Programming optics for quantum & neuromorphic computing

AI Hardware Program Fall Research Update

Dirk Englund*

* ∈ {RLE, MTL, EECS}@MIT Quantum Photonics Laboratory @ MIT & Collaborators ∈ {MIT, MIT-LL, Harvard, Brookhaven National Laboratory, Sandia National Laboratory, U-Arizona, Stanford, BBN Technologies, MITRE Corp. U. Vienna, Air Force Research Laboratory@Rome, Army Research Laboratory, Oak Ridge NL, U-Maryland,QuEra Computing, TU Berlin, QuEra}



Designing Systems to Tackle Power & Data Movement in DNNs

ISSCC 2014 / SESSION 1 / PLENARY / 1.1

1.1 Computing's Energy Problem (and what we can do about it)

Mark Horowitz

Departments of Electrical Engineering and Computer Science, Stanford University, Stanford, CA

1. Introduction

Technology scaling has decreased the cost of computing to the point where it can be included in almost anything. As a result, we now live in a world surrounded by computing devices. They power our sacrehes on Google, connect to our friends on Facebook, answer our questions to Siri, and serve us our entertainment on Youtube; they are in our homes everywhere, in all our appliances (i recently had to reboot my refrigerator), cars, workplaces, and even in the cards we send to each other. We have become so accustomed to computing becoming taster, cheaper, and lower power, we simply assume it will continue. Aiready, smartphone capabilities are being embedded in eye glasses [1] and smart watches [2].

CONTRIBUTED

PAPER

While scaling computing performance I tors have made scaling increasingly diff power to become the principal constainreviews how computing became power field scaling (3) proke down, and explain change to fix our problems. The rest of to to addressing this computing-energycit will take more than parallelism to get scaling computing performance is to cr are better matched to the task and each tools that allow application experts to cr creating these tools is challenging, they is optimized computing!

2 Processor Scaling

processors became power constrained and leakage current grew, it became apparent that one could dramatically reduce the power dissipation, and improve the performance yield of a processor if each processor chip could specify the supply voltage that was required for it to operate at the desired performance. This would allow a chip fabricated with high-leakage, lower-average-V,, transistors, to run at a lower supply voltage, reducing both the dynamic and leakage power, for overall power optimization. Correspondingly, processors with higher V_m transistors, and lower leakage could run at a higher supply voltage while still operating within the total power budget, enabling these transistors to operate at the desired speed. While this has been good for processor specification. it has made it much more difficult to track how the average supply voltages have been scaling over the past decade. Thus, the numbers in the voltage plot in Figure 1.1.4 are the peak allowable supply voltages, and do not represent the average voltages used. From limited data, the actual operating supply voltages seem to remain in the 0.9 to 1.1 volt range for peak performance. But, the recent move to 3-D channel structures with reduced leakage currents, has enabled about a 100 to 200mV decrease in operating voltage.

Like most chips today, processors used to run at a fixed supply voltage, and this voltage depended on the fabrication technology that was used. But, as

2 Technology to the Becaus?

Efficient Processing of Deep Neural Networks: A Tutorial and Survey

This article provides a comprehensive tutorial and survey coverage of the recent advances toward enabling efficient processing of deep neural networks.

By VIVIENNE SZE[©], Senior Member IEEE, YU-HSIN CHEN, Student Member IEEE, TIEN-JU YANG, Student Member IEEE, AND JOEL S. EMER, Fellow IEEE

ABSTRACT | Deep neural networks (DNNs) are currently widely used for many artificial intelligence (A) applications including computer vision, speech recognition, and robotics. While DNNs deliver state-of-the-art accuracy on many AI tasks, it comes at the cost of high computational complexity. Accordingly, techniques that enable efficient processing of DNNs to improve energy efficiency and throughput without sacrificing application accuracy or increasing hardware cost are critical to the wide denoments of DNNs in Accuracy. This article aline for sample between various hardware architectures and platforms; be able to evaluate the utility of various DNN design techniques for efficient processing; and understand recent implementation trends and opportunities.

