

10-301/601: Introduction to Machine Learning

Lecture 19 – Pre-training, Fine-tuning & In-Context Learning

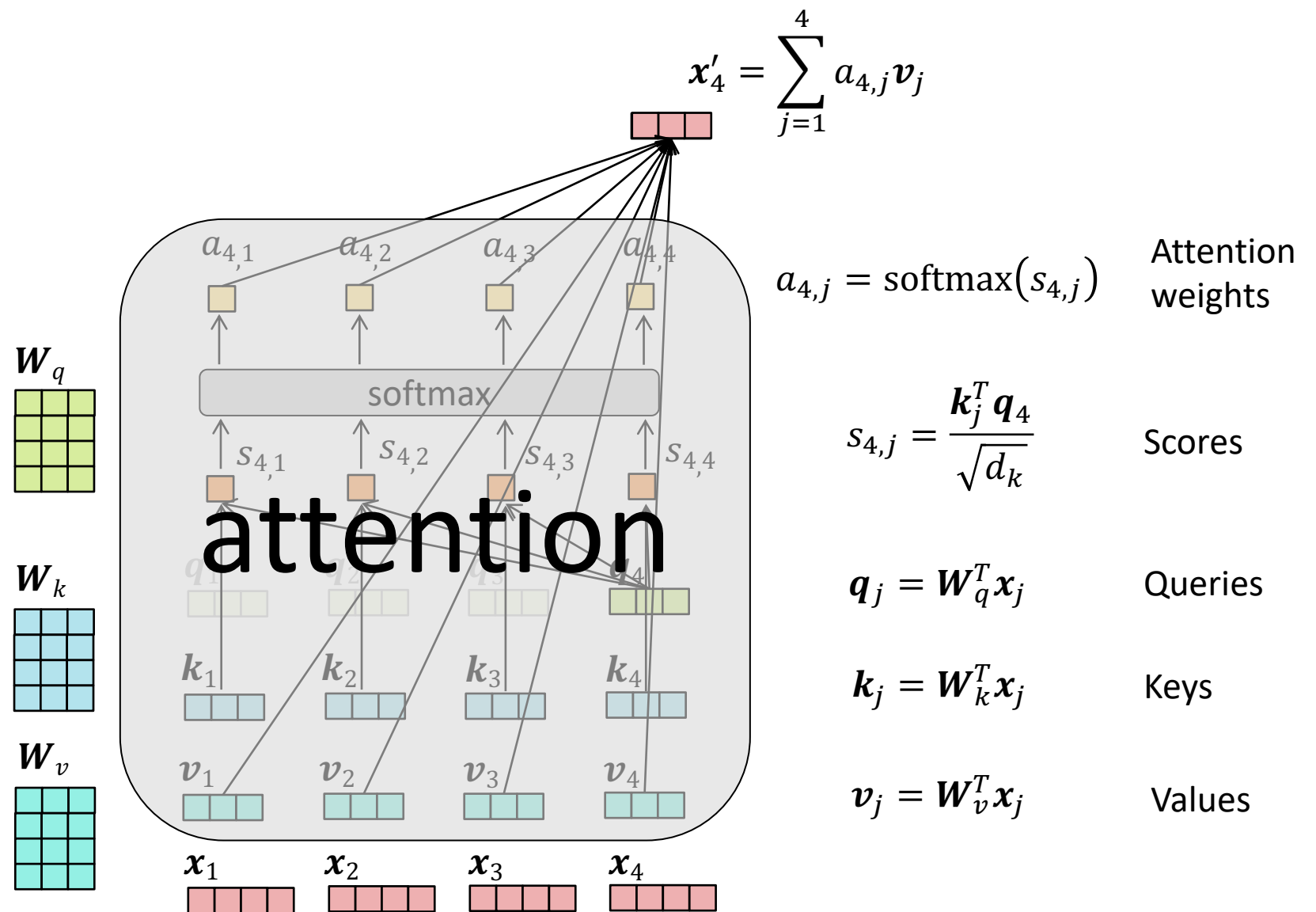
Henry Chai & Matt Gormley

11/6/23

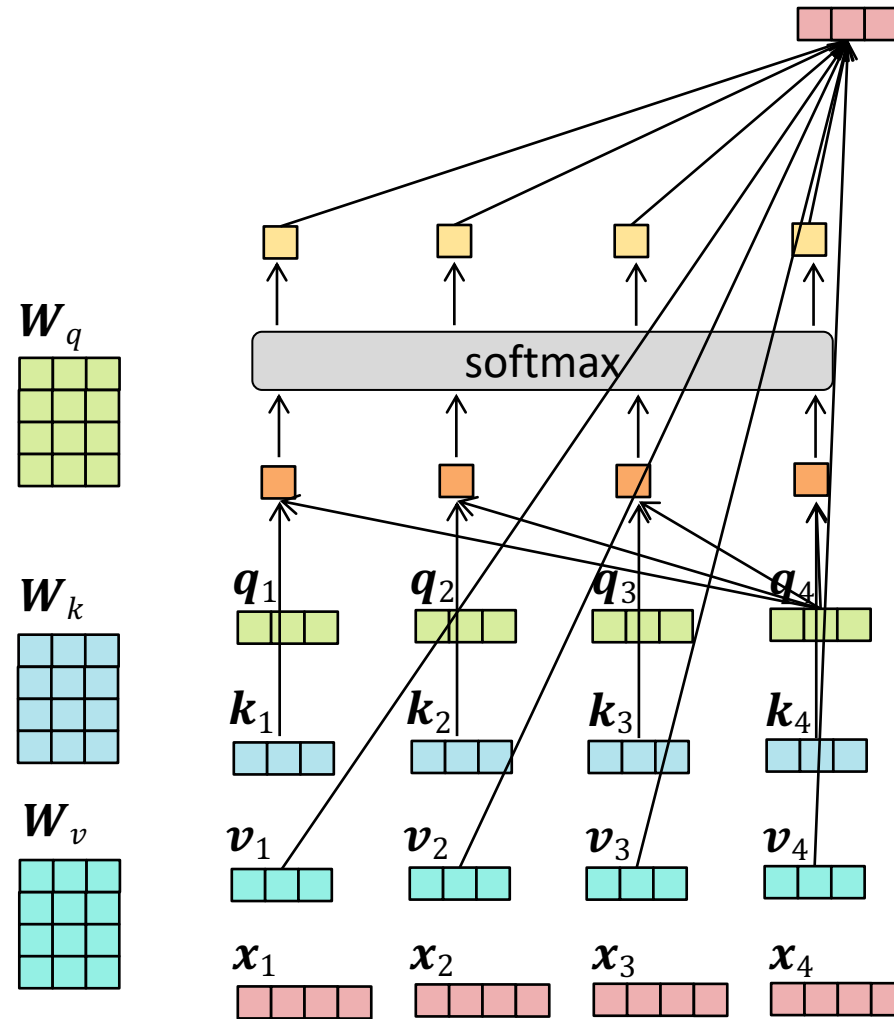
Front Matter

- Announcements:
 - Exam 2 on 11/9 (Thursday!)
 - All topics from Lecture 8 - 16 are in-scope
 - Exam 1 content may be referenced but will not be the primary focus of any question
 - No electronic devices (you won't need them!)
 - You may bring one letter-size sheet of notes; you can put *whatever* you want on *both sides*

