

CSC413/2516 Lecture 6: Interpretability and Large-scale Generative Models

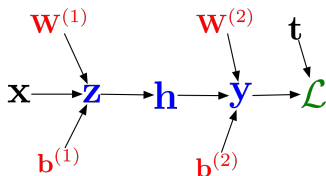
Jimmy Ba and Bo Wang

Overview

- We've seen convolutional neural networks are very successful computer vision models.
 - But, how do we know the network has learnt useful patterns from the training set?
- The interpretation of deep learning models is a challenge due to their size, complexity, and often opaque internal state.
- In this lecture, we discuss a few some tools to help understand the behavior of ML models.

Overview

- Recall the computation graph:



- From this graph, you could compute $\partial\mathcal{L}/\partial\mathbf{x}$, but we never made use of this.
- Basic idea: $\partial\mathcal{L}/\partial\mathbf{x}$ contains the model's sensitivity wrt changes of its input. It could be useful for interpreting or breaking the model!

Overview

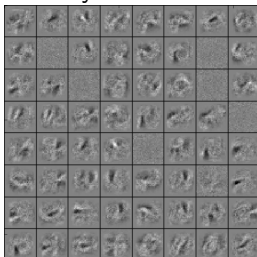
Use cases of input gradients:

- Visualizing what learned features represent
 - Visualizing image gradients to give us per-image feature visualization
 - Optimizing an image to maximize activations
- Adversarial inputs
- “Deep Dream”

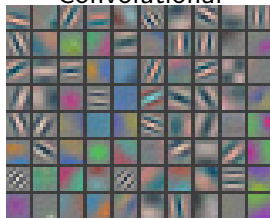
Feature Visualization

- Recall: we can understand what first-layer features are doing by visualizing the weight matrices.

Fully connected



Convolutional



- What to do with the higher-level features?

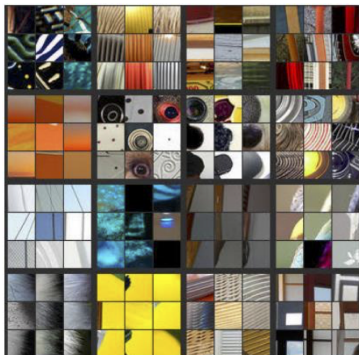
Feature Visualization

- One way to formalize: pick the images in the training set which activate a unit most strongly.
- Here's the visualization for layer 1:



Feature Visualization

- Layer 3:



Feature Visualization

- Layer 4:



Feature Visualization

- Layer 5:



Feature Visualization

- Higher layers seem to pick up more abstract, high-level information.
 - Problem: can't tell what the unit is actually responding to in the image!

Overview

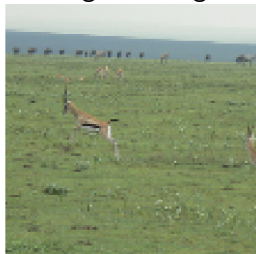
Use cases of input gradients:

- Visualizing what learned features represent
 - **Visualizing image gradients to give us per-image feature visualization**
 - Optimizing an image to maximize activations
- Adversarial inputs
- “Deep Dream”

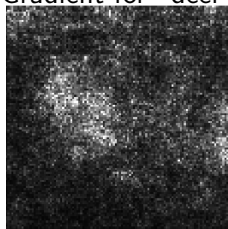
Feature Visualization

- Input gradients can be noisy and hard to interpret.
- Take a good object recognition conv net (Alex Net) and compute the gradient of $\log p(y = \text{“deer”} | x)$:

Original image



Gradient for “deer”



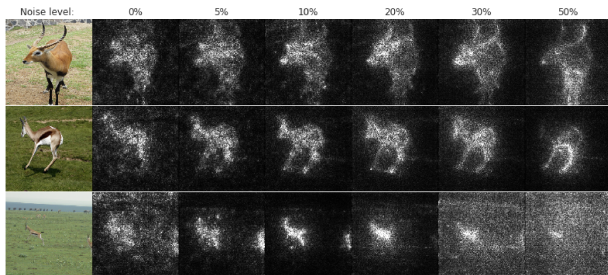
- This is partially due to the steepest directions are local sensitivity wrt the current input pixels. It is difficult to pick out important global features from one instance of the local changes.