KEYWORDS | ASIC; computer architecture; convolutional neural networks; dataflow processing; deep learning; deep neural networks; energy-efficient accelerators; low power; machine learning; spatial architectures; VLSI Takes more energy to fetch longer distances in memory

- Off-chip >> on-chip
- Large RAM > Small RAM ($E_{bit} \propto L \propto (N_{bit})^{0.5}$)
- Goal: max. use of small register memories → "local data reuse"



V. Sze, Proc. IEEE 105(12), 2295 (2017)

M. Horowitz, ISSCC 2014, pp.10-14

Systolic Array architecture is optimized for GEMM



Two key limiters to the TPU:

- Power → Heat dissipation
- Interconnects / Data Movement

In-Datacenter Performance Analysis of a Tensor Processing Unit

Norman P. Jouppi, Cliff Young, Nishant Patil, David Patterson, Gaurav Agrawal, Raminder Bajwa,

.. Google team: >cs>arXiv:1704.04760



Figure 1. TPU Block Diagram. The main computation is the yellow Matrix Multiply unit. Its inputs are the blue Weight FIFO and the blue Unified Buffer and its output is the blue Accumulators. The yellow Activation Unit performs the nonlinear functions on the Accumulators, which go to the Unified Buffer.

Google TPU cloud

offerings:





Cloud TPU v2 180 teraflops 64 G8 High Bandwidth Memory (HBM)



Cloud TPU v2 Pod 11.5 petaflops 4 TB HBM 2-D toroidal mesh network



Cloud TPU v3 420 teraflops 128 GB HBM



Cloud TPU v3 Pod 100+ petaflops 32 TB HBM 2-D toroidal mesh network

How does this matter for actual DNNs?

Start with a DNN accelerator design (such as Eyeriss)

Not just a systolic array!

Top-Level Cor

Filter

Ifmap

Ofmap

RLC

Link Clock Core Clock

Has a balanced hierarchy of memory to handle many tasks movement costs, not raw processing!

Find the optimal dataflow

Trying to minimize data Despite all the optimizations, chips are dominated by data ita must

'weight-stationa**rv**"

, not be

Moore's Law: limited gains for data movement This is why optics is interesting for future DNN hardware A compiler picks the best data-flow after a parameter



(weight-stationary)





V. Sze, Proc. IEEE 105(12), 2295 (2017)

Off-Chip

DRAM

64

bits



Taken from L Bernstein*, A Sludds*, R Hamerly, V Sze, J Emer, D Englund, Sci Reports (2021) 11:3144



A Renaissance in Physics of Computing ?

Hardware

Programmable photonic integrated circuits (PICs)

Review: G Wetzstein et al, Nature 588 (2020) Review: W Bogaerts et al, <u>Nature 586, 207–216</u>

Example:



Matrix algebra in programmable photonic integrated circuits (PICs)



Programmable optical matrix transformation: 176 phase shift degrees of freedom (8 bit)

- Y. Shen*, N. C. Harris*, et al [with M Soljacic, MIT], Nature Photon 11 (2017). * equal authors
- See also D.A.B.M Miller, "Sorting out Light" , Science 347 (2017)
- Reviews:
 - N Harris et al, Optica 5 (2018)
 - W Bogaerts, D Pérez, J Capmany, D A. B. Miller, J Poon, D Englund, F Morichetti & A Melloni, Nature 586 (2020)
 - o G Wetzstein, A Ozcan, S Gigan , S Fan, D Englund, M Soljačić, C Denz, D A. B. Miller, D Psaltis, Nature 588 (2020)

Thermal phase shifters burn ~ 10 mW/each, but we have piezoelectric alternative (M Dong et al -Nature Photonics).. and nonvolatile phase shifters as by Pernice group are super interesting if we could program precisely

> OPSIS Foundry Collaborators: Michael Hochberg, Michael Fanto, Paul Alsing (AFRL), Stefan Preble (RIT), Philip Walther (U. Vienna)

"Weight-stationary architecture": flying output **y**=(stationary **W**).(flying **x**)



Y Shen*, N C. Harris*, S Skirlo, M Prabhu, T Baehr-Jones, M Hochberg, X Sun, S Zhao, H Larochelle, D Englund, and M Soljacić, *Nature Photonics* 11 (2017). *equal authors

Scalability of analog systems?



C. Babbage "difference engine" (1832)



A CONTINUOUS INTEGRAPH.*

BY

V. BUSH, F. D. GAGE, and H. R. STEWART. Electrical Engineering Department, Massachusetts Institute of Technology.

ABSTRACT.

