Qualcomm
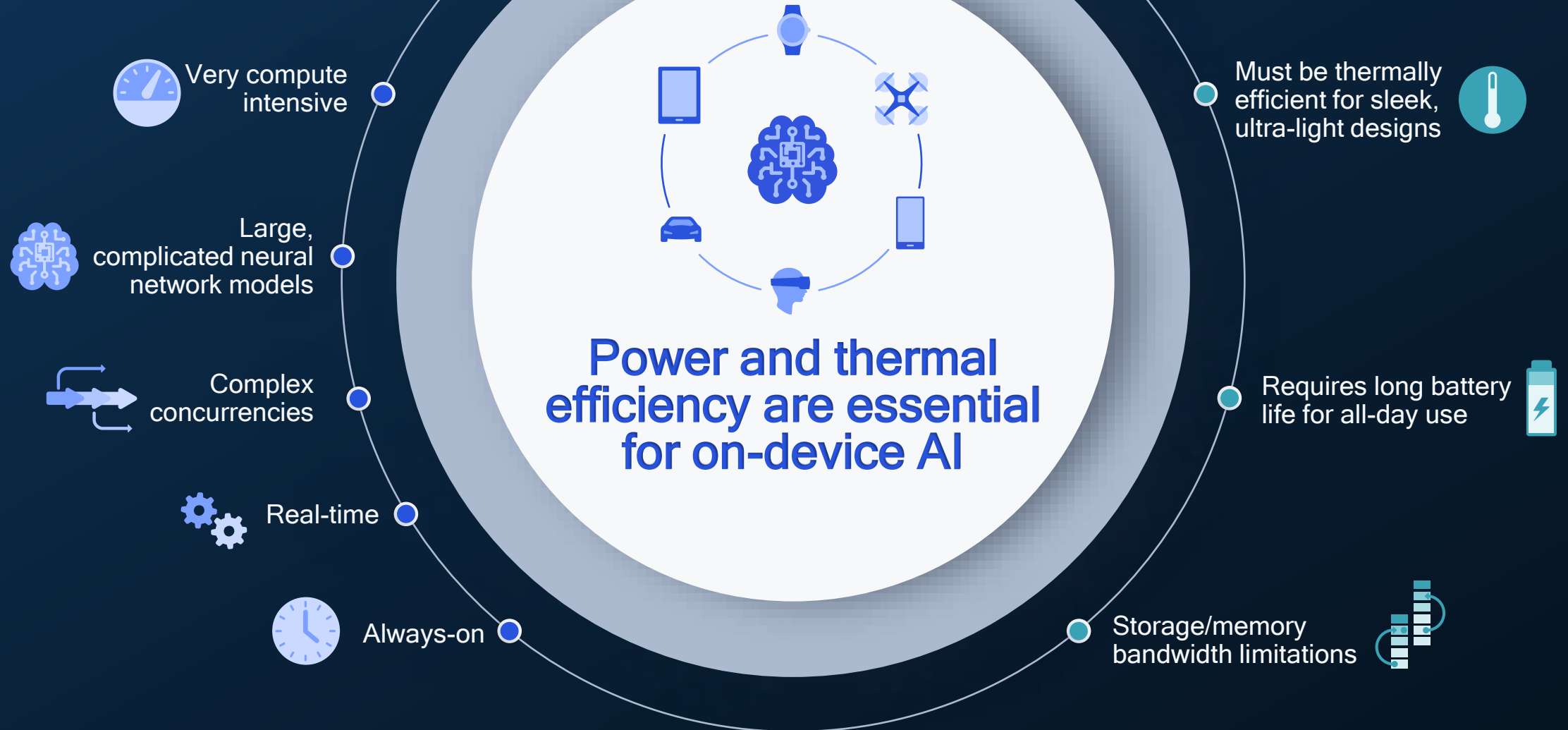
Qualcomm
AI research

# AI Model Efficiency Toolkit (AIMET)

Chirag Patel
Principal Eng./Mgr.
Qualcomm AI Research

# The challenge of AI workloads

## Constrained mobile environment

Very compute intensive

Large, complicated neural network models

Complex concurrencies

Real-time

Always-on

**Power and thermal efficiency are essential for on-device AI**

Must be thermally efficient for sleek, ultra-light designs

Requires long battery life for all-day use

Storage/memory bandwidth limitations

# Leading machine learning research for edge AI

across the entire spectrum of topics

Deep generative models

Causality and system-2

G-CNN

Reinforcement learning

Self-supervised learning

## Fundamental research

Bayesian distributed learning

AI for wireless and RF sensing

Graph and kernel optimization

Deep learning for 3D/geometry

Federated learning

Audio and video compression

Qualcomm
AI research

Compute-in-memory

Video recognition and prediction

Model quantization, compression, and NAS

Fingerprint

AI for chip design

## Platform research

## Applied research

Energy-efficient perception

HW-SW co-design

Deep learning for graphics

On-device learning

Power management

Visual quality improvement

AI Model Efficiency Toolkit (AIMET)

Voice UI

## Model quantization

Invented the best techniques for fast deployment of 8-bit quantization

Best power-efficiency toolkit in the industry

## On-device learning

Invented continuous learning techniques for SOTA on-device voice-UI

First demonstration of 30% improvement to keyword spotting

## Federated learning

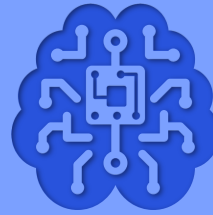Invented methods for combining differential privacy and compression

First end-to-end research software framework deployable on mobile

## Video semantic segmentation

Top the Cityscape leaderboard with loss function innovation for boundary-awareness

