



CS 760: Machine Learning

Transfer learning

Misha Khodak

University of Wisconsin-Madison

8 December 2025

Announcements

- **Exam:** 5:05 – 7:05 PM on Dec 16th in Sewell 5208
 - covers all course topics, with emphasis on second half
 - otherwise same policies as midterm
- What should we review on Wednesday?
 1. neural network basics
 2. CNNs
 3. RNNs, Transformers, LMs
 4. generative models
 5. learning theory
 6. kernel methods
 7. reinforcement learning
 8. data-efficient learning

Outline

- **Transfer learning**

- motivations, basic approaches, self-supervised learning

- **Learning across multiple tasks**

- setups, multi-task methods, meta-learning methods

- **Foundation models**

- overview, fine-tuning, in-context learning

Outline

- **Transfer learning**

- motivations, basic approaches, self-supervised learning

- **Learning across multiple tasks**

- setups, multi-task methods, meta-learning methods

- **Foundation models**

- overview, fine-tuning, in-context learning

Dealing with low-data scenarios

Numerous approaches (too many to cover in detail)

- which one to take is highly application-dependent
- can construct a basic taxonomy:

less-than-full supervision

- do more with less (labeled data)
- last week's lecture

semi-
supervised
learning

active
learning

weakly-
supervised
learning

transfer learning

- do more with more (o.o.d. data)
- today's lecture

multi-task
learning

meta
learning

foundation
models

Transfer learning

We typically assume labeled points $(x_1, y_1), \dots, (x_n, y_n) \sim D$ drawn i.i.d. from the **target distribution** D

What if:

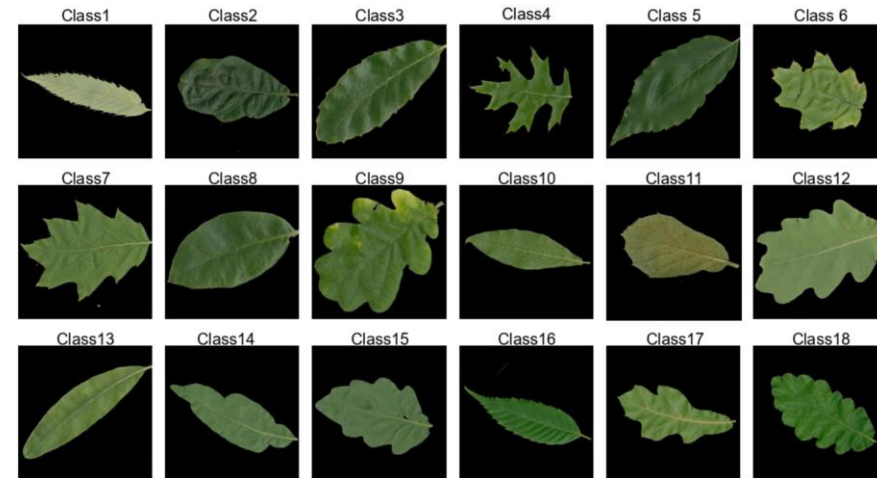
- n is too small to learn a sufficiently expressive model
- but we have access to more data $(x'_1, y'_1), \dots, (x'_N, y'_N) \sim D'$ from a **related distribution** D' ?

Using data from a related distribution to improve performance on the target distribution is **transfer learning**

Canonical example: ImageNet

standard vision pipeline:

1. collect a bunch of data for your target task
2. download a large CNN (e.g. a big ResNet) trained on ImageNet and **replace its classification layer**
3. then
 - I. either pass its **features** to a simpler model
 - II. or **fine-tune** it directly on the task



Arun et al. *J. Phytopathology*.

