bottlenecks analysis and numerical mismatches

# TVM: ML Compiler Framework



Wide range  
of ML frameworks  
and deployment  
targets



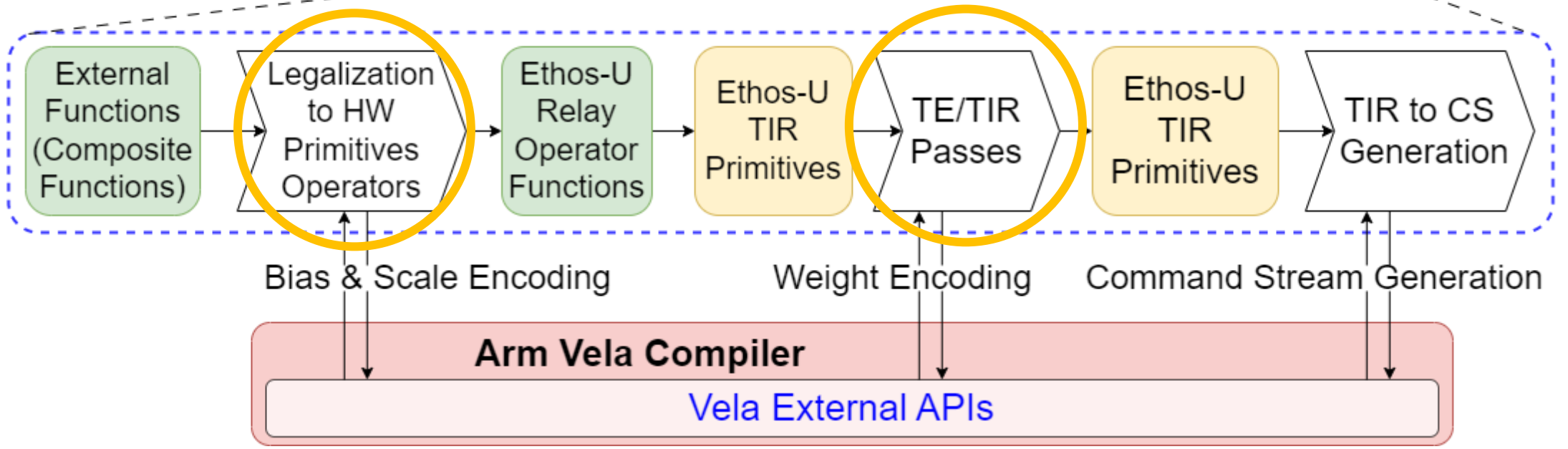
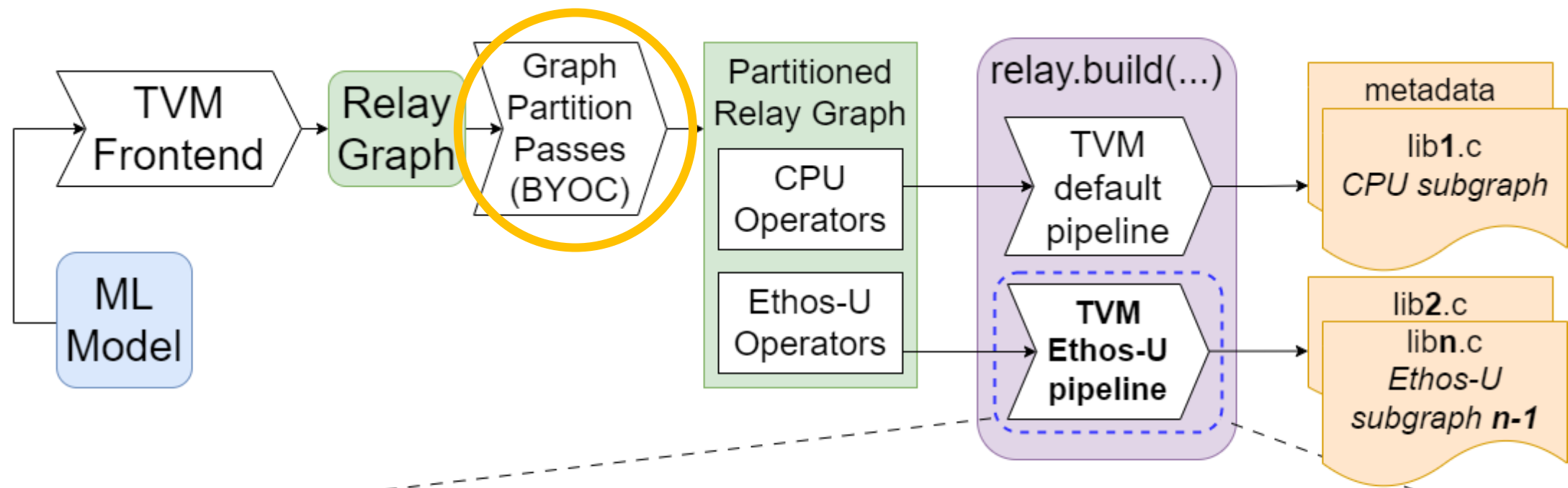
Open-source  
project, large  
community



