# AI FIRSTS FROM QUALCOMM AI RESEARCH

## THE GROUP HAS PIONEERED NEW TECHNOLOGIES FROM ON-DEVICE LEARNING TO WIRELESS AI AND HAS AMBITIOUS PLANS FOR THE NEXT SET OF "FIRSTS".

Whenever people take photos or speak to a digital assistant using a mobile phone, they often don't realize that they just took advantage of Artificial Intelligence (AI). If they think of AI at all, it is typically in the context of Autonomous Vehicles or perhaps Facebook's (Meta's) massive data centers. While AI is becoming ubiquitous and distributed across edge devices and cloud servers, many challenges remain to realize the connected intelligent edge vision CEO Cristiano Amon has for AI to enable automated perception, reasoning, and action on edge devices.

For AI to enable the levels of automation and personalization Mr. Amon and his team envisions, AI hardware and software must become much smaller, faster, more efficient, lower power, and able to learn continuously at the edge in the real world. This provides the perfect complement to remote processing in the cloud, whose reach has been further advanced through Qualcomm's 5G technology. These are not just engineering challenges; progress and breakthroughs must come from fundamental scientific, applied, and platform research. This is where Qualcomm AI Research comes in, built with Qualcomm's core culture of innovation and evolved from over a decade of investment in R&D focused on machine learning and AI.

## AI FIRSTS FROM QUALCOMM AI RESEARCH

Qualcomm AI Research recently held a webinar where Jilei Hou, VP of engineering and the head of AI R&D at Qualcomm AI Research, outlined "firsts" in eight research areas for which they are justifiably proud. Their novel AI research and full-stack AI optimizations have pushed the AI industry forward and enabled first-ever proof-of-concept demonstrations on commercial mobile devices.

Their AI firsts include 8-bit model quantization leading to the best power-efficiency toolkit in the industry (the AI Model Efficiency Toolkit, or AIMET), on-device learning demonstrating a 30% improvement in keyword detection, federated learning with an end-to-end software framework, video semantic segmentation in real-time, group equivariant CNNs, AI for wireless with the first weakly supervised method for passive RF sensing, video super-resolution at 4K at 100 FPS on mobile, and neural video compression with the first real-time HD decoding on mobile. You can check out the rest of Qualcomm AI Research's firsts by downloading this presentation.

Artificial Intelligence at Qualcomm, like most everything else in the company, begins with mobile. And that heritage sets a high bar, since mobile devices are battery-powered, lightweight consumer devices with limited processing, memory, and I/O. The ongoing research underway at Qualcomm AI Research is centered around three primary domains: fundamental research, applied research, and platform research. The latter covers topics such as power efficiency and on-device learning, which are essential in mobile but have also benefited Qualcomm in other vertical markets like IOT, XR, and automotive. But efficiency doesn't stop with hardware; Qualcomm AI Research is also advancing model design, compression, quantization, algorithms, and software tools to attack the problem through full-stack research.



Efficiency comes in part from using less data in the model and quantization. If one can run a model with 4-bit integer math, it will be up to 64 times more efficient than using 32-bit floating point. But all that efficiency cannot come at the cost of reduced accuracy. To that point, Qualcomm AI Research has demonstrated a state-of-the-art 8-bit transformer model with less than 1% accuracy degradation.

Today, some AI processing has to be offloaded from the mobile device to a cloud service, but the trend is toward more and more processing at the edge devices. Faster response times, improved privacy, better personalization, and improved understanding within the context of the request are all benefits of on-device learning.

## LOOKING TO THE FUTURE

The researchers at Qualcomm AI Research are pursuing additional new "firsts" that could result in significant value and differentiation for customers and partners. Targeting improvements in AI efficiency on mobile devices and the rest of the connected intelligent edge can pave the way for improved personalization and automation of many tasks, enabling new use cases and improving the overall user experience.



## CONCLUSIONS

Increased AI processing on edge devices is both inevitable and extremely valuable to improve user experience and device functionality. The challenges remain daunting, given the limited amount of processing and memory capacity on mobile platforms. Qualcomm AI Research is rising to the challenge, and these 8 "AI Firsts" demonstrate that the organization is already making significant progress, with more innovations in sight.

## IMPORTANT INFORMATION ABOUT THIS PAPER

**AUTHOR:** Karl Freund, Founder, and Principal Analyst at Cambrian-AI Research

**INQUIRIES:**

Contact us if you would like to discuss this report, and Cambrian-AI Research will respond promptly.

## CITATIONS

This paper can be cited by accredited press and analysts but must be mentioned in the context, displaying the author's name, author's title, and "Cambrian-AI Research." Non-press and non-analysts must receive prior written permission from Cambrian-AI Research for any citations.

## LICENSING

This document, including any supporting materials, is owned by Cambrian-AI Research. This publication may not be reproduced, distributed, or shared in any form without Cambrian-AI Research's prior written permission.

## DISCLOSURES

This document was developed with Qualcomm Technologies, Inc. (QTI) funding and support. Although the paper may utilize publicly available material from various vendors, including QTI, it does not necessarily reflect the positions of such vendors on the issues addressed in this document.

## DISCLAIMER