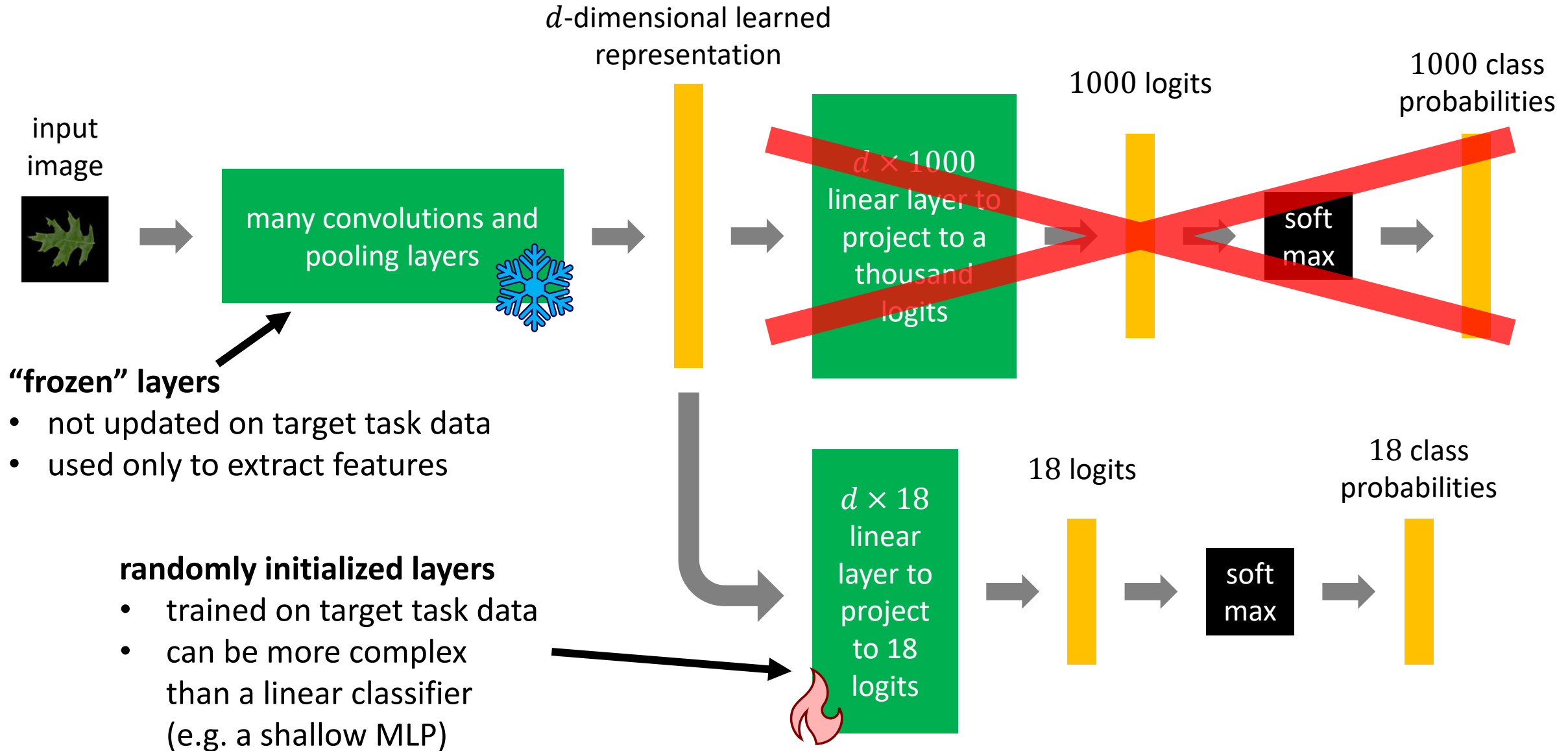
a few
datapoints
for a few
classes



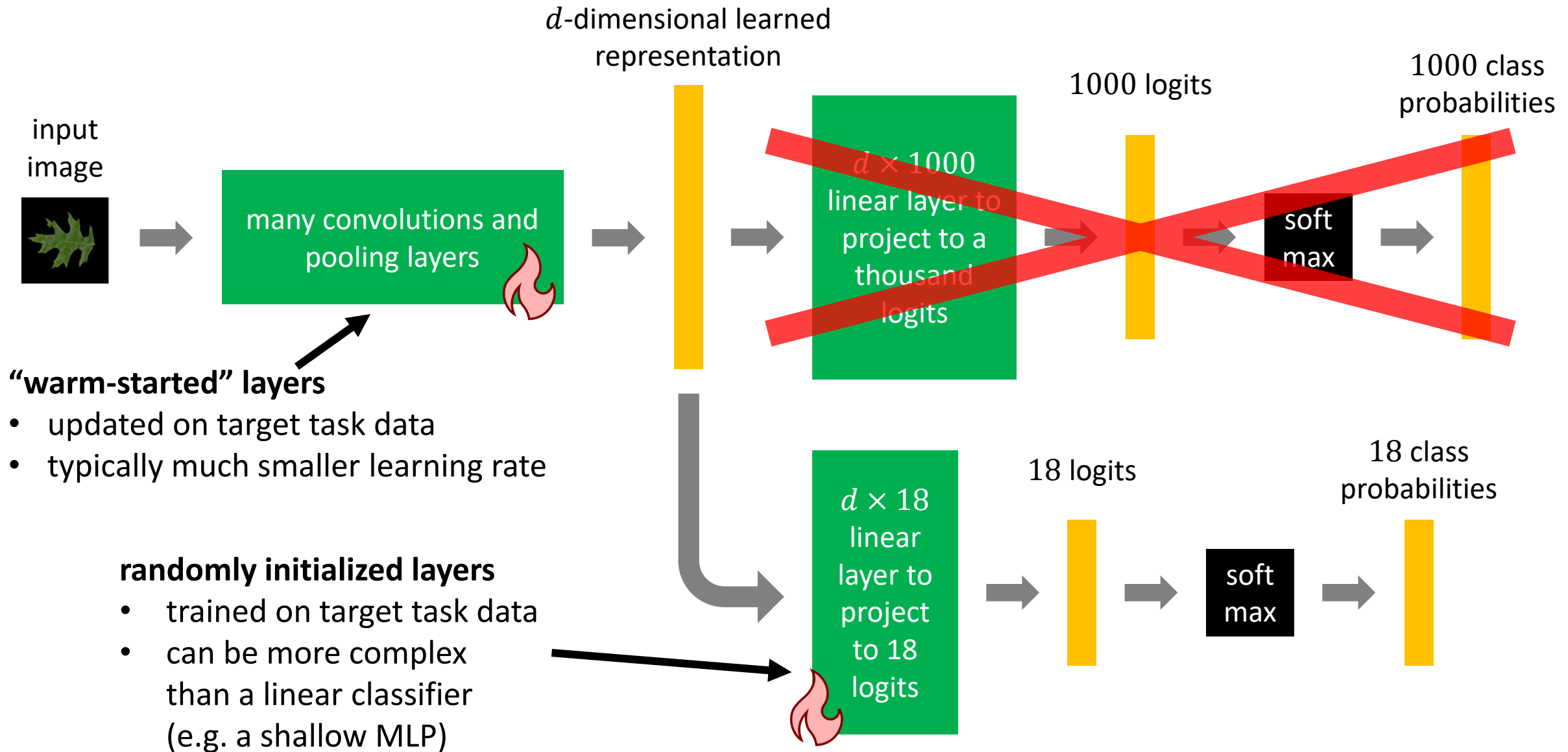
Kaggle

thousands of
datapoints
for each of a
thousand
classes

Approach I: feature extraction



Approach II: fine-tuning



Transfer learning

- Transfer learning has been hugely successful
- Numerous other potential approaches
- Big remaining question: **what if the related data lacks labels?**
 - we chop off the classification layers anyway, so we just need to extract some **representation** of the data
 - can do so using classical unsupervised learning (PCA, etc.)
 - or we can do it with **self-supervised learning (SSL)**

Self Supervision: Basic Idea

- Use domain-specific properties of the inputs (x) to create pseudo-labels (y) corresponding to “**pretext tasks**”
- Ex: predict stuff you already know

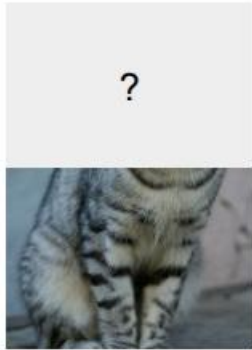
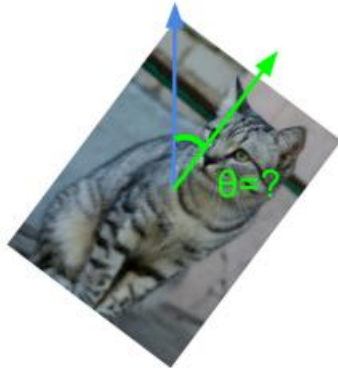


image completion
Stanford CS 231n



rotation prediction



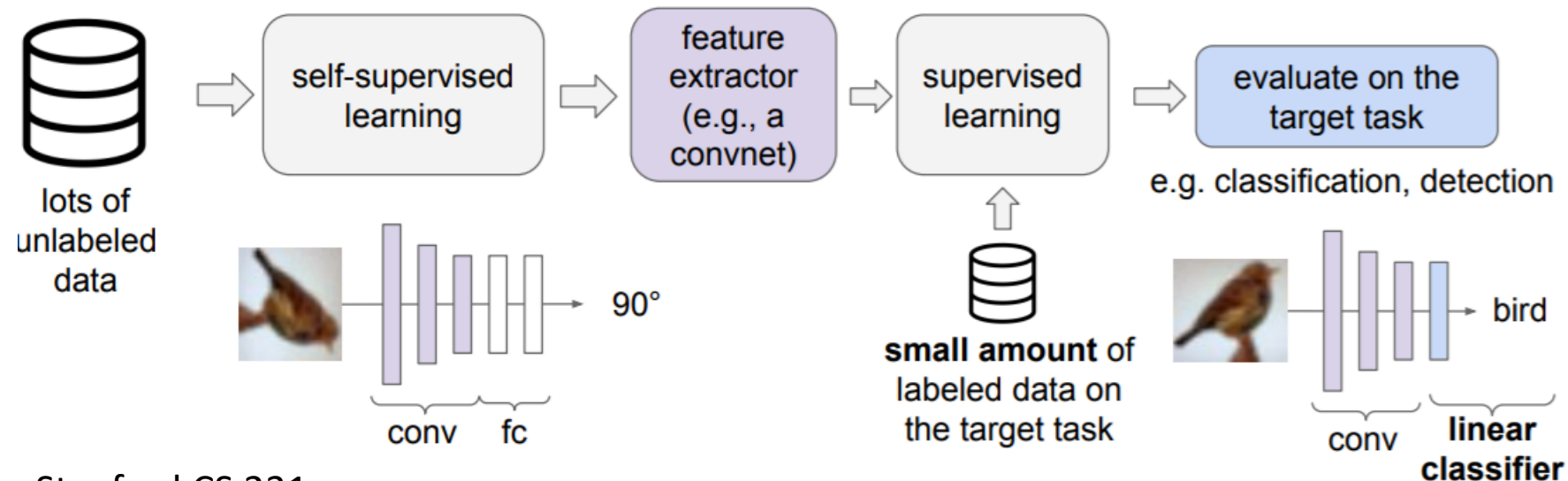
“jigsaw puzzle”