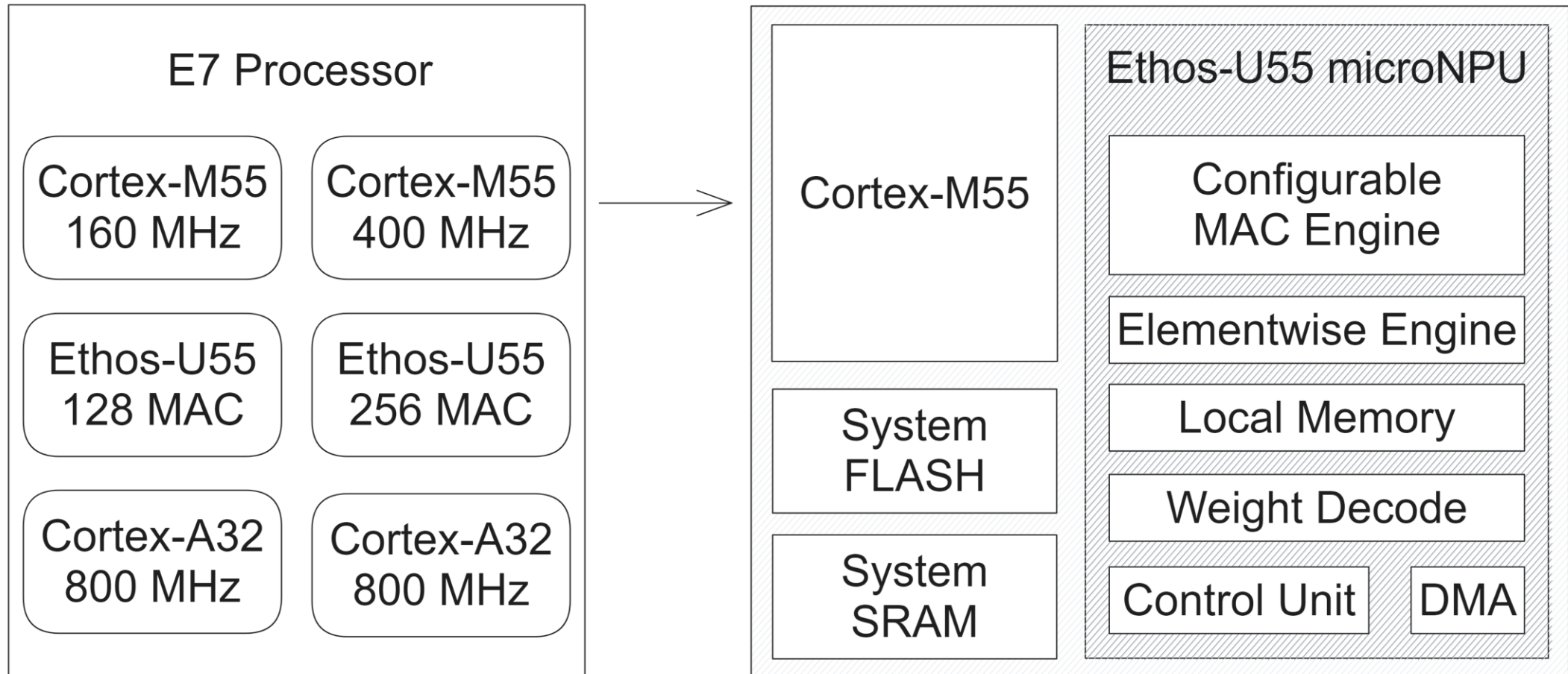
Integrated with  
ARM<sup>®</sup> Vela  
Compiler for  
acceleration on  
Ethos<sup>™</sup>-U55 NPU



