

Efficient Learning

Danna Gurari

University of Colorado Boulder

Fall 2022



Review

- Last week:
 - Motivation
 - Key idea: knowledge distillation
 - Knowledge distillation for CNNs (vision problems)
 - Knowledge distillation for Transformers (language problems)
- Assignments (Canvas):
 - Final project presentations due Monday
 - Note: provide video URLs with YouTube or Vimeo
- Questions?

Today's Topics

- Motivation
- Efficient learning: curriculum learning
- Efficient learning: active learning
- Efficient learning: other considerations
- Faculty course questionnaire

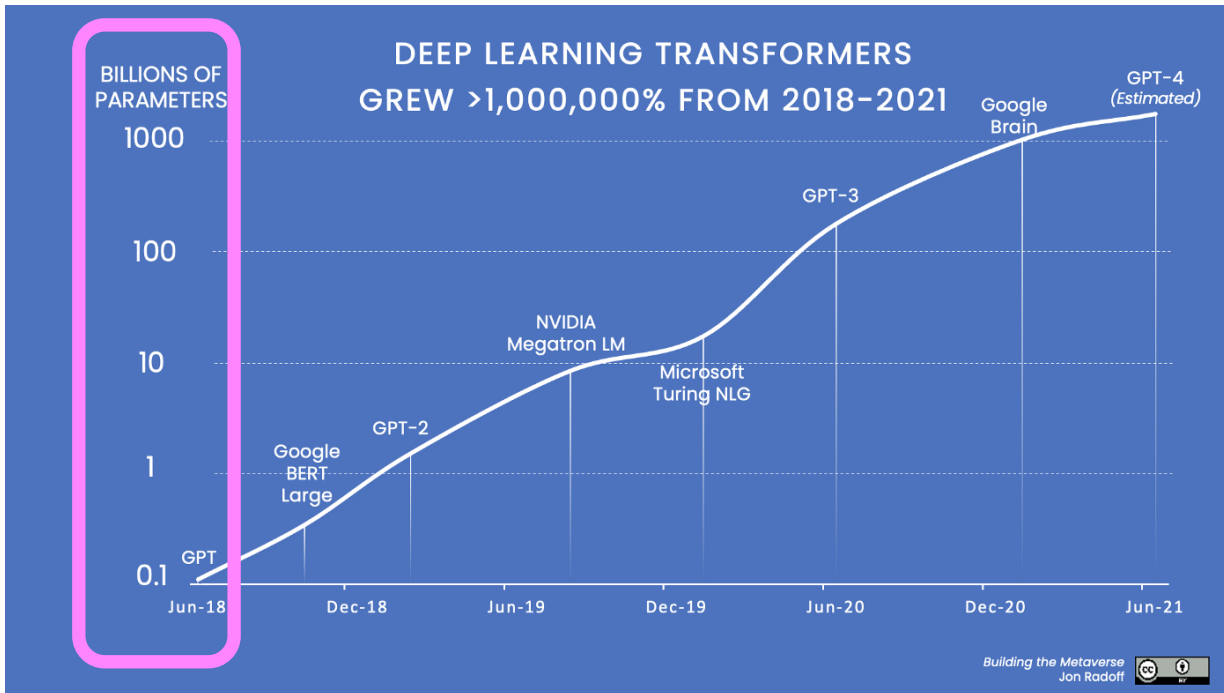
Today's Topics

- Motivation
- Efficient learning: curriculum learning
- Efficient learning: active learning
- Efficient learning: other considerations
- Faculty course questionnaire

Trend: Parameter-Heavy Models

Language – pretrained transformers

Vision – ImageNet classification



Architecture	Year	Top-1 Accuracy	# Parameters
DenseNet-169	2017	76.2%	14M
Inception-v3	2016	78.8%	24M
Inception-resnet-v2	2017	80.1%	56M
PolyNet	2017	81.3%	92M
SENet	2018	82.7%	146M
GPipe	2018	84.3%	557M
ResNeXt-101 32x48d	2019	85.4%	829M

<https://medium.com/building-the-metaverse/the-metaverse-and-artificial-intelligence-ai-577343895411>

Trend: Parameter-Heavy Models Are Often Predicated on Extensive Training

(Measured on Nvidia A100)

Models	#Params (M)	Training Time (GPU Hours)
ResNet-50	26	31
ResNet-101	45	44
BERT-Base	108	84
Turing-NLG 17B	17,000	TBA
GPT-3 175B	175,000	3,100,000

On a single GPU, it would take 335 years to train GPT-3

Why Is Extensive Training Is Costly?

- Time-consuming
- Expensive
- Increased environmental impact from extra computations

Extensive Training Is Costly; e.g., Training BERT Cost:

~\$80k-\$1.6m:

THE COST OF TRAINING NLP MODELS A CONCISE OVERVIEW

Or Sharir
AI21 Labs
ors@ai21.com

Barak Peleg
AI21 Labs
barakp@ai21.com

Yoav Shoham
AI21 Labs
yoavs@ai21.com

On the Dangers of Stochastic Parrots: Can Language Models Be Too Big?

Emily M. Bender*
ebender@uw.edu
University of Washington
Seattle, WA, USA

Timnit Gebru*
timnit@blackinai.org
Black in AI
Palo Alto, CA, USA

Angelina McMillan-Major
aymm@uw.edu
University of Washington
Seattle, WA, USA

Shmargaret Shmitchell
shmargaret.shmitchell@gmail.com
The Aether

as much energy as a
trans-American flight:



Boss: What did you do last month?

You: Trained the model for one epoch.



Boss: Umm, fine, what is your plan for next month?

You: Train... train the model for one more epoch?



Today's Topics

- Motivation
- **Efficient learning: curriculum learning**
- Efficient learning: active learning
- Efficient learning: other considerations
- Faculty course questionnaire

How to teach machines to learn more efficiently?



Intuition: How to Teach a Child Math?

Random Order of Examples

Meaningful Order of Examples

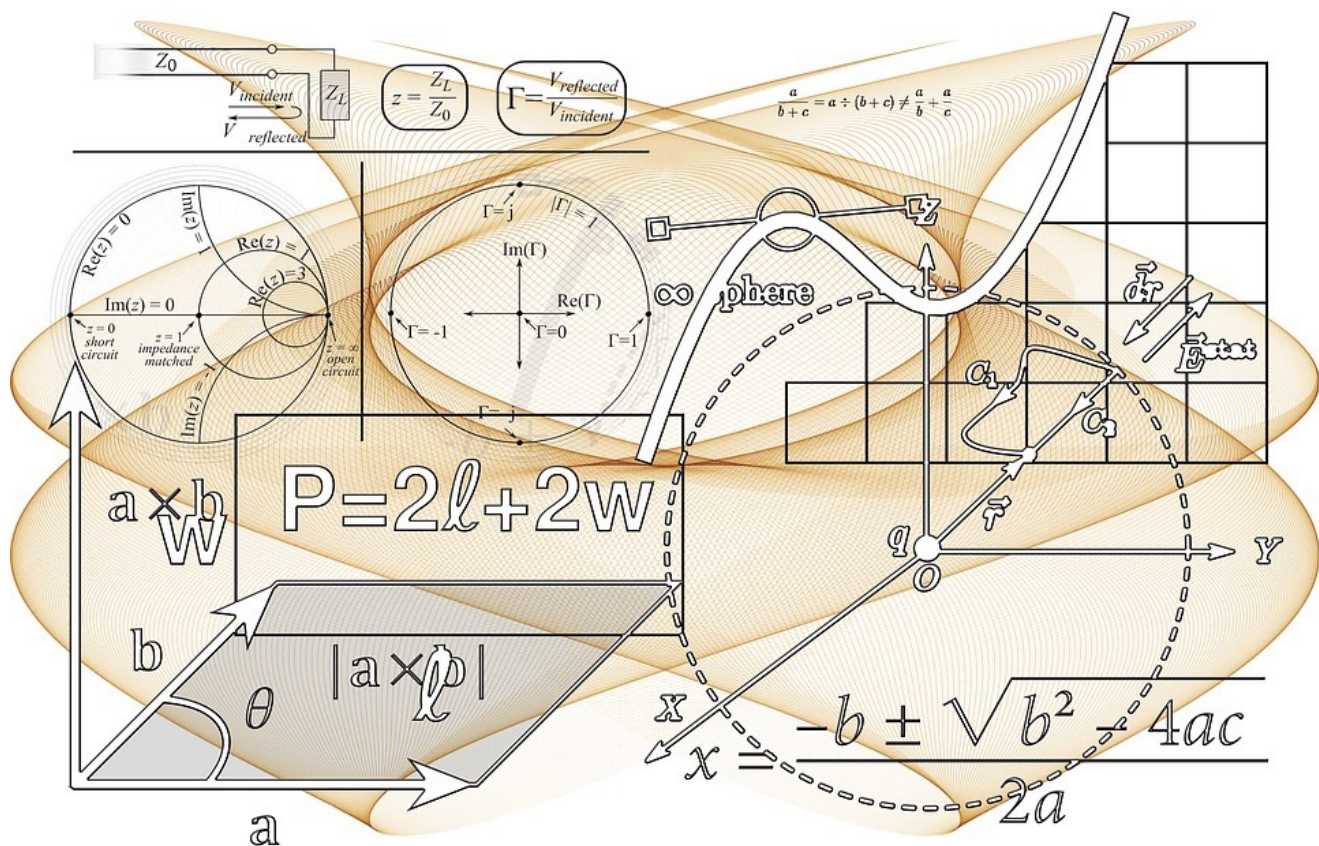


Table of Contents	
Letter from Dinah Zike	ix
Introduction to Foldables	1
Why Use Foldables in Mathematics?	1
Foldable Basics	1
Choosing the Appropriate Foldable	2
Folding Instructions	3
Using Foldables	3
1-Page Foldables	4
Book Book	4
Factor Book	7
Area Book	8
Two-Digit Book	9
2-Page Foldables	10
Factor Book	11
Area Book	12
3-Page Foldables	13
Table Book	13
Area Book	14
Area Book	15
Area Book	16
4-Page Foldables	17
Area Book	17
Area Book	18
Area Book	19
Area Book	20
Area Book	21
Area Book	22
Area Book	23
Area Book	24
Any Number of Pages	25
Area Book	25
Folding the Foldables	26
Folding the Foldables	27
Folding the Foldables	28
Area Book	29
Area Book	30
Area Book	31
Area Book	32
Area Book	33
Area Book	34
Area Book	35
Area Book	36
Area Book	37
Area Book	38
Area Book	39
Area Book	40
Area Book	41
Area Book	42
Area Book	43
Area Book	44
Area Book	45
Area Book	46
Area Book	47
Area Book	48
Area Book	49
Area Book	50
Area Book	51
Area Book	52
Area Book	53
Area Book	54
Area Book	55
Area Book	56
Area Book	57
Area Book	58
Area Book	59
Area Book	60
Area Book	61
Area Book	62
Area Book	63
Area Book	64
Area Book	65
Area Book	66
Area Book	67
Area Book	68
Area Book	69
Area Book	70
Area Book	71
Area Book	72
Area Book	73
Area Book	74
Area Book	75
Area Book	76
Area Book	77
Area Book	78
Area Book	79
Area Book	80
Area Book	81
Area Book	82
Area Book	83
Area Book	84
Area Book	85
Area Book	86
Area Book	87
Area Book	88
Area Book	89
Area Book	90
Area Book	91
Area Book	92
Area Book	93
Area Book	94
Area Book	95
Area Book	96
Area Book	97
Area Book	98
Area Book	99
Area Book	100