Recall: Scaled Dot- Product Attention



Scaled Dot-Product Attention: Matrix Form

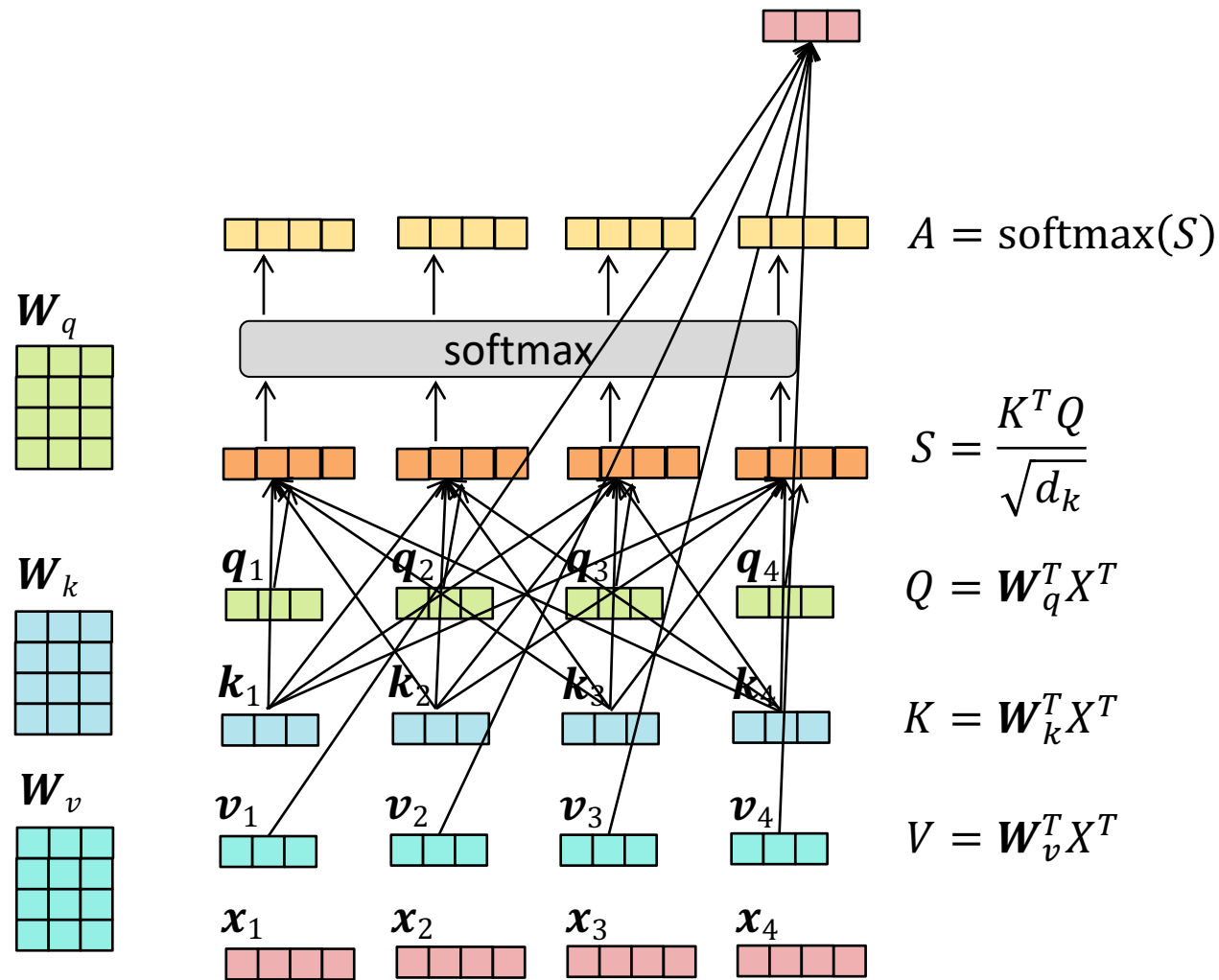


$$Q = [q_1, \dots, q_N] = W_q^T [x_1, \dots, x_N]$$

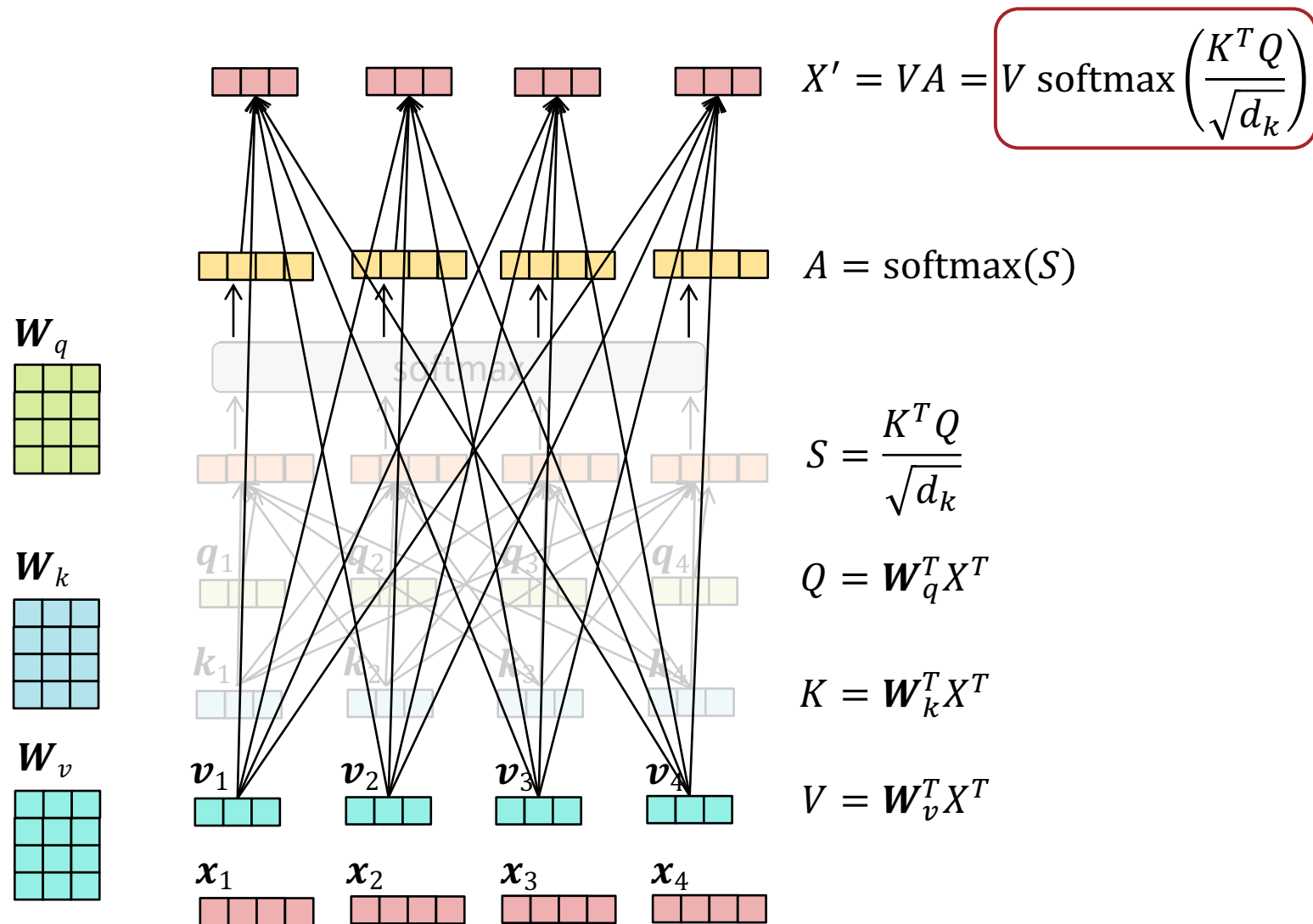
$$K = [k_1, \dots, k_N] = W_k^T [x_1, \dots, x_N]$$

$$V = [v_1, \dots, v_N] = W_v^T [x_1, \dots, x_N]$$

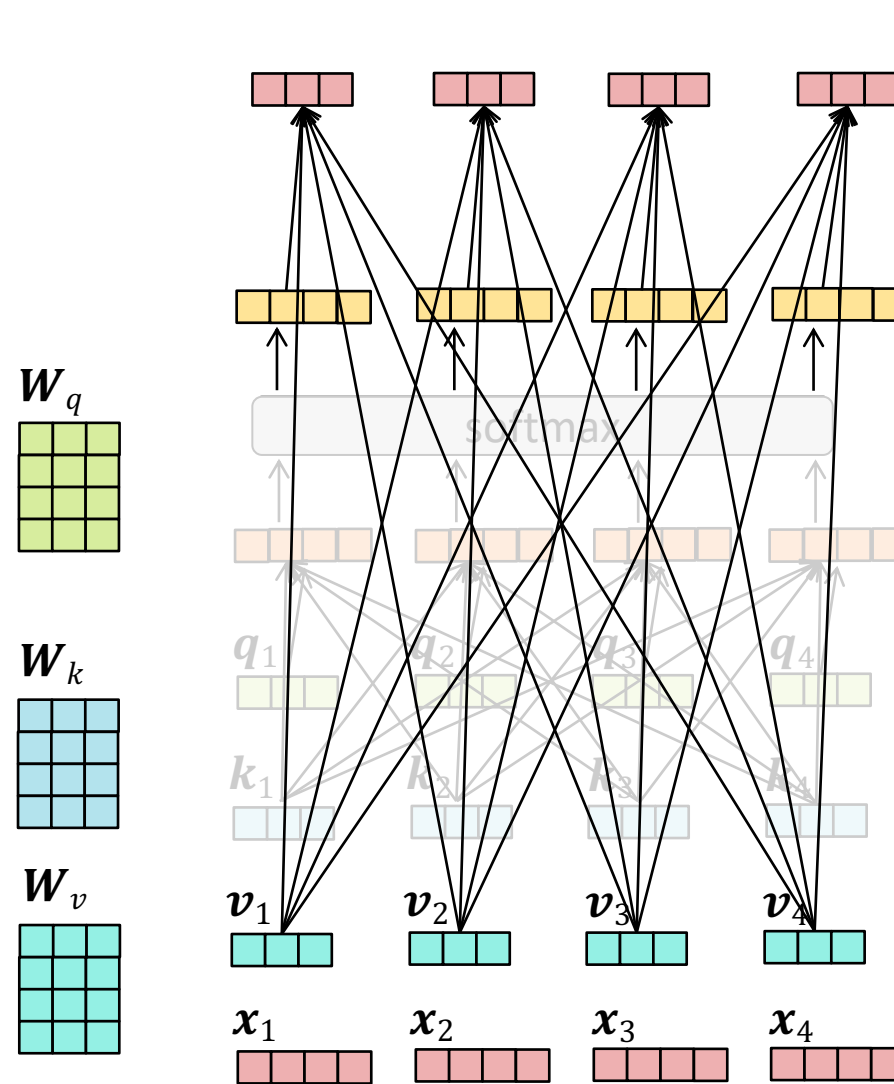
Scaled Dot-Product Attention: Matrix Form



Scaled Dot-Product Attention: Matrix Form



Decoding



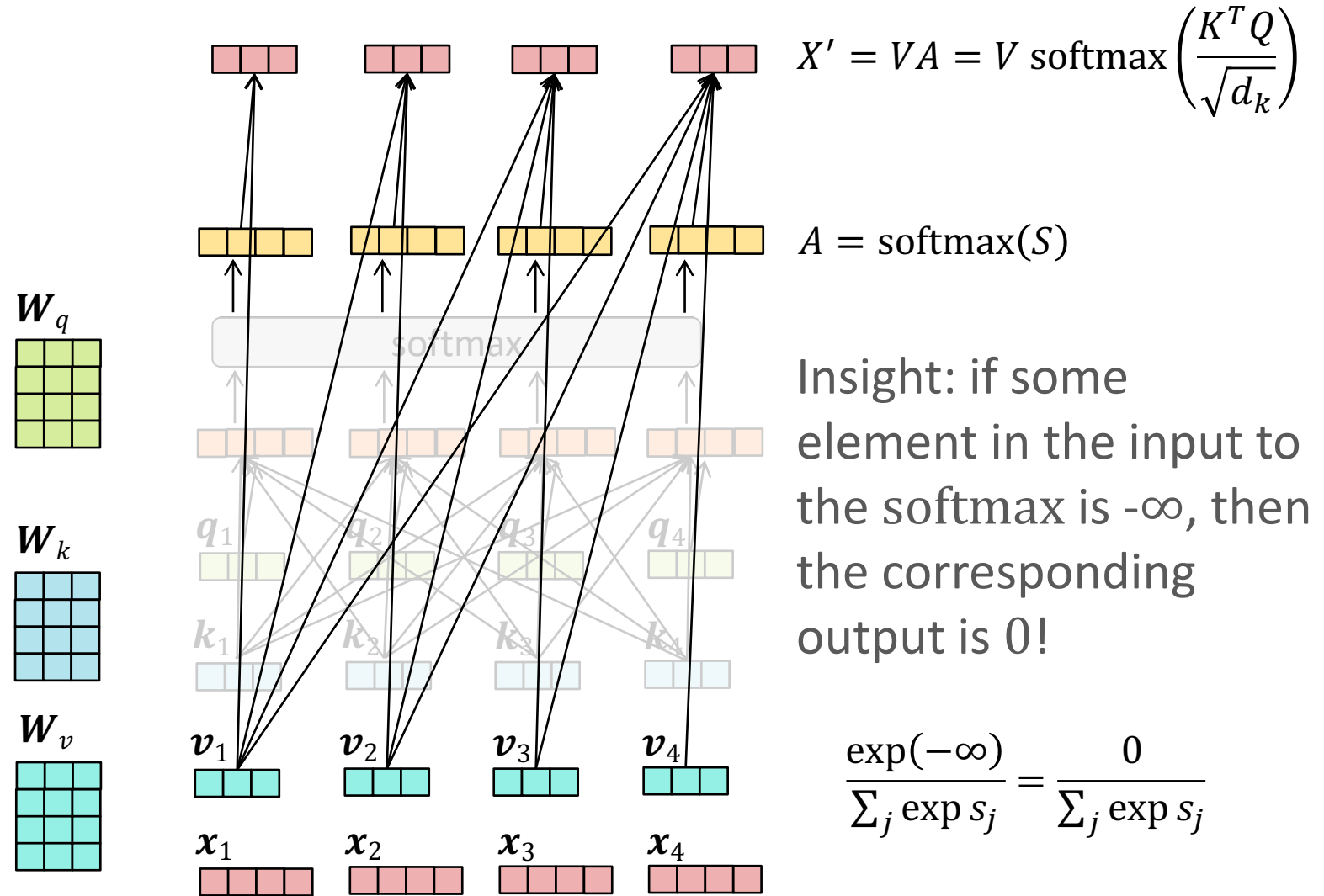
$$X' = VA = V \operatorname{softmax} \left(\frac{K^T Q}{\sqrt{d_k}} \right)$$

$$A = \operatorname{softmax}(S)$$

- Suppose we're training our transformer to predict the next token(s) given the input...
- ... then attending to tokens that come after the current token is cheating!

