

tinyML[®] EMEA

Enabling Ultra-low Power Machine Learning at the Edge

tinyML EMEA Technical Forum 2021 Proceedings

June 7 – 10, 2021

Virtual Event



www.tinyML.org

The Model- Efficiency Pipeline

Enabling deep learning
inference at the edge

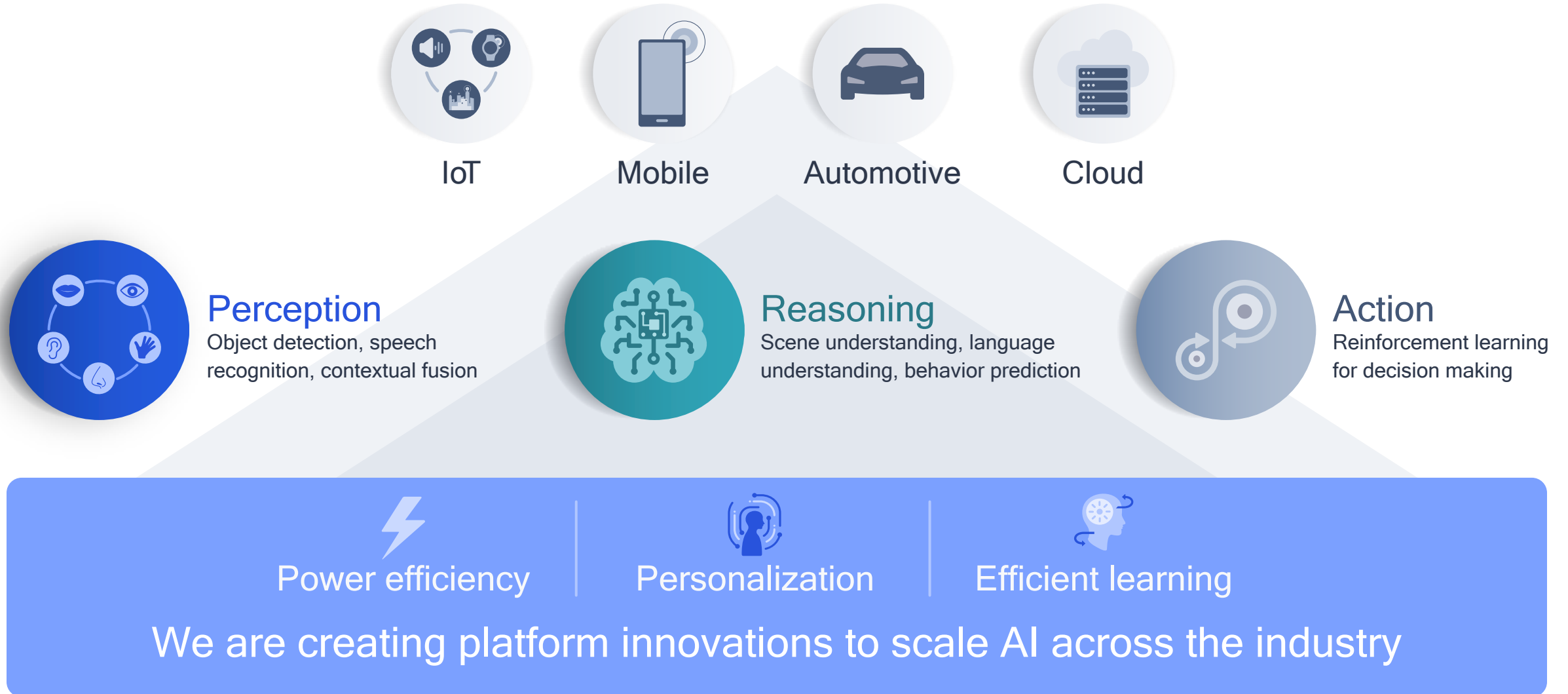
Bert Moons,
Engineer, Senior
Qualcomm Technologies Netherlands B.V.



Qualcomm AI Research



Advancing research to make AI ubiquitous



Qualcomm Research Netherlands

qualcomm.com/careers



search for Amsterdam



Agenda

- Energy-Efficient machine learning and the computational budget gap
- The Model-Efficiency Pipeline reduces the cost of on-device inference



Qualcomm Innovation Center, Inc. open sources through AI Model Efficiency Toolkit (AIMET)

- What's next in energy-efficient AI

Smartphone



Smart homes



Video conferencing



Autonomous vehicles



Smart factories



Extended reality



Smart cities



Video monitoring

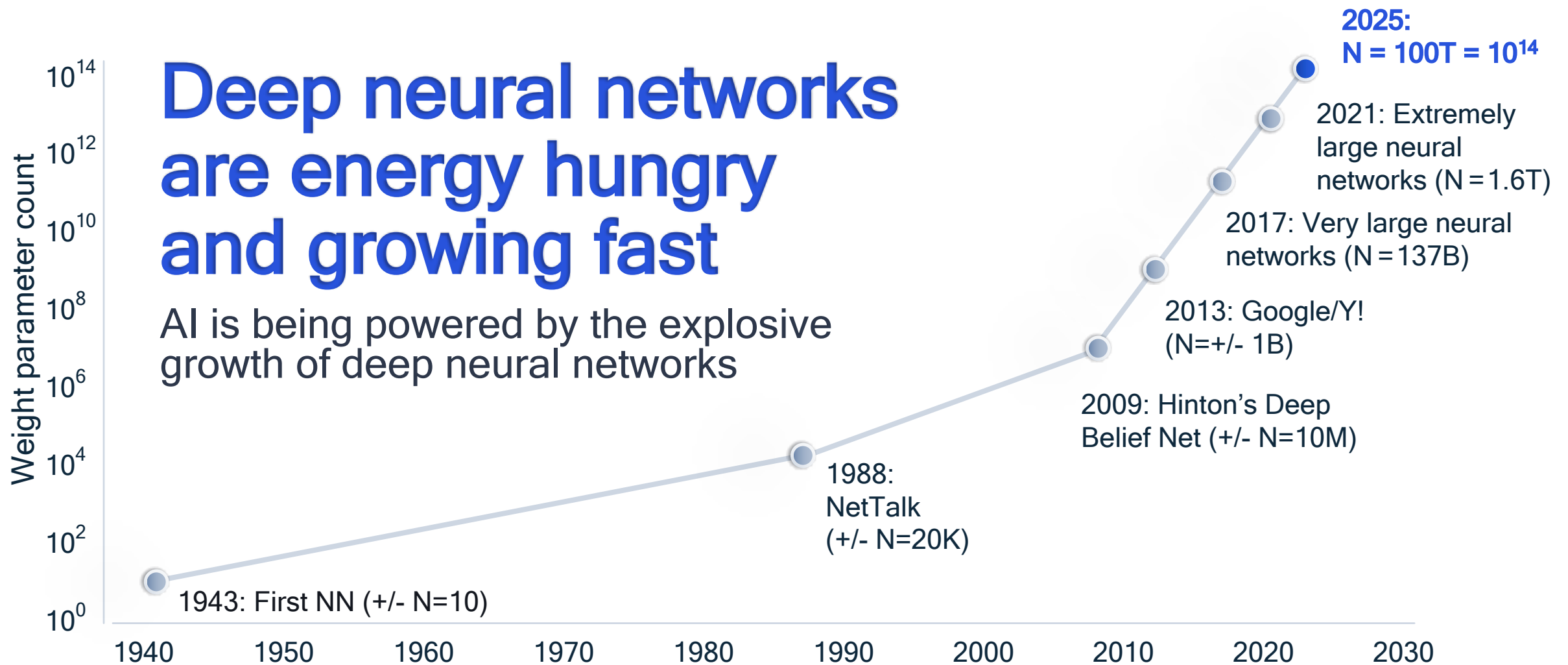


AI is being used all around us

increasing productivity, enhancing collaboration,
and transforming industries

AI video analysis is on the rise

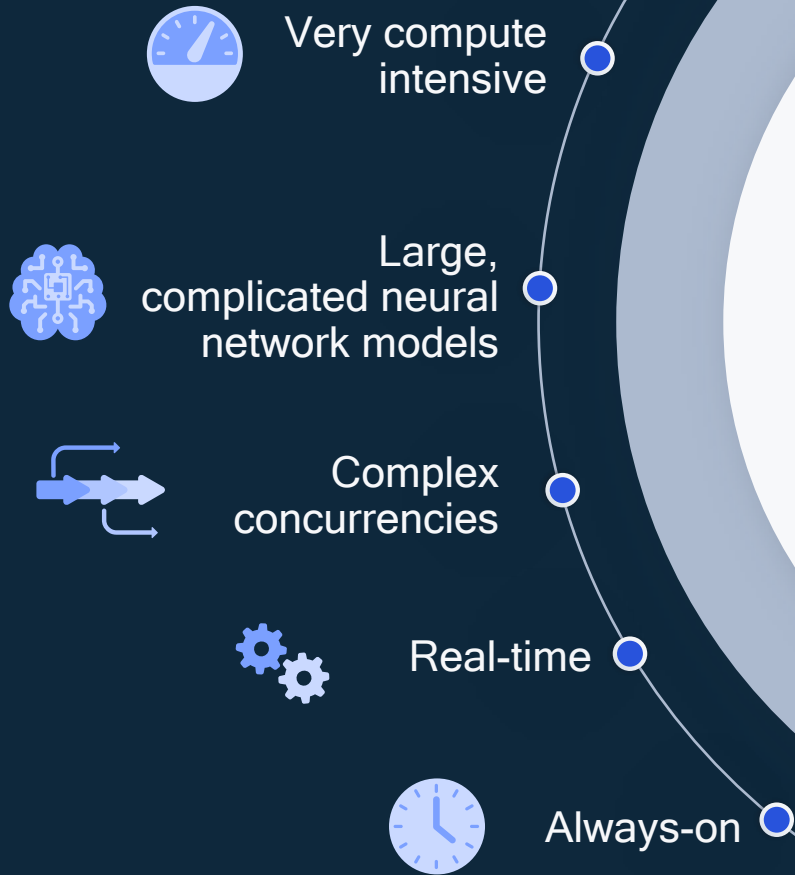
Trend toward more cameras, higher resolution,
and increased frame rate across devices



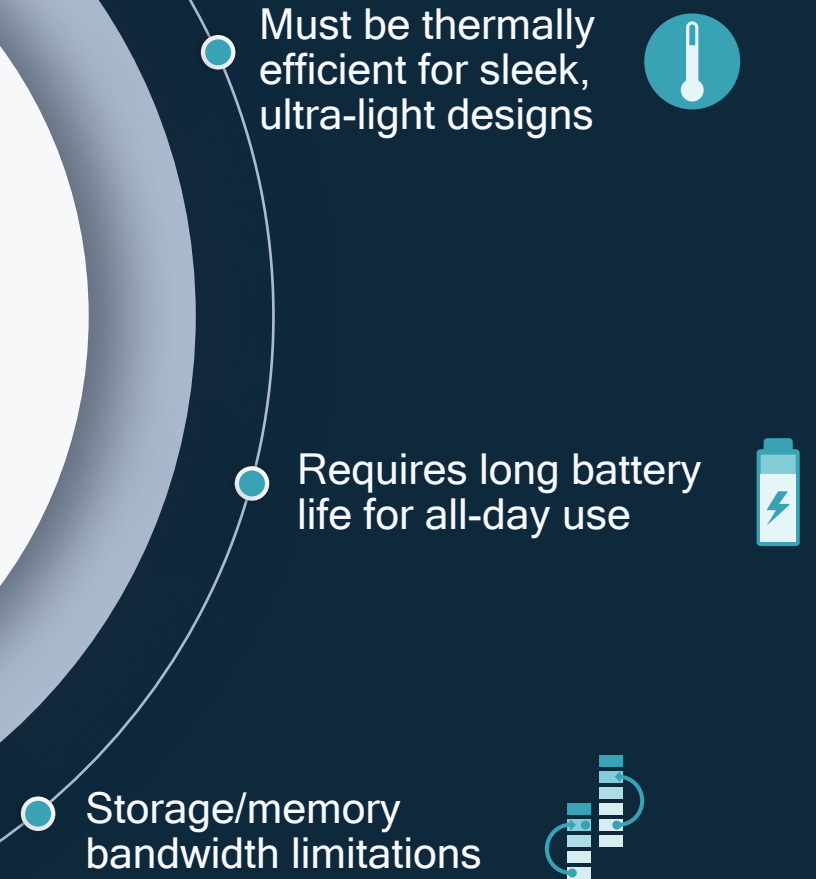
2025

Increasingly large and complex neural networks for Natural Language Processing, Image and Video Processing

The challenge of AI workloads



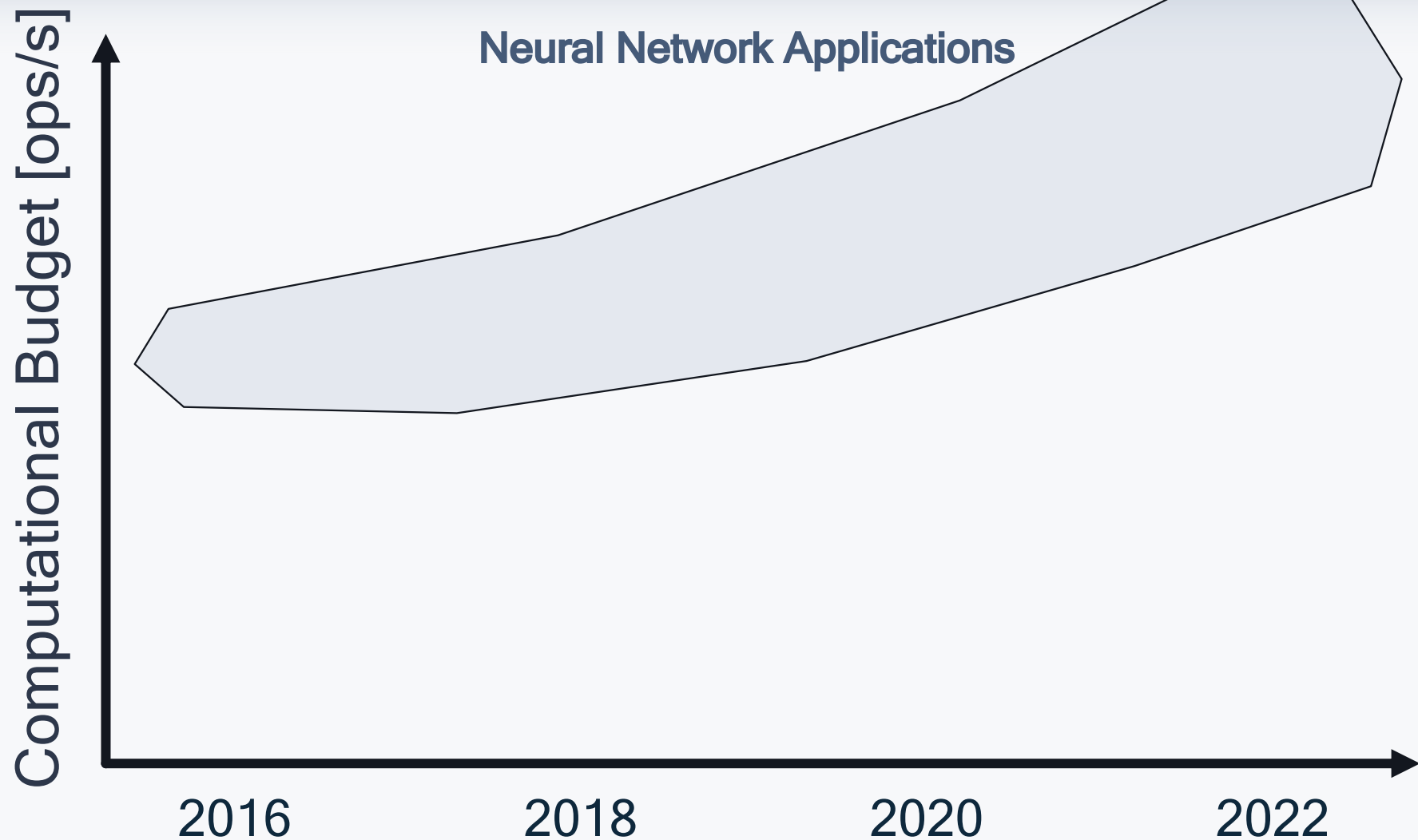
Constrained mobile environment



Power and thermal efficiency are essential for on-device AI

The Deep Learning Budget Gap

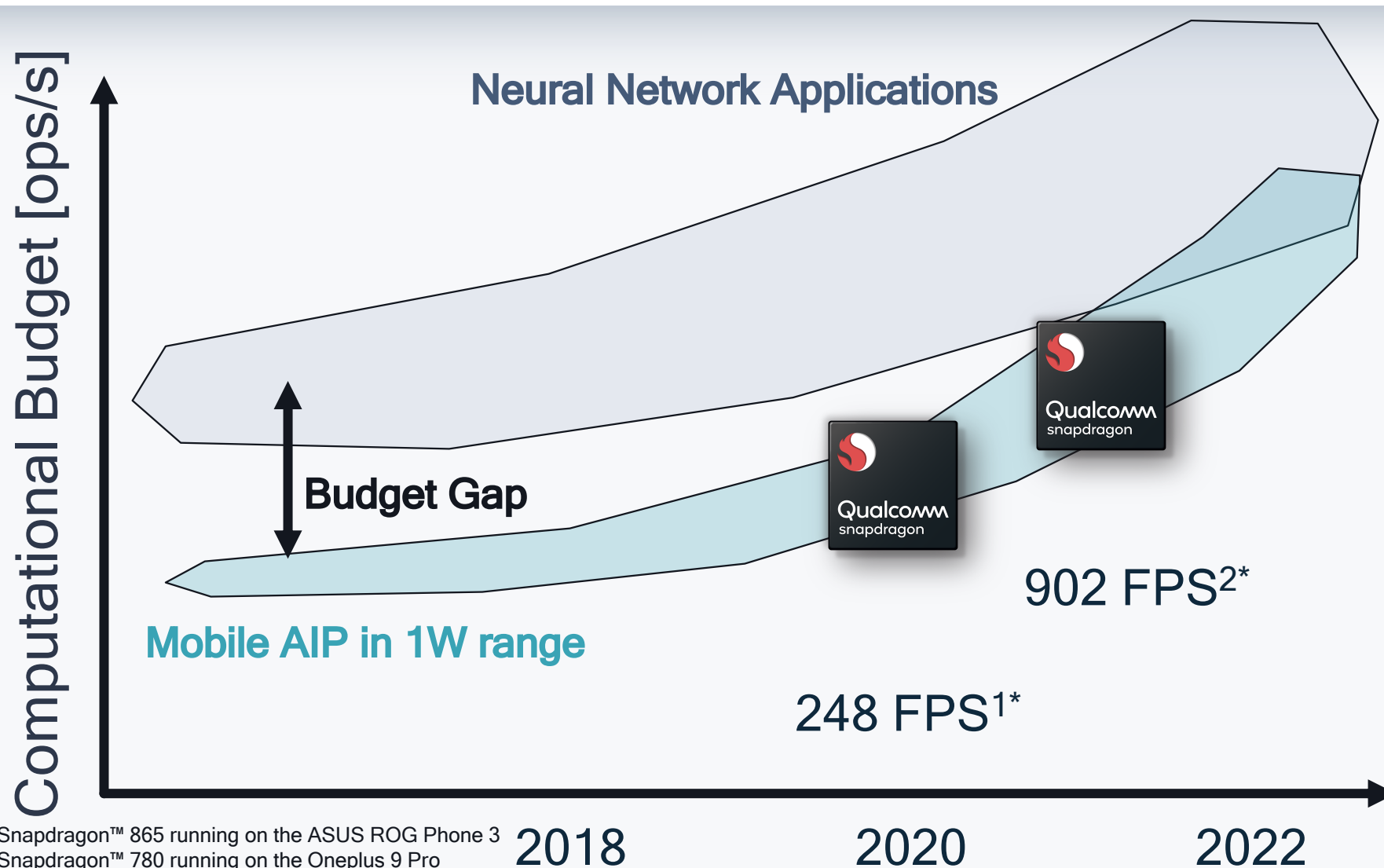
Trend 1:
Increasingly complex
Neural Networks:
Image, NLP, video,
ensembles, higher
resolution, ...



The Deep Learning Budget Gap

Trend 1:
Increasingly complex
Neural Networks:
Image, NLP, video,
ensembles, higher
resolution, ...