First real-time SS at FHD on mobile

# AI Firsts

## Brought to you by Qualcomm AI Research

## Group equivariant CNN

Pioneer for rotational equivariance; best paper at ICLR'18

First G-CNN segmentation for health on mobile

## AI for wireless

Invented neural augmentation to enhance physical layer algorithms

First weakly supervised method for real-world passive RF sensing

## Video super resolution

Full stack optimization for visual quality improvement at 4K resolution

First 4K SR at 100+ FPS on mobile

## Neural video compression

Invented instance-adaptive for SOTA performance & new deployment scenarios

First real-time HD decoding on mobile

# Holistic model efficiency research

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

## Quantization
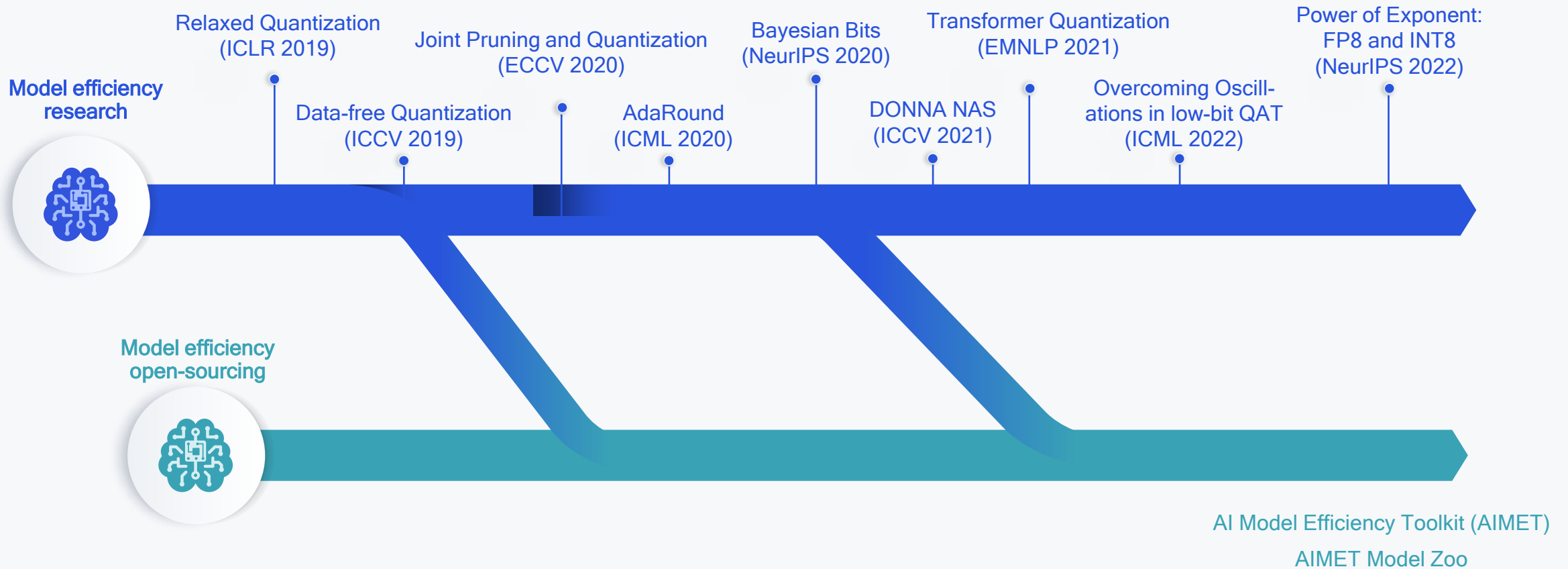Learning to reduce bit-precision while keeping desired accuracy

## Compilation
Learning to compile AI models for efficient hardware execution

## Conditional compute
Learning to execute only parts of a large inference model based on the input

## Neural architecture search
Learning to design smaller neural networks that are on par or outperform hand-designed architectures on real hardware

**Model efficiency research**

- Relaxed Quantization (ICLR 2019)
- Data-free Quantization (ICCV 2019)
- Joint Pruning and Quantization (ECCV 2020)
- AdaRound (ICML 2020)
- Bayesian Bits (NeurIPS 2020)
- DONNA NAS (ICCV 2021)
- Transformer Quantization (EMNLP 2021)
- Overcoming Oscillations in low-bit QAT (ICML 2022)
- Power of Exponent: FP8 and INT8 (NeurIPS 2022)

**Model efficiency open-sourcing**

AI Model Efficiency Toolkit (AIMET)

