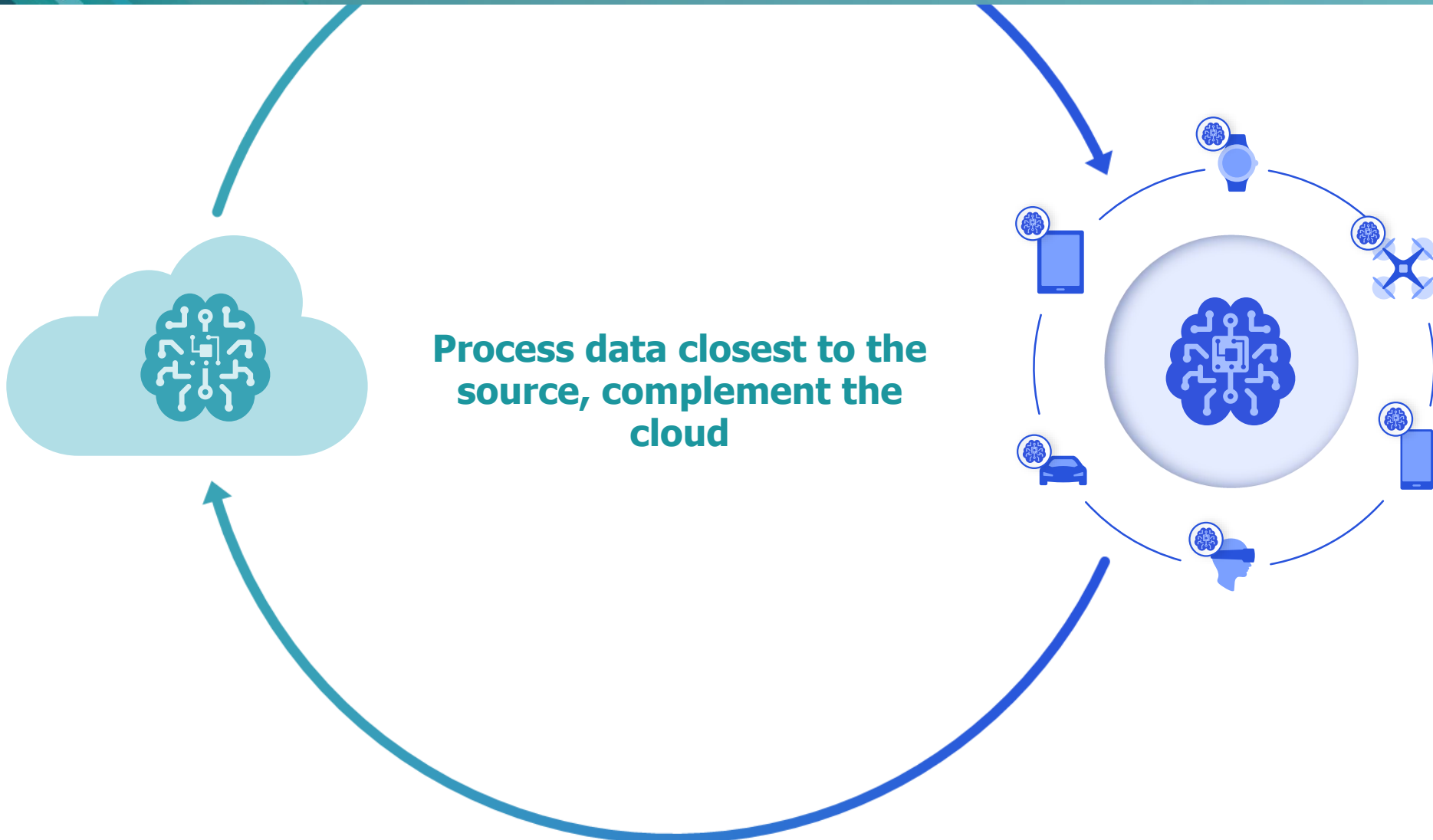




# **The Future of AI is Here Today: Deep Dive into Qualcomm's On-Device AI Offering**

Vinesh Sukumar  
Sr. Director, Product Management (AI/ML)  
Qualcomm Technologies, Inc.

# Center of Gravity Moving to the Edge...



Process data closest to the source, complement the cloud

## Historically

Privacy

---

Reliability

---

Low latency

---

Efficient use of network bandwidth

## Increased Demand

Autonomy

---

Personalization

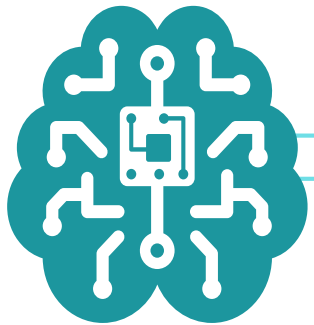
---

Efficiency

---

Security

# At Qualcomm – AI Deployed Across Various Technologies & Verticals



**MORE THAN  
A DECADE  
OF AI R&D**

**Enhance**

**TECHNOLOGIES**

- Camera
- Video
- Voice
- Audio
- AR/VR
- Modem

**Superior solutions**

**Create**

**QUALCOMM® AI ENGINE**

- Application layer
- Software layer
- Hardware layer

**Enabling developer ecosystems**

**BU VERTICALS**

## Windows 11 – PC Computing

Devices will need a neural processing unit (NPU) to use these new Windows 11 features, which means they'll show up first on Lenovo's new ThinkPad X13s, which is powered by Qualcomm's Snapdragon 8cx Gen 3 compute platform.

## Smart Camera – IoT Markets

### Qualcomm debuts smart camera processor at ISC West

Qualcomm showcased a smart camera processor named QCS7230 at the International Security Conference and Exposition (ISC West) that expands its Vision Intelligence Platform portfolio to cities, enterprises, and public spaces.

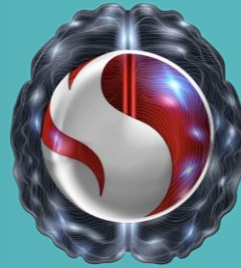
## ADAS Markets

Ferrari and Qualcomm team up for tech projects for road, racing cars



Snapdragon, Qualcomm QCS7230, and Qualcomm AI Engine are products of Qualcomm Technologies, Inc. and/or its subsidiaries.