Masking

Idea: we can effectively delete or “mask” some of these arrows by selectively setting attention weights to 0



Insight: if some element in the input to the softmax is $-\infty$, then the corresponding output is 0!

$$\frac{\exp(-\infty)}{\sum_j \exp s_j} = \frac{0}{\sum_j \exp s_j}$$

Masked Multi-headed Attention: Matrix Form

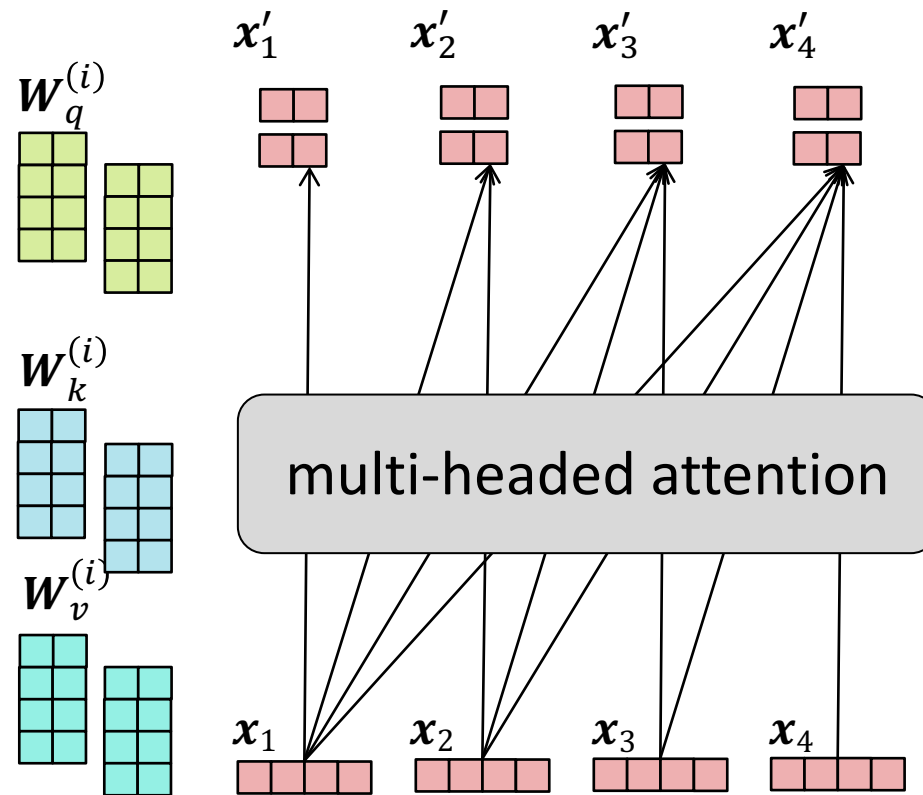
$$X' = \text{concat}_i \left\{ V^{(i)} \text{softmax} \left(\frac{K^{(i)T} Q^{(i)}}{\sqrt{d_k}} + M \right) \right\}$$

where

$$Q^{(i)} = W_q^{(i)T} X^T$$

$$K^{(i)} = W_k^{(i)T} X^T$$

$$V^{(i)} = W_v^{(i)T} X^T$$



Practical Considerations

1. Where on earth do tokens come from?
 - Example: “Henry is giving a lecture on transformers”
 - Word-based tokenization:
[“henry”, “is”, “giving”, “a”, “lecture”, “on”, “transformers”]

2. How can we handle variable-length sequences?

Practical Considerations

1. Where on earth do tokens come from?
 - Example: “Henry is givin’ a lectrue on transformers”
 - Word-based tokenization:
[“henry”, “is”, “?”, “a”, “?”, “on”, “transformers”]
 - Can have difficulty trading off between vocabulary size and computational tractability
 - Similar words e.g., “transformers” and “transformer” can get mapped to completely disparate representations
 - Typos will typically be out-of-vocabulary (OOV)
2. How can we handle variable-length sequences?

Practical Considerations

1. Where on earth do tokens come from?
 - Example: “Henry is givin’ a lectrue on transformers”
 - Character-based tokenization:
[“h”, “e”, “n”, “r”, “y”, “i”, “s”, “g”, “i”, “v”, “i”, “n”, “ ’ ”, ...]
 - Much smaller vocabularies but a lot of semantic meaning is lost...
 - Sequences will be much longer than word-based tokenization, potentially causing computational issues
 - Can do well on logographic languages e.g., Kanji 漢字
2. How can we handle variable-length sequences?

Practical Considerations

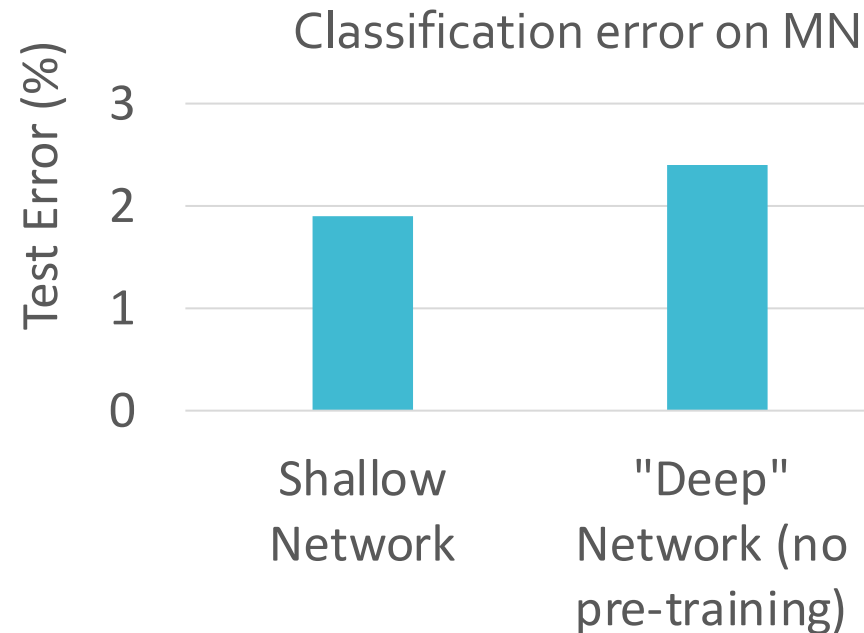
1. Where on earth do tokens come from?
 - Example: “Henry is givin’ a lectrue on transformers”
 - Subword tokenization:
[“henry”, “is”, “giv”, “##in”, “ ’”, “a”, “lect” “##re”, “on”, “transform”, “##ers”]
 - Split long or rare words into smaller, semantically meaningful components or *subwords*
2. How can we handle variable-length sequences?
 - Artificially make all sequences the same length by
 - Padding: adding special *pad tokens* to short sequences
 - Truncating: using only the first few tokens for long sequences

Recall: Mini-batch Stochastic Gradient Descent...

- Input: training dataset $\mathcal{D} = \{(\mathbf{x}^{(i)}, y^{(i)})\}_{i=1}^N$,
step size γ , and batch size B
 1. Randomly initialize the parameters $\boldsymbol{\theta}^{(0)}$ and set $t = 0$
 2. While TERMINATION CRITERION is not satisfied
 - a. Randomly sample B data points from \mathcal{D} , $\{(\mathbf{x}^{(b)}, y^{(b)})\}_{b=1}^B$
 - b. Compute the gradient of the loss w.r.t. the sampled *batch*,
$$\nabla J^{(B)}(\boldsymbol{\theta}^{(t)})$$
 - c. Update $\boldsymbol{\theta}$: $\boldsymbol{\theta}^{(t+1)} \leftarrow \boldsymbol{\theta}^{(t)} - \gamma \nabla J^{(B)}(\boldsymbol{\theta}^{(t)})$
 - d. Increment t : $t \leftarrow t + 1$
- Output: $\boldsymbol{\theta}^{(t)}$

Reality

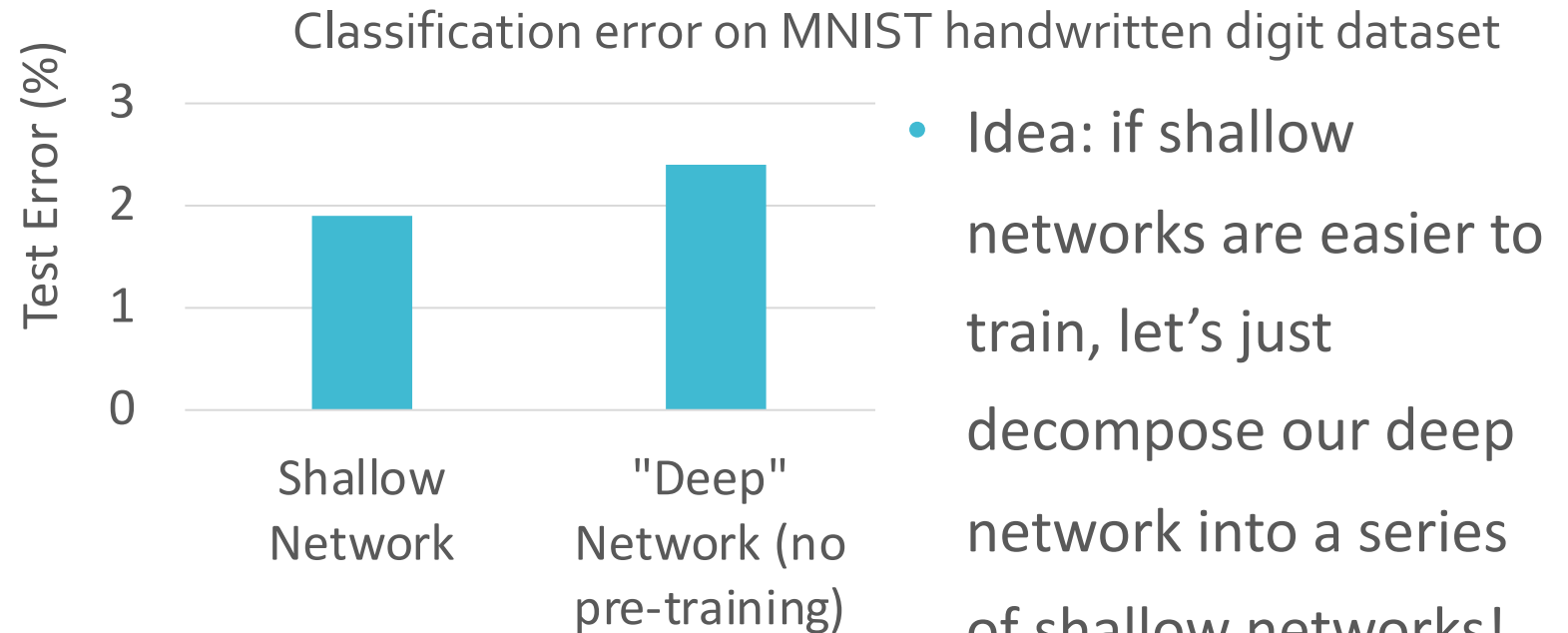
- You have some niche task that you want to apply machine learning to e.g., predicting the author of children's books
- You have a tiny labelled dataset to train with
- You fit a massive deep learning model to the dataset
- Surprise, surprise: it overfits and your test error is super high



- “gradient-based optimization starting from random initialization appears to often get stuck in poor solutions for such deep networks.”

