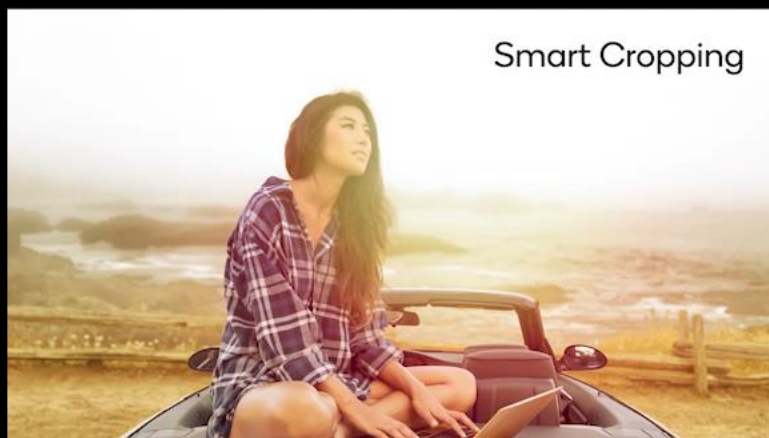
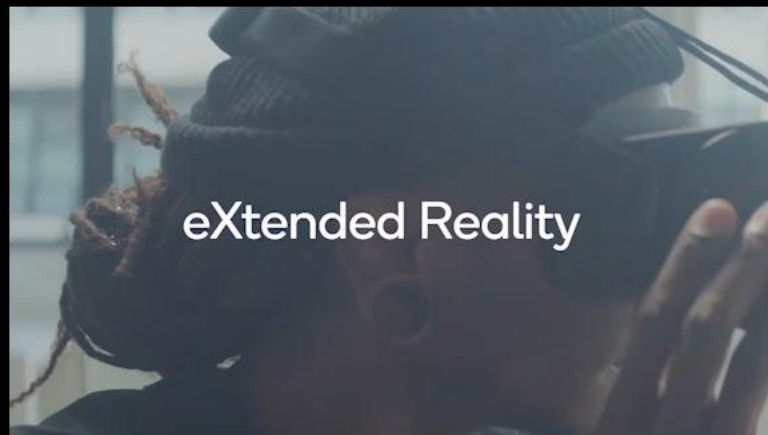




# What we need to Transform Lives and Industries with On-Device AI, Cloud and 5G

Ziad Asghar, Vice President, Product Management at Qualcomm Technologies Inc.