# AI Applications



Snapdragon  
smart

Qualcomm

# AI Applications: Across Various Segments

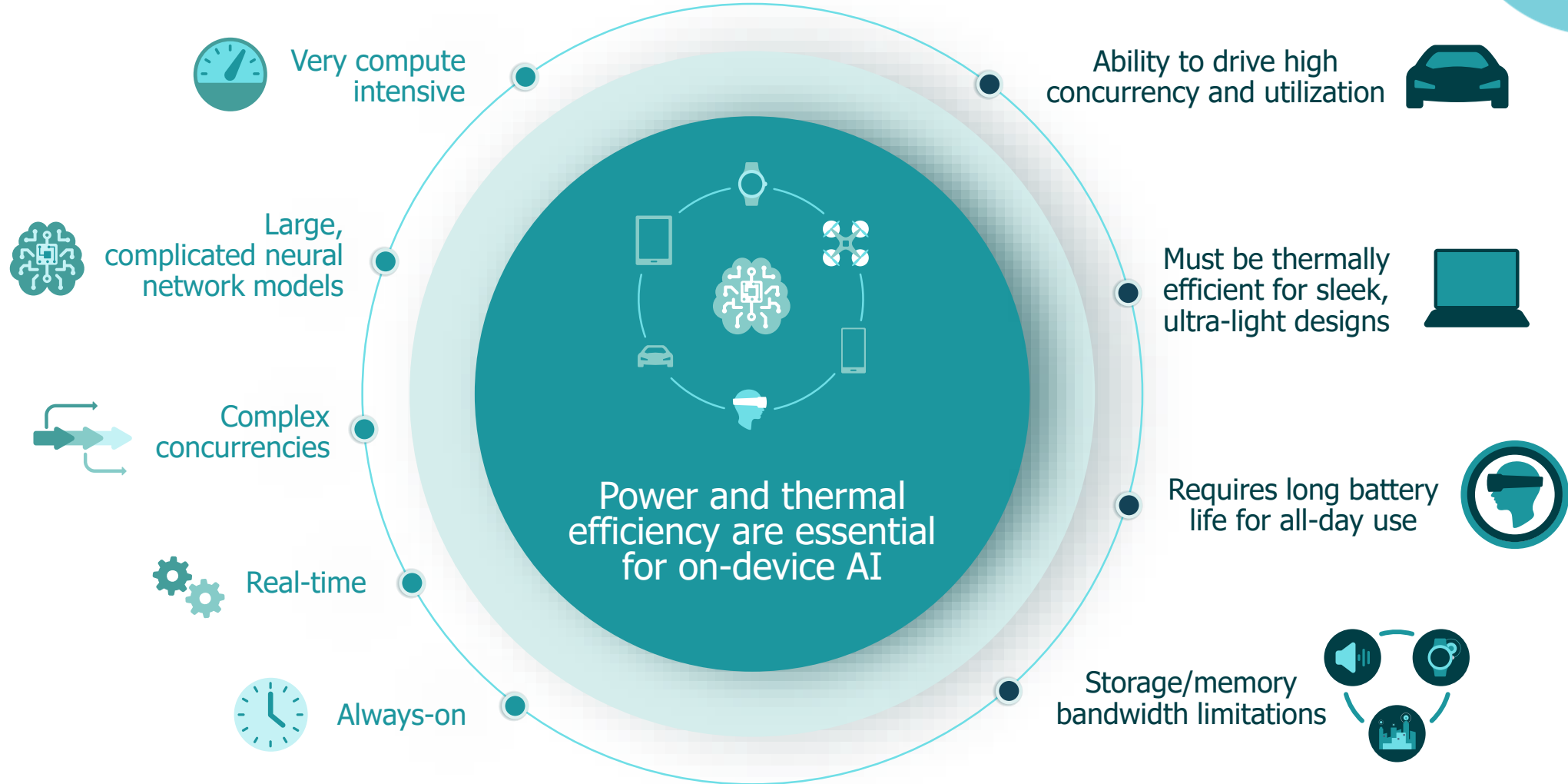
Expanding beyond modalities of computer vision to linguistics, communication, commerce and language understanding



Mobile	IoT	Compute	Cloud	Auto
<p><b>AI Assisted Imaging</b></p> <ul style="list-style-type: none"> <li>AI 3A</li> <li>Scene-based Camera Selection</li> </ul> <p><b>Image Understanding</b></p> <ul style="list-style-type: none"> <li>Face Detection / Tracking / Features</li> <li>Object Detection / Tracking</li> <li>Body Detection / Tracking / Pose</li> <li>Human Segmentation</li> <li>Sky Segmentation</li> <li>Multi-Class Segmentation</li> <li>Depth Estimation</li> </ul> <p><b>Beautify / Augment</b></p> <ul style="list-style-type: none"> <li>Scene-based Image Enhancement</li> </ul> <p><b>Image Processing</b></p> <ul style="list-style-type: none"> <li>AI based NR or Image SR</li> <li>Scene-based Camera Selection</li> </ul> <p><b>Audio</b></p> <ul style="list-style-type: none"> <li>Real time language</li> <li>Natural language processing (NLP)</li> </ul> <p><b>Modem</b></p> <ul style="list-style-type: none"> <li>Parameter optimization</li> <li>Robust sequence predictions</li> </ul>	<p><b>Robotics</b></p> <ul style="list-style-type: none"> <li>Autonomous navigation</li> <li>Obstacle Avoidance</li> <li>Picking and Sorting</li> </ul>	<p><b>Productivity</b></p> <ul style="list-style-type: none"> <li>Background based noise cancellation on Audio (inbound and outbound)</li> <li>Segmentation/Blur/Super Resolution on Video</li> <li>Voice activation without keywords</li> <li>Face tracking</li> <li>Smart photo categorization</li> </ul>	<p><b>Data Centers</b></p> <ul style="list-style-type: none"> <li>Natural language processing</li> <li>Computer vision</li> <li>Recommendation system</li> </ul>	<p><b>IVI</b></p> <ul style="list-style-type: none"> <li>Occupancy monitoring system (OMS)</li> <li>Driver monitoring system (DMS)</li> <li>Surround perception</li> <li>Audio Command &amp; Control</li> </ul>
	<p><b>Retail</b></p> <ul style="list-style-type: none"> <li>Visitor/Face/Gesture Recognition</li> <li>Object/People Detection and Counting</li> <li>Barcode decoding</li> <li>Empty shelf detection</li> <li>Dwell time</li> </ul>	<p><b>Privacy &amp; Security</b></p> <ul style="list-style-type: none"> <li>Automatic screen unlock and login</li> <li>Privacy alert</li> <li>Guard mode</li> </ul>		
	<p><b>Transportation</b></p> <ul style="list-style-type: none"> <li>License plate recognition</li> <li>Face and facial landmark detection</li> <li>Drowsiness detection</li> </ul>	<p><b>Content Creation &amp; Gaming</b></p> <ul style="list-style-type: none"> <li>Gaming with gesture control</li> <li>Gaming with voice commands</li> <li>Intelligent highlight videos</li> <li>Game play improvement</li> </ul>	<p><b>ADAS (Up to L4)</b></p> <ul style="list-style-type: none"> <li>Highway driving assist                             <ul style="list-style-type: none"> <li>Front collision warning</li> <li>lane departure,</li> <li>Traffic jam assist</li> <li>Auto lane change</li> <li>Auto lane merge</li> <li>Traffic light recognition</li> <li>Construction zones</li> <li>Urban autonomous driving</li> </ul> </li> <li>Parking assist                             <ul style="list-style-type: none"> <li>Person detection,</li> <li>Perception</li> <li>Valet parking</li> </ul> </li> <li>Driver monitoring</li> </ul>	
	<p><b>Smart Devices</b></p> <ul style="list-style-type: none"> <li>Object/People detection</li> <li>Speaker detection</li> <li>Gun shot detection</li> </ul>	<p><b>Performance &amp; Efficiency</b></p> <ul style="list-style-type: none"> <li>Power and Screen optimization</li> </ul>		
	<p><b>Smart Buildings</b></p> <ul style="list-style-type: none"> <li>People Tracking</li> <li>Access Control</li> </ul>			
		<p><b>Manufacturing/Logistics</b></p> <ul style="list-style-type: none"> <li>Predictive maintenance</li> <li>Energy management with Asset demand</li> </ul>		