AIMET Model Zoo

# Driving the industry toward integer inference and power-efficient AI

## Leading model efficiency research and fast commercialization

# AIMET makes AI models small

State-of-the-art quantization and compression techniques from Qualcomm AI Research

**Trained**
AI model

**AIMET**

**Optimized**
AI model

**Deployed**
AI model

Compression

Quantization

TensorFlow or PyTorch

Github: https://github.com/quic/aimet

**Features:**

State-of-the-art network compression tools

State-of-the-art quantization tools

Support for both TensorFlow and PyTorch

Benchmarks and tests for many models

Developed by professional software developers

# AIMET

Providing advanced model efficiency features and benefits

## Benefits

- Lower power
- Lower memory bandwidth
- Maintains model accuracy
- Lower storage
- Higher performance
- Simple ease of use

## Features

### Quantization

State-of-the-art INT8 and INT4 performance

Quantization simulation

Quantization-aware training (QAT)

Post-training quantization (PTQ) methods:
- Data-Free Quantization
- Adaptive Rounding (AdaRound),
- Automatic Mixed Precision (AMP)
- AutoQuant

### Compression

Efficient tensor decomposition and removal of redundant channels in convolution layers

Spatial singular value decomposition (SVD)

Channel pruning

### Visualization

Analysis tools for drawing insights for quantization and compression

Weight ranges

Per-layer compression sensitivity

# AdaRound: Adaptive Rounding for Better Quantization
## ICML'20 paper

Rounding-to-the-nearest is not optimal



Round to nearest

-127          128

?