Trend 2:
Faster, more efficient
hardware platforms
close the Budget Gap



1: Qualcomm® Hexagon™ 698 DSP in the Qualcomm® Snapdragon™ 865 running on the ASUS ROG Phone 3

2: Qualcomm® Hexagon™ 780 DSP in the Qualcomm® Snapdragon™ 780 running on the OnePlus 9 Pro

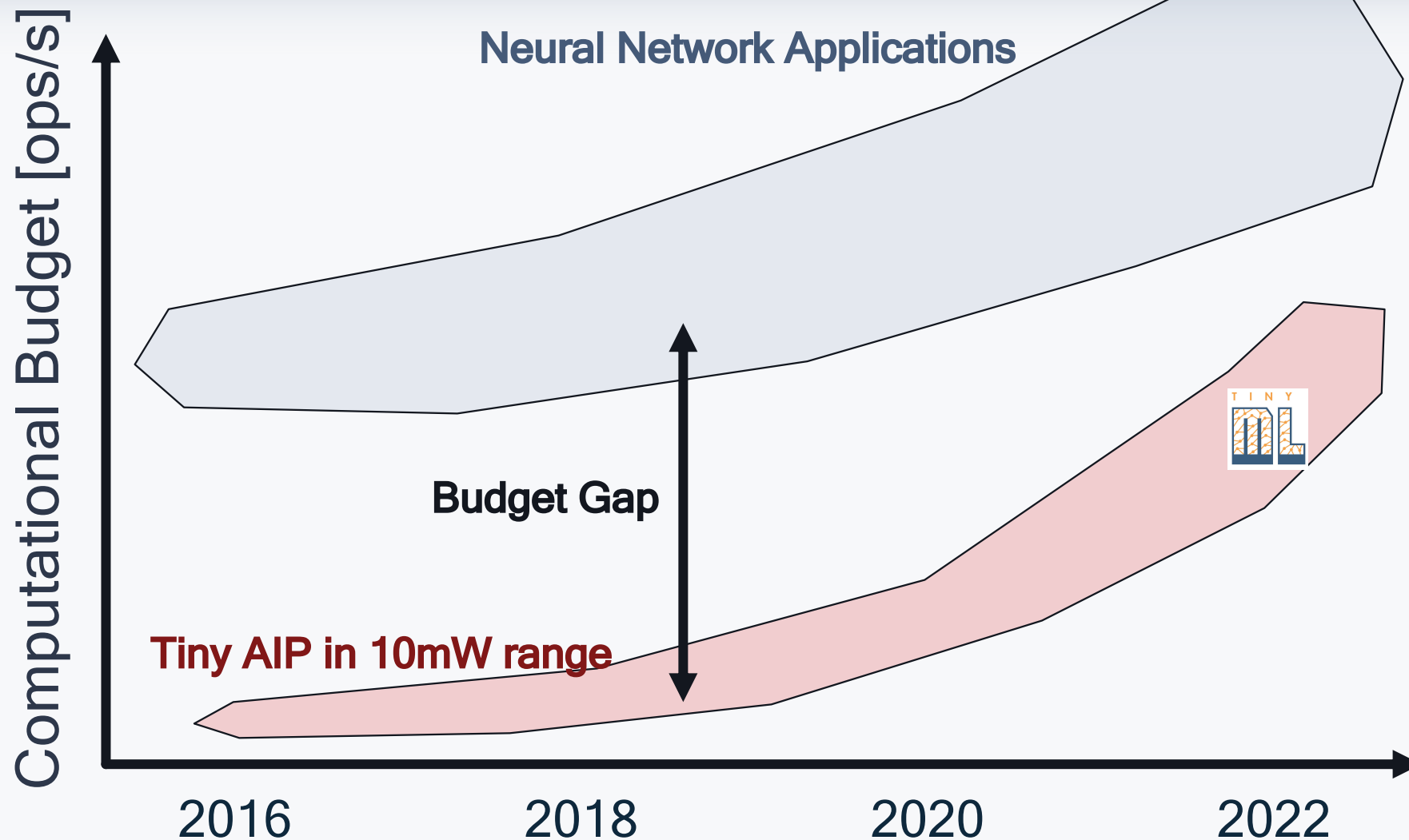
Qualcomm Hexagon and Qualcomm Snapdragon are products of Qualcomm Technologies, Inc. and/or its subsidiaries.

*: MobileNetEdge on mlperf <https://mlcommons.org/en/inference-mobile-10/>

The Deep Learning Budget Gap

Trend 1:
Increasingly complex
Neural Networks:
Image, NLP, video,
ensembles, higher
resolution, ...

Trend 2:
Faster, more efficient
hardware platforms
close the Budget Gap

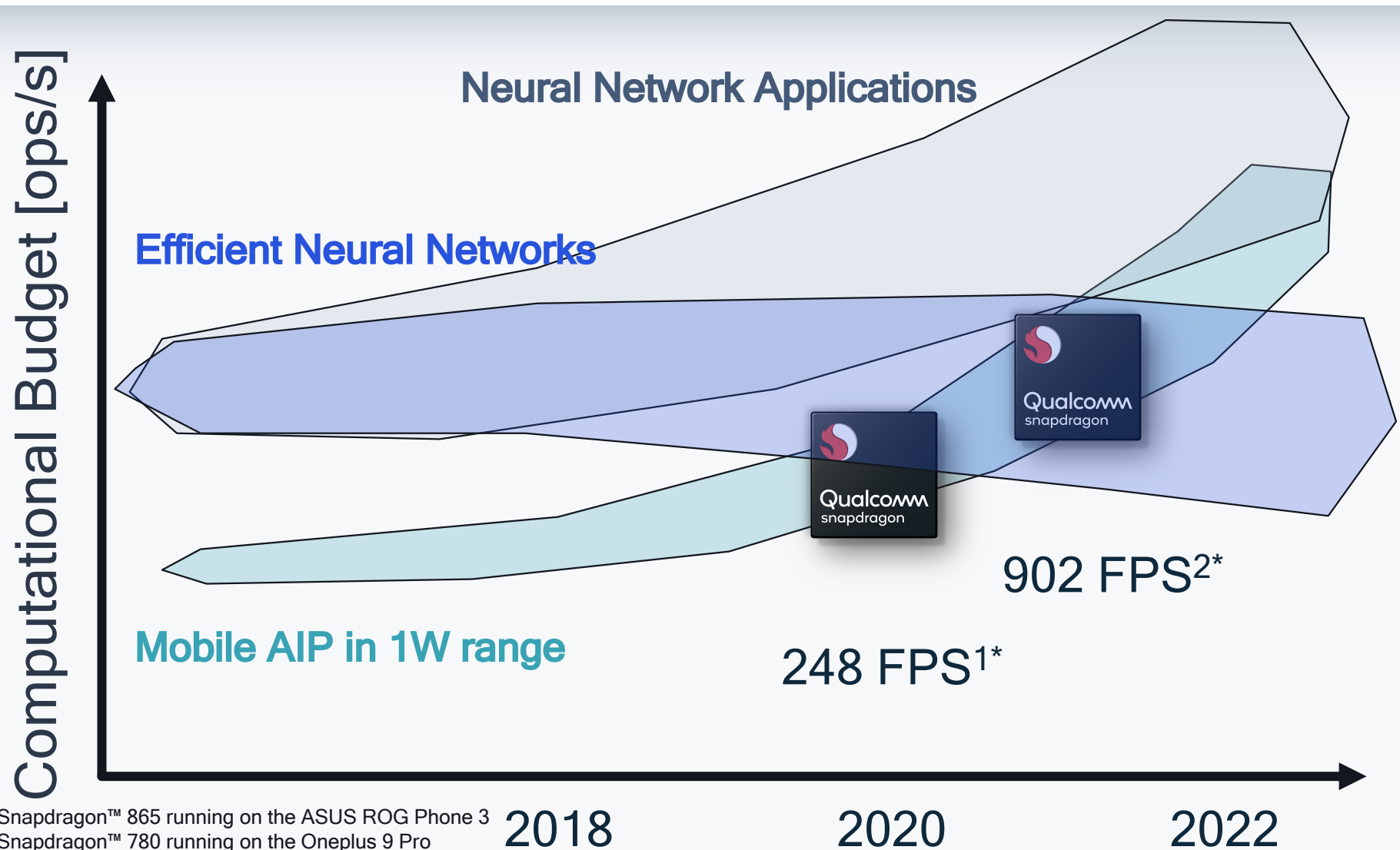


The Deep Learning Budget Gap

Trend 1:
Increasingly complex Neural Networks:
Image, NLP, video, ensembles, higher resolution, ...

Trend 2:
Faster, more efficient hardware platforms close the Budget Gap

Trend 3:
Faster, optimized Neural Networks and Applications close the Budget Gap



1: Qualcomm® Hexagon™ 698 DSP in the Qualcomm® Snapdragon™ 865 running on the ASUS ROG Phone 3

2: Qualcomm® Hexagon™ 780 DSP in the Qualcomm® Snapdragon™ 780 running on the OnePlus 9 Pro

Qualcomm Hexagon and Qualcomm Snapdragon are products of Qualcomm Technologies, Inc. and/or its subsidiaries.

*: MobileNetEdge on mlperf <https://mlcommons.org/en/inference-mobile-10/>

Trend 3: The Model-Efficiency Pipeline

The Model-Efficiency Pipeline

Multiple axes to shrink
AI models and run them
efficiently on hardware

Neural
architecture
search

Accurate
Quantization

Pruning and
Model
Compression

Neural Architecture Search: automated design of on-device optimal networks

Training networks
from scratch is
expensive!

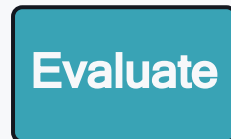
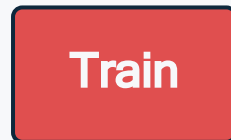
>2 GPU months to train a single SotA network on ImageNet
~4k USD per network using commercial cloud services

Neural Architecture Search: automated design of on-device optimal networks

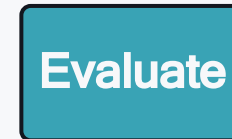
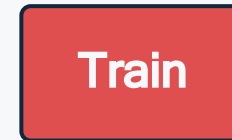
Training networks from scratch is expensive!

Manual network design requires training many networks from scratch for every device

>2 GPU months to train a single SotA network on ImageNet
~4k USD per network using commercial cloud services

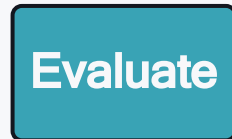
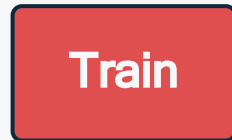
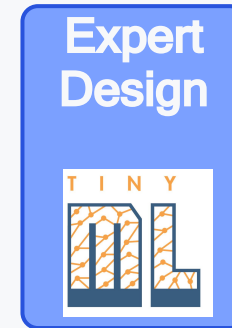


Spec A, Platform A



Spec B, Platform B

...



Spec C, Platform C

Neural Architecture Search: automated design of on-device optimal networks

Training networks from scratch is expensive!

Manual network design requires training many networks from scratch for every device

Solution

>2 GPU months to train a single SotA network on ImageNet
~4k USD per network using commercial cloud services



Train

Evaluate

Spec A, Platform A



Train

Evaluate

Spec B, Platform B



Train

Evaluate

Spec C, Platform C

Cheap, scalable Neural Architecture Search reduces design and training costs of networks optimized for specific devices

Existing NAS solutions do not address all the challenges



Lack diverse search

Hard to search in diverse spaces, with different block-types, attention, and activations
Repeated training for every new scenario



High cost

Brute force search is expensive
>40,000 epochs per platform



Do not scale

Repeated training for every device
>40,000 epochs per platform



Unreliable hardware models

Requires differentiable cost-functions
Repeated training phase for every new device

Introducing new AI research

DONNA

Distilling Optimal Neural
Network Architectures

Efficient NAS with hardware-aware
optimization

Finds pareto-optimal architectures
in terms of accuracy-latency at low
cost



Diverse search to find the best models

Supports diverse spaces with different cell-
types, attention, and activation functions
(ReLU, Swish, etc.)



Low cost

Low start-up cost equivalent to training
2-10 networks from scratch



Scalable

Scales to many hardware devices
at minimal cost



Reliable hardware measurements

Uses direct hardware measurements instead
of a potentially inaccurate hardware model

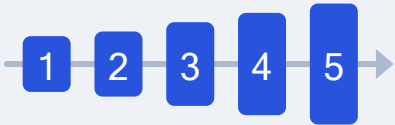
Bert Moons et al, "Distilling Optimal Neural Networks: rapid search in diverse spaces, Arxiv20

A

**Define reference
and search
space once**

Define backbone:

- Fixed channels
- Head and Stem



Varying parameters:

- Kernel Size
- Expansion Factors
- Network depth
- Network width
- Attention/activation
- Different efficient layer types

Define reference architecture and search-space once

A diverse search space is essential for finding optimal architectures with higher accuracy

Select reference architecture

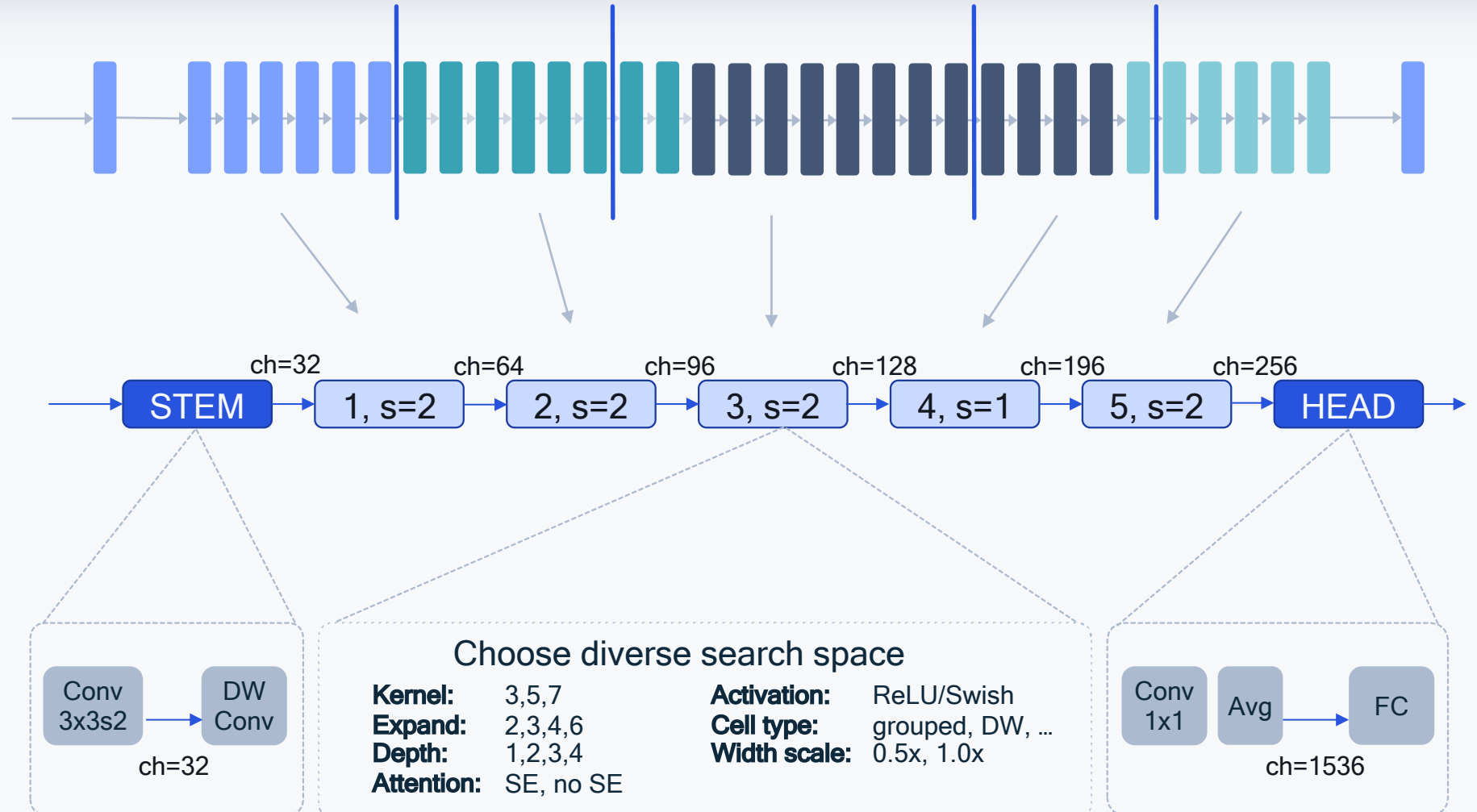
The largest model in the search-space

Chop the NN into blocks

Fix the STEM, HEAD, # blocks, strides, # channels at block-edge

Choose search space

Diverse factorized hierarchical search space, including variable cell-types, kernel-size, expansion-rate, depth, # channels, activation, attention

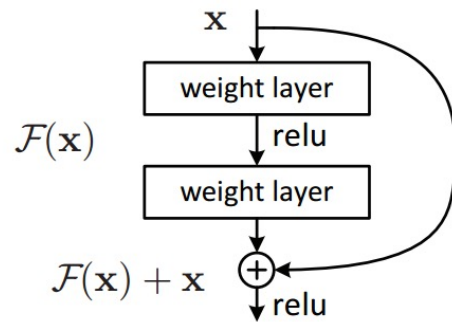
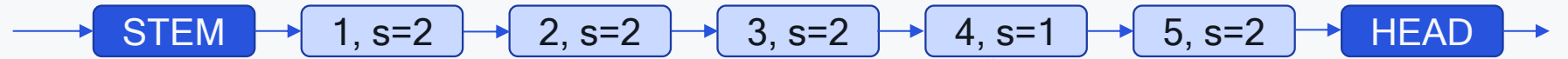


Define reference architecture and search-space once

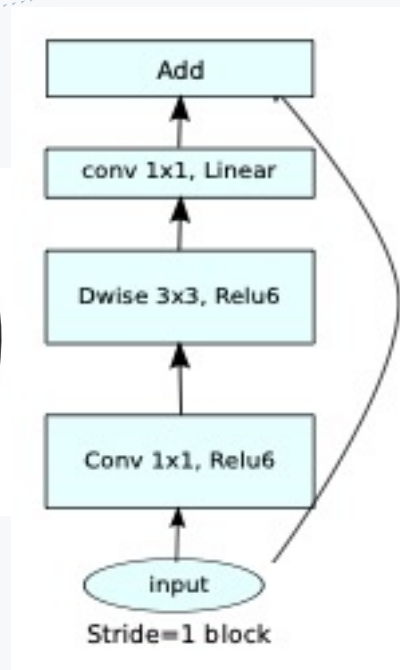
Some example blocks in the shared search space: BasicBlocks, ShiftNets, MobileConv, Squeeze-and-Excitation

Choose search space

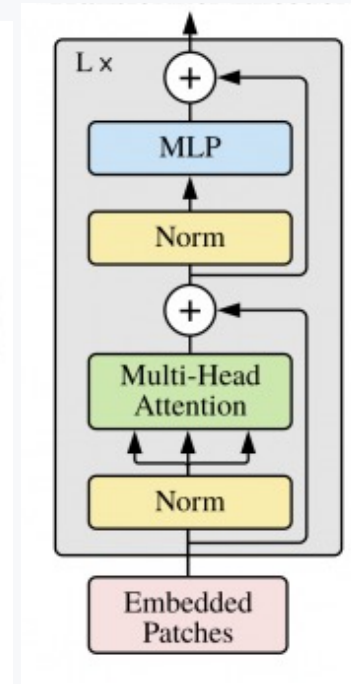
Examples of variable cell types that can be combined in a single search space



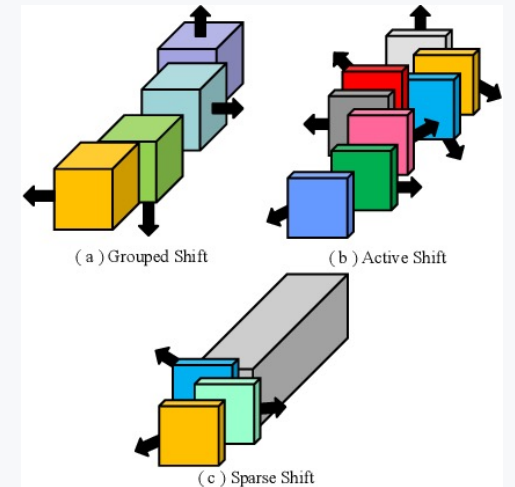
ResNet-Style
BasicBlock [1]



MobileNet-Style
Inverted Bottleneck [2]



Vision-Transformer [3]



ShiftNets [4]

[1] K. He, "Deep Residual Learning for image recognition", CVPR16

[2] M. Sandler, "MobileNetV2: Inverted Residuals and Linear Bottlenecks", CVPR18

[3] A. Dosovitskiy, "An image is Worth 16x16 words: transformers for image recognition at scale", ICLR21