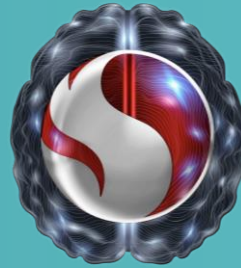


# Challenges of AI Applications



Lays the foundation for Qualcomm AI Hardware

# AI Hardware



Snapdragon  
smart

Qualcomm

# Vision: Drive Leadership Capability Across All Markets

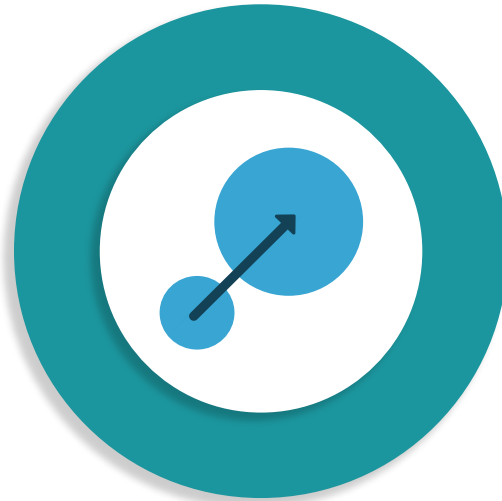


## Performance



Invest in performance (Inf/Sec) and Power efficiency (Inf/Sec/W)

## Scalability



Leverage existing AI HW engines to scale across various TDP points

## Innovation



Feature innovation to drive leadership (Datatypes, Sparsity, Streaming...)

## Co-Design



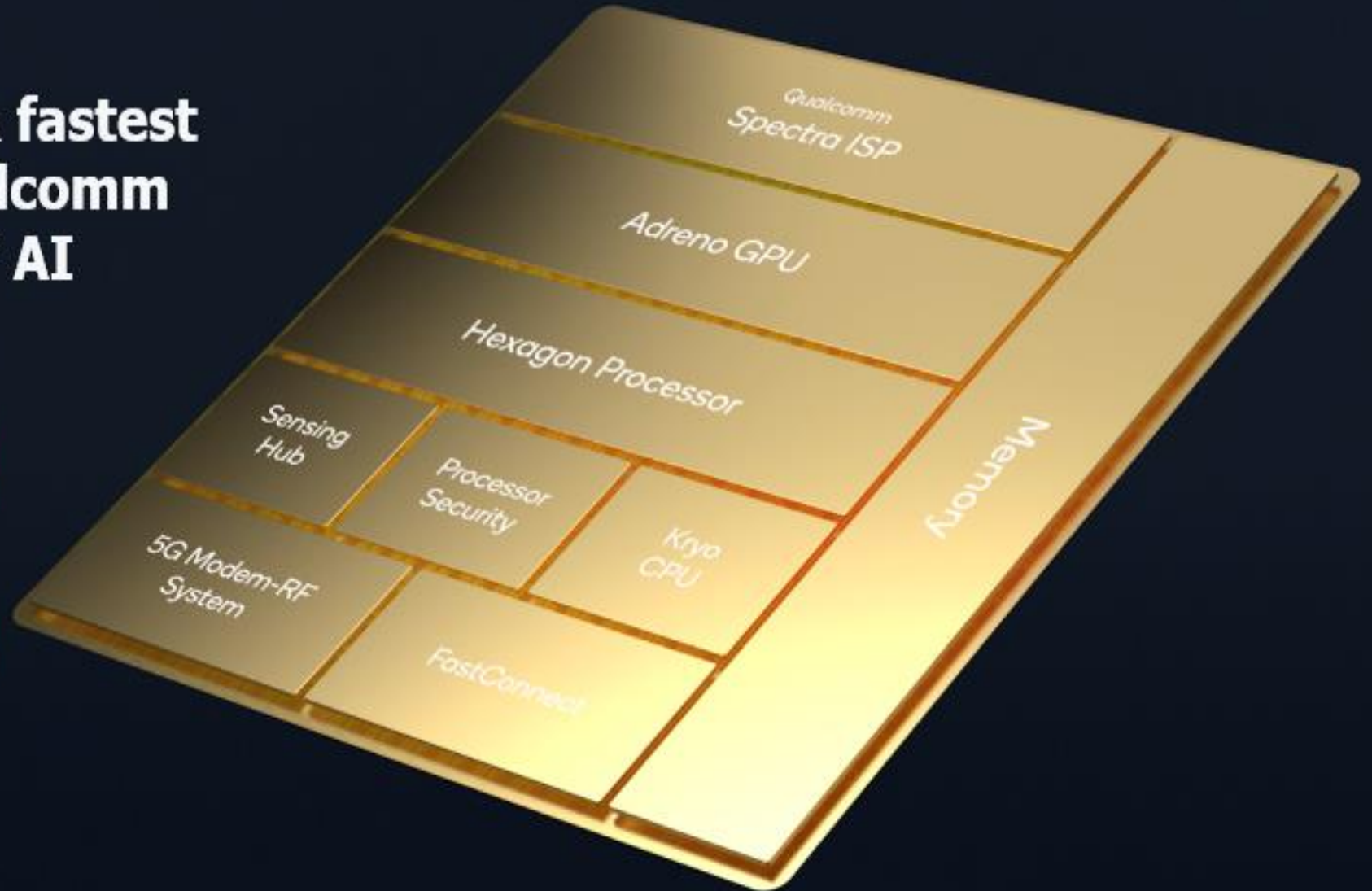
HW & SW co-design to make programmability easier (NAS, Compatibility)



