

An aerial photograph of a coastal landscape. In the foreground, there are green, rolling hills and a small peninsula with some buildings. The water is a deep blue, and a small white sailboat is visible in the distance. In the background, there are large, rugged mountains with some greenery and some rocky peaks. The sky is blue with scattered white clouds.

Qualcomm

# Ziad Asghar

Vice President, Product Management, AI & Strategy

Qualcomm Technologies, Inc.



**AI is in every aspect of our lives**



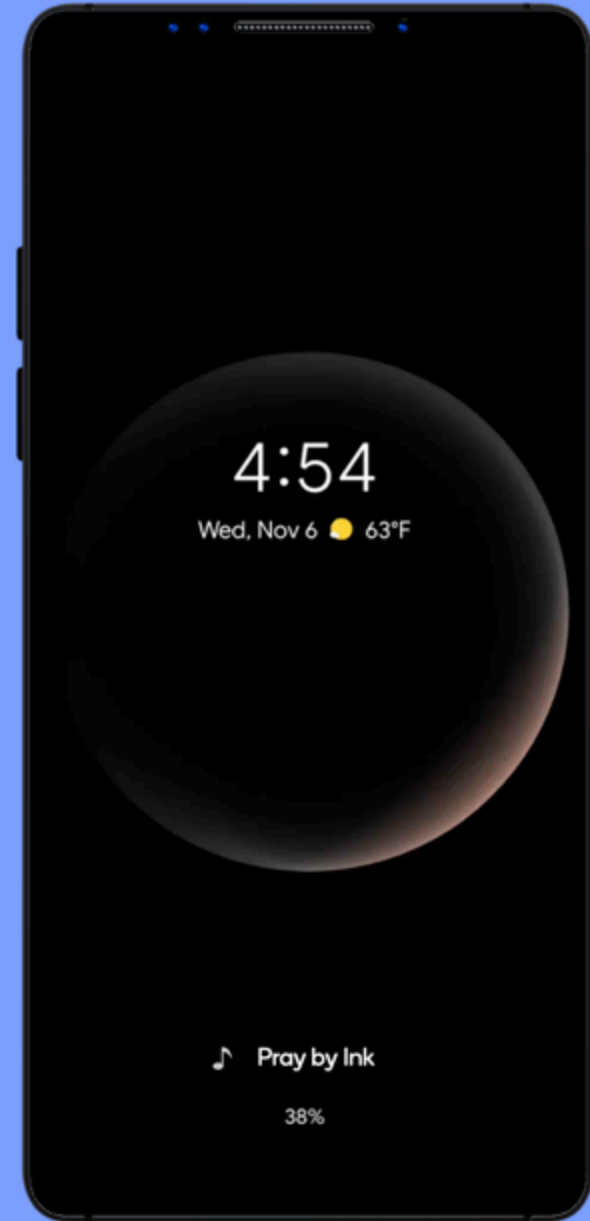


# Conventional music detection

# Music detection with on-device AI

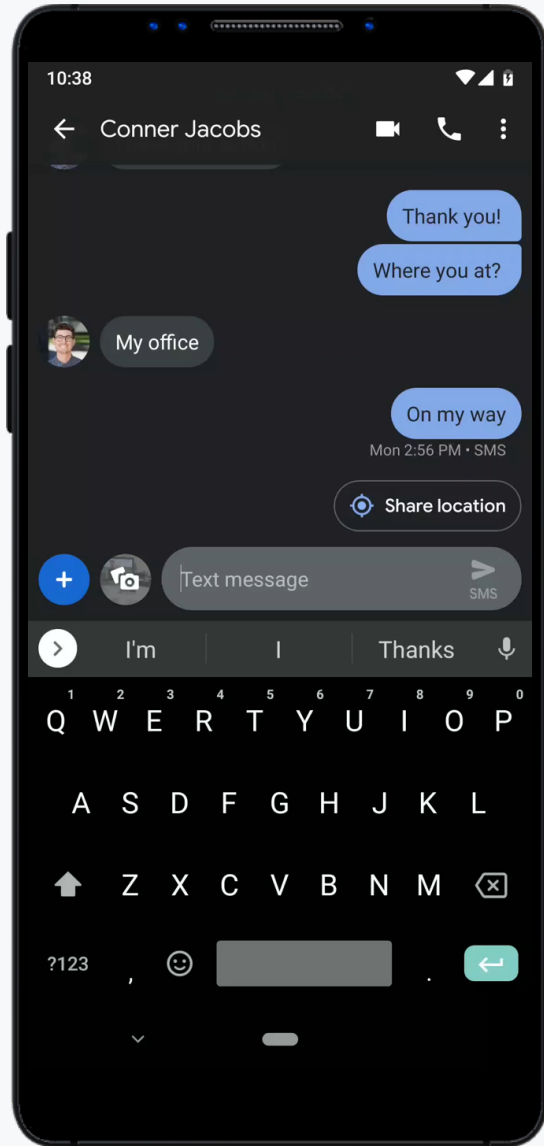
---

Always-on neural nets  
at ultra-low power

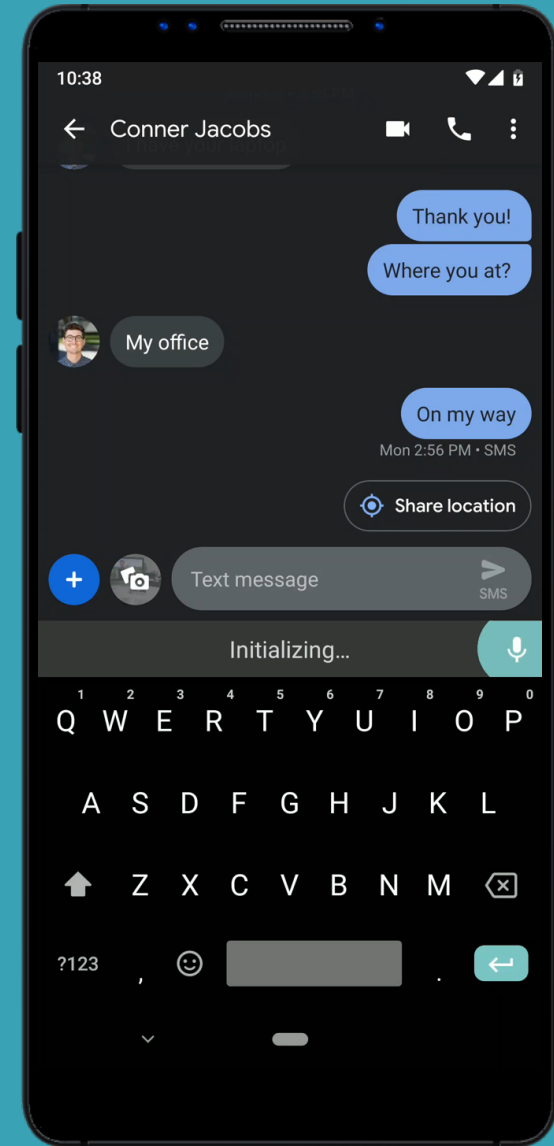


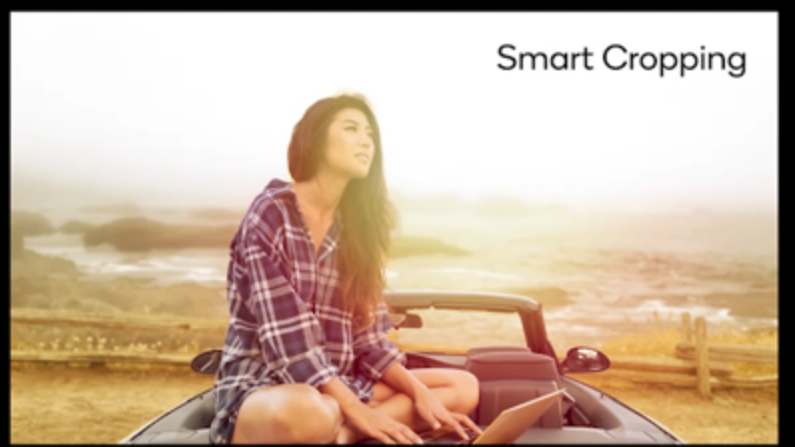
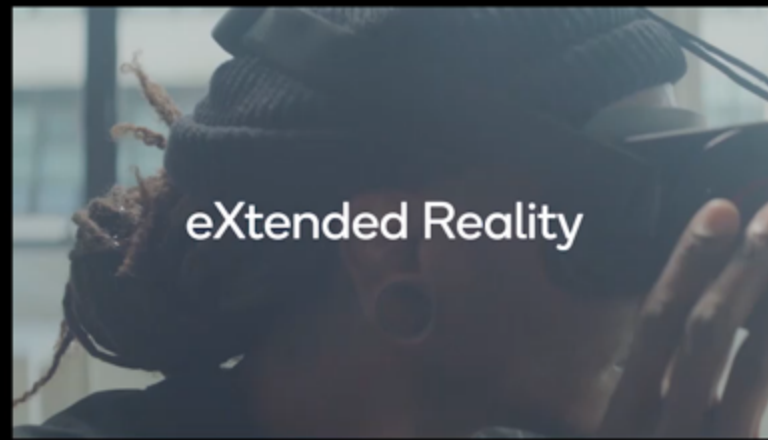


