



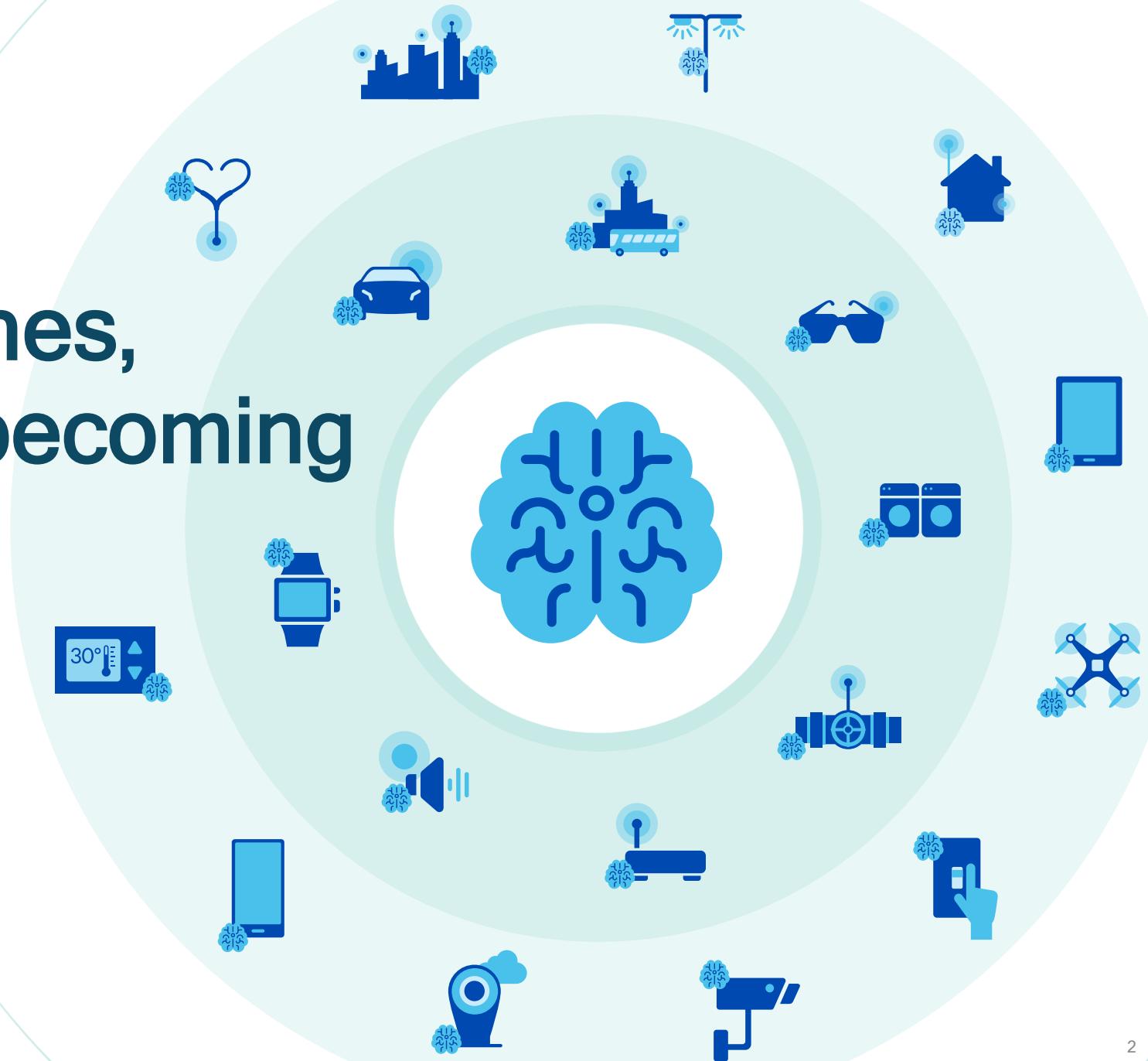
Mobile Innovation – The Key to AI @ The Edge

Gary Brotman
Director, AI Product Management
Qualcomm Technologies, Inc.
October 11, 2017

@GaryZero



Devices, machines, and things are becoming more intelligent



“Embedded and connected systems will reach annual volumes of 24 billion by 2021, with half of all the compute architectures shipping in 2021 supporting and processing AI.”