### Object Detection

| Configuration | mAP |
|---|---|
| Floating point | 82.20 |
| Nearest Rounding – 8-bit weights, 8-bit activations | 49.85 |
| AdaRound – 8-bit weights, 8-bit activations | 81.21 |

mAP: Mean Average Precision

### Semantic Seg. (Deeplabv3)

| Configuration | mIOU |
|---|---|
| Floating point | 72.94 |
| Nearest Rounding – 4-bit weights, 8-bit activations | 6.09 |
| AdaRound – 4-bit weights, 8-bit activations | 70.86 |

mIOU: Mean Intersection Over Union

AdaRound optimizes the network weights without model fine-tuning

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

# AdaRound Results

- Poor baseline INT8 quantization performance

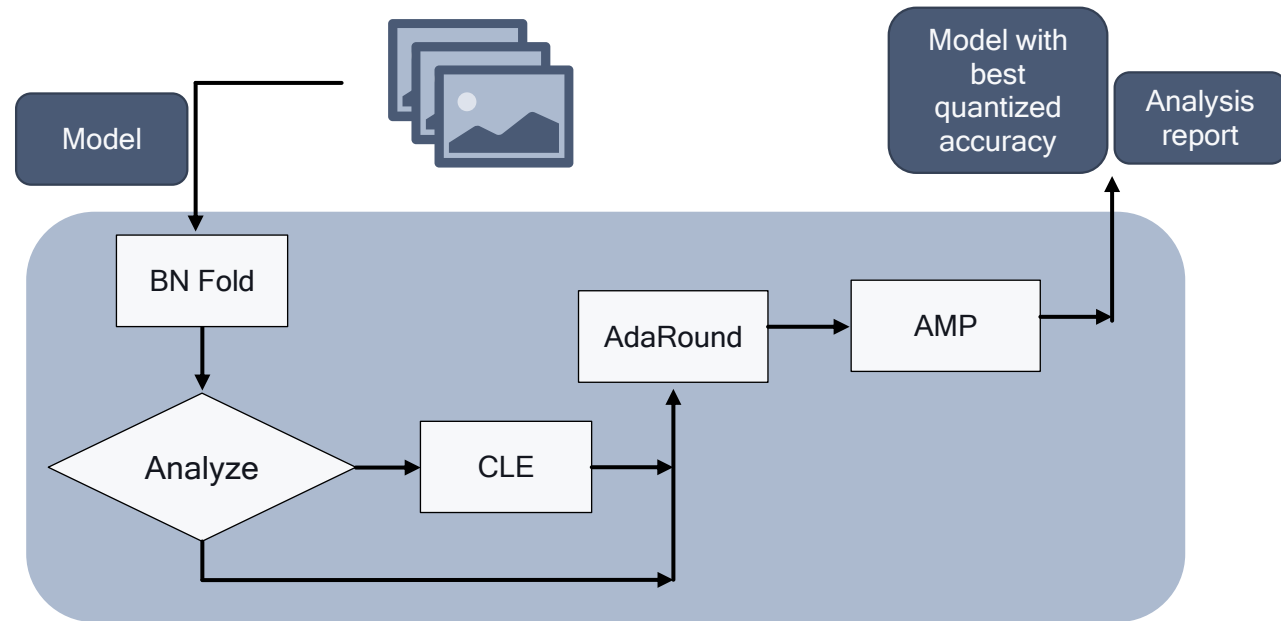- AdaRound performance within 1% of FP32

INT8, Baseline

INT8, AdaRound

# AutoQuant simplifies post-training quantization

- Analyzes the model

- Applies the best sequence of already existing post-training quantization (PTQ) features