Feature Visualization

- **SmoothGrad** is a method to estimate a global “saliency” map.
- Do the backward pass on a few noisy version of the input images, then average their input gradients.

$$S_{deer} = \frac{1}{N} \sum_{i=1}^N \frac{\partial \mathcal{L}_{deer}}{\partial x} (x + \epsilon_i), \quad \epsilon \sim \mathcal{N}(0, \sigma^2 I)$$

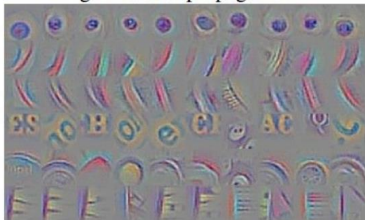
- We want to average out the local sensitivity effect from slightly perturbed input images.
- Results



Cautionary Tales of Image Gradients

This looks very convincing!

guided backpropagation



corresponding image crops



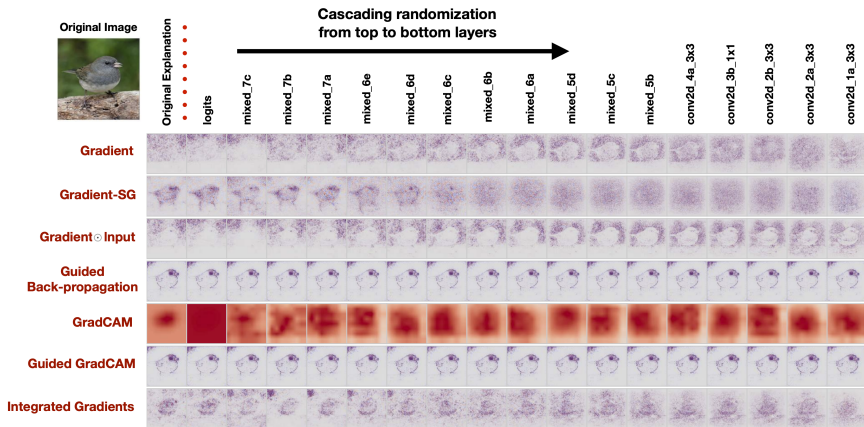
guided backpropagation



corresponding image crops

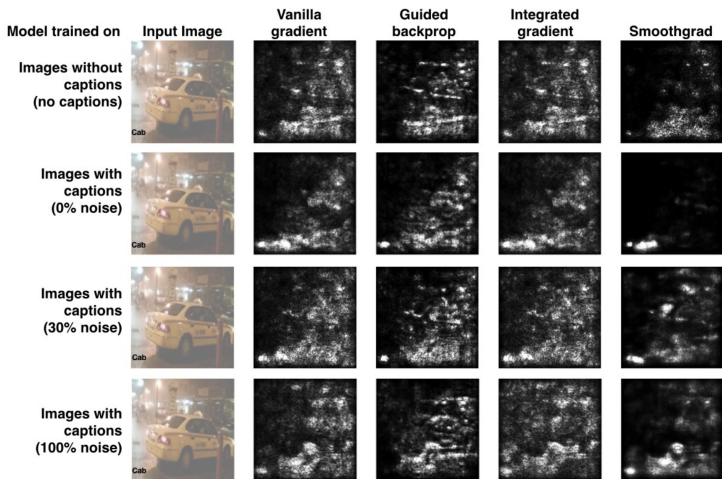


Cautionary Tales of Image Gradients



Sanity check for saliency maps, <http://papers.nips.cc/paper/8160-sanity-checks-for-saliency-maps.pdf>

Cautionary Tales of Image Gradients



Testing with Concept Activation Vectors, <https://arxiv.org/pdf/1711.11279.pdf>

Cautionary Tales of Image Gradients

Understanding the predictions from a trained neural network remains a challenging yet important problem when machine learning models:

- A **spurious correlation** is the unexpected correlation in training data that may not exist in test. E.g. the text “Cab” in the previous slides.
- How can we tell our model picked up any **spurious correlation**?
 - Visualizing image gradients may give clues about ConvNets.
 - What about Neural Language Models?
 - What about medical applications? Expert knowledge is required to discern spurious or not.

Overview

Use cases of input gradients:

- Visualizing what learned features represent
 - Visualizing image gradients to give us per-image feature visualization
 - **Optimizing an image to maximize activations**
- Adversarial inputs
- “Deep Dream”

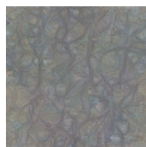
Gradient Ascent on Images

- Can do gradient ascent on an image to maximize the activation of a given neuron.
- Requires a few tricks to make this work; see <https://distill.pub/2017/feature-visualization/>

Starting from random noise, we optimize an image to activate a particular neuron (layer mixed4a, unit 11).