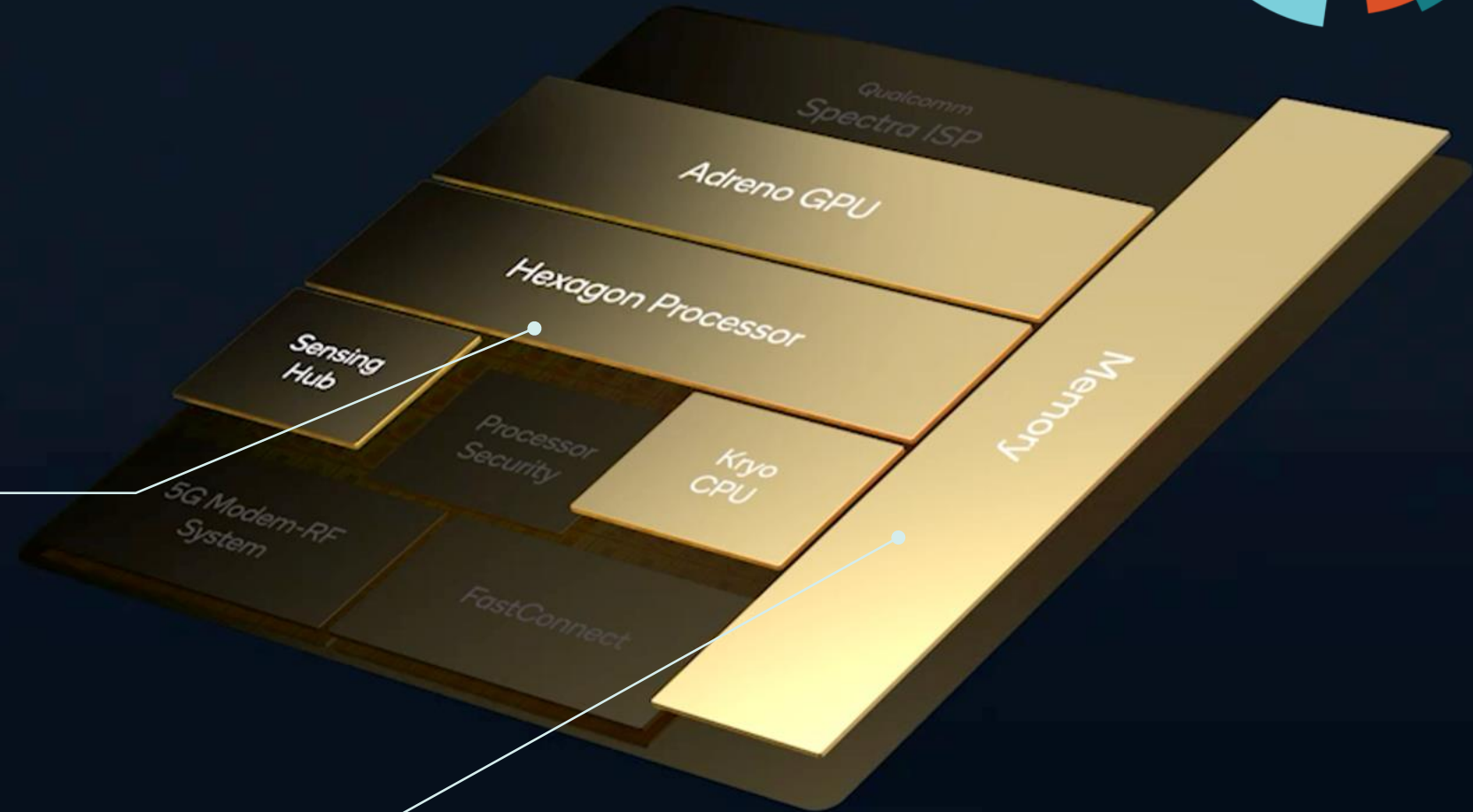
# 7<sup>th</sup> Generation : Introducing Qualcomm AI Engine for All Verticals



**The most powerful & fastest  
AI Engine from Qualcomm  
to run "Modern" AI**



# 7th Generation : Qualcomm AI Engine



**2X**

Computational performance

**2X**

Large-shared memory

# Scaling – Dedicated Investments in AI HW Engines



Using Mobile Design as an Anchor Point



Investment in sparsity modules for supporting Auto ADAS usages



Optimize design for sustained low power with the ability to enable high concurrency



Investment in dedicated datatypes for higher performance and TCO (Total Cost of Ownership)

# Performance Leadership – Using MLPerf™ 2.0



## Qualcomm® 8C GEN 1 Mobile



**Xiaomi Mi12 Platform**

Scaling from Mobile to Cloud..