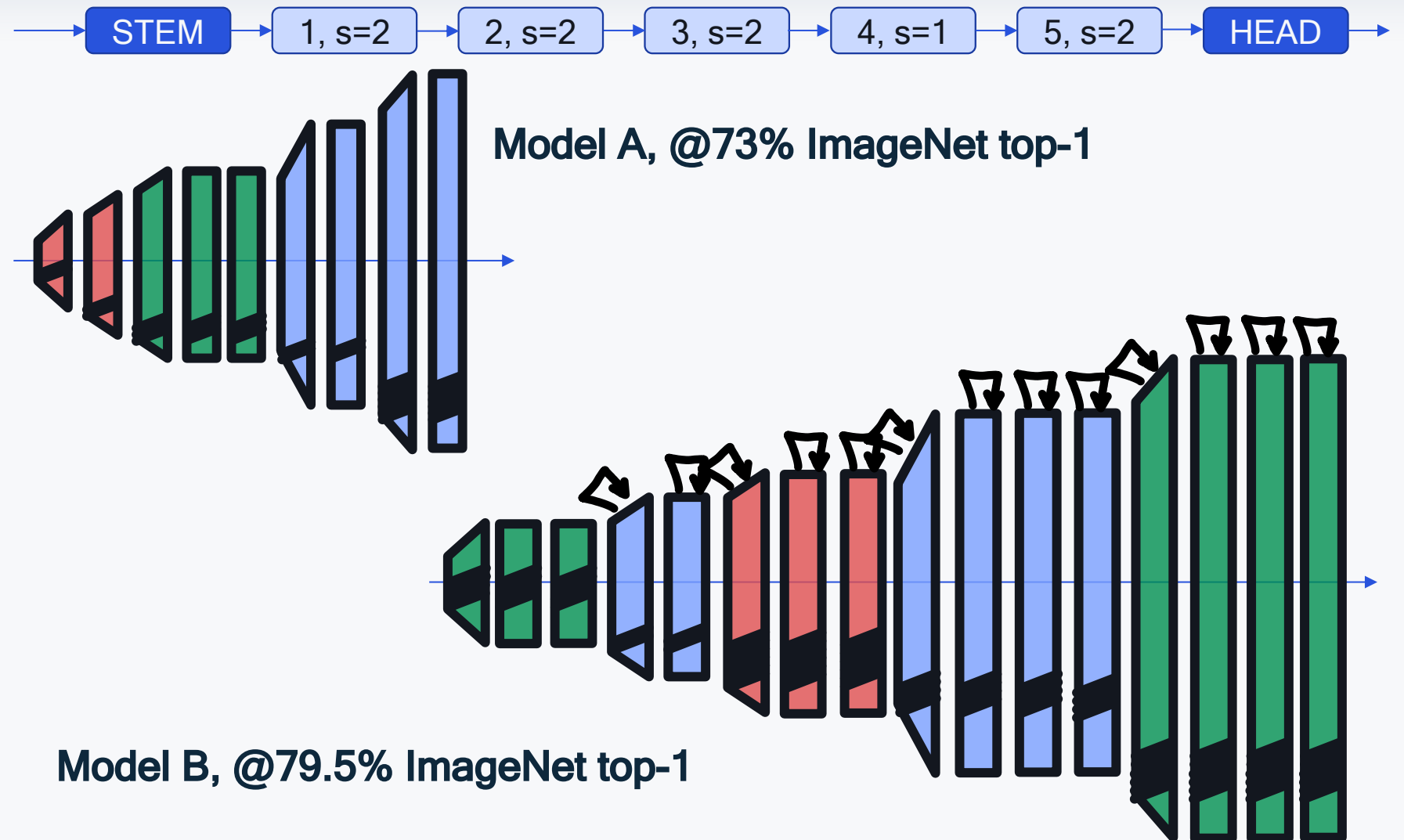
[4] W. Chen, "All you need is a few shifts: Designing efficient convolutional neural networks for image classification", CVPR19

Define reference architecture and search-space once

Some example blocks in the shared search space: BasicBlocks, ShiftNets, MobileConv, Squeeze-and-Excitation

Choose search space

Two example models, pareto-optimal on a desktop GPU



DONNA 4-step process

Objective: Build accuracy model of search space once, then deploy to many scenarios

Bert Moons et al, "Distilling Optimal Neural Networks: rapid search in diverse spaces, Arxiv20

A

Define reference and search space once

Define backbone:

- Fixed channels
- Head and Stem



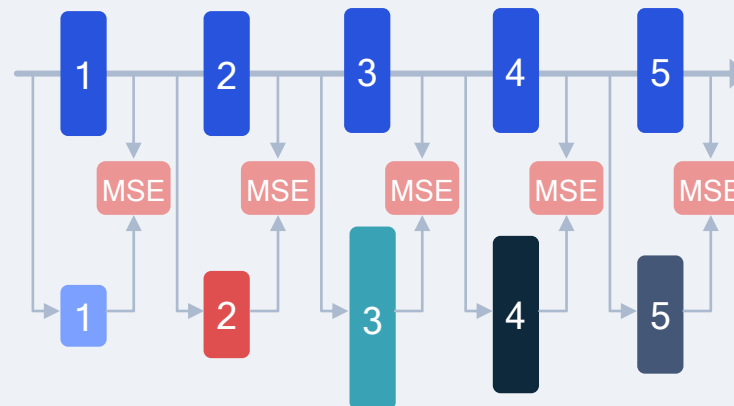
Varying parameters:

- Kernel Size
- Expansion Factors
- Network depth
- Network width
- Attention/activation
- Different efficient layer types

B

Build accuracy model via Knowledge Distillation (KD) once

Approximate ideal projections of a reference model through KD



Use quality of blockwise approximations to build accuracy model



Build accuracy predictor via Blockwise Knowledge Distillation once

Low-cost hardware-agnostic training phase

Block library

Pretrain all blocks in search-space through blockwise knowledge distillation

Block pretrained weights



Block quality metrics



Fast block training
Trivial parallelized training
Broad search space

Architecture library

Quickly finetune a representative set of architectures



Finetuned architectures

Finetune sampled networks
Fast network training
Only 20-30 NN required

Accuracy predictor

Fit linear regression model

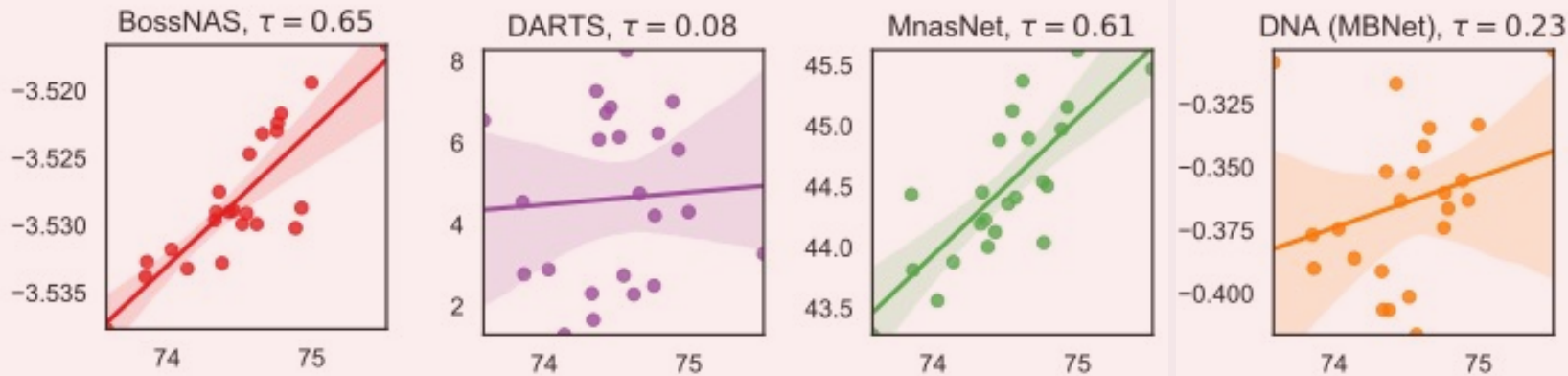


Linear Regression Model
Accurate predictions
Up to 10x improved ranking vs DARTS

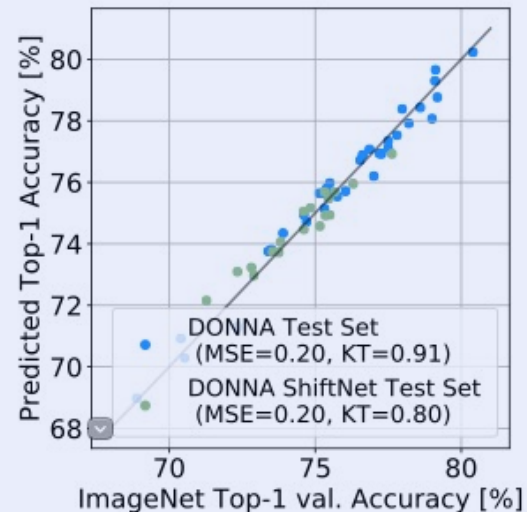
Build accuracy predictor via BKD once

Low-cost hardware-agnostic training phase

State-of-the-art references achieve up to 0.65 KT ranking*



DONNA achieves up to 0.91 KT on basic test sets and reliably extends to test sets with previously unseen cell-types: 0.8KT



Accuracy predictor

Fit linear regression model



Linear Regression Model
Accurate predictions
Up to 10x improved ranking vs DARTS

DONNA 4-step process

Objective: Build accuracy model of search space once, then deploy to many scenarios

Bert Moons et al, "Distilling Optimal Neural Networks: rapid search in diverse spaces, Arxiv20

A

Define reference and search space once

Define backbone:

- Fixed channels
- Head and Stem



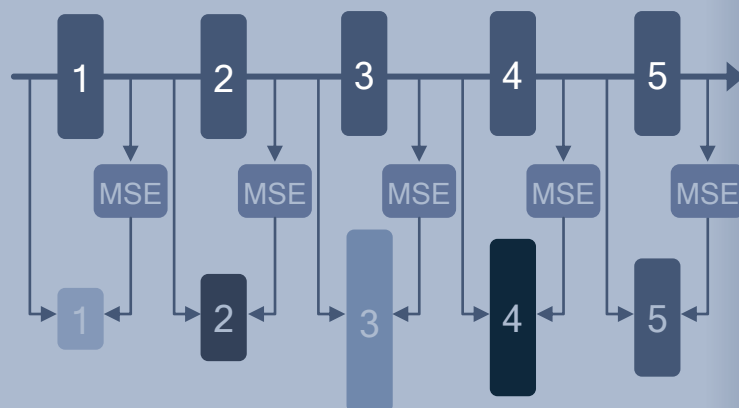
Varying parameters:

- Kernel Size
- Expansion Factors
- Network depth
- Network width
- Attention/activation
- Different efficient layer types

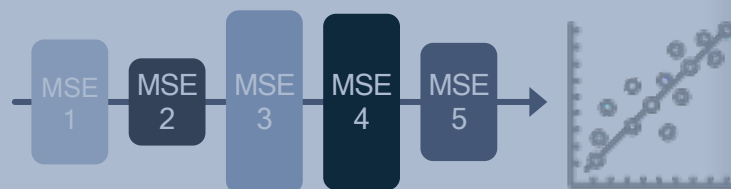
B

Build accuracy model via Knowledge Distillation (KD) once

Approximate ideal projections of a reference model through KD

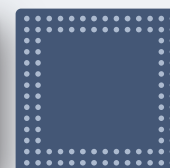


Use quality of blockwise approximations to build accuracy model



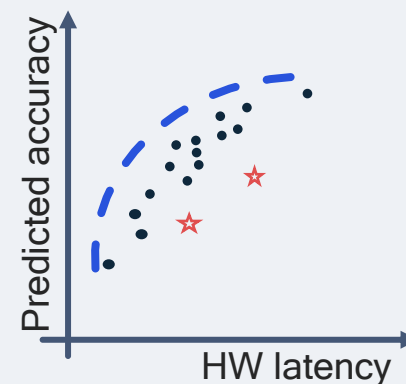
C

Evolutionary search in 24h



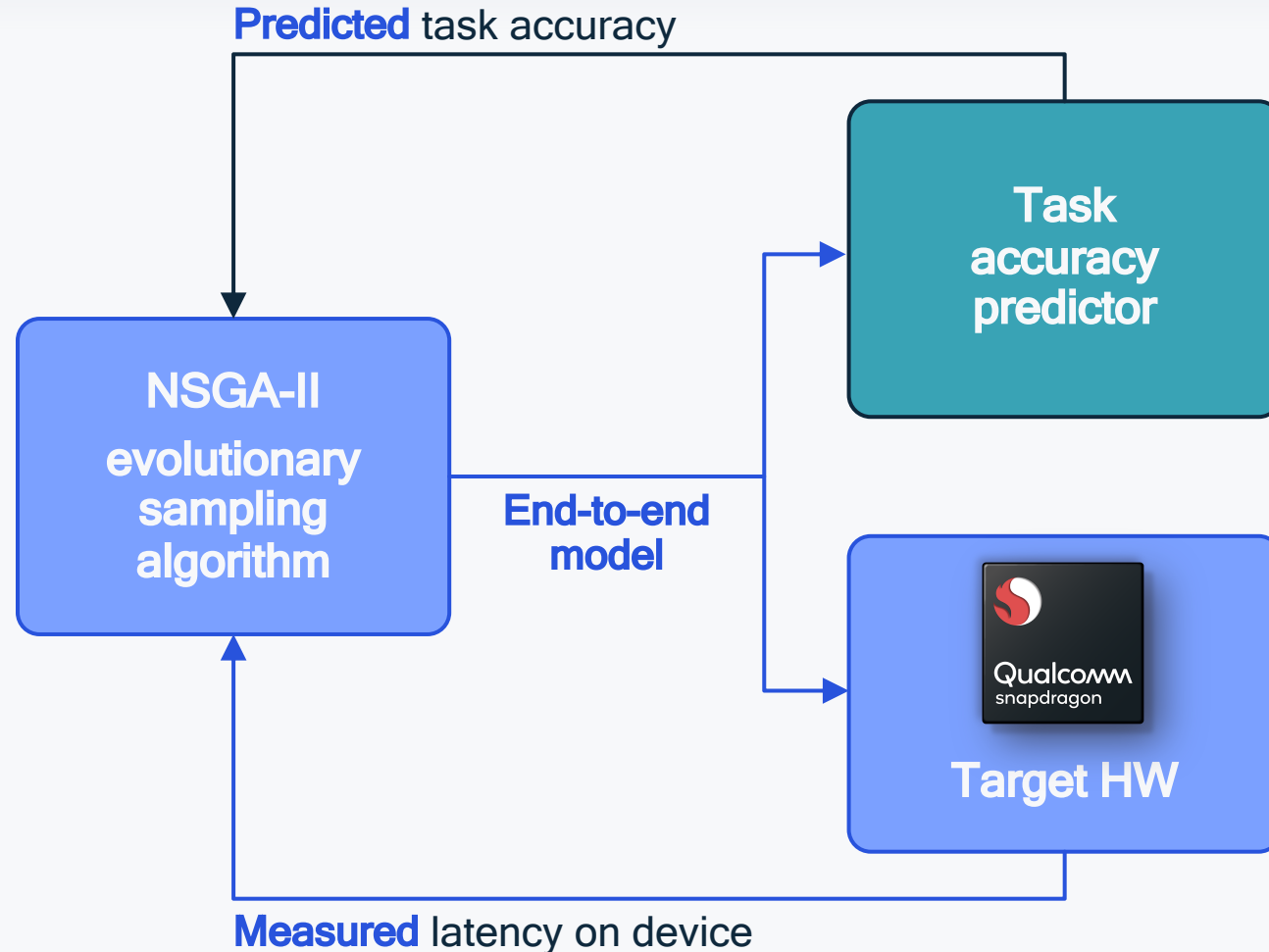
different compiler versions,
different image sizes

Scenario-specific search



Evolutionary search with real hardware measurements

Scenario-specific search allows users to select optimal architectures for real-life deployments



Quick turnaround time

Results in +/- 1 day using one measurement device

Accurate scenario-specific search

Captures all intricacies of the hardware platform and software – e.g. run-time version or devices

DONNA 4-step process

Objective: Build accuracy model of search space once, then deploy to many scenarios

Bert Moons et al, "Distilling Optimal Neural Networks: rapid search in diverse spaces, Arxiv20

A Define reference and search space once

Define backbone:

- Fixed channels
- Head and Stem

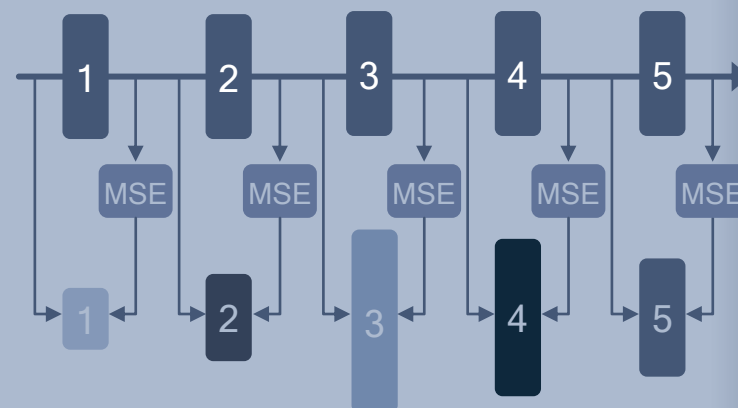


Varying parameters:

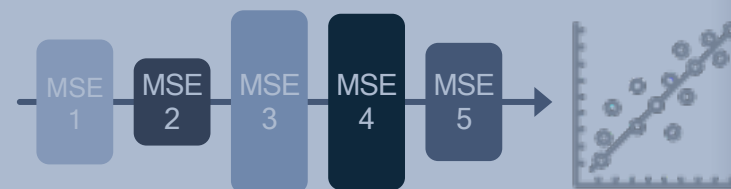
- Kernel Size
- Expansion Factors
- Network depth
- Network width
- Attention/activation
- Different efficient layer types

B Build accuracy model via Knowledge Distillation (KD) once

Approximate ideal projections of a reference model through KD



Use quality of blockwise approximations to build accuracy model

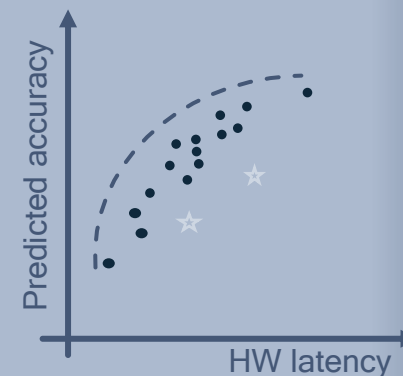


C Evolutionary search in 24h

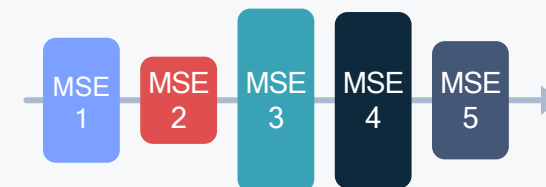


different compiler versions,
different image sizes

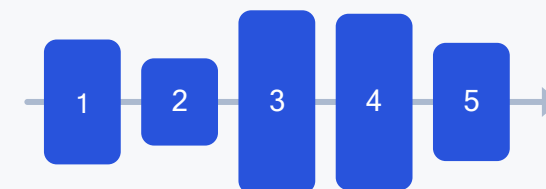
Scenario-
specific
search



D Sample and finetune

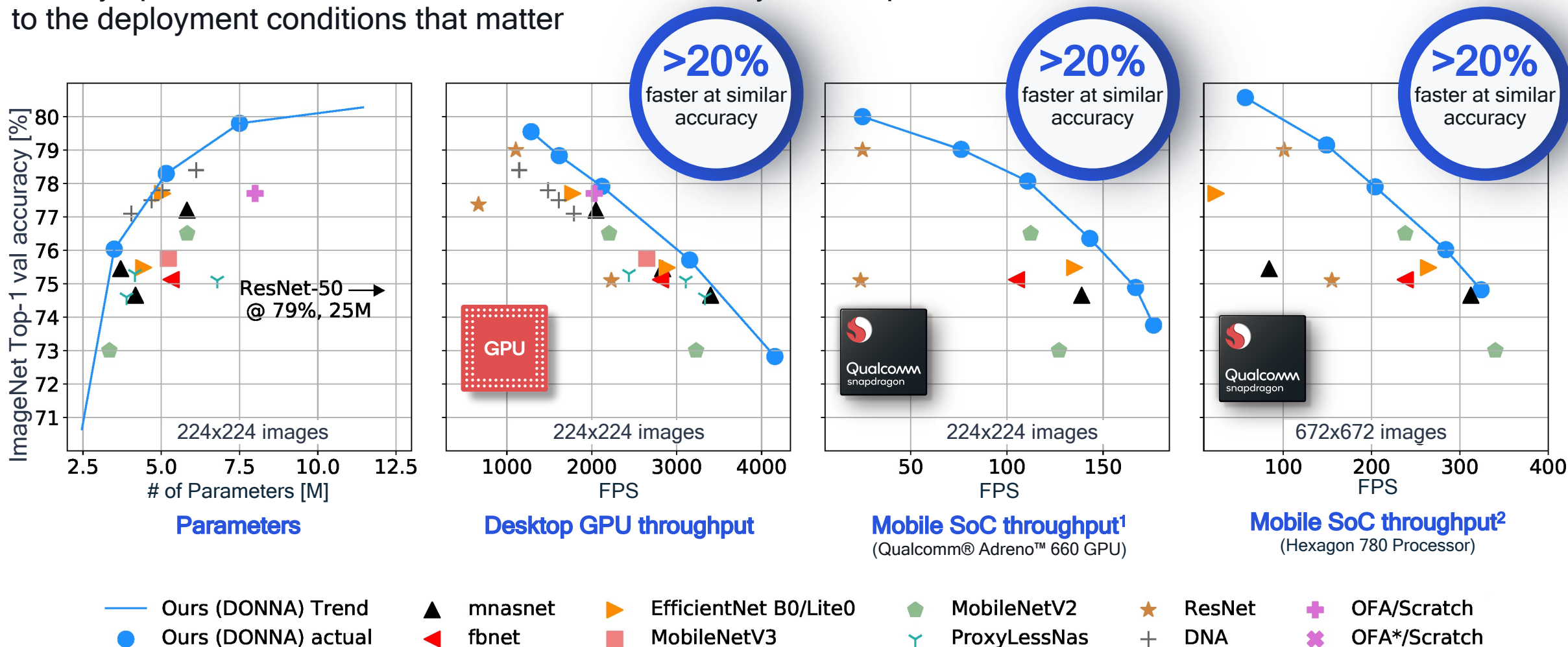


Use KD-initialized
blocks from step B
to finetune any
network in the
search space in
**15-50 epochs
instead of 450**



DONNA finds state-of-the-art networks for on-device scenarios

Quickly optimize and make tradeoffs in model accuracy with respect to the deployment conditions that matter



Qualcomm Adreno 660 GPU in the Snapdragon 888 running on the Samsung Galaxy S21. 2: Qualcomm Hexagon 780 Processor in the Snapdragon 888 running on the Samsung Galaxy S21.