colorization

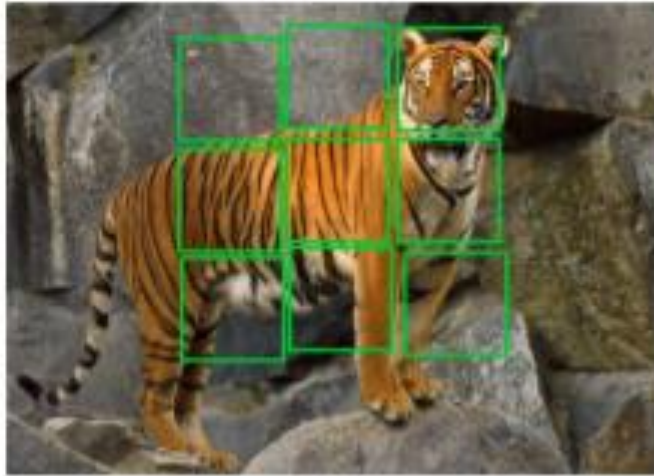
Self Supervision: Using the Representations

- Don't care specifically about our performance on pretext task
- Use the learned network as a feature extractor
- Once we have labels for a particular task, train on a small amount of data



Self Supervision: Pretext Tasks

- Lots of options for pretext tasks
 - Predict rotations
 - Coloring
 - Fill in missing portions of the image
 - Solve puzzles



(a)



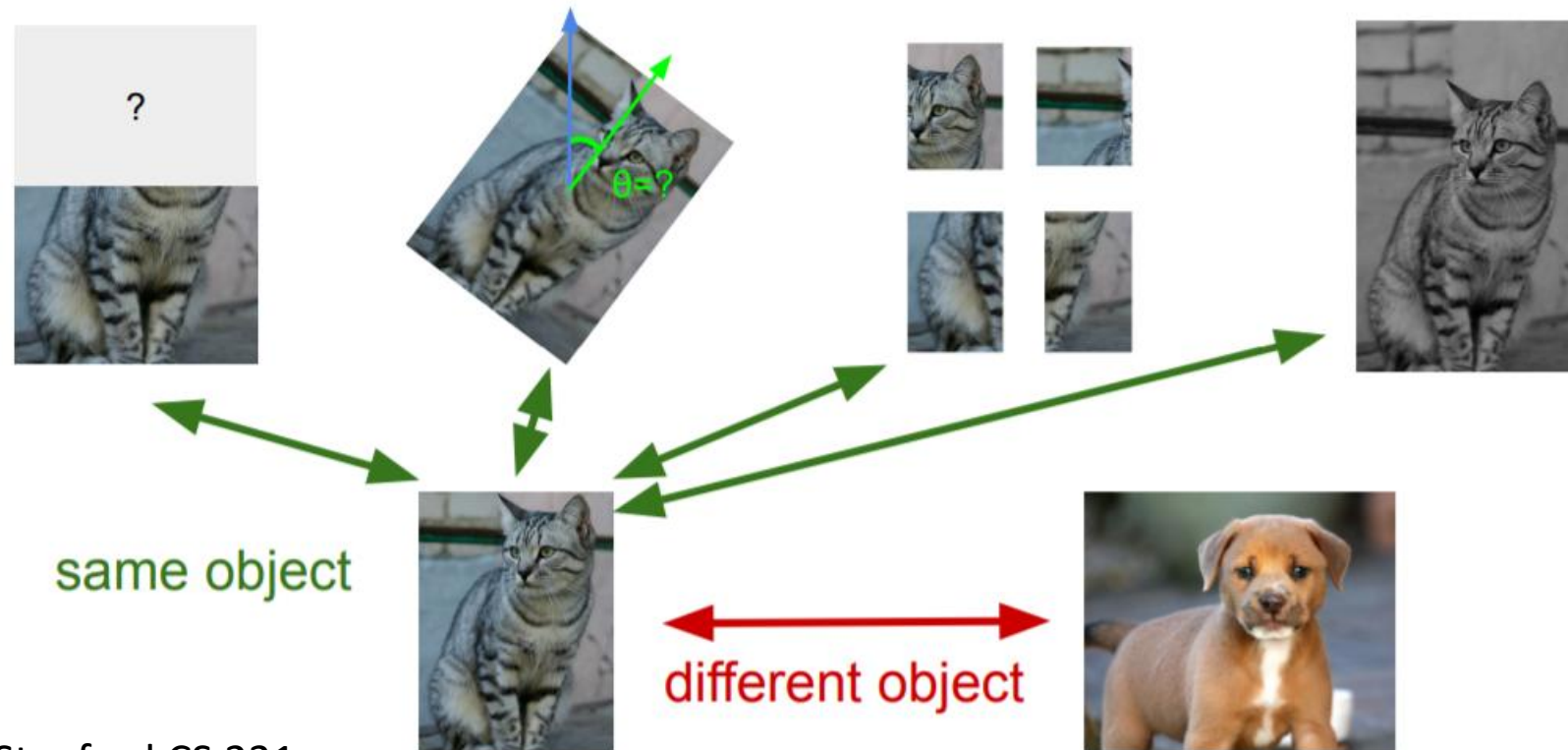
(b)



(c)

Contrastive Learning: Basics

- Type of SSL where we learn representations such that:
 - transformed versions of single sample are similar
 - different samples are different



Contrastive Learning: Motivation

- Goal:
 - Keep together related representations, push unrelated apart.
 - The InfoNCE loss function:

Van den Oord et al., 2018

$$L = -E_X \left[\log \frac{\exp(s(f(x), f(x^+)))}{\exp(s(f(x), f(x^+))) + \sum_{j=1}^{k-1} \exp(s(f(x), f(x_j^-)))} \right]$$



x



x^+

Positive sample:
keep close



Negative sample:
keep far



x



x_1^-



x_2^-



x_3^-

...