# Cloud speech-to-text



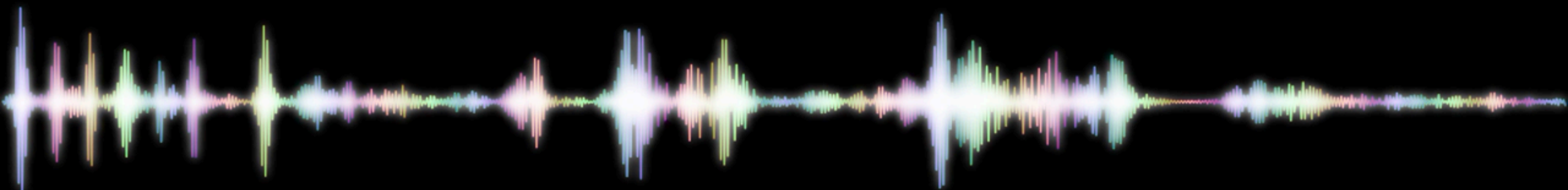
# On device speech-to-text







**AI is happening as we “speak”**





Software

---

# Hardware

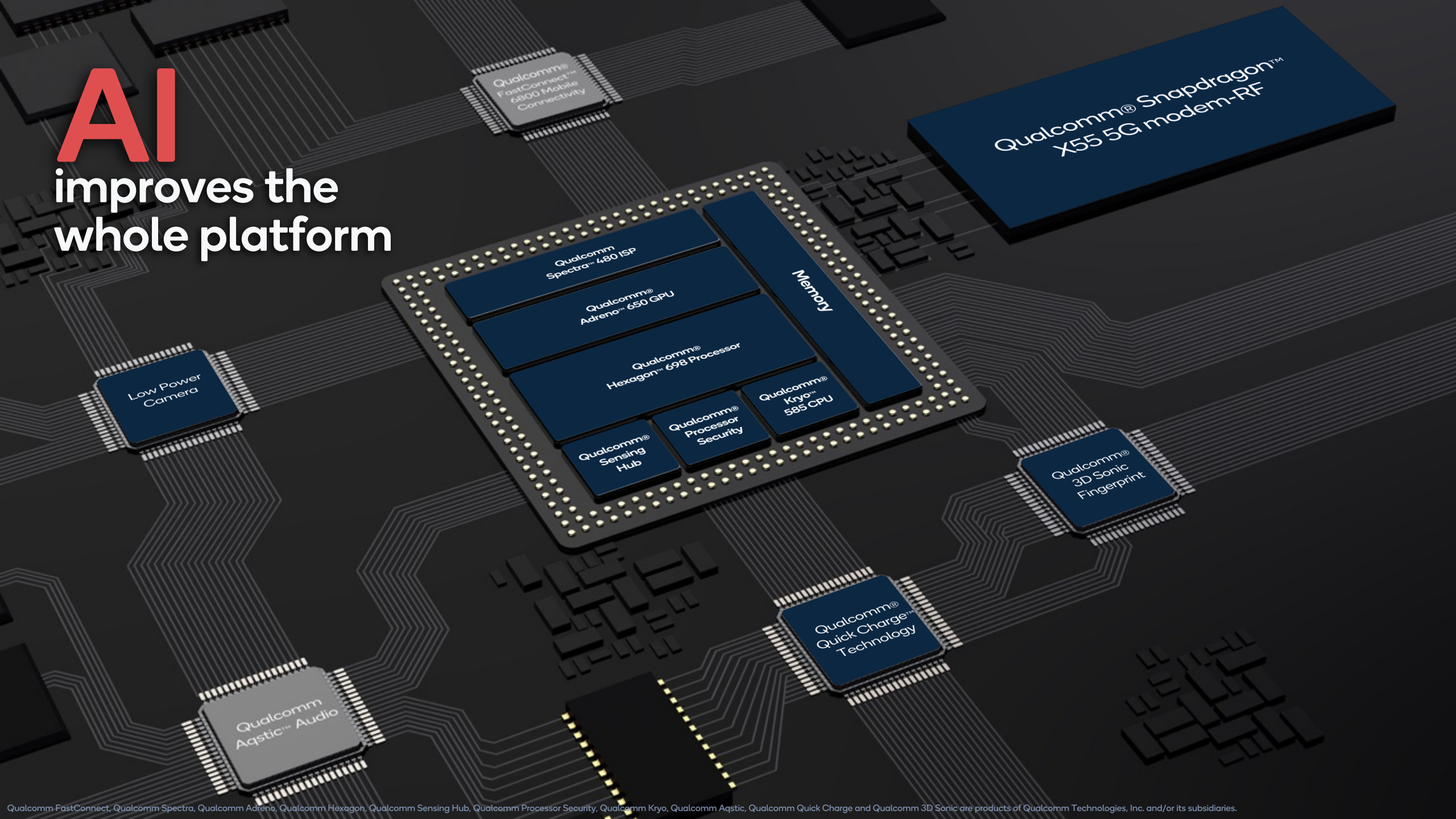
---

Developer tools



# AI

improves the  
whole platform



# 5<sup>th</sup> Generation Qualcomm® AI Engine





## Adreno 650

New AI mixed precision instructions

2x higher TOPS\*

16-bit and 32-bit FP

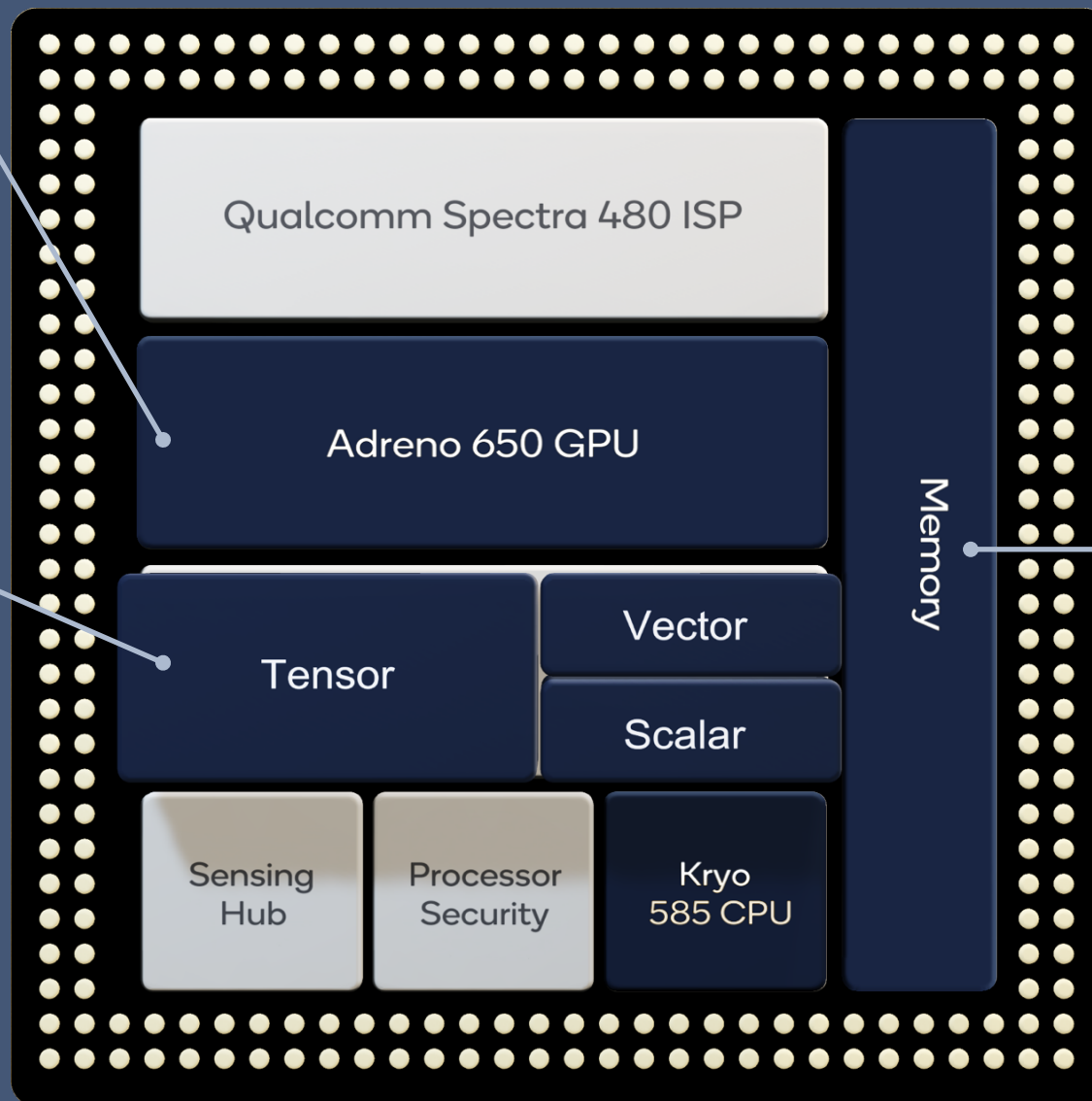
## Hexagon 698

New Tensor Accelerator

- 4x higher TOPS\*
- Up to 35% power savings\*
- 8-bit and 16-bit INT

Deep Learning Bandwidth Compression

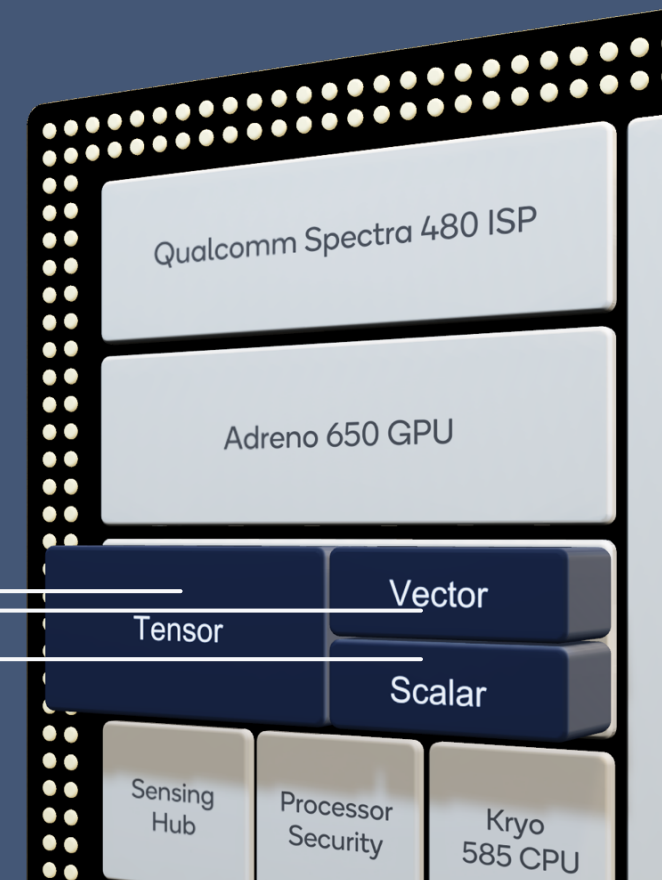
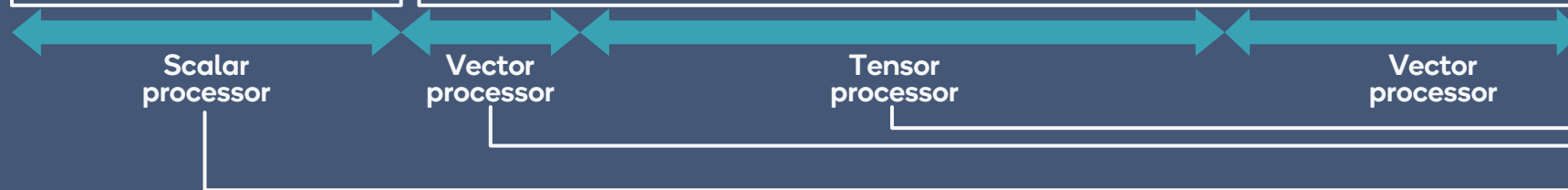
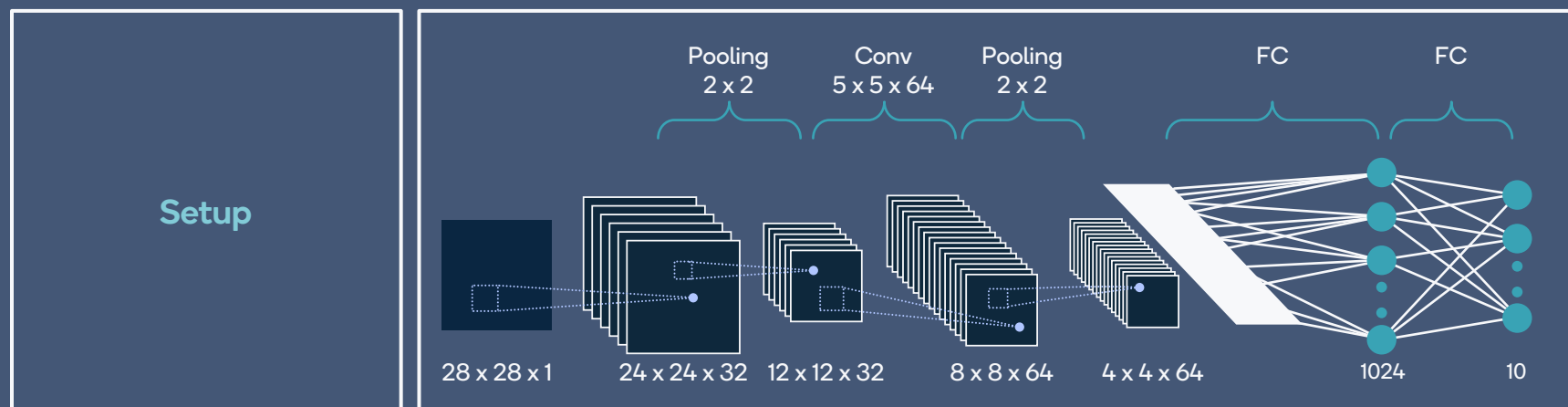
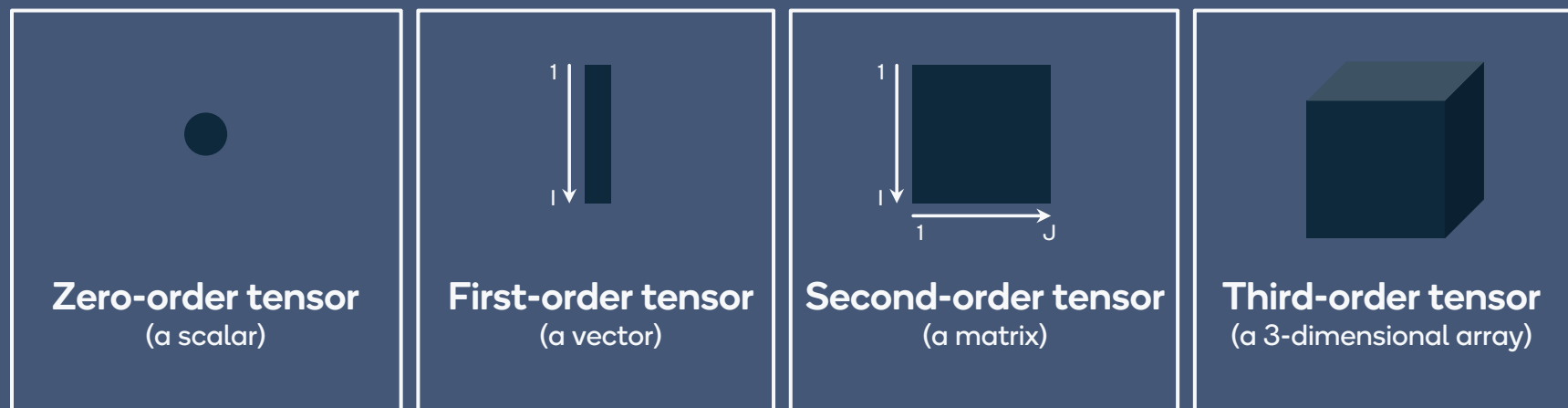
- Up to 50% lossless Compression
- Frees up bandwidth for other parts of the SoC
- Saves power due to reduced memory transfers



