



NeurIPS 2019 recap

René Donner, contextflow

The background image shows an aerial view of the Vancouver skyline during sunset. The city is built on a peninsula, with a dense cluster of skyscrapers on the right and a more spread-out urban area on the left. In the foreground, a large white building with a distinctive scalloped roof, resembling a sail or a series of tents, sits on a pier extending into the water. Several small white water taxis are docked at the pier. The water is a deep blue-green color. The sky is clear with a few wispy clouds.

Overview

Tutorials

Papers

Workshops & Hardware

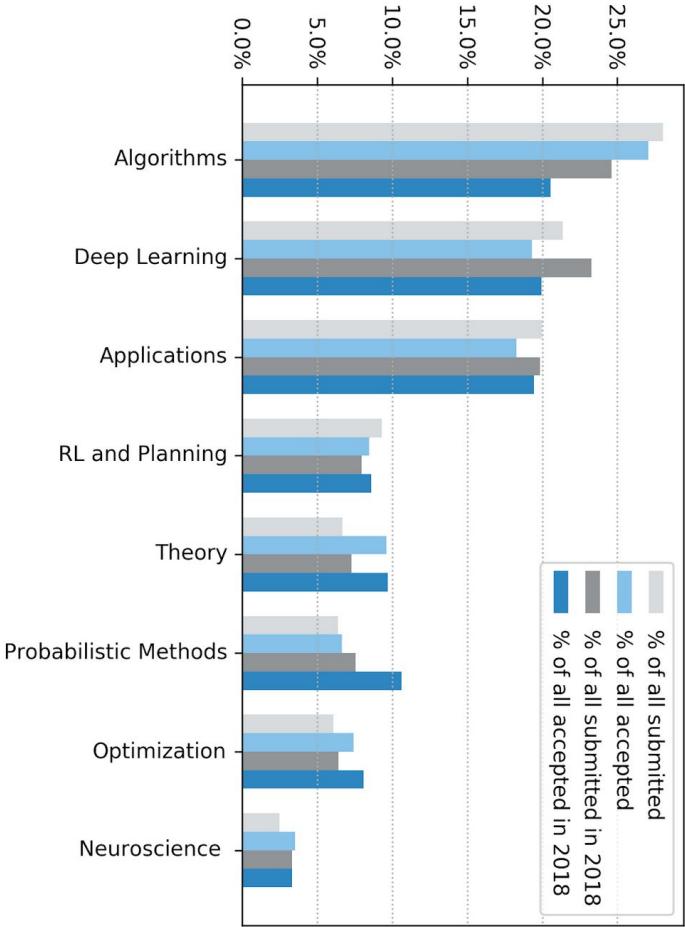
Overview

13.500 attendees - ticket lottery
1000 papers at main conference
430 papers at workshops

1 day of industrial talks / expo
1 day tutorials
3 day orals + posters
2 days workshops

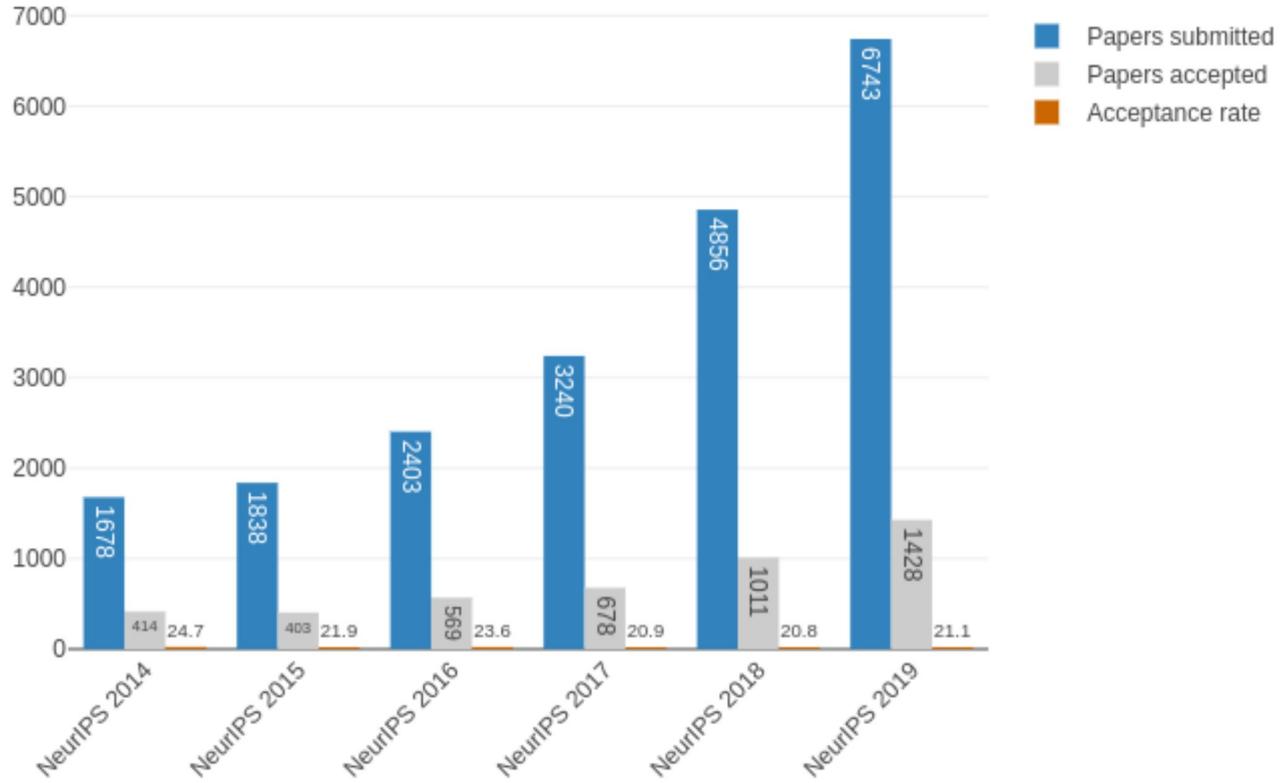


Overview



Statistics of acceptance rate NeurIPS

Overview



Welcome to NeurIPS

Poster Session is Full

Please try again later

NeurIPS 2019

Videos

NeurIPS | 2019

Thirty-third Conference on Neural Information Processing Systems

Year (2019) ▾

Help ▾

My Registrations

Profile ▾

Contact NeurIPS

Sponsor Info

Publications

Future Meetings

Diversity & Inclusion

New in ML

Code of Conduct

About Us

Press

News

Dates Schedule ▾ Submit ▾ Attend ▾ Organizers ▾

Vancouver Convention Center, Vancouver CANADA

Sun Dec 8th through Sat the 14th
(Sunday is a full day industry expo)

Mobile-friendly Schedule

Expo (Sun Dec 8th)

[Talks](#) » [Demonstrations](#) » [Workshops](#) » [Sun](#) »

[Expo Brochure PDF](#) »

Main Meeting

[Mon](#) » [Tue](#) » [Wed](#) » [Thu](#) » [Fri](#) » [Sat](#) »

[Tutorials](#) » [Orals](#) » [Posters](#) » [Demos](#) »

[Workshops](#) » [Socials](#) »

Schedule

Expo (Sun Dec 8th)

[Talks](#) » [Demonstrations](#) » [Workshops](#) » [Sun](#) »

[Expo Brochure PDF](#) »

Main Meeting

[Mon](#) » [Tue](#) » [Wed](#) » [Thu](#) » [Fri](#) » [Sat](#) »

[Tutorials](#) » [Orals](#) » [Posters](#) » [Demos](#) » [Workshops](#) »

[Socials](#) » [Everything](#) »

[Diversity and Inclusion Information and Schedules](#) »

[Live Streams](#) » [Pre-recorded Videos](#) »

Reproducibility Challenge
Paper Discussion Forum

The live streams will be available as an archive immediately after the stream finishes. Everything but the expo, socials, and poster sessions are streamed.

Registration

- General Admission Lottery

Sponsors

[View NeurIPS 2019 sponsors](#) »

Schedule

Thu Dec 12th 10:45 AM -- 12:45 PM @ East Exhibition Hall B + C #84

Poster

PC-Fairness: A Unified Framework for Measuring Causality-based Fairness

Yongkai Wu · Lu Zhang · Xintao Wu · Hanghang Tong

In Applications -- Fairness, Accountability, and Transparency

[Paper »](#)

[Poster »](#)

[Slides »](#)

Thu Dec 12th 10:45 AM -- 12:45 PM @ East Exhibition Hall B + C #85

Poster

This Looks Like That: Deep Learning for Interpretable Image Recognition

Chaofan Chen · Oscar Li · Daniel Tao · Alina Barnett · Cynthia Rudin · Jonathan K Sul In Applications -- Fairness, Accountability, and Transparency

[Paper »](#)

[3 min Video »](#)

Thu Dec 12th 10:45 AM -- 12:45 PM @ East Exhibition Hall B + C #86

Poster

Towards Automatic Concept-based Explanations

Amirata Ghorbani · James Wexler · James Zou · Been Kim

In Applications -- Fairness, Accountability, and Transparency

[Paper »](#)

[Slides »](#)

Thu Dec 12th 10:45 AM -- 12:45 PM @ East Exhibition Hall B + C #87

Poster

Adversarial Training and Robustness for Multiple Perturbations

Florian Tramer · Dan Boneh

In Applications -- Privacy, Anonymity, and Security

[Paper »](#)

[Slides »](#)

[3 min Video »](#)

Thu Dec 12th 10:45 AM -- 12:45 PM @ East Exhibition Hall B + C #88

Poster

Poster

An aerial photograph of the Vancouver skyline during sunset. The city is built on a peninsula, with a dense cluster of skyscrapers on the right and a more spread-out residential area on the left. In the foreground, a large white geodesic dome, likely the Canada Place convention center, sits on a pier extending into the dark blue water. Several small white water taxis are docked at a platform further out. The sky is a clear, pale blue.

Overview

Tutorials

Papers

Workshops & Hardware

9 Tutorials

Mon Dec 9th 11:15 AM -- 01:15 PM @ West Exhibition Hall C + B3

Tutorial

Efficient Processing of Deep Neural Network: from Algorithms to Hardware Architectures

Vivienne Sze

Mon Dec 9th 02:45 -- 04:45 PM @ West Exhibition Hall A

Tutorial

Representation Learning and Fairness

Moustapha Cisse · Sanmi Koyejo

Mon Dec 9th 05:45 -- 06:35 PM @ West Exhibition Hall C + B3

Invited Talk

How to Know

Celeste Kidd

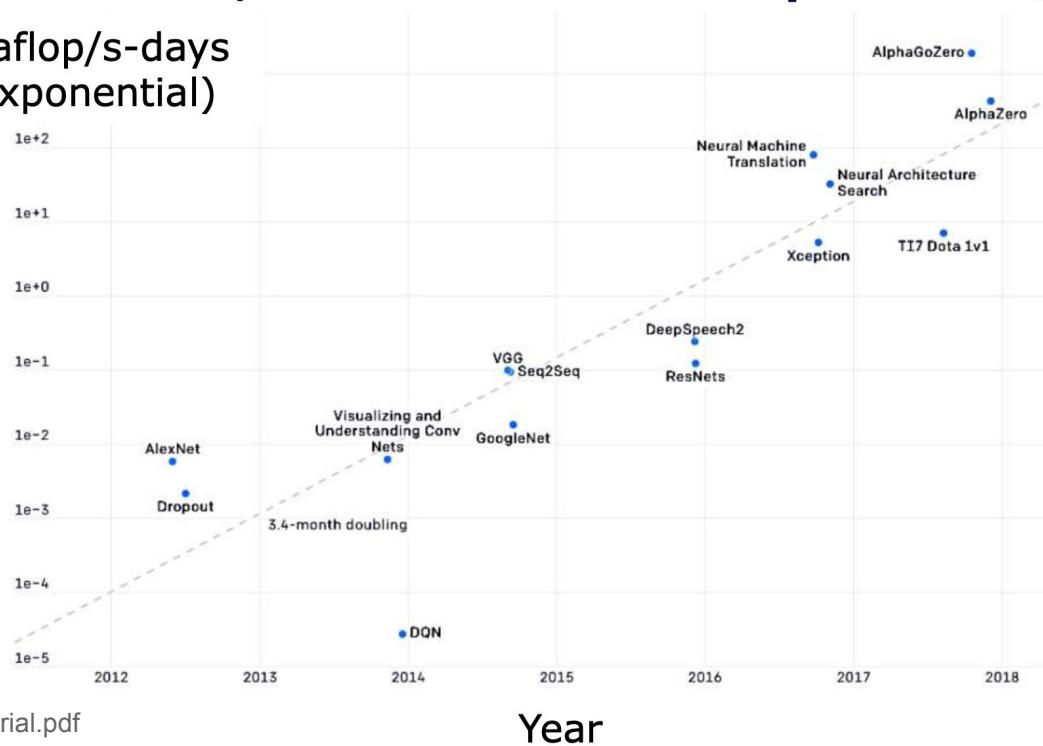
[PDF »](#)