Fine-grained  
control over model  
compilation,  
deployment and  
execution

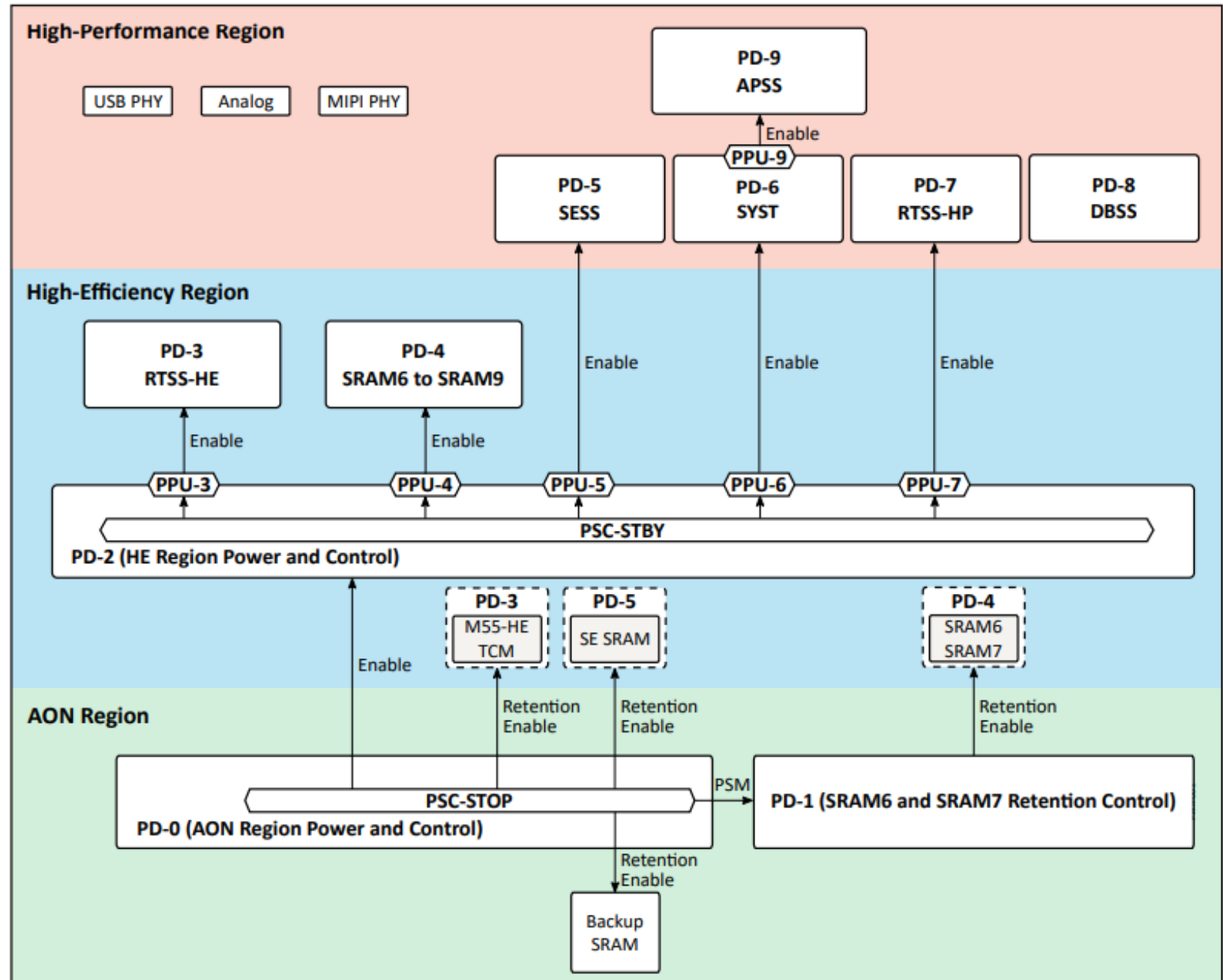
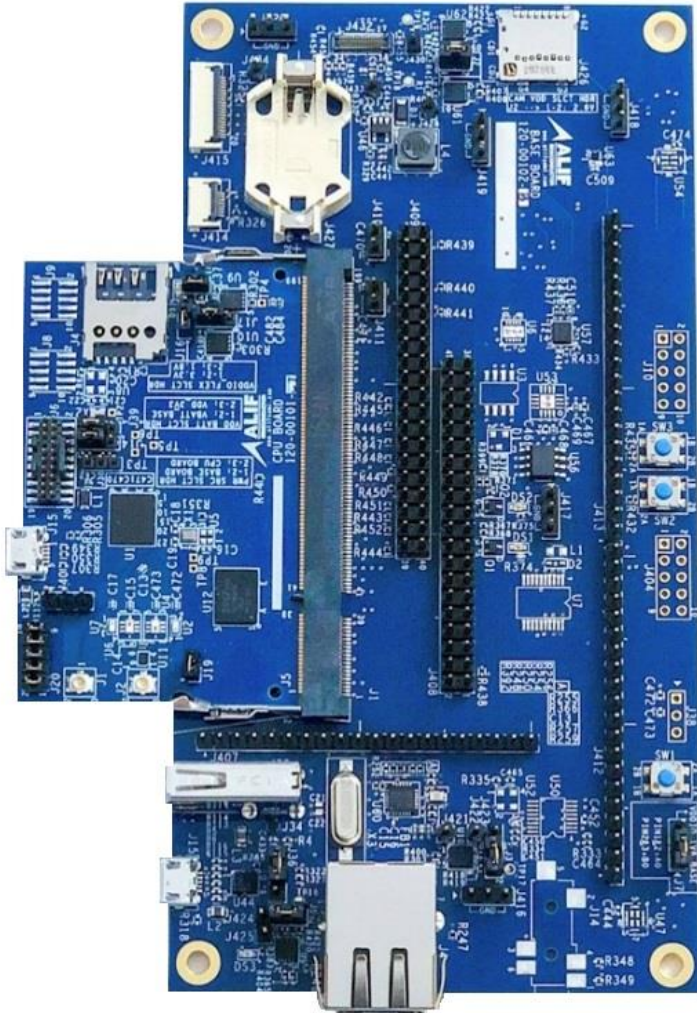