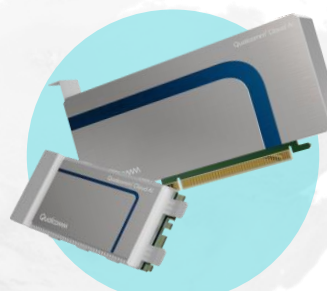
RESNET50  
2221  
Inf/sec

MOSAIC  
752  
Inf/sec

BERT  
101  
Inf/sec

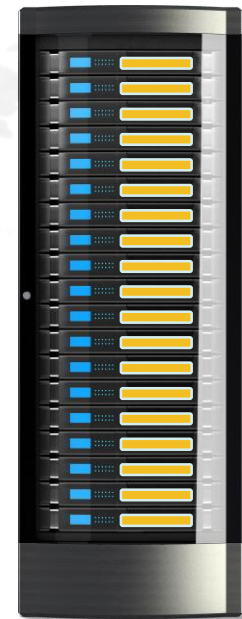
8C GEN 1 achieves an average of about 26% better latency than Exynos devices across various categories

## Qualcomm® Cloud AI 100



16  
75W cards  
in one server  
RESNET50  
371,473  
Inf/sec

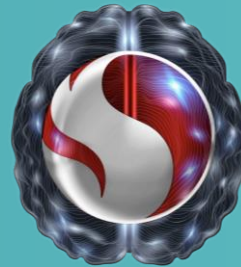
20 server  
units in one  
server rack  
7.4+M  
Inf/sec



**Gigabyte Platform**

Cloud AI 100 servers achieve 3.7x higher rack-level ResNet-50 inference performance than Nvidia A100 servers

# AI Software



Snapdragon  
smart

Qualcomm



# Vision: Accelerate AI Innovation and Solution Deployment

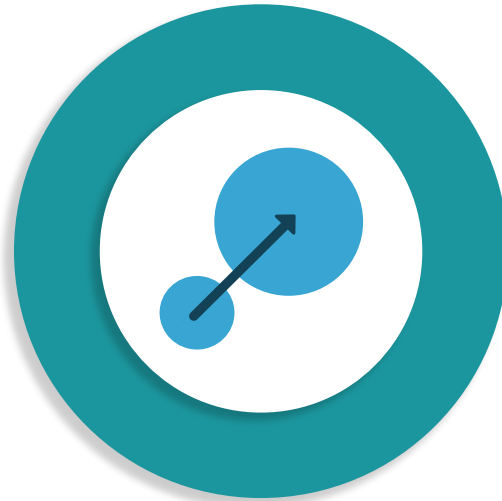


## Performance



Accelerate “out of box” operator functionality and performance

## Scalability



Ability to have programming consistency from Cloud to Edge

## Tools



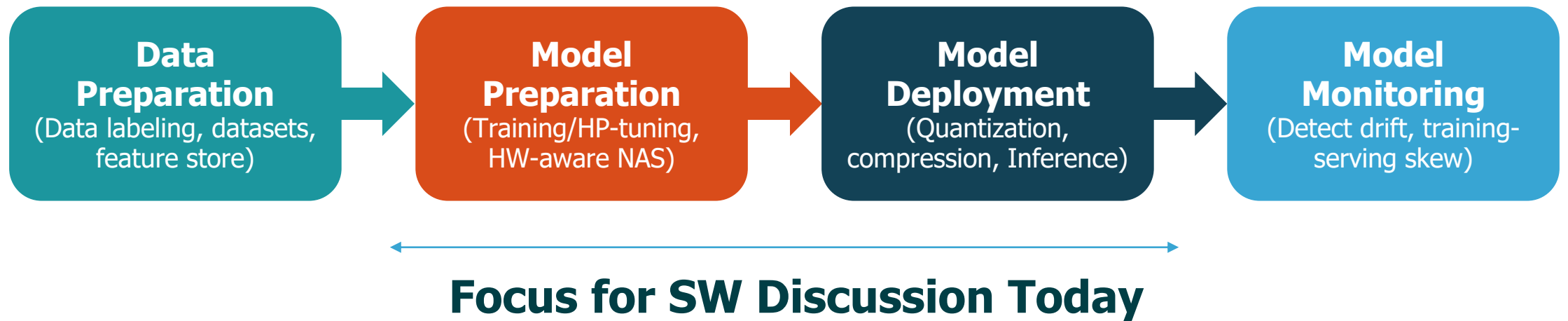
Accelerate AI Solution deployment with investment in Tools