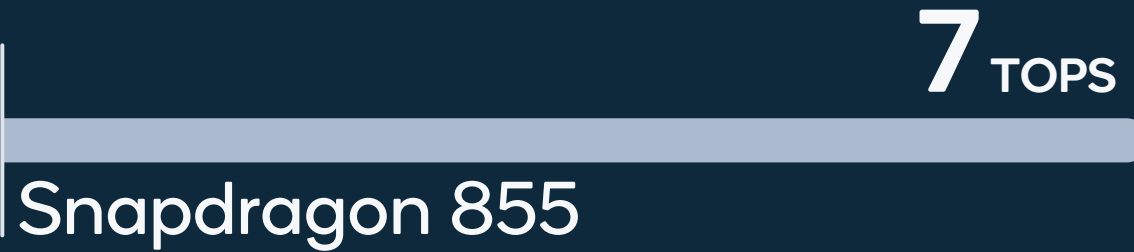
## LP-DDR5 memory

30% more bandwidth\*

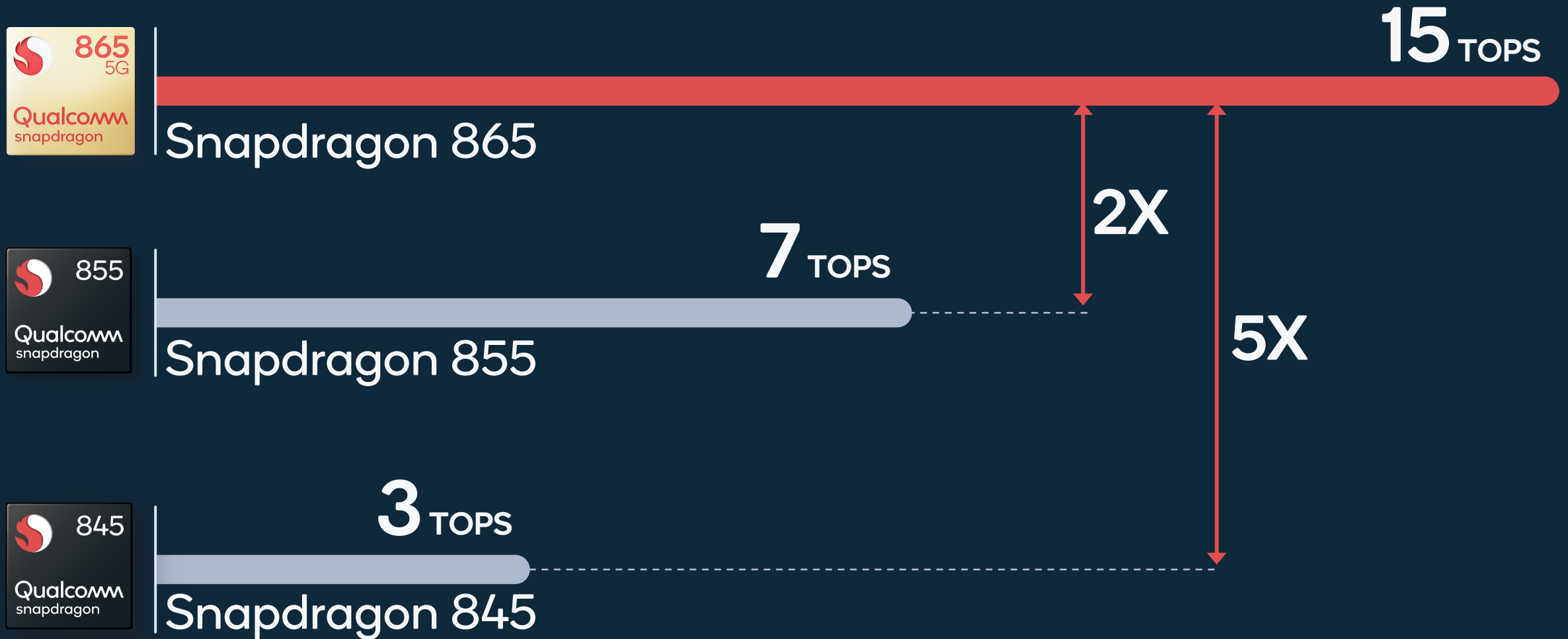
Improved AI processing



# Trillion operations per second

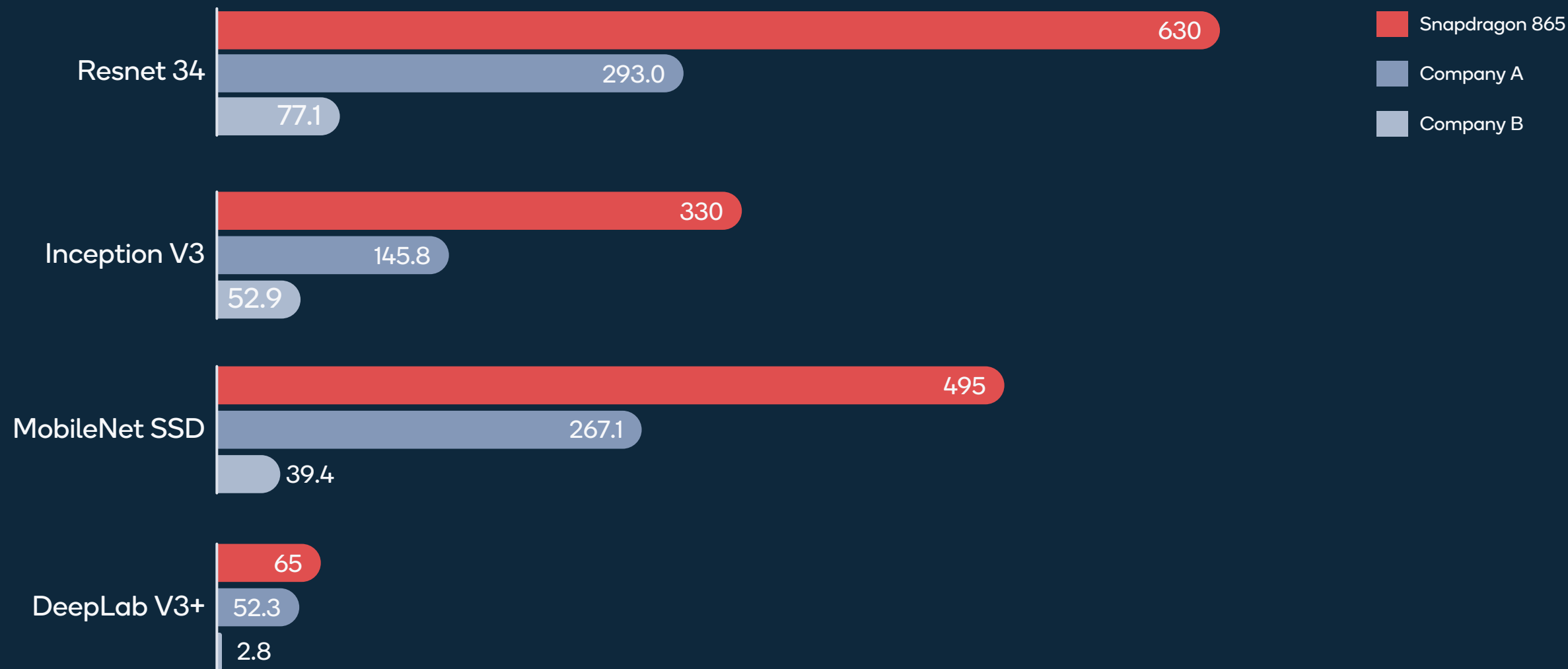


# Trillion operations per second



# Peak performance on classification networks

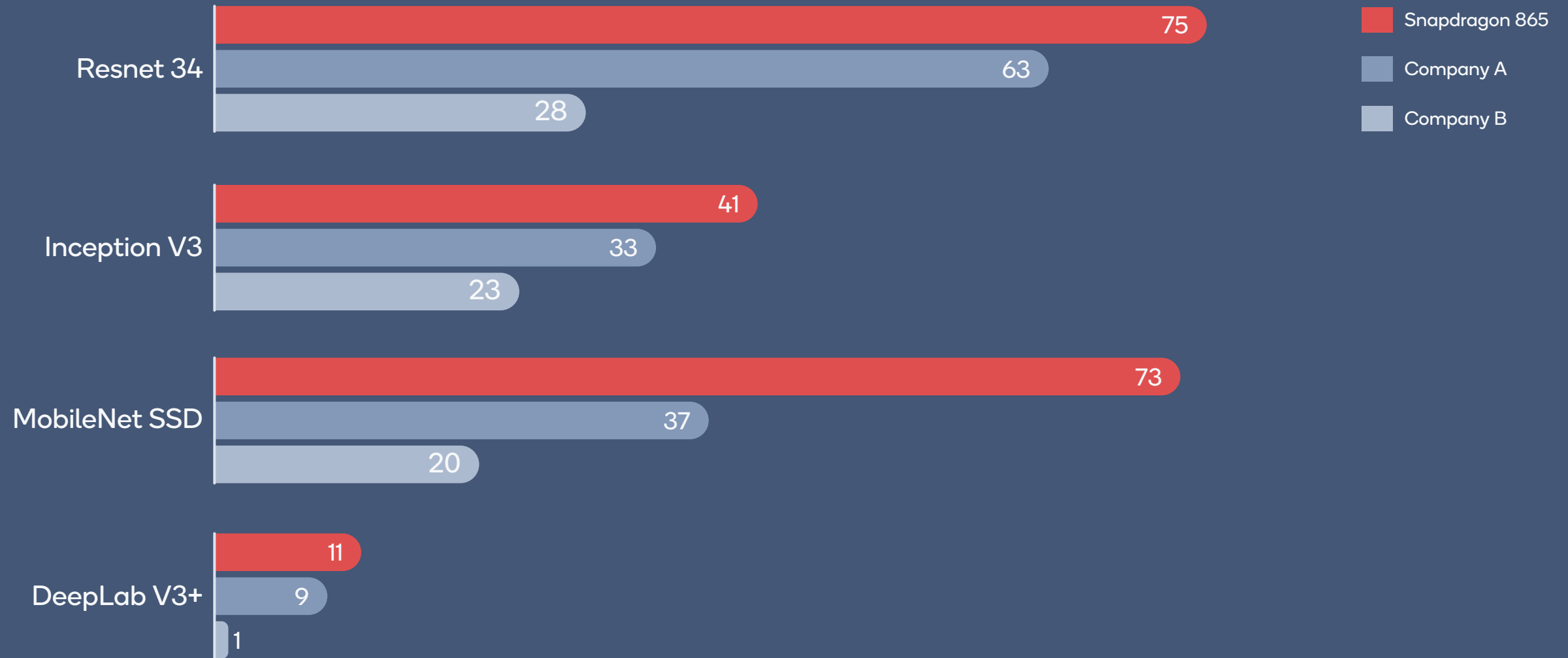
(Inf/sec)





# Power consumption

(Inf/Watt)



# Qualcomm Sensing Hub

<1mA

Qualcomm® Hexagon™  
698 Processor

Audio

Video

Sensors

Qualcomm®  
Processor  
Security

Qualcomm®  
Kryo™  
585 CPU

Memory

Safety



Health & Fitness



Security



Context