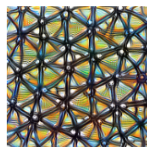
Step 0



Step 4



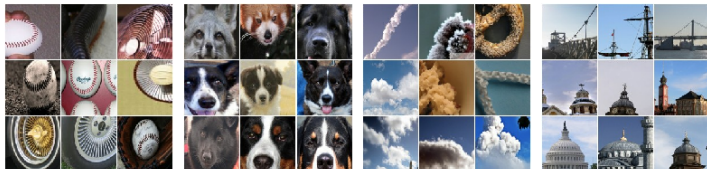
Step 48



Step 2048

Gradient Ascent on Images

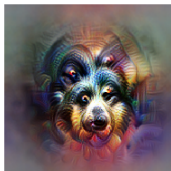
Dataset Examples show us what neurons respond to in practice



Optimization isolates the causes of behavior from mere correlations. A neuron may not be detecting what you initially thought.



Baseball—or stripes?
mixed4a, Unit 6



Animal faces—or snouts?
mixed4a, Unit 240



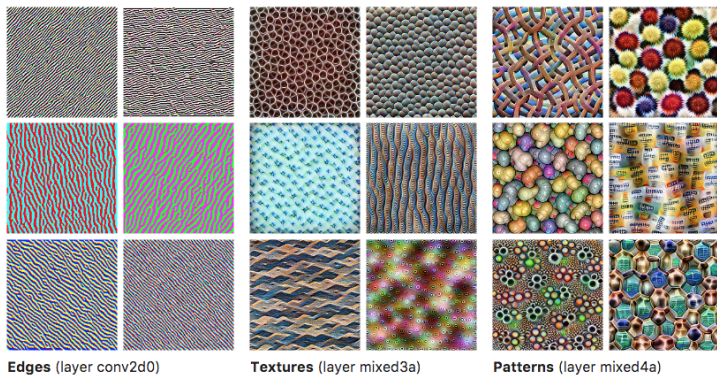
Clouds—or fluffiness?
mixed4a, Unit 453



Buildings—or sky?
mixed4a, Unit 492

Gradient Ascent on Images

- Higher layers in the network often learn higher-level, more interpretable representations



<https://distill.pub/2017/feature-visualization/>

Gradient Ascent on Images

- Higher layers in the network often learn higher-level, more interpretable representations



Parts (layers mixed4b & mixed4c)

Objects (layers mixed4d & mixed4e)

<https://distill.pub/2017/feature-visualization/>

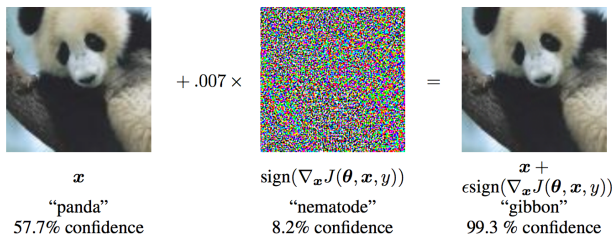
Overview

Use cases of input gradients:

- Visualizing what learned features represent
 - Visualizing image gradients to give us per-image feature visualization
 - Optimizing an image to maximize activations
- **Adversarial inputs**
- “Deep Dream”

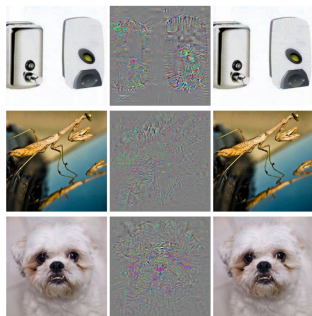
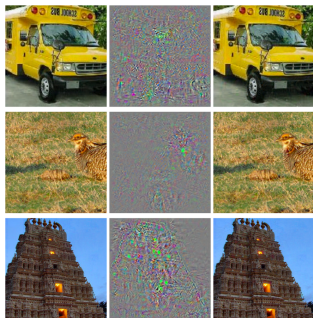
Adversarial Examples

- One of the most surprising findings about neural nets has been the existence of **adversarial inputs**, i.e. inputs optimized to fool an algorithm.
- Given an image for one category (e.g. “cat”), compute the image gradient to maximize the network’s output unit for a different category (e.g. “dog”)
 - Perturb the image very slightly in this direction, and chances are, the network will think it’s a dog!
 - Works slightly better if you take the sign of the entries in the gradient; this is called the **fast gradient sign method**.