Intuition: How to Teach a Child To Read



Random Order of Examples



Meaningful Order of Examples



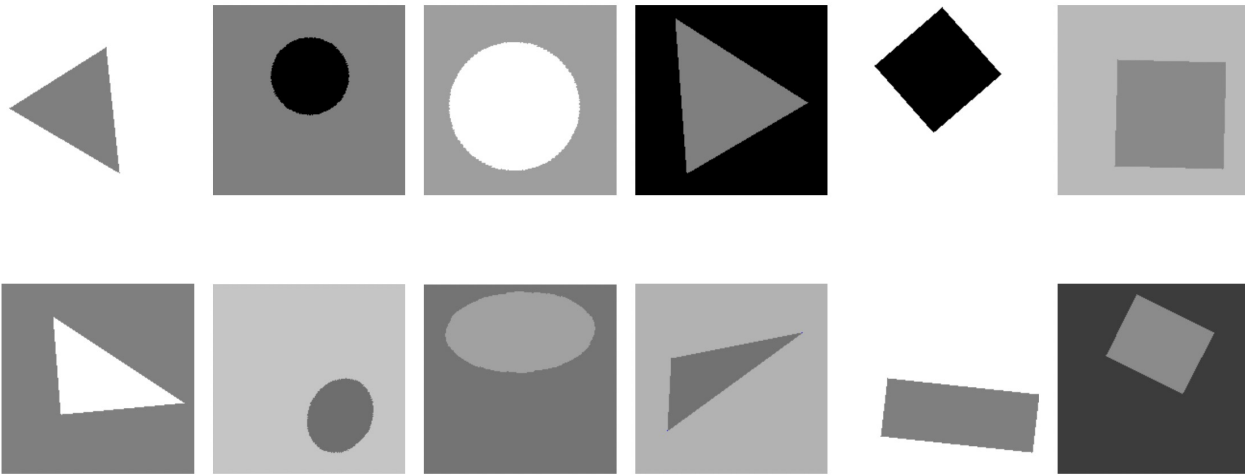
Idea: Teach Machines As We Teach Humans

Curriculum

Train with simpler examples first and progressively harder examples over time

Tasks

1. Classify each shape as rectangle, ellipse, or triangle

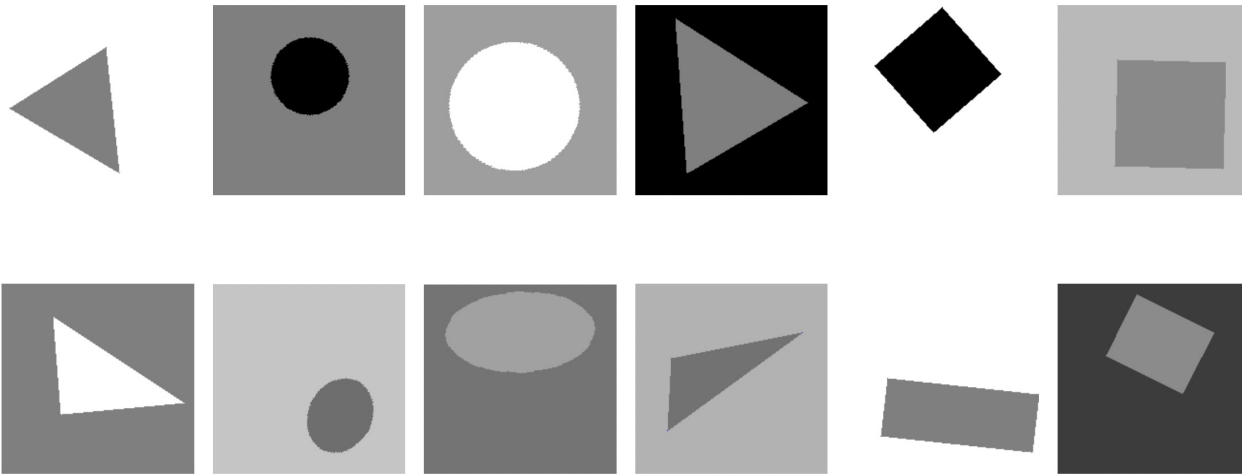


2. Predict the next word

Background music from a

Shape Prediction: Curriculum Learning

1. Classify each shape as rectangle, ellipse, or triangle



Architecture: 3-layer neural network

Easy (Basic): less shape variability
(squares, circles, and equilateral triangles);
10,000 examples

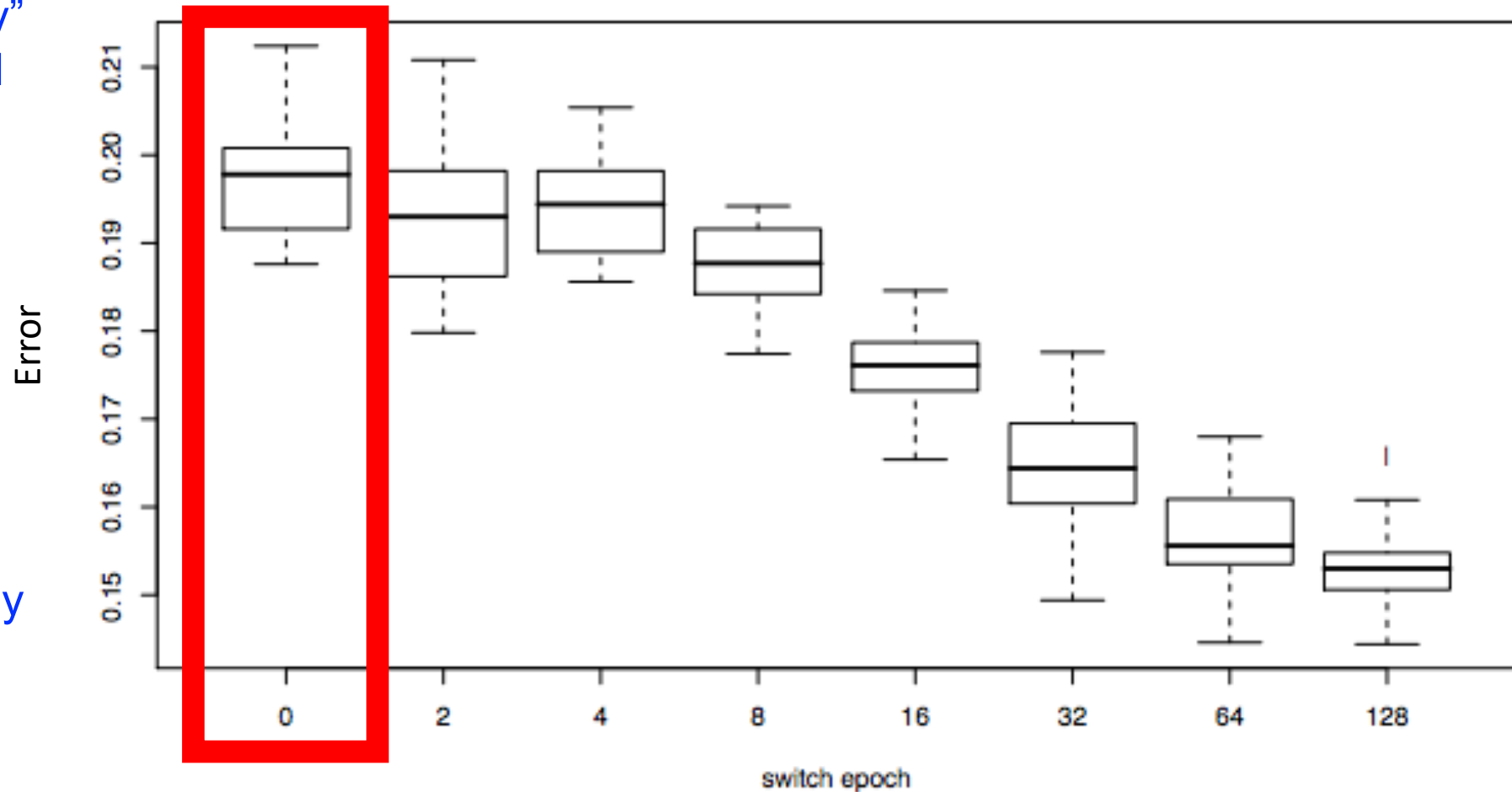
Hard (Geom): more shape variability
(rectangles, ellipses, and triangles);
10,000 examples

Shape Prediction: Curriculum Learning

Results of training on “easy” examples for n epochs and then training on “hard” examples until 256 epochs (20 random initializations).

What are benefits of curriculum learning?

How many epochs should the algorithm train with easy examples before switching to difficult examples?



No curriculum

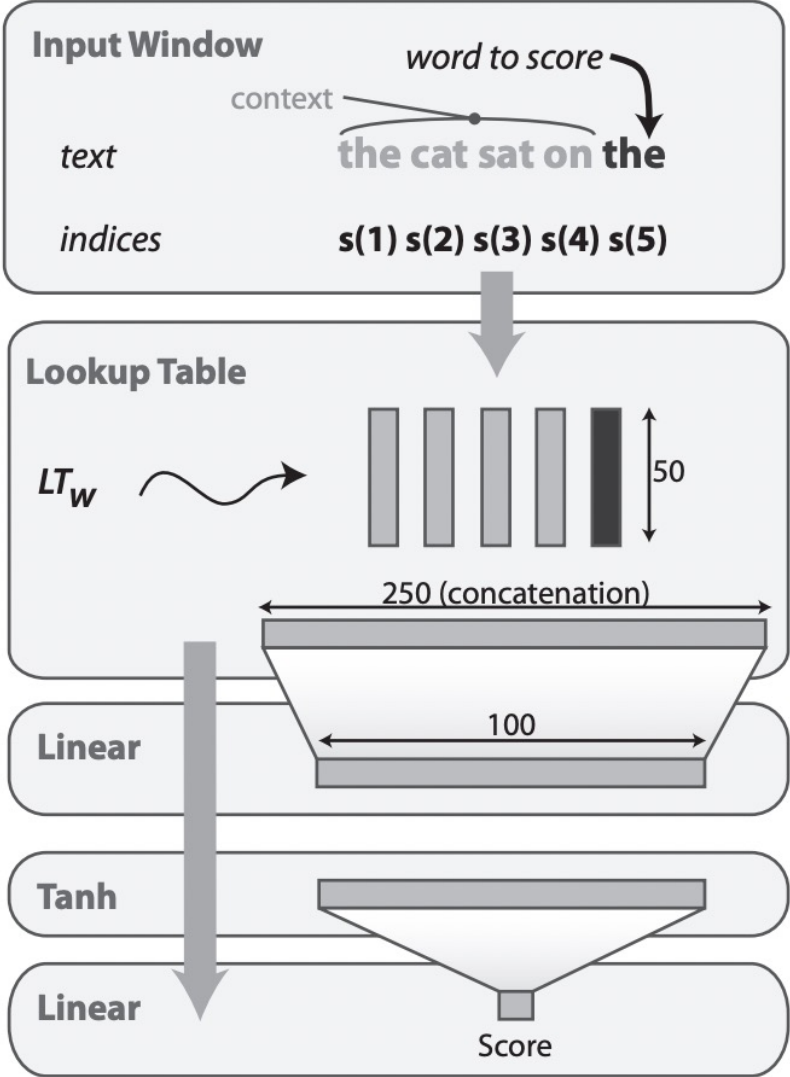
Next Word Prediction: Curriculum Learning