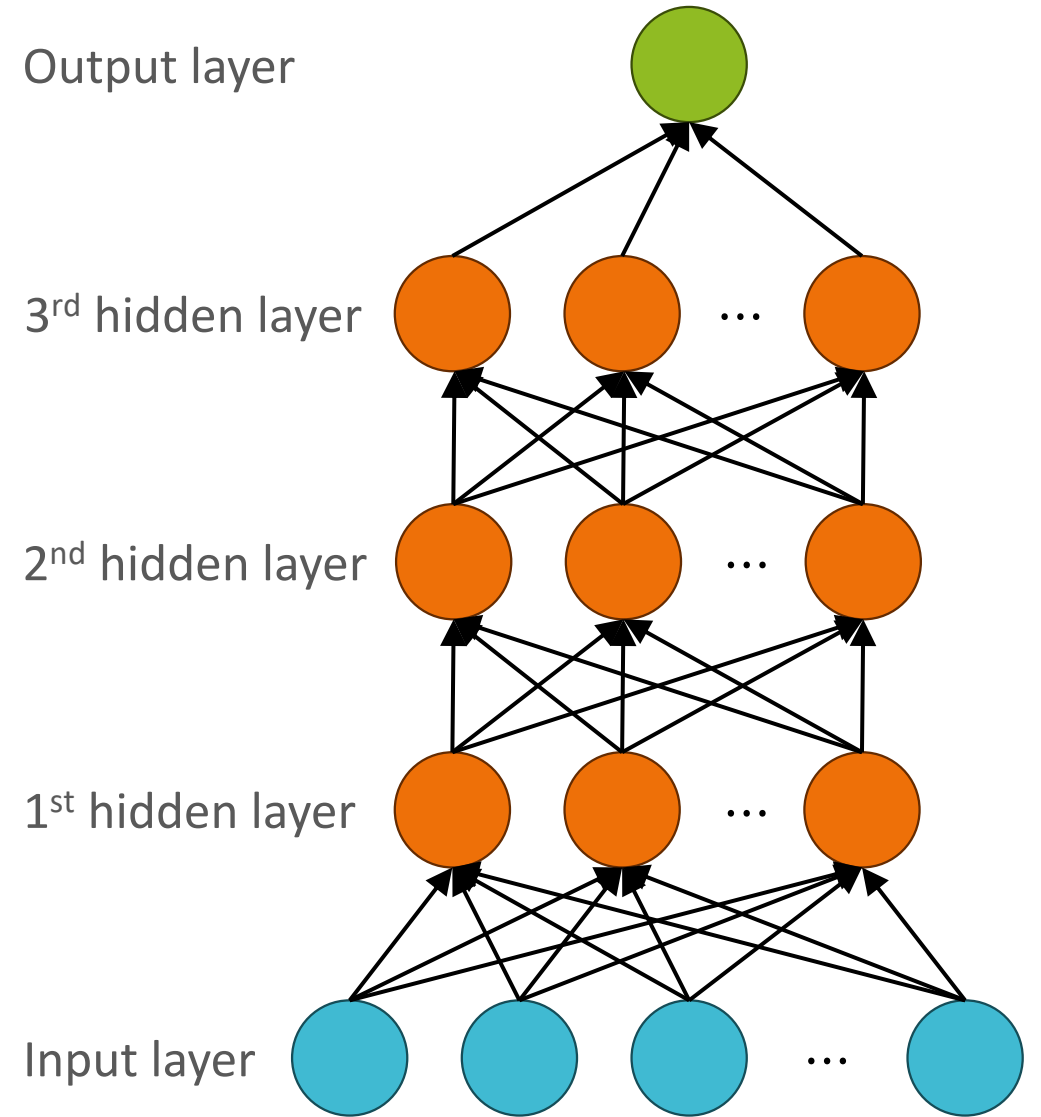
Reality

- You have some niche task that you want to apply machine learning to e.g., predicting the author of children's books
- You have a tiny labelled dataset to train with
- You fit a massive deep learning model to the dataset
- Surprise, surprise: it overfits and your test error is super high



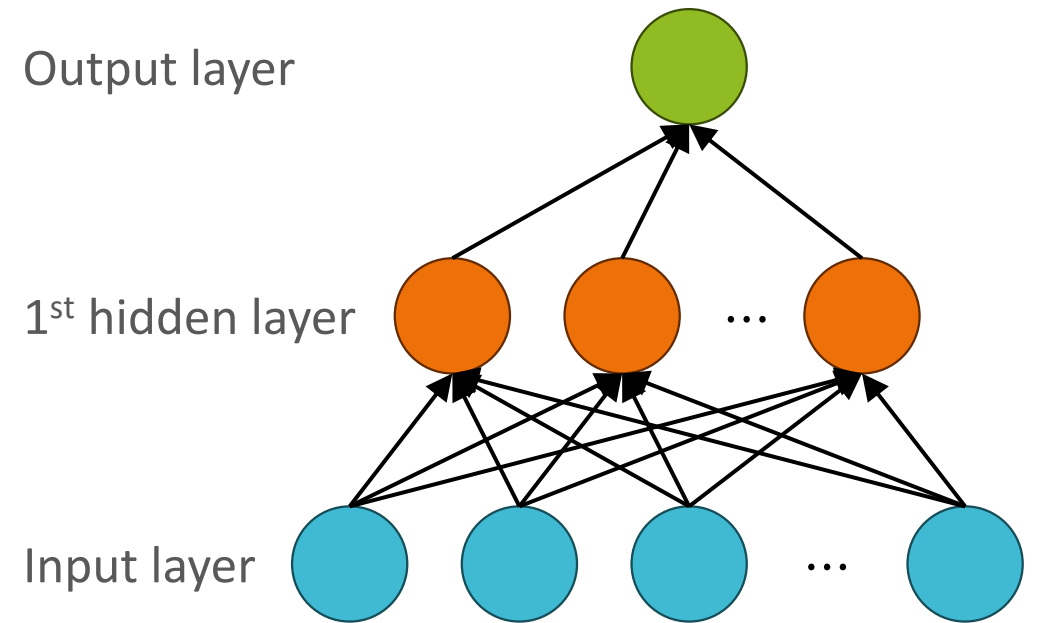
Pre-training (Bengio et al., 2006)

- Train each layer of the network iteratively using the training dataset
- Start at the input layer and move towards the output layer
- Once a layer has been trained, fix its weights and use those to train subsequent layers



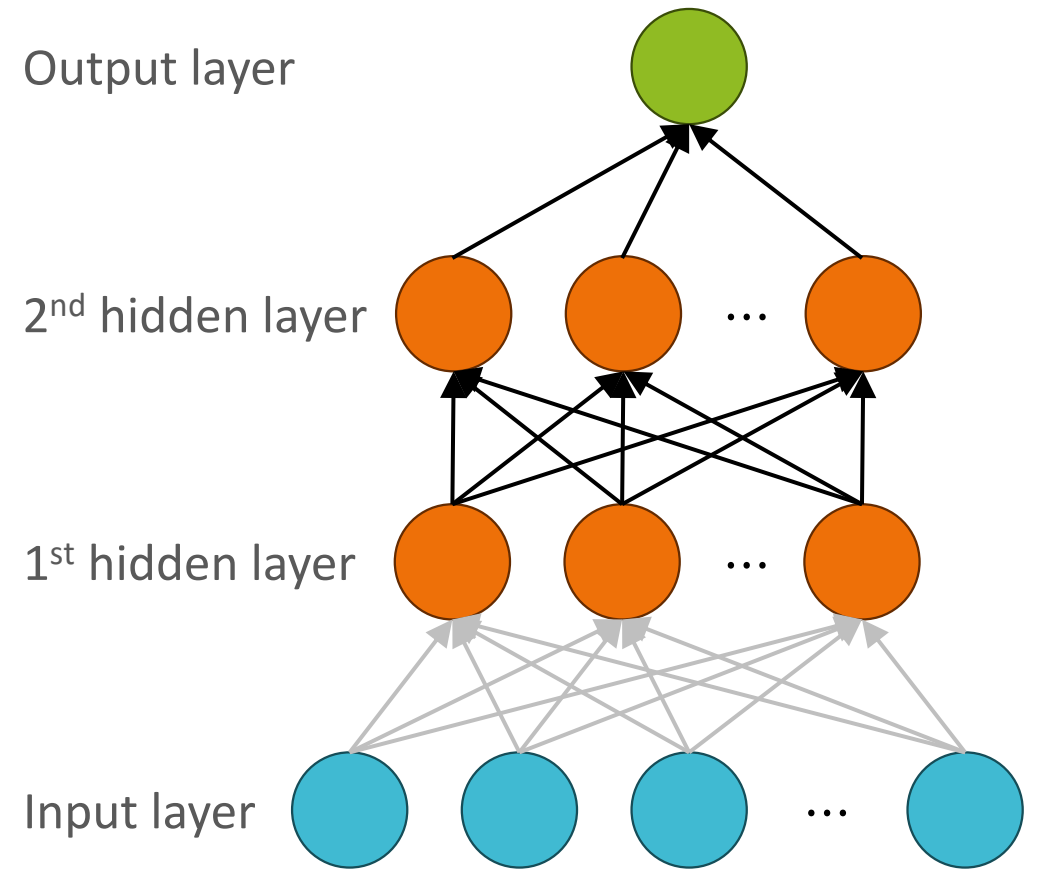
Pre-training (Bengio et al., 2006)

- Train each layer of the network iteratively using the training dataset
- Start at the input layer and move towards the output layer
- Once a layer has been trained, fix its weights and use those to train subsequent layers



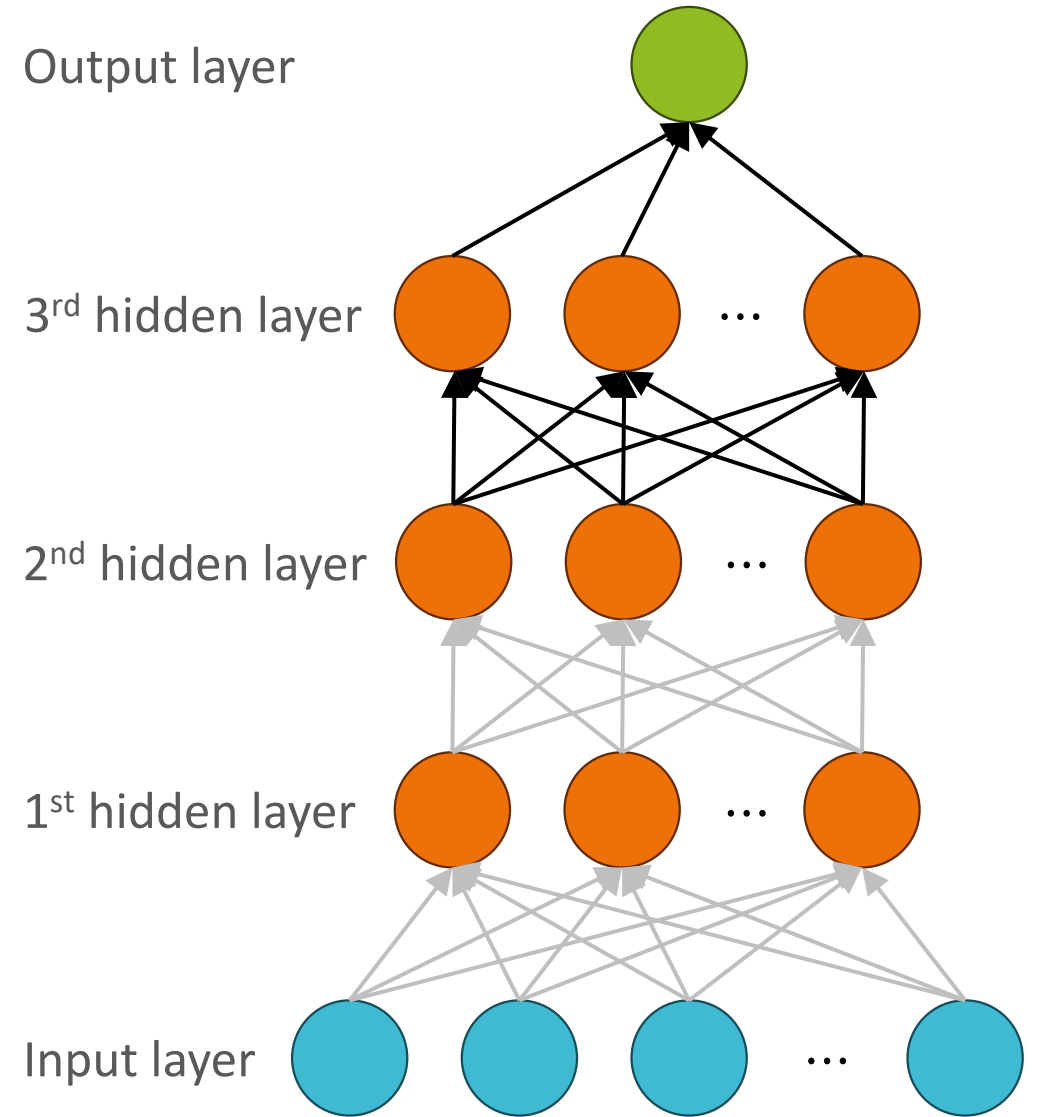
Pre-training (Bengio et al., 2006)

- Train each layer of the network iteratively using the training dataset
- Start at the input layer and move towards the output layer
- Once a layer has been trained, fix its weights and use those to train subsequent layers