Adversarial Examples

- The following adversarial examples are misclassified as ostriches. (Middle = perturbation $\times 10$.)



Adversarial Examples

- 2013: ha ha, how cute!
 - The paper which introduced adversarial examples was titled “Intriguing Properties of Neural Networks.”

Adversarial Examples

- 2013: ha ha, how cute!
 - The paper which introduced adversarial examples was titled “Intriguing Properties of Neural Networks.”
- 2018: serious security threat
 - Nobody has found a reliable method yet to defend against them.
 - 7 of 8 proposed defenses accepted to ICLR 2018 were cracked within days.
 - Adversarial examples transfer to different networks trained on a totally separate training set!
 - You don't need access to the original network; you can train up a new network to match its predictions, and then construct adversarial examples for that.
 - Attack carried out against proprietary classification networks accessed using prediction APIs (MetaMind, Amazon, Google)

Adversarial Examples

- You can print out an adversarial image and take a picture of it, and it still works!



- Can someone paint over a stop sign to fool a self-driving car?

Adversarial Examples

- An adversarial example in the physical world (network thinks it's a gun, from a variety of viewing angles!)



Overview

Use cases of input gradients:

- Visualizing what learned features represent
 - Visualizing image gradients to give us per-image feature visualization
 - Optimizing an image to maximize activations
- Adversarial inputs
- **“Deep Dream”**

Deep Dream

- Start with an image, and run a conv net on it.
- Pick a layer in the network.
- Change the image such that units which were already highly activated get activated even more strongly. “Rich get richer.”
 - I.e., set $\bar{h} = h$, and then do backprop.
 - Aside: this is a situation where you'd pass in something other than 1 to `backward_pass` in autograd.
- Repeat.
- This will accentuate whatever features of an image already kind of resemble the object.

Deep Dream



Deep Dream



"Admiral Dog!"



"The Pig-Snail"



"The Camel-Bird"



"The Dog-Fish"

Deep Dream



After the break: **Large-scale (Generative) Models**

Large-scale (Generative) Models

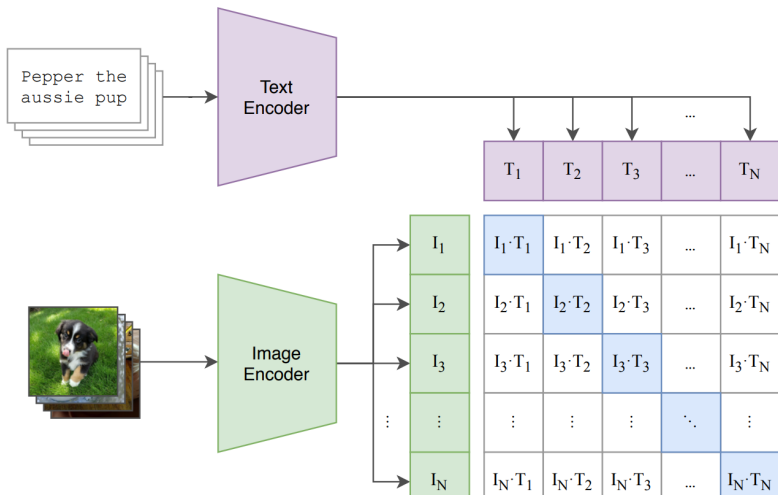
- Let's now take a look at the recent advances in internet-scale deep learning models.
 - We are already familiar with LLMs, e.g ChatGPT, from the assignment.
- This lecture hopes to shed some light on the engineering choices that contributed to the success of these models.

Confusing Terminology

- **Large Language Models (LLMs)** refers to large-scale neural networks (usually $> 10B$ parameters) trained on language modeling objective.
 - **Pre-trained Language Models (PLMs)** is the special case of LLMs often used in NLP community.
- **Foundation Models** is the umbrella name given to a wide range of pre-trained internet-scale models beyond just language models, such as vision, speech, robotics.