- Returns the best accuracy model with analysis report
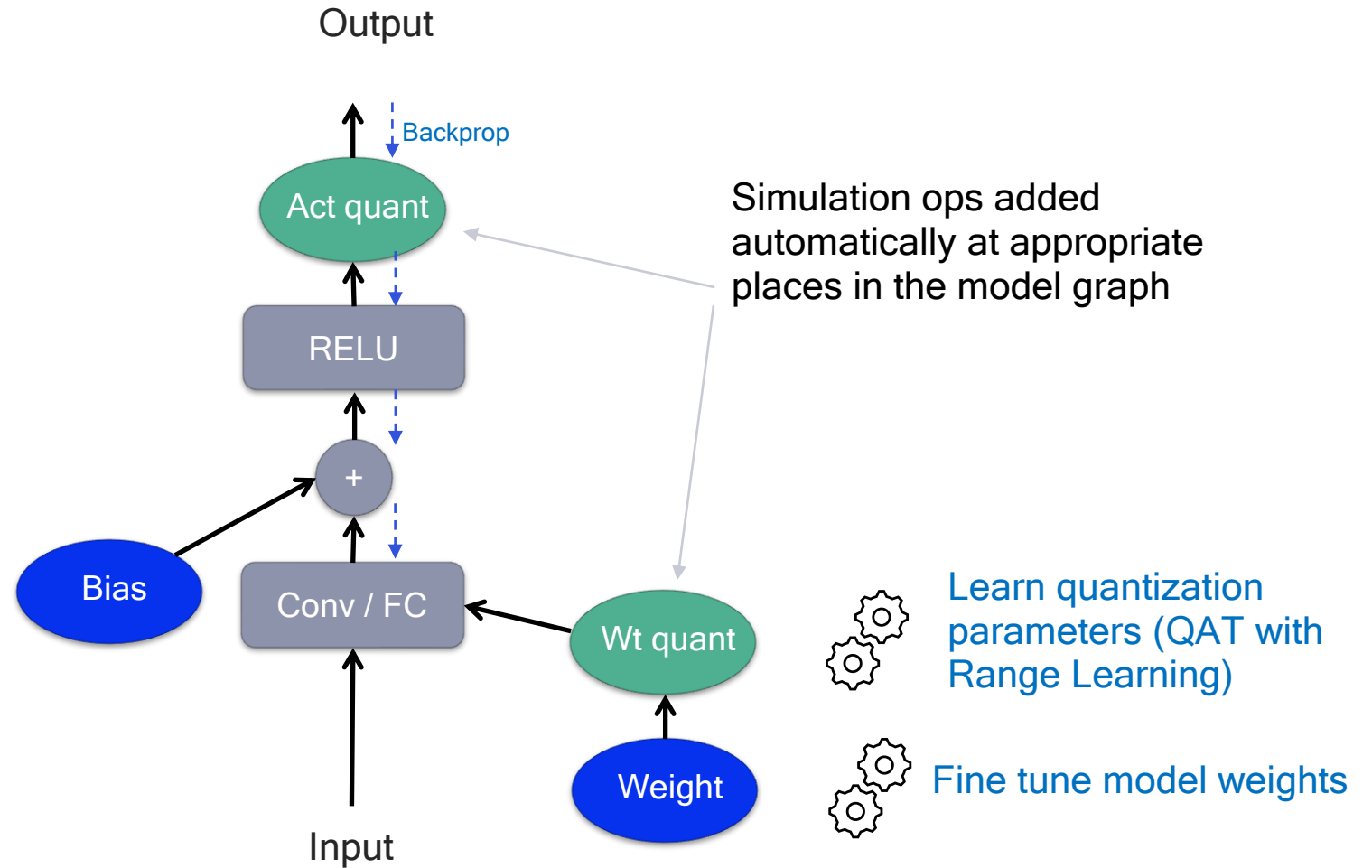
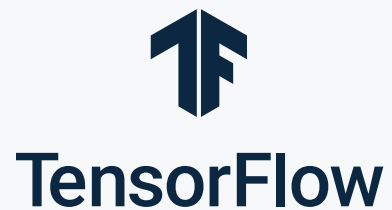- A simple, blackbox, push-button solution