Efficient Processing of Deep Neural Network: from Algorithms to Hardware Architectures

Vivienne Sze

AlexNet to AlphaGo Zero: A 300,000x Increase in Compute

Petaflop/s-days
(exponential)



Efficient Processing of Deep Neural Network: from Algorithms to Hardware Architectures

Vivienne Sze

Common carbon footprint benchmarks

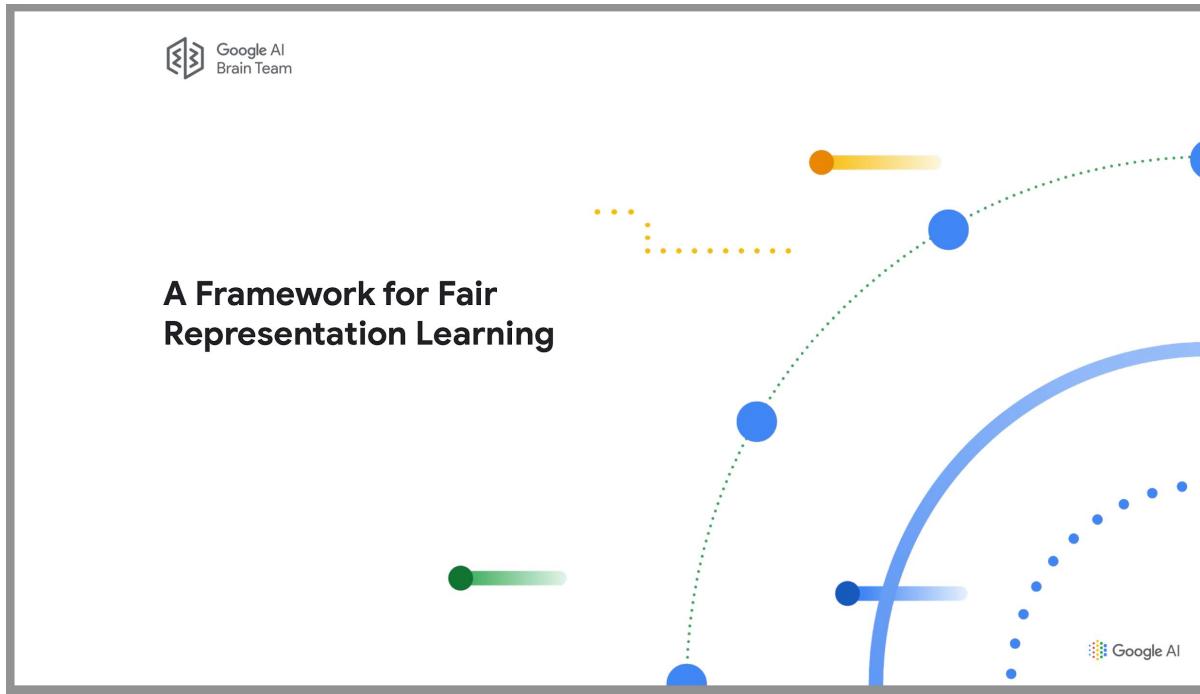
in lbs of CO₂ equivalent



Chart: MIT Technology Review · [Strubell, ACL 2019]

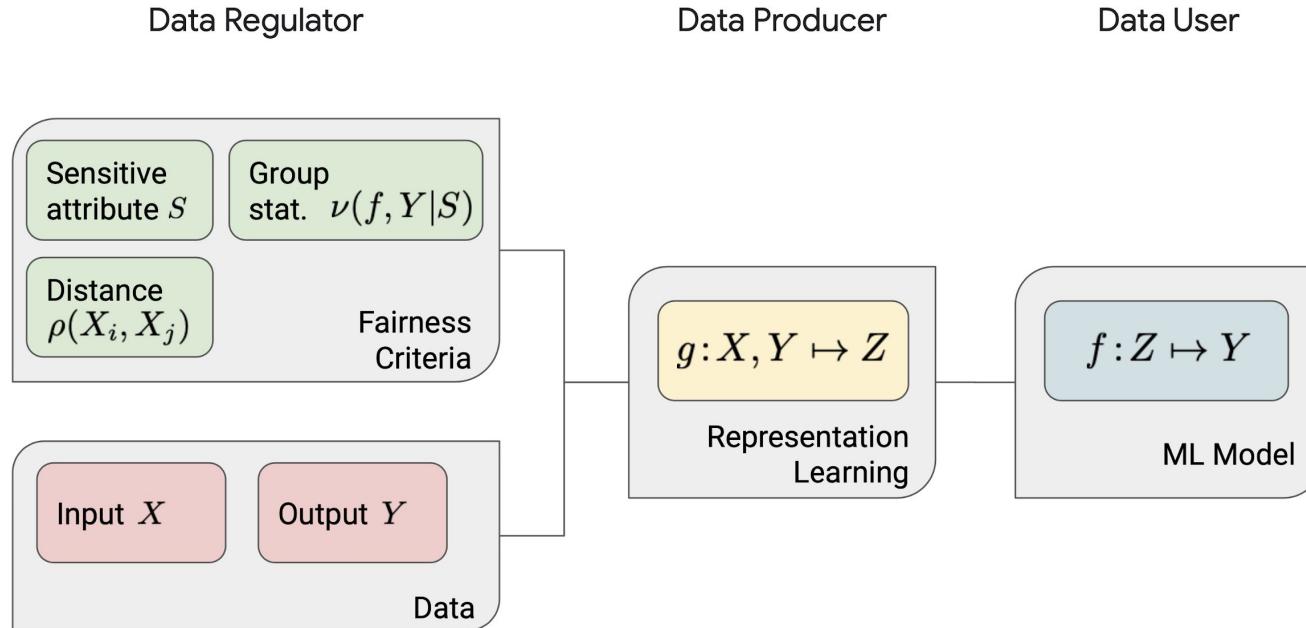
Representation Learning and Fairness

Moustapha Cisse · Sanmi Koyejo



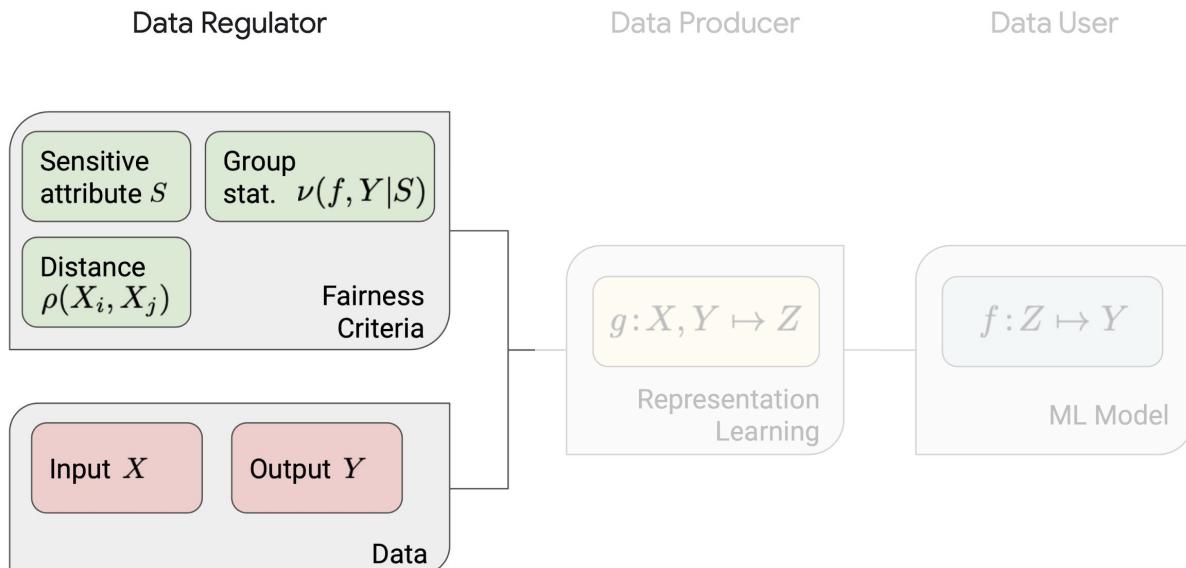
Representation Learning and Fairness

Moustapha Cisse · Sanmi Koyejo



Representation Learning and Fairness

Moustapha Cisse · Sanmi Koyejo



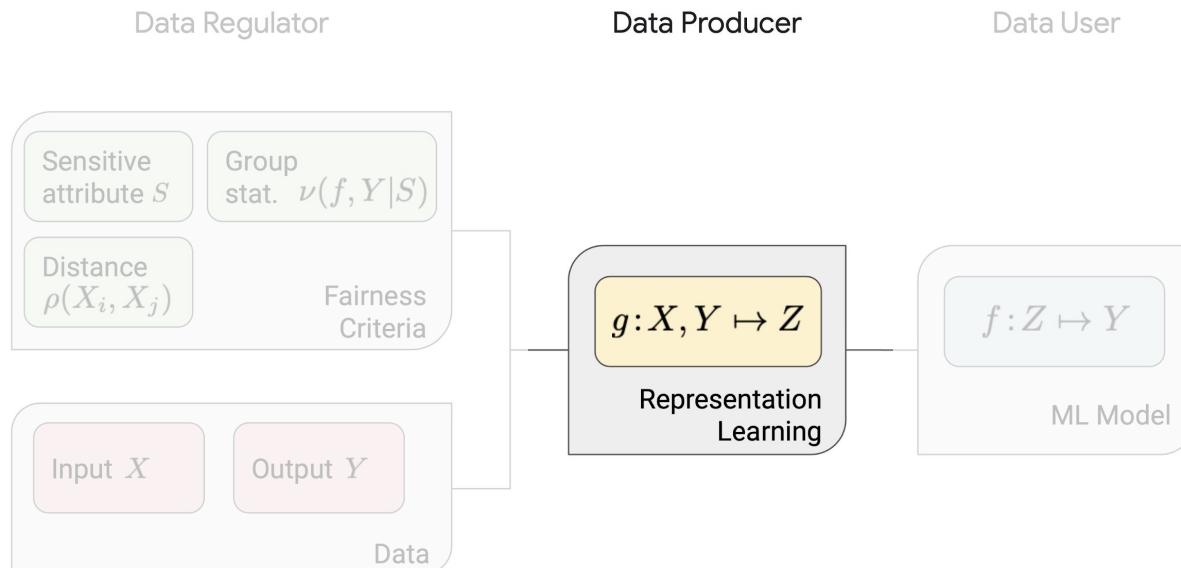
Data Regulator

Determines fairness criteria,
determines data source(s),
audits results

- INPUT: Data
 - OUTPUT: Fairness criteria
- AUDITING
- INPUT: Models
 - OUTPUT: Satisfactory?