## Innovation

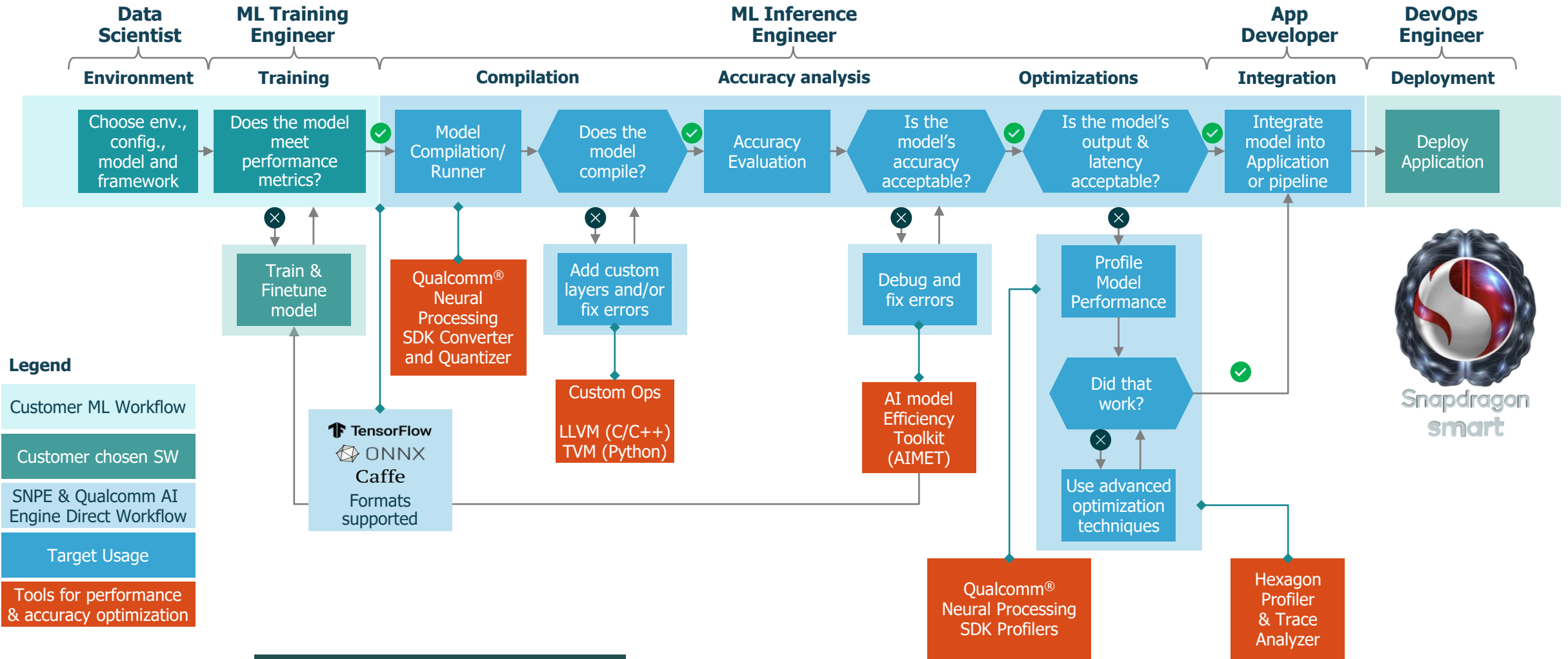


Innovation to drive product leadership (Pre-emption, DFS, Multi Chaining)

# Qualcomm AI SW MLOps Cycle (Inception to Monitoring)



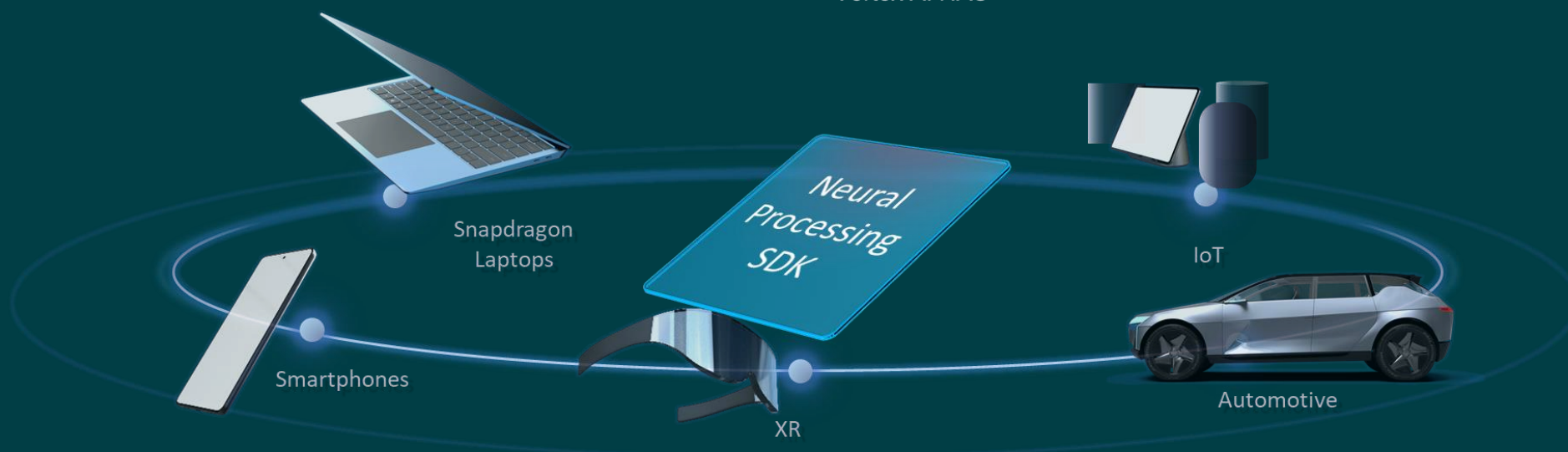
# AI Software Workflow





# STEP : 1 | NAS : Model Optimization Mapped to HW Intrinsic

## Partnering with Google Vertex AI Team - Integrated into Qualcomm Software Stack



### How →

#### Search Space

Space of allowable architectures (Structure, operations, connectivity)

#### Search Algorithm

Sampling populations of good architecture candidates

#### Evaluation Strategy

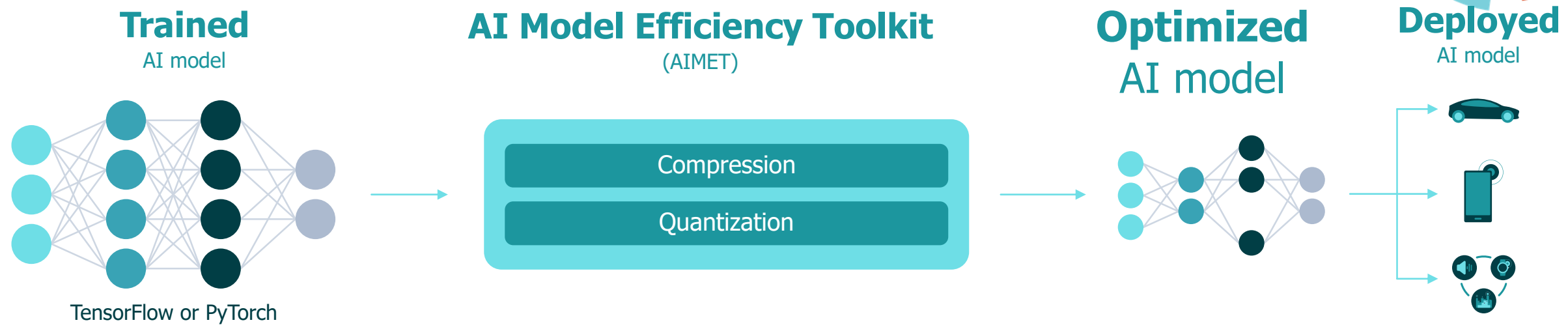
Estimate performance of sampled architecture

### Results →

- 8-10 % Accuracy Improvements
- ~20 to 30% Latency Reduction

Note : DL architecture choices influence results

# STEP : 2 | AIMET : Quantize & Compress AI Models

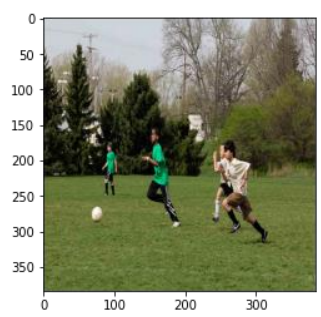


**Why →**

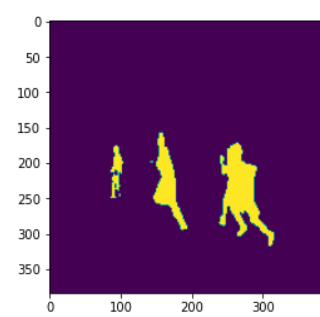
Automated way of enabling reduction in precision of weights and activation while maintaining accuracy

State of the art network compression and quantization tools for various DL architectures (CNN, BERT, GAN's..)