# Why Ethos-U and Alif Ensemble SoC





# Our Experience with Alif E5



# Pooling with high strides

Benefits:

Inference speed-up: +25%

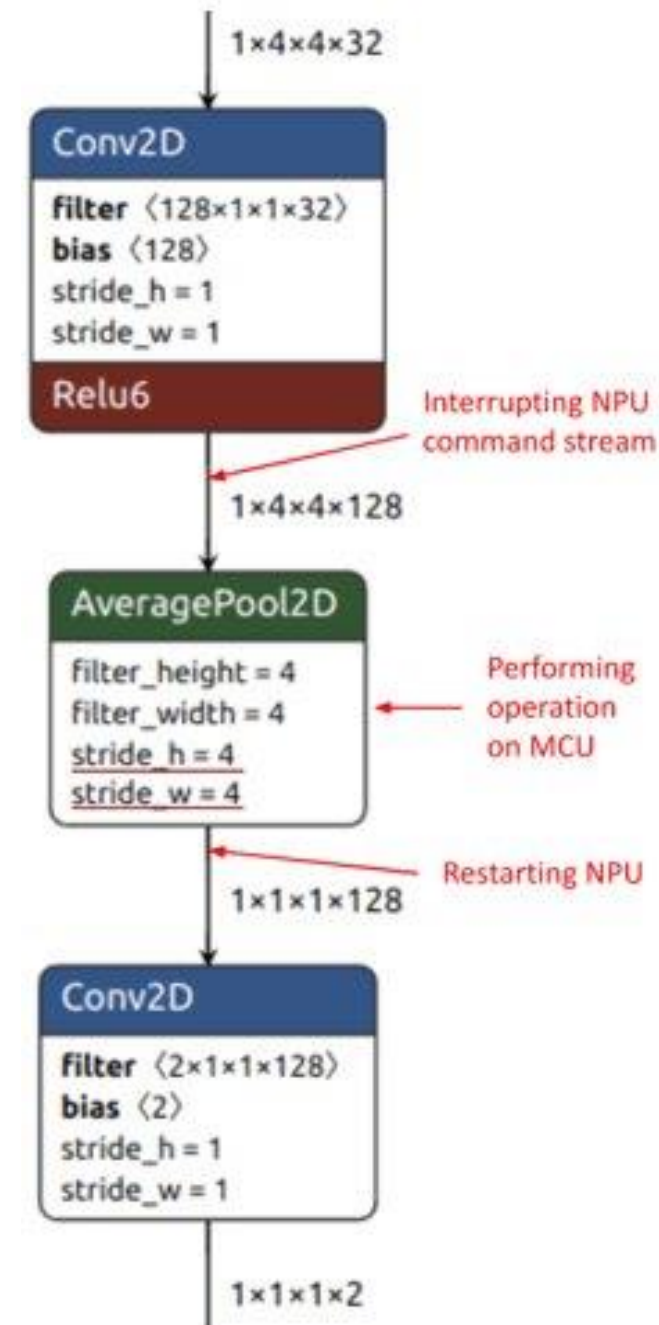
ARM ML Zoo Models affected: ~17%

```
cmsis-nn <- cmsis-nn.qnn_avg_pool2d
cmsis-nn <- %213 = cast(%212, dtype="int32")
cmsis-nn <- %214 = nn.avg_pool2d(%213, pool_size=[4, 4], strides=[4, 4], padding=[0, 0, 0, 0], layout="NHWC")
cmsis-nn <- %215 = cast(%214, dtype="int8")
```

```
strides = params.strides
if params.strides[0] > 3 or params.strides[1] > 3:
    strides = [1, 1]
```