“Heterogeneous architectures will be instrumental in enabling developers to program once and leverage the most efficient solution for AI workloads.”



Mobile Drives Innovation at the Edge



Mobile computing



Smart homes



Industrial IoT



Wearables



Rapid replacement cycles



Smart cities



Automotive



Healthcare



Networking



Extended reality

Superior scale

Integrated and optimized technologies

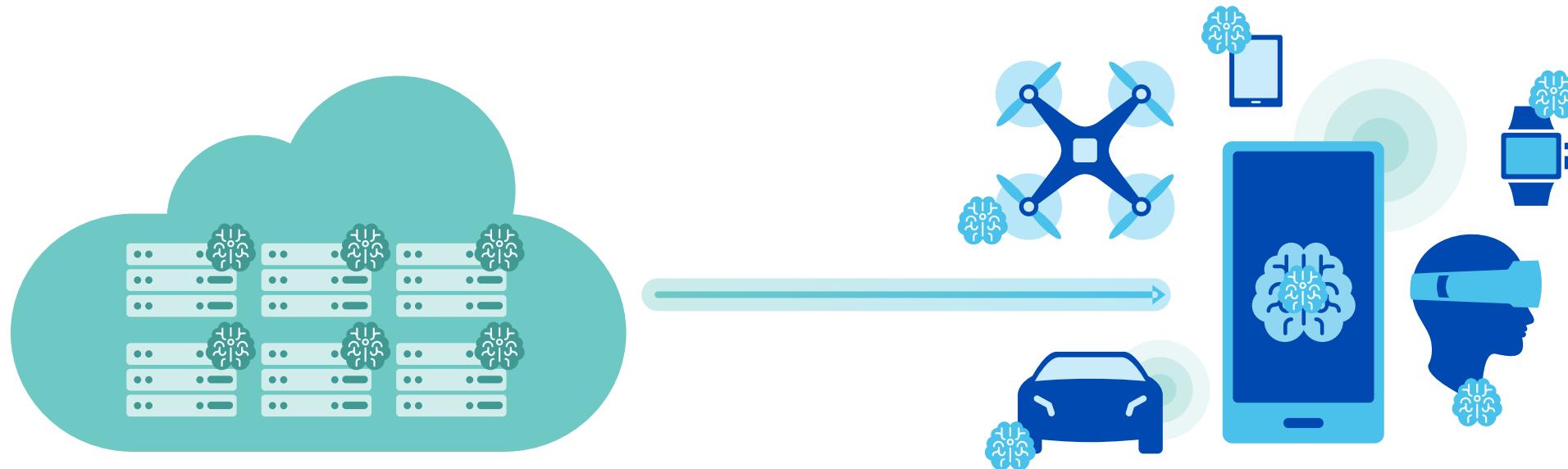
Mobile - THE most pervasive AI platform

>8.5 Billion

Cumulative smartphone unit shipments
forecast between 2017-2021



