October 2019

@qualcomm_tech

Amsterdam

Qualcom

The Next Generation Convolutions

Max Welling VP Technology <u>Qualcomm Technologies Netherlands</u>



Convolutional Neural Networks

- Each neuron is an image / feature map / channel
- Each synapse / edge is a filter
- The filter slides over the image
- At every location it multiplies the image values with the filter values and adds them up



Convolutions

- Convolution is arguably the technology that is responsible for much of the deep learning success.
- Can we generalize convolutions?

- Convolutions are expensive to run on hardware.
- Can we build radical new hardware to improve the efficiency of computing convolutions?

Generalized Convolutions

- Graph convolutions
- Group convolutions
- Gauge convolutions
- Quantum convolutions

- In memory convolutions
 Analogue calculations
 Brain inspired distributed memory
 - Breaks with Von Neumann architecture

Computers are consuming an increasing amount of energy

By 2025, the data center sector could be using 20% of all available electricity in the world¹

A cloud provider used the equivalent energy consumption of ~366,000 US households in 2014²

Bitcoin mining in 2017 used the same energy as did all of Ireland³

Given the economic potential of AI, these numbers will only be increasing

1. Andrae, Anders (2017) Total Consumer Power Consumption Forecast; 2. The Verge (2014); 3.The Guardian (2017)





2025 Will we have reached the capacity of the human brain? Energy efficiency of a brain is 100,000x better than current hardware

Value created by AI must exceed the cost to run the service



Economic feasibility per transaction may require cost as low as a micro-dollar (1/10,000th of a cent)

We need more AI per Joule!



Santiago Ramón y Cajal Father of neuroscience

• Discovered the neuron



Von Neumann architectures are not ideal for deep learning



We need architectures that look more like the brain!

In Memory Computation

- Inputs (neural activations) are voltages (word lines)
- Weights are conductances (synapses)
- Outputs are currents and compute Y = W * X (bit lines)

- Computation is analogue (uses less power)
- Memory is where computation is performed (like a brain)
- Deep learning can compensate for extra noise (think dropout)





Albert Einstein Inventor of General Relativity



Today's deep learning Traditional CNNs

Produce state-of-the art results but... do not generalize input like rotations



(Convolutional neural networks would need to be retrained with new rotated images to determine new set of parameters—like filter weights)

Tomorrow's deep learning Gauge Equivariant CNNs

No matter how you rotate or move the object, the generalized model will still identify the object





Rotated objects and images applicable to drones, robots, cars, fisheye-lens cameras. VR, AR,..



(Generalized CNNs (G-CNN): Gauge equivariant CNN, Group, and Steerable CNN pioneered by Qualcomm AI Research do not need to be retrained)

Convolutional NNs on curved surfaces

Gauge Equivariant CNNs for deep learning on curved spaces

- In the future "images" may become 3D (e.g. with depth)
- Can we do deep learning on 2D surfaces like we do on images?
- Yes, but we need to use the same theory Einstein used in General Relativity



- Kernels rotate when we transport them over a curved surface
- To convolve at p, we need to transport the signal to p.
- Every point has its own frame, so we must match these frames
- We study how results change under changes of local frames (a.k.a. gauge transformations)

12





Niels Bohr, Erwin Schrödinger

Inventors quantum mechanics (among others)



The Crazy world of Quantum mechanics

$$P(C) = P(A)P(A \rightarrow C) + P(B)P(B \rightarrow C)$$



Dutch book argument: quantum mechanics is a consistent theory of probabilities

No Measurement at A,B

$P(C) \neq P(A)P(A \rightarrow C) + P(B)P(B \rightarrow C)$



Quantum statistics is consistent but weird: when not measuring the photon it is in both paths simultaneously!

Quantum computation

A qubit is more like a sphere than a bit.

• Qubit:



• Entanglement:



Entangled qubits span all possible combinations:

"00", "10", "01", "11"

Quantum inspired convolutional kernels

Can we leverage the strange world of quantum statistics to design new quantum convolutions?



• Quantum statistics opens new doors to deep learning, even classically! (but with the benefit that it runs efficiently on a quantum computer) Convolutions are the core of Deep Learning's success

AI algorithms and hardware need to be energy efficient: perform convolutions in memory!

We can generalize convolutions to curved spaces. Applications in e.g. medical imaging.

Quantum mechanics provides a new theory of probability, leading to "quantum convolutions".



Qualcom

Thank you!

Follow us on: **f** 🎔 in

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.

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

Qualcomm is a trademark of Qualcomm Incorporated, registered in the United States and other countries. 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.