Architecture:
context size
set to 5

Easy: 5,000 most
frequent words

Hard: additional 5,000
words at each epoch
until 20,000 words

**Examples with words
not in the vocab were
discarded from training**

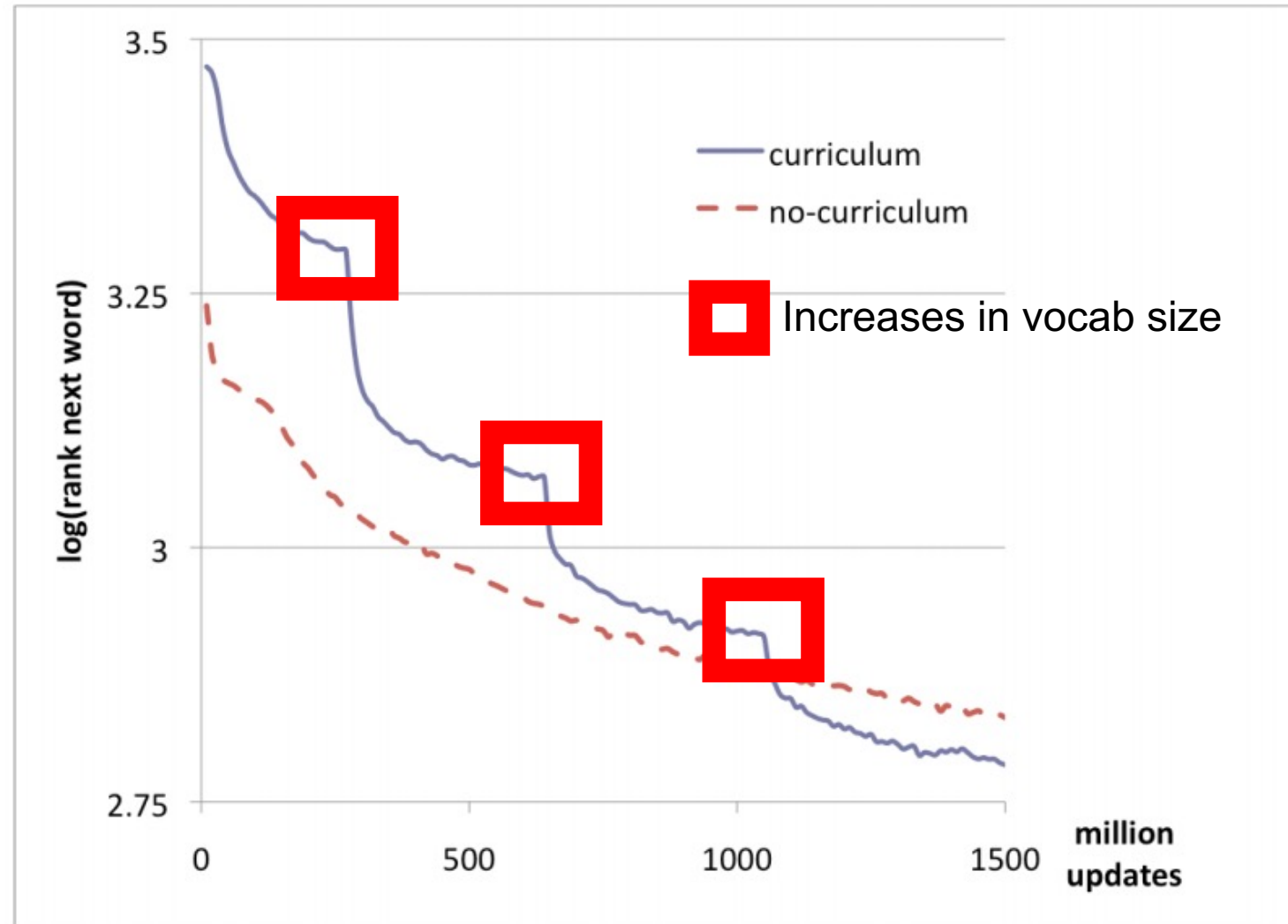


2. Predict the next word

Background music from a _____

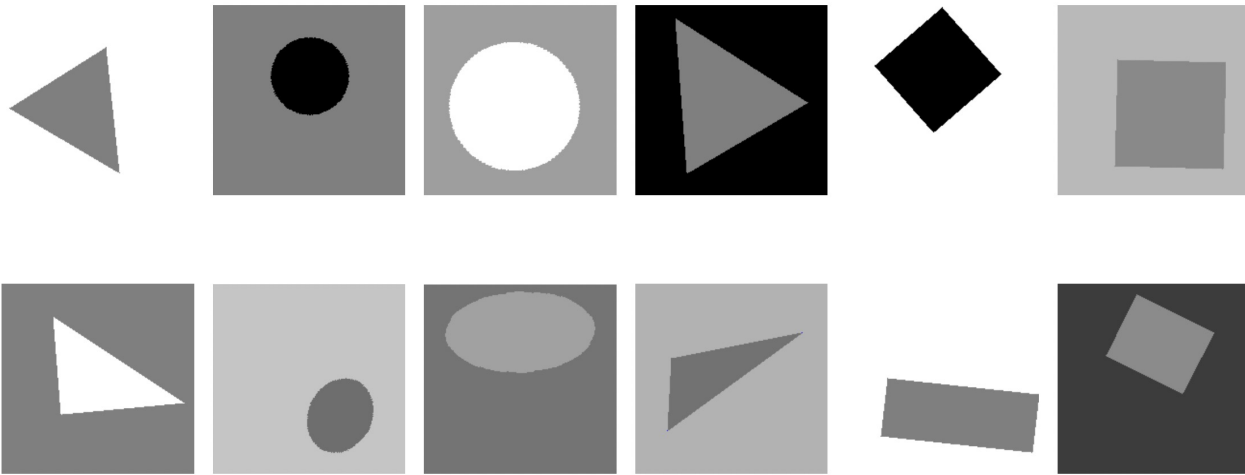
Next Word Prediction: Curriculum Learning

What are benefits of curriculum learning?



Summary: Curriculum Learning is a Form of Transfer Learning that Accelerates Optimization

1. Classify each shape as rectangle, ellipse, or triangle



2. Predict the next word

Background music from a _____

Key Questions for Curriculum Learning; e.g., for Visual Question Answering



Is my monitor on?



Hi there can you please
tell me what flavor this is?



Does this picture look
scary?



Which side of the
room is the toilet on?

Questions

1. What criteria should be used to order examples?
2. How would you update the training data (and how often)?

Today's Topics

- Motivation
- Efficient learning: curriculum learning
- **Efficient learning: active learning**
- Efficient learning: other considerations
- Faculty course questionnaire

How to teach machines with minimal human supervision?



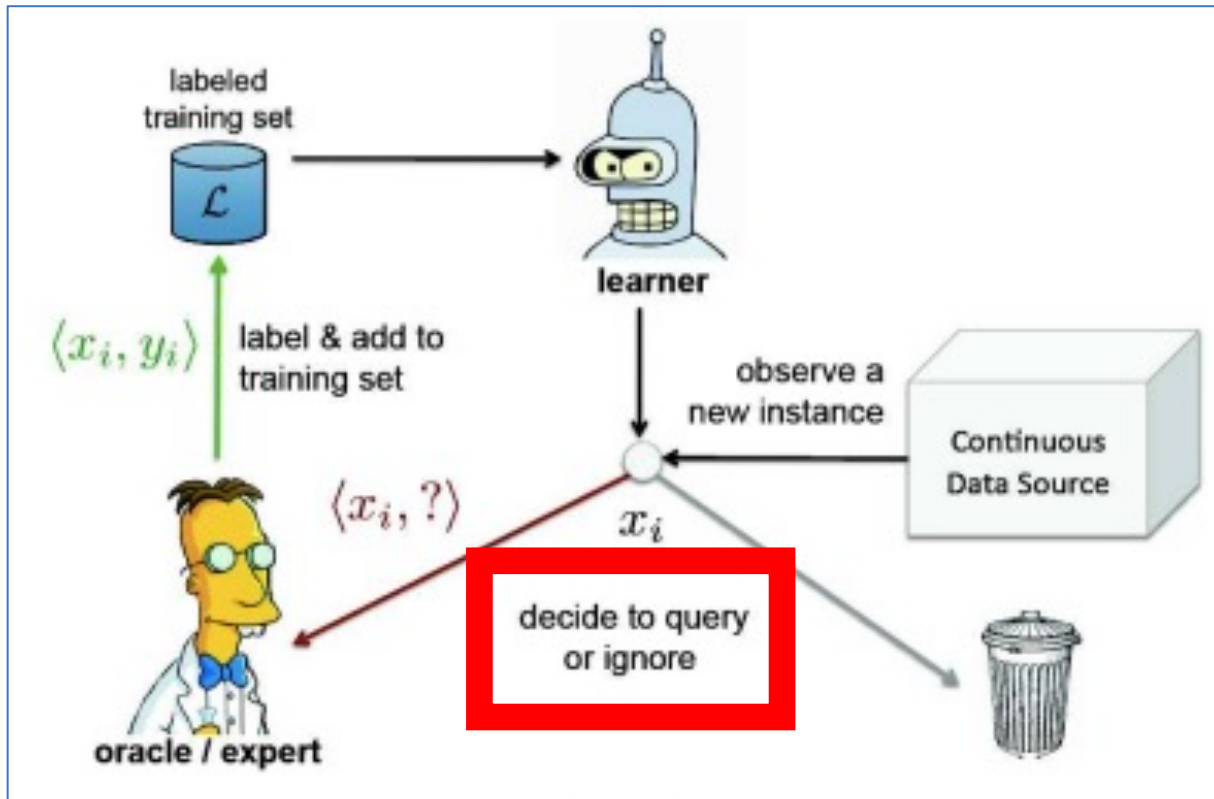
e.g., limited access to
(expert) annotators