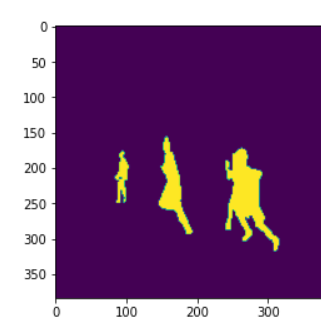
**Results →**



Original Image



FP16 Segmented Map

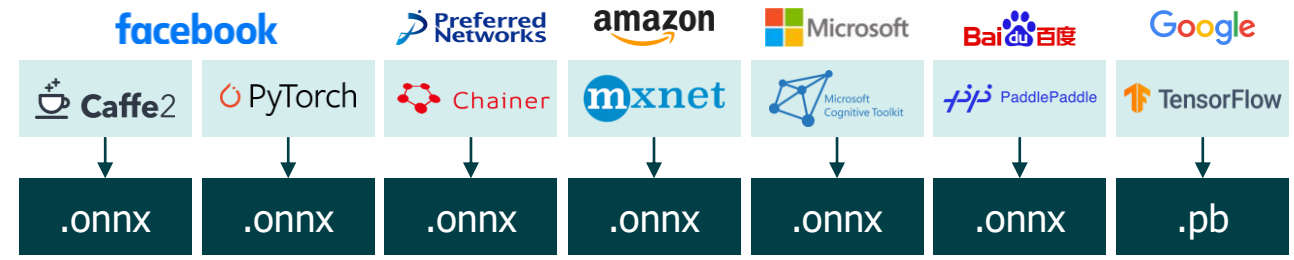


Quantized Segmented Map



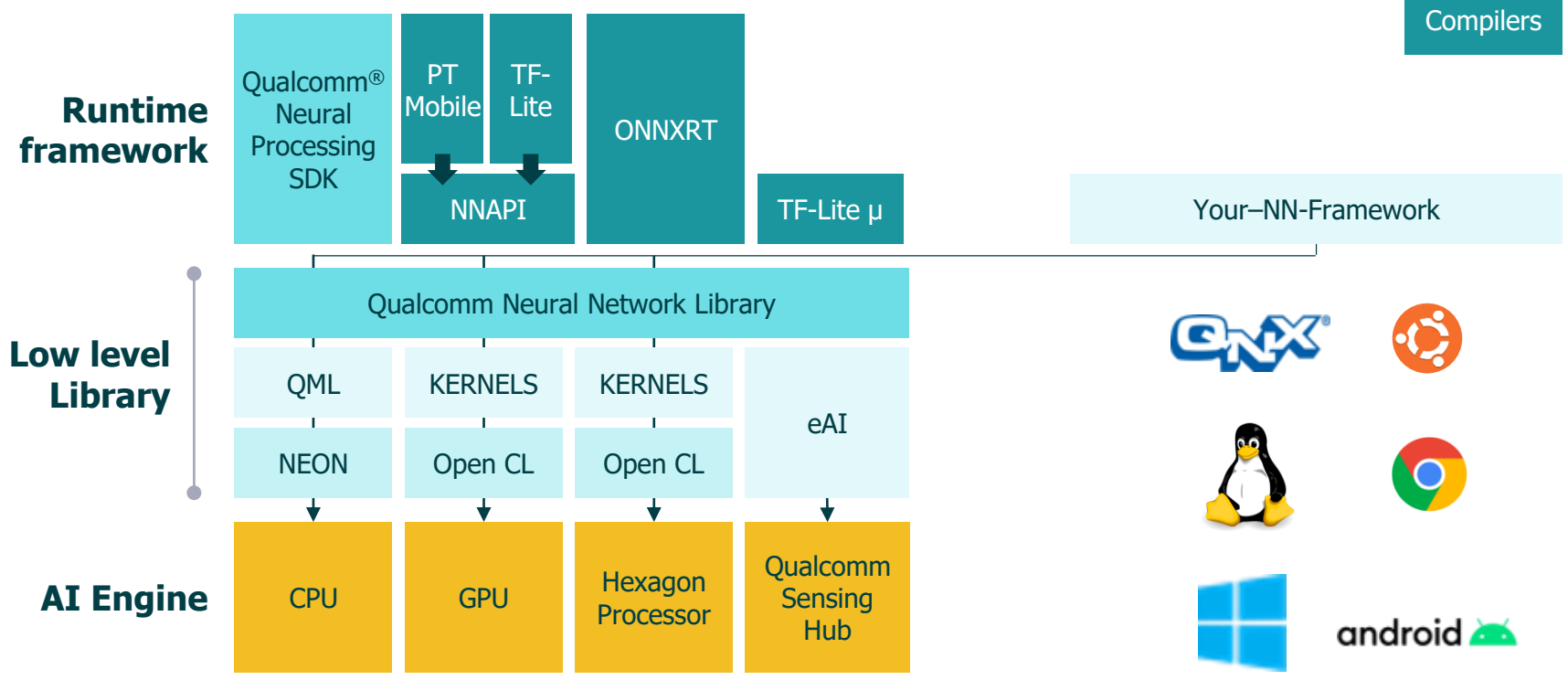
# STEP : 3

# Run Time: Qualcomm AI Software Stack for Performance and Scalability Support - Application Deployment



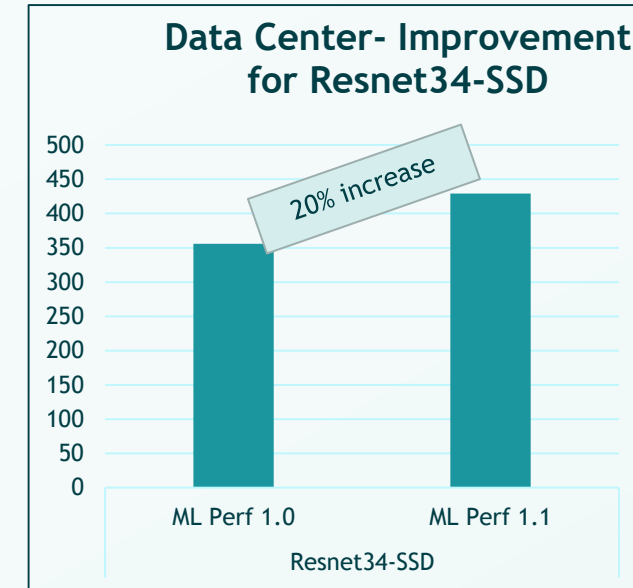
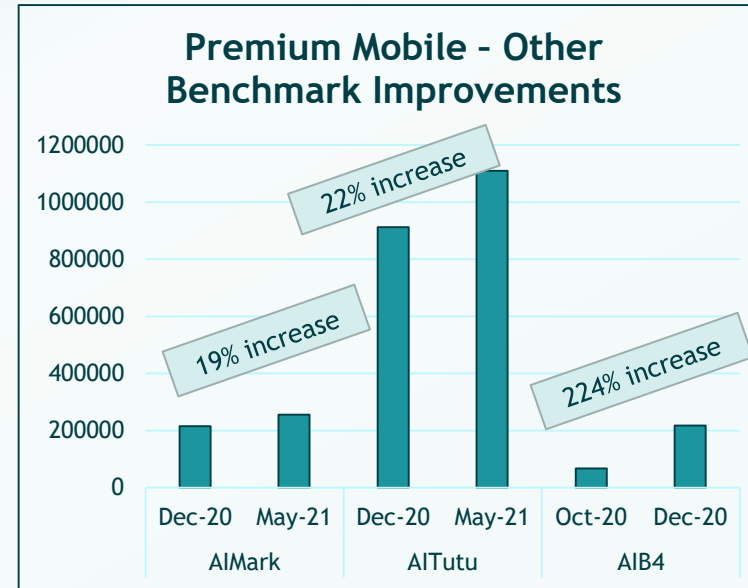
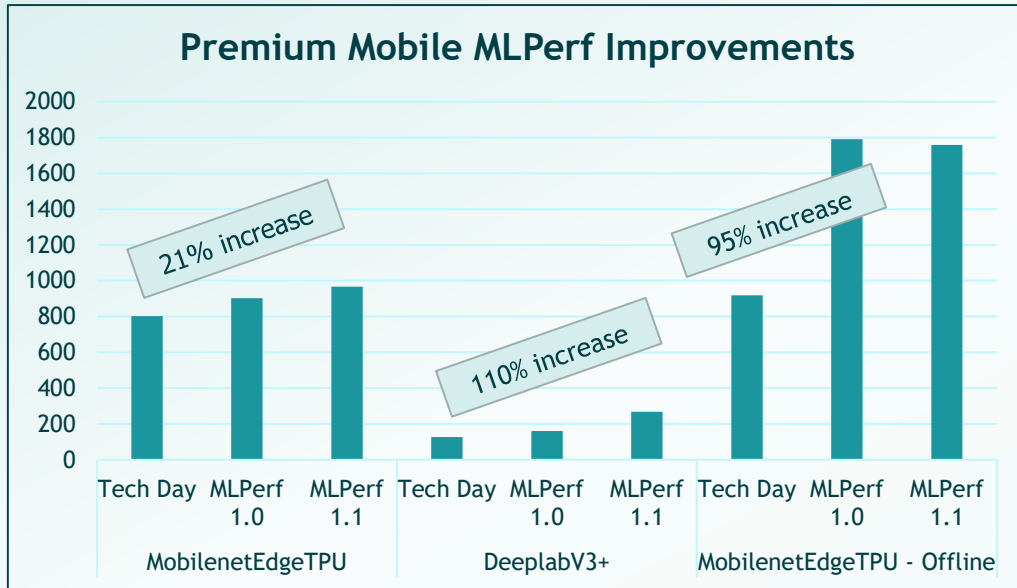
- Profiler
- Debugger
- Visualizer
- Compilers

Performance — Scalability



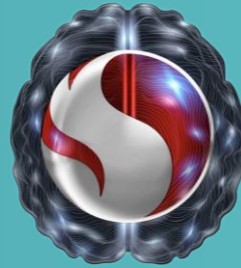
General purpose compute engines — AI engines

# STEP : 4 | Continuous Improvements in SW Stack - For Performance Leadership



- Meaningful improvements after initial chip shipment ; achieved through software investment / innovation
- **Consequences -**
  - Additional gains across the entire SoC portfolio for downstream chips and adjacent BUs from a common investment
  - Platform software that can be delivered incrementally to already released products can continue to improve
  - Improvements directly leveraged into next generation chips
  - Improvements in current generation can lead to relative performance gap compression with next-gen devices

# Conclusions



Snapdragon  
smart

- AI Applications expanding beyond modalities of computer vision to linguistics, communication, commerce and language understanding