Representation Learning and Fairness

Moustapha Cisse · Sanmi Koyejo



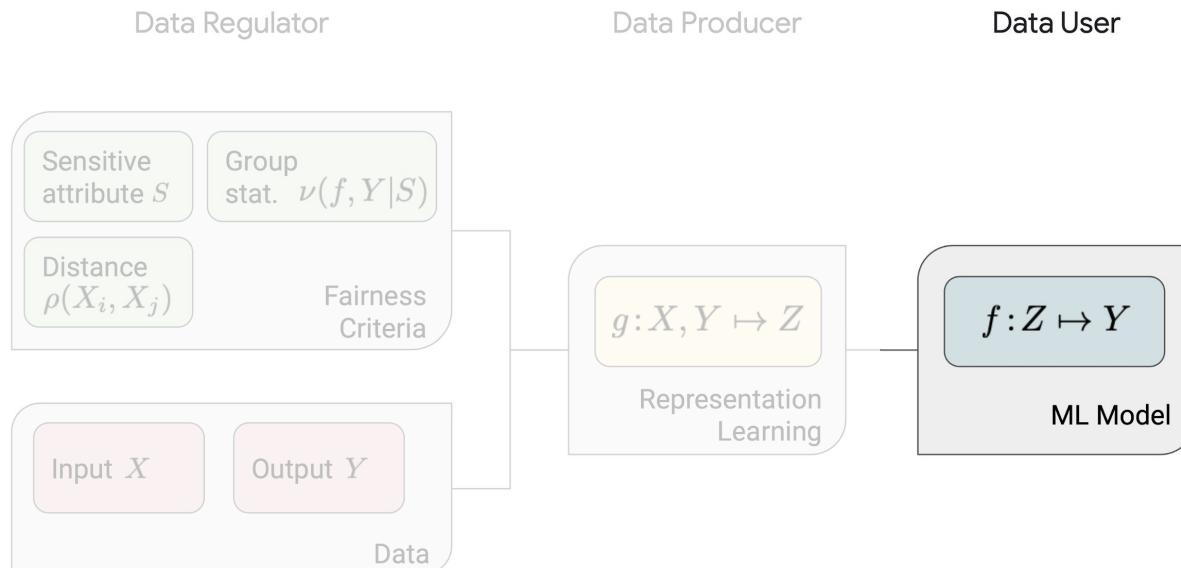
Data Producer

Computes the fair representation given data regulator criteria

- INPUT: Fairness criteria
- OUTPUT: Representation

Representation Learning and Fairness

Moustapha Cisse · Sanmi Koyejo



Data User

Computes ML model given
sanitized data

- INPUT: Sanitized data
- OUTPUT: ML model

How to Know

Celeste Kidd

[PDF »](#)

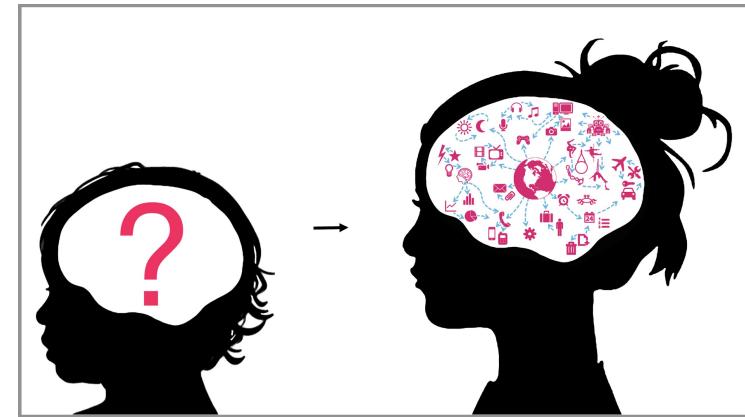


How to know

Celeste Kidd (celestekidd@berkeley.edu)
Psychology, University of California, Berkeley

 [@celestekidd](#)

NeurIPS, Vancouver, 9 Dec 2019



An aerial photograph of the Vancouver skyline during sunset. The city is built on a peninsula, with a dense cluster of skyscrapers on the right and a more spread-out residential area on the left. In the foreground, a large white building with a distinctive scalloped roof (the Canada Place convention center) sits on a pier extending into the dark blue water. A green-roofed building is visible behind it. On the water, several white water taxis are docked at a floating platform. The sky is a clear, pale blue.

Overview

Tutorials

Papers

Workshops & Hardware

SuperGLUE: A Stickier Benchmark for General-Purpose Language Understanding Systems

Alex Wang*
New York University

Yada Pruksachatkun*
New York University

Nikita Nangia*
New York University

Amanpreet Singh*
Facebook AI Research

Julian Michael
University of Washington

Felix Hill
DeepMind

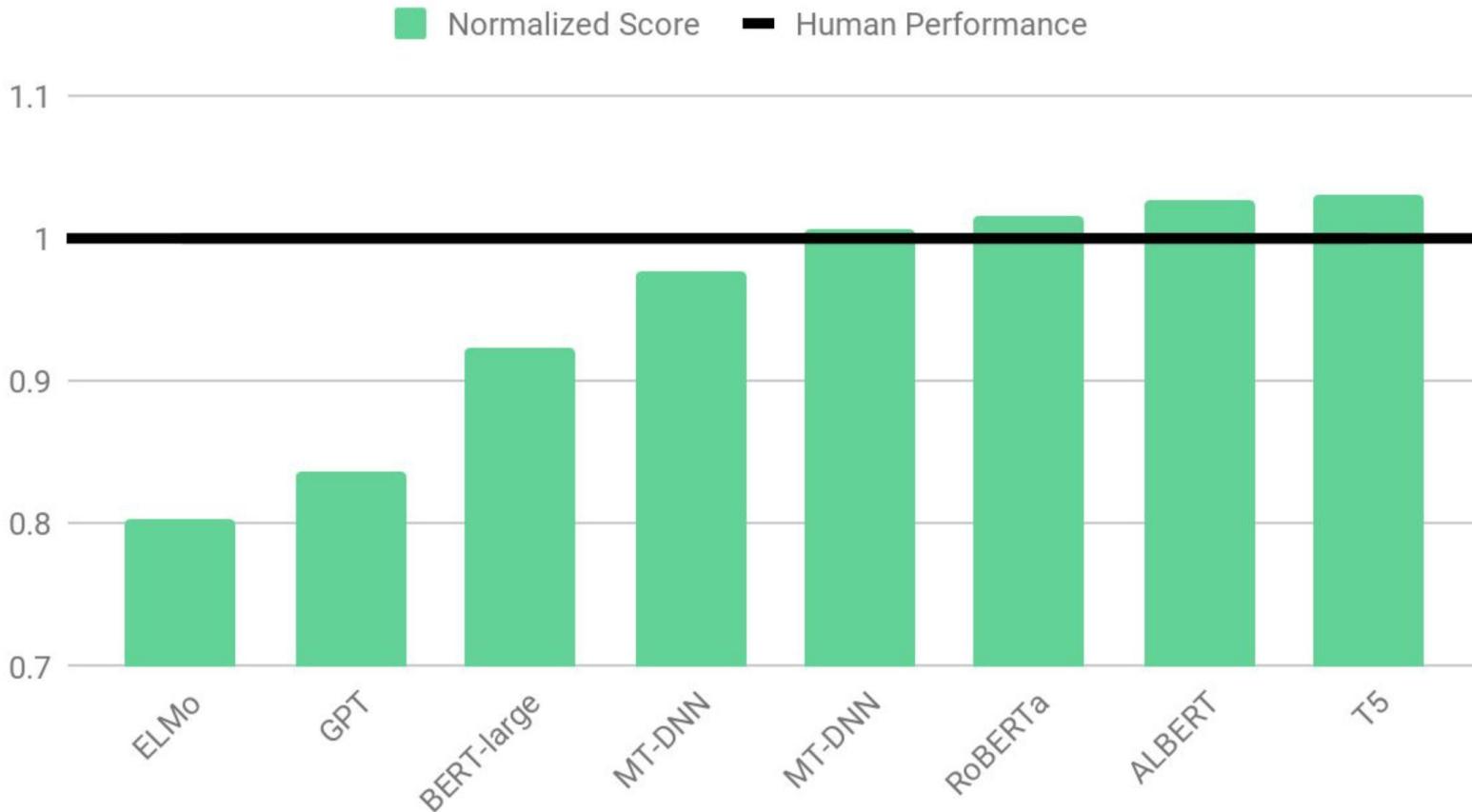
Omer Levy
Facebook AI Research

Samuel R. Bowman
New York University

First Attempt: GLUE

- Benchmark of 9 sentence- and sentence-pair classification tasks
 - Different tasks (sentiment analysis, paraphrase detection, etc.), genre, amount of data
 - Evaluate system on all nine tasks; overall score is average across tasks
- Released May 2018





SuperGLUE

- New benchmark of 8 NLU tasks
- Also:
 - Additional diagnostics
 - Rules updates
 - Starter code
- Tasks were selected from an open call to the NLP community
 - Screen each proposed task to be easy for humans, hard for machines
 - Emphasized tasks with little training data
 - More diverse set of task formats, e.g. QA, coreference
- Released May 2019



Passage: Barq's – Barq's is an American soft drink. Its brand of root beer is notable for having caffeine. Barq's, created by Edward Barq and bottled since the turn of the 20th century, is owned by the Barq family but bottled by the Coca-Cola Company. It was known as Barq's Famous Olde Tyme Root Beer until 2012.

Question: is barq's root beer a pepsi product **Answer:** No

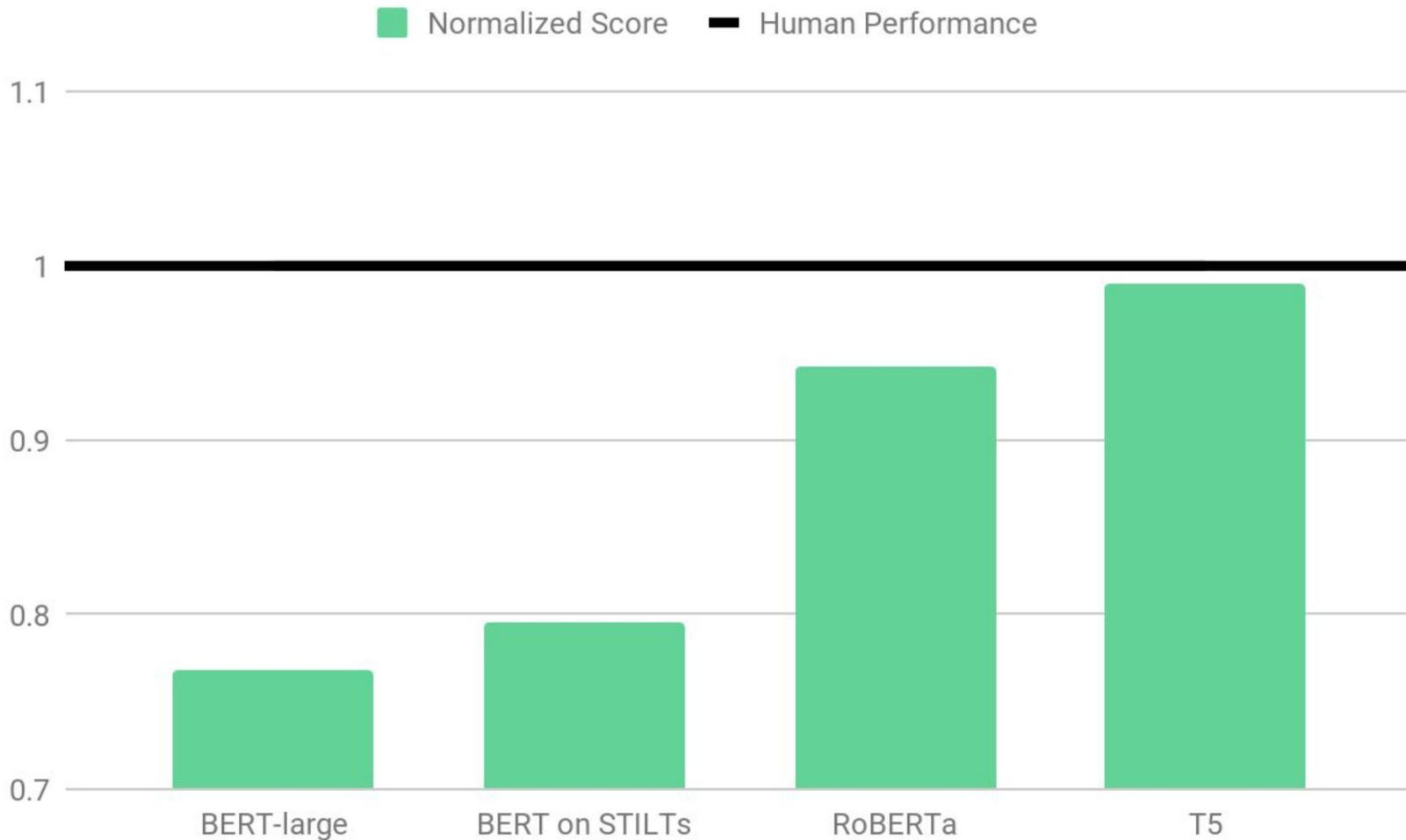
Premise: My body cast a shadow over the grass. **Question:** What's the CAUSE for this?

Alternative 1: The sun was rising. **Alternative 2:** The grass was cut.