Qualcomm

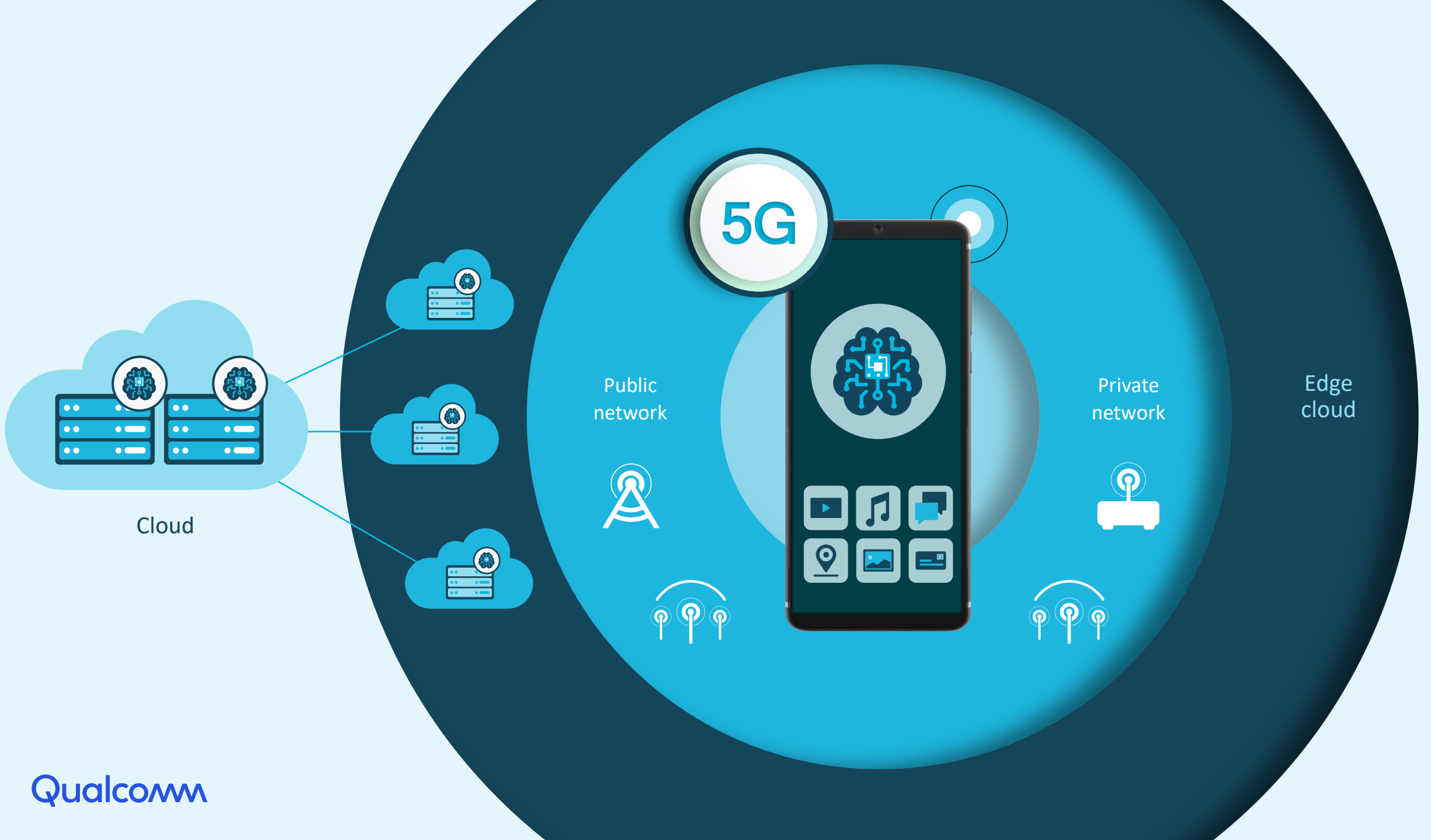




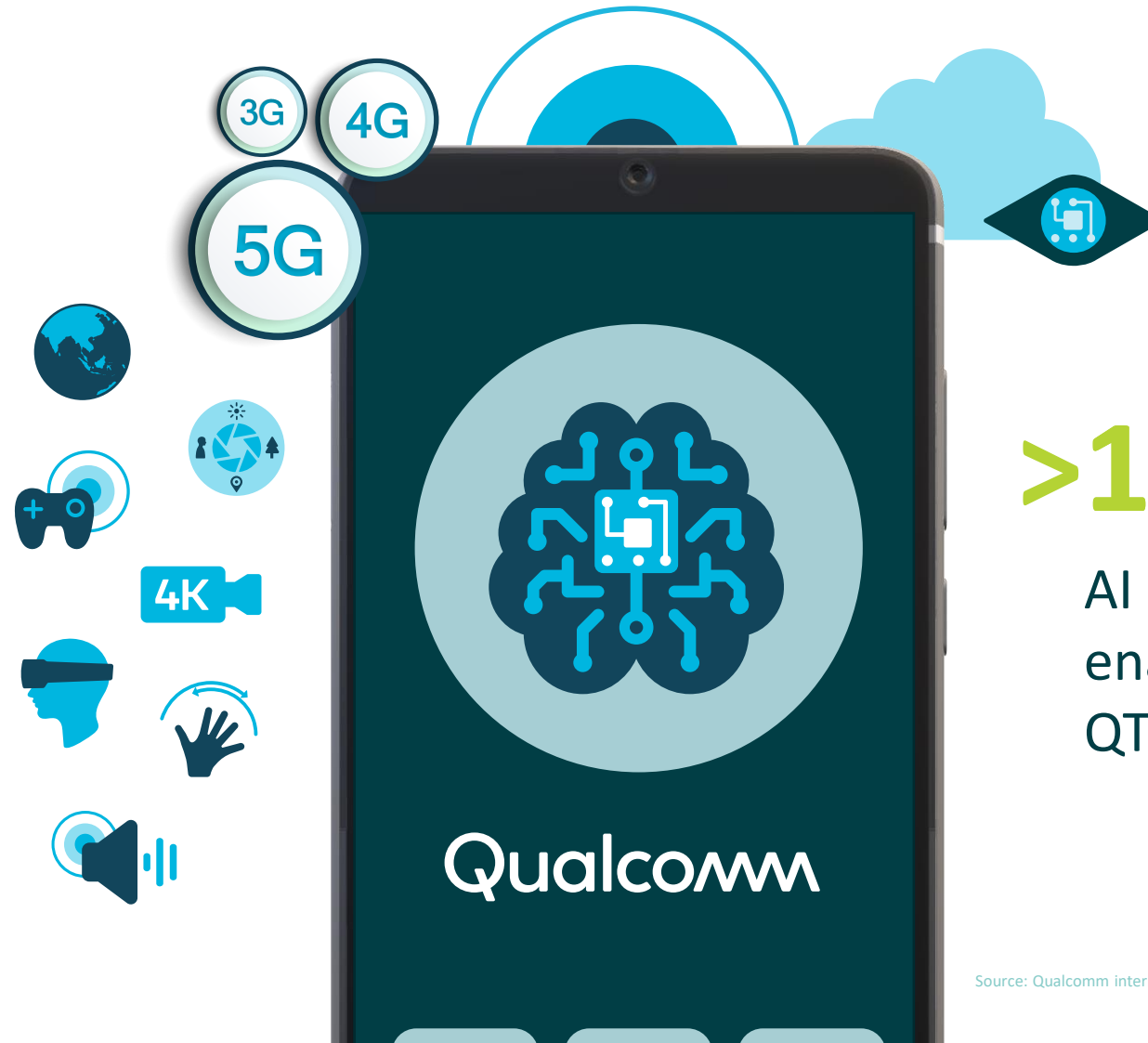








# Mobile — THE most pervasive AI platform



>105 Billion

AI capable devices  
enabled with  
QTI technology

# Snapdragon AI Milestones

2016

**Snapdragon 820**  
first introduction to the Qualcomm® Hexagon™ Vector eXtensions for more powerful AI processing

2017

**3<sup>rd</sup> generation Qualcomm® AI Engine with Snapdragon 845**  
Enabling an AI-based voice assistant

2018

**4<sup>th</sup> generation Qualcomm® AI Engine with Snapdragon 855**  
including the first Tensor Accelerator, bringing 7 TOPs and enabling AI based single-mic noise cancellation

2019

**5<sup>th</sup> generation Qualcomm® AI Engine with Snapdragon 865**  
featuring the first on-device real time voice translation powered by 15 TOPS

2020

**6<sup>th</sup> generation Qualcomm® AI Engine with Snapdragon 888**  
re-engineered Hexagon 780 Processor features a fused AI-accelerator architecture and brings the total Qualcomm AI Engine performance up to an astonishing 26 TOPS