Qualcomm Adreno is a product of Qualcomm Technologies, Inc. and/or its subsidiaries.

DONNA efficiently finds optimal models over diverse scenarios

Cost of training
is a handful of
architectures*

FS_e : From-Scratch-Equivalent training cost, 450 epochs

Method	Granularity	Macro-diversity	Search-cost / scenario 1 scenario, 10 models/scenario [FS_e]	Search-cost / scenario ∞ scenarios, 10 models/ scenario [FS_e]
OFA	Layer-level	Fixed	$2.7+10\times[0.05-0.15]$	0.5 - 1.5
DNA	Layer-level	Fixed	$1.5+10\times 1$	10
MNasNet	Block-level	Variable	$90+10\times 1$	100
This Work	Block-level	Variable	$9+10\times 0.1$	1

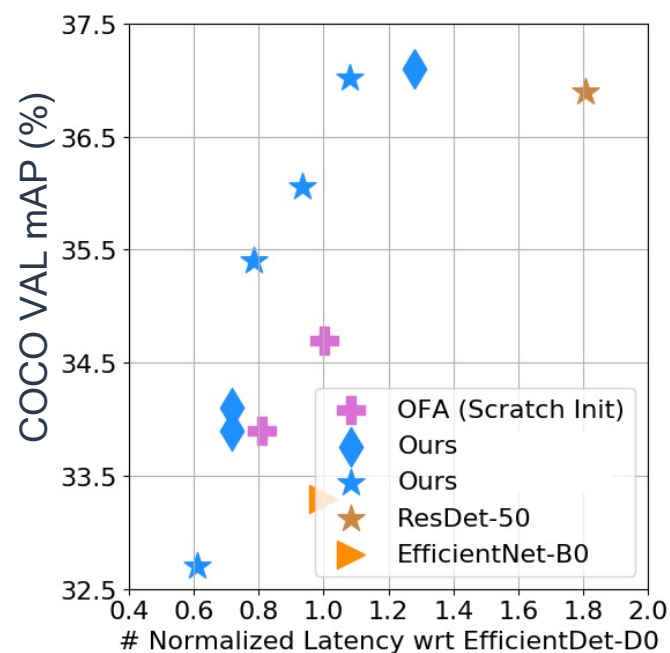
■ Good
 ■ OK
 ■ Not good

DONNA == MnasNet-level diversity at 100x lower cost

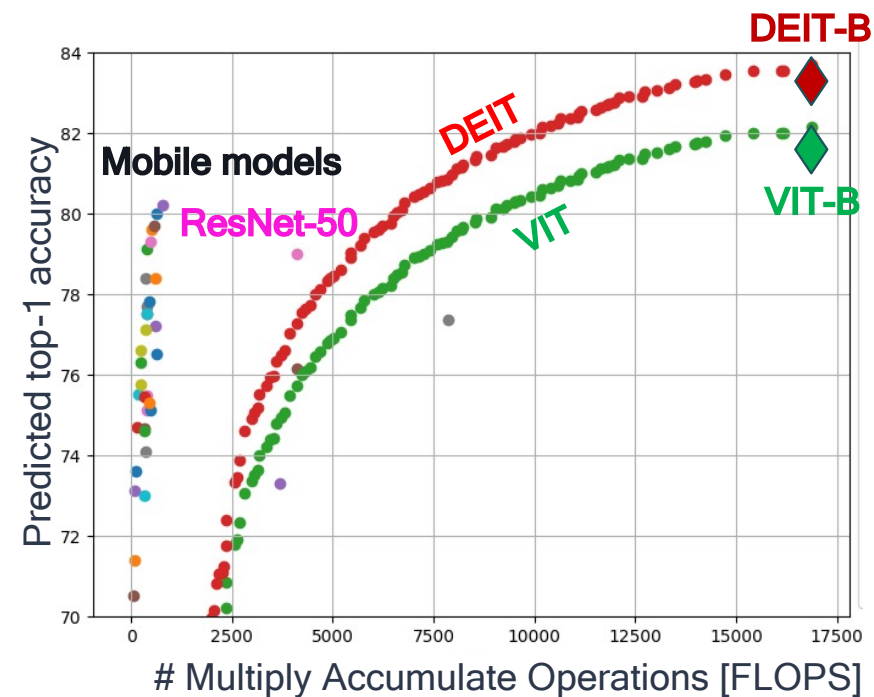
OFA = Han Cai, et al, "Once For All: Train One Network and Specialize it for Efficient deployment", ICLR2020
DNA = Changlin Li, "Blockwisely Supervised Neural Architecture Search with Knowledge Distillation", CVPR20
MNasNet = MingXing Tan, et al, "MNasNet: Platform-Aware Neural Architecture Search for Mobile", CVPR19
This work = Bert Moons et al, "Distilling Optimal Neural Networks: rapid search in diverse spaces, Arxiv20"

DONNA applies directly to downstream tasks and non-CNN neural architectures without conceptual code changes

Object Detection



Vision Transformers



The Model-Efficiency Pipeline

Multiple axes to shrink
AI models and run them
efficiently on hardware

Neural
architecture
search

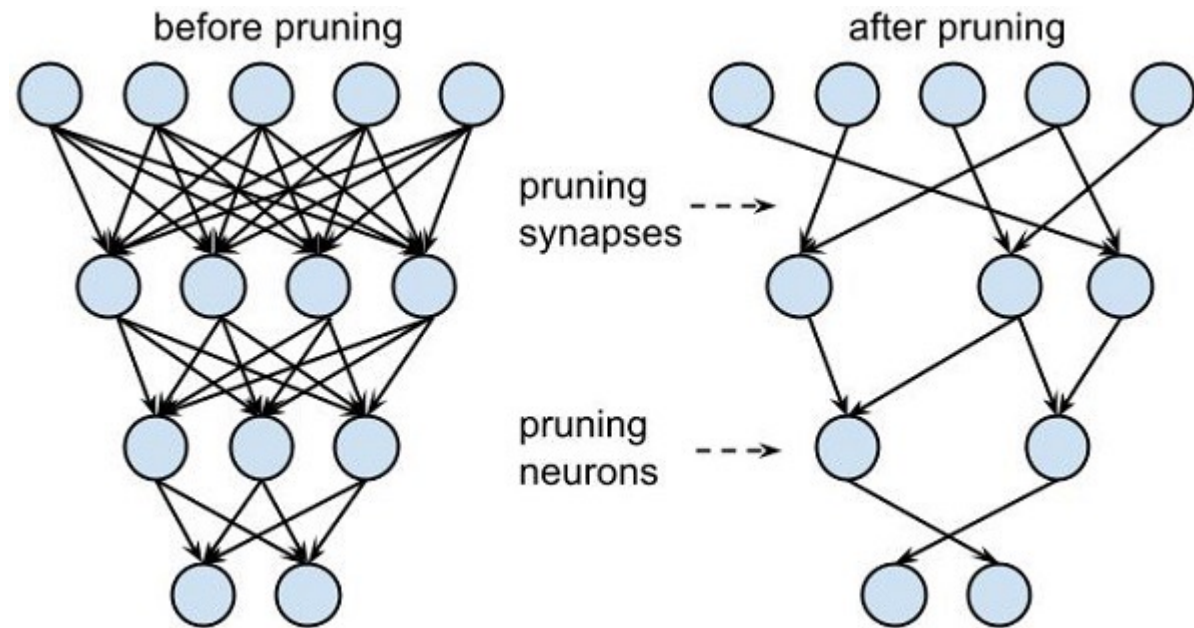
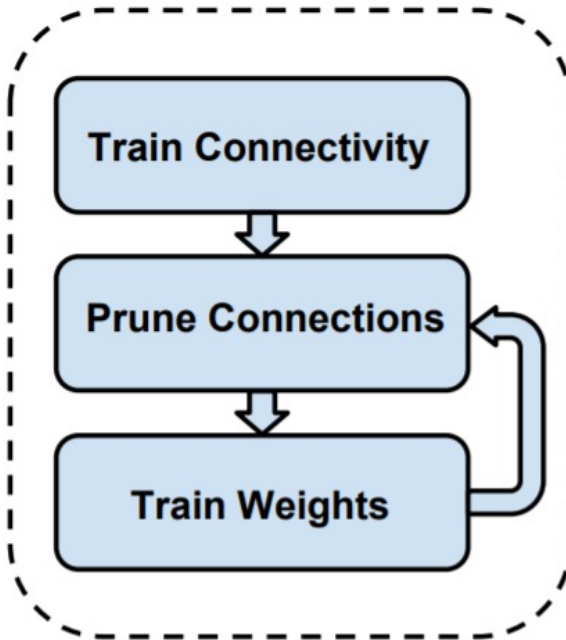
Accurate
Quantization

Pruning and
Model
Compression

Unstructured Pruning of Neural Networks

Pruning removes unnecessary connections in the neural network. Unstructured pruning is non-trivial to accelerate on parallel hardware.

Pruning: less number of weights



Structured compression through low rank approximations

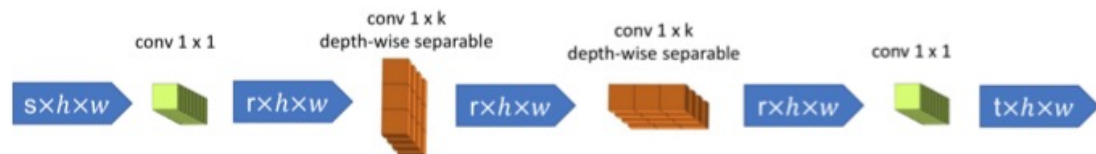
Structured mathematical decompositions (SVD, CP, Tucker-II, Tensor-train,...) are easier to accelerate on parallel hardware



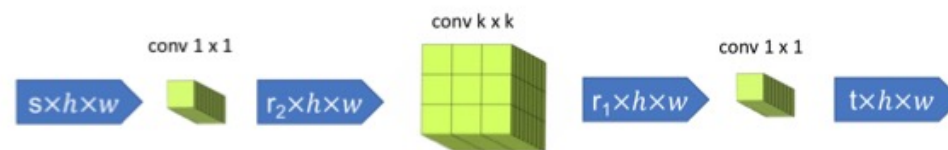
(a) Weight SVD



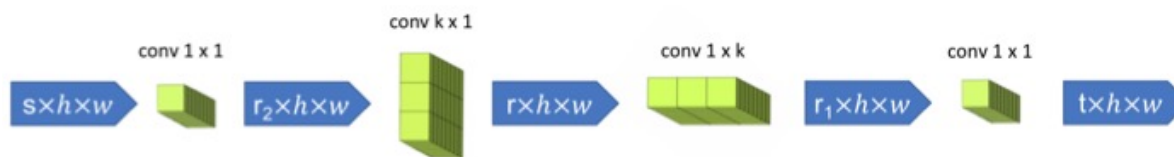
(b) Spatial SVD



(c) CP-decomposition



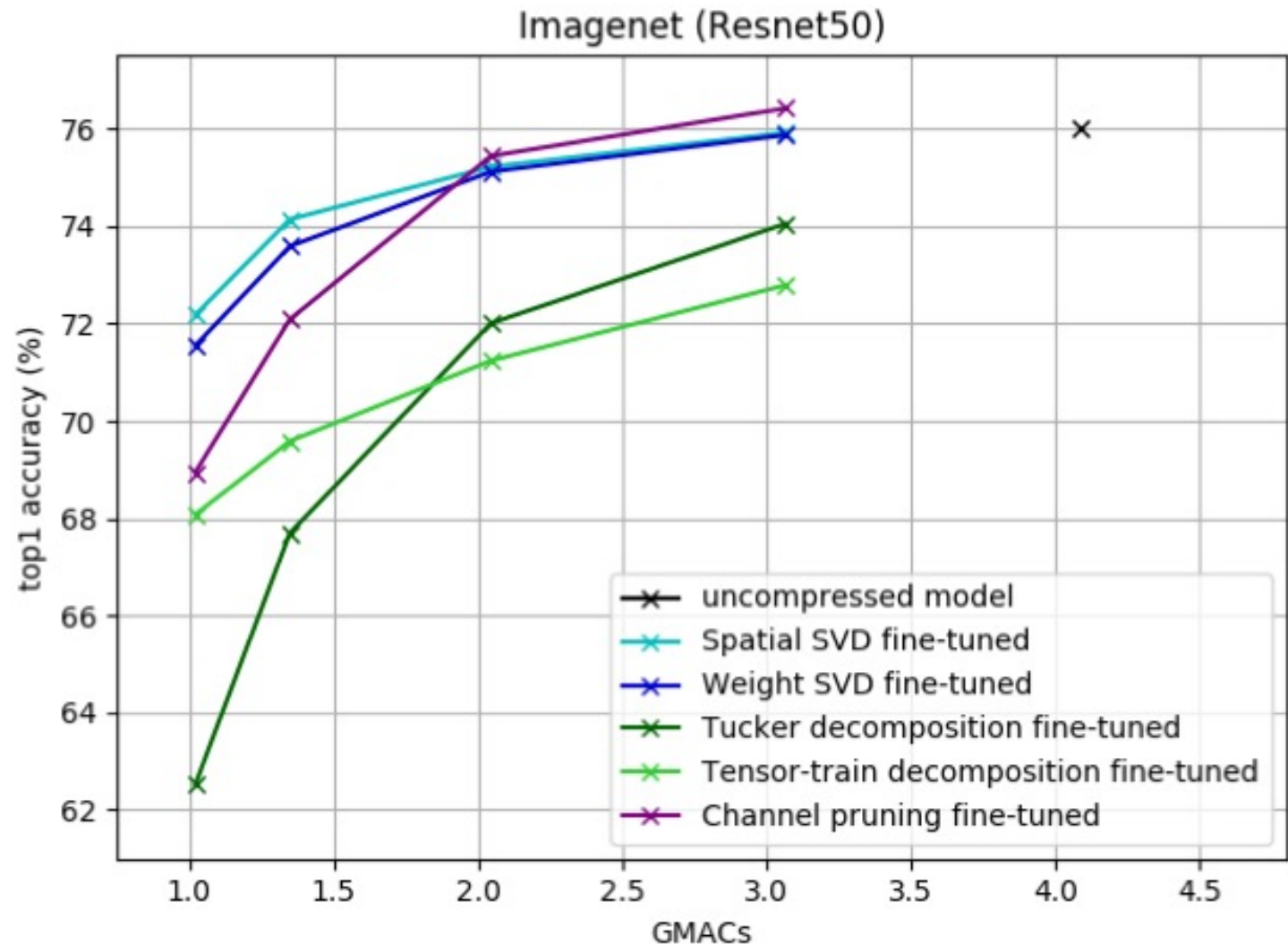
(d) Tucker decomposition



(e) Tensor-train decomposition

Structured compression through low rank approximations

- (Structured) Channel Pruning and Spatial-SVD typically work best.
- 50% compression @0.3% accuracy loss for ResNet-50



The Model-Efficiency Pipeline

Multiple axes to shrink
AI models and run them
efficiently on hardware

Neural
architecture
search

Accurate
Quantization

Pruning and
Model
Compression

What is neural network quantization?

For any given trained neural network:

- Store weights in n bits
- Compute calculations in n bits

Benefits

- Reduced memory usage
- Reduced energy usage
- Lower latency

Quantization example



Pushing the limits of what's possible with quantization

Data-free quantization

Baseline training-free method with equalization and bias correction

- ✓ No training
- ✓ Data free

AdaRound

Outperform rounding to nearest

- ✓ No training
- ✓ Minimal unlabeled data

Bayesian bits

Automated mixed-precision

- ✓ Training required
- ✓ Training data required
- ✓ Jointly learns bit-width precision and pruning

SOTA 8-bit results

<1%

Accuracy drop for MobileNet V2 against FP32 model

Data-Free Quantization Through Weight Equalization and Bias Correction (Nagel, van Baalen, et al., ICCV 2019)

SOTA 4-bit weight results

<2.5%

Accuracy drop for MobileNet V2 against FP32 model

Up or Down? Adaptive Rounding for Post-Training Quantization (Nagel, Amjad, et al., ICML 2020)

SOTA mixed-precision results

<1%

Accuracy drop for MobileNet V2 against FP32 model for mixed precision model with **computational complexity equivalent to a 4-bit weight model**

Bayesian Bits: Unifying Quantization and Pruning van Baalen, Louizos, et al., NeurIPS 2020)

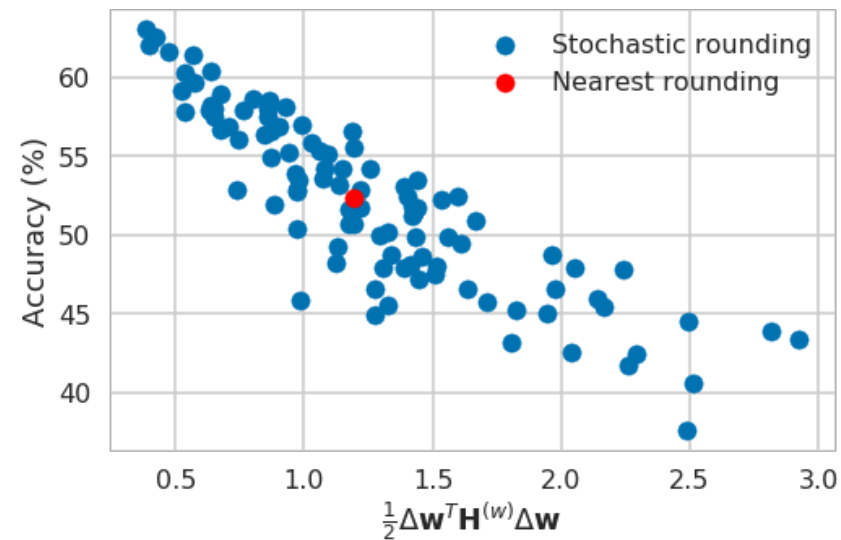
Adaround

Up or Down? Adaptive Rounding for Post-Training Quantization (Nagel, Amjad, et al., ICML 2020)

- Traditional post-training weight quantization uses **rounding to nearest**:
- However, rounding-to-nearest is not optimal

Rounding Method	Accuracy (%)
Nearest	52.29
Floor / Ceil	00.10
Stochastic	52.06±5.52
Stochastic (best)	63.06

4-bit weight quantization of 1st layer of Resnet18
,tested on ImageNet.



Up or Down?

How can we systematically find the best rounding choice?