Correct Alternative: 1

Paragraph: Susan wanted to have a birthday party. She called all of her friends. She has five friends. Her mom said that Susan can invite them all to the party. Her first friend could not go to the party because she was sick. Her second friend was going out of town. Her third friend was not so sure if her parents would let her. The fourth friend said maybe. The fifth friend could go to the party for sure. Susan was a little sad. On the day of the party, all five friends showed up. Each friend had a present for Susan. Susan was happy and sent each friend a thank you card the next week

Question: Did Susan's sick friend recover? **Candidate answers:** Yes, she recovered (T), No (F), Yes (T), No, she didn't recover (F), Yes, she was at Susan's party (T)



CondConv: Conditionally Parameterized Convolutions for Efficient Inference

Brandon Yang*

Google Brain

bcyang@google.com

Gabriel Bender

Google Brain

gbender@google.com

Quoc V. Le

Google Brain

qvl@google.com

Jiquan Ngiam

Google Brain

jngiam@google.com

CondConv

Fundamental assumption:

Convolutional kernels should be shared for all examples in a dataset

Conditionally parameterized convolutions (CondConv)

each individual example is processed with different weights!

Enables to increase the size and capacity of a network

while maintaining efficient inference

CondConv

Fundamental assumption:

Convolutional kernels should be shared for all examples in a dataset

Conditionally parameterized convolutions (CondConv)

each individual example is processed with different weights!

Enables to increase the size and capacity of a network

while maintaining efficient inference

CondConv

Fundamental assumption:

Convolutional kernels should be shared for all examples in a dataset

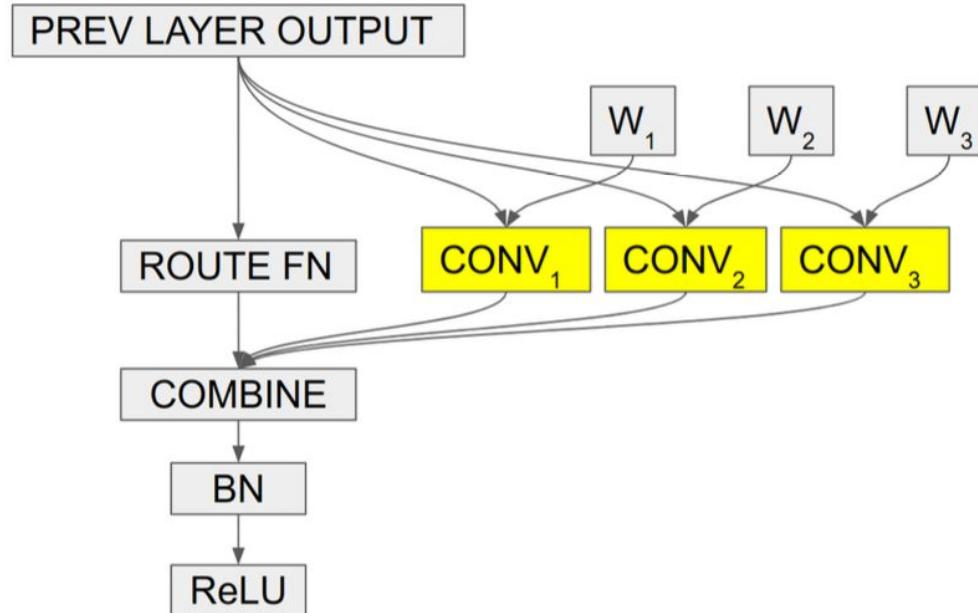
Conditionally parameterized convolutions (CondConv)

each individual example is processed with different weights!

Enables to increase the size and capacity of a network

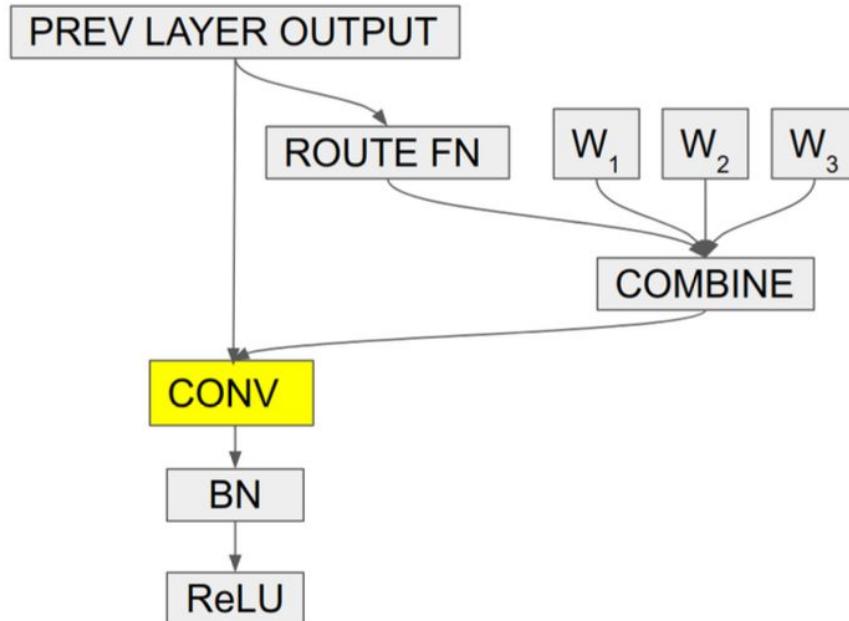
while maintaining efficient inference

CondConv



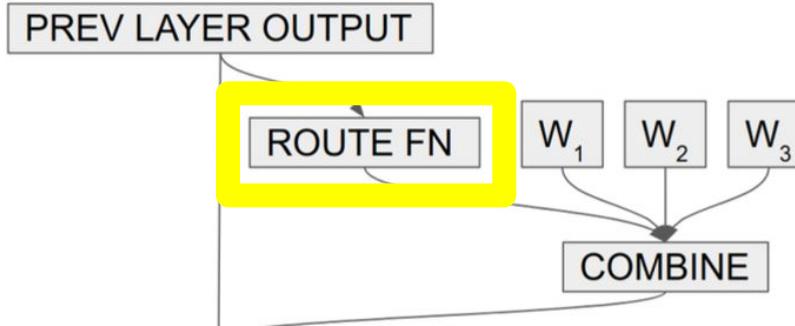
(b) Mixture of Experts: $\alpha_1(W_1 * x) + \dots + \alpha_n(W_n * x)$

CondConv



(a) CondConv: $(\alpha_1 W_1 + \dots + \alpha_n W_n) * x$

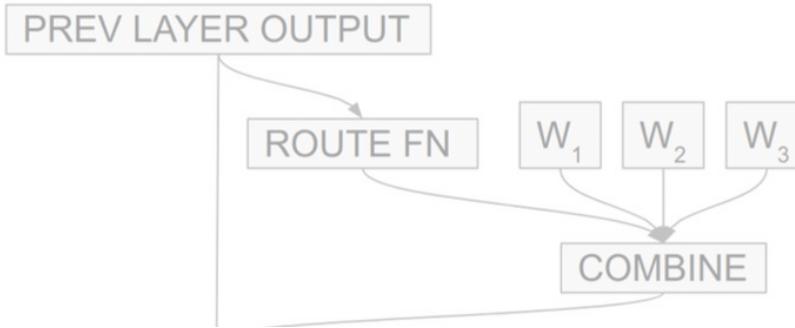
CondConv



We wish to design a per-example routing function that is computationally efficient, able to meaningfully differentiate between input examples, and is easily interpretable. We compute the example-dependent routing weights $\alpha_i = r_i(x)$ from the layer input in three steps: global average pooling, fully-connected layer, Sigmoid activation.

$$r(x) = \text{Sigmoid}(\text{GlobalAveragePool}(x) R)$$

CondConv

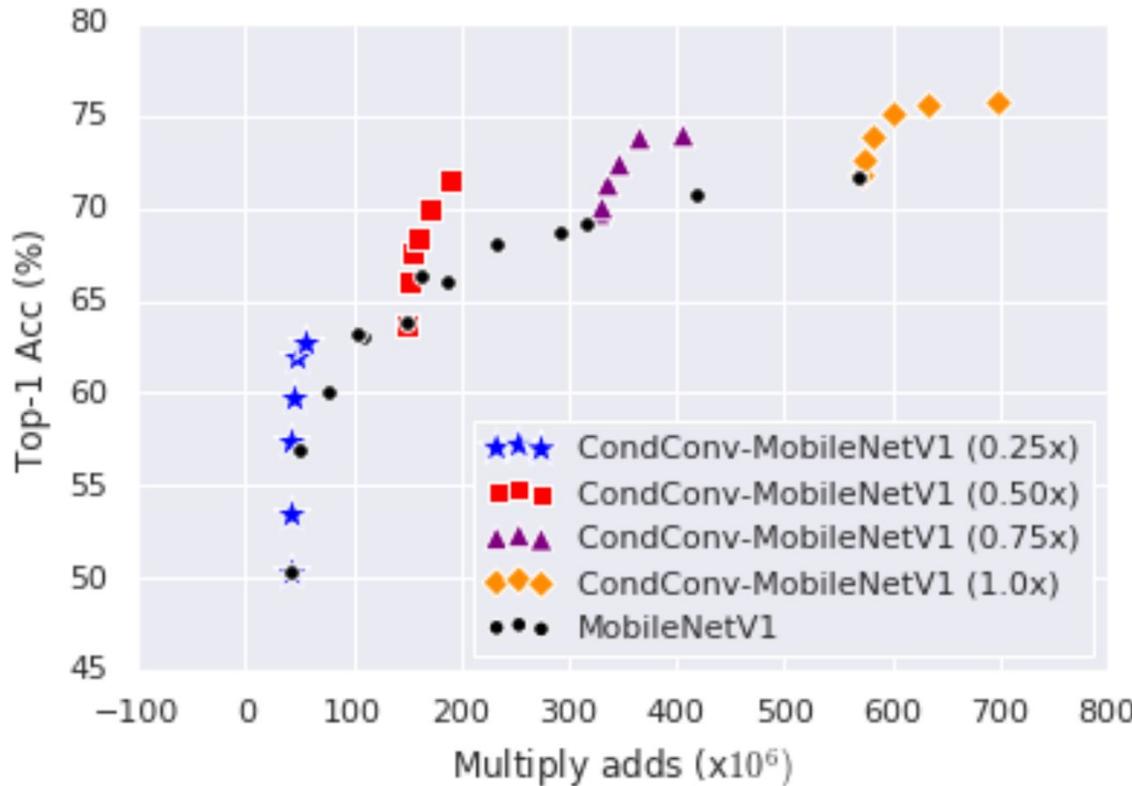


We wish to design a per-example routing function that is computationally efficient, able to meaningfully differentiate between input examples, and is easily interpretable. We compute the example-dependent routing weights $\alpha_i = r_i(x)$ from the layer input in three steps: global average pooling, fully-connected layer, Sigmoid activation.

$$r(x) = \text{Sigmoid}(\text{GlobalAveragePool}(x) R)$$

Finally, we experiment with the *Softmax* activation function to compute routing weights. The baseline's *Sigmoid* significantly outperforms *Softmax*, which suggests that multiple experts are often useful for a single example.