**26 TOPS**





# Key Ingredients for AI Leadership

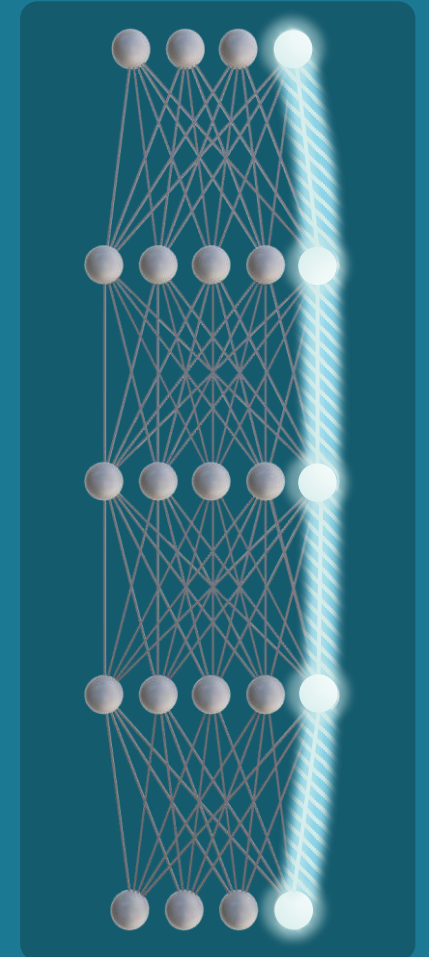
Hardware



+ Software



+ Tools



## Connectivity

3rd Gen Snapdragon  
X60 Modem-RF System

3rd Gen mmWave

5G Carrier Aggregation\*

Qualcomm® FastConnect™ 6900  
6GHz, 4K QAM, BT 5.2

## Performance

Kryo 680  
X1  
Architecture,  
25% faster

Adreno 660  
35% faster  
graphics  
rendering

5nm  
Process  
Technology

## Security

Hypervisor

CAI Compliant camera

## AI

6th Gen  
AI Engine  
Hexagon 680  
Fused AI  
Accelerators

3x Performance per watt  
26 TOPS  
2nd Gen  
Qualcomm  
Sensing Hub



## Gaming

3rd Gen Qualcomm®  
Snapdragon Elite Gaming™

Ultra smooth gaming:

Qualcomm® Game Quick Touch,

Desktop level features:

Qualcomm® Variable Rate Shading

## Camera

2.7 Gigapixels  
Per Second:  
35% Faster

Burst Capture  
120 photos in  
1 second at 12MP

Triple ISP:

Triple  
Concurrency

4K HDR  
video capture with  
computational HDR

Triple Parallel  
Processing

low light architecture  
photo capture

Triple Capture:

Triple 4K HDR  
Video Capture

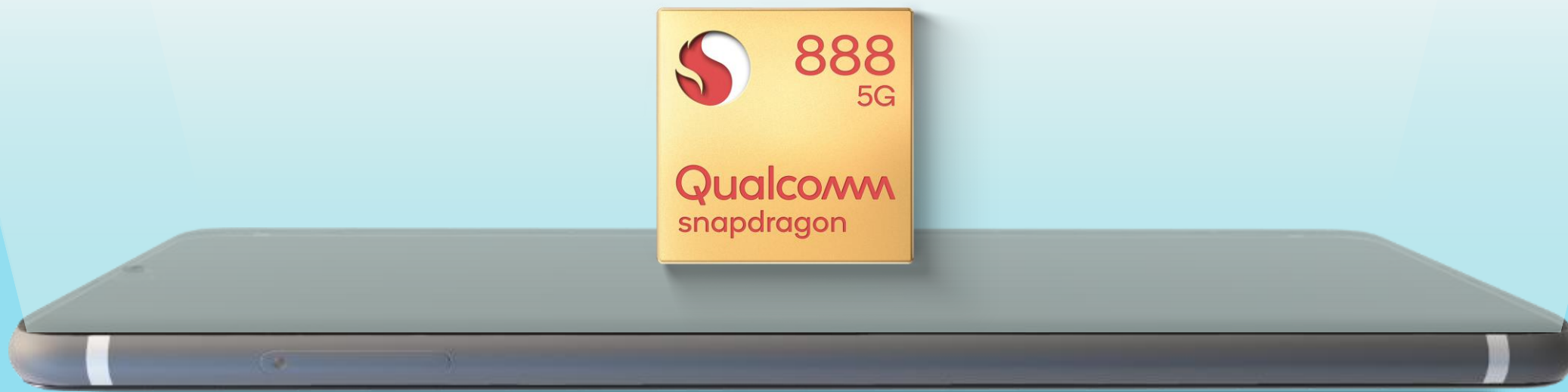
10-bit HDR  
HEIF photo  
capture

Triple 28MP  
Photo Capture

\* across TDD and FDD

Qualcomm Snapdragon Elite Gaming, Qualcomm Game Smoother,  
Qualcomm Game Performance Engine and Qualcomm Game Color Plus  
are product of Qualcomm Technologies, Inc. and/or its subsidiaries.

# 6<sup>th</sup> Generation Qualcomm AI Engine





# Hexagon 780 Processor



Qualcomm Spectra, Qualcomm Sensing Hub, Qualcomm Processor Security, Qualcomm Kryo and Qualcomm FastConnect are products of Qualcomm Technologies, Inc. and/or its subsidiaries.



# Fused AI accelerator



Qualcomm Spectra, Qualcomm Sensing Hub, Qualcomm Processor Security, Qualcomm Kryo and Qualcomm FastConnect are products of Qualcomm Technologies, Inc. and/or its subsidiaries.



## 43% faster AI performance

New instructions:

4-input mixed precision  
dot product

Wave Matrix Multiply for  
16/32-bit floating point

## Fused AI accelerator:

Up to  
**3X**  
Performance  
per watt

Tensor  
**2X**  
compute  
capacity

Scalar  
**50%**  
performance  
improvement

## Shared Memory:

**16X**  
dedicated  
memory

Up to  
**1000x**  
hand off time improvement  
in certain use cases

Qualcomm

Compared to previous generation

Qualcomm Spectra™  
580 ISP

Qualcomm® Adreno™  
660 GPU

Scalar

Tensor

Vector

Qualcomm®  
Sensing  
Hub

Qualcomm®  
Processor  
Security

Qualcomm®  
Kryo™  
680 CPU

Qualcomm®  
Snapdragon™  
X60 5G Modem-RF  
System

Qualcomm®  
FastConnect™  
6000

# Power Efficiency

Efficiency: (inf/s/W)