Low-power camera



Hardware

---

**Software**

---

Developer tools

An aerial photograph of a coastal landscape. In the foreground, there are green, rolling hills and a small peninsula with some buildings. The water is a deep blue, and a small white sailboat is visible in the distance. The background features more rugged, brownish hills under a blue sky with scattered white clouds.

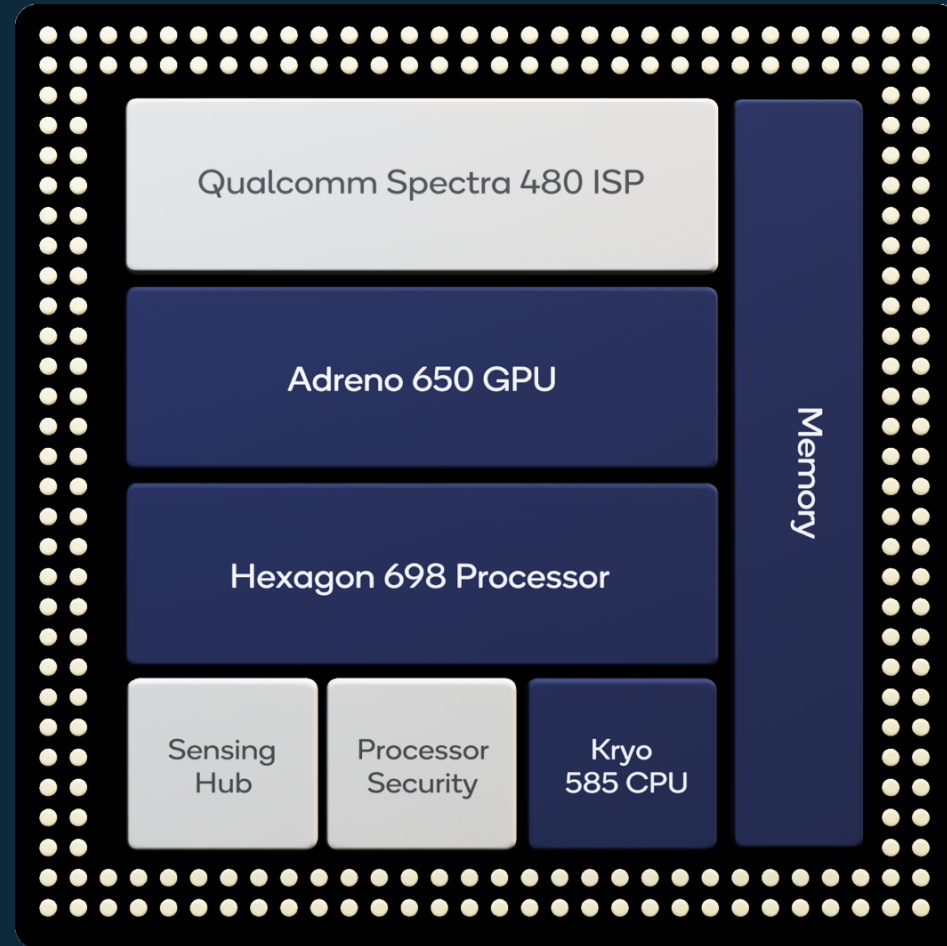
Qualcomm

# Jeff Gehlhaar

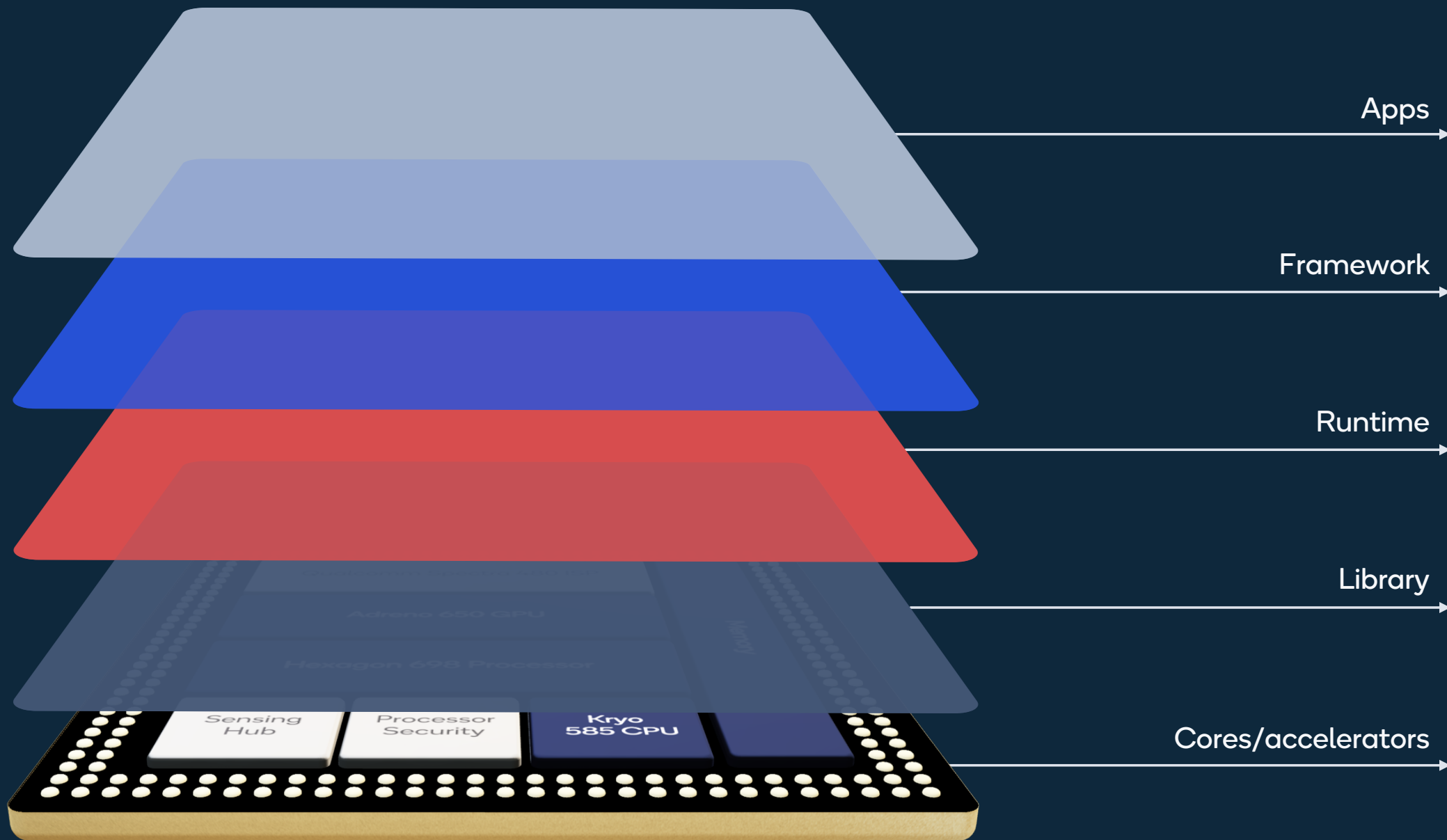
Vice President, Technology, AI Software

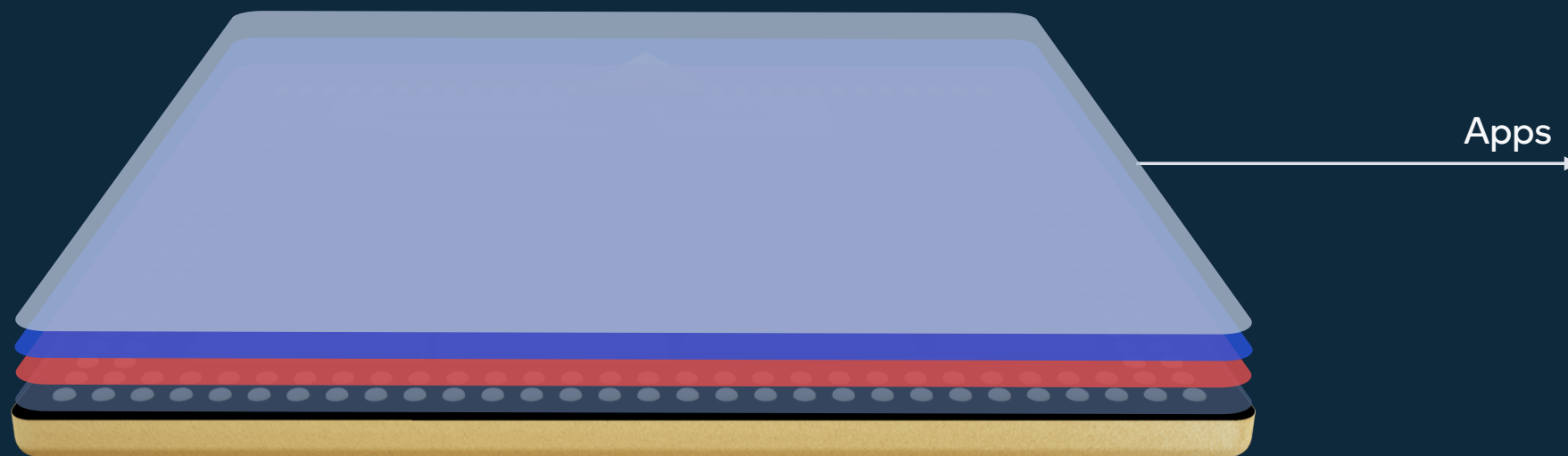
Qualcomm Technologies, Inc.