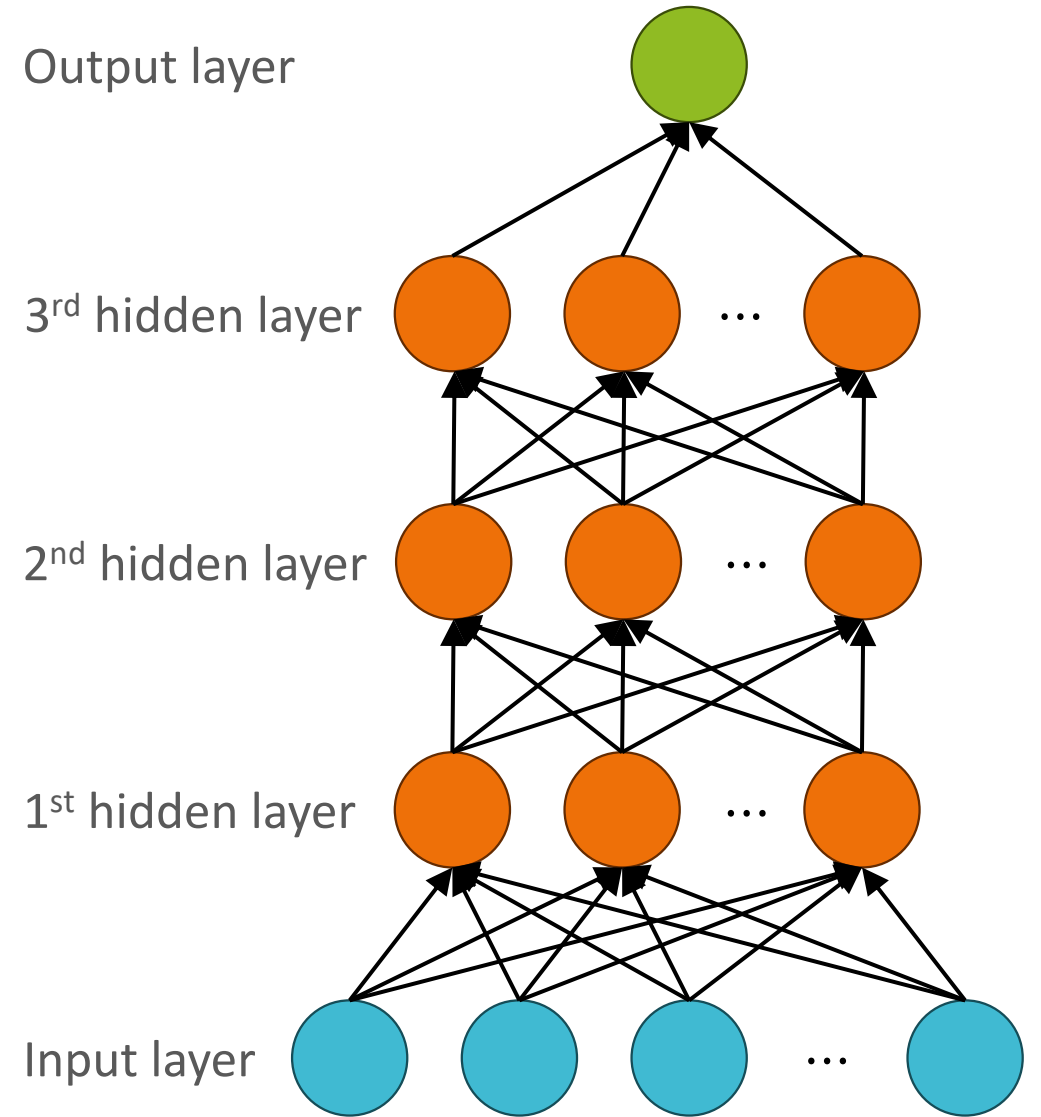
Pre-training (Bengio et al., 2006)

- Train each layer of the network iteratively using the training dataset
- Start at the input layer and move towards the output layer
- Once a layer has been trained, fix its weights and use those to train subsequent layers



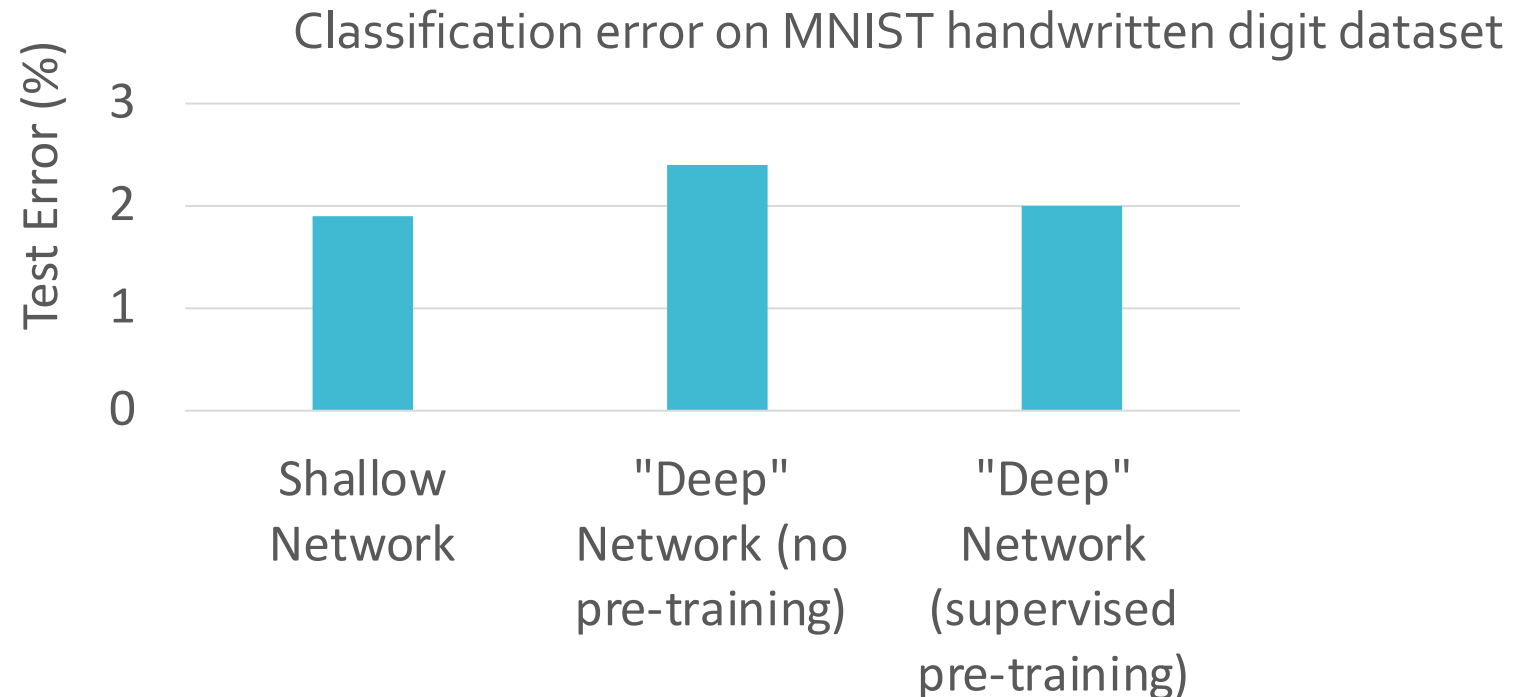
Fine-tuning (Bengio et al., 2006)

- Train each layer of the network iteratively using the training dataset
- Use the pre-trained weights as an initialization and *fine-tune* the entire network e.g., via SGD with the training dataset



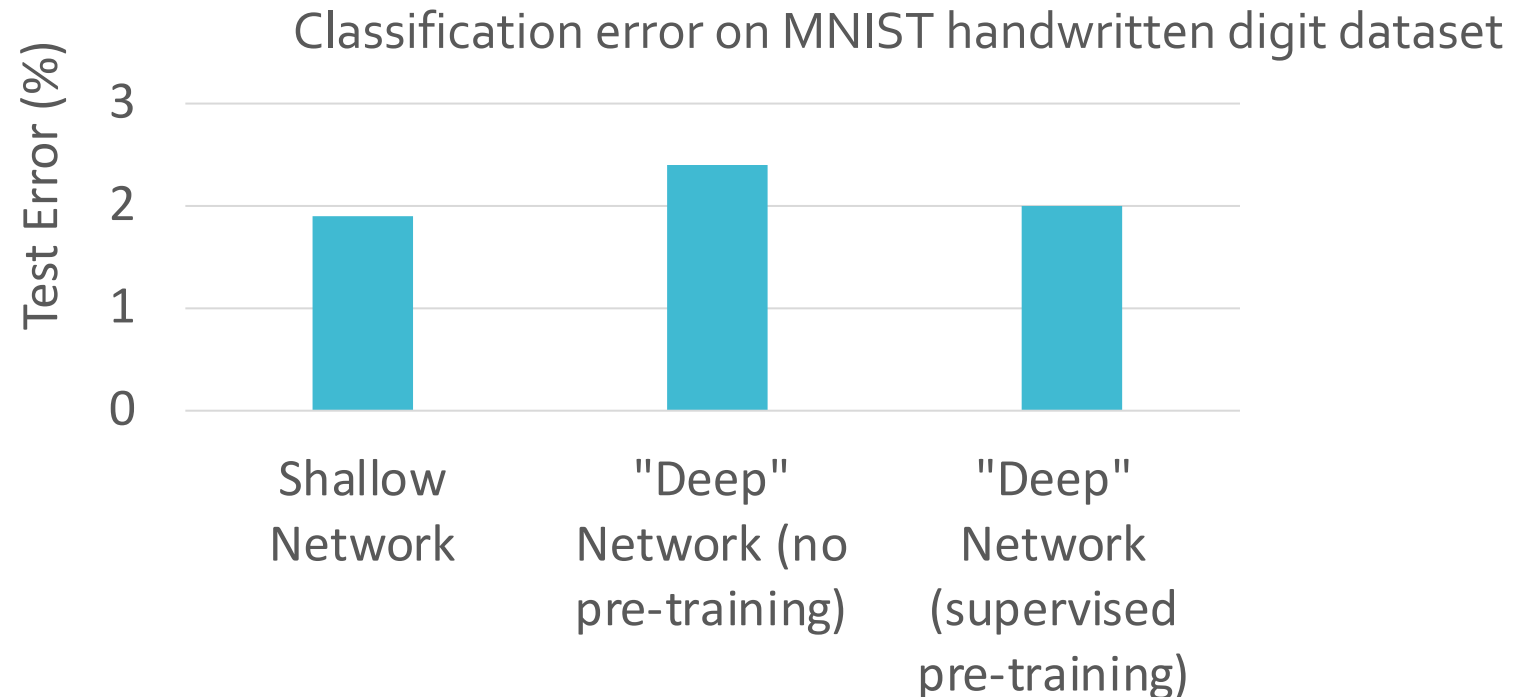
Supervised Pre-training (Bengio et al., 2006)

- Train each layer of the network iteratively using the training dataset
- Use the pre-trained weights as an initialization and *fine-tune* the entire network e.g., via SGD with the training dataset