Resnet 34

Snapdragon 888

227

Company A

150

Company B

137



MobileNet SSD

215

120

102

Deeplabv3

66

39

38

Inception V3

110

80

74

# 2<sup>nd</sup> Gen Qualcomm Sensing Hub <1mA

Qualcomm

Qualcomm® Hexagon™  
780 Processor

Qualcomm®  
Sensing  
Hub

Qualcomm®  
Processor  
Security

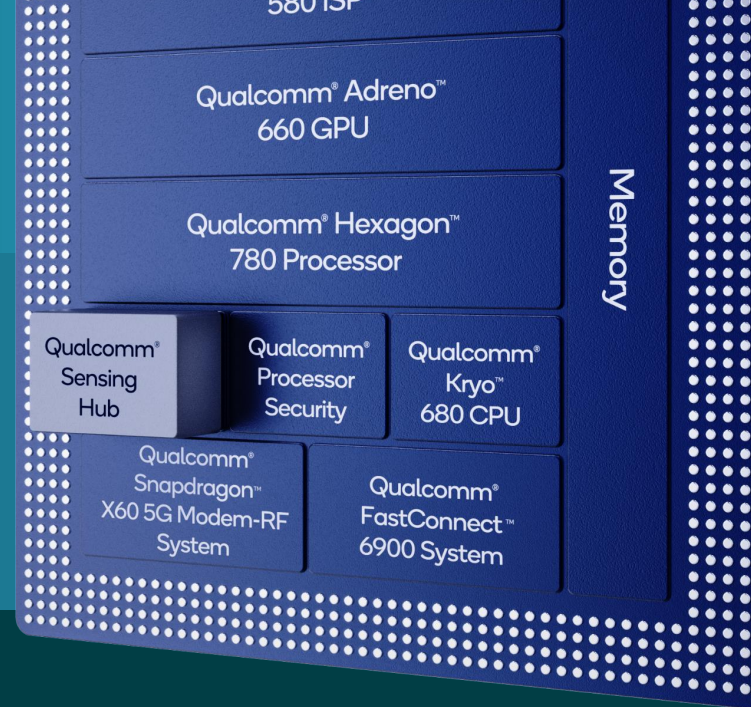
Qualcomm®  
Kryo  
680 Core

Qualcomm®  
Snapdragon™  
X60 5G Modem-RF  
System

Qualcomm®  
FastConnect  
6900 System



# AI performance



With 2<sup>nd</sup> Gen  
Qualcomm Sensing Hub



80%

1st Gen  
Qualcomm Sensing Hub



5X

Task offload from  
Hexagon Processor

# Sensing Hub – Use case



Tools + Compilers

AIMET

TVM

Models

ResNet

SSD

MobileNet

Mobile  
BERT

VDSR

DeepLab

Applications



Qualcomm®  
Neural Processing SDK



Android Neural  
Networks API

Frameworks

# Qualcomm AI software stack

Supporting every  
AI software layer  
from applications  
to the metal



Runtime

Qualcomm® AI  
Engine Direct

NNAPI

SDKs

# AIMET + Tools

>4x

Increase in  
performance  
per watt with  
quantization

## AI Model Efficiency Toolkit

3x

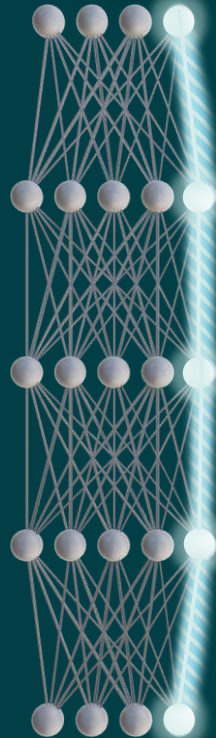
Compression with  
less than 1% loss  
in accuracy\*

Improved/ Robust  
quantization  
for INT16,8,4

Quantization  
aware training  
with range learning

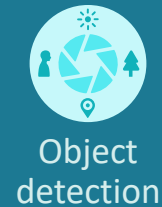
Mix precision support

Opensource



AI Model Efficiency Toolkit is a product  
of Qualcomm Innovation Center, Inc.

## AIMET Model Zoo



```
1 def quantized_add_generic(size, s, b, a_offset, b_offset, a_mult, b_mult, output, target, ctx):
2     # Construct the TVM computation.
3     A_offset = tvm.var('A_offset', dtype='uint8')
4     B_offset = tvm.var('B_offset', dtype='uint8')
5     A_mult = tvm.var('A_mult', dtype='uint16')
6     B_mult = tvm.var('B_mult', dtype='uint16')
7
8     N = tvm.var('N')
9     A = tvm.placeholder((N,), name='A', dtype='uint8')
10    B = tvm.placeholder((N,), name='B', dtype='uint8')
11
12    C = tvm.compute((A.shape),
13                    lambda i: ((A[i].astype('int32') - A_offset.astype('int32')) * A_mult.astype('int32')) +
14                               ((B[i].astype('int32') - B_offset.astype('int32')) * B_mult.astype('int32')), name='C'))
15
16    # Create the schedule.
17    s = tvm.create_schedule(C.op);
18    px, x = s[C].split(s[C].op.axis[0], nparts=1)
19    s[C].bind(px, tvm.thread_axis("pipeline"))
20
21    # Construct the callable object "func" corresponding to the computation.
22    func = tvm.build(s, [A, B, C, N, A_offset, B_offset, A_mult, B_mult], target, name='qadd_tvm')
23
24    func(tvm.nd.array(s, ctx=ctx), tvm.nd.array(b, ctx=ctx), output,
25         size, a_offset, b_offset, a_mult, b_mult)
```