- With evolution of AI Applications across many verticals, continued push for innovation around latency, heterogeneity, concurrency and user personalization
  - This is putting a lot of emphasis for the need of custom HW modules in several BU verticals in Qualcomm
- Qualcomm silicon continues to show leadership in performance and energy efficiency in industry leading benchmarks
- High investment in software continues to accelerate AI solution deployment across all verticals



## Qualcomm Mobile AI

[Mobile AI | On-Device AI | Qualcomm®](#)

## Qualcomm & Google NAS

[Qualcomm Technologies and Google Cloud Announce Collaboration on Neural Architecture Search for the Connected Intelligent Edge | Qualcomm](#)

Vinesh Sukumar

Senior Director, Product Management – AI/ML  
vinesuku@qti.qualcomm.com

## 2022 Embedded Vision Summit

- *"Powering the Intelligent Connected Edge and the Future of On-Device AI"* **Ziad Asghar May 18 9:30 - 10:00 AM PT**
- *"A Practical Guide to Getting the DNN Accuracy You Need and the Performance You Deserve"* **Felix Baum May 18 2:40 - 3:10 PM PT**
- *"Tools for Creating Next-Gen Computer Vision Apps on Snapdragon"* **Judd Heape May 18 10:50 - 11:20 AM PT**
- *"Seamless Deployment of Multimedia and Machine Learning Applications at the Edge"* **Megha Daga May 17 2:40 - 3:10 PM PT**



**Thank You**

Qualcomm