# AIMET Quantization Aware Training

Simulate quantization noise in the forward pass and fine-tune for improved robustness
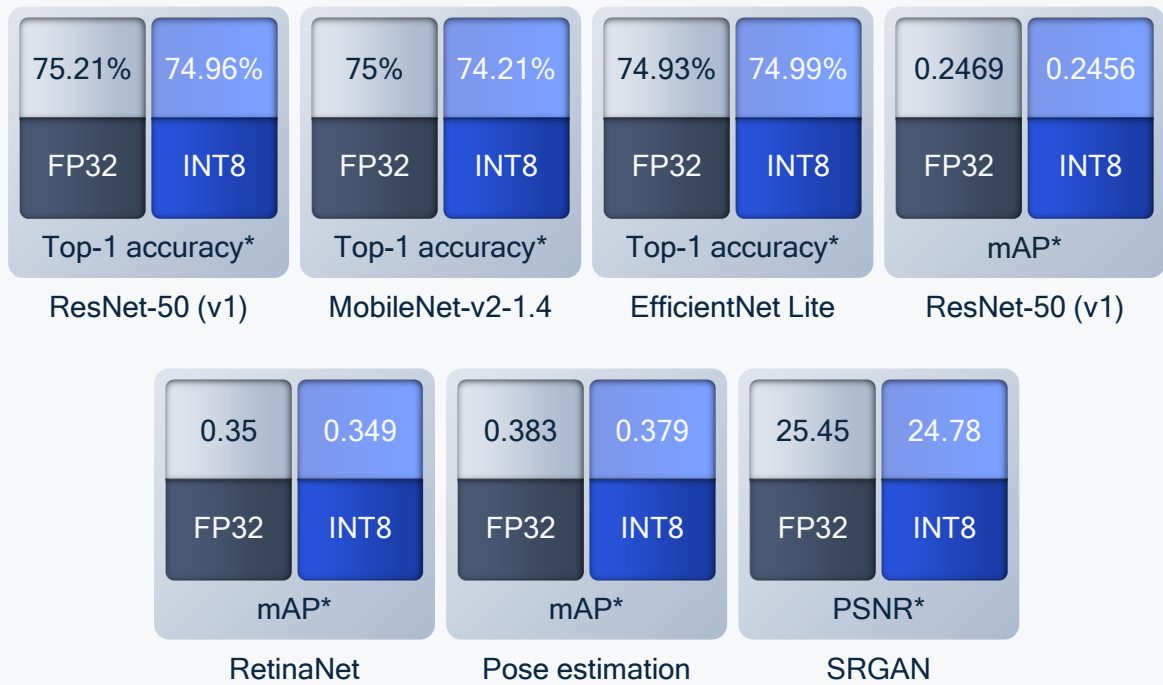
Provides accurate prediction of on-target performance by HW/run-time awareness

INT8 performance typically within 0.5-1% of FP32 performance

Output

Backprop

Act quant

RELU

+

Bias

Conv / FC

Wt quant

Weight

Input

Simulation ops added automatically at appropriate places in the model graph

Learn quantization parameters (QAT with Range Learning)

Fine tune model weights

**TensorFlow**

| | | | |
|---|---|---|---|
| 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* — ResNet-50 (v1) |
| 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 | |

**<1% Loss in accuracy***

**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 | 32.75 / 32.69 (FP32 / INT8) PSNR — ABPN |

# AIMET Model Zoo includes popular quantized AI models

## Accuracy is maintained for INT8 and INT4 models – less than 1% loss*

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

13

# Transformer Quantization

| Model | FP32 | INT8 |
|-------|------|------|
| BERT-base-uncased | 82.73 (GLUE) | 82.53 |
| DistilBERT-base-uncased | 80.35 (GLUE) | 79.81 |
| mobileBERT | 81.24 (GLUE) | 81.27 |
| VIT (vision transformer) | 81.30 | 81.50 |