CondConv



This Looks Like That: Deep Learning for Interpretable Image Recognition

Chaofan Chen*

Duke University

cfchen@cs.duke.edu

Oscar Li*

Duke University

oscarli@alumni.duke.edu

Chaofan Tao

Duke University

chaofan.tao@duke.edu

Alina Jade Barnett

Duke University

abarnett@cs.duke.edu

Jonathan Su

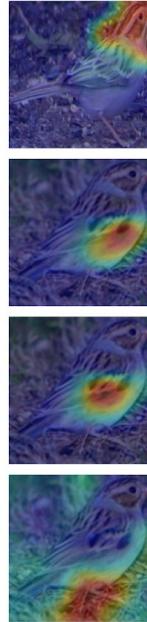
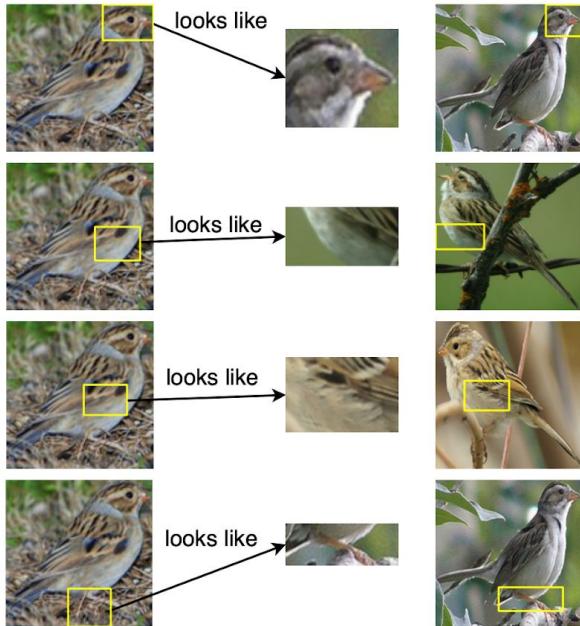
MIT Lincoln Laboratory[†]

su@ll.mit.edu

Cynthia Rudin

Duke University

cynthia@cs.duke.edu



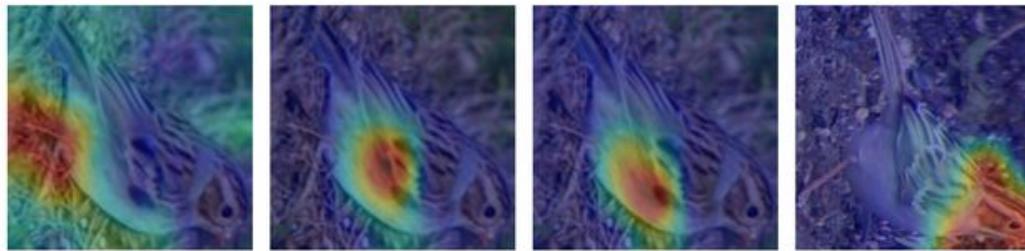
Leftmost: a test image of a clay-colored sparrow
Second column: same test image, each with a bounding box generated by our model -- the content within the bounding box is considered by our model to look similar to the prototypical part (same row, third column) learned by our algorithm

Third column: prototypical parts learned by our algorithm

Fourth column: source images of the prototypical parts in the third column

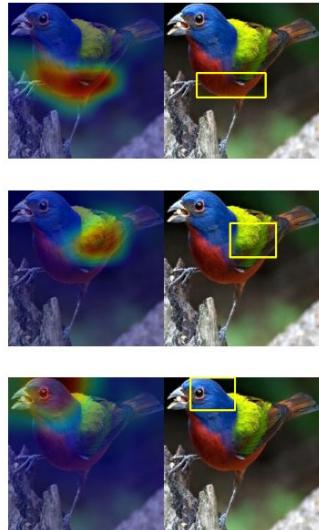
Rightmost column: activation maps indicating how similar each prototypical part resembles part of the test bird

Our work relates to (but contrasts with) those that perform *posthoc* interpretability analysis for a trained convolutional neural network (CNN). In posthoc analysis, one interprets a trained CNN by fitting explanations to how it performs classification. Examples of posthoc analysis techniques include activation maximization [6, 13, 23, 43, 31, 37, 47], deconvolution [48], and saliency visualization [37, 41, 40, 35]. All of these posthoc visualization methods do not explain the reasoning process of how a network *actually* makes its decisions. In contrast, our network has a built-in case-based reasoning process, and the explanations generated by our network are actually used during classification and are not created posthoc.

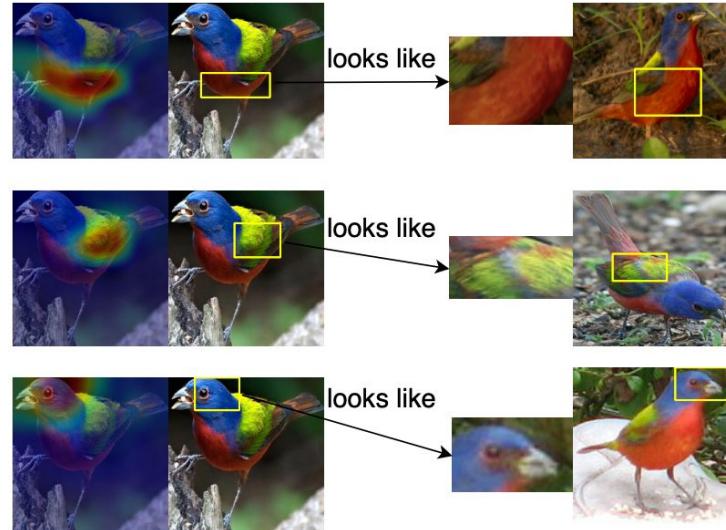




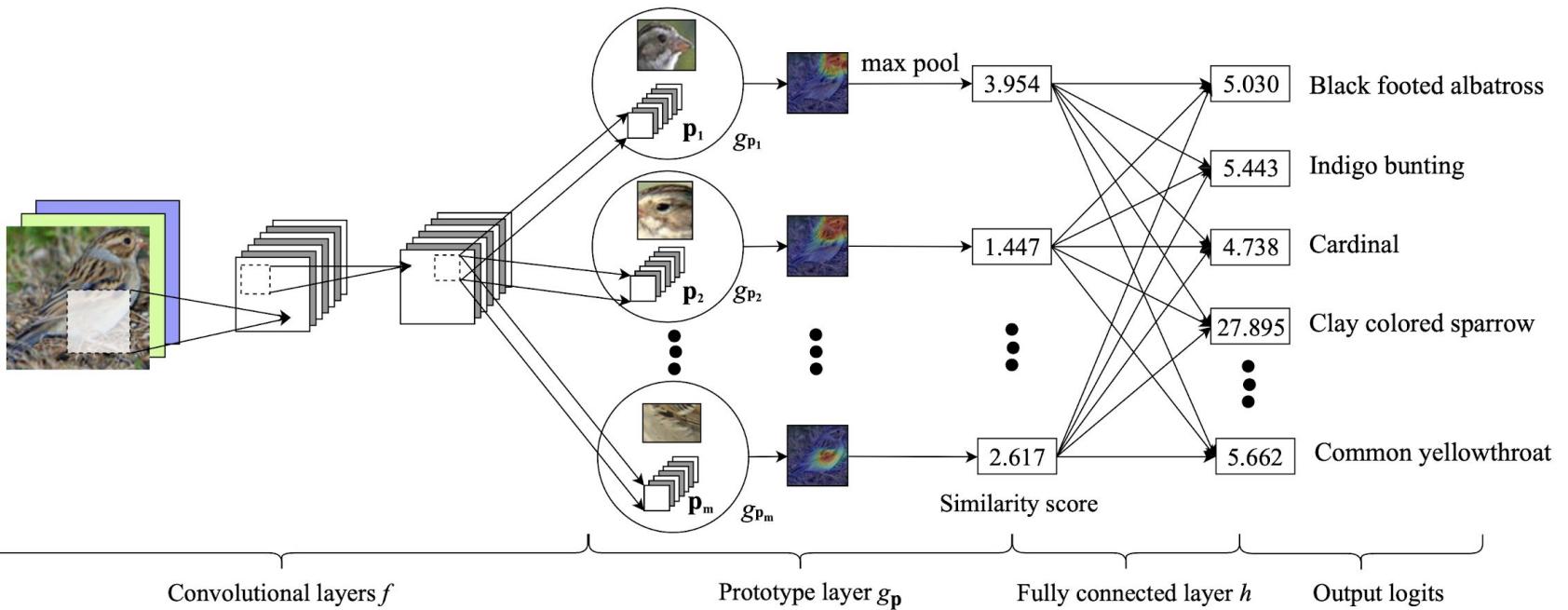
(a) Object attention
(class activation map)

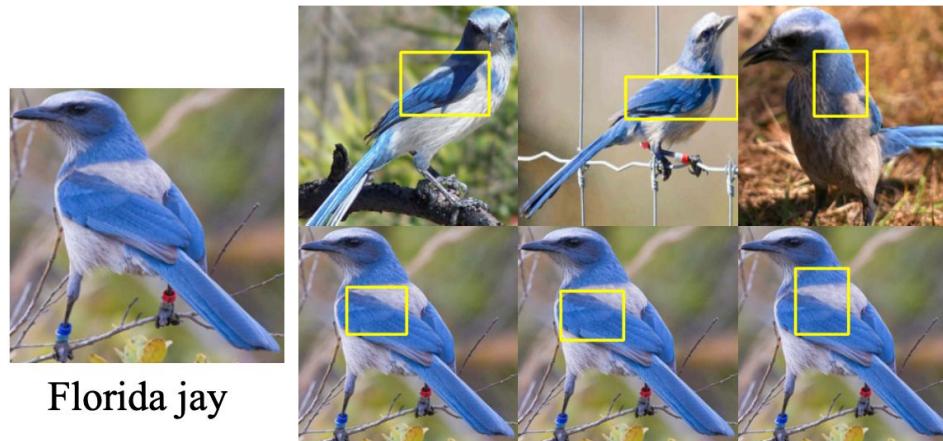


(b) Part attention
(attention-based models)

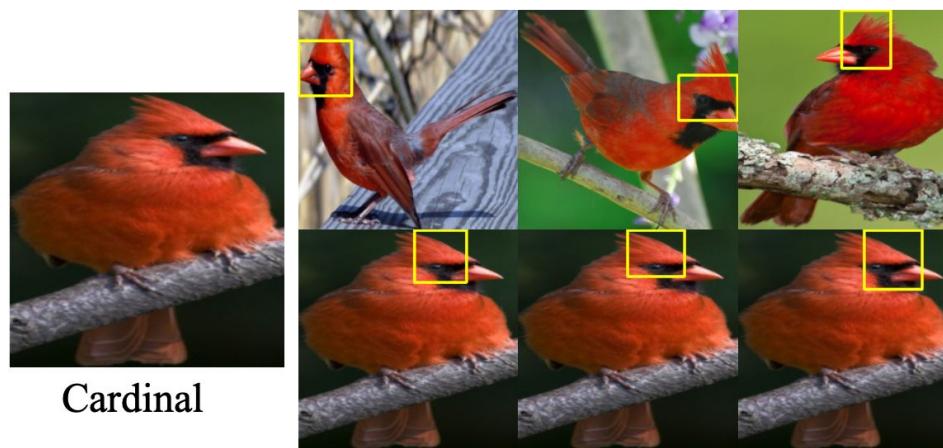


(c) Part attention + comparison with learned
prototypical parts (our model)





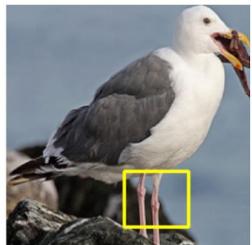
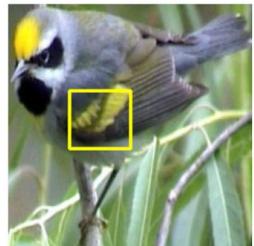
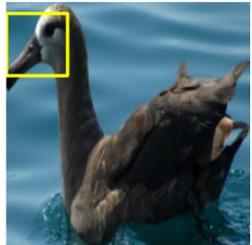
Florida jay



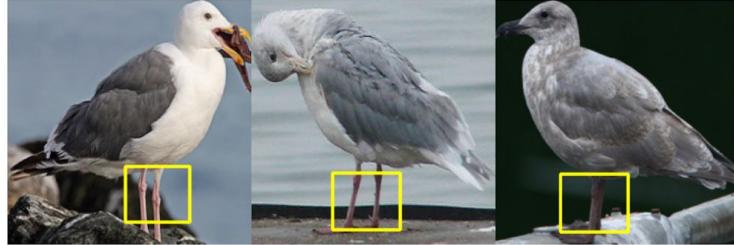
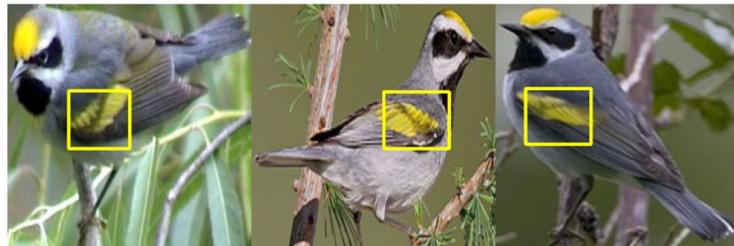
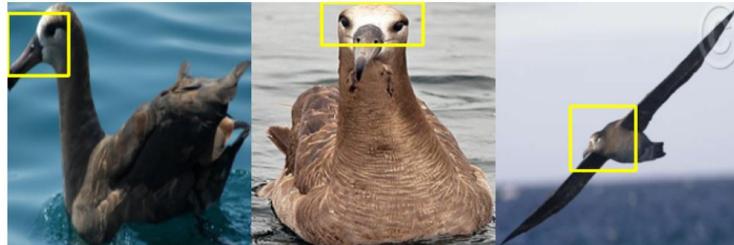
Cardinal

(a) nearest prototypes of two test images

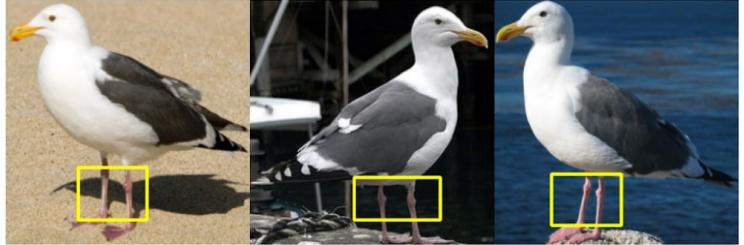
Prototype
(in bounding box)



Nearest training patches
(in bounding box)



Nearest test patches
(in bounding box)



(b) nearest image patches to prototypes

Table 1: Top: Accuracy comparison on cropped bird images of CUB-200-2011
 Bottom: Comparison of our model with other deep models

Base	ProtoPNet	Baseline	Base	ProtoPNet	Baseline
VGG16	76.1 ± 0.2	74.6 ± 0.2	VGG19	78.0 ± 0.2	75.1 ± 0.4
Res34	79.2 ± 0.1	82.3 ± 0.3	Res152	78.0 ± 0.3	81.5 ± 0.4
Dense121	80.2 ± 0.2	80.5 ± 0.1	Dense161	80.1 ± 0.3	82.2 ± 0.2

An aerial photograph of the Vancouver skyline during sunset. The city is built on a peninsula, with a dense cluster of skyscrapers on the right and a more spread-out residential area on the left. In the foreground, a large white geodesic dome, likely the Canada Place convention center, sits on a pier extending into the dark blue water. Several small white water taxis are docked at a platform in the lower right corner. The sky is a clear, pale blue.