AdaRound: learning to round

- Minimize per-layer L2 loss of output features

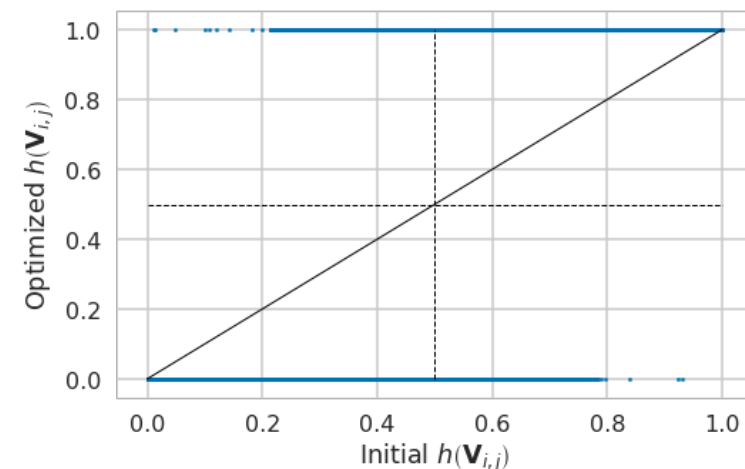
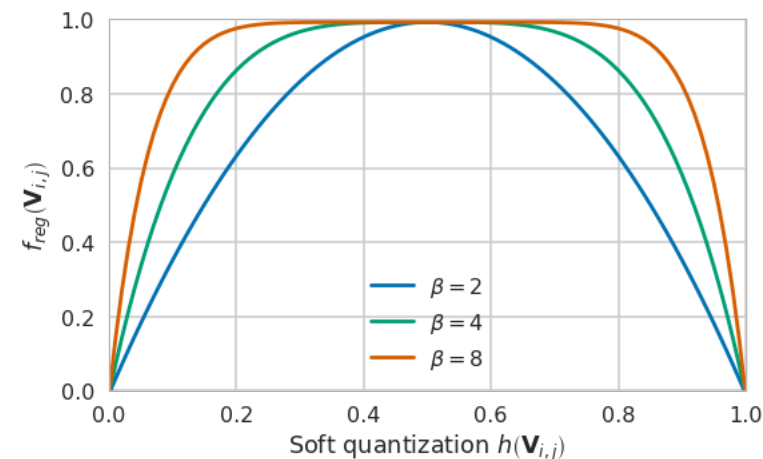
$$\arg \min_{\mathbf{V}} \left\| \mathbf{W}\mathbf{x} - \widetilde{\mathbf{W}}\mathbf{x} \right\|_F^2 + \lambda f_{reg}(\mathbf{V})$$

Regularizer forces \mathbf{V} to be 0 or 1

$$\widetilde{\mathbf{W}} = s \cdot \text{clip} \left(\left\lfloor \frac{\mathbf{W}}{s} \right\rfloor + h(\mathbf{V}), n, p \right)$$

round down + learned value between [0,1]

- Regularization: $f_{reg}(\mathbf{V}) = \sum_{i,j} 1 - |2h(\mathbf{V}_{i,j}) - 1|^\beta$



Comparison to literature

Setting a new SOTA for 4-bit post-training weight quantization

Optimization	#bits W/A	Resnet18	Resnet50	InceptionV3	MobilenetV2
Full precision	32/32	69.68	76.07	77.40	71.72
DFQ (Nagel et al., 2019)	8/8	69.7	-	-	71.2
DFQ (our impl.)	4/8	38.98	52.84	-	46.57
Bias corr (Banner et al., 2019)	4*/8	67.4	74.8	59.5	-
AdaRound w/ act quant	4/8	68.55±0.01	75.01±0.05	75.72±0.09	69.25±0.06[†]

Table 7. Comparison among different post-training quantization strategies in the literature. We report results for various models in terms of ImageNet validation accuracy (%). *Uses per channel quantization. [†]Using CLE (Nagel et al., 2019) as preprocessing.

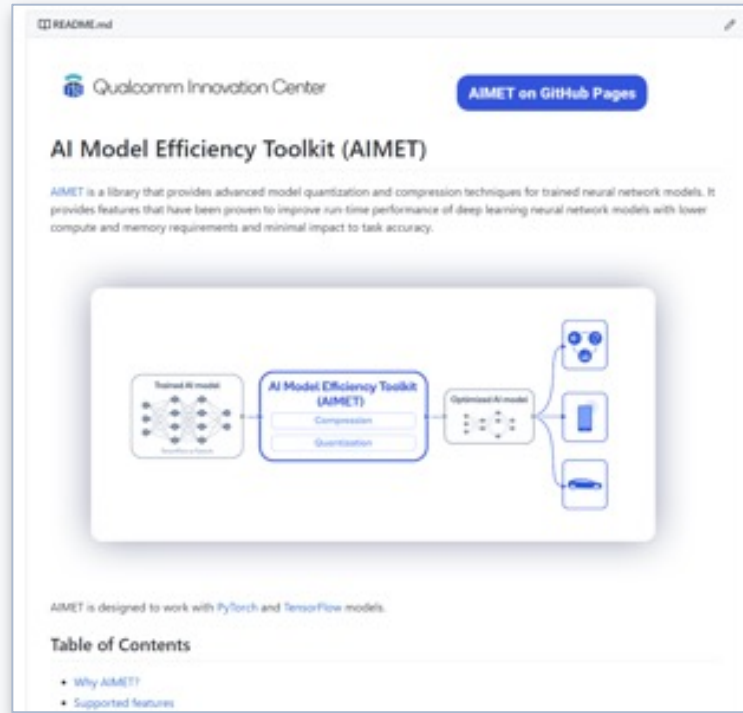
Tools are open-sourced through AIMET

github.com/quic/aimet

github.com/quic/aimet-model-zoo

AIMET

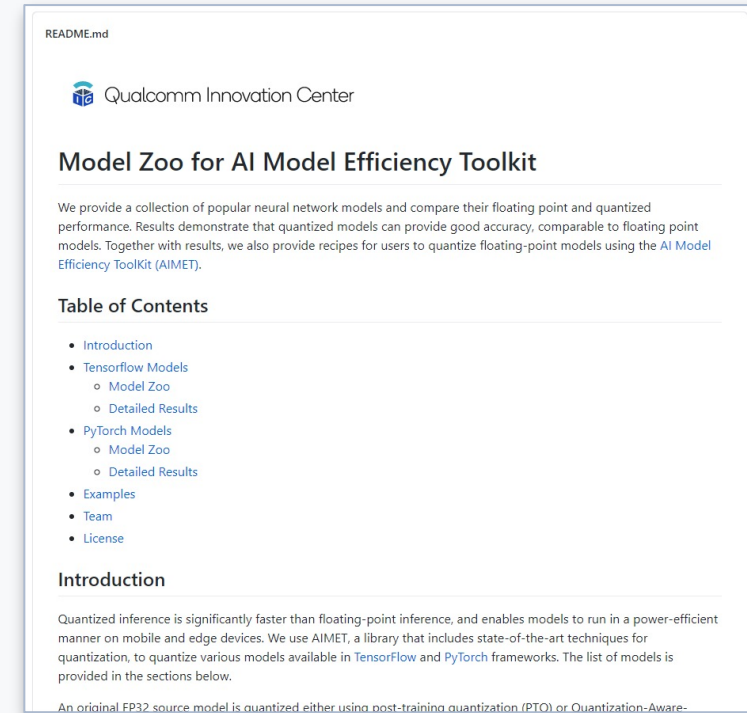
State-of-the-art quantization and compression techniques



github.com/quic/aimet

AIMET Model Zoo

Accurate pre-trained 8-bit quantized models



github.com/quic/aimet-model-zoo

Join our open-source projects

AIMET Model Zoo includes popular quantized AI models

Accuracy is maintained for INT8 models – less than 1% loss*

Tensorflow

<1%
Loss in
accuracy*

75.21% 74.96%
FP32 INT8

Top-1 accuracy*

ResNet-50
(v1)

75% 74.21%
FP32 INT8

Top-1 accuracy*

MobileNet-
v2-1.4

74.93% 74.99%
FP32 INT8

Top-1 accuracy*

EfficientNet
Lite

0.2469 0.2456
FP32 INT8

mAP*

SSD
MobileNet-v2

0.35 0.349
FP32 INT8

mAP*

RetinaNet

0.383 0.379
FP32 INT8

mAP*

Pose
estimation

25.45 24.78
FP32 INT8

PSNR*

SRGAN

Pytorch

71.67% 71.14%
FP32 INT8

Top-1 accuracy*

MobileNetV2

75.42% 74.44%
FP32 INT8

Top-1 accuracy*

EfficientNet-
lite0

72.62% 72.22%
FP32 INT8

mIoU*

DeepLabV3+

68.7% 68.6%
FP32 INT8

mAP*

MobileNetV2-
SSD-Lite

0.364 0.359
FP32 INT8

mAP*

Pose
estimation

25.51 25.5
FP32 INT8

PSNR

SRGAN

9.92% 10.22%
FP32 INT8

WER*

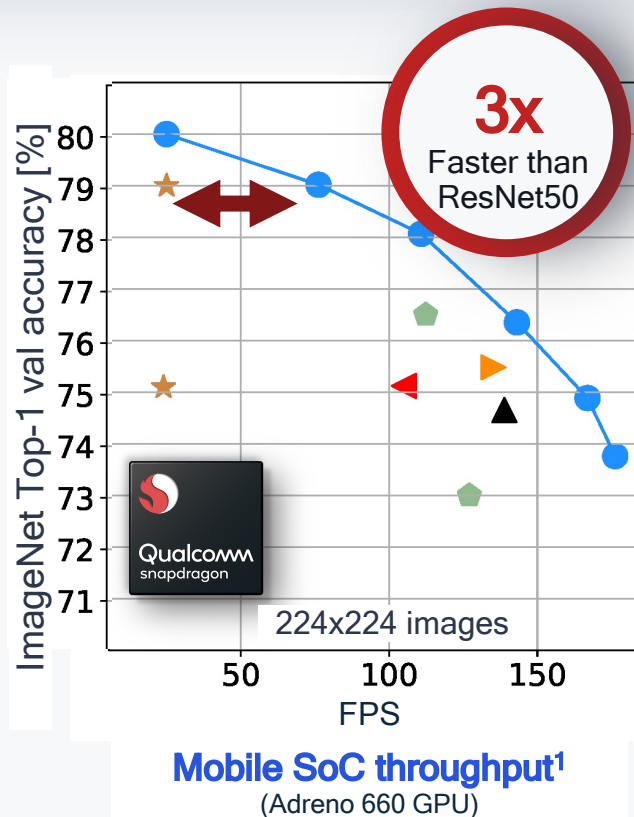
DeepSpeech2

*: Comparison between FP32 model and INT8 model quantized with AIMET.
For further details, check out: <https://github.com/quic/aimet-model-zoo/>

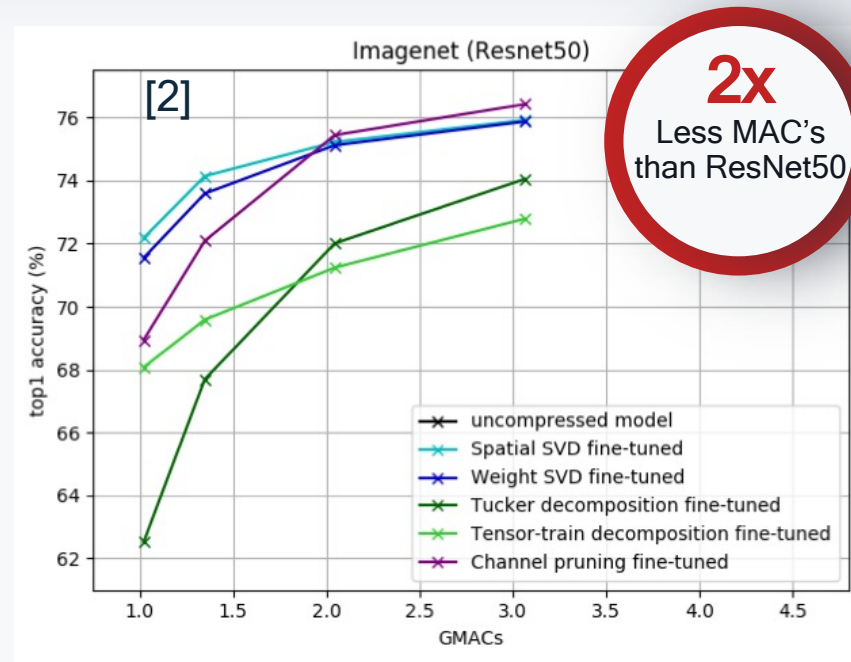
What's next in efficient on-device AI

Ultimately limited gains from NAS, compression, uniform quantization

Current tools optimize existing architectures, leading to 1-3x gains over standard networks on device



NAS



Compression

8-bit Integer

up to
16X

4-bit Integer

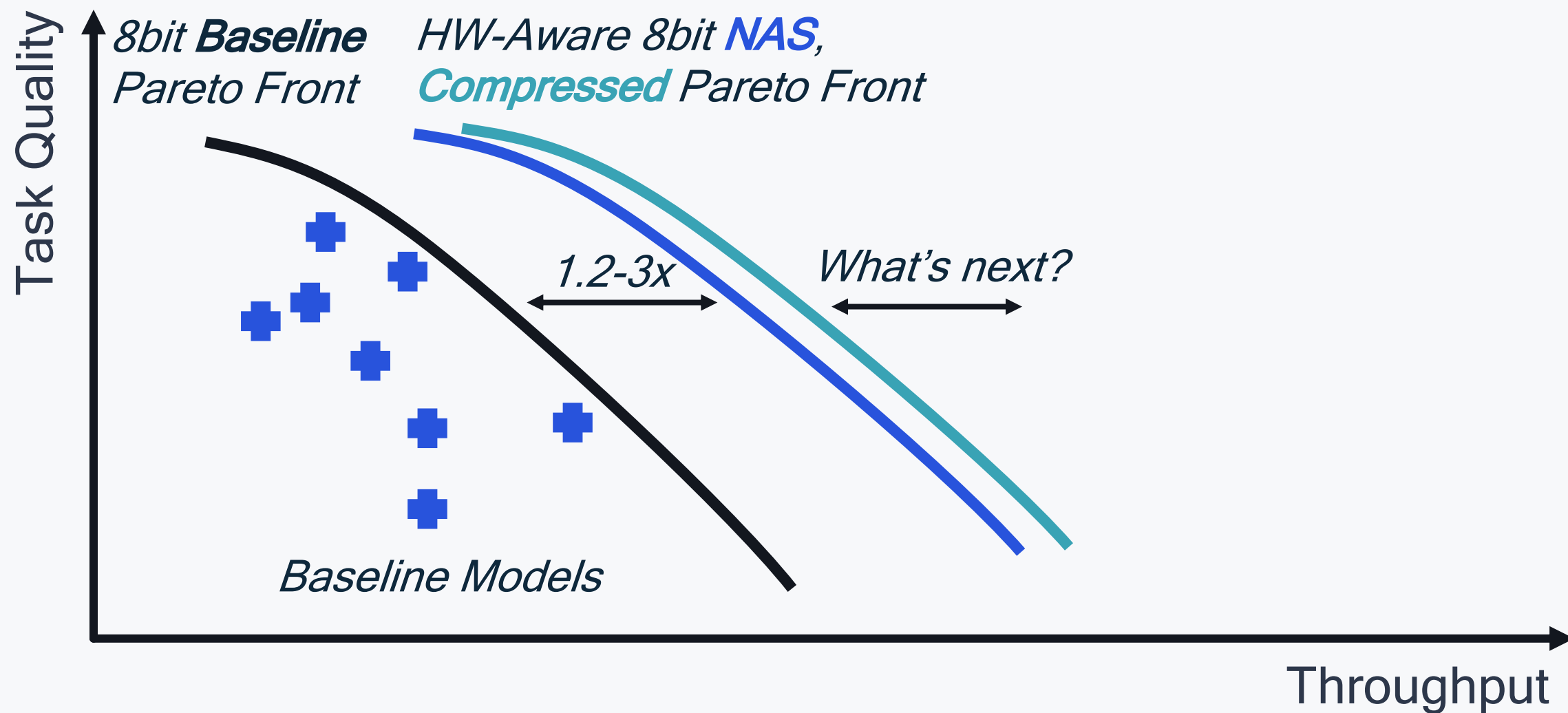
up to
64X

4-8b established

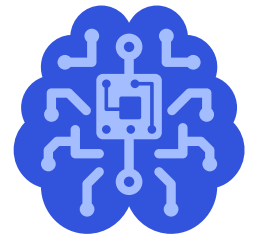
¹ Qualcomm Adreno 660 GPU in the Snapdragon 888 running on the Samsung Galaxy S21

[2] Andrey Kuzmin, et al, "Taxonomy and Evaluation of Structured Compression of Convolutional Neural Networks", Arxiv 2019

What's next in efficient AI models?

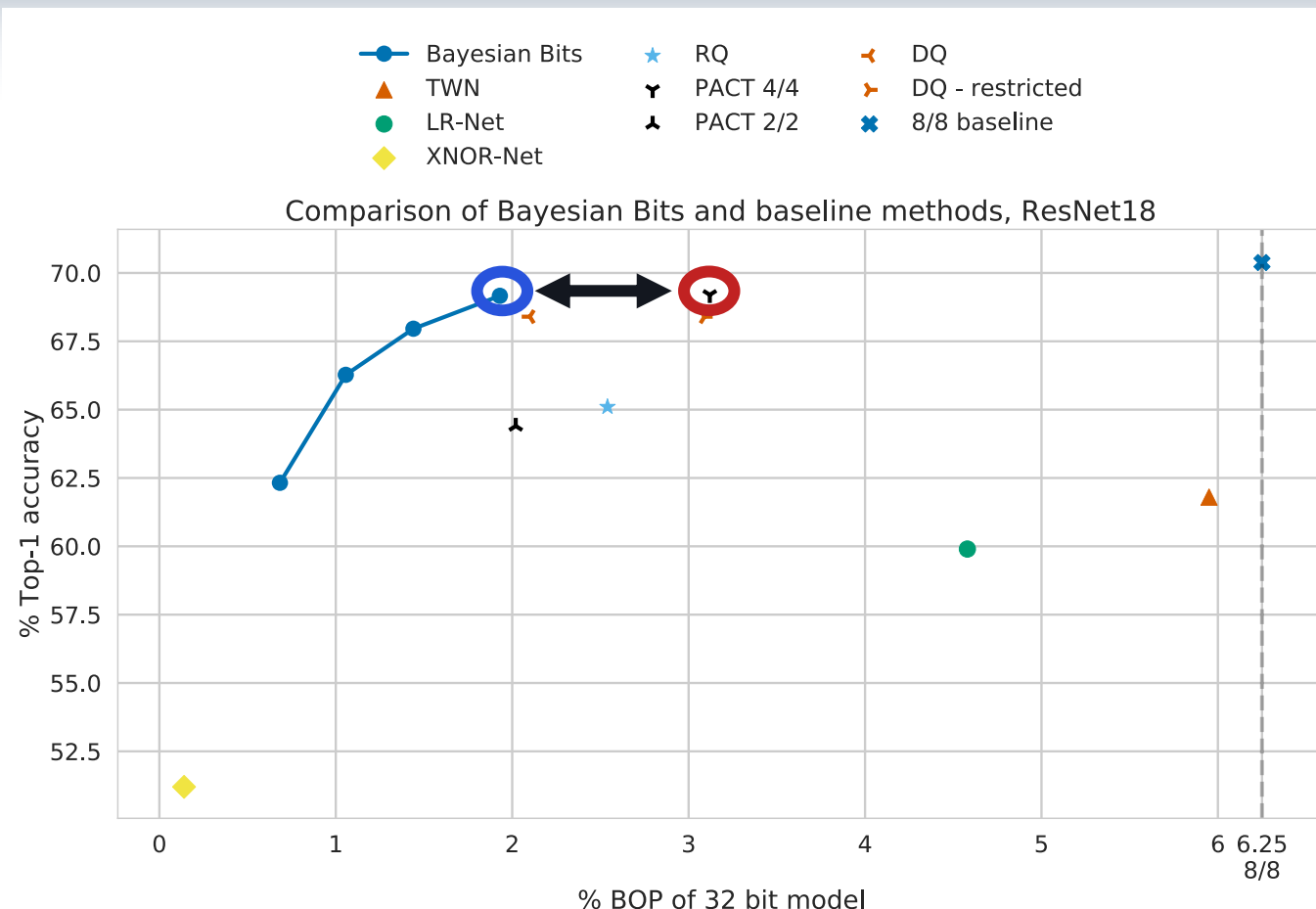


Mixed-Precision Quantized NAS



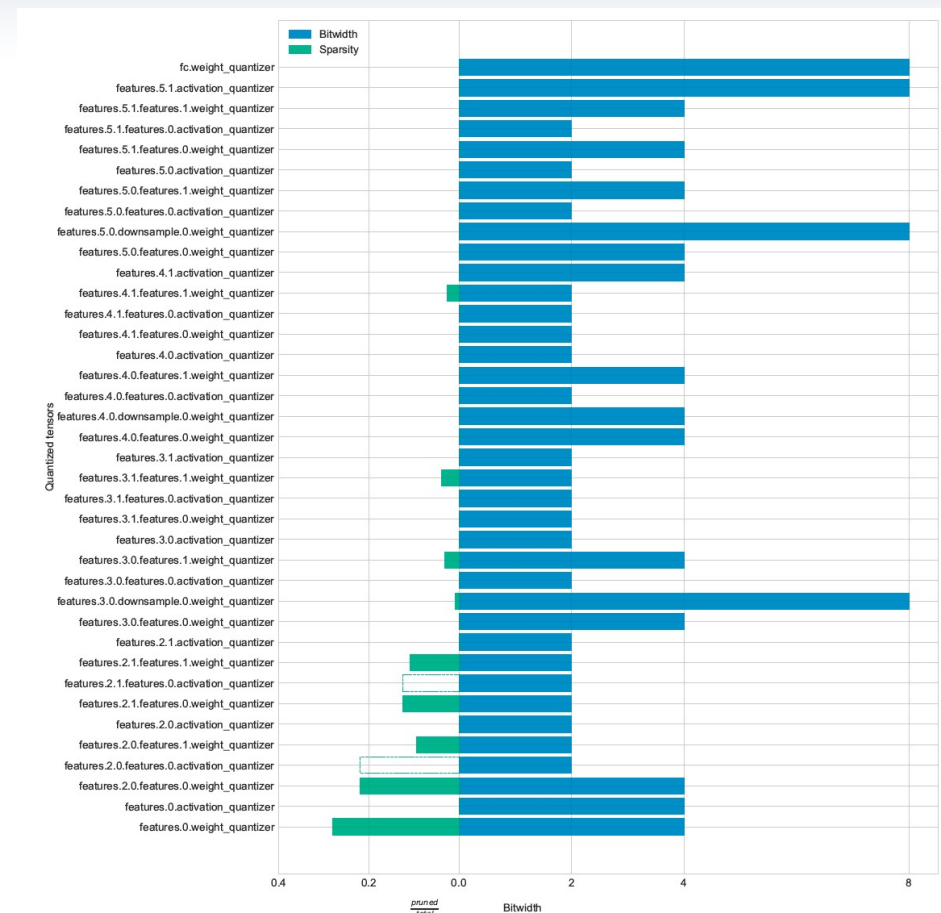
Mixed Precision outperforms uniform quantization

Bayesian Bits: Neural Networks can be optimized for mixed-precision.



