# Computer Vision for Embedded Systems

Yung-Hsiang Lu
Purdue University
yunglu@purdue.edu

**PURDUE**
UNIVERSITY.

# **Evaluate Computer Vision**

# *Evaluating Computer Vision*

For many people, the only metric is the accuracy using a specific dataset. Even this leaves many questions:

- Which dataset is used?
- Why is this dataset chosen?
- How is accuracy defined?
- What other methods are compared?

Yung-Hsiang Lu, Purdue University

**Artificial intelligence** / Machine learning

# Training a single AI model can emit as much carbon as five cars in their lifetimes

Deep learning has a terrible carbon footprint.

by **Karen Hao**

June 6, 2019

Yung-Hsiang Lu, Purdue University

| Consumption | $CO_2e$ (lbs) |
| --- | --- |
| Air travel, 1 passenger, NY↔SF | 1984 |
| Human life, avg, 1 year | 11,023 |
| American life, avg, 1 year | 36,156 |
| Car, avg incl. fuel, 1 lifetime | 126,000 ⇦ |

| Training one model (GPU) | |
| --- | --- |
| NLP pipeline (parsing, SRL) | 39 |
| w/ tuning & experimentation | 78,468 |
| Transformer (big) | 192 |
| w/ neural architecture search | 626,155 ⇦ |

Emma Strubell, Ananya Ganesh, Andrew McCallum, "Energy and Policy Considerations for Deep Learning in NLP" 2019

# IMAGE RECOGNITION

**16X** Model

8 layers
1.4 GFLOP
~16% Error

152 layers
22.6 GFLOP
~3.5% error

**2012** AlexNet

**2015** ResNet

Microsoft

# SPEECH RECOGNITION

**10X** Training Ops

80 GFLOP
7,000 hrs of Data
~8% Error

465 GFLOP
12,000 hrs of Data
~5% Error

**2014** Deep Speech 1

**2015** Deep Speech 2

Baidu

**(Training)**

Source: cs231n.stanford.edu/slides/2017/cs231n_2017_lecture15.pdf

# *Transmit all data from cameras to to servers?*

- latency: wireless signals travel at 3.33 microseconds/km
- data rates:

  - Bluetooth up to 3Mb/s, up to 10 meters

  - Wifi (802.11ax) up to 2.4 Gb/s, 70 meters (indoors), 240 outdoors

  - 5G up to 20Gbps, 500 meters
- power: omnidirectional antenna - power proportional to the square of distance. directional antenna can be much more efficient
- privacy: who owns the servers? is data encrypted?
- Homomorphic encryption is not ready yet.

Yung-Hsiang Lu, Purdue University

# *ATT and T Mobile Coverage*

Yung-Hsiang Lu, Purdue University

# *Precision and Recall*

source: wikipedia

# *Factors and metrics for performance*

- Accuracy: precision, recall, top-3, top-5, hierarchical
- Execution time: per image (or video frame)
- FPS: frames per second
- FLOPS: floating-point operations
- Memory: to store machine learning model and to process data
- Resolution: number of pixels (width x height)
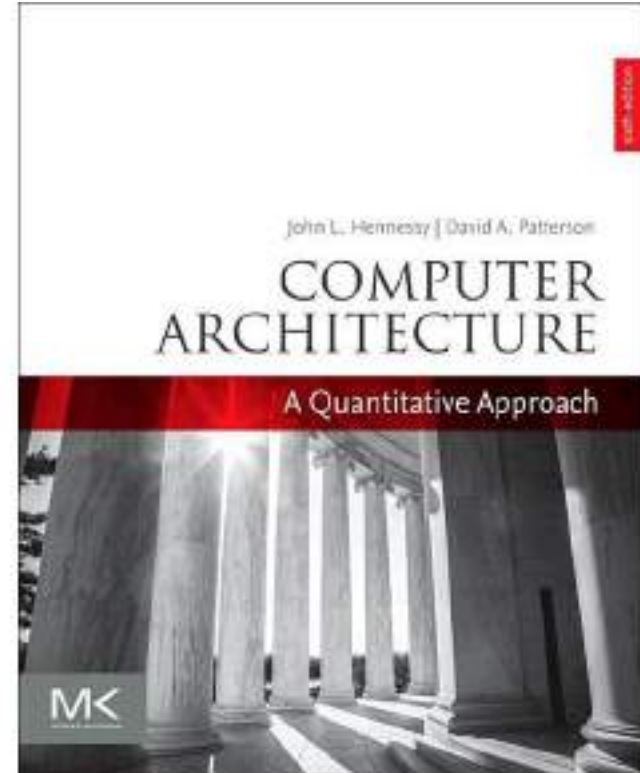