Here we already know that IFM.shape == kernel.shape

```
ethos-u <- ethos-u.avgpool2d
ethos-u <- %213 = cast(%212, dtype="int32")
ethos-u <- %214 = nn.avg_pool2d(%213, pool_size=[4, 4], strides=[4, 4], padding=[0, 0, 0, 0], layout="NHWC")
ethos-u <- %215 = cast(%214, dtype="int8")
```



# Padding over channel axis

Benefits:

Inference speed-up: 250% - 400%

ARM ML Zoo Models affected: ~10%

```
# pad channels before
if params.ch_padding[0] > 0:
    identity1 = ethosu_ops.ethosu_identity(pad_values, ...)
    concat_args.append(identity1)

identity2 = ethosu_ops.ethosu_identity(ifm.tensor, ...)
concat_args.append(identity2)

# pad channels after
if params.ch_padding[1] > 0:
    identity3 = ethosu_ops.ethosu_identity(pad_values, ...)
    concat_args.append(identity3)

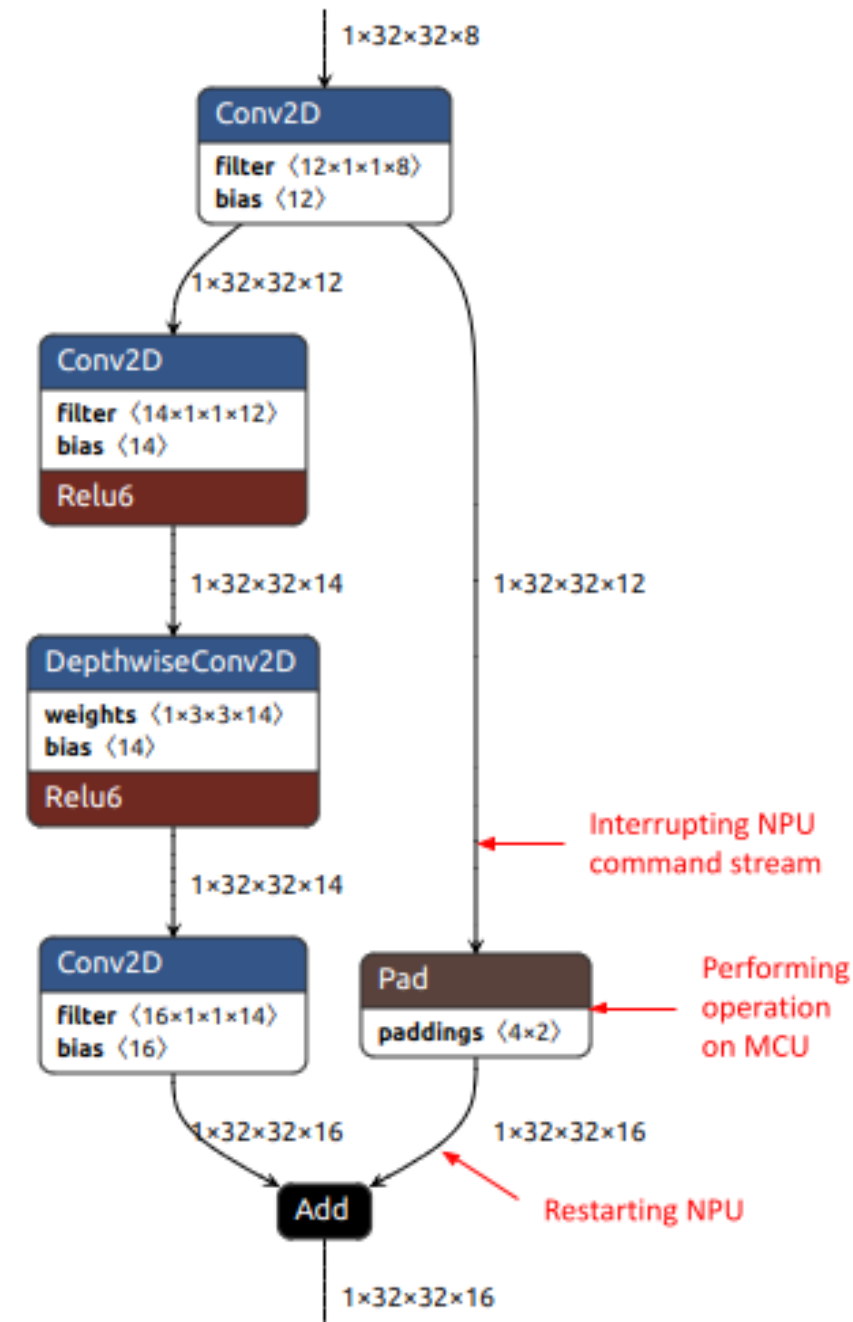
relay.op.concatenate(relay.Tuple(concat_args), axis=3)
```