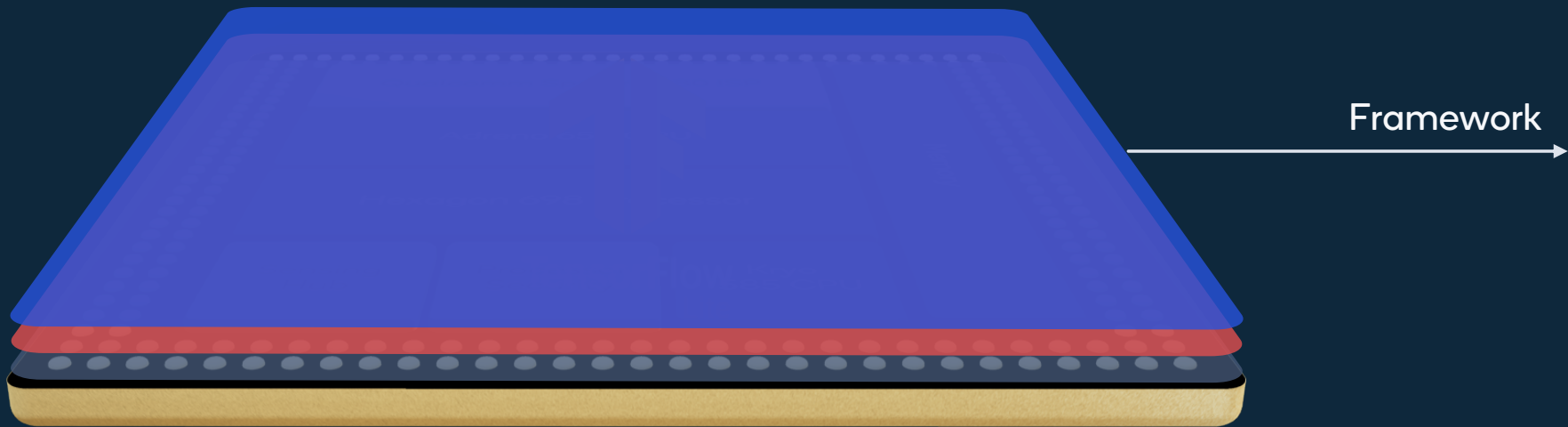
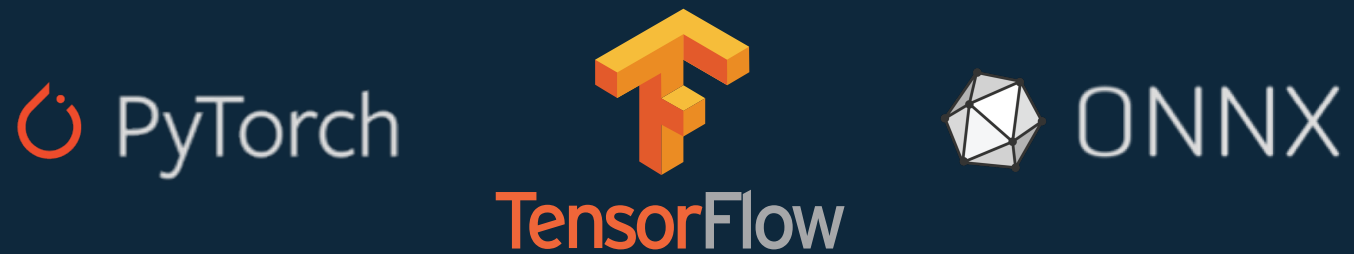
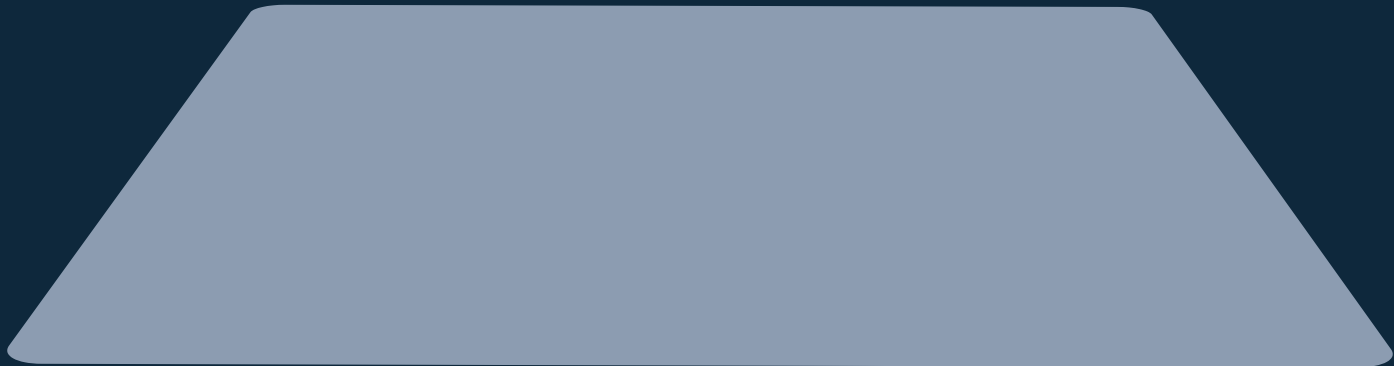














Qualcomm  
Neural Processing SDK



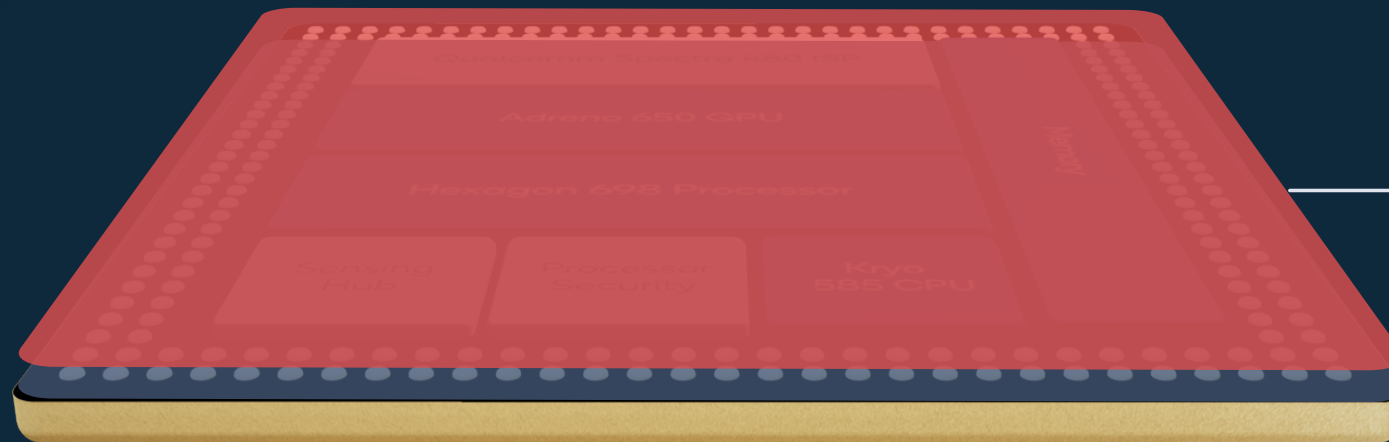
Runtime



Qualcomm  
Neural Processing SDK



Android Neural  
Networks API



Runtime





## Qualcomm Neural Processing SDK

Improved  
performance

Data free  
quantization

Support for more  
network models

Graph analysis

**3-5x**

Speed improvements  
on Hexagon  
processor



## Android Neural Networks API

3x number of  
accelerated  
operators

Shipping  
with  
Android 10

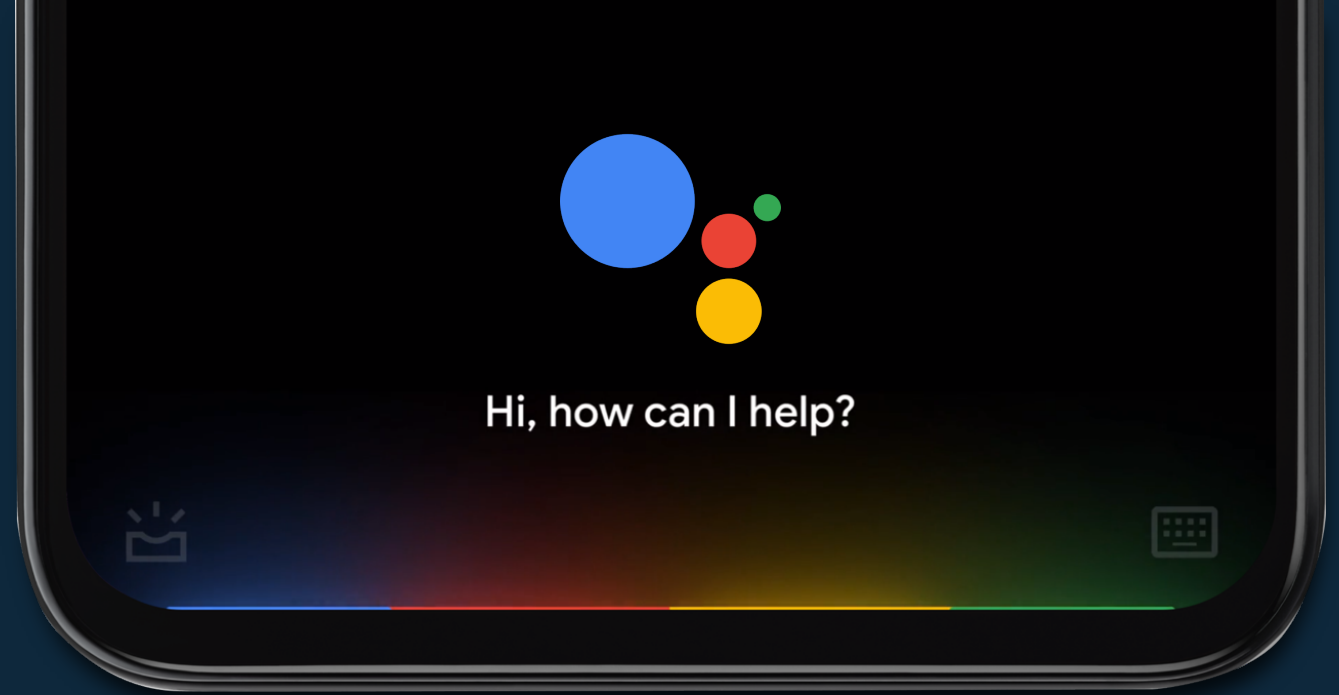
Optimized for  
Google speech  
recognition and  
Google Lens