Overview
Tutorials
Papers
Workshops & Hardware

- Black in AI (BAI) Affinity Workshop
Women in Machine Learning (WiML) Affinity Workshop
LatinX in AI (LAI) Affinity Workshop
New In Machine Learning
Queer in AI (QAI) Affinity Workshop
CiML 2019: Machine Learning Competitions for All
Solving inverse problems with deep networks: New architectures, theoretical foundations, and applications
Retrospectives: A Venue for Self-Reflection in ML Research
Biological and Artificial Reinforcement Learning
Robust AI in Financial Services: Data, Fairness, Explainability, Trustworthiness, and Privacy
Bayesian Deep Learning
- MLSys: Workshop on Systems for ML**
- Information Theory and Machine Learning
KR2ML - Knowledge Representation and Reasoning Meets Machine Learning
Machine Learning for Health (ML4H): What makes machine learning in medicine different?
Competition Track Day 1
Workshop on Human-Centric Machine Learning
Safety and Robustness in Decision-making
Visually Grounded Interaction and Language
Machine Learning for the Developing World (ML4D): Challenges and Risks
Minding the Gap: Between Fairness and Ethics
Shared Visual Representations in Human and Machine Intelligence
Graph Representation Learning
Beyond first order methods in machine learning systems
Perception as generative reasoning: structure, causality, probability
EMC2: Energy Efficient Machine Learning and Cognitive Computing (5th edition)
Optimal Transport for Machine Learning
AI for Humanitarian Assistance and Disaster Response
Workshop on Federated Learning for Data Privacy and Confidentiality
Learning Meaningful Representations of Life
Learning with Rich Experience: Integration of Learning Paradigms
Meta-Learning
Real Neurons & Hidden Units: future directions at the intersection of neuroscience and AI
Fair ML in Healthcare
- Document Intelligence
Deep Reinforcement Learning
Privacy in Machine Learning (PriML)
Robot Learning: Control and Interaction in the Real World
Tackling Climate Change with ML
Medical Imaging meets NeurIPS
Learning Transferable Skills
Machine Learning with Guarantees
Machine Learning and the Physical Sciences
Emergent Communication: Towards Natural Language
Context and Compositionality in Biological and Artificial Neural Systems
Sets and Partitions
The third Conversational AI workshop – today's practice and tomorrow's potential
Program Transformations for ML
- ML For Systems**
- NeurIPS Workshop on Machine Learning for Creativity and Design 3.0
Learning with Temporal Point Processes
Machine Learning for Autonomous Driving
Joint Workshop on AI for Social Good
“Do the right thing”: machine learning and causal inference for improved decision making
The Optimization Foundations of Reinforcement Learning
Science meets Engineering of Deep Learning
Competition Track Day 2
Bridging Game Theory and Deep Learning

An aerial photograph of the Vancouver skyline during sunset. The city is built on a peninsula, with a dense cluster of skyscrapers on the right and a more spread-out residential area on the left. In the foreground, a large white geodesic dome, likely the Canada Place convention center, sits on a pier extending into the dark blue water. Several small white water taxis are docked at a platform further out. The sky is a clear, pale blue.

Overview
Tutorials
Papers
Workshops & Hardware

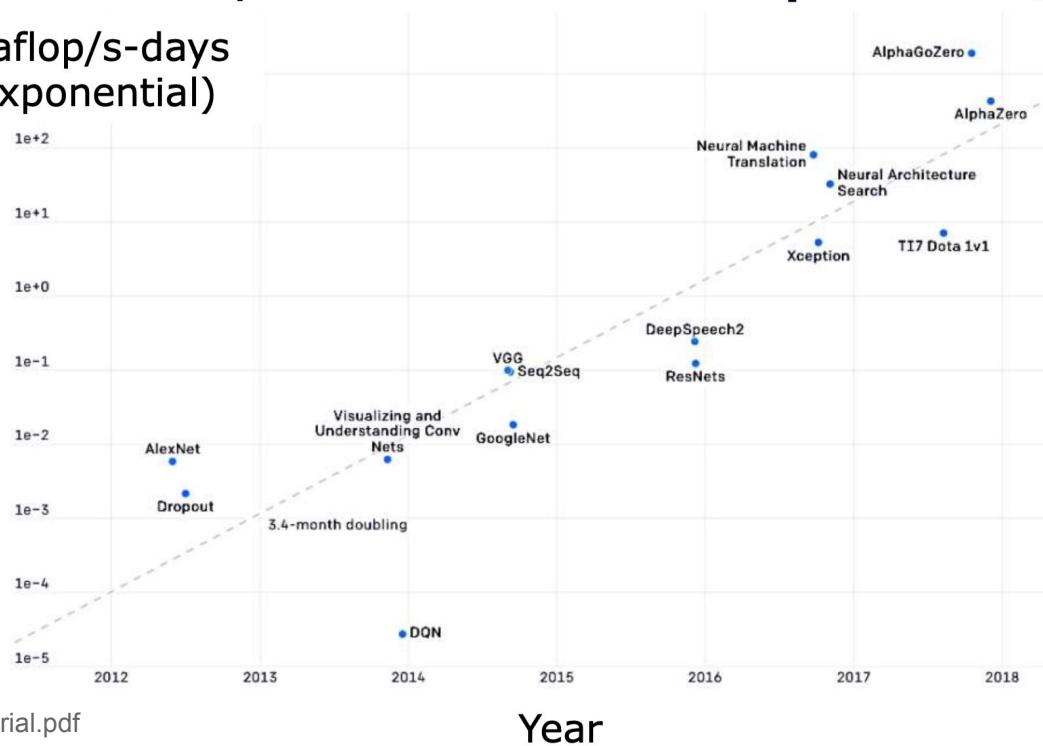
Hardware

Efficient Processing of Deep Neural Network: from Algorithms to Hardware Architectures

Vivienne Sze

AlexNet to AlphaGo Zero: A 300,000x Increase in Compute

Petaflop/s-days
(exponential)

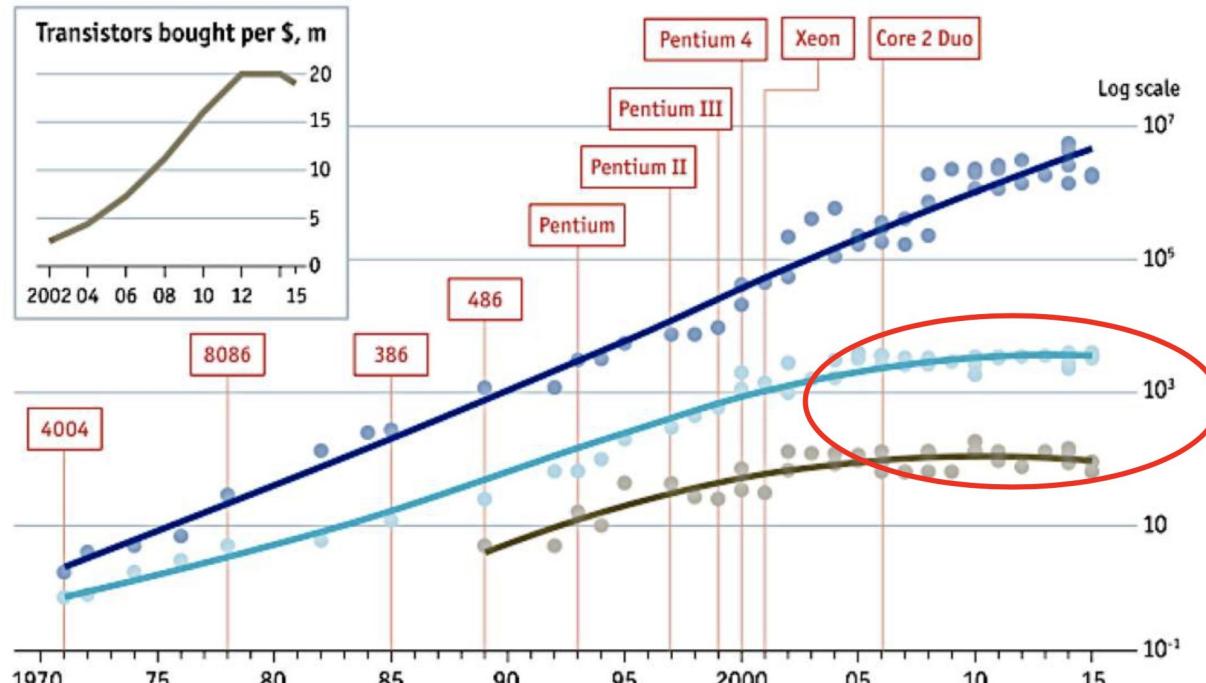


Future Hardware

Stuttering

● Transistors per chip, '000 ● Clock speed (max), MHz ● Thermal design power*, w

Chip introduction
dates, selected



Sources: Intel; press reports; Bob Colwell; Linley Group; IB Consulting; *The Economist*

*Maximum safe power consumption

Future Hardware

Intel AI: Intel® Nervana™ NNP: domain specific architectures for inference & training

Graphcore: Innovative approaches in training large scale language models

Cerebras Systems: Accelerating deep learning at wafer scale

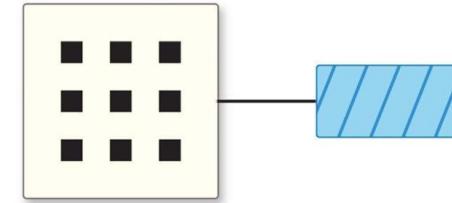
Habana Labs: Hands-on Workshop:
Implementing High-Performance AI
Workloads with Habana AI Processors

Traditional Memory Architectures not Optimized for DL

In neural networks, weights and activations are local, with low reuse

Traditional memory designs are punished

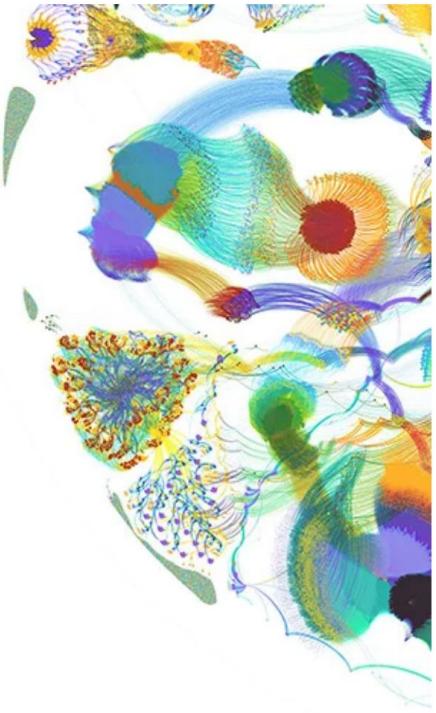
- Central shared memory is slow & far away
- Requires high data reuse (caching)
- Fundamental operation (matrix*vector) has low data reuse