Create a memory area with padding values "before"

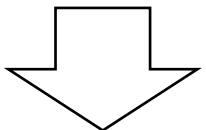
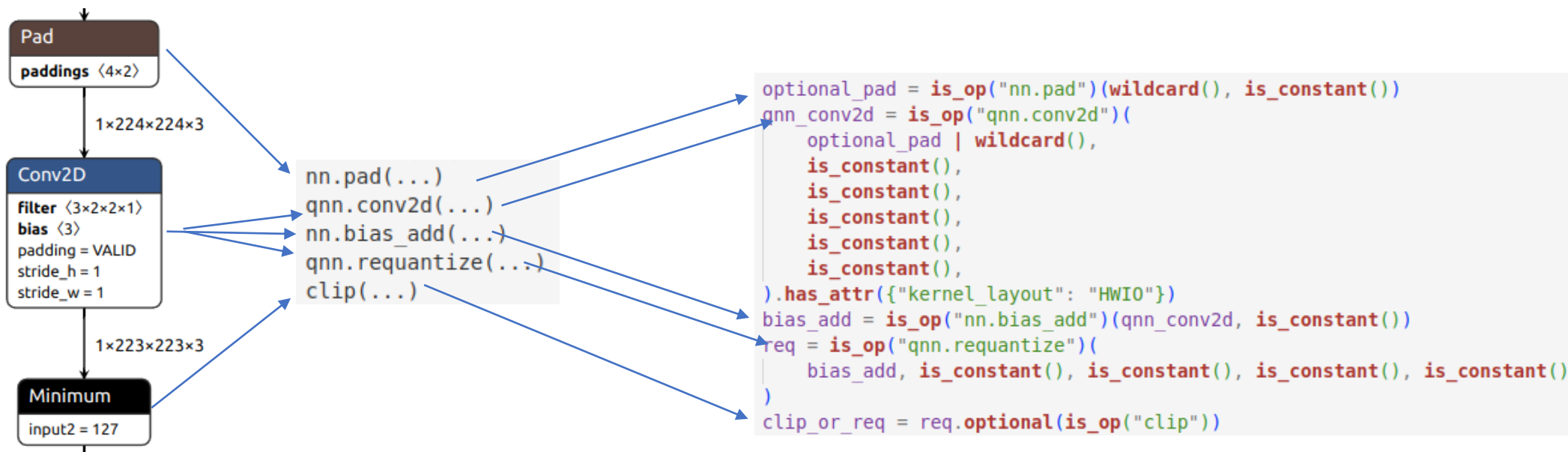
Our actual channel data

Create a memory area with padding values "after"

Concatenate everything together over channel axis



# Understanding TVM's patterns



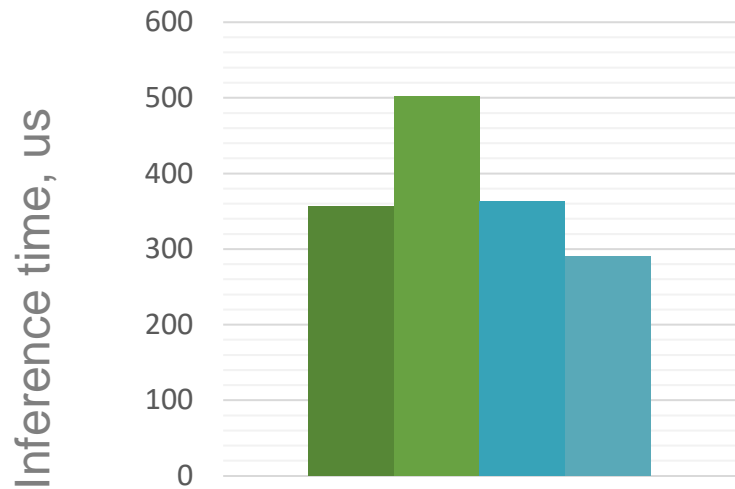
```

ethos-u <- %0 = nn.pad(%data, -128
ethos-u <- %1 = qnn.conv2d(%0, %v_
ethos-u <- %2 = nn.bias_add(%1, %v
ethos-u <- %3 = qnn.requantize(%2,
ethos-u <- %4 = clip(%3, a_min=-12
    
```