During training, the network automatically finds the optimal trade-off between network complexity and accuracy

Mart Van Baalen, et al "Bayesian Bits: Unifying Quantization and Pruning", Neurips 2020

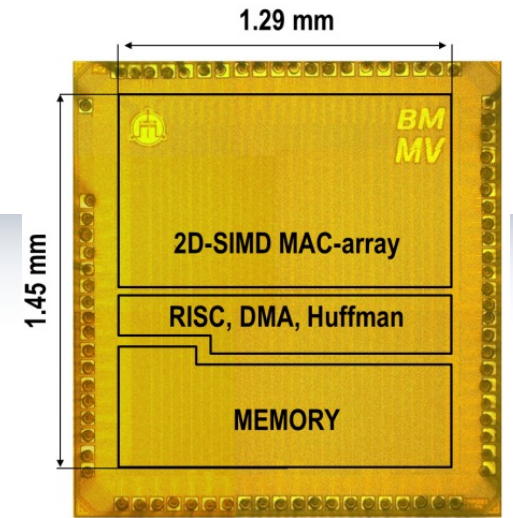


The result: Some layers are fine with 8 bits, while others are fine with 2 bits. And some layers are pruned (green)

Many Ways to gain from mixed-precision

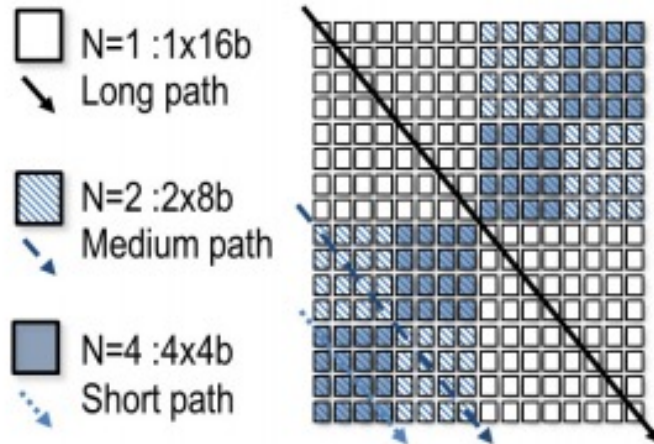
An academic system level example

- DVAFS: DVAS + subword parallelism op/J **10x** @ 2b vs 8b



Dynamic Voltage Accuracy Frequency Scaling exploits varying precision

A 16b DVAFS multiplier example



Modulate precision at constant throughput:

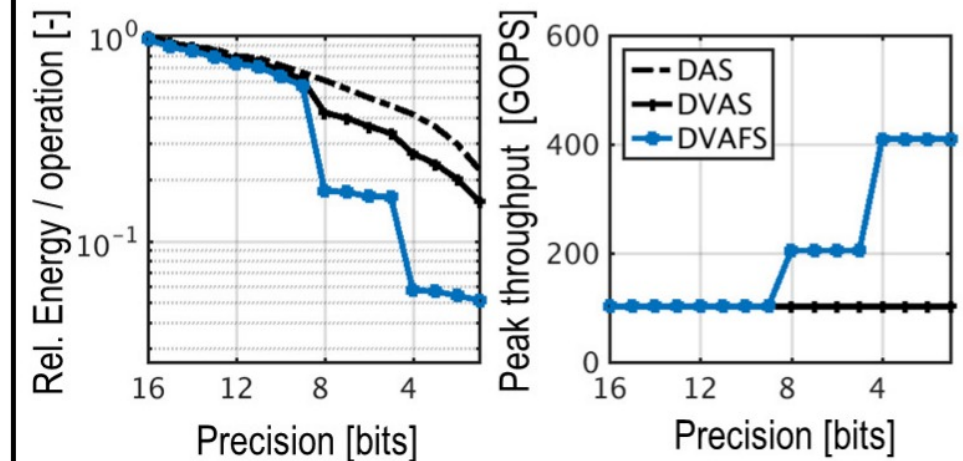
- Dynamic Accuracy Scaling

$$P_{DAS} = \frac{\alpha}{k_1} C f V^2$$
- Dynamic Voltage Accuracy Scaling

$$P_{DVAS} = \frac{\alpha}{k_1} C f \left(\frac{V}{k_2}\right)^2$$
- Dynamic Voltage Accuracy Freq. Scaling

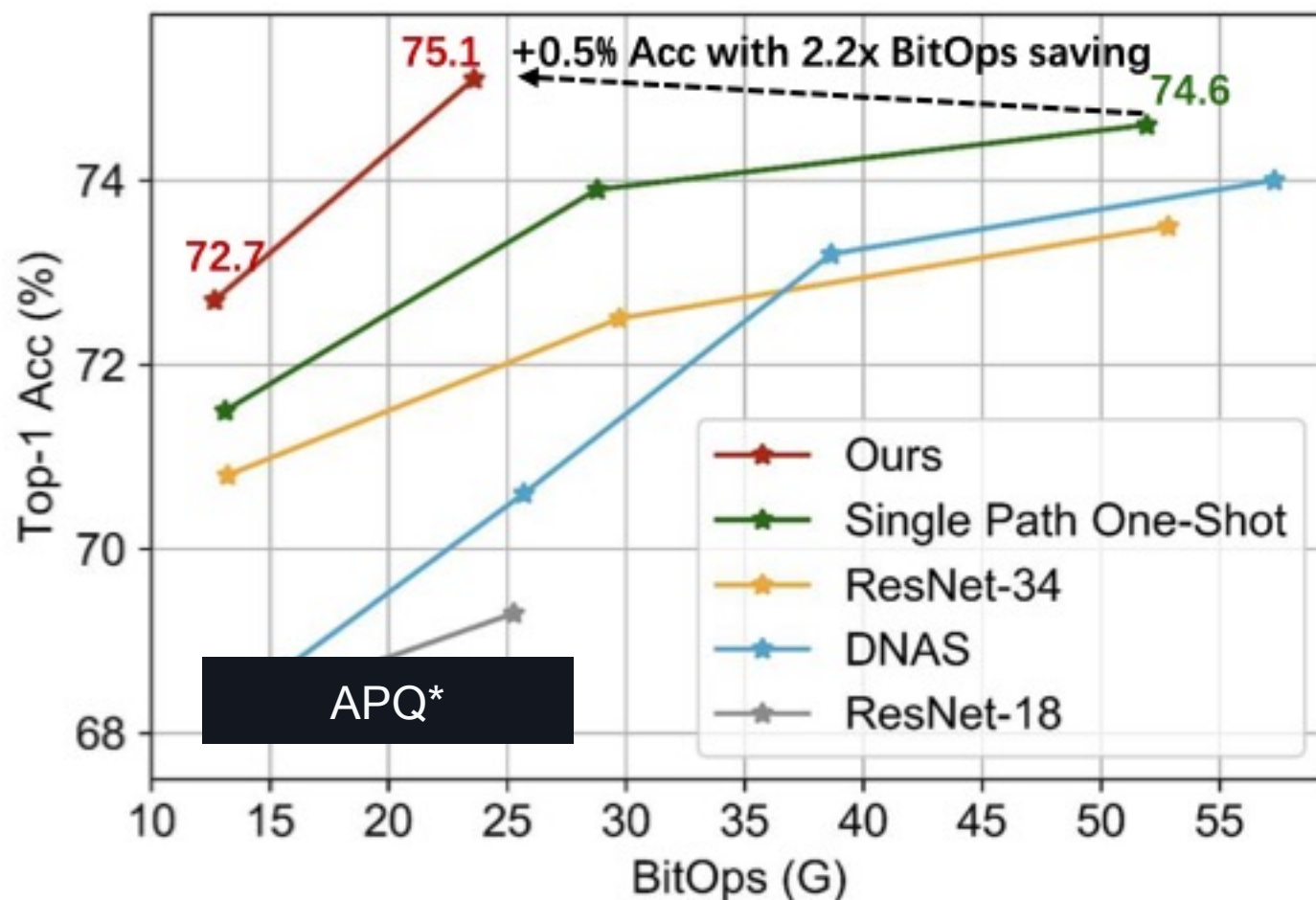
$$P_{DVAFS} = \frac{\alpha}{k_3} C \frac{f}{N} \left(\frac{V}{k_4}\right)^2$$

Measured comparison in this chip @ 200 MHz

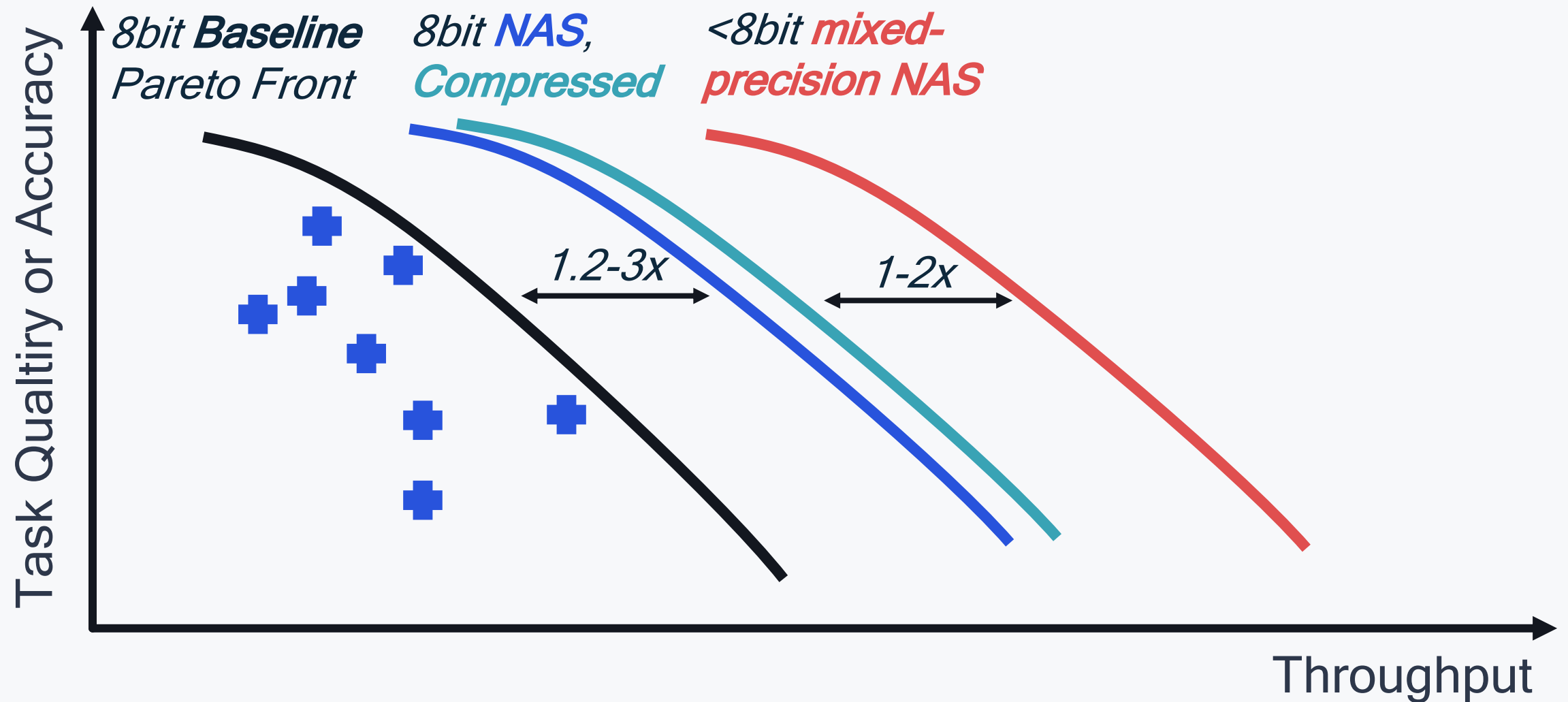


Mixed Precision Quantized NAS

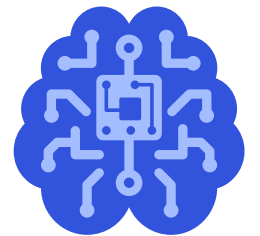
- APQ builds on top of OFA
- +/- 1%, or 2.2x BOPS gains expected through joint NAS and Quantization



What's next in efficient AI models?



Conditional networks



Conditional computing as a complementary technique to NAS

Input-dependent network architectures spend less time on easier samples

Classification: some samples are easier than others



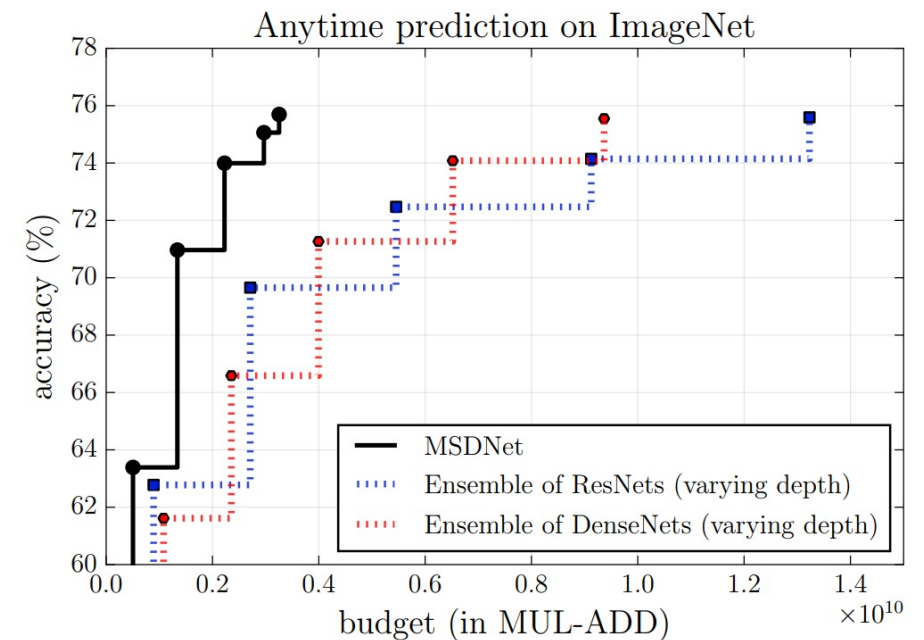
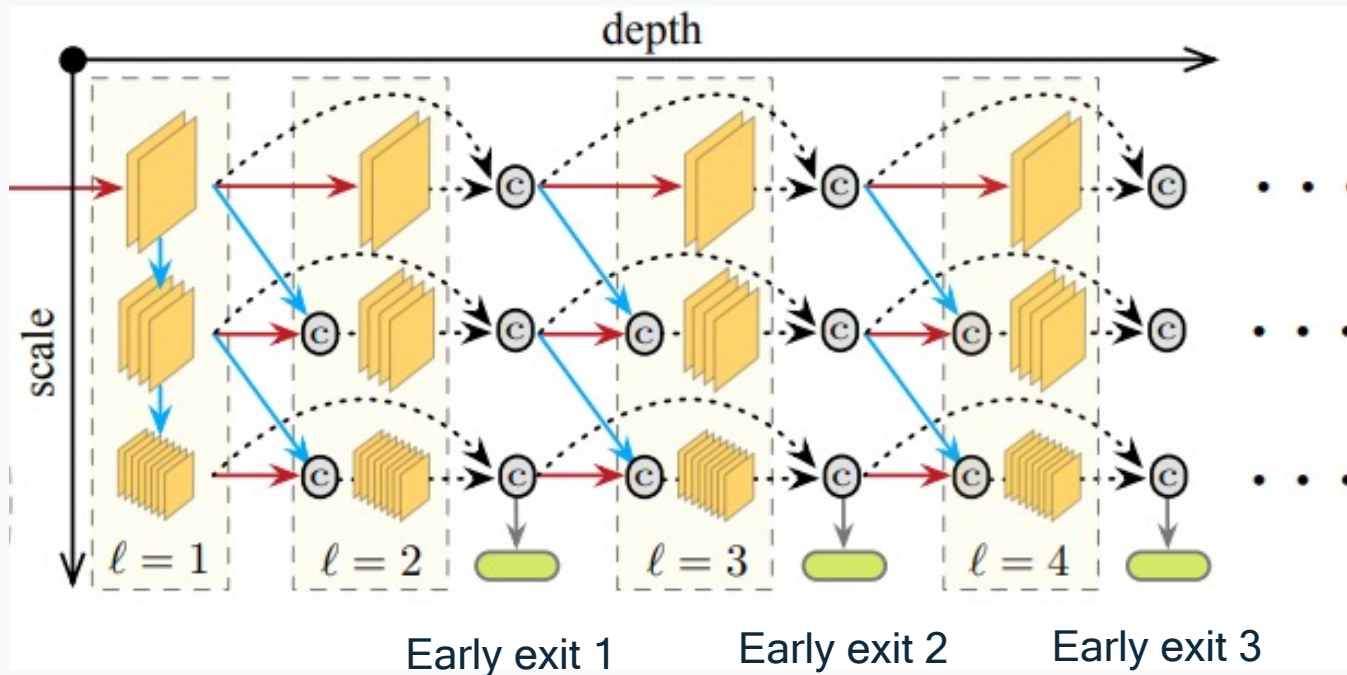
Copyright Pixel Addict and Doyle (CC BY-ND 2.0), no changes made

found through Huang, G, et al, "Multi-Scale Dense Networks for Resource Efficient Image Classification", ICLR2018