Wide adoption  
from developers



Android Neural  
Networks API

# Moving ASR for Google Assistant from CPU to Hexagon processor



**3x**  
Power  
savings

**30%**  
Lower  
latency

Current operator count

160+

operators

User defined operators



operators

OpenCL

Hexagon  
SDK

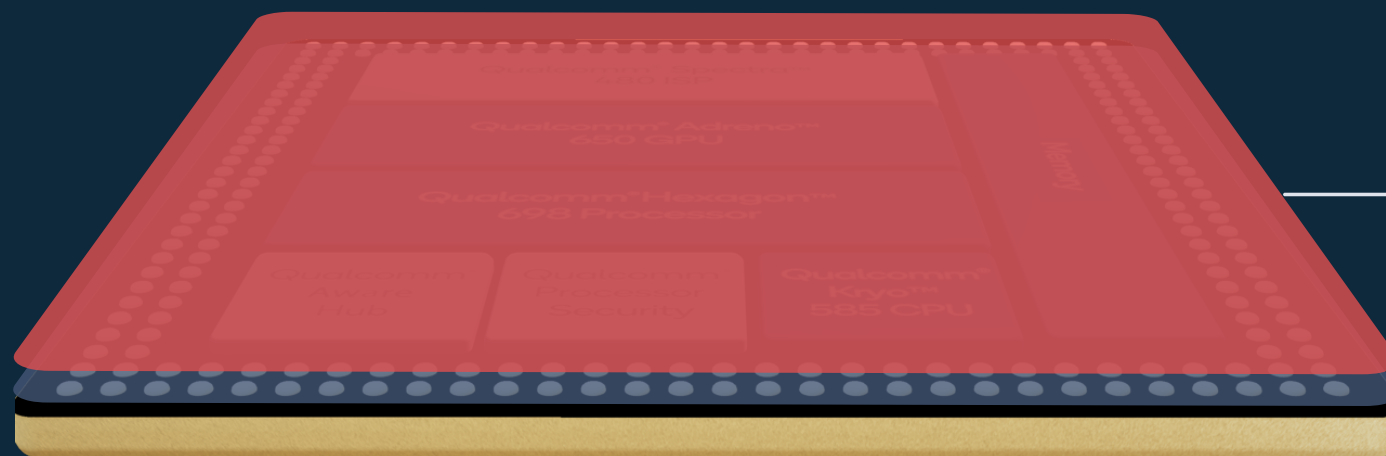




Qualcomm  
Neural Processing SDK



Android Neural  
Networks API



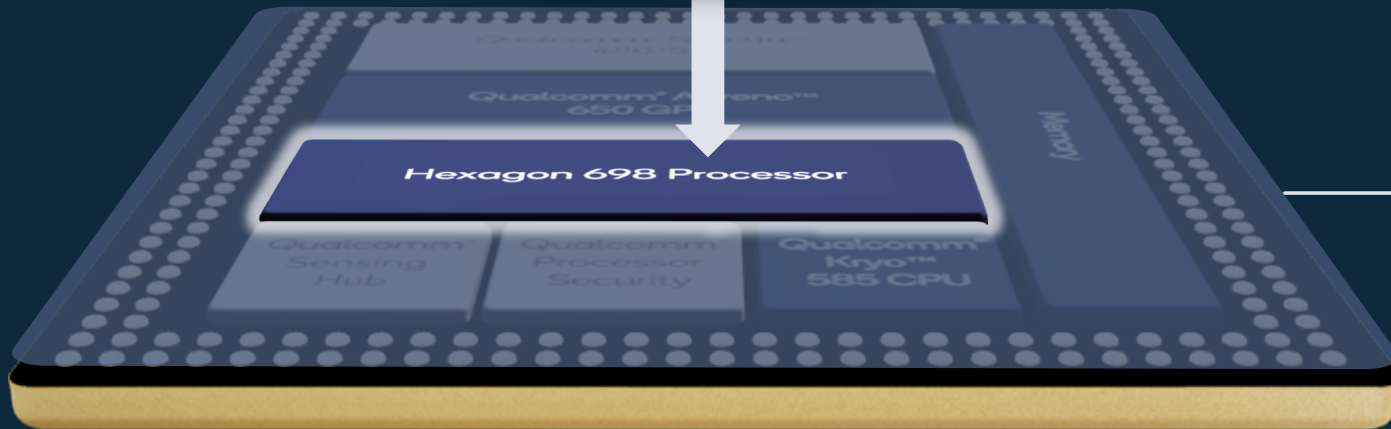
Runtime





TensorFlow  
Lite

Hexagon  
NN Direct



Library



Hexagon  
NN Direct



Library →

Hardware

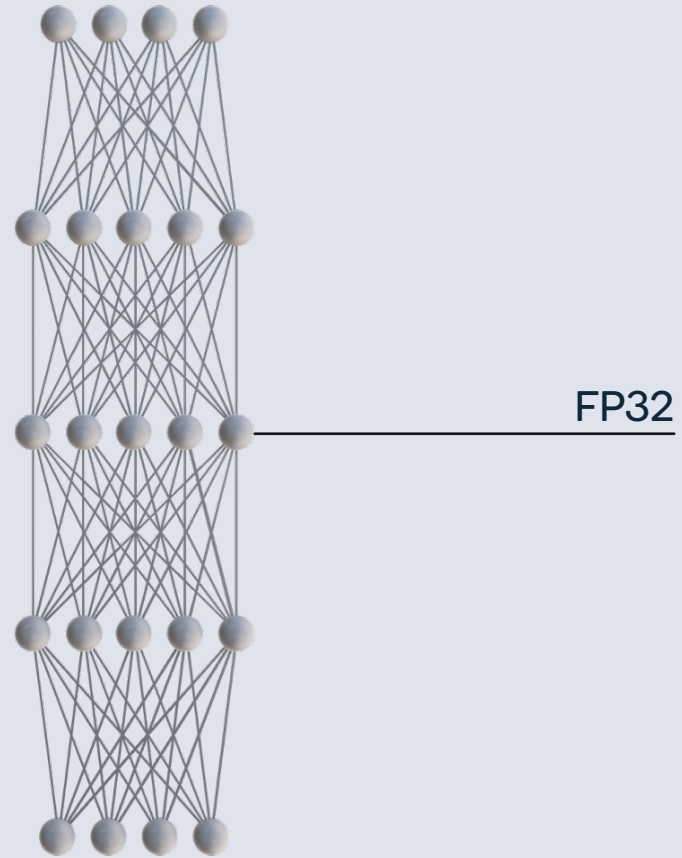
---

# Developer tools

---

Software

# New Qualcomm® AI Model Efficiency Toolkit



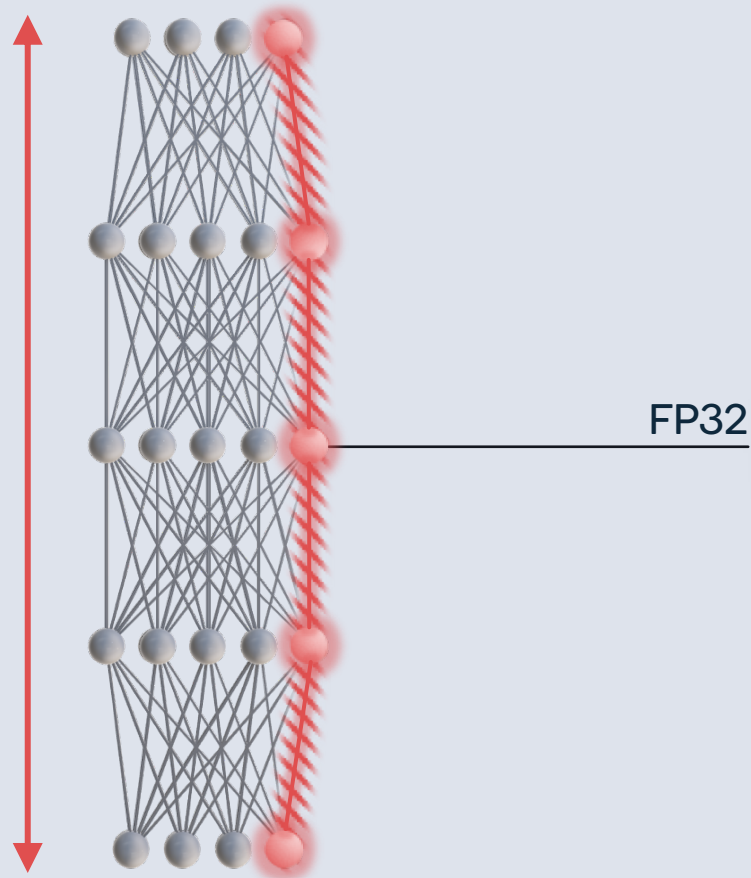


## New Qualcomm® AI Model Efficiency Toolkit

### Model compression

Spatial SVD

Bayesian compression



**3x**  
Compression with  
less than 1% loss  
in accuracy\*



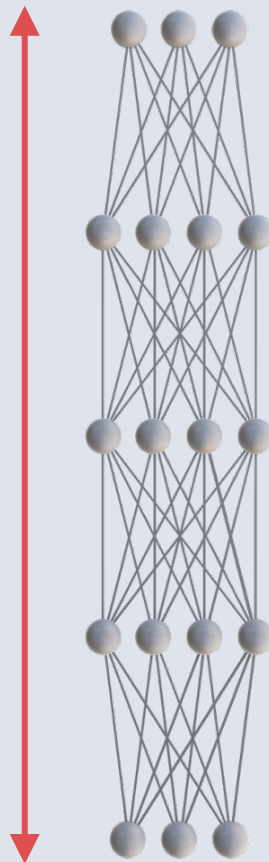
\*Comparison between baseline and compression with both Bayesian compression and spatial SVD. Example uses ResNet18 as baseline.