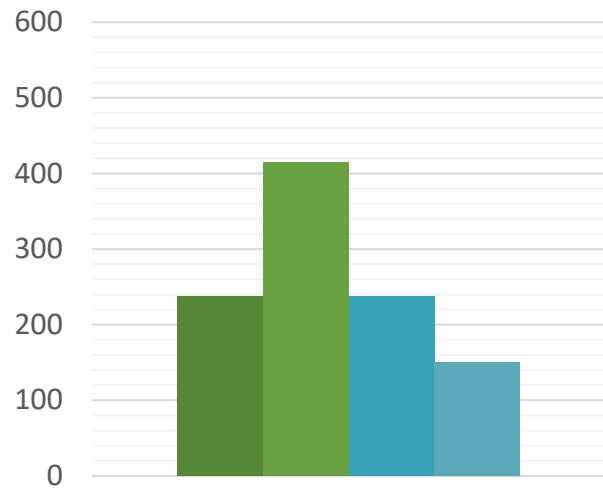


```

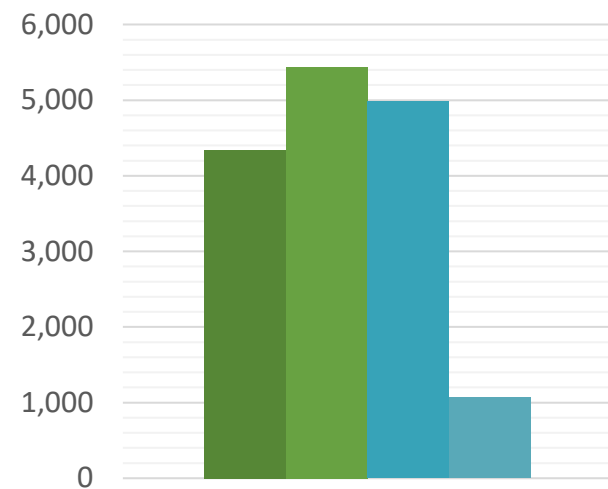
Conv2D name=None:
  IFM: h=1,w=296,c=39, region=3, NHWC, INT8, size=11544, scale: 0.17129258811473846, zero: 4
    Stride y/x/c: 1/39/1, tiles: w0=296, h0=1, h1=0, base=['0x0', '0x0', '0x0', '0x0']
    name=None
  OFM: h=1,w=148,c=250, region=1, NHCWB16, INT8, size=37000, scale: 0.24689388275146484, zero: 3
    Stride y/x/c: 1/16/2368, tiles: w0=148, h0=1, h1=0, base=['0x0', '0x0', '0x0', '0x0']
    name=None
  Kernel: w=48, h=1, stride=(2, 1), dilation=(1, 1)
  NpuPadding(top=0, left=23, bottom=0, right=23)
  Weights: (region=0, address=0x1f6090, length=271472)
  Scales: (region=0, address=0x1b810, length=2512)
  Activation: TABLE_LOOKUP, min=None, max=None, lut index=0
  NpuBlockTraversal.PART_KERNEL_FIRST
  Block config: h=2,w=60,c=64, NpuResamplingMode.NONE, NpuRoundingMode.TFL
    
```



ResNet-8

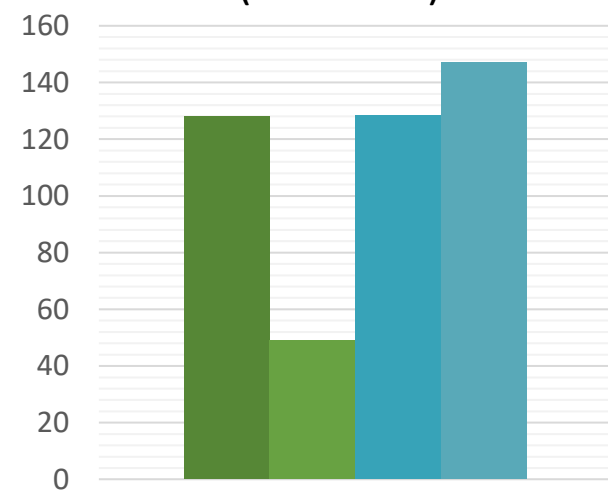
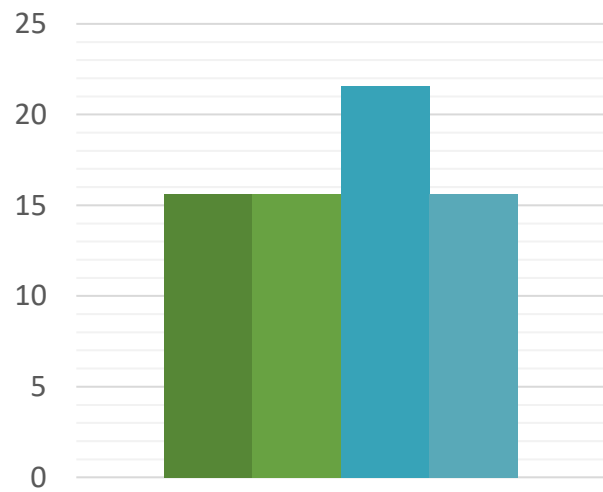
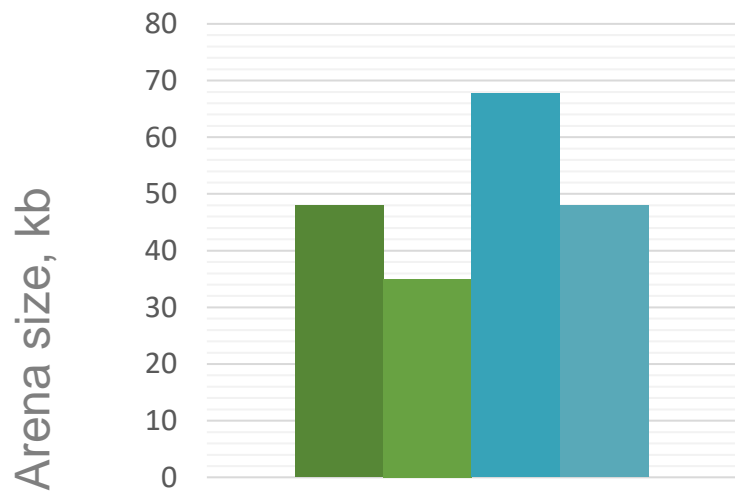