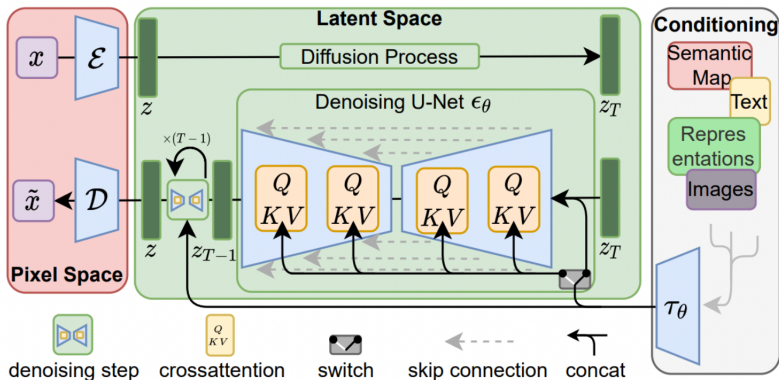
Examples of Large-scale Models

- **Contrastive Language-Image Pre-training (CLIP)** is a large-scale image-text model that learns from 400 millions of training examples.



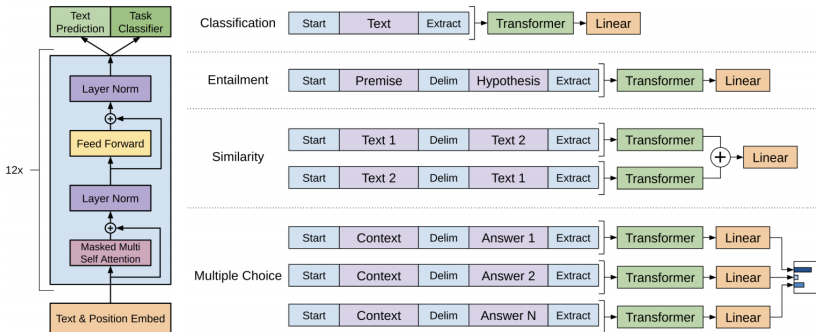
Examples of Large-scale Models

- **Stable Diffusion** is a large-scale image generator uses a pre-trained CLIP model.



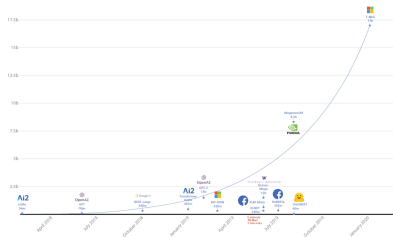
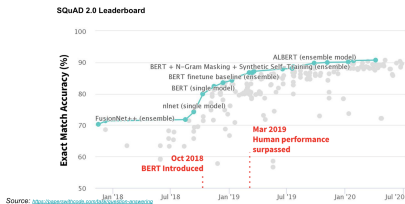
Examples of Large-scale Models

- Generative Pre-trained Transformer (GPT) is a large-scale language model trained on 100 billions of internet text.



Performance

- We had an inflection point in NLP around the end of 2018.



Performance

- We had an inflection point in NLP around the end of 2018.

Translation Benchmarks



Q&A Benchmarks



Document Classification Benchmarks



Source: <https://paperswithcode.com/>

Performance

- Unlike the previous fine-tuning approaches, pre-trained LLMs can be used directly on new tasks.

Standard Prompting

Model Input

Q: Roger has 5 tennis balls. He buys 2 more cans of tennis balls. Each can has 3 tennis balls. How many tennis balls does he have now?

A: The answer is 11.

Q: The cafeteria had 23 apples. If they used 20 to make lunch and bought 6 more, how many apples do they have?

Model Output

A: The answer is 27. ❌

Performance

- **Chain-of-thought** prompting addresses a fundamental limit in LLMs: decouple the compute and the answer length.

Chain-of-Thought Prompting

Model Input

Q: Roger has 5 tennis balls. He buys 2 more cans of tennis balls. Each can has 3 tennis balls. How many tennis balls does he have now?

A: Roger started with 5 balls. 2 cans of 3 tennis balls each is 6 tennis balls. $5 + 6 = 11$. The answer is 11.

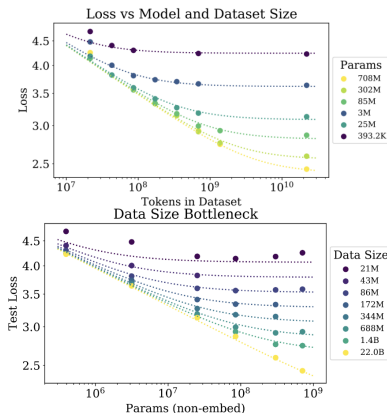
Q: The cafeteria had 23 apples. If they used 20 to make lunch and bought 6 more, how many apples do they have?

Model Output

A: The cafeteria had 23 apples originally. They used 20 to make lunch. So they had $23 - 20 = 3$. They bought 6 more apples, so they have $3 + 6 = 9$. The answer is 9. ✓

Scaling Law Curves

- These models are so expensive to train now. Researchers realized they need well-thoughtout experiments.