"Skynet processing at 60 Teraflops a second."
- Movie Terminator 3

Yung-Hsiang Lu, Purdue University

# *Tuning Parameters (depend on applications)*

- Resolution: how many pixels are needed?
- Frame rate: do you really need 30 frames per second?
- Accuracy: estimating crowd or recognizing faces for secure areas?
- General or special purpose?
- Layers of neural networks
- Size of convolution filters

Yung-Hsiang Lu, Purdue University

# *Measuring performance can be complex*

GFLOPS/second does not directly translate to performance (execution time)

- Integer operations
- Pipeline processors
- Memory hierarchy
- Thermal throttling
- ...

John L. Hennessy | David A. Patterson

COMPUTER ARCHITECTURE

A Quantitative Approach

MK

Yung-Hsiang Lu, Purdue University
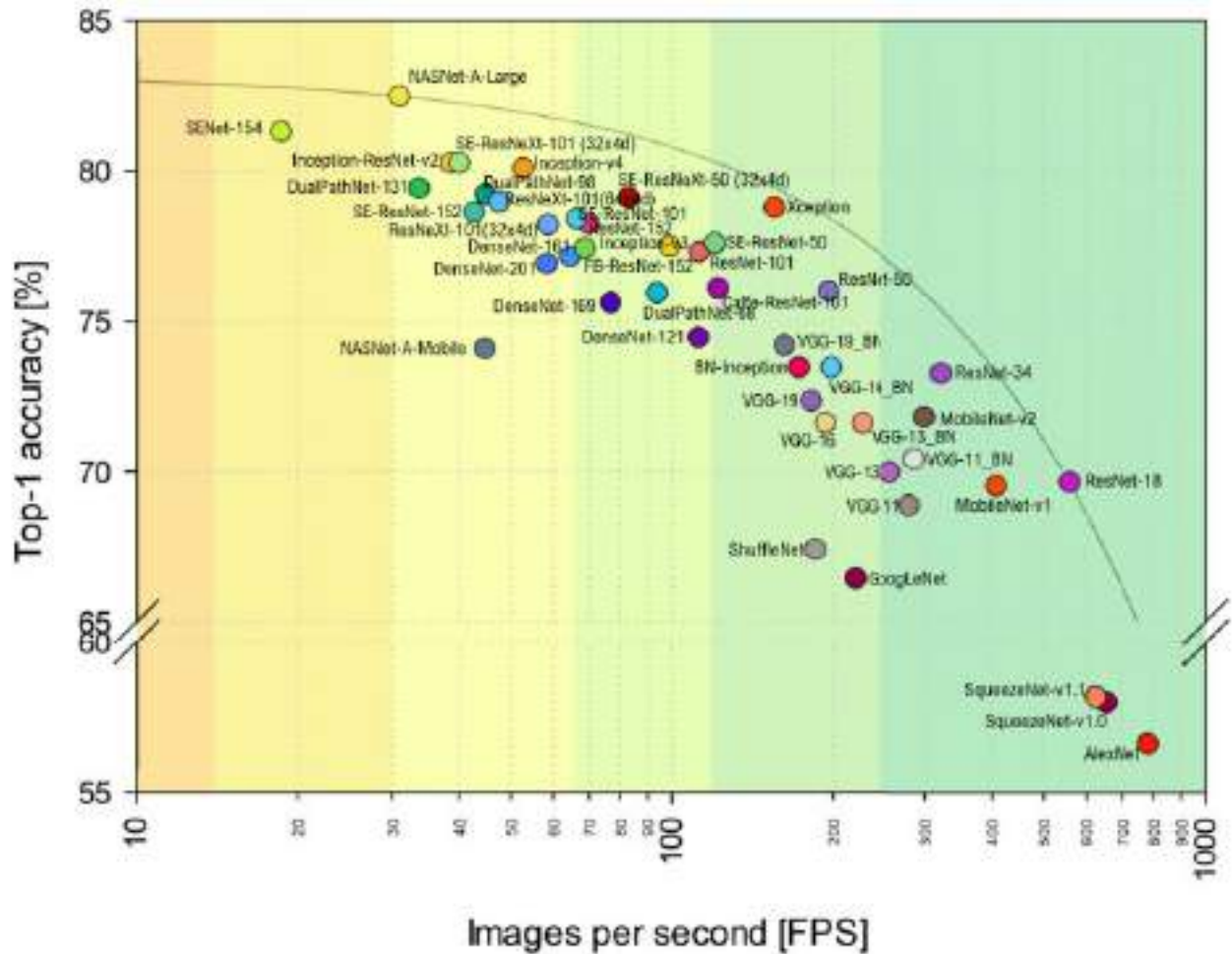
# Analysis of Neural Networks