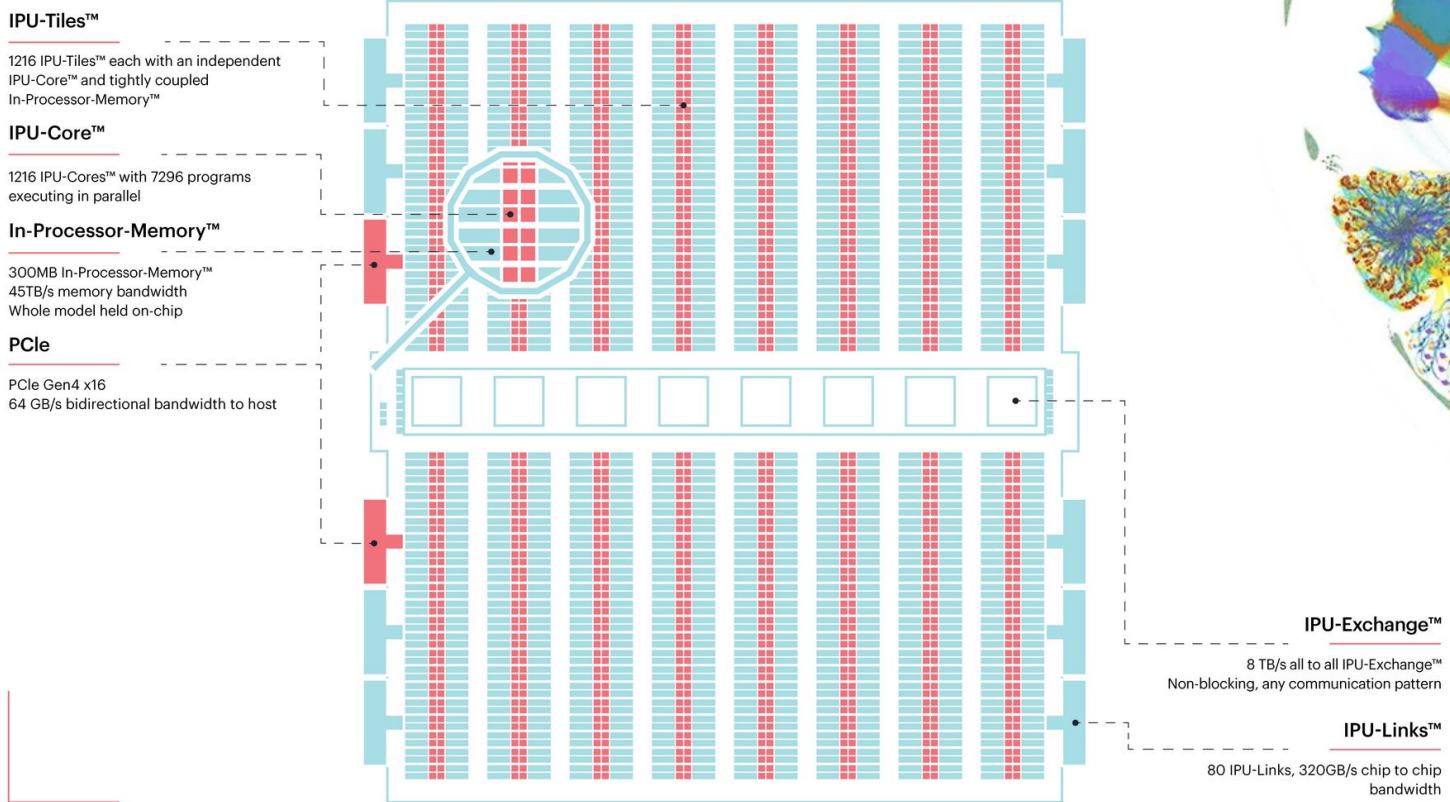
Memory separate from cores

■ Core / Memory

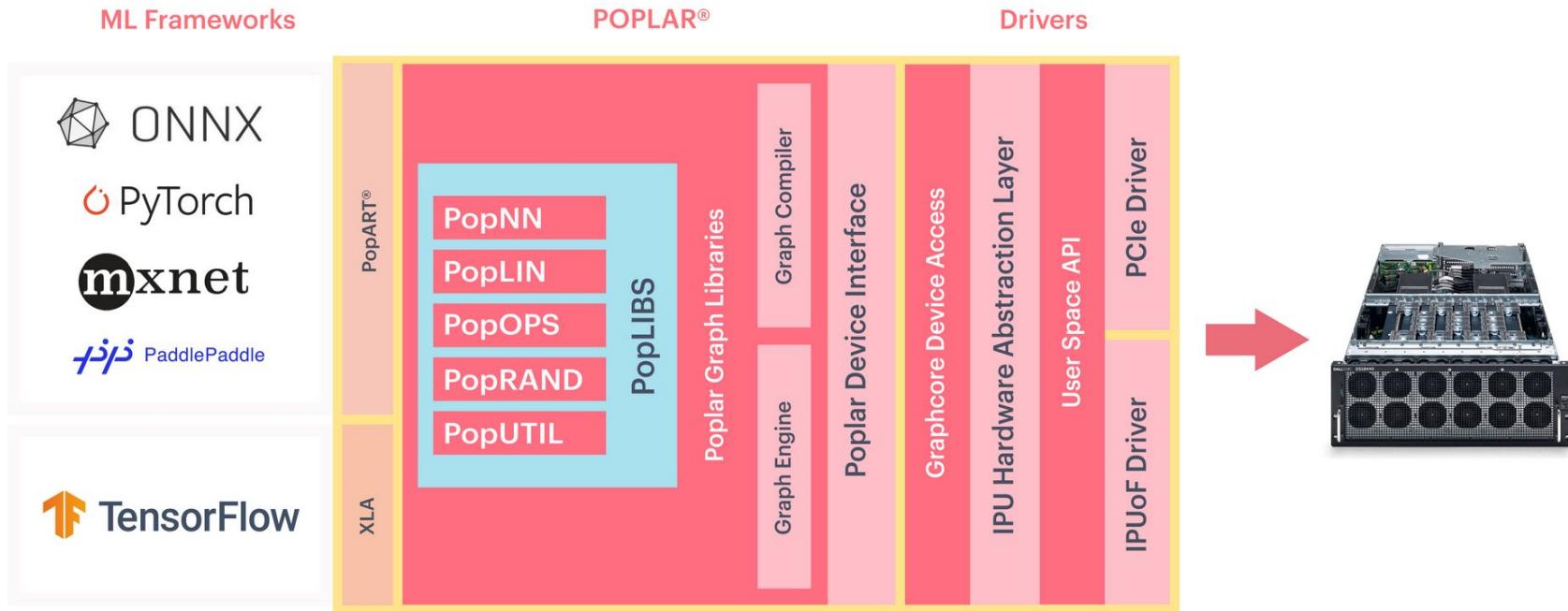
Graphcore



Graphcore



Graphcore



Graphcore



MICROSOFT AZURE IPU CLOUD PREVIEW IS HERE

The world's first Graphcore IPU cloud service is now available on Microsoft Azure, letting innovators around the world create new breakthroughs in machine intelligence.

[Get started →](#)

A promotional graphic for the Microsoft Azure Graphcore IPU Cloud Preview. It features a dark, abstract background with a complex, glowing network of nodes and connections, similar to a neural network or a cloud. In the center, there is a red rectangular callout box. Inside the box, the text "MICROSOFT AZURE IPU CLOUD PREVIEW IS HERE" is displayed in white, bold, uppercase letters. Below this, a smaller text block reads "The world's first Graphcore IPU cloud service is now available on Microsoft Azure, letting innovators around the world create new breakthroughs in machine intelligence." At the bottom of the red box is a white button with the text "Get started" and a right-pointing arrow. To the right of the arrow is a small white circle containing a red downward-pointing arrow.

Cerebras



Hotchips 2019

4:45	ML Training	Session Chair: Cliff	
PM		Young	
4:45	A Scalable unified architecture for Neural Network computing from Nano-level to high performance computing	Liao Heng	Huawei
5:15	Deep Learning Training at Scale – Spring Crest Deep Learning Accelerator	Andrew Yang, Nitin Garegrat, Connie Miao & Karthik Vaidyanathan	Intel
PM			
5:45	Wafer Scale Deep Learning	Sean Lie	Cerebras
PM			
6:15	Habana Labs Approach to Scaling AI Training	Eitan Medina	Habana
PM			

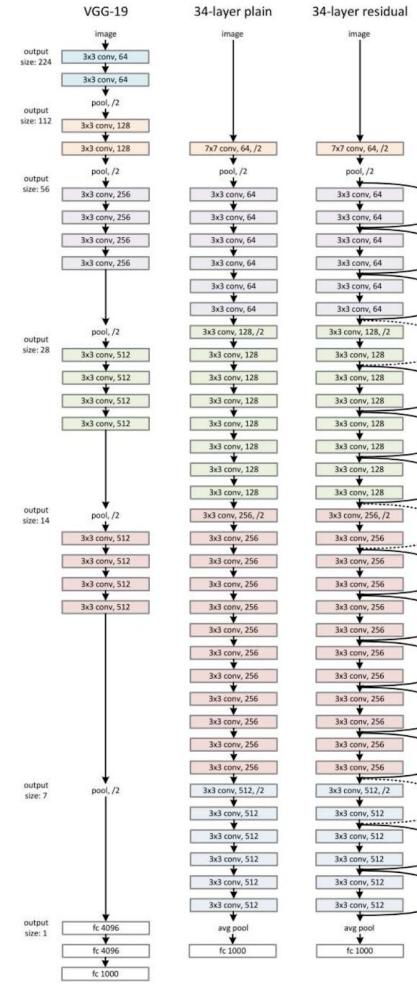
Deep Learning Training is Hard

Size:

- Billions-trillions of ops per sample
- Millions-billions of samples per training
- Peta-exa scale compute

Shape:

- **Fine-grained:** A lot of parallelism; presents opportunity to accelerate
- **Coarse-grained:** Inherently serial



Legacy Technologies: Brute Force Parallelism

Fine-grained

- Dense vector processors (e.g. GPUs)
- Limited when compute not large uniform blocks

Coarse-grained

- Scale out clustering (e.g. PCIe, Ethernet, IB, NVLink)
- Run multiple instances of the same model (data parallel)
- Limited by inherent serial nature of problem

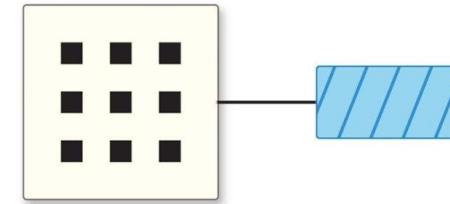
Result: scaling is limited and costly

Traditional Memory Architectures not Optimized for DL

In neural networks, weights and activations are local, with low reuse

Traditional memory designs are punished

- Central shared memory is slow & far away
- Requires high data reuse (caching)
- Fundamental operation (matrix*vector) has low data reuse



Memory separate from cores

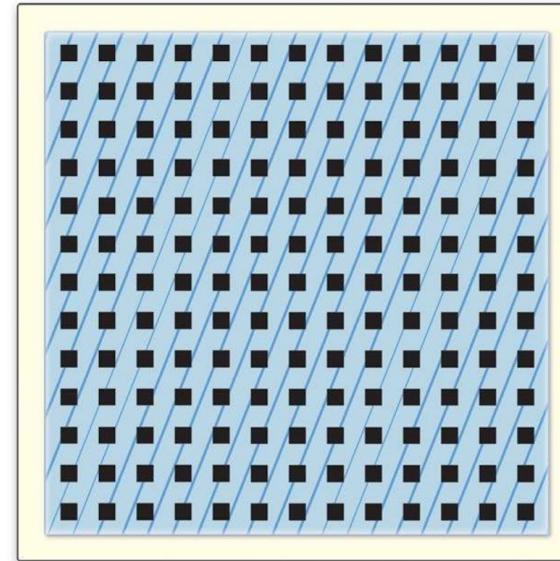
■ Core / Memory

A Memory Architecture that *is* Optimized for DL

In neural networks, weights and activations are local, with low data reuse

**The right answer is distributed,
high performance, on-chip memory**

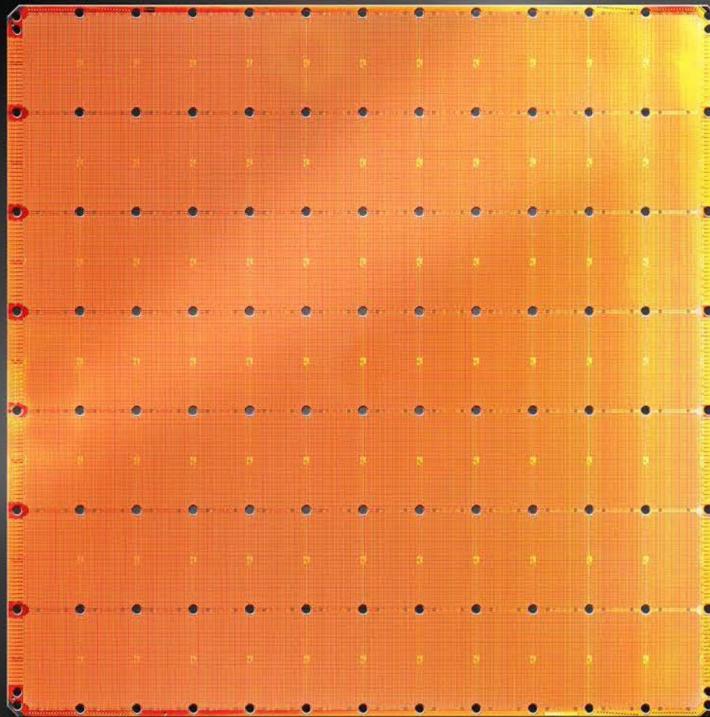
- All memory is fully distributed along with compute datapaths
- Datapath has full performance from memory