Self-supervised learning: Summary

Procedure:

- **pretrain** a network to do well on a pretext task
- **transfer** the network to your target task



Most well-known example: predict-the-next-word

Difference with regular (un)supervised training not obvious:

- sometimes pretext tasks are useful (e.g. autocomplete)
- sometimes unsupervised methods are implicitly SSL (e.g. GloVe)

Outline

- **Transfer learning**

- motivations, basic approaches, self-supervised learning

- **Learning across multiple tasks**

- setups, multi-task methods, meta-learning methods

- **Foundation models**

- overview, fine-tuning, in-context learning

Transfer learning from **multiple tasks**

What if instead of one related task with lots of data we have **many related tasks with similar amounts of data?**

Many setups:

- multi-task learning
- meta-learning
- continual learning
- lifelong learning
- ...

$$(x_{1,1}, y_{1,1}), \dots, (x_{1,n_1}, y_{1,n_1}) \sim D_1$$

$$\vdots$$

$$(x_{t,1}, y_{t,1}), \dots, (x_{t,n_t}, y_{t,n_t}) \sim D_t$$

$$\vdots$$

We'll cover two of them: **multi-task** and **meta-learning**

Multi-task learning

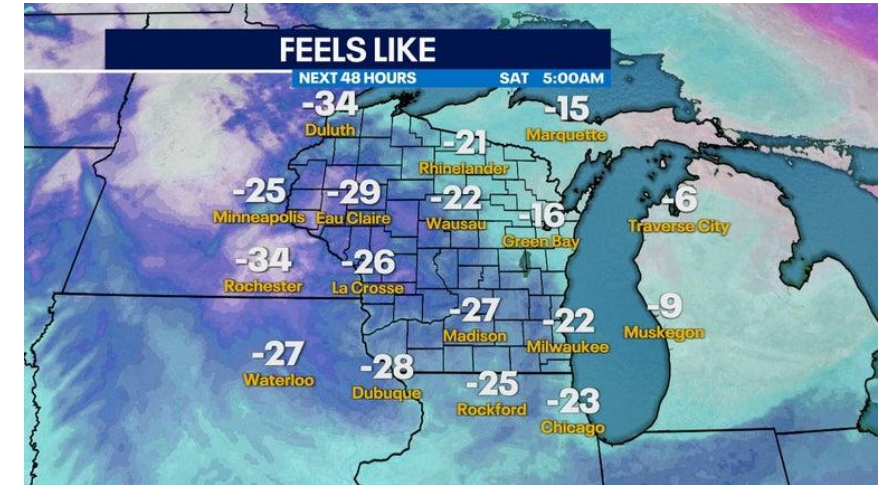
Setup: **fixed number of related tasks**

Examples:

- predict the weather in nearby cities
- diagnose patients in different hospitals

Key challenges:

- how to encode task-relationships?
- how to avoid conflicting tasks?



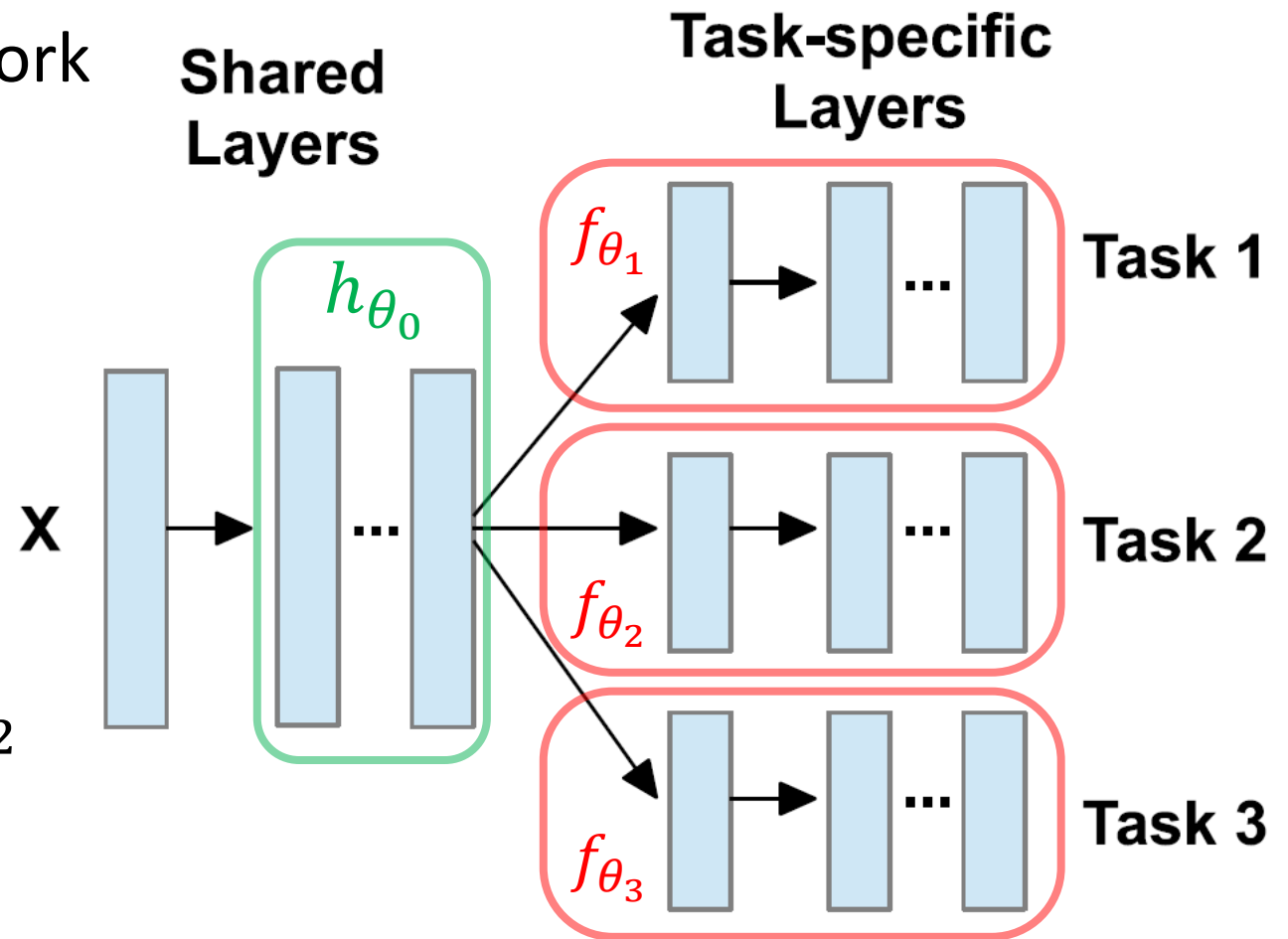
Fox6 (don't worry, this isn't the actual forecast...yet)



One common approach: **Layer-sharing**

- jointly train a multi-output network
- assumes existence of a good **shared representation** h_{θ_0}
- example objective:

$$\sum_{t=1}^T \sum_{i=1}^{n_t} \left(y_{t,i} - f_{\theta_t} \left(h_{\theta_0}(x_{t,i}) \right) \right)^2$$



Another common approach: **Regularization**

- jointly train separate networks
- regularize parameters to be closer together
- example objective:

$$\sum_{t=1}^T \sum_{i=1}^{n_t} \left(y_{t,i} - f_{\theta_t}(x_{t,i}) \right)^2 + \sum_{t=1}^T \sum_{u=t+1}^T \lambda_{t,u} \| \theta_t - \theta_u \|^2$$

- allows hand-encoding of task-relationships via the regularization strengths $\lambda_{t,u}$