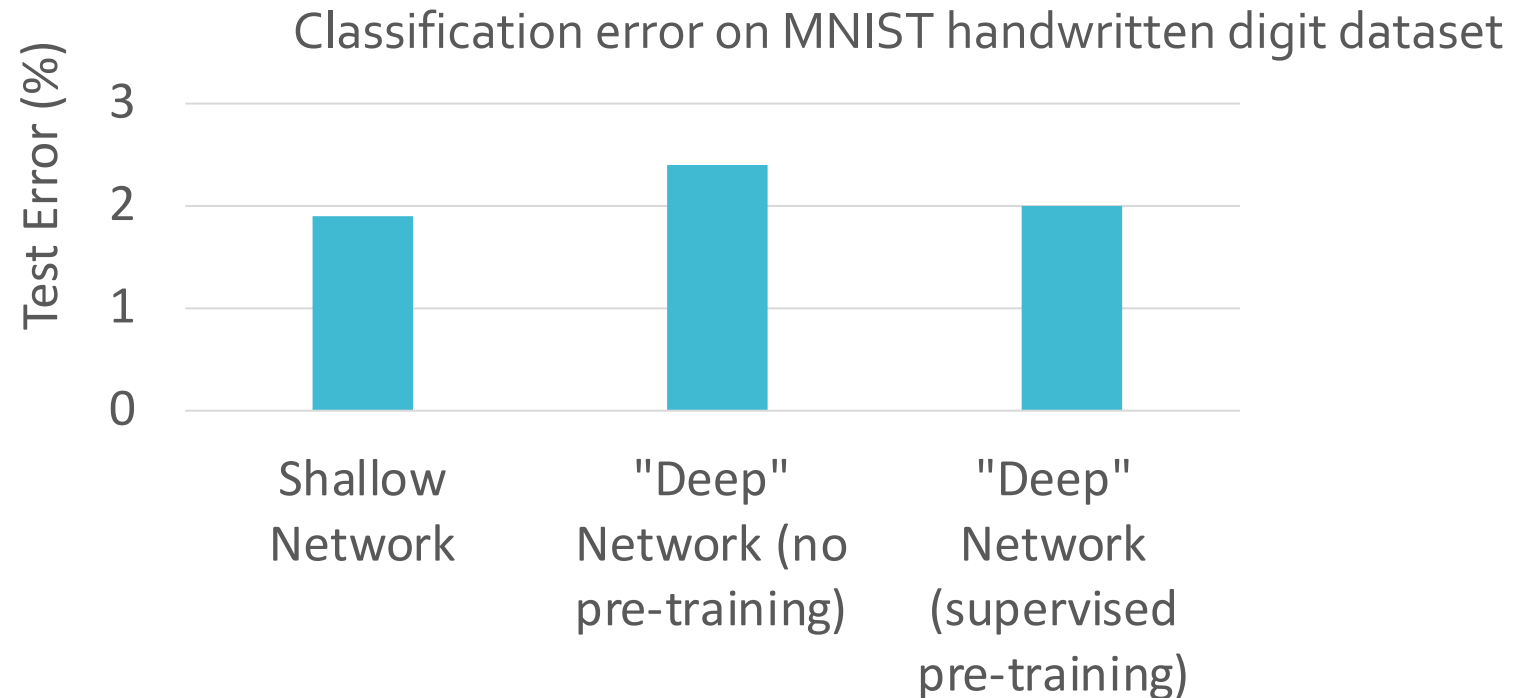
Supervised Pre-training (Bengio et al., 2006)

- Train each layer of the network iteratively using the training dataset *to predict the labels*
- Use the pre-trained weights as an initialization and *fine-tune* the entire network e.g., via SGD with the training dataset



Unsupervised Pre-training (Bengio et al., 2006)

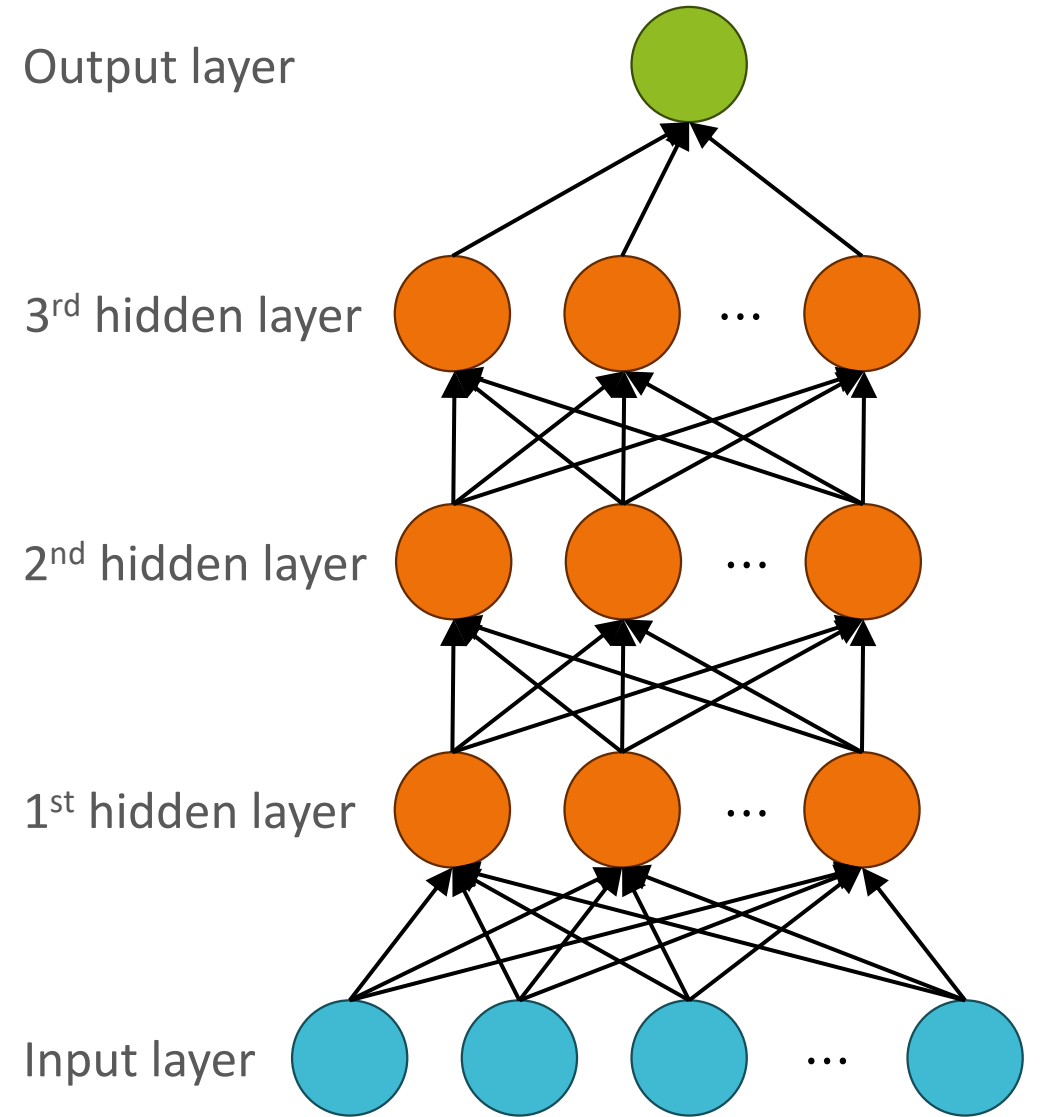
- Train each layer of the network iteratively using the training dataset to learn *useful* representations
- Idea: a good representation is one preserves a lot of information and could be used to recreate the inputs



Unsupervised Pre-training (Bengio et al., 2006)

- Train each layer of the network iteratively using the training dataset by minimizing the *reconstruction error*

$$\|\mathbf{x} - h(\mathbf{x})\|_2$$

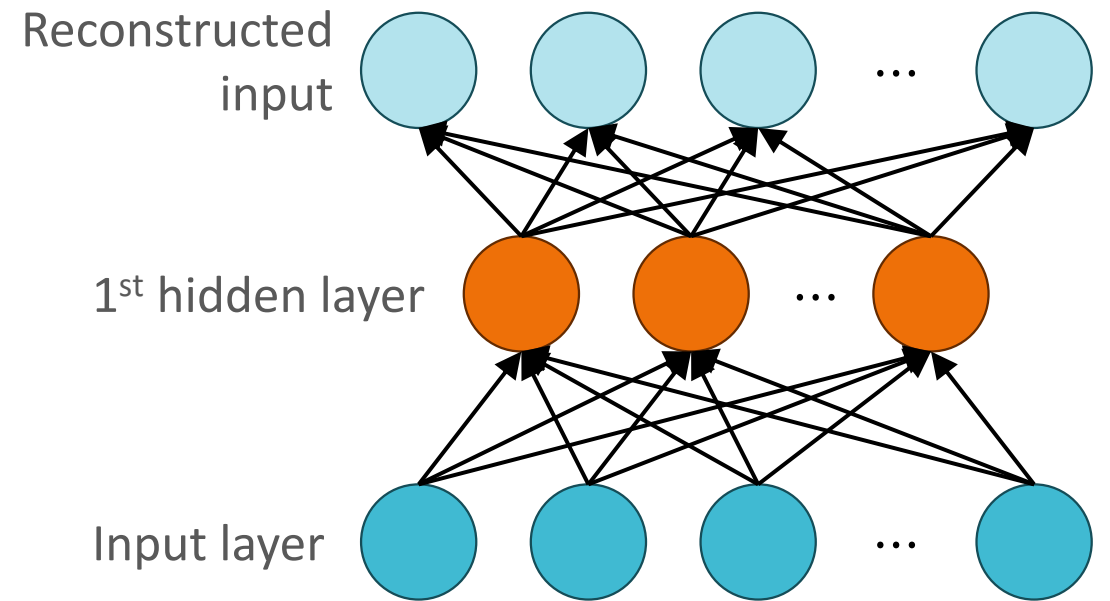


Unsupervised Pre-training (Bengio et al., 2006)

- Train each layer of the network iteratively using the training dataset by minimizing the *reconstruction error*

$$\|x - h(x)\|_2$$

- This architecture/objective defines an *autoencoder*

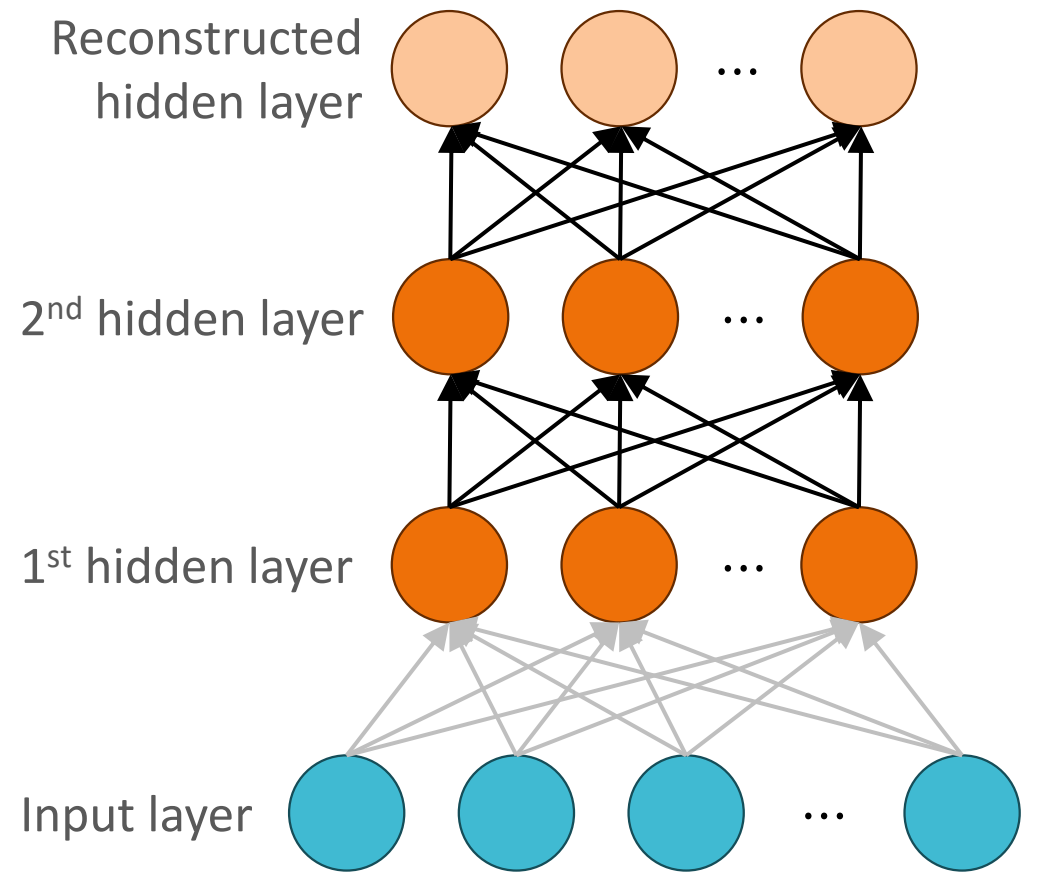


Unsupervised Pre-training (Bengio et al., 2006)