# *Analysis of Neural Networks*

- Which neural networks are analyzed?
- What metrics are used?
- How do these networks perform?
- What patterns can be observed?

Yung-Hsiang Lu, Purdue University

Benchmark Analysis of Representative Deep Neural Network Architectures
10.1109/ACCESS.2018.2877890

Top-1 accuracy density [%/M-params]

# Compare Networks by Training

| Application | Model | Number of Layers | Dominant Layer | Implementations | Dataset |
|---|---|---|---|---|---|
| Image classification | ResNet-50 [56] Inception-v3 [80] | 50 (152 max) 42 | CONV | TensorFlow, MXNet, CNTK | ImageNet1K [73] |
| Machine translation | Seq2Seq [79] Transformer [82] | 5 12 | LSTM Attention | TensorFlow, MXNet TensorFlow | IWSLT15 [21] WMT-14 [18] |
| Object detection | Faster R-CNN [71] | 101[a] | CONV | TensorFlow, MXNet | Pascal VOC 2007 [37] |
| Speech recognition | Deep Speech 2 [13] | 9[b] | RNN | MXNet | LibriSpeech [64] |
| Adversarial learning | WGAN [40] | 14+14[c] | CONV | TensorFlow | Downsampled ImageNet [29] |
| Deep reinforcement learning | A3C [62] | 4 | CONV | MXNet | Atari 2600 |

| Dataset | Number of Samples | Size | Special |
|---|---|---|---|
| ImageNet1K | 1.2million | 3x256x256 per image | N/A |
| IWSLT15 | 133k | 20-30 words long per sentence | vocabulary size of 17188 (English to Vietnamese) |
| WMT-14 | 4.5million | up to 50 words (most sentences) | vocabulary size of 37000 (English to German) |
| Pascal VOC 2007 | 5011[d] | around 500x350 | 12608 annotated objects |
| LibriSpeech | 280k | 1000 hours[e] | N/A |
| Downsampled ImageNet | 1.2million | 3x64x64 per image | N/A |
| Atari 2600 | N/A | 4x84x84 per image | N/A |

Fig. 5: GPU compute utilization for different models on multiple mini-batch sizes.