Meta-learning

Setup:

- **meta-training** dataset of related tasks
- **at meta-test time** we get a new dataset $(x_1, y_1), \dots, (x_n, y_n) \sim D$
- **our goal:** low expected error on unseen examples $(x, y) \sim D$

$$(x_{1,1}, y_{1,1}), \dots, (x_{1,n_1}, y_{1,n_1}) \sim D_1$$

\vdots

$$(x_{T,1}, y_{T,1}), \dots, (x_{T,n_T}, y_{T,n_T}) \sim D_T$$

Applications:

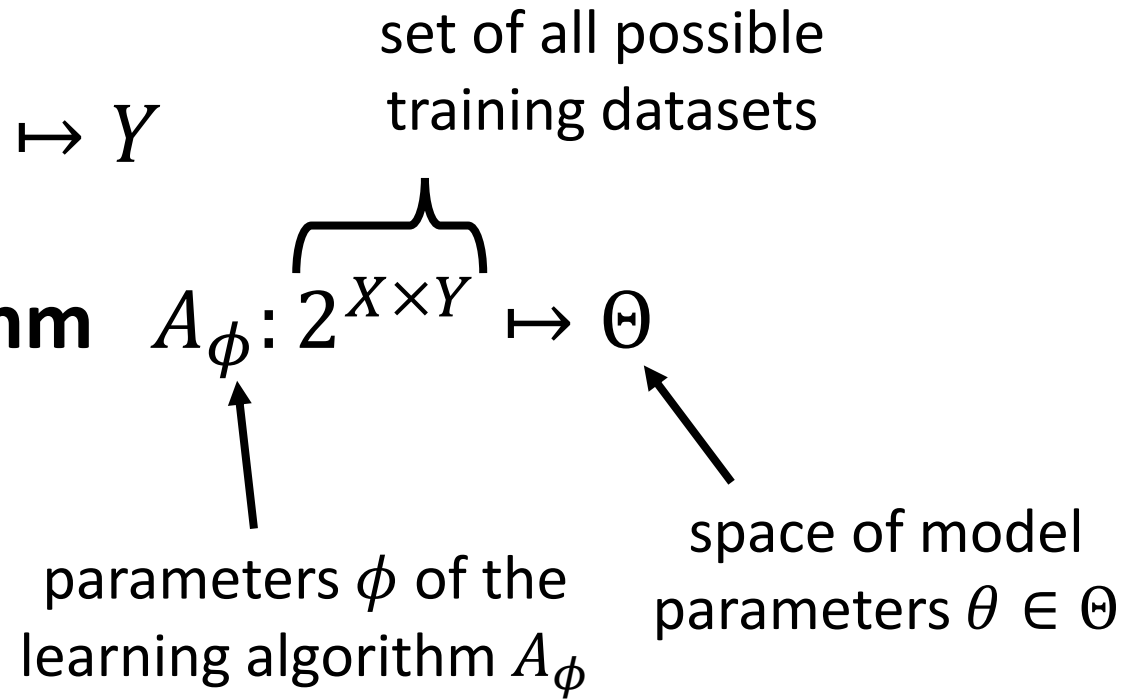
- auto-complete for new cellphone users (federated learning)
- image classification with limited labels (few-shot learning)
- robots in related environments (meta-RL)

Why is it called meta-learning?

- no longer learning a model $f_{\theta}: X \mapsto Y$

- we are learning a **learning algorithm**

- thus also called **learning-to-learn**




Example: meta-learn an initialization ϕ for gradient descent

Meta-learning (one step) gradient descent

MAML approach: $A_\phi(\{(x_i, y_i)\}_{i=1}^n) = \phi - \alpha \nabla_\phi \sum_{i=1}^n \ell(y_i, f_\phi(x_i))$

- meta-training objective:

$$\operatorname{argmin}_{\phi} \sum_{t=1}^T \sum_{i=n_t/2+1}^{n_t} \ell \left(y_{t,i}, f_{A_\phi(\{(x_{t,i}, y_{t,i})\}_{i=1}^{n_t/2})}(x_{t,i}) \right)$$


- at meta-test time:

- $\theta \leftarrow \phi - \alpha \nabla_\phi \sum_{i=1}^n \ell_i(y_i, f_\phi(x_i))$
- make predictions using $f_\theta(x)$

Outline

- **Transfer learning**

- motivations, basic approaches, self-supervised learning

- **Learning across multiple tasks**

- setups, multi-task methods, meta-learning methods

- **Foundation models**

- overview, fine-tuning, in-context learning

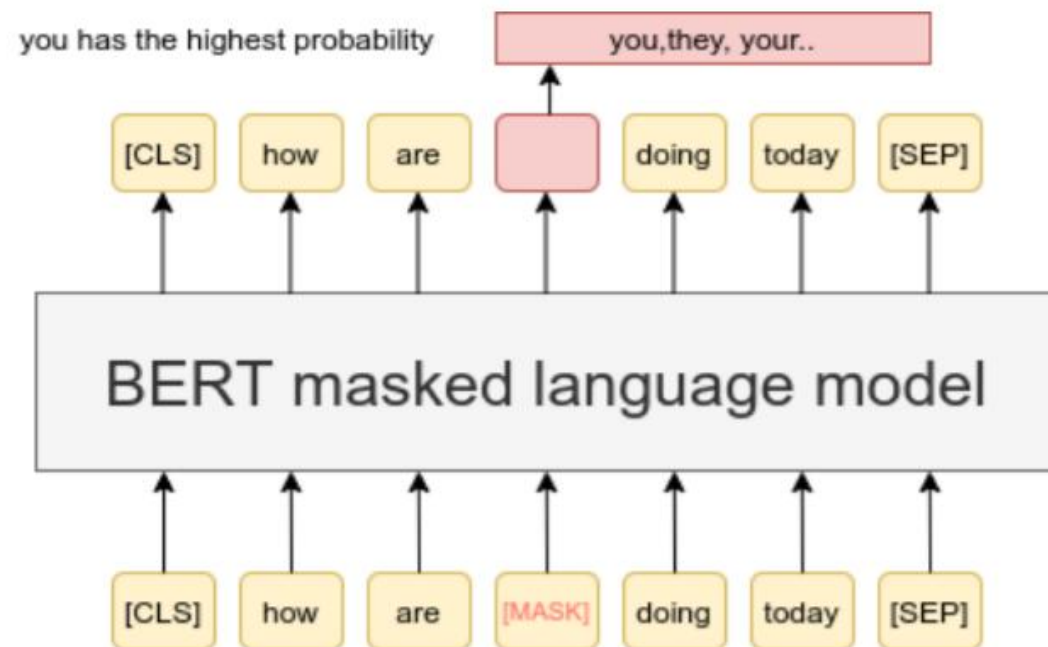
What is a foundation model?