- Train each layer of the network iteratively using the training dataset by minimizing the *reconstruction error*

$$\|\mathbf{x} - h(\mathbf{x})\|_2$$

- This architecture/objective defines an *autoencoder*

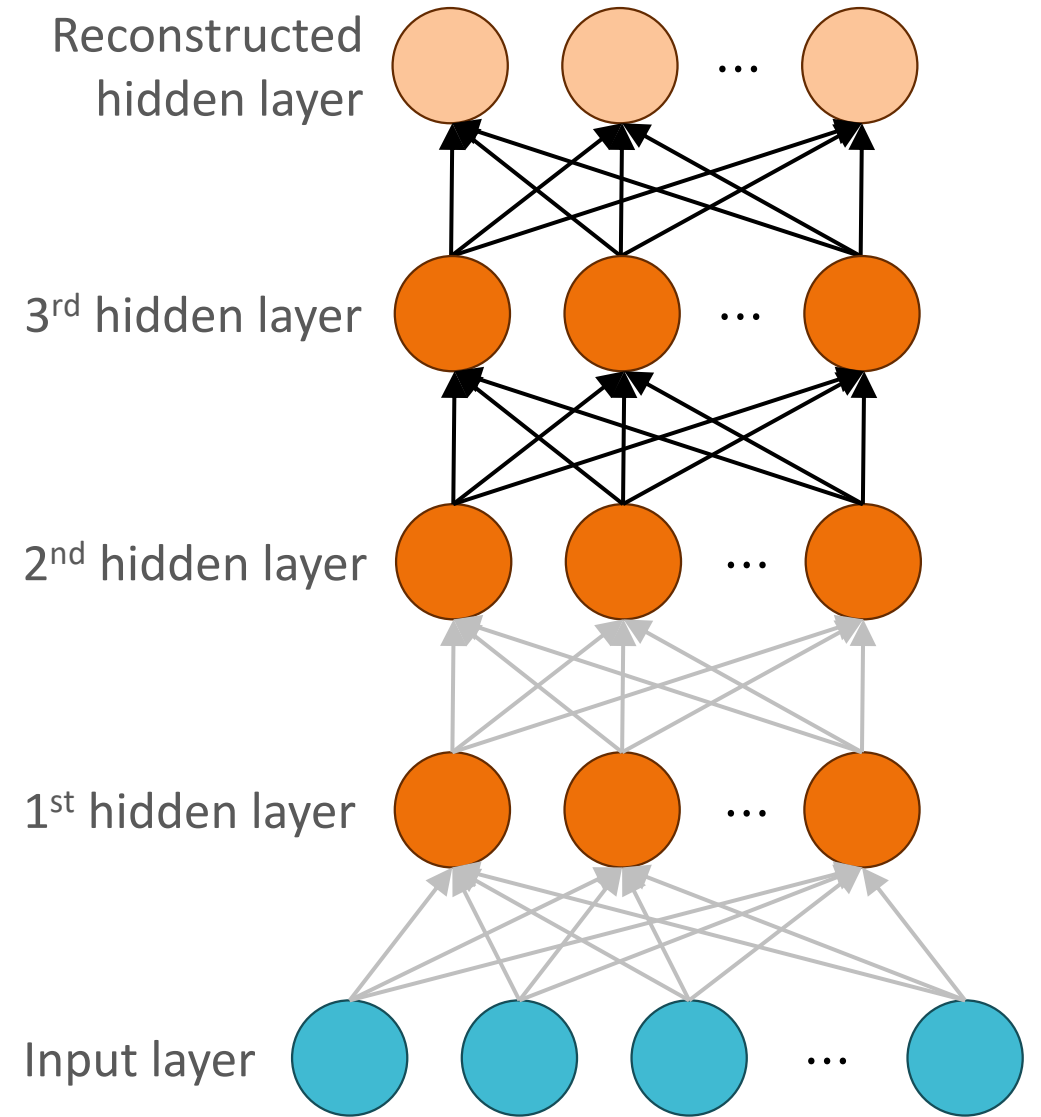


Unsupervised Pre-training (Bengio et al., 2006)

- Train each layer of the network iteratively using the training dataset by minimizing the *reconstruction error*

$$\|x - h(x)\|_2$$

- This architecture/objective defines an *autoencoder*

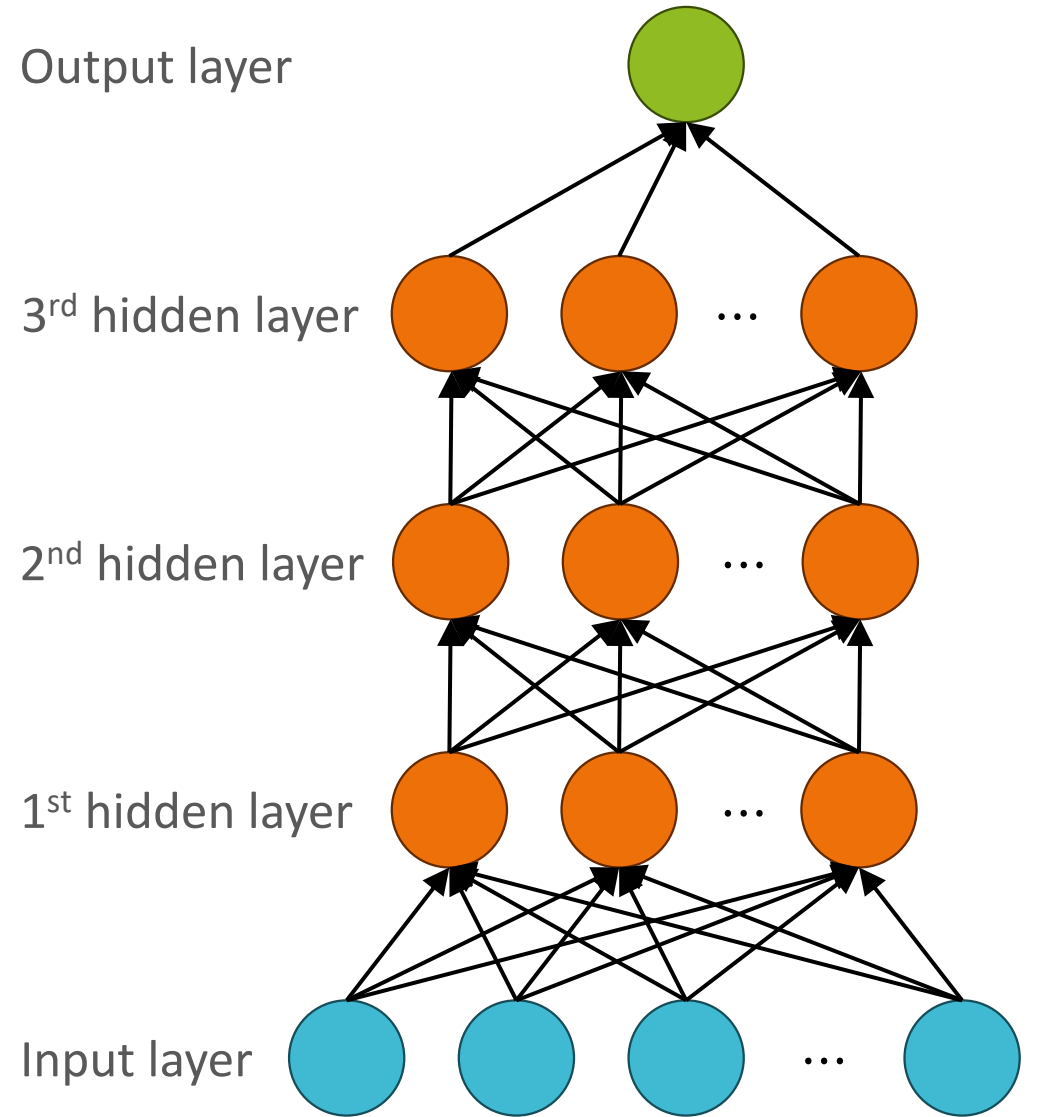


Fine-tuning (Bengio et al., 2006)

- Train each layer of the network iteratively using the training dataset by minimizing the *reconstruction error*

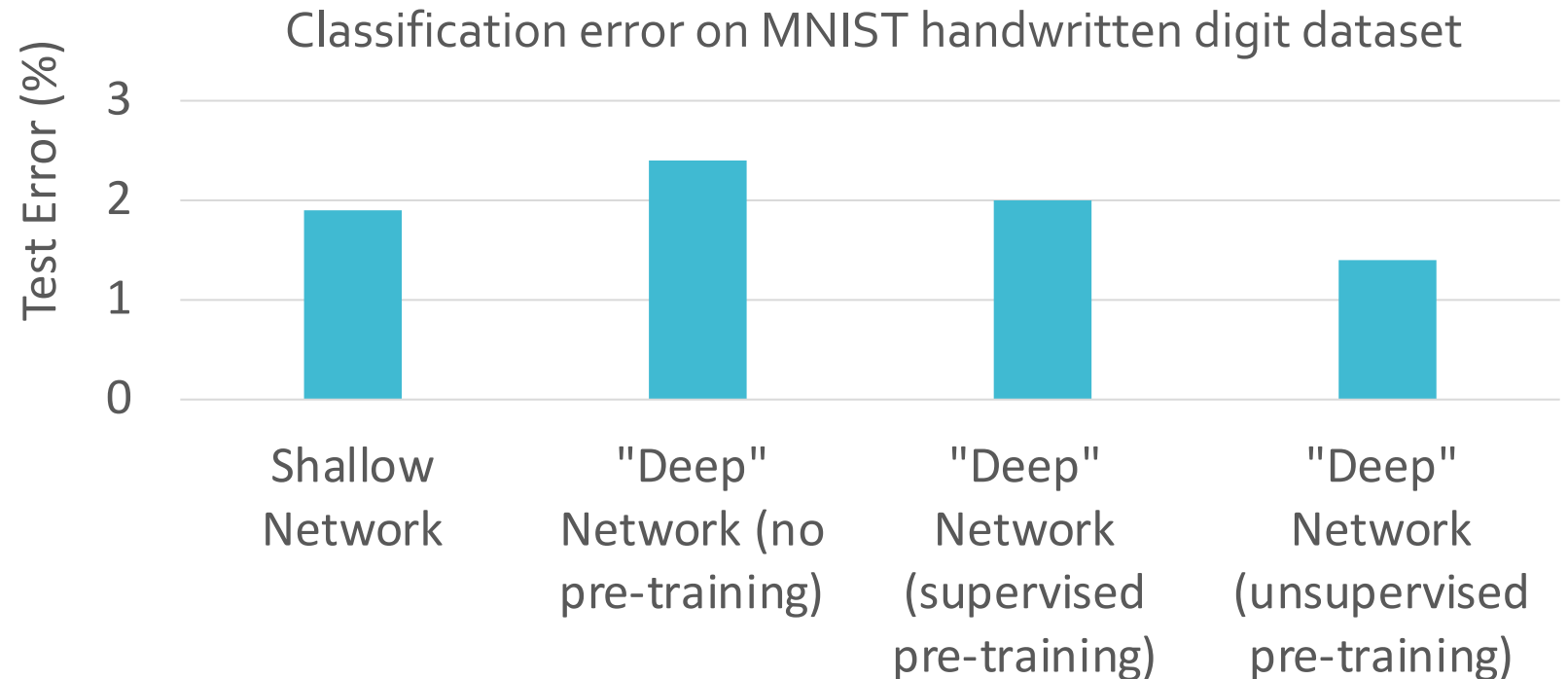
$$\|\mathbf{x} - h(\mathbf{x})\|_2$$

- When fine-tuning, we're effectively swapping out the last layer and fitting all the weights to the training dataset