Custom operators efficiently written in

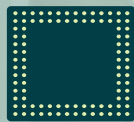


# AI Performance Continuum

IoT

0.5-1+  
TOPS





0.5-1+  
TOPS

# AI for IoT

Home  
Industrial/Enterprise  
Smart cities

Object Detection

Access Control

Semantic Filtering

Classification

Picking and Sorting

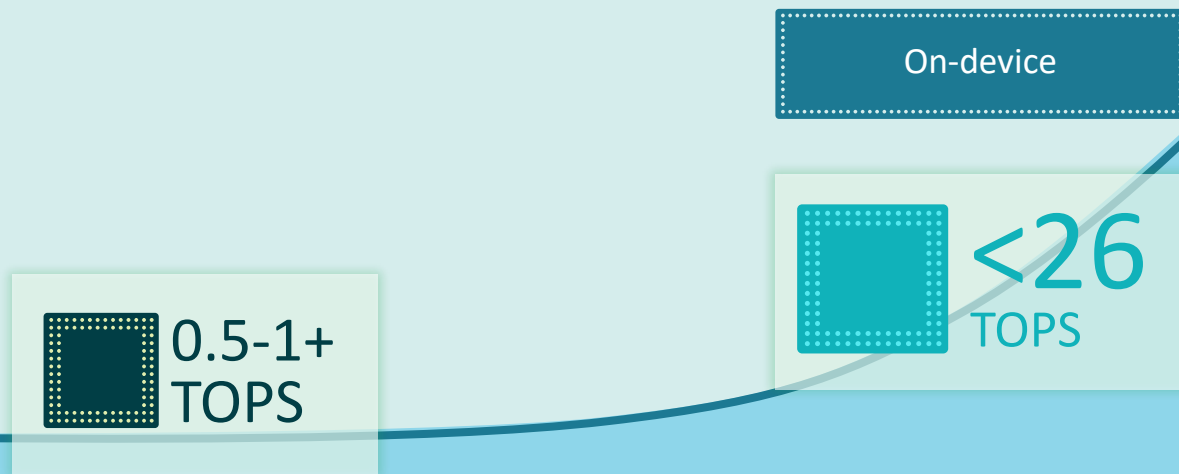
Super Resolution

<10 Concurrencies

Qualcomm



# AI Performance Continuum







# On-Device AI

Mobile

Always connected PCs

XR

Segmentation

De-noise

Semantic Filtering

Classification

Detection

Super Resolution

10 – 15 Concurrencies

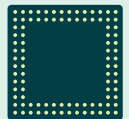
# AI Performance Continuum



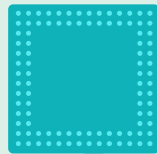
Qualcomm® Snapdragon Ride™ Platform  
Automotive ADAS

Qualcomm  
Cloud AI 100

400  
TOPS



0.5-1+  
TOPS



>26  
TOPS



Qualcomm  
Cloud AI 100

400  
TOPS

# AI for Auto

## Snapdragon Ride

Qualcomm  
Cloud AI 100

800  
TOPS

Occupancy

Voice Command

Lane Merge

Traffic Assist

Monitoring

Surround System

Pedestrian Assist

Lane Change

8 – 10 Concurrencies

10 – 12 Concurrencies

10 TOPS

System Performance

800+ TOPS



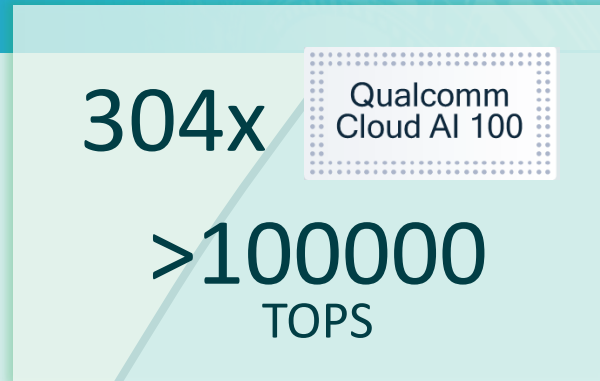
L1/L2



L2+L3



L4/L5



Qualcomm  
Cloud AI 100

400  
TOPS

# Qualcomm Cloud AI 100

Industrial/Enterprise  
Infrastructure  
Auto

Recommendation

Linguistics

Communication

Classification

Pedestrian Assist

Intersection Assist

3 – 5 Concurrencies

400 TOPS

Performance

125 POPS

Qualcomm Cloud AI 100 is a product of Qualcomm Technologies, Inc. and/or its subsidiaries.