(c) Seq2Seq



(f) Deep Speech 2



(c) Seq2Seq


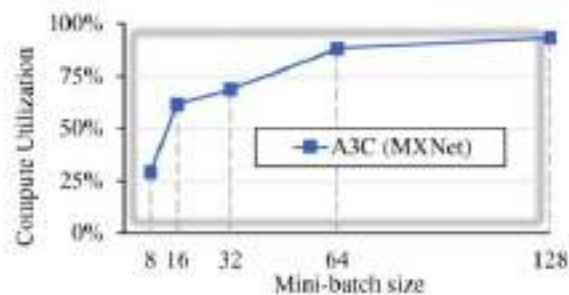
(f) Deep Speech 2

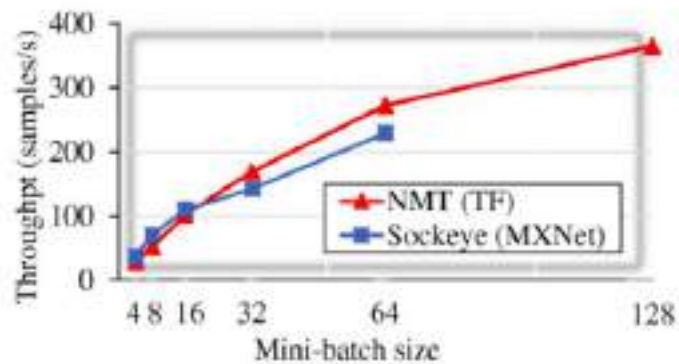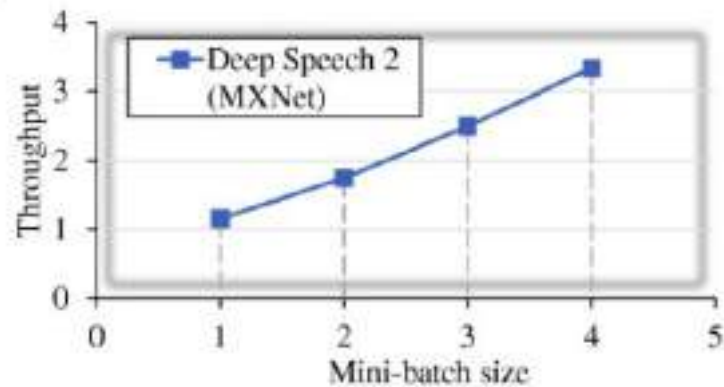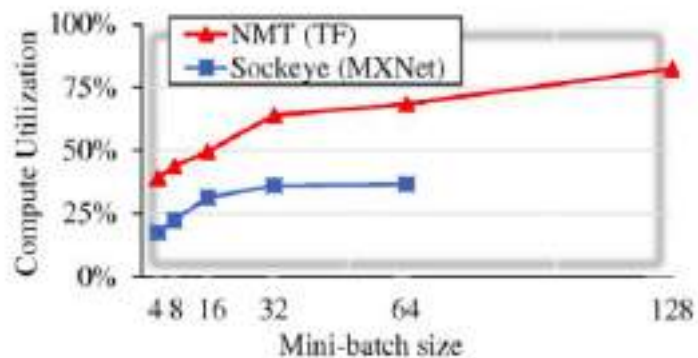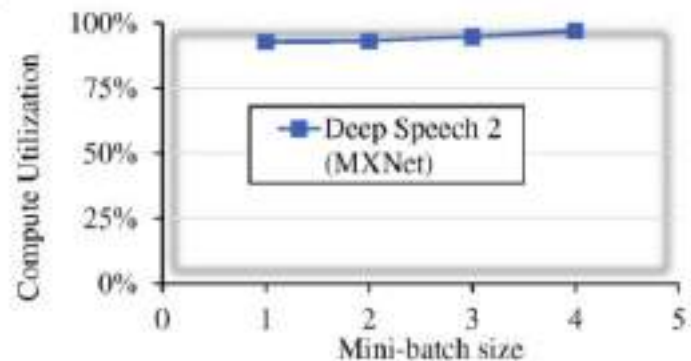