e.g., limited funding

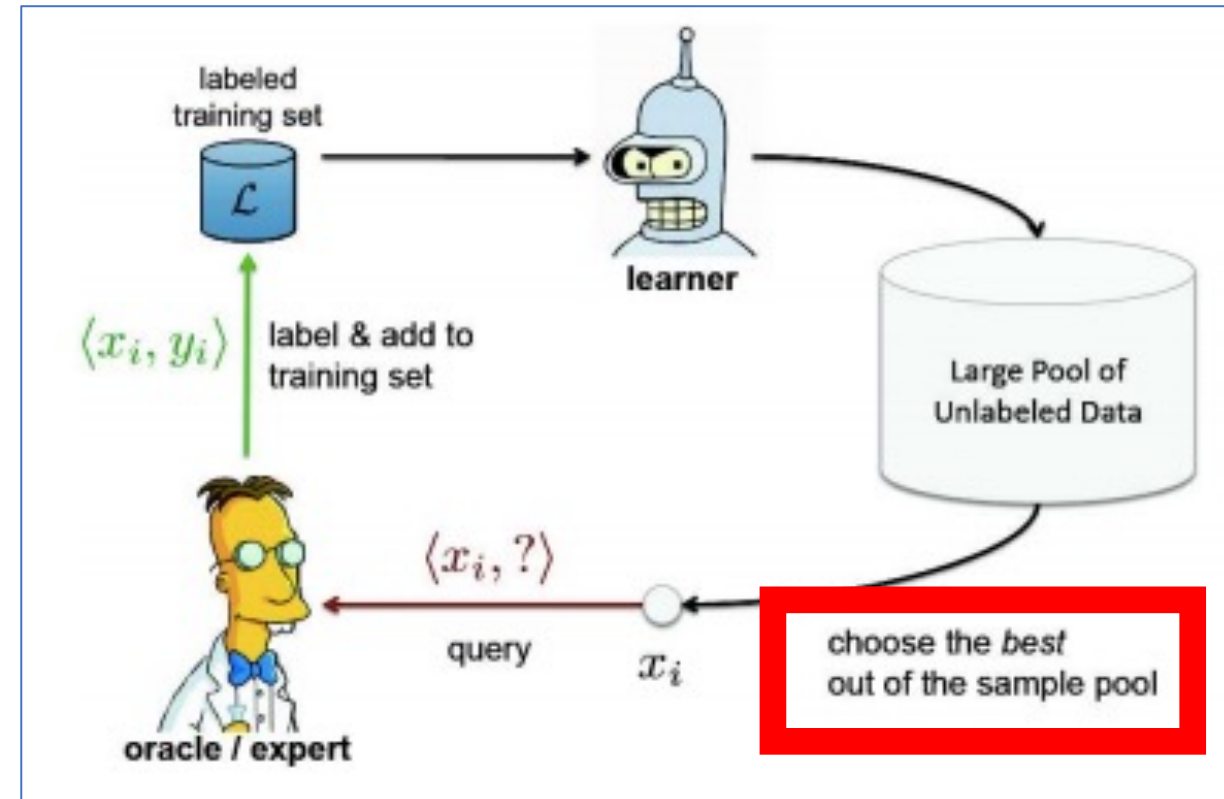
Idea: Choose Most Informative Data to Label

Stream-Based



Consider one example at a time

Pool-Based

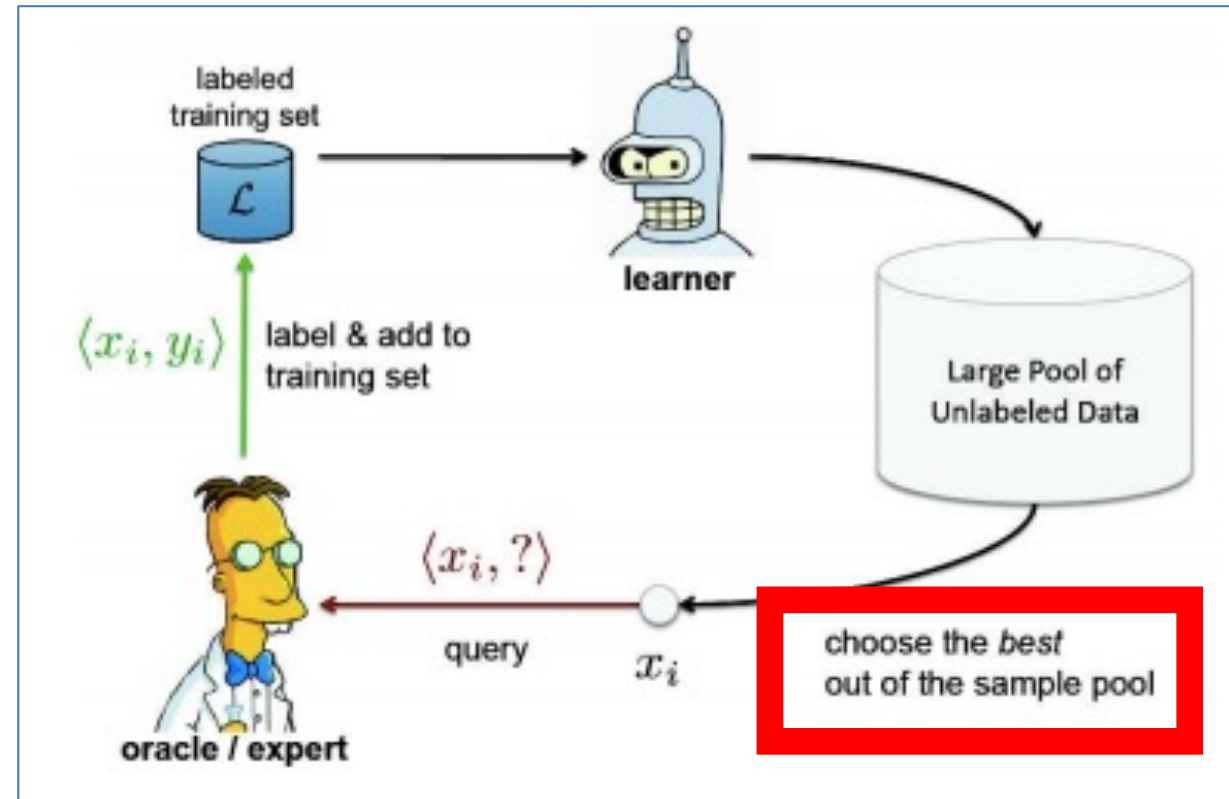


Consider many examples at a time

Active Learning for Neural Networks: Status Quo

Iteratively add more labelled training examples after n epochs; different from curriculum learning because labels need to be collected for the added data

Pool-Based



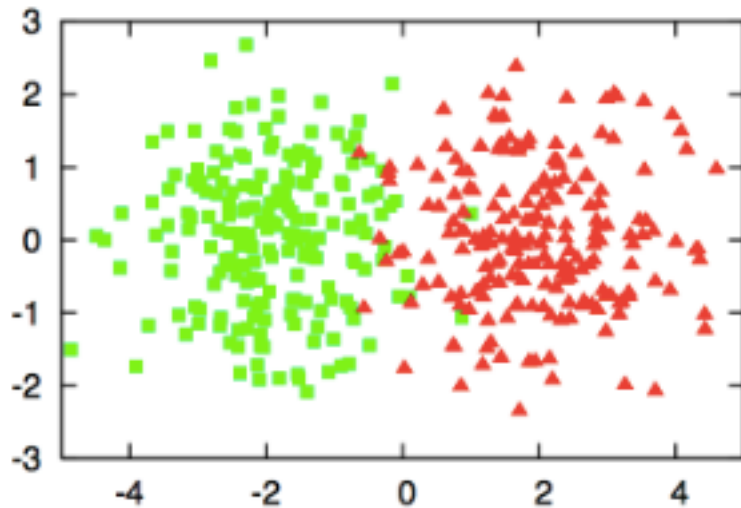
Consider many examples at a time

What approach might be effective in identifying the most informative data to label for training neural networks?

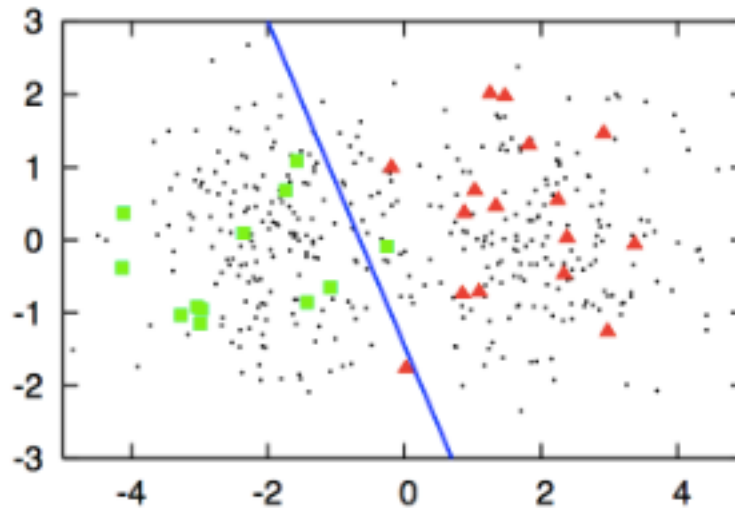
Common Approach: Uncertainty Sampling

Query instance(s) the classifier is most uncertain about.