New Qualcomm® AI  
Model Efficiency Toolkit

Model compression

Data free quantization

Quantization aware training



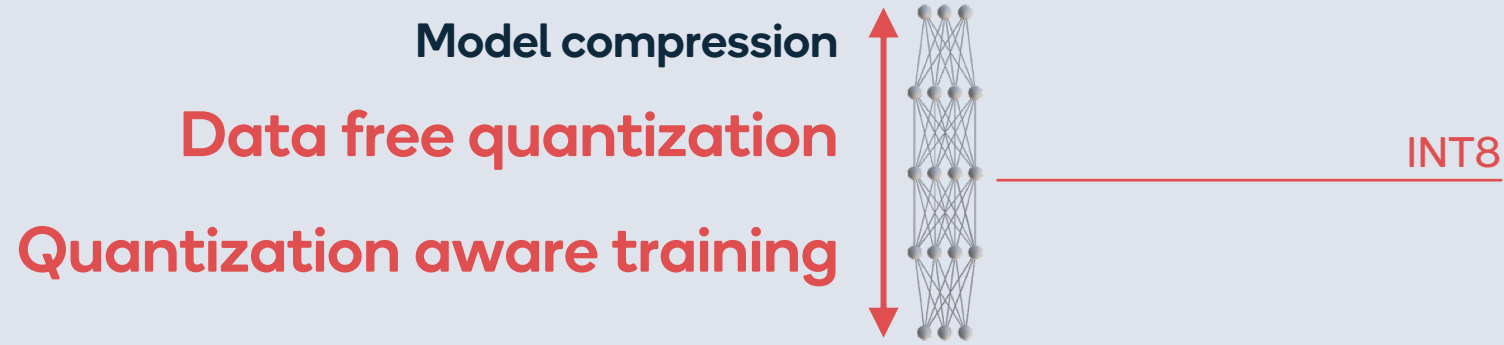
FP32

**3x**  
Compression with  
less than 1% loss  
in accuracy\*



\*Comparison between baseline and compression with both Bayesian compression and spatial SVD. Example uses ResNet18 as baseline.

## New Qualcomm® AI Model Efficiency Toolkit



**>4x**

Increase in  
performance  
per watt with  
quantization

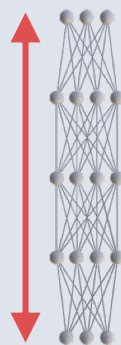


3x

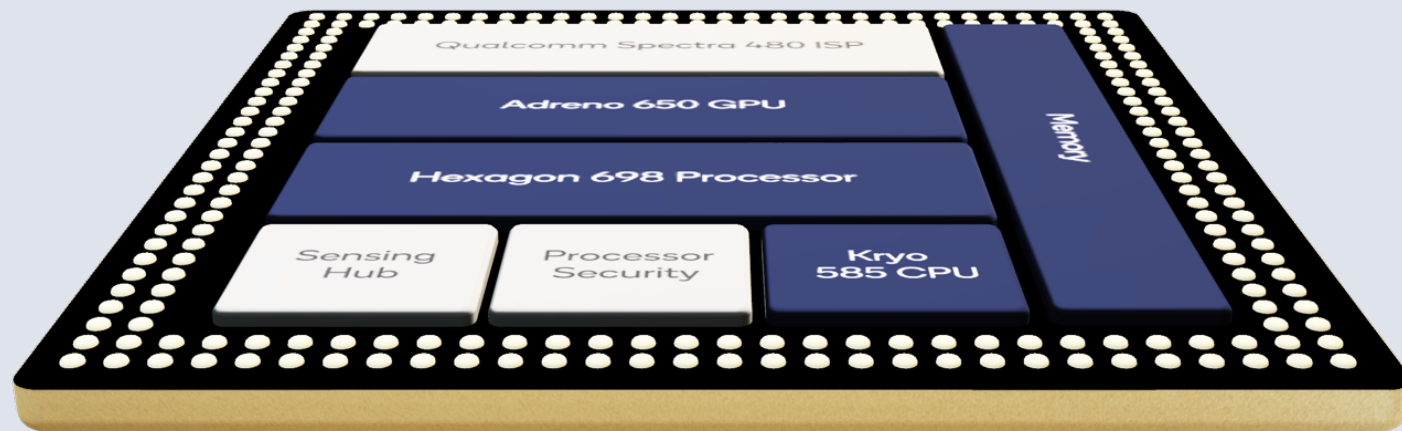
Compression with  
less than 1% loss  
in accuracy\*