Intelligence is moving to the device

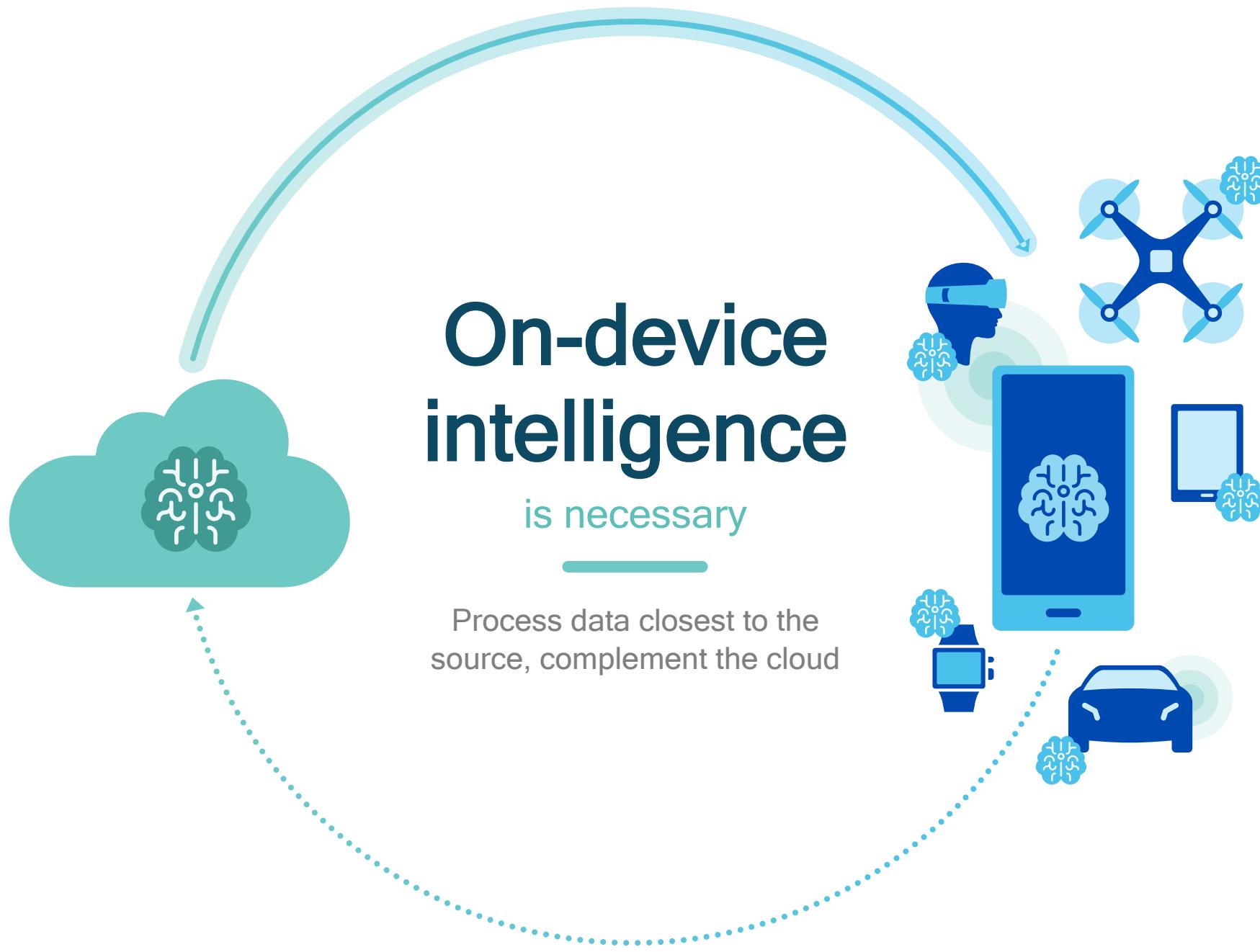


Server/Cloud

Training
Execution/Inference

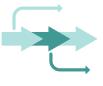
Devices

Execution/Inference
Training (emerging)



Power and thermal efficiency are essential for on-device AI

The challenge of AI workloads

-  Very compute intensive
-  Large, complicated neural network models
-  Complex concurrencies
-  Always-on
-  Real-time



Constrained environments

-  Must be thermally efficient for sleek, ultra-light designs
-  Requires long battery life for all-day use
-  Storage / Memory bandwidth limitations