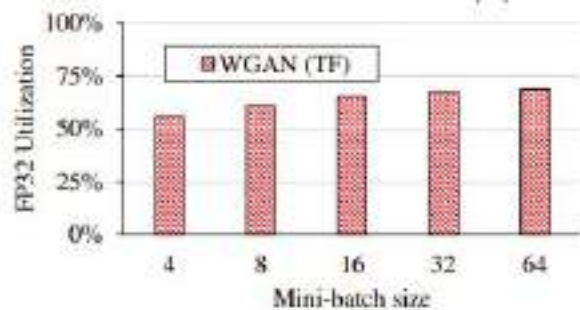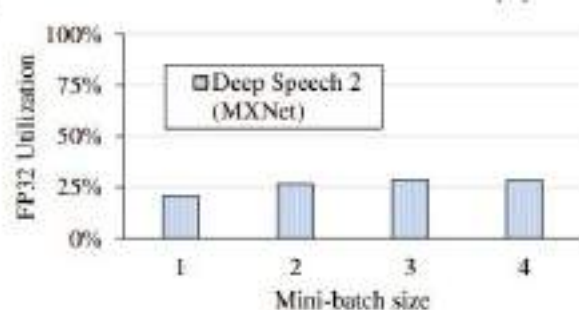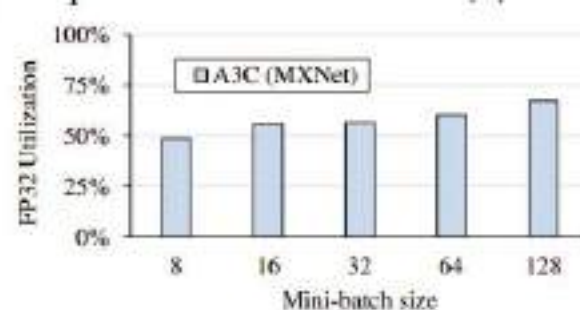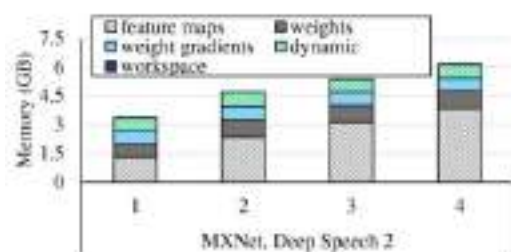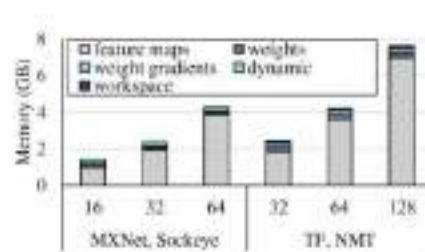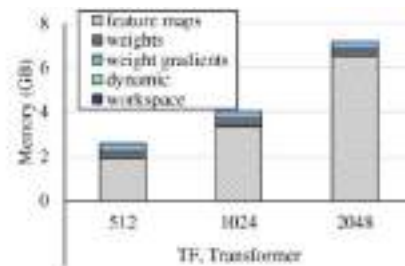(e) WGAN　　　　　(f) Deep Speech 2　　　　　(g) A3C