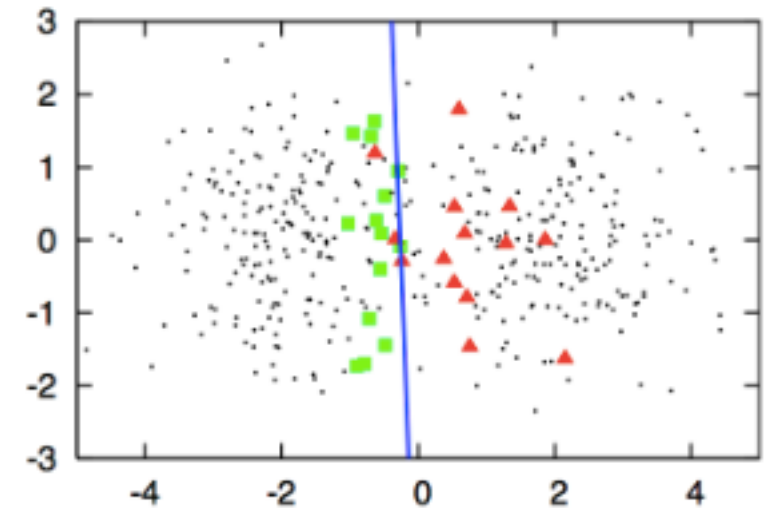
True Representation
(Assume Labels Are
Not Known)



Passive Learner
(Random Selection)



Active Learner
(Uncertainty Sampling)



e.g., Uncertainty Estimation for Neural Networks **Using Robustness Testing**

Use model's predictions on random augmentations of the input to measure consistency/uncertainty; e.g.,

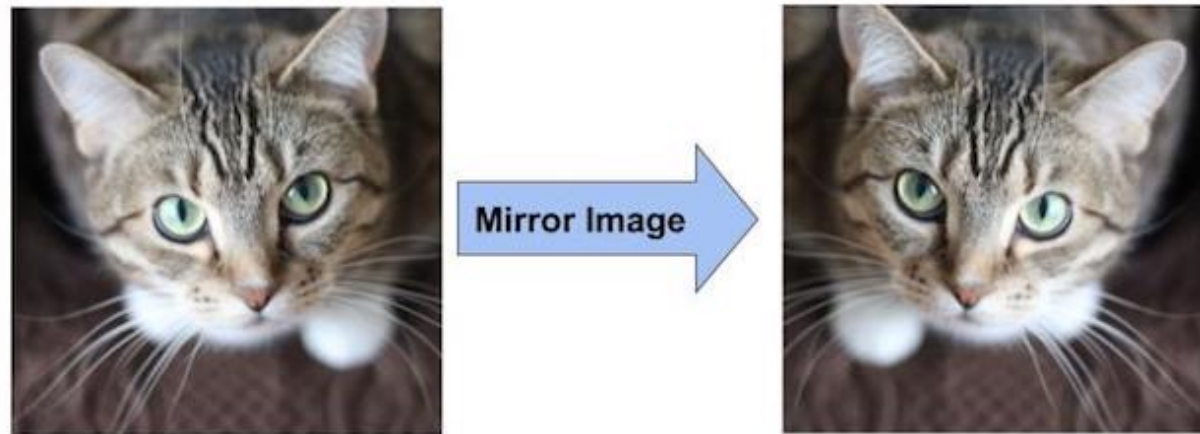
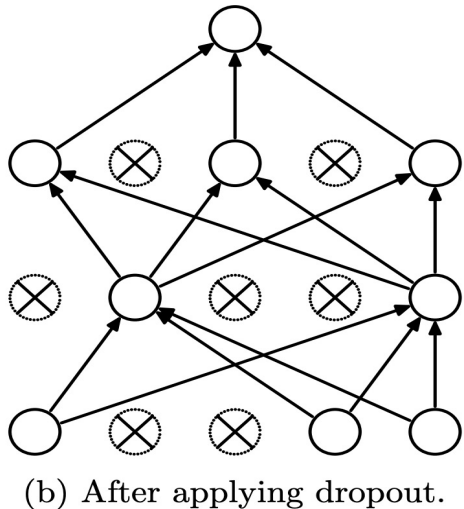


Figure Source: <https://learnopencv.com/understanding-alexnet/>

e.g., Uncertainty Estimation for Neural Networks Using Ensembles (Two Approaches)

1. Dropout with different masks at inference time



2. Multiple neural networks

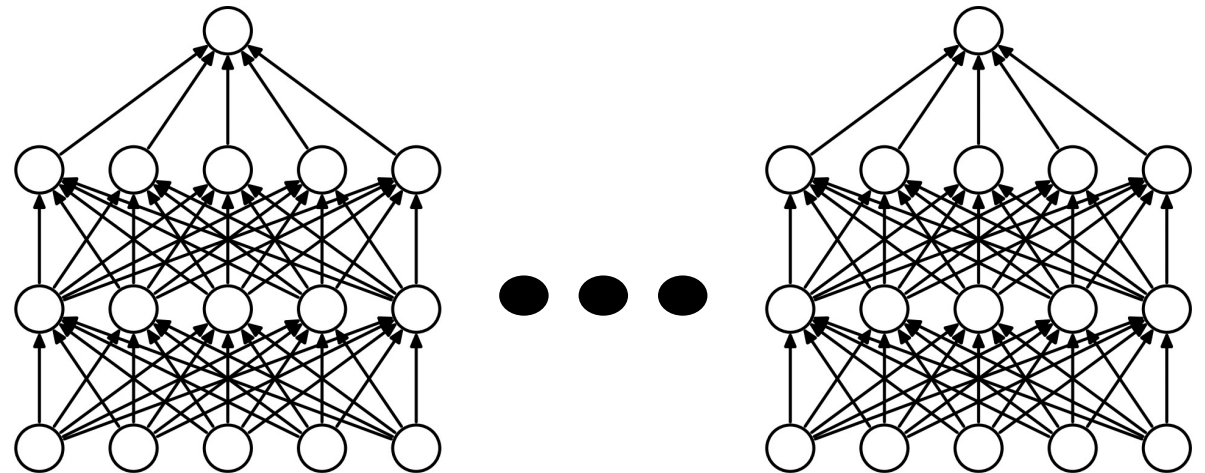


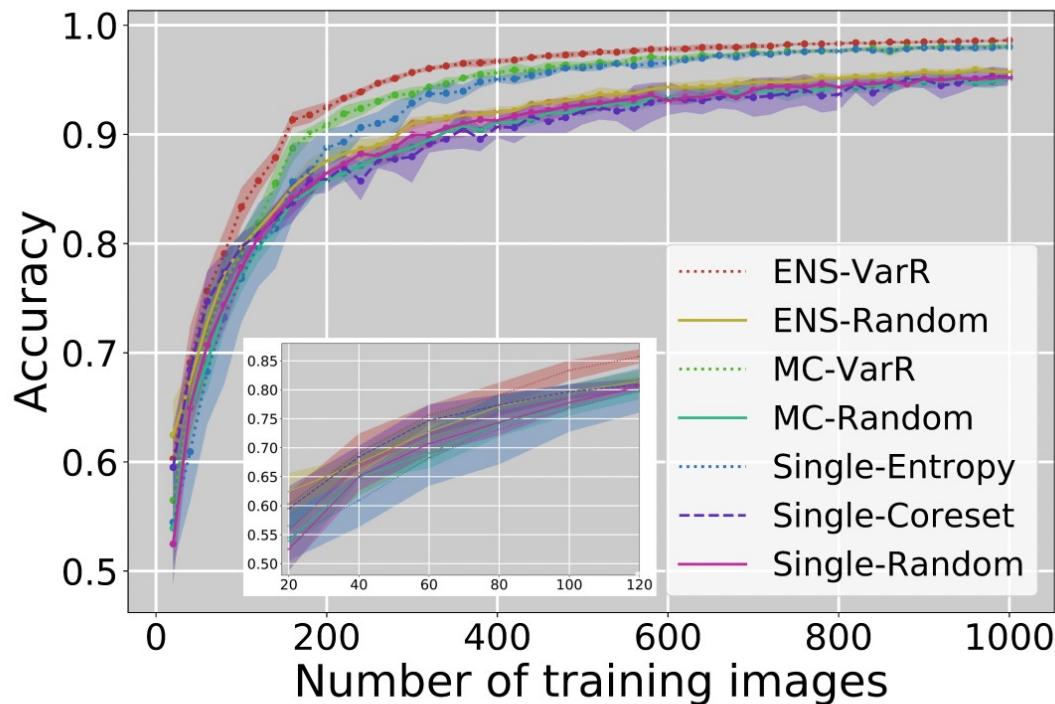
Figure Source: Srivastava et al. Dropout: A Simple Way to Prevent Neural Networks from Overfitting. Journal of Machine Learning Research. 2014

Predicted softmax probabilities used to estimate uncertainty (e.g., entropy across softmax values), with average taken across all ensemble's softmax distributions

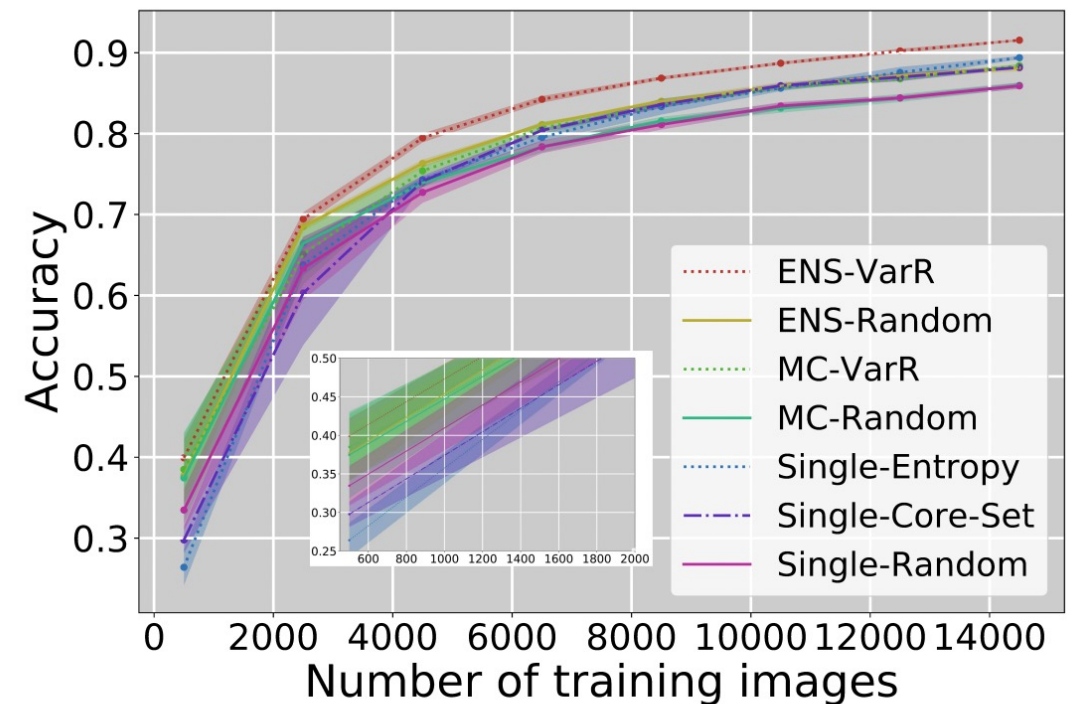
Beluch et al. The power of ensembles for active learning in image classification. CVPR 2018

e.g., Uncertainty Estimation for Neural Networks Using Ensembles (Two Approaches)

Active learning methods lead to **faster learning** and **reduced human annotation effort** than passive (random) learning for two image classification datasets