Making on-device intelligence a reality

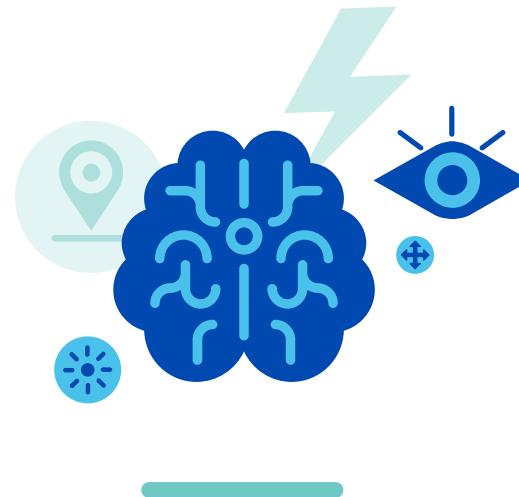
Focusing on high performance HW/SW and optimized network design



Efficient hardware

Developing heterogeneous compute to run demanding neural networks at low power and within thermal limits

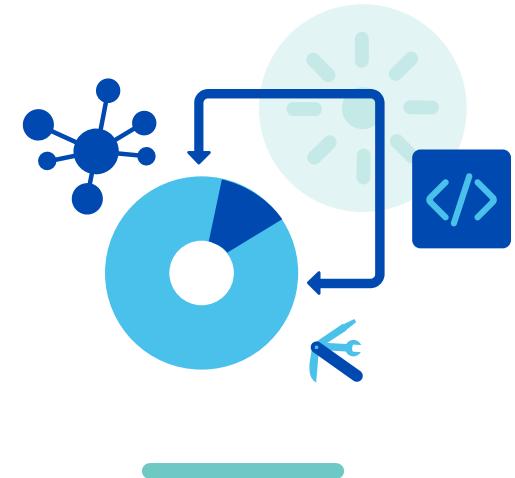
Selecting the right compute block for the right task



Algorithmic advancements

Algorithmic research that benefits from state-of-the-art deep neural networks

Optimization for space and runtime efficiency

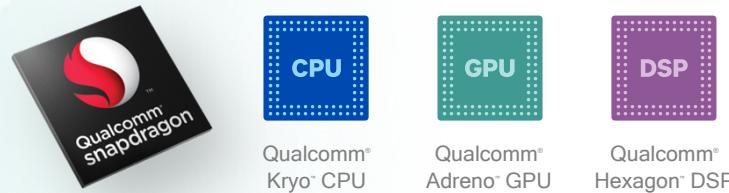
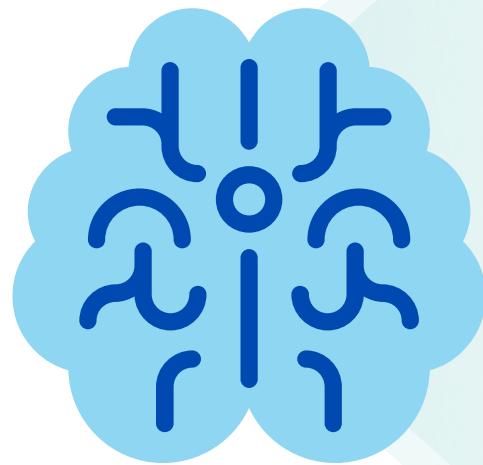


Software tools

Software accelerated run-time for deep learning

SDK/development frameworks

Snapdragon Neural Processing Engine



Available at: developer.qualcomm.com

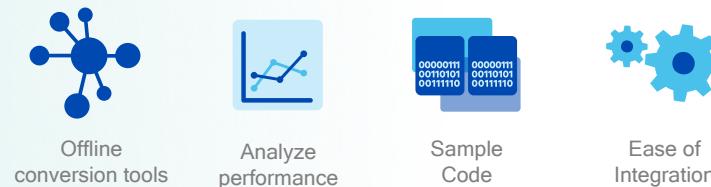
Efficient execution on Snapdragon

- Takes advantage of Snapdragon heterogeneous computing capabilities
- Runtime and libraries accelerate deep neural net processing on all engines: CPU, GPU, and DSP with HVX