1. take a **massive neural network**
 - older / specialized models had 100M+ params
 - latest models have 1-100 billion or more
2. **pretrain** it on Internet-scale data
3. (optionally) **post-train** on large-scaled supervised data
4. use it for transfer learning for many different tasks

Early history

2017: BERT model (340M)

- Transformer trained on masked language modeling (pretext task)
- “solved” transfer learning for language

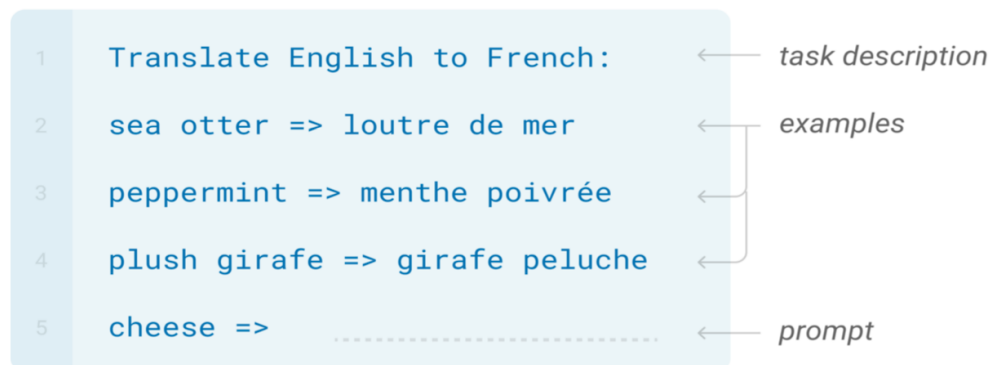


2017-present: GPT series

- Transformer trained on next-word prediction
- first observation of **in-context** learning capabilities in GPT-3
- ChatGPT post-trained on GPT-3.5

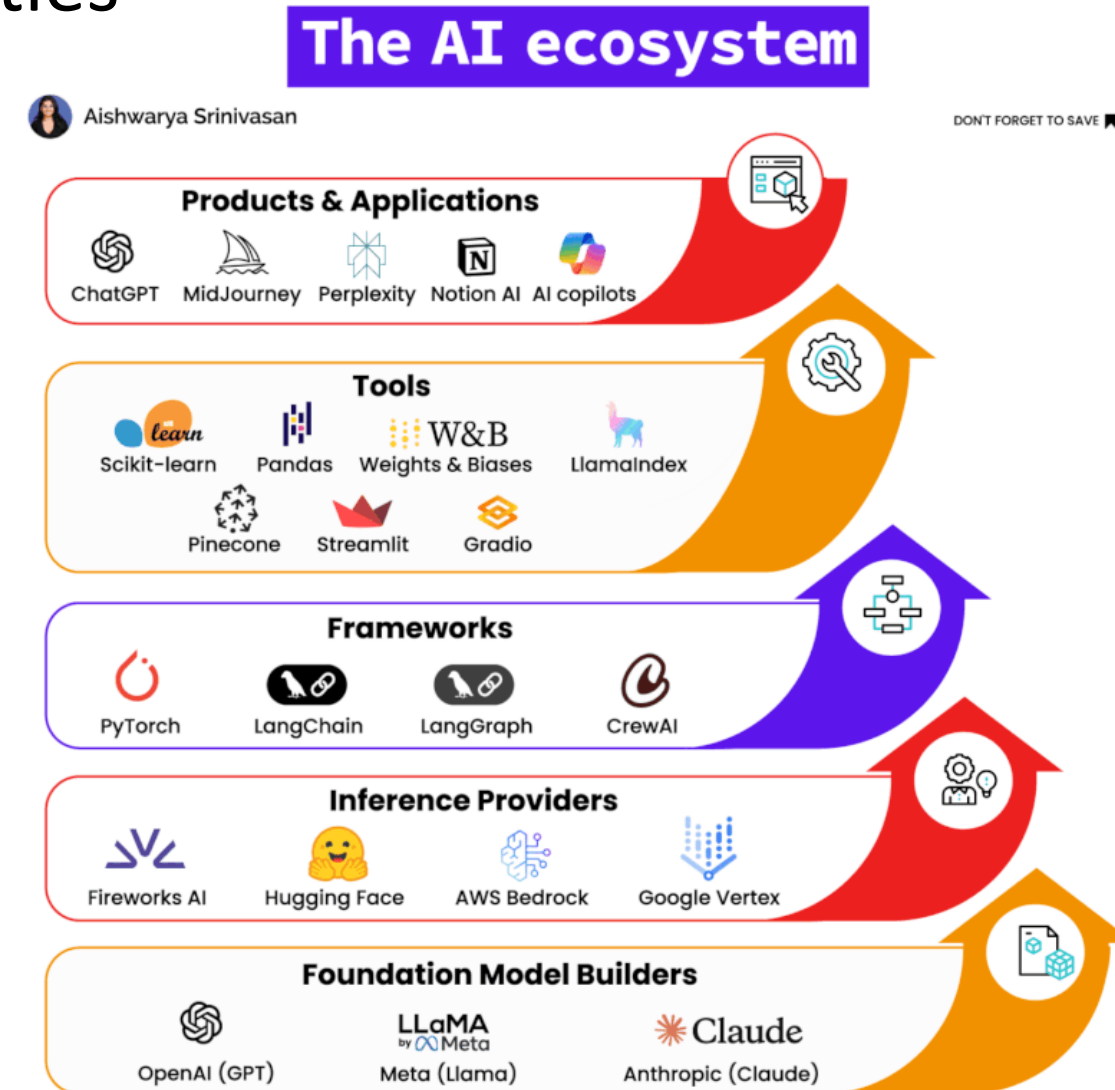
Few-shot

In addition to the task description, the model sees a few examples of the task. No gradient updates are performed.



Post-ChatGPT

- many models with varying capabilities
- closed-source models typically outperform open-source models
- new challenges:
 - **massive compute costs**
 - privacy, security, safety
- new opportunities:
 - **in-context learning**
 - reasoning



Challenge: Compute costs

pretraining FMs limited to large orgs

- one training run requires 100s of GPUs
- need many training runs (to tune) and engineers (to manage training)

even fine-tuning is hard:

- SGD on GPT-3 (175B) uses 1.2TB VRAM
- NVIDIA GPUs max out below 200GB
- what can we do?

Model	Microarchitecture	Launch	Core	Core clock (MHz)	Shaders			Memory				
					Core config [c]	Base clock (MHz)	Max boost clock (MHz) [c]	Bus type	Bus width (bit)	Size (GB)	Clock (MT/s)	Bandwidth (GB/s)
A100 GPU accelerator (PCIe card) [442][443]	Hopper	May 14, 2020 [444]	1× GA100-883AA-A1	—	10212.432 :160:432:0 (108)	765	1410	HBM2	5120	40 or 80	1215	1555
H100 GPU accelerator (PCIe card) [445]		March 22, 2022 [446]	1× GH100 [447]	—	14592.456 :24:456:0 (114)	1065	1755 CUDA 1620 TC	HBM2E	5120	80	1000	2039
H100 GPU accelerator (SXM card)				—	1065	1980 CUDA 1830 TC	HBM3	5120	64 or 80 or 96	1500	3352	
H200 GPU accelerator (PCIe card) [448]		November 18, 2024 [449]	1× GH100	—	16896.528 :24:528:0 (132)	1365	1785	HBM3E	5120	141	1313	3360
H200 GPU accelerator (SXM card)				—	1590	1980	HBM3E	5120	141	1313	3360	
H800 GPU accelerator (SXM card)		March 21, 2023 [450]	1× GH100	—	1095	1755	HBM3	5120	80	1313	3360	
L40 GPU accelerator [451]	Ada Lovelace	October 13, 2022	1× AD102 [452]	—	18176.568 :192:568:142 (142)	735	2490	GDDR6	384	48	2250	864
L4 GPU accelerator [453][454]		March 21, 2023 [455]	1x AD104 [456]	—	7424.240 :80:240:0 (60)	795	2040	GDDR6	192	24	1563	300
B100 GPU accelerator [457]	Blackwell	November 2024	2× GB102	—	2× 16896.528 :24:528:0 (132)	1665	1837	HBM3E	2× 4096	2× 96	2000	2× 4100
B200 GPU accelerator [459]		2024	2× GB100	—		1665	1837	HBM3E	2× 4096	2× 96	2000	2× 4100