Memory uniformly distributed across cores

■ Core ■ Memory

Cerebras Wafer Scale Engine



Cerebras WSE

1.2 Trillion Transistors
46,225 mm² Silicon

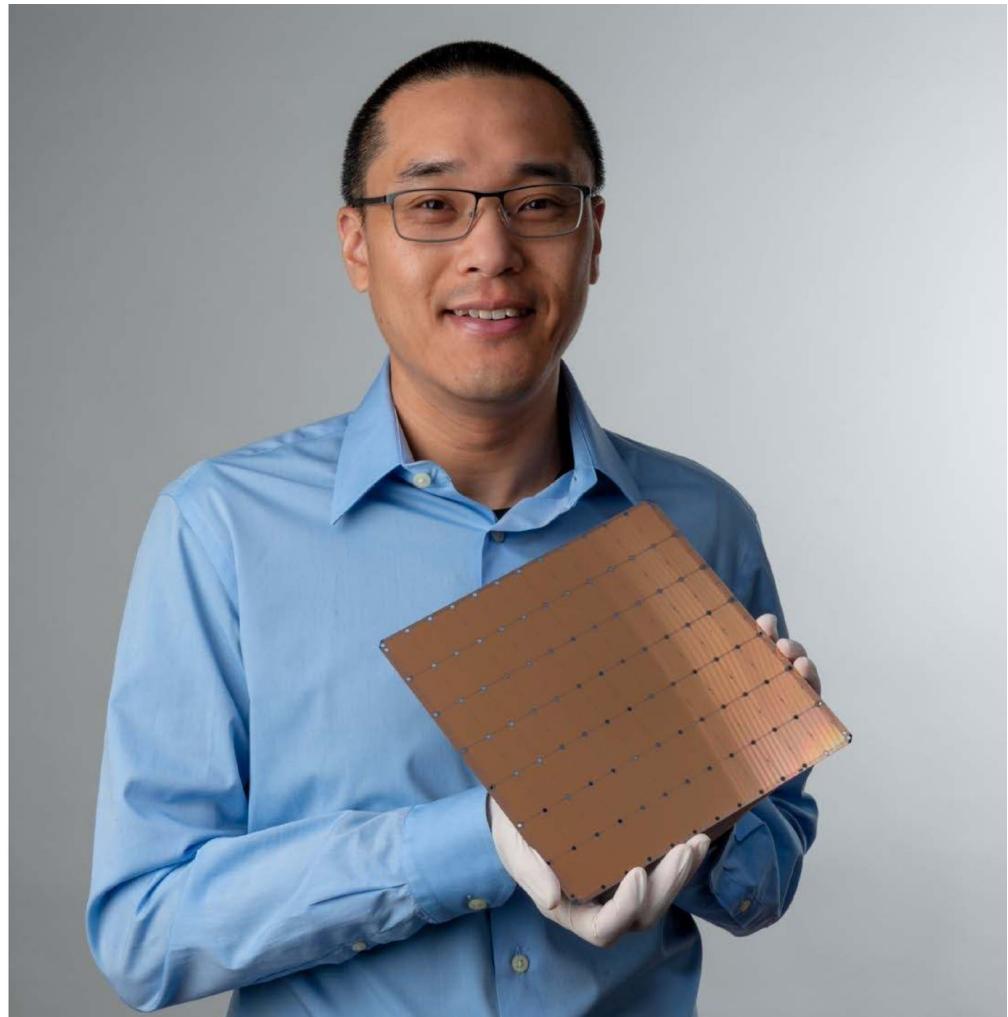


Largest GPU

21.1 Billion Transistors
815 mm² Silicon

Largest Chip Ever Built

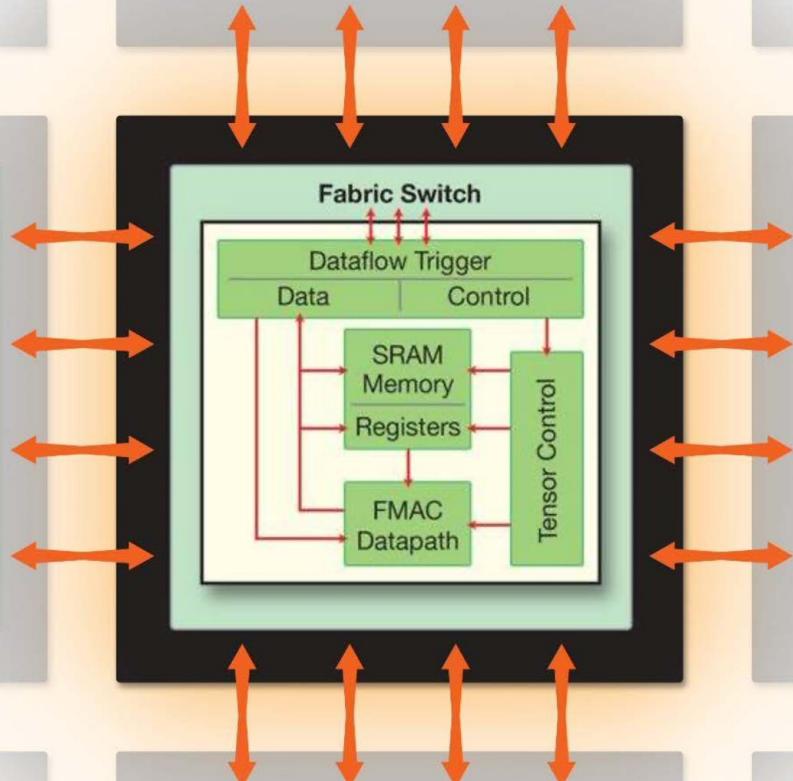
- 46,225 mm² silicon
- 1.2 trillion transistors
- 400,000 AI optimized cores
- 18 Gigabytes of On-chip Memory
- 9 PByte/s memory bandwidth
- 100 Pbit/s fabric bandwidth
- TSMC 16nm process



Flexible Cores Optimized for Tensor Operations

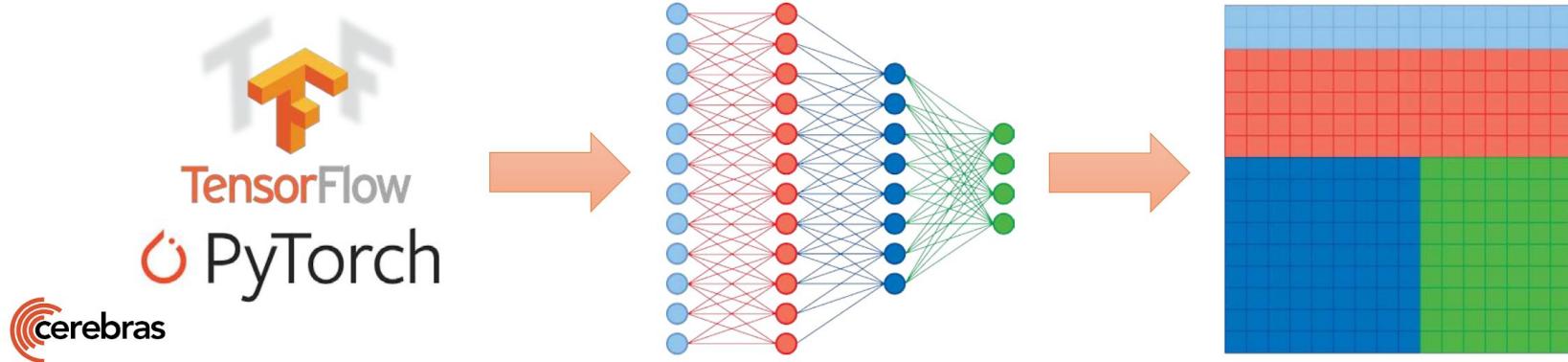
Key to supporting rapidly evolving NN architectures

- Fully programmable compute core
- Full array of general instructions with ML extensions
- Flexible **general ops** for control processing
 - e.g. arithmetic, logical, load/store, branch
- Optimized **tensor ops** for data processing
 - Tensors as first class operands
 - e.g. `fmac [z] = [z], [w], a`
3D 3D 2D scalar



Programming the Wafer Scale Engine

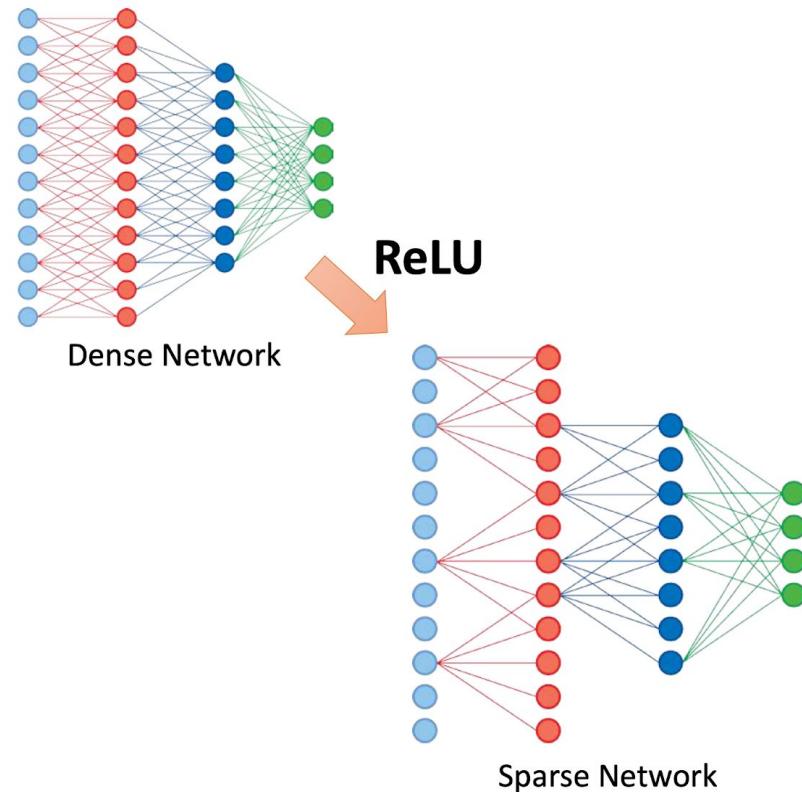
- Neural network models expressed in common ML frameworks
- Cerebras interface to framework extracts the neural network
- Performs placement and routing to map neural network layers to fabric
- The entire wafer operates on the single neural network



Sparse Compute Engine for Neural Networks

NN operations like nonlinearities naturally create fine-grained sparsity

- Native, sparse processing enables higher efficiency and performance
- Dataflow scheduling in hardware
 - Triggered by data
 - Filters out sparse zero data
 - Skips unnecessary processing
- Fine-grained execution datapaths
 - Small cores with independent instructions
 - Maximizes utilization
 - Efficiently processes dynamic, non-uniform work





**Cluster-scale Deep
Learning compute in a
single system**

15 Rack Units

Fits in a standard datacenter rack

1.2 Terabits/sec

System IO over 12x standard 100 GbE

20 kW

Maximum power draw

[Explore our product](#)



It's working,
running customer workloads.

Stay tuned...



NeurIPS 2019 recap

René Donner, contextflow