Model framework/network support

- Convolutional neural networks and LSTMs
- Support for Caffe/Caffe2, TensorFlow, and user/developer defined layers



Optimization/debugging tools

- Offline network conversion tools
- Debug and analyze network performance
- API and SDK documentation with sample code
- Ease of integration into customer applications

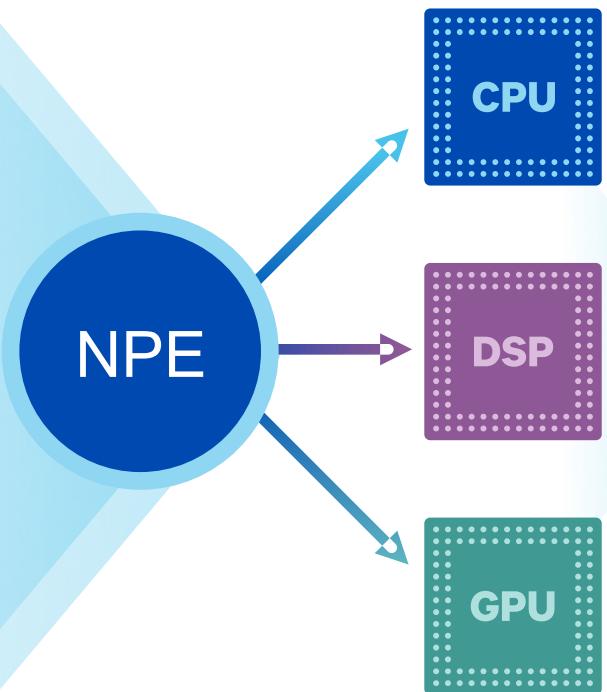
Software accelerated runtime for the execution of deep neural networks on device