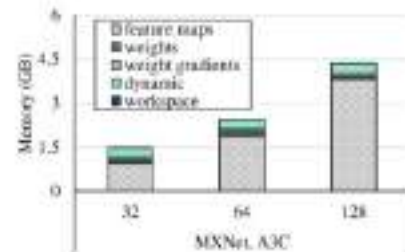Fig. 6: GPU FP32 utilization for different models on multiple mini-batch sizes.



(d) Deep Speech 2　　　(e) Seq2Seq　　　(f) Transformer　　　(g) A3C
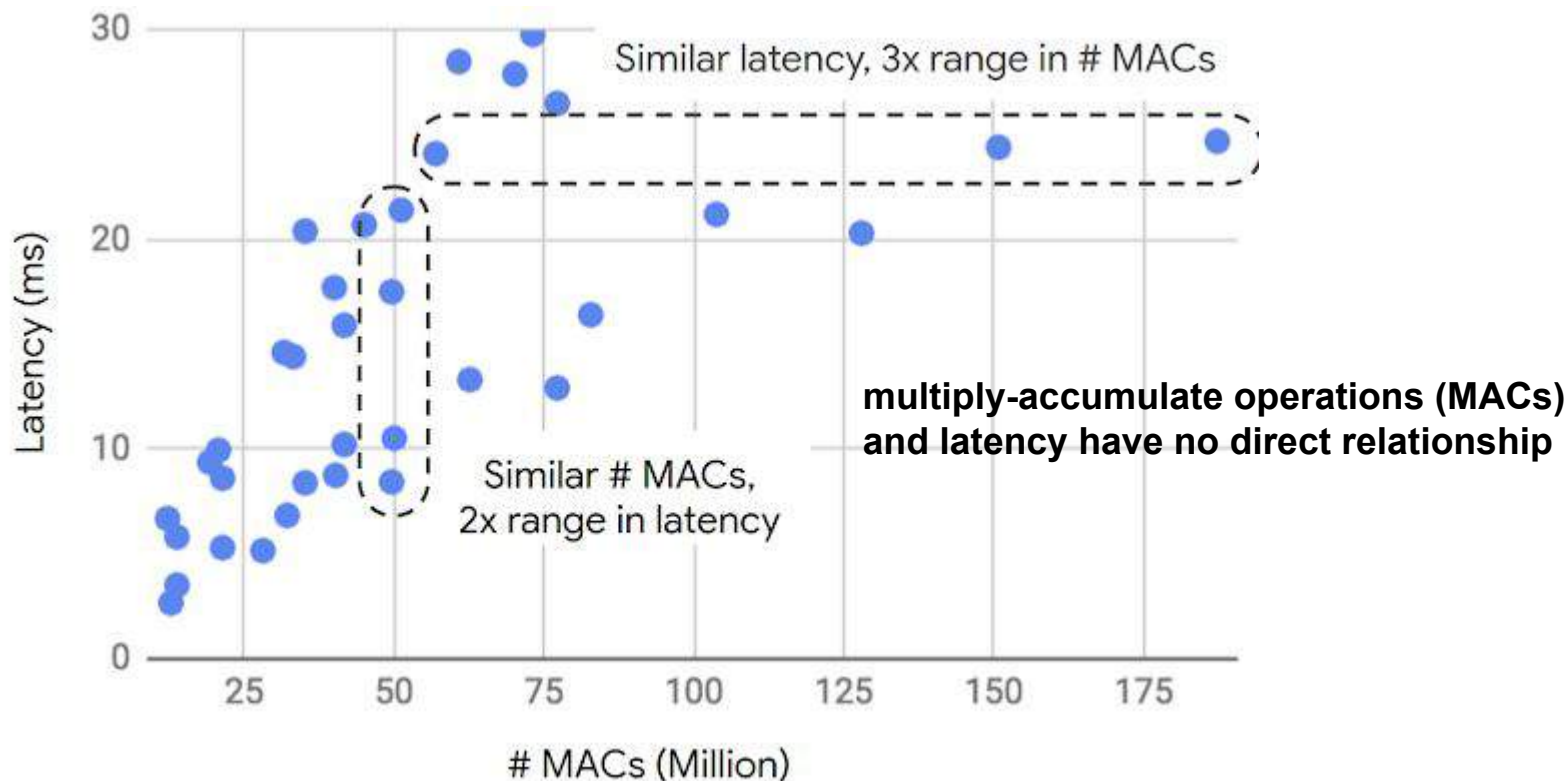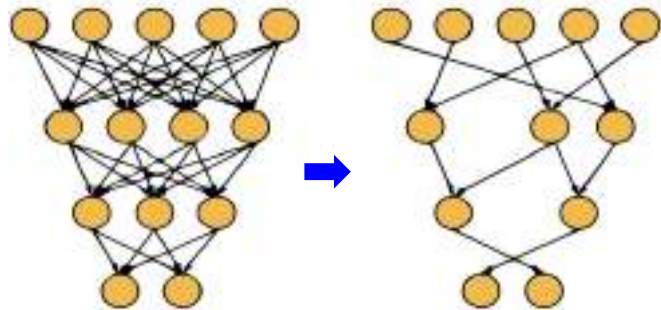
Fig. 8: GPU memory usage breakdown for different models on multiple mini-batch sizes.