Conditional computing as a complementary technique to NAS

Input-dependent network architectures spend less time on easier samples

Early exiting, some samples are easier than others

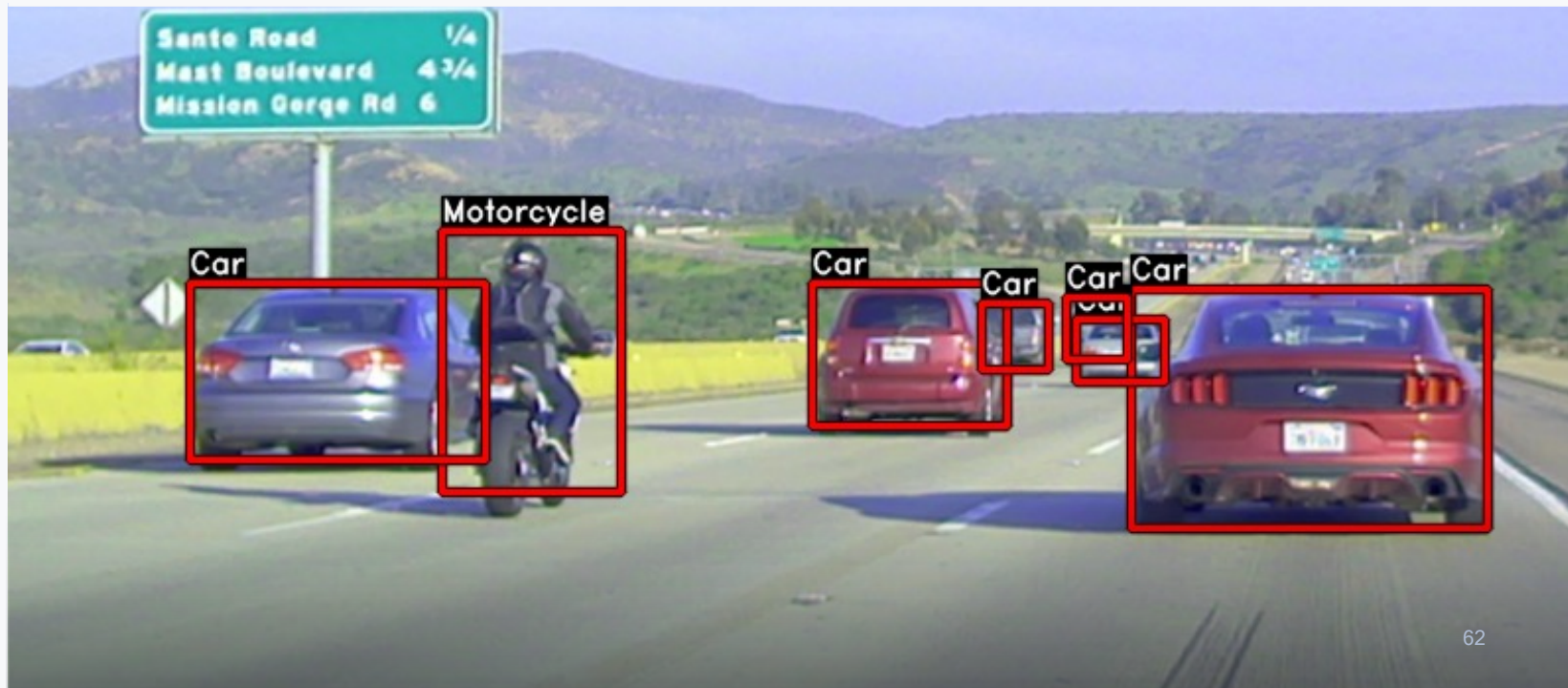


Conditional computing as a complementary technique to NAS

Input-dependent network architectures spend less time on easier samples

Segmentation: **backgrounds** are **abundant** and **easy** to recognize

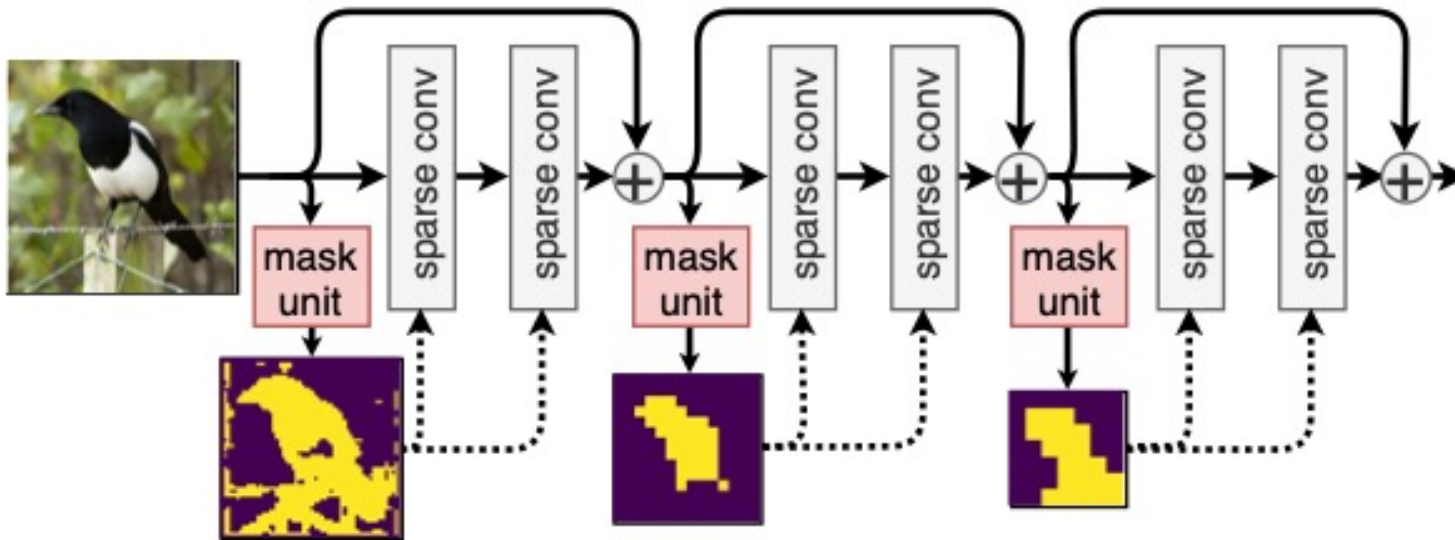
Detection: **Objects of interest** are relatively **rare**



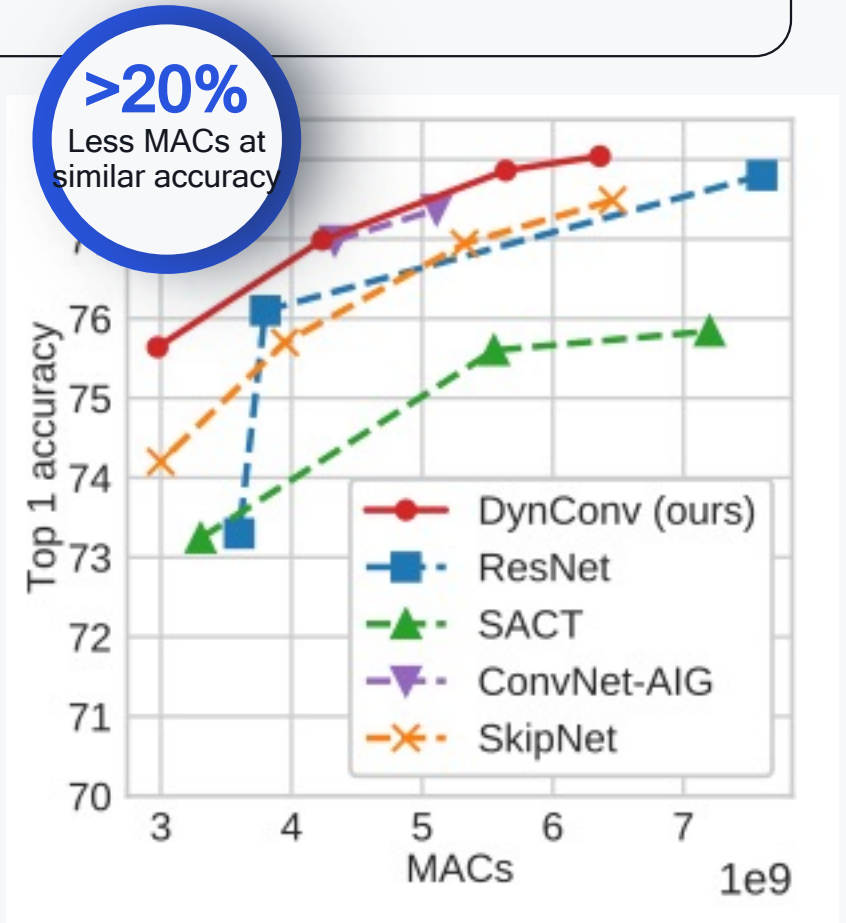
Conditional computing as a complementary technique to NAS

Input-dependent network architectures spend less time on easier samples

Dynamic Convolutions: exploiting spatial sparsity



Thomas Verelst, et al, "Dynamic Convolutions: Exploiting spatial sparsity for faster inference", CVPR20



Conditional computing as a complementary technique to NAS

Input-dependent network architectures spend less time on easier samples

Video: Subsequent frames are correlated

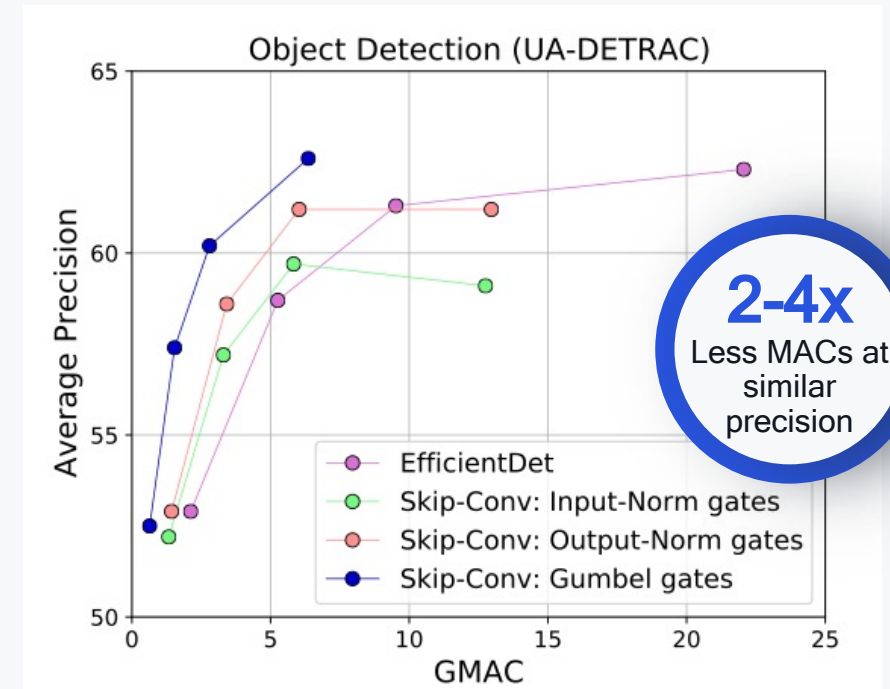
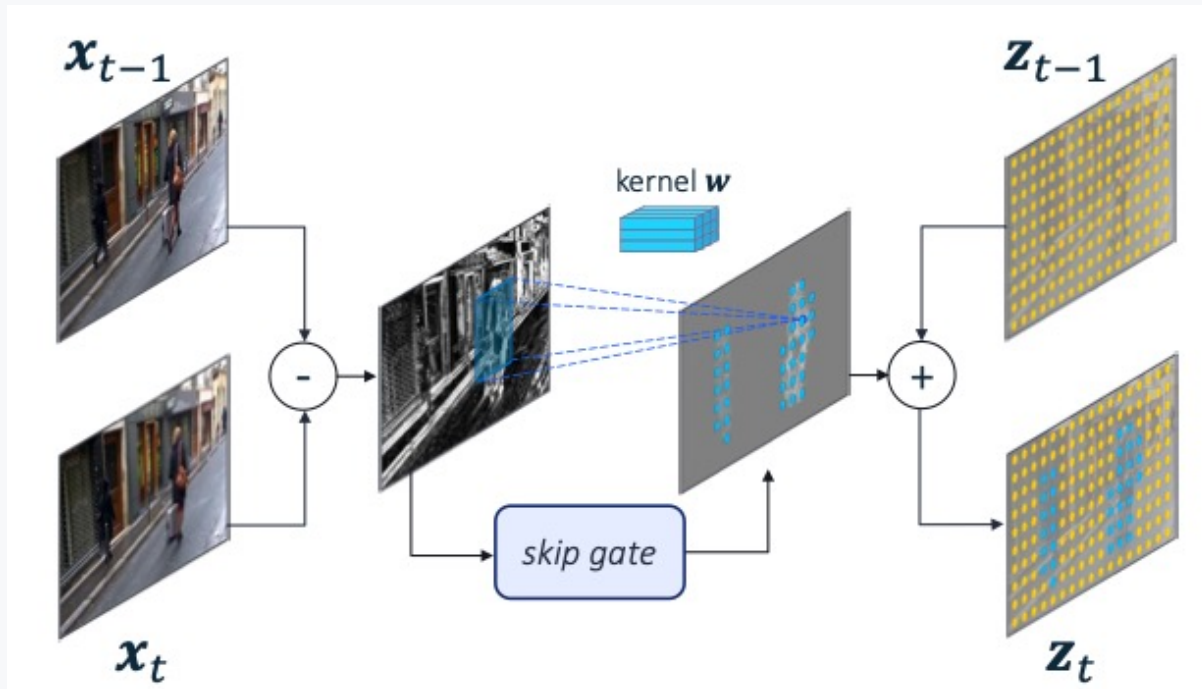


*Artur Andrzej, https://commons.wikimedia.org/wiki/File:Gdańsk_skrzyżowanie_ulic_Grunwaldzkiej_i_Słowackiego.jpg
(Creative Commons CC0 1.0 Universal Public Domain Dedication), no changes made*

Conditional computing as a complementary technique to NAS

Input-dependent network architectures spend less time on easier samples

Temporal Skip-Convolutions in video segmentation/detection



Amirhossein Habibian, et al, "Skip-Convolutions for Efficient Video Processing", CVPR21

Conditional computing as a complementary technique to NAS

Input-dependent network architectures spend less time on easier samples

Conditional Early exiting in video/action recognition

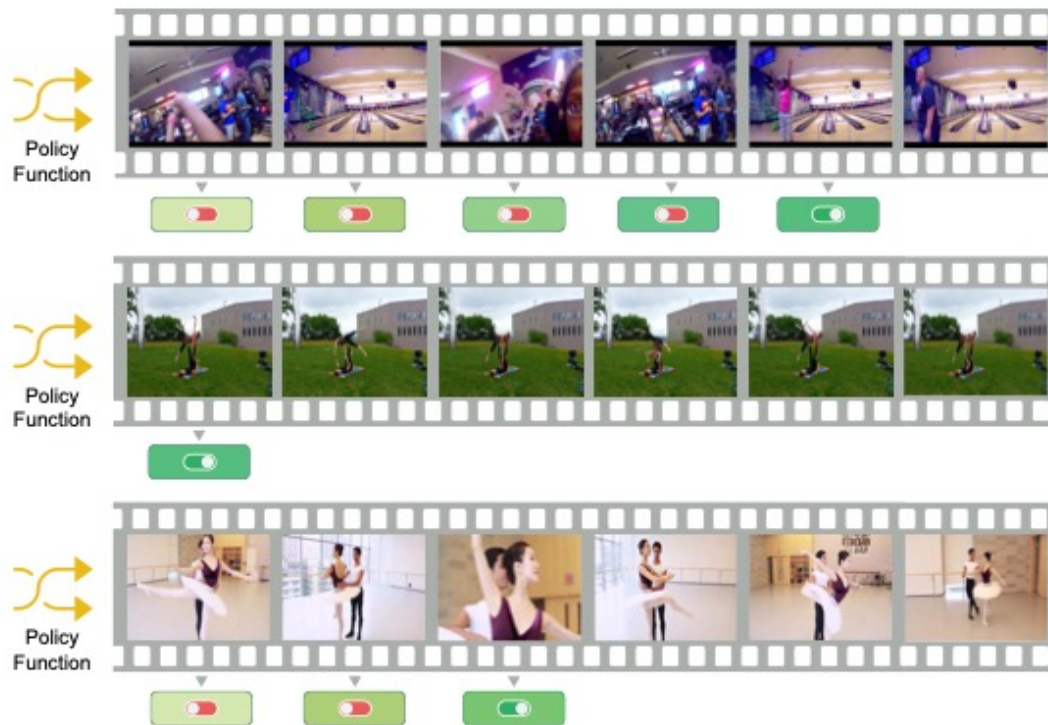


Figure 1: **Efficient video recognition by early exiting.**

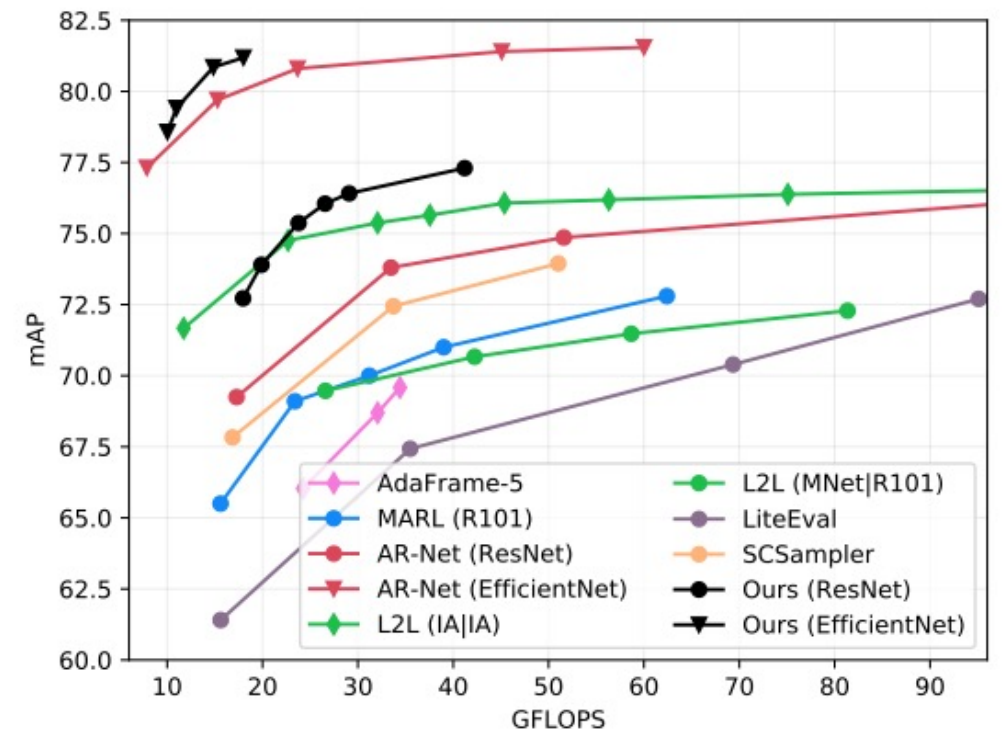
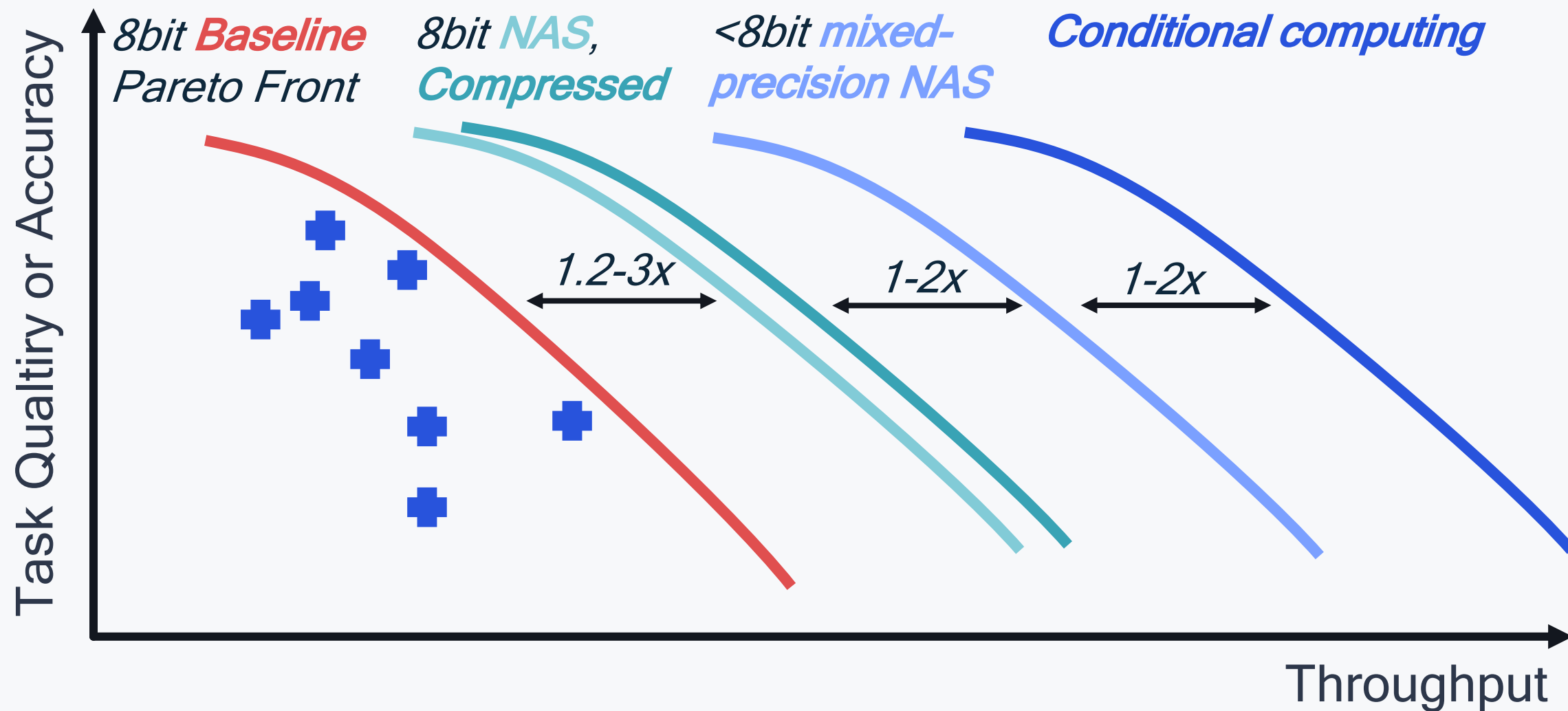
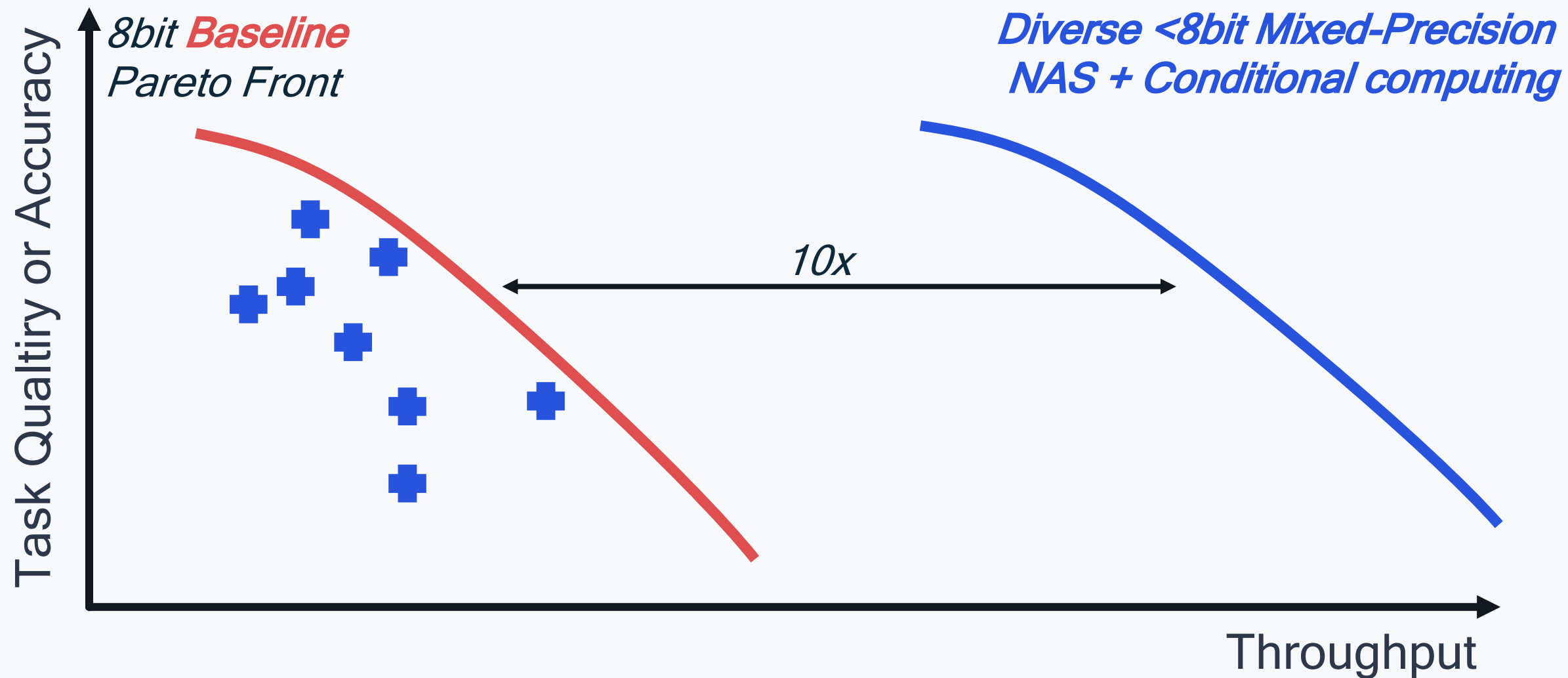


Figure 6: **Accuracy vs. efficiency curves on ActivityNet.**

What's next in efficient AI models?