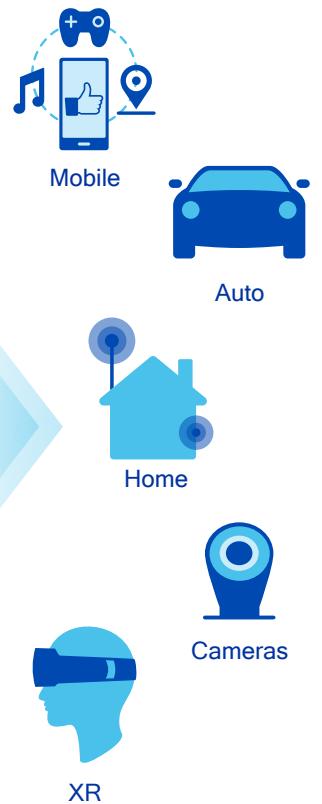
Caffe



GoogleNet/
Inception
SSD
Alexnet
ResNet
SqueezeNet
Faster - RCNN
MobileNet



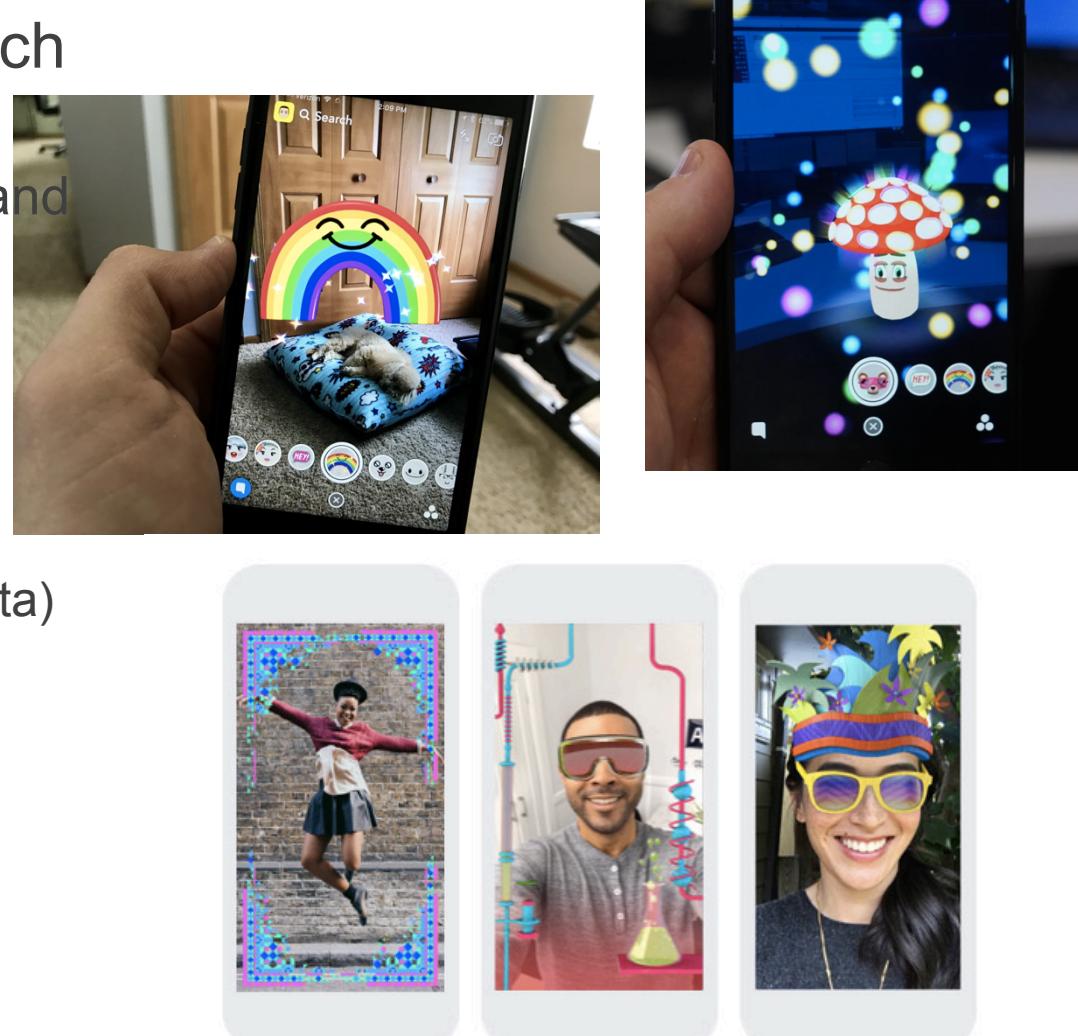
- Object classification
- Face detection
- Scene segmentation
- Natural language understanding
- Speaker recognition
- Security/Authentication
- Resource management



Social Media Apps: Camera + AR + AI = Consumer Delight

AI workloads moving from cloud to the edge

- Leading social media apps using AR & AI tech in abundance
 - Camera at the center – key to consumer delight and communication
 - Augmented reality features powered by AI
 - Style Transfer & Filters
 - Frames & Masks
 - Photo and live videos, including 360
 - Contextual awareness (e.g. location/sensor metadata)
- Progression of high level tasks
 - Detection, Perception, Understanding, Prediction
- Why on device?
 - Lower latency, privacy, connection independent