Similar latency, 3x range in # MACs

multiply-accumulate operations (MACs)
and latency have no direct relationship

Similar # MACs,
2x range in latency

Introducing the CVPR 2018 On-Device Visual Intelligence Challenge
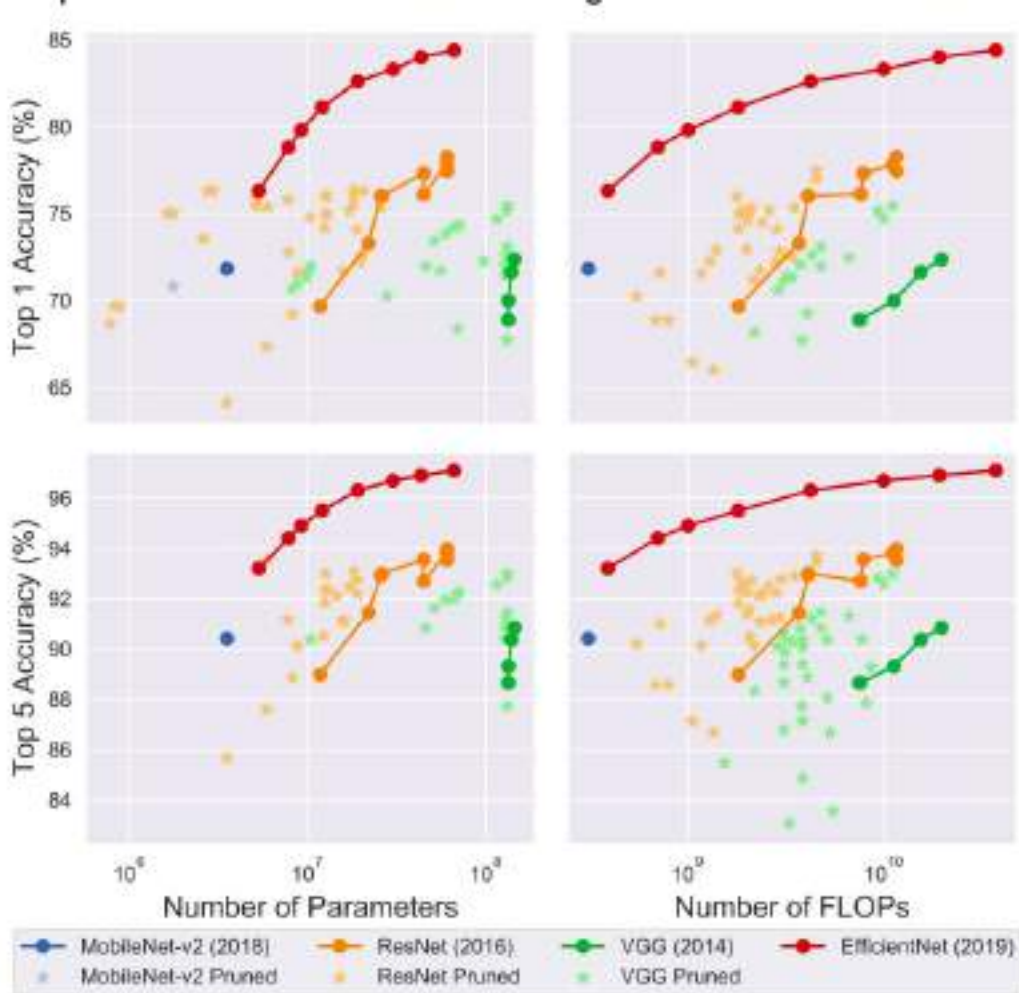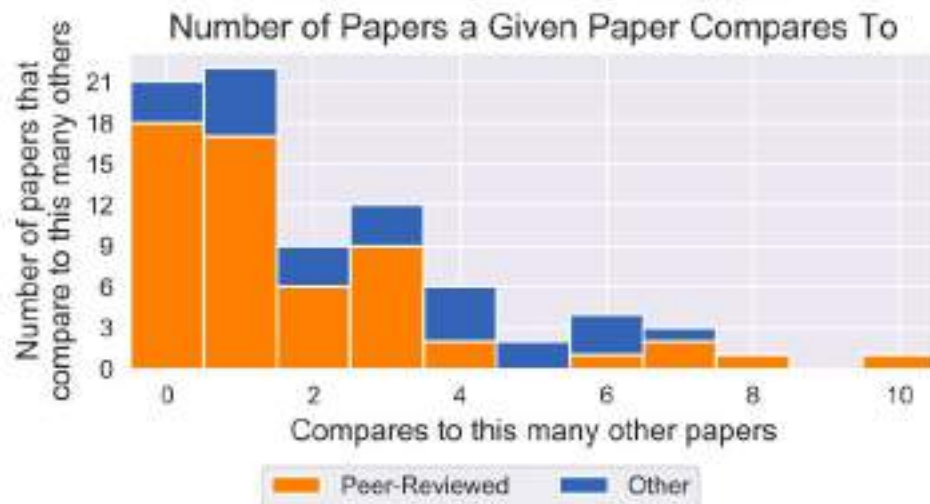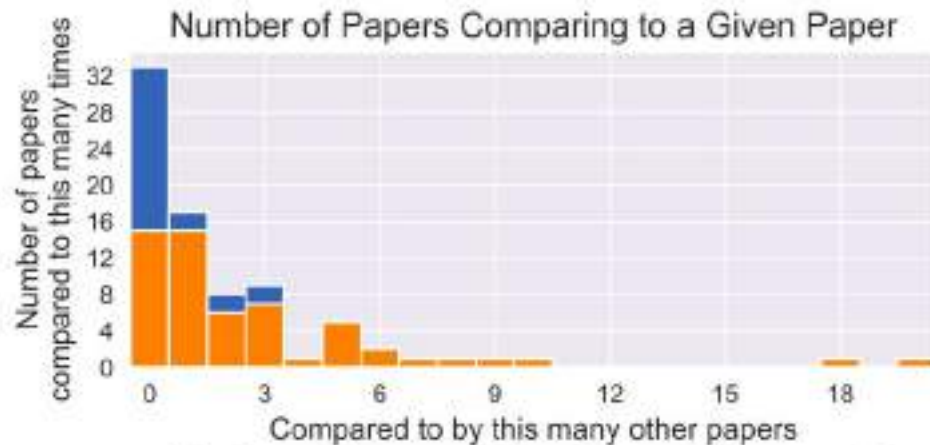Friday, April 20, 2018, Google AI Blog

# Network Pruning

$$f\left(\sum_{k=1}^{n} w_k \times i_k + b\right)$$

**What is the state of neural network pruning?**
Davis Blalock, Jose Javier Gonzalez Ortiz, Jonathan Frankle, John Guttag

Number of Papers Comparing to a Given Paper

Number of Papers a Given Paper Compares To

Peer-Reviewed    Other

# *Reproducibility Challenge*

Reproducing results in research papers can be hard:
- lack of source code
- comparable hardware
- software environment, library and right versions
- undocumented parameters

Yung-Hsiang Lu, Purdue University

# *Summary*

- Computer vision can be evaluated in many different ways, including performance.
- Performance can be defined in different ways, such as execution time.
- Many factors affect performance, such as the sizes of the networks, but the relationships are not straight lines.
- Training time is affected by the sizes of mini batches.

Yung-Hsiang Lu, Purdue University