Qualcomm Cloud AI  
100

400  
TOPS

1 card



Gigabyte  
platform

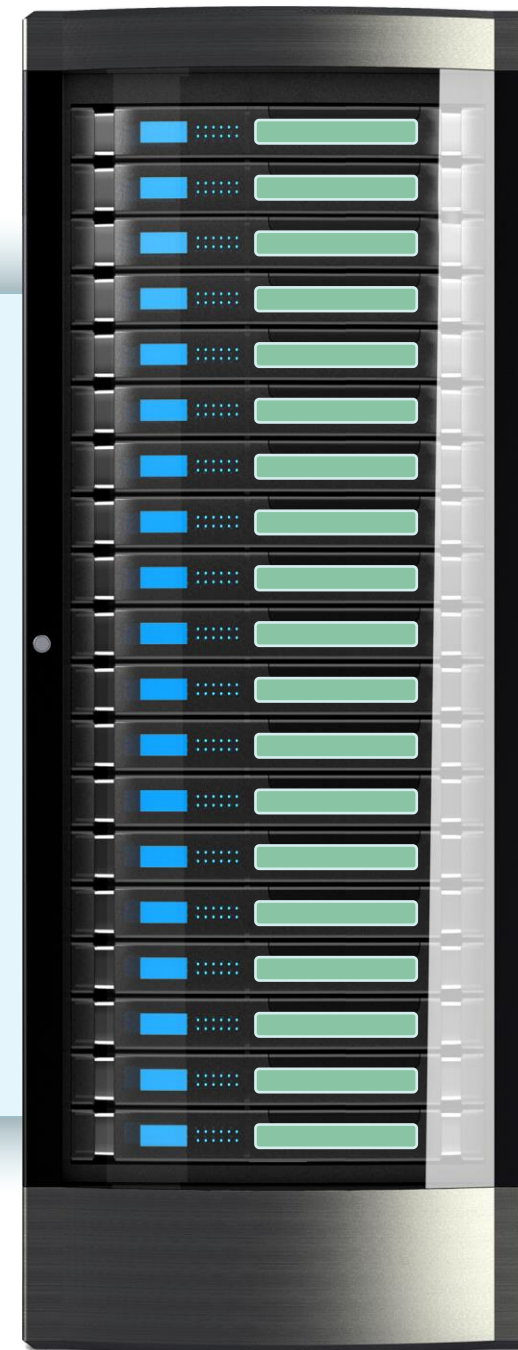
6.4  
POPS

16 cards  
in one  
server

Full  
server rack

125+  
POPS

20 server  
units in  
one server  
rack

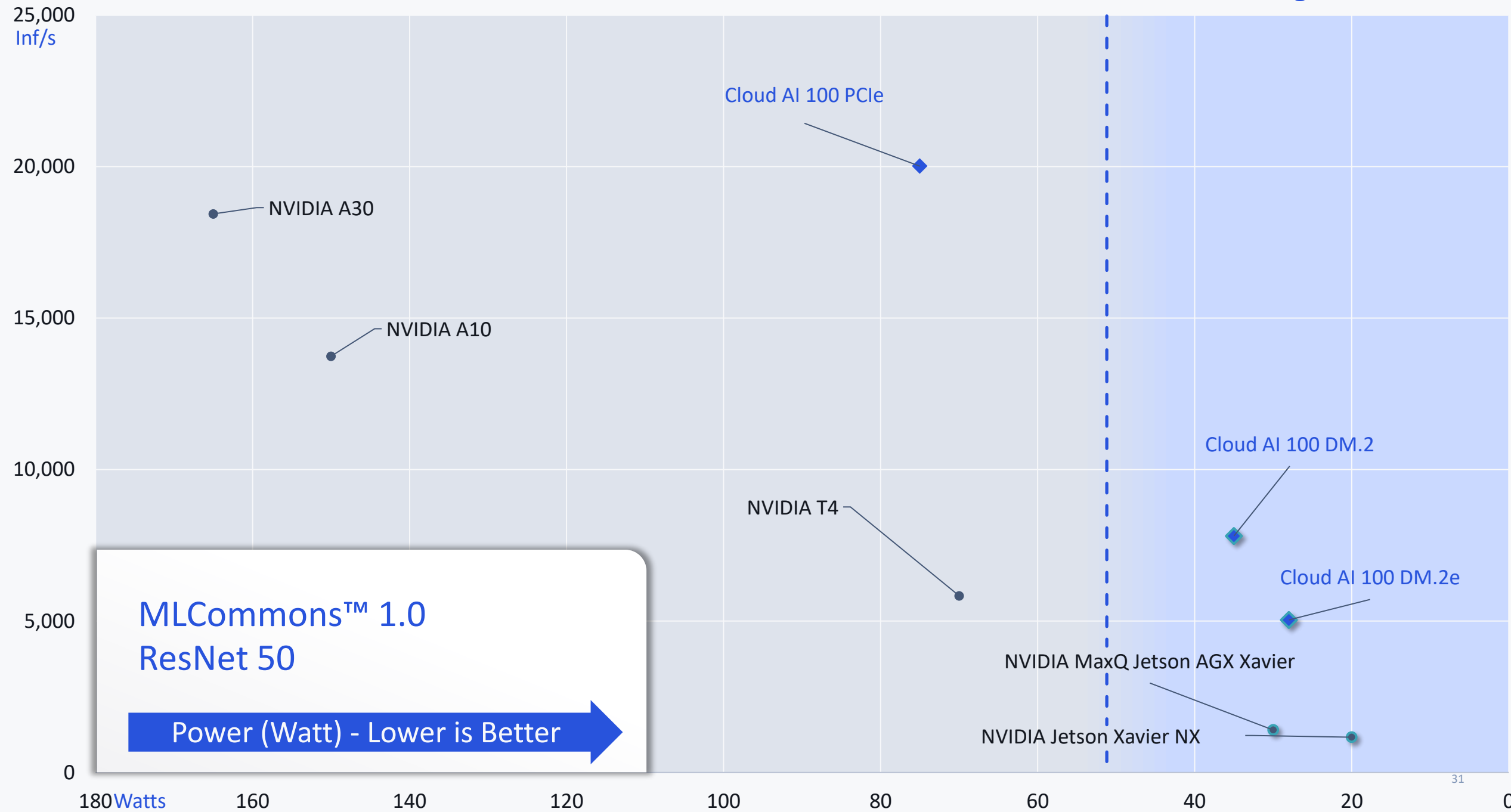


1 Peta OPS = 1000 TOPS

# Offline Performance

>50 Watt - Cloud

<50 Watt - Edge



# The future of AI

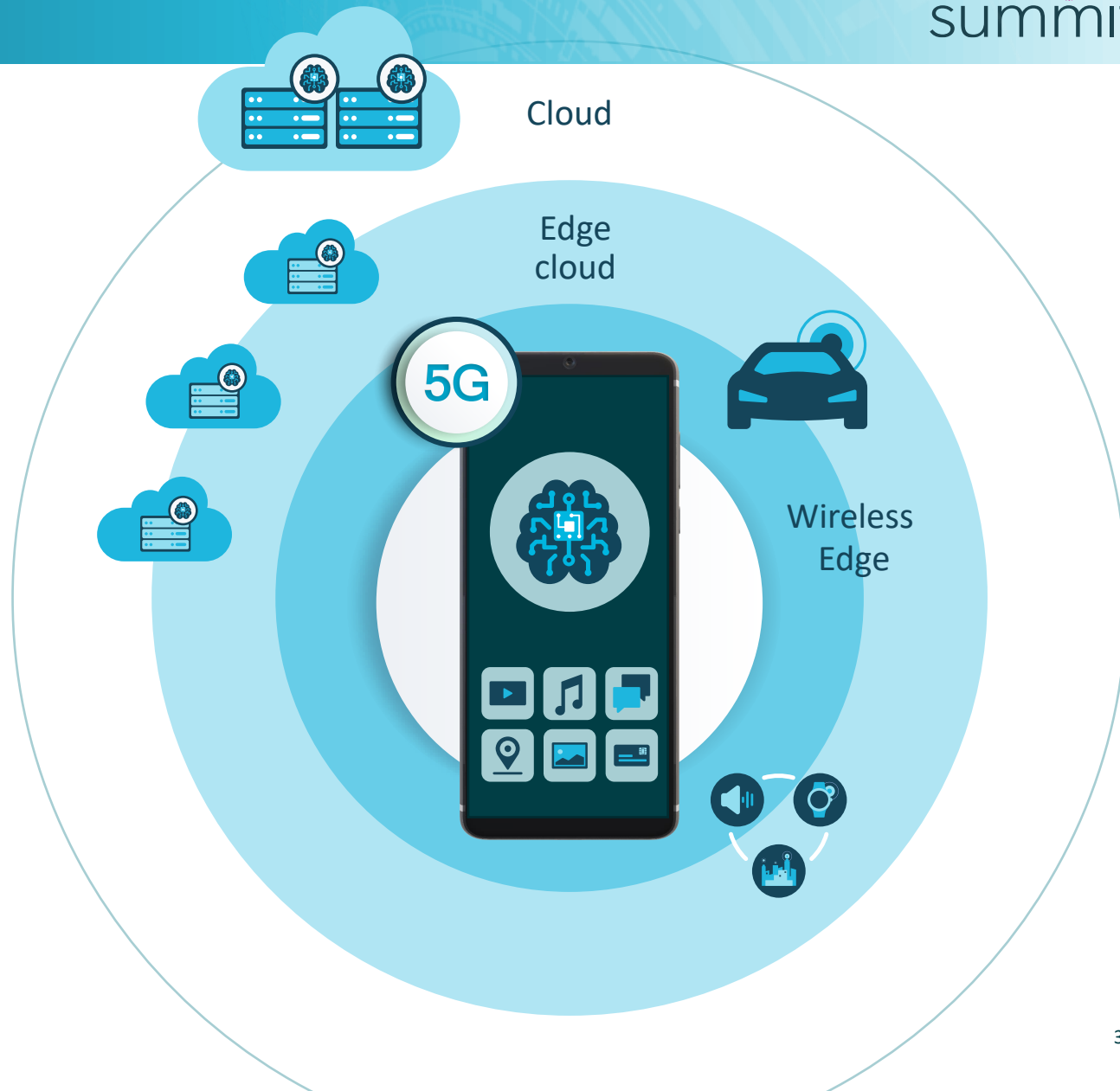
2021  
embedded  
**VISION**  
summit



Qualcomm



# We're creating a future of distributed intelligence







# Thank You

- Qualcomm AI page:

<https://www.qualcomm.com/invention/artificial-intelligence>

- Qualcomm AI research:

[https://www.qualcomm.com/invention/artificial-intelligence/ai-research?cmpid=fofyus193556&gclid=CjwKCAjw19z6BRAYEiwAmo64LfQjU8vqH8TxqKTM2PZQp8JibXrjev85wLfKFknJnS\\_b494yZ7e\\_WhoCPQkQAvD\\_BwE](https://www.qualcomm.com/invention/artificial-intelligence/ai-research?cmpid=fofyus193556&gclid=CjwKCAjw19z6BRAYEiwAmo64LfQjU8vqH8TxqKTM2PZQp8JibXrjev85wLfKFknJnS_b494yZ7e_WhoCPQkQAvD_BwE)

- Qualcomm® Platform Solution Ecosystem:

<https://www.qualcomm.com/support/qan/platform-solutions-ecosystem>

- GitHub AI Model Efficiency Toolkit (AIMET):

<https://github.com/quic/aimet>

- Qualcomm Mobile AI page:

<https://www.qualcomm.com/products/smartphones/mobile-ai>

- Qualcomm Mobile AI blog:

<https://www.qualcomm.com/news/onq/2020/12/02/exploring-ai-capabilities-qualcomm-snapdragon-888-mobile-platform>

- Qualcomm Cloud AI 100 blog:

<https://www.qualcomm.com/news/onq/2021/03/15/qualcomm-cloud-ai-100-amd-epyc-7003-series-processor-and-gigabyte-server>