What's next in efficient AI models?



Overview

- Energy-Efficient machine learning and the computational budget gap
- The Model-Efficiency Pipeline reduces the cost of on-device inference



Qualcomm Innovation Center, Inc. open sources through AI Model Efficiency Toolkit (AIMET)

- What's next in energy-efficient AI

Questions?

Connect with Us



www.qualcomm.com/ai



www.qualcomm.com/news/onq



[@QCOMResearch](https://twitter.com/QCOMResearch)







<https://www.youtube.com/qualcomm?>



<http://www.slideshare.net/qualcommwirelessevolution>



Thank you

Follow us on:    

For more information, visit us at:

www.qualcomm.com & www.qualcomm.com/blog

Nothing in these materials is an offer to sell any of the components or devices referenced herein.

©2018-2021 Qualcomm Technologies, Inc. and/or its affiliated companies. All Rights Reserved.

Qualcomm, Snapdragon, Adreno, and Hexagon are trademarks or registered trademark of Qualcomm Incorporated. Other products and brand names may be trademarks or registered trademarks of their respective owners.

References in this presentation to “Qualcomm” may mean Qualcomm Incorporated, Qualcomm Technologies, Inc., and/or other subsidiaries or business units within the Qualcomm corporate structure, as applicable. Qualcomm Incorporated includes our licensing business, QTL, and the vast majority of our patent portfolio. Qualcomm Technologies, Inc., a subsidiary of Qualcomm Incorporated, operates, along with its subsidiaries, substantially all of our engineering, research and development functions, and substantially all of our products and services businesses, including our QCT semiconductor business.

Premier Sponsor



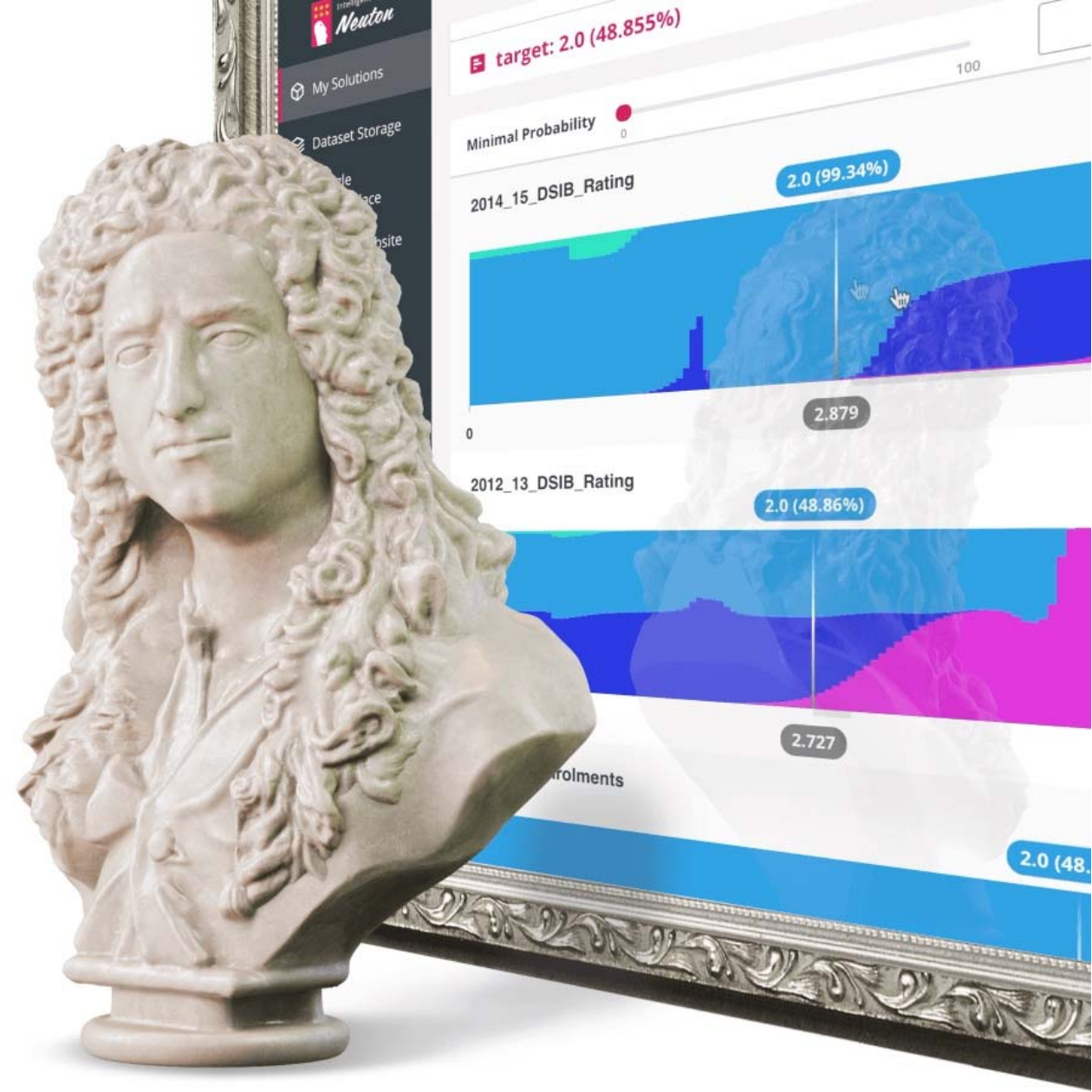
Automated TinyML

Zero-code SaaS solution

**Create tiny models, ready for embedding,
in just a few clicks!**

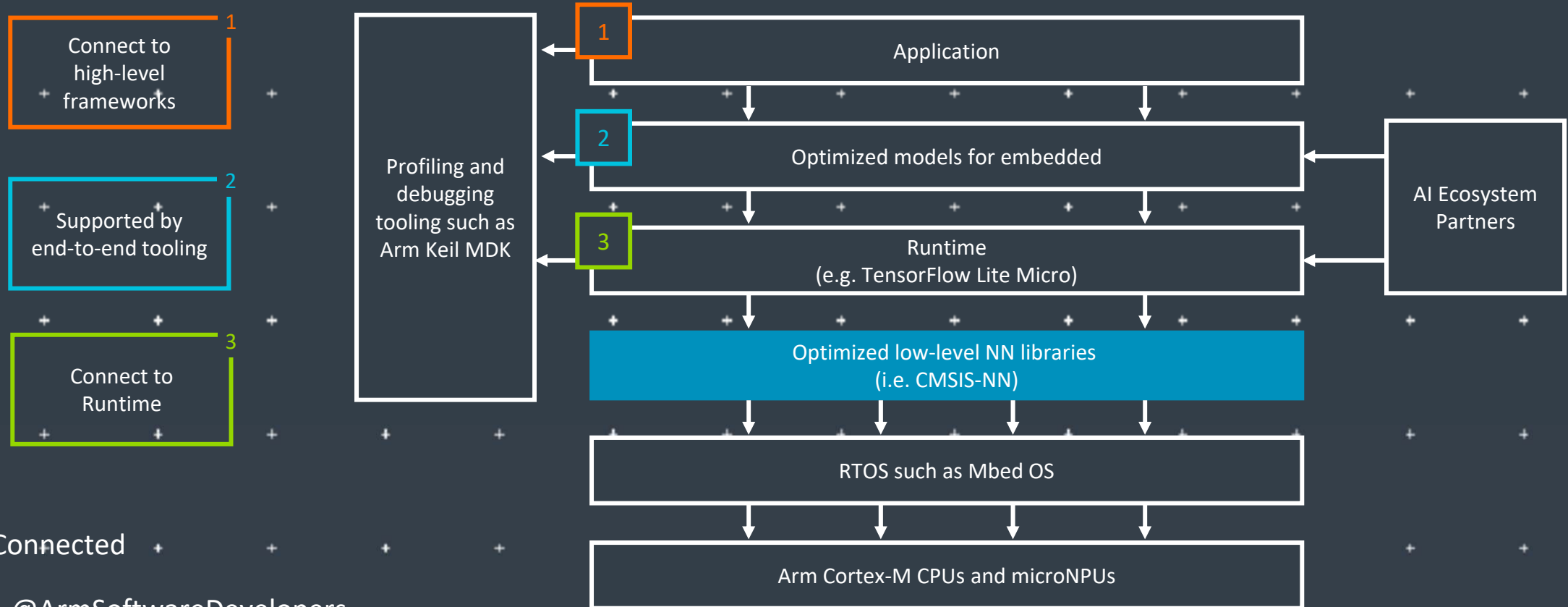
Compare the benchmarks of our compact models to those of TensorFlow and other leading neural network frameworks.

Build Fast. Build Once. Never Compromise.



Executive Sponsors

Arm: The Software and Hardware Foundation for tinyML



Stay Connected



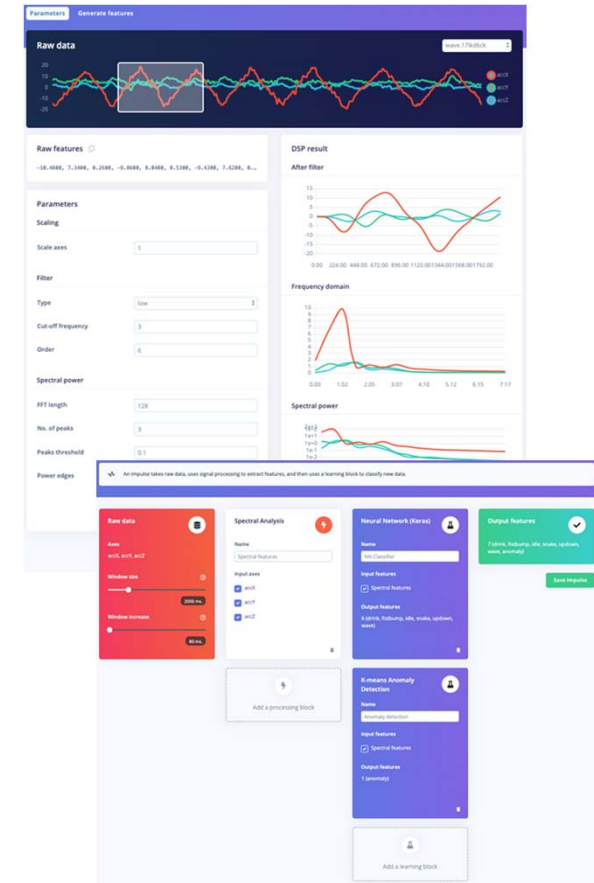
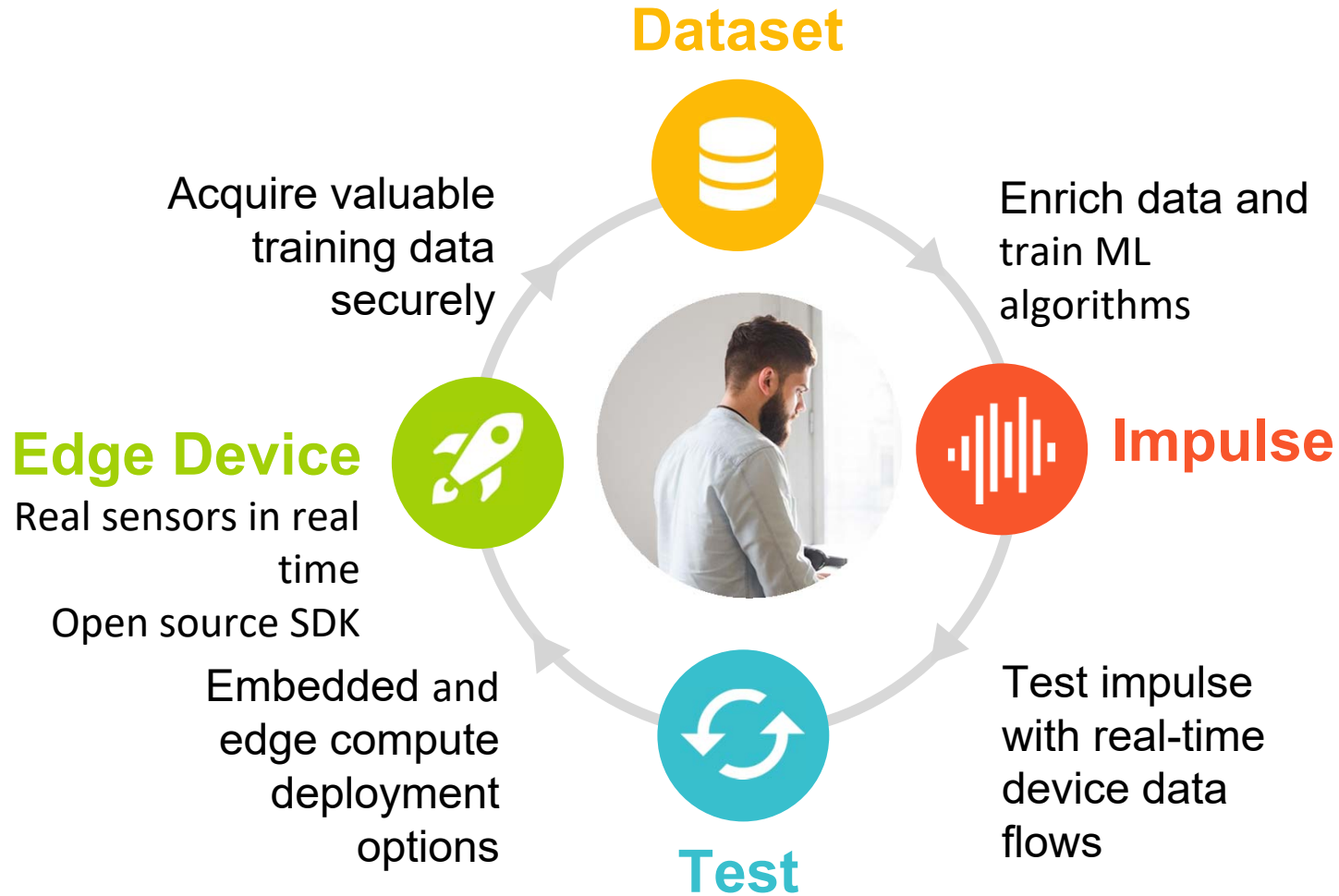
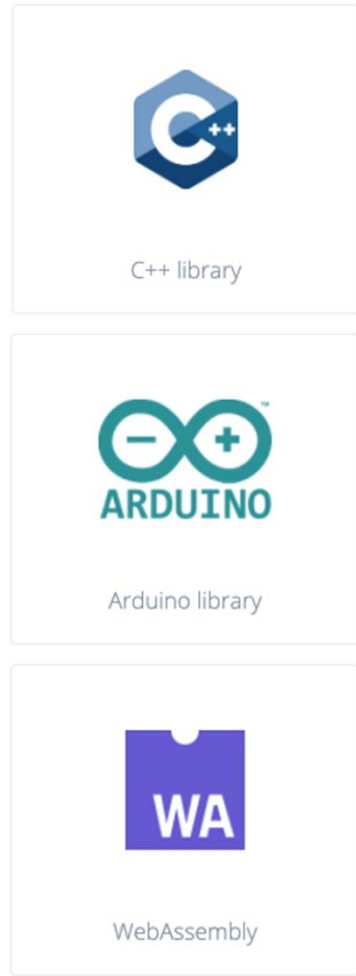
@ArmSoftwareDevelopers



@ArmSoftwareDev

Resources: developer.arm.com/solutions/machine-learning-on-arm

TinyML for all developers



www.edgeimpulse.com

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



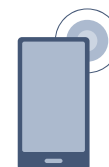
Cloud



IoT/IIoT



Automotive



Mobile

SYNTIANT

[Syntiant Corp.](#) is moving artificial intelligence and machine learning from the cloud to edge devices. Syntiant's chip solutions merge deep learning with semiconductor design to produce ultra-low-power, high performance, deep neural network processors. These network processors enable always-on applications in battery-powered devices, such as smartphones, smart speakers, earbuds, hearing aids, and laptops. Syntiant's Neural Decision Processors™ offer wake word, command word, and event detection in a chip for always-on voice and sensor applications.

Founded in 2017 and headquartered in Irvine, California, the company is backed by Amazon, Applied Materials, Atlantic Bridge Capital, Bosch, Intel Capital, Microsoft, Motorola, and others. Syntiant was recently named a [CES® 2021 Best of Innovation Awards Honoree](#), [shipped over 10M units worldwide](#), and [unveiled the NDP120](#) part of the NDP10x family of inference engines for low-power applications.

www.syntiant.com



@Syntiantcorp

Platinum Sponsors



Part of your life. Part of tomorrow.

www.infineon.com



Reality AI[®]

Add Advanced Sensing to your Product with Edge AI / TinyML

<https://reality.ai>



info@reality.ai



[@SensorAI](https://twitter.com/SensorAI)



[Reality AI](#)

Pre-built Edge AI sensing modules, plus tools to build your own

Reality AI solutions

Prebuilt sound recognition models for
indoor and outdoor use cases

Solution for industrial anomaly detection

Pre-built automotive solution that lets cars
“see with sound”

Reality AI Tools[®] software

Build prototypes, then turn them into
real products

Explain ML models and relate the function
to the physics

Optimize the hardware, including
sensor selection and placement

Gold Sponsors



LatentAI

Adaptive AI for the Intelligent Edge

[Latentai.com](https://latent.ai)



Build Smart IoT Sensor Devices From Data

SensiML pioneered TinyML software tools that auto generate AI code for the intelligent edge.

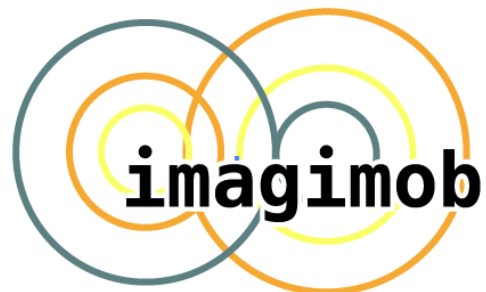
- End-to-end AI workflow
- Multi-user auto-labeling of time-series data
- Code transparency and customization at each step in the pipeline

We enable the creation of production-grade smart sensor devices.



sensiml.com

Silver Sponsors



Copyright Notice

The presentation(s) in this publication comprise the proceedings of tinyML® EMEA Technical Forum 2021. The content reflects the opinion of the authors and their respective companies. This version of the presentation may differ from the version that was presented at tinyML EMEA. 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