Parameter-efficient fine-tuning (PEFT)

Most popular approach: **LoRA**

1. take an FM with **pretrained weight** matrices $\mathbf{W}_1, \dots, \mathbf{W}_N$
2. for each matrix $\mathbf{W}_i \in \mathbb{R}^{d \times k}$:
 - set $r \ll \min\{d, k\}$ and initialize **fine-tuning weights**:
 - $\mathbf{B}_i \in \mathbb{R}^{d \times r}$ to $\mathbf{B}_i = 0$
 - $\mathbf{A}_i \in \mathbb{R}^{r \times k}$ to $\mathbf{A}_i \sim \text{Gaussian}$
 - replace \mathbf{W}_i by $\mathbf{W}_i + \mathbf{B}_i \mathbf{A}_i$
3. fine-tune on target task but
 - freeze \mathbf{W}_i
 - update \mathbf{B}_i and \mathbf{A}_i

$$\begin{aligned} f_i(\mathbf{x}) &= \mathbf{W}_i \mathbf{x} = \begin{array}{c} \text{green box } \mathbf{W}_i \\ \text{yellow box } \mathbf{x} \end{array} \\ &\Downarrow \\ &= \mathbf{W}_i \mathbf{x} + \mathbf{B}_i \mathbf{A}_i \mathbf{x} \\ &= \begin{array}{c} \text{green box } \mathbf{W}_i \\ \text{yellow box } \mathbf{x} \end{array} + \begin{array}{c} \text{green box } \mathbf{B}_i \\ \text{green box } \mathbf{A}_i \end{array} \begin{array}{c} \text{yellow box } \mathbf{x} \end{array} \end{aligned}$$

The diagram illustrates the LoRA fine-tuning process. It shows the transformation of the original weight matrix \mathbf{W}_i into a sum of the original matrix and a product of two smaller matrices, \mathbf{B}_i and \mathbf{A}_i . The original matrix \mathbf{W}_i is represented by a green box, and the input vector \mathbf{x} is represented by a yellow box. The fine-tuning weights \mathbf{B}_i and \mathbf{A}_i are also represented by green boxes. The input vector \mathbf{x} is shown again as a yellow box. The diagram uses a large grey arrow to indicate the transformation. The final expression shows the original matrix \mathbf{W}_i (green box) and input \mathbf{x} (yellow box) being added to the product of \mathbf{B}_i (green box) and \mathbf{A}_i (green box) and \mathbf{x} (yellow box). The \mathbf{W}_i box is marked with a blue snowflake icon, indicating it is frozen. The \mathbf{B}_i and \mathbf{A}_i boxes are marked with red flame icons, indicating they are updated.

How does LoRA save memory?

- original weights $\mathbf{W}_i \in \mathbb{R}^{d \times k}$ have dk trainable params
- new weights $\mathbf{B}_i \in \mathbb{R}^{d \times r}$ and $\mathbf{A}_i \in \mathbb{R}^{r \times k}$ have $(d + k)r$
- typical values in GPT-3 175B:
 - $d \approx k \approx 10^4$
 - $r \leq 10$
- $\geq 10^4$ x fewer trainable params!
- 3x less fine-tuning VRAM

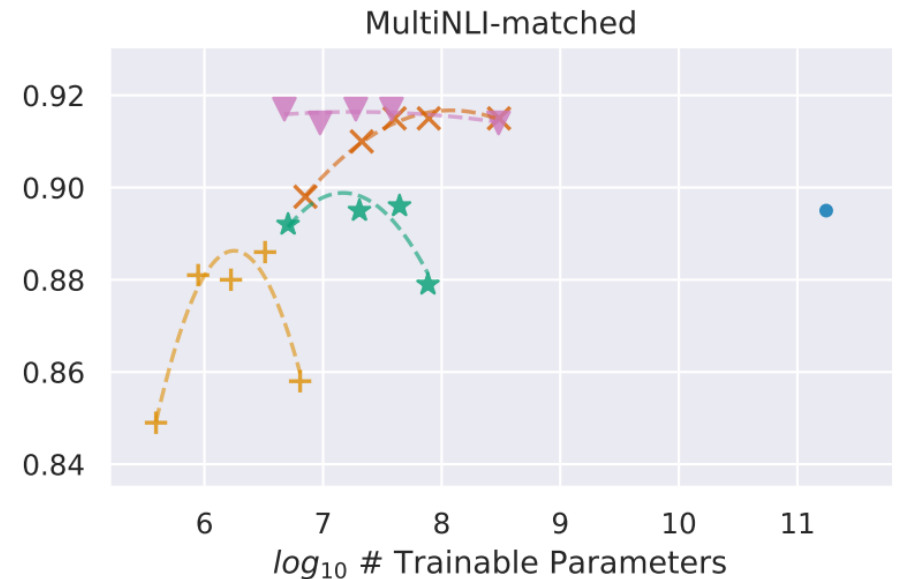
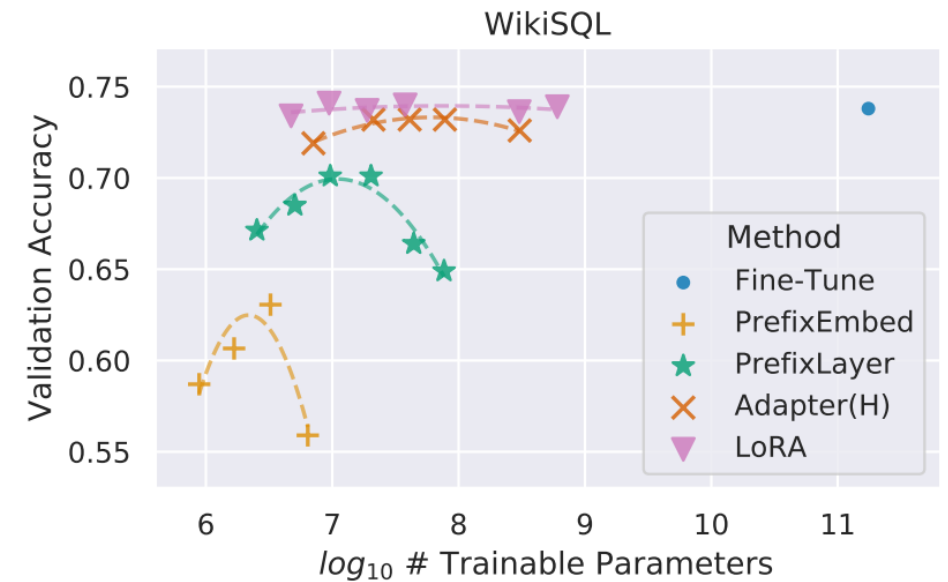
$$\begin{aligned} f_i(\mathbf{x}) &= \mathbf{W}_i \mathbf{x} = \begin{array}{|c|} \hline \mathbf{W}_i \\ \hline \end{array} \begin{array}{|c|} \hline \mathbf{x} \\ \hline \end{array} \\ &\Downarrow \\ &= \mathbf{W}_i \mathbf{x} + \mathbf{B}_i \mathbf{A}_i \mathbf{x} \\ &= \begin{array}{|c|} \hline \mathbf{W}_i \\ \hline \end{array} \begin{array}{|c|} \hline \mathbf{x} \\ \hline \end{array} + \begin{array}{|c|} \hline \mathbf{B}_i \\ \hline \end{array} \begin{array}{|c|} \hline \mathbf{A}_i \\ \hline \end{array} \begin{array}{|c|} \hline \mathbf{x} \\ \hline \end{array} \end{aligned}$$