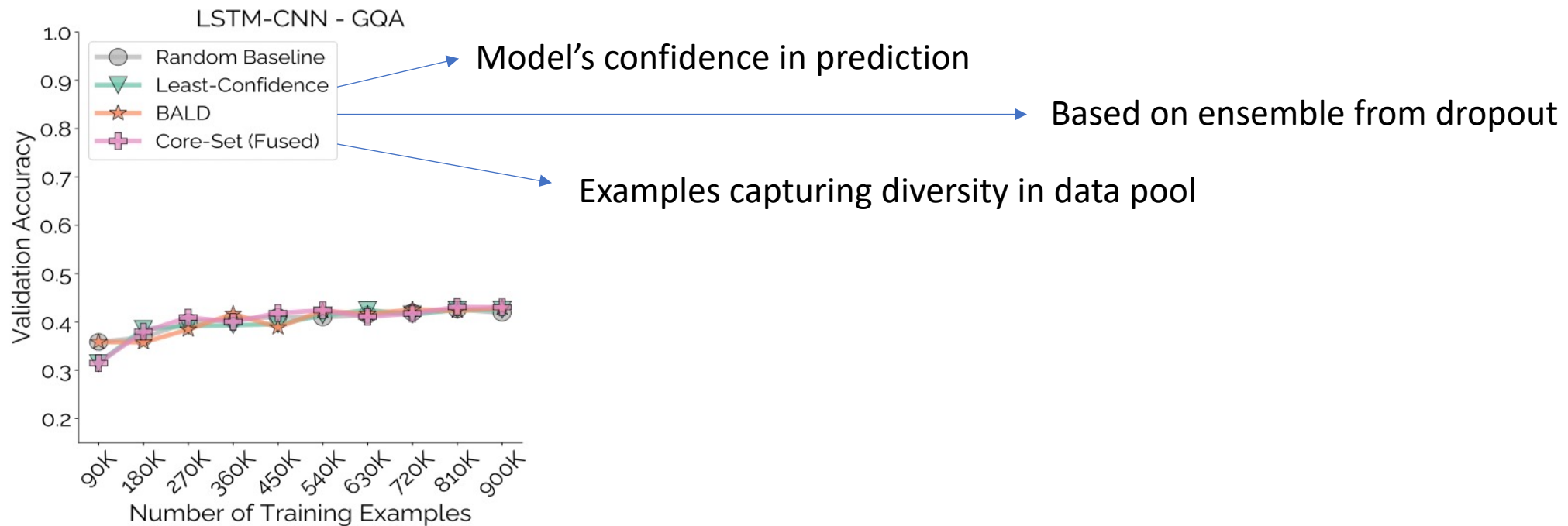
(a) MNIST on S-CNN



(b) CIFAR-10 on DenseNet

Common AL Techniques Have Mixed Results

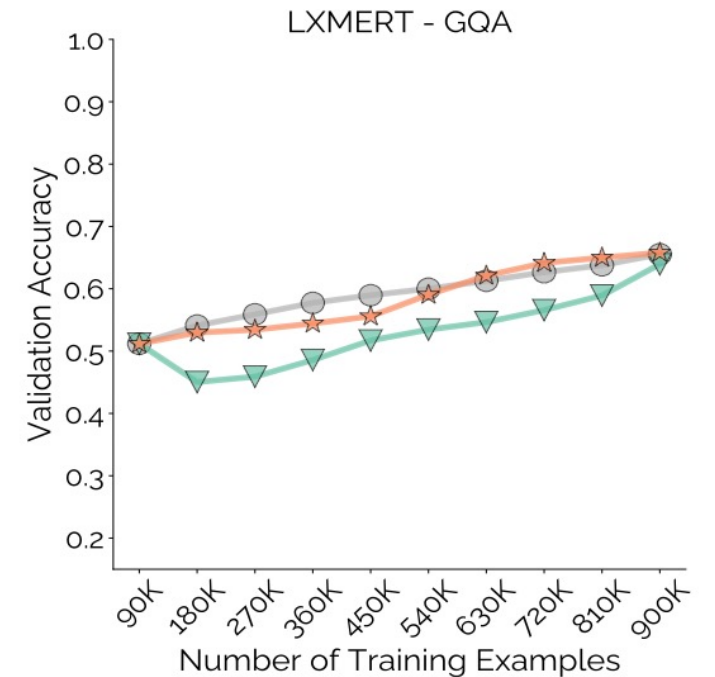
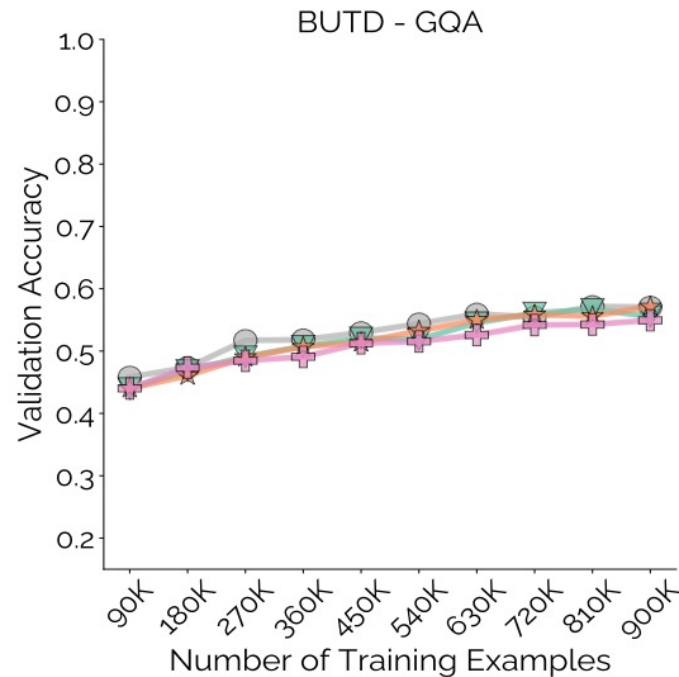
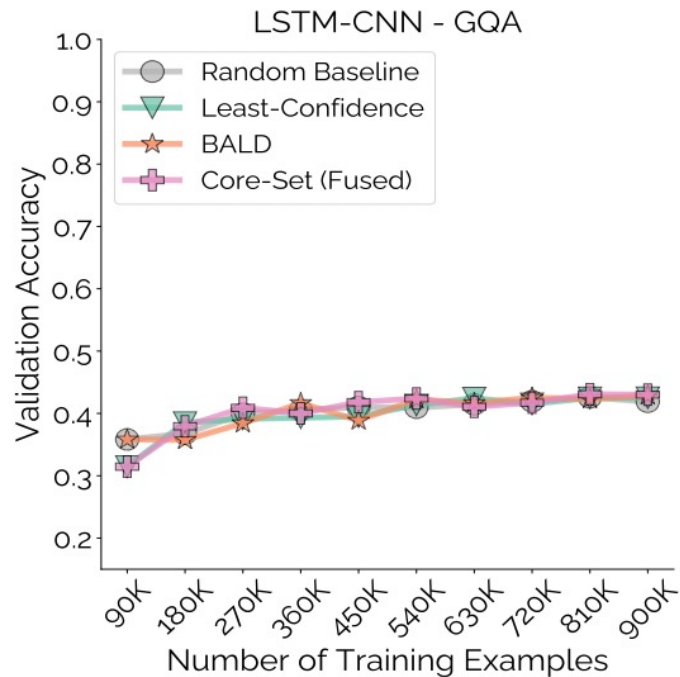
- **Successes:** image classification, object detection
- **Failure: VQA** (e.g., AL methods label 10% of overall pool per iteration; initial model trained on 10% of pool)



Karamcheti et al. Mind your outliers! Investigating the negative impact of outliers on active learning for visual question answering. Association for Computational Linguistics (ACL) 2021

Common AL Techniques Have Mixed Results

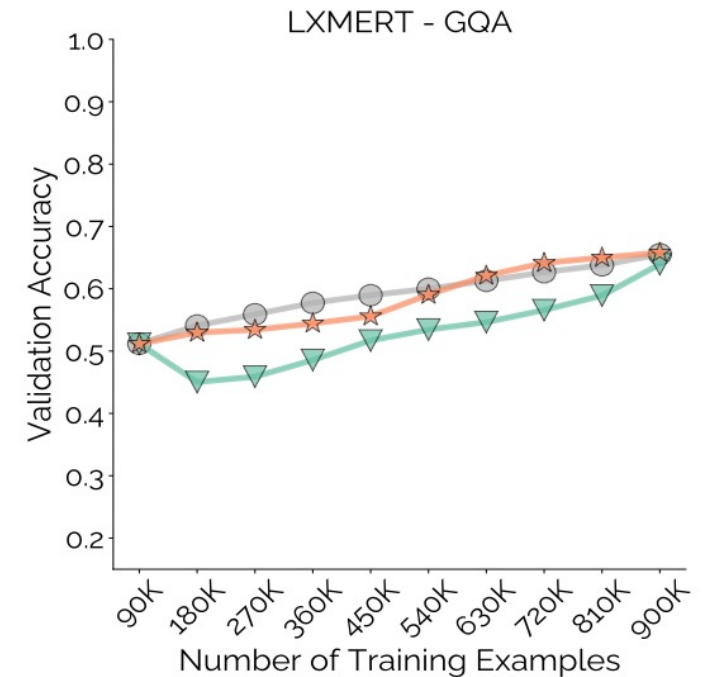
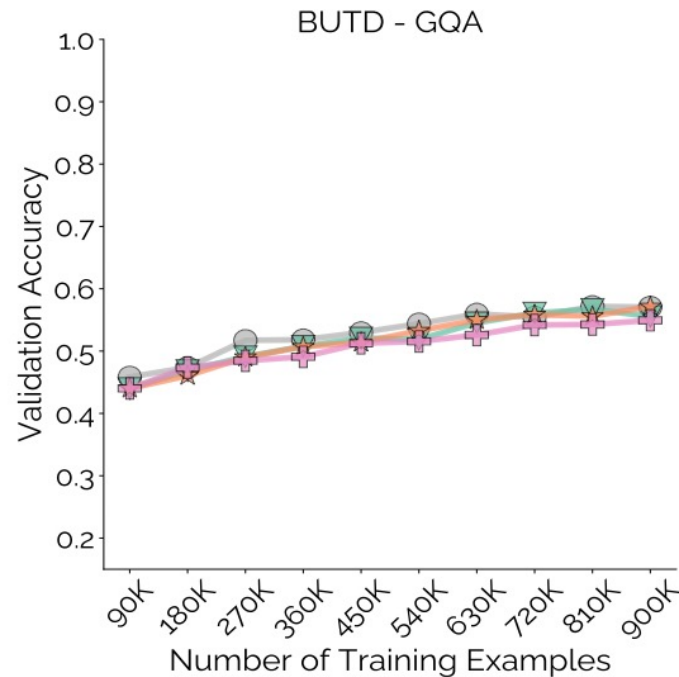
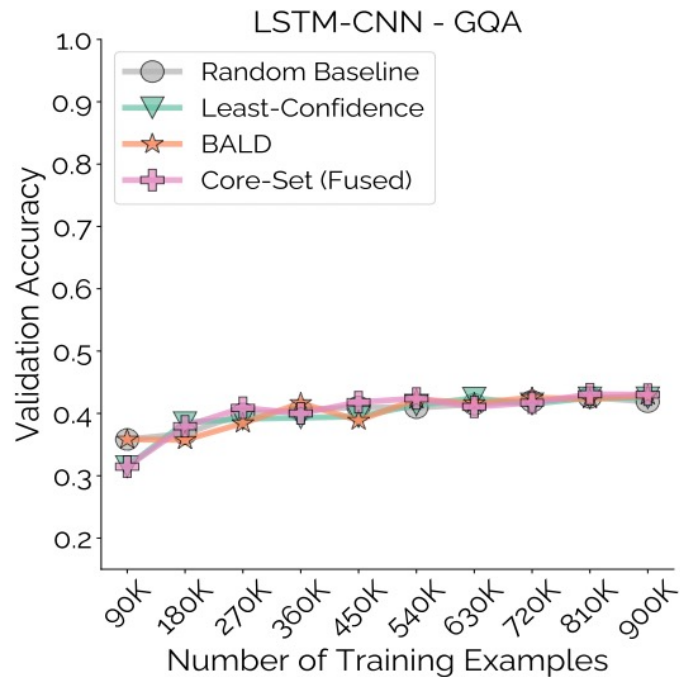
How do the 3 AL methods compare to random selection for the 3 VQA models?