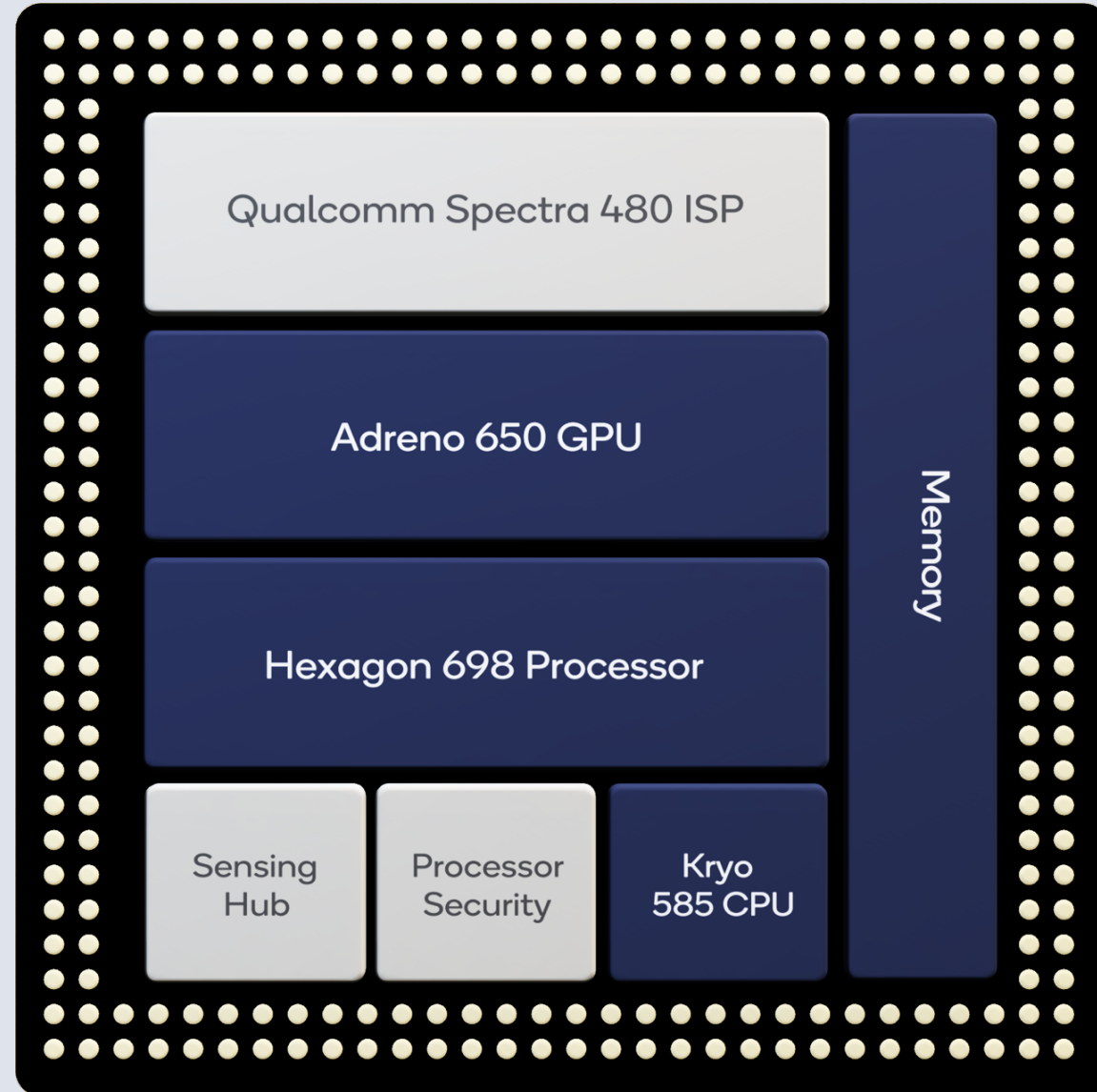
>4x

Increase in  
performance  
per watt with  
quantization



INT8







戲

文

很

早

一

興

反

反

謝

的

嗎

人

身

報

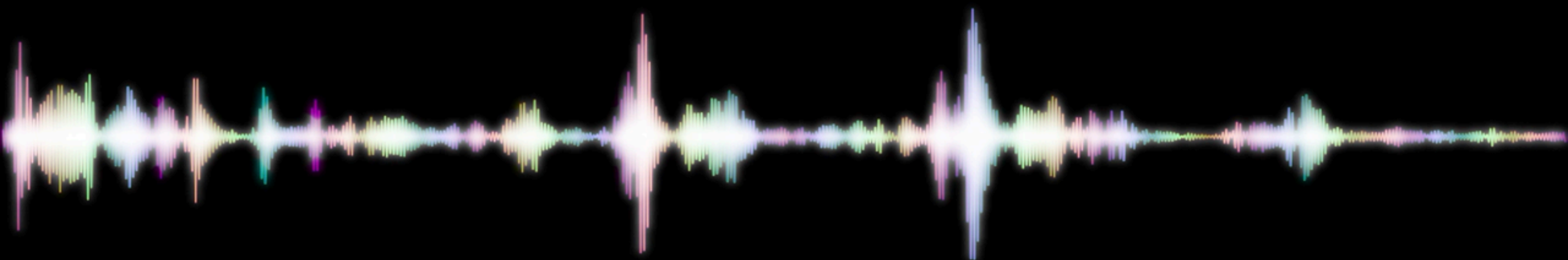
海

文

南



**Snapdragon 865 is a leading mobile AI platform**

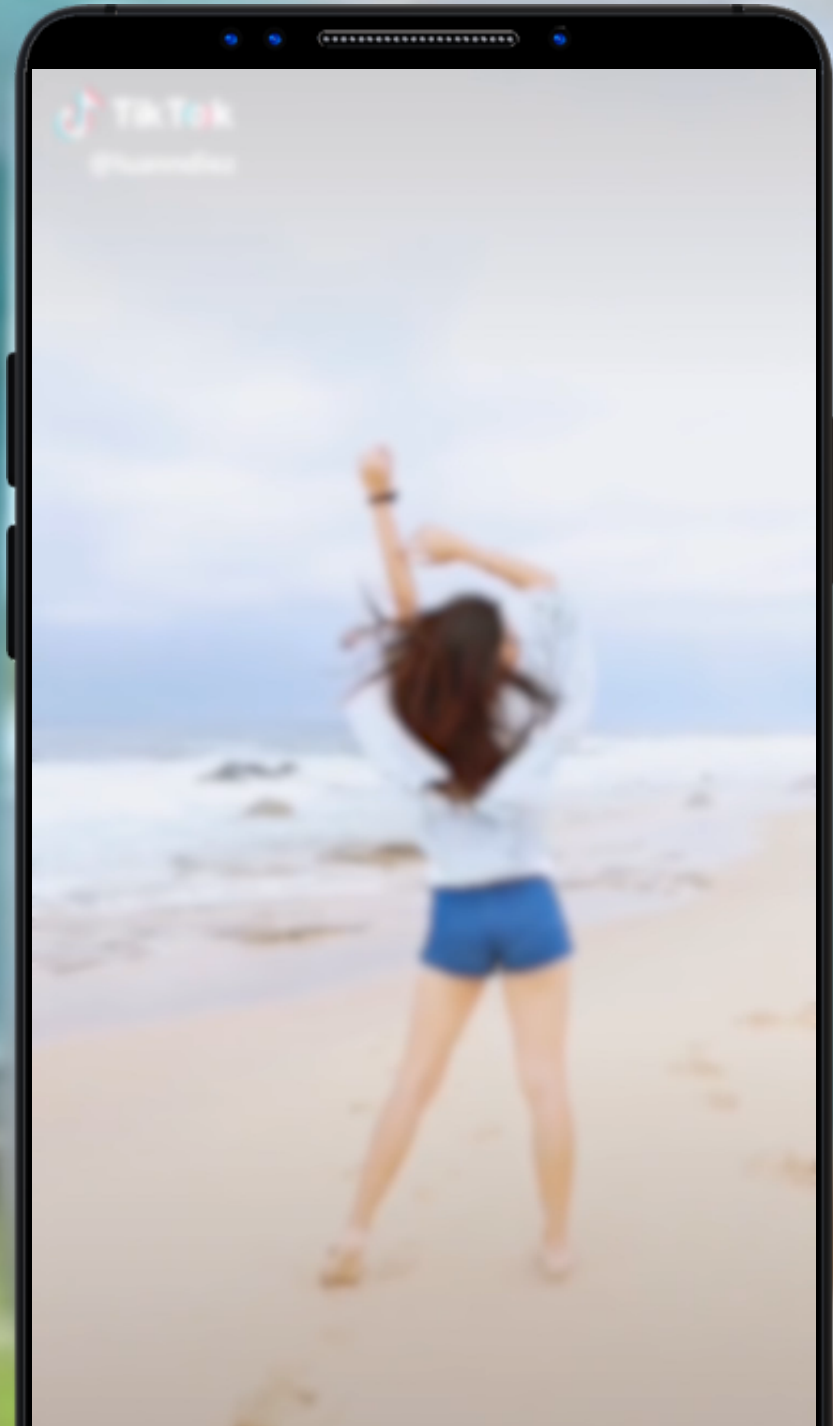










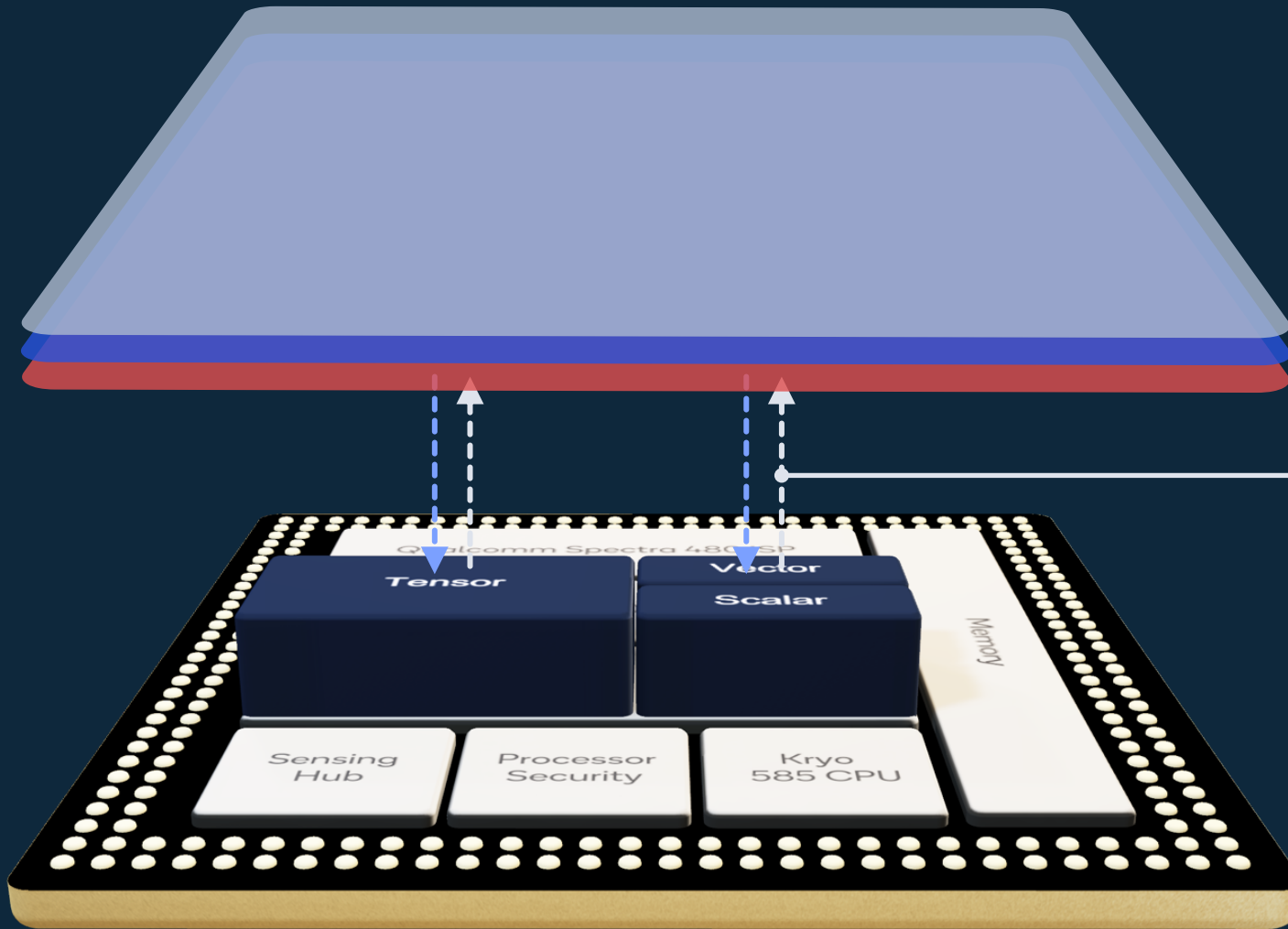


 ByteDance





# ByteDance



Qualcomm  
Neural Processing SDK

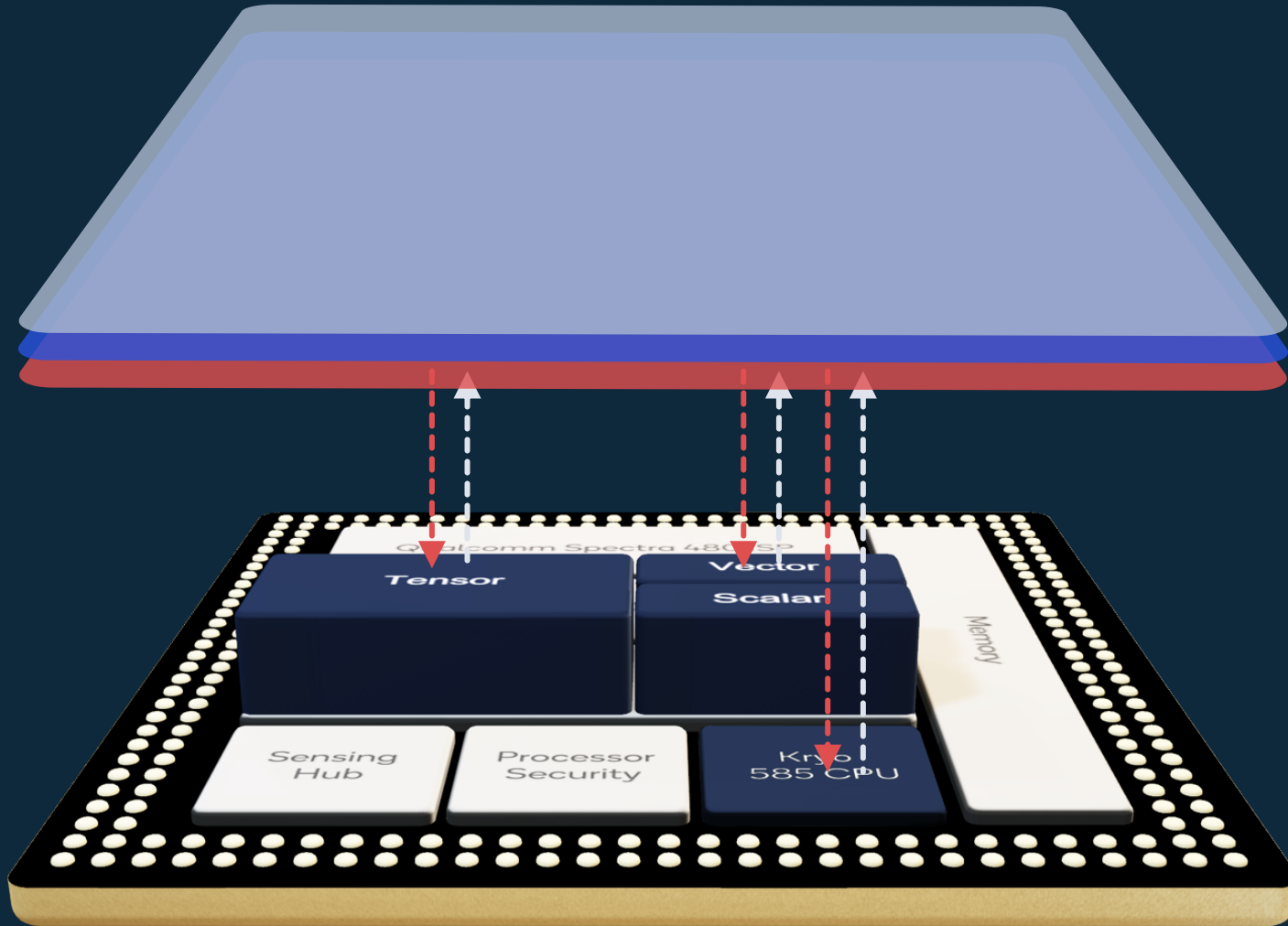
Hexagon  
NN Direct

**Super resolution—**  
Convolutional Neural Network (CNN)

loom.ai



# loom.ai



Qualcomm  
Neural Processing SDK

**Face detection –**  
Convolutional Neural Network (CNN)

**Landmark detection –**  
Convolutional Neural Network (CNN)

**Expression detection –**  
Convolutional Neural Network (CNN)

**Replace face with avatar –**  
Convolutional Neural Network (CNN)



## Gaming



## XR



## Contextual Awareness



## Adreno 650

New AI mixed  
precision  
instructions

2x higher TOPS  
16-bit and 32-bit FP

## Hexagon 698

4x higher TOPS  
Up to 35%  
power savings

Deep learning  
bandwidth  
compression

## LP-DDR5 Memory

30% more bandwidth

## 5th gen Qualcomm AI engine

15 TOPS

# AI highlights



## Qualcomm Neural Processing SDK

New features  
and improvements

## NNAPI support

3x number of  
accelerated operators  
Optimized for Google ASR and  
Google Lens

## Hexagon NN Direct

Provide developers direct  
access to Hexagon

## AI model efficiency toolkit

Data free  
quantization

Quantization  
aware training

Model  
compression





# Thank you

Follow us on:    

For more information, visit us at:

[www.qualcomm.com](http://www.qualcomm.com) & [www.qualcomm.com/blog](http://www.qualcomm.com/blog)

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

©2018-2019 Qualcomm Technologies, Inc. and/or its affiliated companies. All Rights Reserved.

Qualcomm, Snapdragon, Qualcomm Spectra, Adreno, Hexagon, Kryo and Qualcomm Aqstic are trademarks of Qualcomm Incorporated, registered in the United States and other countries. Quick Charge and FastConnect are trademarks 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 Qualcomm’s licensing business, QTL, and the vast majority of its patent portfolio. Qualcomm Technologies, Inc., a wholly-owned subsidiary of Qualcomm Incorporated, operates, along with its subsidiaries, substantially all of Qualcomm’s engineering, research and development functions, and substantially all of its product and services businesses, including its semiconductor business, QCT.