The diagram illustrates the LoRA decomposition of a weight matrix \mathbf{W}_i into two smaller matrices \mathbf{B}_i and \mathbf{A}_i . The original matrix multiplication $\mathbf{W}_i \mathbf{x}$ is shown as a large green square \mathbf{W}_i multiplied by a yellow vertical rectangle \mathbf{x} . A large grey arrow points down to the decomposed form $\mathbf{W}_i \mathbf{x} + \mathbf{B}_i \mathbf{A}_i \mathbf{x}$. This is further visualized as the sum of two products: $\mathbf{W}_i \mathbf{x}$ (green square \mathbf{W}_i and yellow rectangle \mathbf{x}) and $\mathbf{B}_i \mathbf{A}_i \mathbf{x}$ (green rectangle \mathbf{B}_i , green rectangle \mathbf{A}_i , and yellow rectangle \mathbf{x}). A blue snowflake icon is placed below the \mathbf{W}_i matrix, and a red flame icon is placed below the \mathbf{A}_i matrix, symbolizing the reduction in memory usage (cold) and the increase in trainable parameters (hot) respectively.

Does LoRA affect accuracy?

Yes, it constrains weights of the fine-tuned model:

- fine-tuned matrices $\mathbf{W}_i + \mathbf{B}_i \mathbf{A}_i$ at most a rank $r \ll \min\{d, k\}$ update away from pretrained matrices \mathbf{W}_i
- LoRA = **L**ow-**R**ank **A**daptation
- in practice do not need large r for good performance
- learning theory intuition?



Opportunity: In-context learning

Observation: the perfect next-word predictor can be **prompted** to answer any question correctly

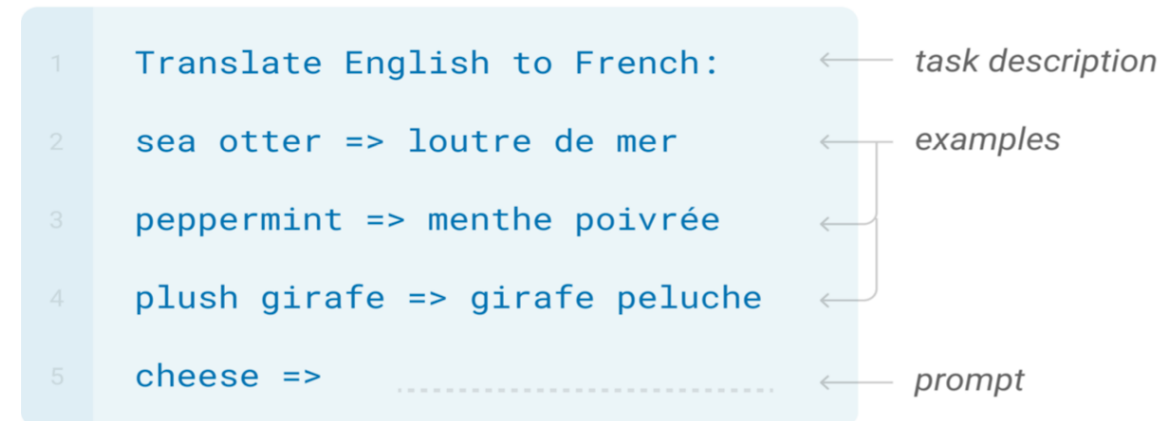
Idea: in-context learning

1. encode task instructions and data as a **context** sequence
2. make the FM generate the remainder of the sequence

Enables learning with target data
without updating the weights at all!

Few-shot

In addition to the task description, the model sees a few examples of the task. No gradient updates are performed.



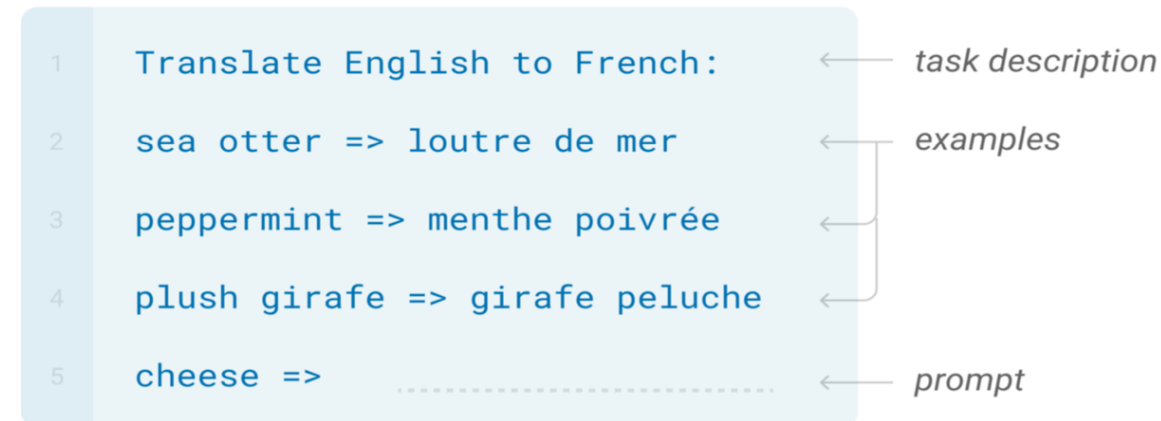
Opportunity: In-context learning

Usefulness:

- handles tasks with diverse input and output structures
- directly incorporates pretraining knowledge
- enables multi-step reasoning

Few-shot

In addition to the task description, the model sees a few examples of the task. No gradient updates are performed.





Thanks Everyone!

Some of the slides in these lectures have been adapted/borrowed from materials developed by Mark Craven, David Page, Jude Shavlik, Tom Mitchell, Nina Balcan, Elad Hazan, Tom Dietterich, Pedro Domingos, Jerry Zhu, Yingyu Liang, Volodymyr Kuleshov, Fred Sala