Karamcheti et al. Mind your outliers! Investigating the negative impact of outliers on active learning for visual question answering. Association for Computational Linguistics (ACL) 2021

Common AL Techniques Have Mixed Results

Why might AL methods perform comparable or worse to random selection?



Karamcheti et al. Mind your outliers! Investigating the negative impact of outliers on active learning for visual question answering. Association for Computational Linguistics (ACL) 2021

Common AL Techniques Have Mixed Results

Why might AL methods perform comparable or worse to random selection?

- Challenging examples to learn are sampled; e.g.,

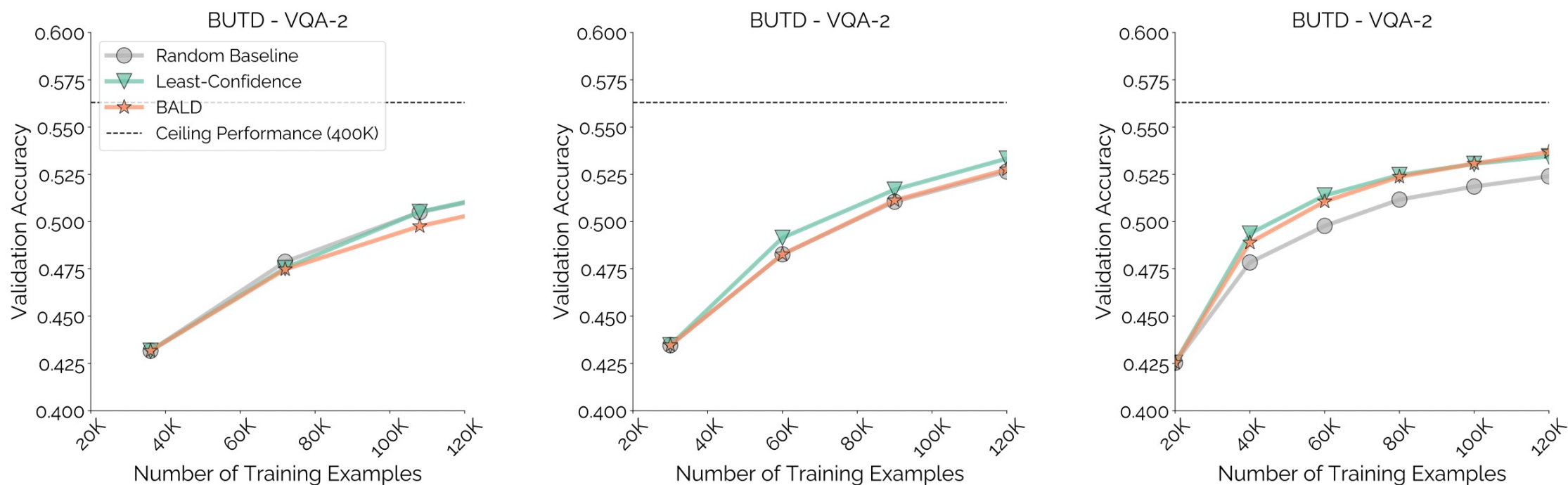
VQA-2		External knowledge: What does the symbol on the blanket mean?		OCR: What is the first word on the black car?
GQA		Underspecification: What is on the shelf?		Multi-hop reasoning: What is the vehicle that is driving down the road the box is on the side of?

Figure 7: Example groups of collective outliers in the VQA-2 and GQA datasets.

Karamcheti et al. Mind your outliers! Investigating the negative impact of outliers on active learning for visual question answering. Association for Computational Linguistics (ACL) 2021

Idea: Remove “Unlearnable” Data from Pool

What is the performance trend for AL approaches compared to random selection when removing “challenging” examples from data pool?



(a) 10% of Dataset Removed

(b) 25% of Dataset Removed

(c) 50% of Dataset Removed

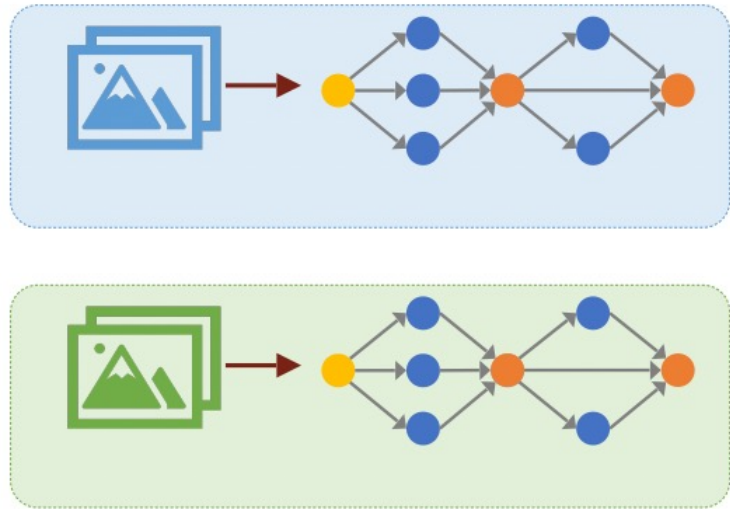
Karamcheti et al. Mind your outliers! Investigating the negative impact of outliers on active learning for visual question answering. Association for Computational Linguistics (ACL) 2021

Today's Topics

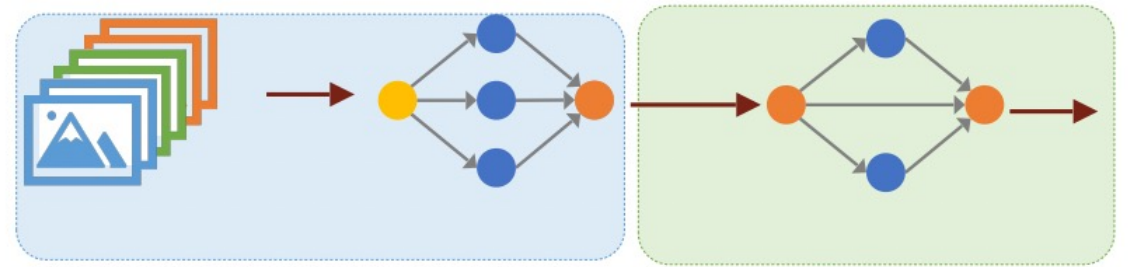
- Motivation
- Efficient learning: curriculum learning
- Efficient learning: active learning
- **Efficient learning: other considerations**
- Faculty course questionnaire

How to teach machines so they learn
(1) faster and (2) with fewer resources?

Distributed Training



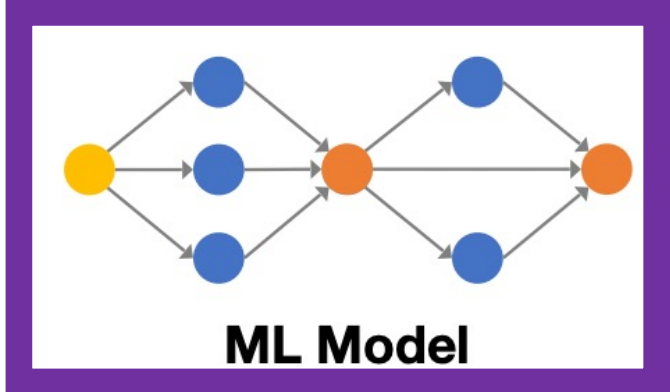
Data Parallelism:



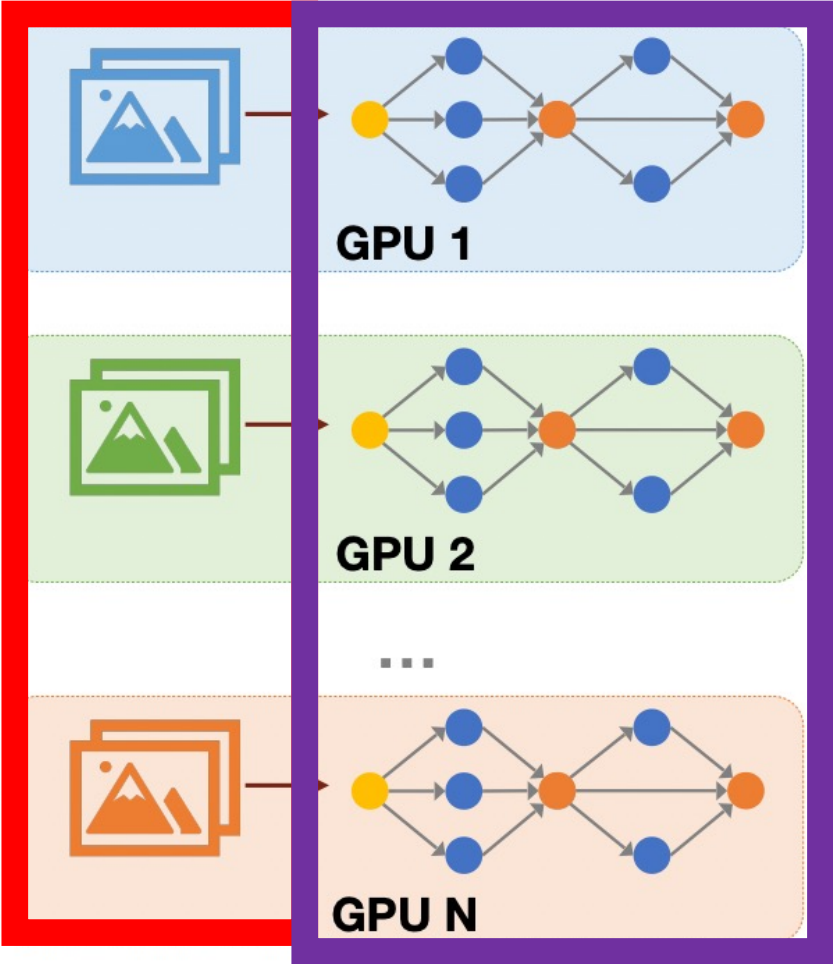
Model Parallelism:

Distributed Training: Data Parallelism

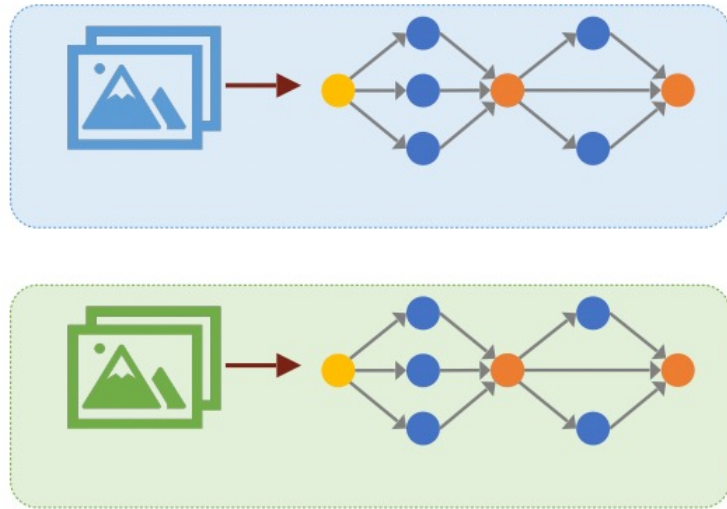
Model copied across GPUs



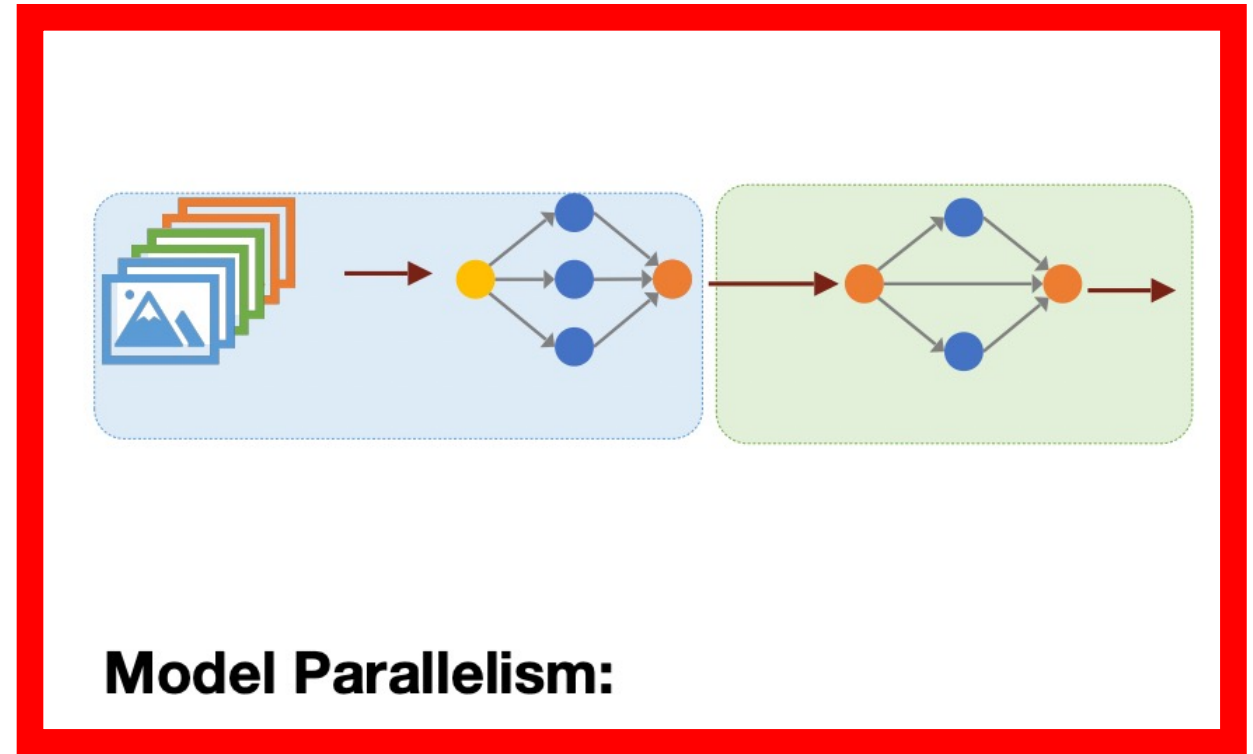
Data split across GPUs



Distributed Training



Data Parallelism:



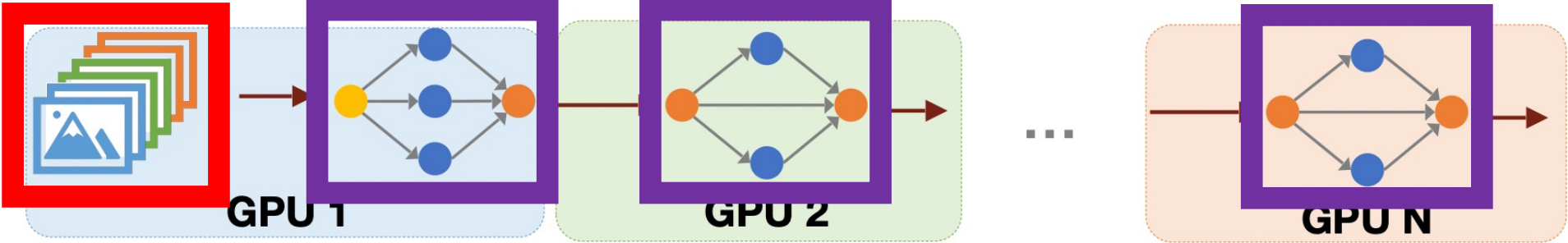
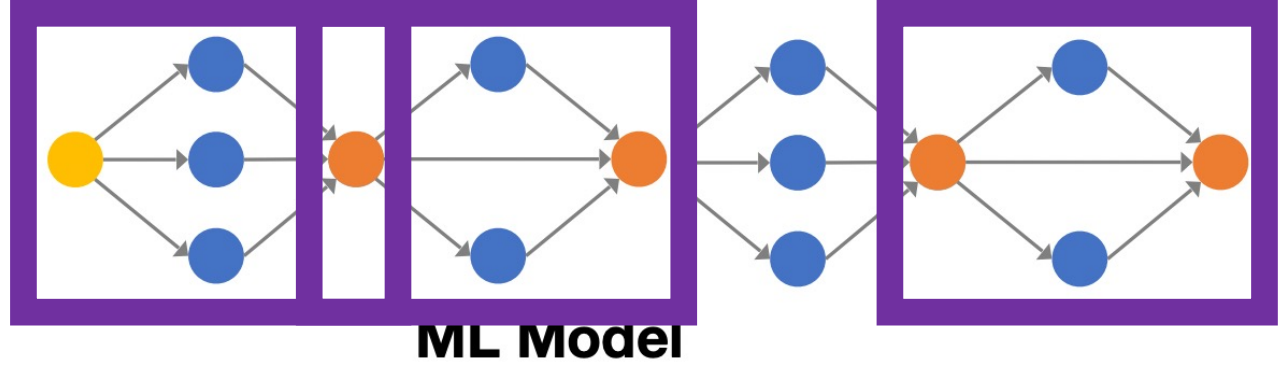
Model Parallelism:

Distributed Training: Model Parallelism

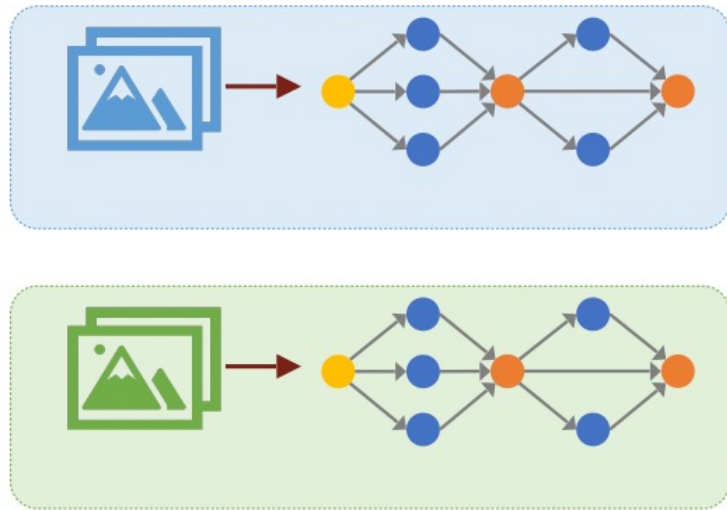


One copy of the data

Model split across GPUs

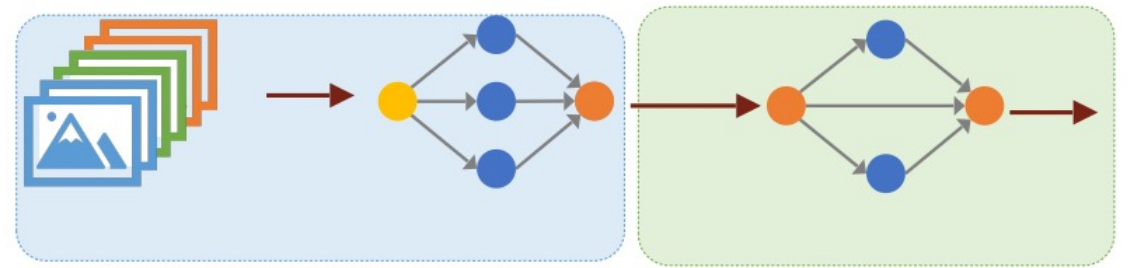


Distributed Training



Data Parallelism:

- Split the data
- Same model across devices
- Easy to parallelize, high utilization
- N copies of model



Model Parallelism:

- Split the model
- Move activations through devices
- Hard to parallelize, load balancing issue
- Single copy of model

How to teach machines so they learn
(1) faster and (2) with fewer resources

On-Device Training: Learn with Limited Memory and Compute



Laptop



Smart phone



Security camera



Snowball



Crop field analysis



Media server



Autonomous cars

DL systems may need to adapt to users' data on-device for reasons such as poor/no internet connection and privacy

Tutorial: <https://hanlab.mit.edu/files/course/slides/MIT-TinyML-Lec15-On-Device-Training-And-Transfer-Learning-I.pdf>

Figure: <https://aws.amazon.com/blogs/machine-learning/demystifying-machine-learning-at-the-edge-through-real-use-cases/>

Today's Topics

- Motivation
- Efficient learning: curriculum learning
- Efficient learning: active learning
- Efficient learning: other considerations
- Faculty course questionnaire: <https://colorado.campuslabs.com/courseeval>

A film strip border with white sprocket holes on a dark grey background. In the center, a soft, circular white glow contains the text "The End" in a white, elegant script font.

The End