DS-CNN



VisWakeWords4  
(micronet)

- VELA Perf.
- VELA Size
- TVM 02.2023
- TVM 06.2023



## — Open Challenges

- Per-layer Analysis: computational and memory usage
- The memory scheduling according to the overall network topology [2]
- Transitioning Network Weights to External Storage [5]
- On Device Learning
- Inference of multiple NNs on heterogeneous computing architectures
- Dig in Ethos-U Platforms specific
- More practice and experience on real applications

# References

[1]: 2023 Edge AI Technology Report  
<https://www.wevolver.com/article/2023-edge-ai-technology-report>

[2]: MCUNet: Tiny Deep Learning on IoT Devices  
<http://tinymml.mit.edu>

[3]: Tiny Reservoir Computing for Extreme Learning of Motor Control  
[https://www.researchgate.net/publication/354752261\\_Tiny\\_Reservoir\\_Computing\\_for\\_Extreme\\_Learning\\_of\\_Motor\\_Control](https://www.researchgate.net/publication/354752261_Tiny_Reservoir_Computing_for_Extreme_Learning_of_Motor_Control)

[4]: Partha Pratim Ray, 2021. A review on TinyML: State-of-the-art and prospects  
<https://www.sciencedirect.com/science/article/pii/S1319157821003335>

[5]: Miao, H. and Lin, F.X., 2021. Enabling Large Neural Networks on Tiny Microcontrollers with Swapping. arXiv preprint arXiv:2101.08744.  
<https://arxiv.org/abs/2101.08744>

[6]: Work With microTVM  
[https://tvm.apache.org/docs/how\\_to/work\\_with\\_microtvm/index.html](https://tvm.apache.org/docs/how_to/work_with_microtvm/index.html)

[7]: Arm Ethos-N Processor Series, Product Brief  
<https://developer.arm.com/-/media/Arm%20Developer%20Community/PDF/AI-ML%20Datasheet%20and%20briefs/Arm%20Ethos-N%20Product%20Brief%20-%20May%202020.pdf>



AI Powered Devices

**arm** AI Partner

Thank you for your attention!

[gozman@grovety.com](mailto:gozman@grovety.com)





# Copyright Notice

This multimedia file is copyright © 2023 by tinyML Foundation. All rights reserved. It may not be duplicated or distributed in any form without prior written approval.

tinyML<sup>®</sup> is a registered trademark of the tinyML Foundation.

[www.tinyml.org](http://www.tinyml.org)



# Copyright Notice

This presentation in this publication was presented as a tinyML® Talks webcast. The content reflects the opinion of the author(s) and their respective companies. The inclusion of presentations in this publication does not constitute an endorsement by tinyML Foundation or the sponsors.

There is no copyright protection claimed by this publication. However, each presentation is the work of the authors and their respective companies and may contain copyrighted material. As such, it is strongly encouraged that any use reflect proper acknowledgement to the appropriate source. Any questions regarding the use of any materials presented should be directed to the author(s) or their companies.

tinyML is a registered trademark of the tinyML Foundation.

**[www.tinyml.org](http://www.tinyml.org)**