A MECHANICAL integraph has been developed which plots continuously the integral of the product of two functions. It uses the principle of the electrical integrating watthour-meter combined with a moving table. Errors of the machine have been reduced to an average of I per cent. for common uses. By cross-connecting the device in a simple mechanical way, it is possible to solve cer-



- A Renaissance in Physics of Computing ?
- Programmable photonic integrated circuits (PICs)

Hardware

EC in hardware compute, training: \rightarrow "error-free optics"

Review: G Wetzstein et al, Nature 588 (2020) Review: W Bogaerts et al, <u>Nature 586, 207–216</u>

S. Bandyopadhyay, Optica 8, 1247, (2021) Hamerly et al., arXiv:2106.03249 (2021) R Hamerly, S Bandyopadhyay, DE, arXiv: 2109.05367 (2021)

Error correction techniques in Mach Zehnder Meshes

Numerical Optimization

Nonlinear optimizer: eg minimize norm $|U - U_0|$

Burgwal et al., OE 25, 28236 (2017) Mower et al., PRA 92, 032322 (2015) Pai et al., PRApp 11, 064044 (2019)

(+) Works for arbitrary meshes (-) Slow

(-) Needs pre-calibration (knowledge of MZI errors)

"Local" EC

Given errors, find (θ, ϕ) that fix them per-MZI.



In-situ Training

Adjoint-related method. Send input & weight gradient, get gradients in internal detectors.



Progressive Methods

Adjust MZIs one-by-one to match certain I/O modes.



Question: Can one construct "perfect" photonics from imperfect components? And can it be done efficiently?



36



Hardware

- A Renaissance in Physics of Computing ?
- Programmable photonic integrated circuits (PICs)
- EC in hardware compute, training: \rightarrow "error-free optics"
- EC in ML in-situ training?
 - Fully integrated on-chip inference

Review: G Wetzstein et al, Nature 588 (2020) Review: W Bogaerts et al, <u>Nature 586, 207–216</u>

S. Bandyopadhyay, Optica 8, 1247,(2021) Hamerly et al., arXiv:2106.03249 (2021) R Hamerly, S Bandyopadhyay, DE, arXiv: 2109.05367 (2021)

arXiv:2203.05466 (2022)

See Saumil's talk coming up

Single chip photonic deep neural network with accelerated training

Saumil Bandyopadhyay,^{1, *} Alexander Sludds,¹ Stefan Krastanov,¹ Ryan Hamerly,^{1, 2} Nicholas Harris,¹ Darius Bunandar,¹ Matthew Streshinsky,³ Michael Hochberg,⁴ and Dirk Englund^{1, †}

¹Research Laboratory of Electronics, MIT, Cambridge, MA 02139, USA
²NTT Research Inc., PHI Laboratories, 940 Stewart Drive, Sunnyvale, CA 94085, USA
³Nokia Corporation, New York, NY, 10016, USA
⁴Luminous Computing Inc., Mountain View, CA, 94041, USA



Hardware

- A Renaissance in Physics of Computing ?
- Programmable photonic integrated circuits (PICs)
- EC in hardware compute, training: \rightarrow "error-free optics"
- EC in ML in-situ training?
- Fully integrated on-chip inference

Review: G Wetzstein et al, Nature 588 (2020) Review: W Bogaerts et al, <u>Nature 586, 207–216</u>

S. Bandyopadhyay, Optica 8, 1247, 2021) Hamerly et al., arXiv:2106.03249 (2021) R Hamerly, S Bandyopadhyay, DE, arXiv: 2109.05367 (2021)

arXiv:2203.05466 (2022)

Computing Across the Internet's Edge: Tops/sec @mW on edge devices

ArXiv:cs.ET cs.LG 2205.09103

(2022); to appear in Science

Hey Siri





Hardware

- A Renaissance in Physics of Computing ?
- Programmable photonic integrated circuits (PICs)
- EC in hardware compute, training: \rightarrow "error-free optics"
- EC in ML in-situ training?
- Fully integrated on-chip inference

e optics" S. Ban Harr

S. Bandyopadhyay, Optica 8, 1247,(2021) Hamerly et al., arXiv:2106.03249 (2021) R Hamerly, S Bandyopadhyay, DE, arXiv: 2109.05367 (2021)