AIMET quantizes transformers with high accuracy, comparable to FP32

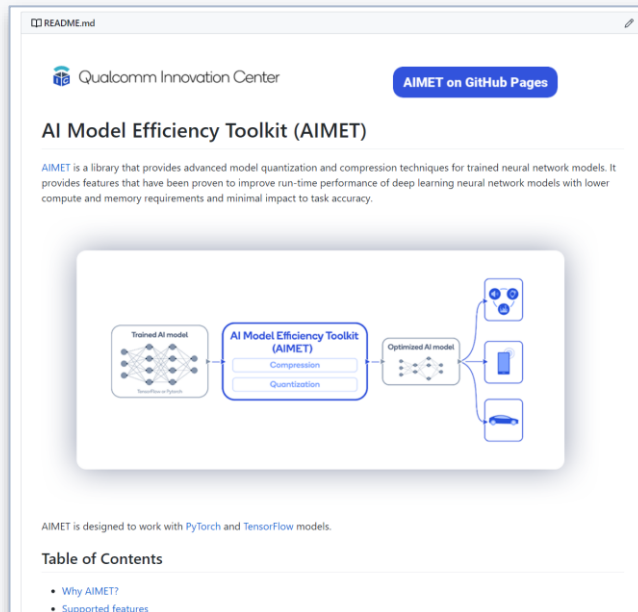# AIMET enables accurate INT4 (4-bit weights, 8-bit activations) for wide range of use cases

| Task | Model | FP32 | INT W4A8 |
|------|-------|------|----------|
| Classification | ResNet50 | 76.10% | 75.4% |
| | ResNet18 | 69.75% | 68.96% |
| | EfficientNet-Lite | 75.31% | 74.33% |
| | Regnext | 78.30% | 77.20% |
| Segmentation | Deeplabv3 (RN-50) | 76.07% | 75.91% |
| Super-resolution | ABPN | 31.97 dB | 31.67 (dB) |
| Pose detection | PoseNet (HRNet-32) | 0.765 | 0.763 |

Low resolution input (540p)

High resolution output 4bit Model (4K)

With better PTQ and QAT techniques, more models will achieve better power efficiency

# 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**

# Quantization whitepaper



**arxiv.org/abs/2201.08442**

# Explore our open-source projects and tools

# Qualcomm AI Stack

**Tools:**

Qualcomm AI Model Studio

- AIMET
- AIMET Model Zoo
- NAS
- Model analyzers

**Infrastructure:**

- Prometheus
- kubernetes
- docker

**AI Frameworks, Runtimes**

### AI Frameworks

TensorFlow  PyTorch  ONNX

### AI Runtimes

Qualcomm® Neural Processing SDK  ONNX RUNTIME  TF Lite Micro  Direct ML  TF Lite

Qualcomm® AI Engine Direct (QNN)

**Developer Libraries and Services**

Math Libraries  Compilers  Virtual platforms

Profilers & Debuggers  Programming Languages  Core Libraries

**System Software**

System Interface  SoC, accelerator drivers  Emulation Support

**OS**

android  Windows  Linux  Zephyr  ubuntu®  CentOS  QNX

**Platforms**

Smartphones  XR  ACPC  IoT  Robotics  Auto  Cloud

# Snapdragon® 8 Gen 2 Mobile Platform
# Qualcomm® AI Engine

## Micro-tile Inferencing

**More Hardware Acceleration**

**2X Tensor Accelerator Performance**

up to **4.35X** Performance improvement*

**60%** More power efficient

Multi-language Translation

World's **1ST** **INT4** support

**HEXAGON DIRECT LINK**

**ArcSoft®** AI Cinematic Mode

**Qualcomm AI Stack** | **UNREAL ENGINE**

Feature updates
Performance Improvements

INT4 Support

AI bot plug-in

AI Frameworks & Runtimes
Developer Libraries and Services
System Software
OS

**Qualcomm** Innovators Developer Kit

**Qualcomm** Sensing Hub

**Dual** AI Processor

Always-sensing Camera

Qualcomm AI Studio

# Thank you