N: #parameters N
D: the dataset size

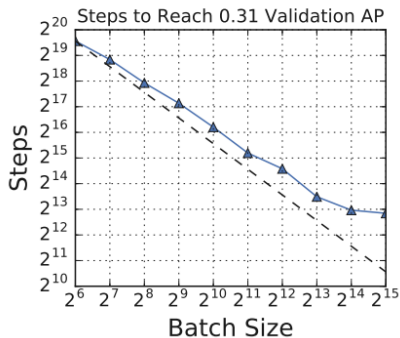
$$L(N, D) = \left[\left(\frac{N_c}{N} \right)^{\frac{\alpha_N}{\alpha_D}} + \frac{D_c}{D} \right]^{\alpha_D}$$

Power-law scaling:

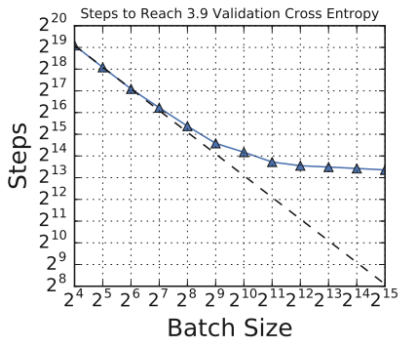
Test loss scales as 1/x to both **parameters** and **dataset size** individually till one of them plateaus.

Scaling Curve in Algorithms

- We have seen batch size scaling before. Let us re-exam its effect.



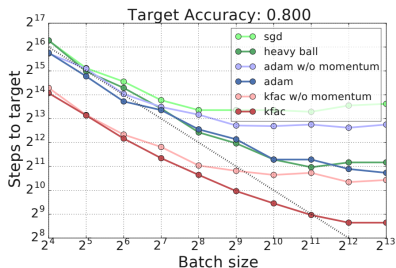
(e) ResNet-50 on Open Images



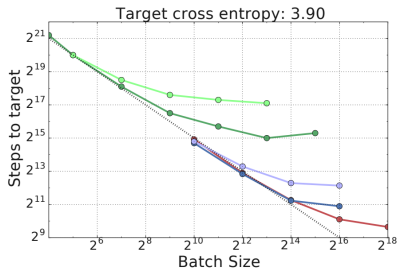
(f) Transformer on LM1B

Scaling Curve in Algorithms

- Not only we see scaling in data and model size, learning algorithms can also be a bottleneck.



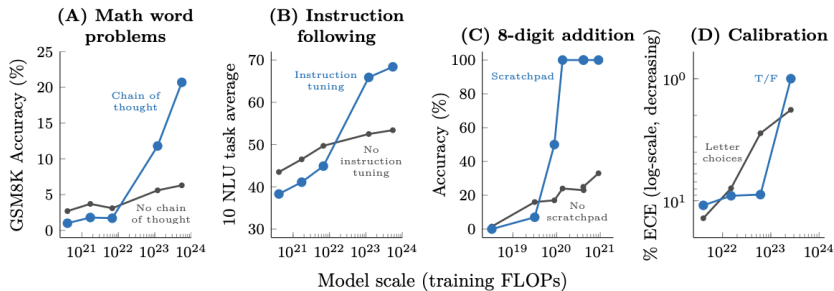
(c) ResNet8 on CIFAR10



(f) Transformer on LM1B

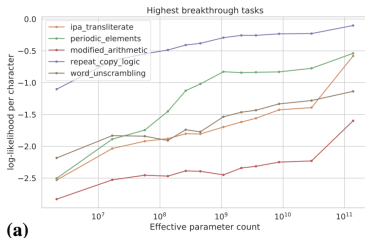
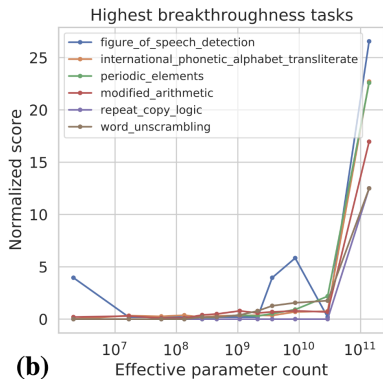
What about the Capabilities?

- Some capabilities can be emergent.



What about the Capabilities?

- Some capabilities can be emergent.



Wider and Deeper

- Understanding the limit of the current scaling is the main focus of deep learning research.