AI will bring XR closer to total immersion

Creating physical presence in real or imagined worlds





AI will make the car of the future possible

Redefining the in-car experience

- Natural user interfaces
- Personalization
- Driver awareness monitoring

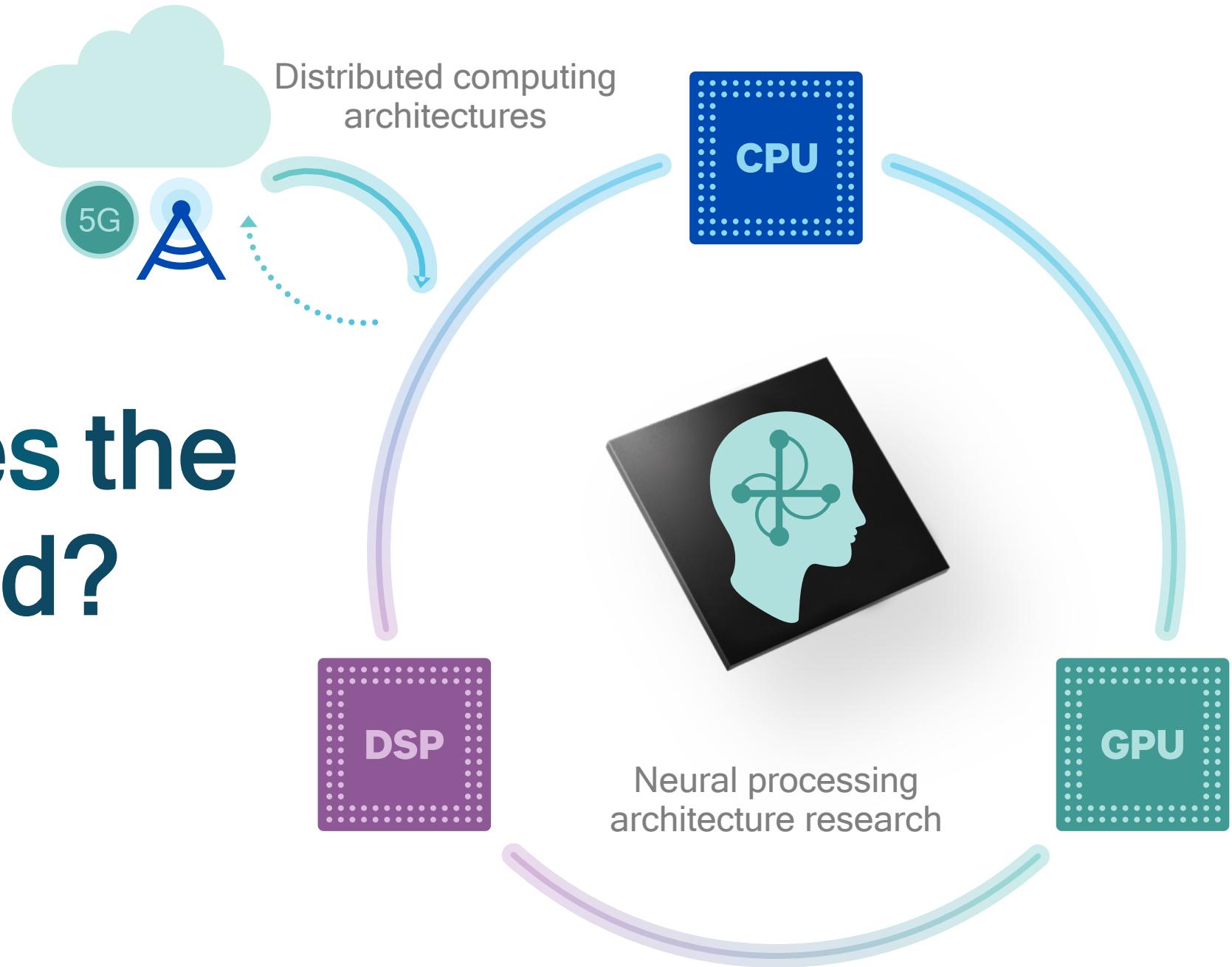
Paving the road to autonomy

- Surround view perception
- Sensor fusion
- Path planning
- Decision making

On-device AI Use Cases



What does the future hold?



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.

©2017 Qualcomm Technologies, Inc. and/or its affiliated companies. All Rights Reserved.

Qualcomm, Snapdragon, MSM, Hexagon, Kryo, and Adreno are trademarks of Qualcomm Incorporated, registered in the United States and other countries. Qualcomm Spectra, Aqstic and All-Ways Aware are trademarks of Qualcomm Incorporated.

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.