Review: G Wetzstein et al, Nature 588 (2020) Review: W Bogaerts et al, <u>Nature 586, 207–216</u>

arXiv:2203.05466 (2022)

Computing Across the Internet's Edge: Tops/sec @mW on edge devices ArX

Applications

Single-shot NNs: 20,000 MACs / 20 fs: Exascale computing

<u>ArXiv:cs.ET cs.LG 2205.09103</u> (2022); to appear in Science

L Bernstein et al, ArXiv:2205.09103 (2022)



Liane Bernstein, Alexander Sludds, Christopher Panuski, Sivan Trajtenberg Mills, Ryan Hamerly, and Dirk Englund, Scalable Ultralow Latency Photonic Tensor Processor, CLEO, IEEE, May 19, 2022, STh5G.7.

Answer:



The arrival of scalable control by analog photonic integrated circuits (APICs)

16-CH APIC² in 200 mm CMOS fabrication process¹



¹Mark Dong et. al . 2021. Nature Phot. (2021) ²Adrian Menssen, Artur Hermans, et. al. 2022 CLEO (2022); to be published (2022)

MIT-Sandia National Laboratory - MITRE - Harvard Separately: QuEra - SNL

16-CH APIC lithium niobate on insulator³



16 programmable projection fields, 16 x 10 GHz bandwidth

³Ian Christen et al, CLEO (2022), to be published

2D Spatial Light Modulator with ns switching rate; Si demonstrator ⁴:

10 B modulators per wafer







⁴ C. Panuski et al, arXiv:2204.10302 (2022)



- A Renaissance in Physics of Computing ?
- Programmable photonic integrated circuits (PICs)
- EC in hardware compute, training: \rightarrow "error-free optics"
- EC in ML in-situ training?
- Fully integrated on-chip inference
- Computing Across the Internet's Edge: *Tops/sec @mW on edge devices*
- Single-shot NNs: 20,000 MACs / 20 fs: Exascale computing

Review: G Wetzstein et al, Nature 588 (2020) Review: W Bogaerts et al, <u>Nature 586, 207–216</u>

S. Bandyopadhyay, Optica 8, 1247,(2021) Hamerly et al., arXiv:2106.03249 (2021) R Hamerly, S Bandyopadhyay, DE, arXiv: 2109.05367 (2021)

arXiv:2203.05466 (2022)

Z. Chen et al, CLEO 2022

ArXiv:cs.ET cs.LG 2205.09103 (2022); to appear in Science

L Bernstein et al, ArXiv:2205.09103 (202220)

Microlaser-based Coherent Scalable Efficient deep Learning ~ *fJ/OP and 25 TeraOP/(mm² · s)* R. Davis et al, CLEO 2022

Quantum optical neural networks: signal processing on quantum states

What could a DNN do with a *quantum* nonlinearity?





A Quantum Optical Neural Network (QONN)





Nonlinearities: atoms, artificial atoms, QDs, molecules, or bulk optical NL^{2,3}

Supervised training¹ \rightarrow

- Quantum optical state compression
- Reinforcement learning
- Black-box quantum simulation
- One-way quantum repeaters

Acknowledgements (ML portion)

Positions available in theory and experiment: http://www.rle.MIT.edu/qp

Chen (--> USC)



MIT Quantum Photonics Group: Dr. Ryan Hamerly (NTT, MIT)

PhD: Saumil Bandyopadhyay, Nick Harris (--> LightMatter), Darius Bunandar (--> LightMatter), Mihika Prabhu, Chris Panuski, Liane Bernstein, Alex Sludds, Ronald Davis









Dr. Nicholas Saumil Prof. Zaiiun Bandyopadhyay Chen Harris Collaborators on this work:

Sandia NL: M Eichenfield

Stephan Reitzenstein

Dr. Zhizhen Zhong

Liane Bernstein



Dr. Stefan Krastanov

Dr. Carlos Herranz





Mihika Prabhu



Hochberg

Dr. Darius

Banundar

Prof. Manva Ghobadi







MIT: Manya Ghobadi, Joel Emer, Vivienne Sze

TU Berlin: T Heuser, N Heermeier, J Lott,

MITRE Corp: Gerry Gilbert, Mark Dong, Gen Clark





Dr. Tom



Postdocs: Stefan Krastanov, Sri Krishna Vadlamani, Zaijun



Dr. Matthew Dr. Ari Novack

Streshinsky