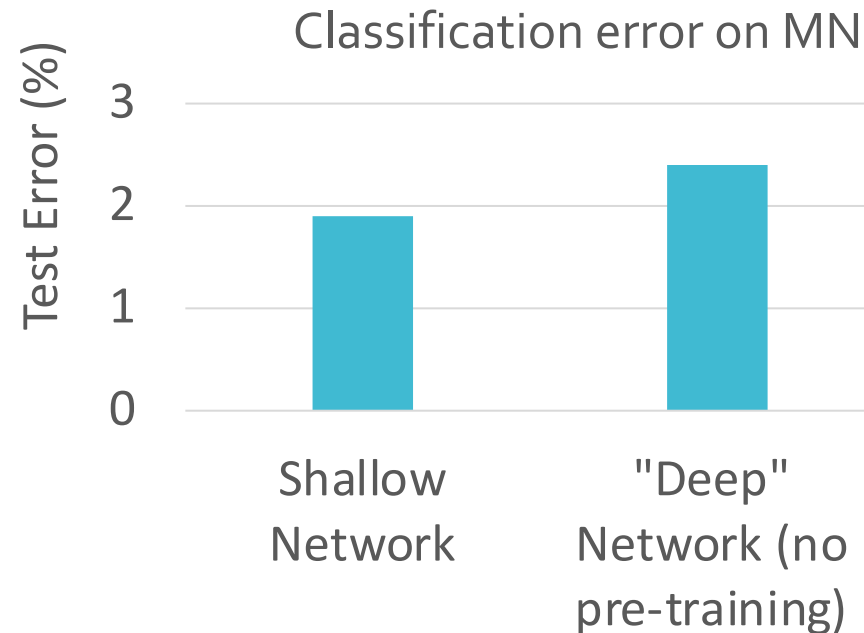
Unsupervised Pre-training (Bengio et al., 2006)

- Train each layer of the network iteratively using the training dataset by minimizing the *reconstruction error*
- Idea: a good representation is one preserves a lot of information and could be used to recreate the inputs



Another dose of Reality

- You have some niche task that you want to apply machine learning to e.g., predicting the author of children's books
- You have a **tiny** labelled dataset to train with
- You fit a **massive** deep learning model to the dataset
- Surprise, surprise: it overfits and your test error is super high



- Problem: what if you don't even have enough data to train a single layer/fine-tune the pre-trained network?

Another dose of Reality

- You have some niche task that you want to apply machine learning to e.g., predicting the author of children's books
- You have a **tiny** labelled dataset to train with
- You fit a **massive** deep learning model to the dataset
- Surprise, surprise: it overfits and your test error is super high
- Key observation: you can pre-train on basically any labelled or unlabelled dataset!
 - Ideally, you want to use a *large* dataset *related* to your goal task

Another dose of Reality

- You have some niche task that you want to apply machine learning to e.g., predicting the author of children's books
- You have a **tiny** labelled dataset to train with
- You fit a **massive** deep learning model to the dataset
- Surprise, surprise: it overfits and your test error is super high
- Key observation: you can pre-train on basically any labelled or unlabelled dataset!
 - GPT-3 pre-training data:

Dataset	Quantity (tokens)	Weight in training mix
Common Crawl (filtered)	410 billion	60%
WebText2	19 billion	22%
Books1	12 billion	8%
Books2	55 billion	8%
Wikipedia	3 billion	3%

Another dose of Reality

- You have some niche task that you want to apply machine learning to e.g., predicting the author of children's books
- You have a **tiny** labelled dataset to train with
- You fit a **massive** deep learning model to the dataset
- Surprise, surprise: it overfits and your test error is super high
- Key observation: you can pre-train on basically any labelled or unlabelled dataset!
- Okay that's great for pre-training and all, but what if
 - A. the concept of labelled data doesn't apply to your task i.e., not every input has a "correct" label e.g., chatbots?
 - B. you don't have enough data to fine-tune your model?

Reinforcement Learning from Human Feedback (RLHF)

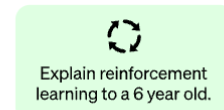
- Insight: for many machine learning tasks, there is no universal ground truth, e.g., there are lots of possible ways to respond to a question or prompt.
- Idea: use human feedback to determine how good or bad some prediction/response is!
- Issue: if the input space is huge (e.g., all possible chat prompts), to train a good model, we might need tons and tons of (potentially expensive) human annotation...
- Idea: use a small number of annotations to learn a “reward” function!

Reinforcement Learning from Human Feedback (RLHF)

Step 1

Collect demonstration data and train a supervised policy.

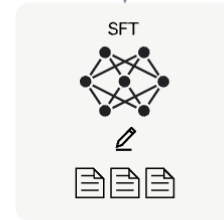
A prompt is sampled from our prompt dataset.



A labeler demonstrates the desired output behavior.



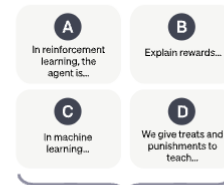
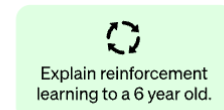
This data is used to fine-tune GPT-3.5 with supervised learning.



Step 2

Collect comparison data and train a reward model.

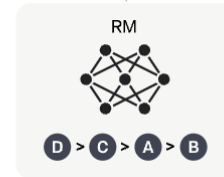
A prompt and several model outputs are sampled.



A labeler ranks the outputs from best to worst.



This data is used to train our reward model.



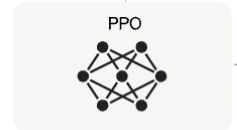
Step 3

Optimize a policy against the reward model using the PPO reinforcement learning algorithm.

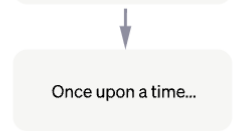
A new prompt is sampled from the dataset.



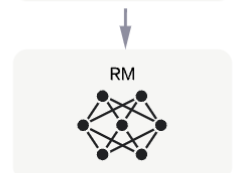
The PPO model is initialized from the supervised policy.



The policy generates an output.



The reward model calculates a reward for the output.



The reward is used to update the policy using PPO.



- RLHF is a special form of fine-tuning, used to fine-tune GPT-3.5 into ChatGPT

In-context Learning

- Problem: given their size, effectively fine-tuning LLMs can require lots of labelled data points.
- Idea: leverage the LLM's context window by passing a few examples to the model as input, *without performing any updates to the parameters*
- Intuition: during training, the LLM is exposed to a *massive* number of examples/tasks and the input conditions the model to “locate” the relevant concepts

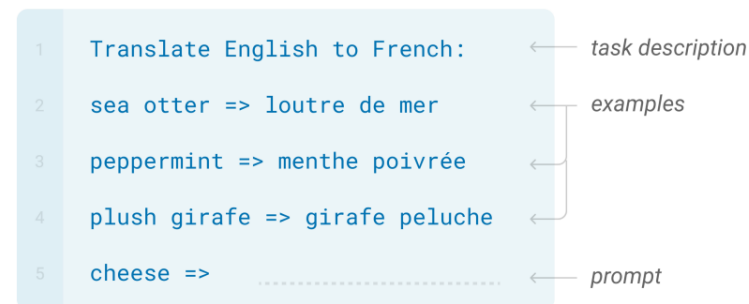
Few-shot, One-shot & Zero-shot (in-context) Learning

- Idea: leverage the LLM's context window by passing a few examples to the model as input, *without performing any updates to the parameters*

The three settings we explore for in-context learning

Few-shot

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



Traditional fine-tuning (not used for GPT-3)

Fine-tuning

The model is trained via repeated gradient updates using a large corpus of example tasks.



Few-shot, One-shot & Zero-shot (in-context) Learning

- Idea: leverage the LLM's context window by passing a ~~few~~ one or more examples to the model as input, *without performing any updates to the parameters*

The three settings we explore for in-context learning

One-shot

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

```
1 Translate English to French: ← task description
2 sea otter => loutre de mer ← example
3 cheese => ..... ← prompt
```

Traditional fine-tuning (not used for GPT-3)

Fine-tuning

The model is trained via repeated gradient updates using a large corpus of example tasks.



Few-shot, One-shot & Zero-shot (in-context) Learning

- Idea: leverage the LLM's context window by passing a ~~few one~~ zero(!) examples to the model as input, *without performing any updates to the parameters*

The three settings we explore for in-context learning

Zero-shot

The model predicts the answer given only a natural language description of the task. No gradient updates are performed.



Traditional fine-tuning (not used for GPT-3)

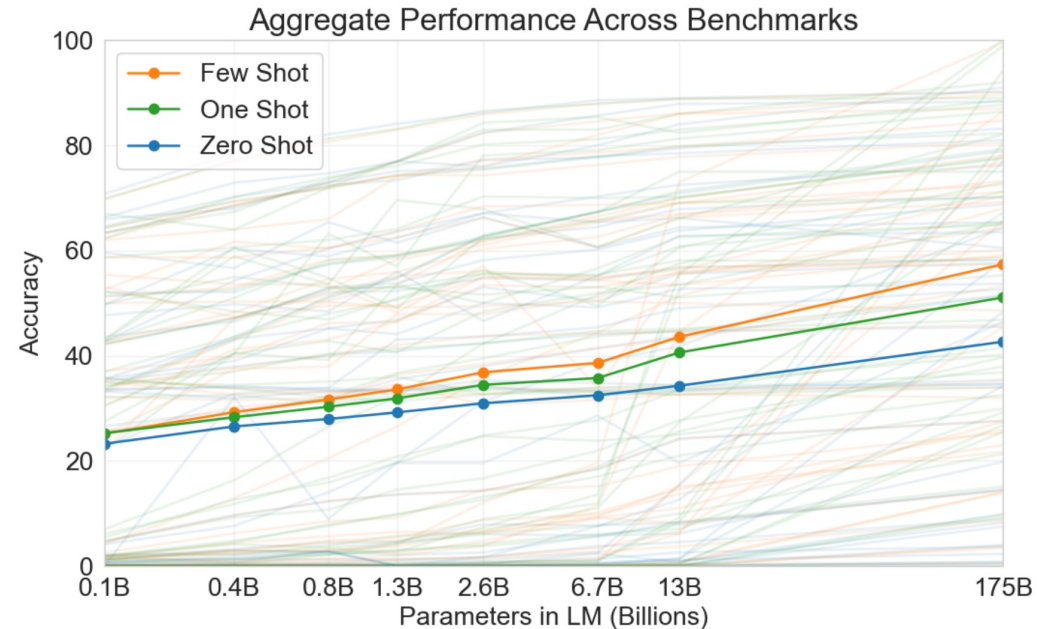
Fine-tuning

The model is trained via repeated gradient updates using a large corpus of example tasks.



Few-shot, One-shot & Zero-shot (in-context) Learning

- Idea: leverage the LLM's context window by passing a ~~few one~~ zero(!) examples to the model as input, *without performing any updates to the parameters*



- Key Takeaway: LLMs can perform well on novel tasks without having to fine-tune the model, sometimes even with just one or zero labelled training data points!

Mini-batch Stochastic Gradient Descent is a lie! just the beginning!

- Input: training dataset $\mathcal{D} = \{(\mathbf{x}^{(i)}, y^{(i)})\}_{i=1}^N$,
step size γ , and batch size B
 1. Randomly initialize the parameters $\boldsymbol{\theta}^{(0)}$ and set $t = 0$
 2. While TERMINATION CRITERION is not satisfied
 - a. Randomly sample B data points from \mathcal{D} , $\{(\mathbf{x}^{(b)}, y^{(b)})\}_{b=1}^B$
 - b. Compute the gradient of the loss w.r.t. the sampled batch,
$$\nabla J^{(B)}(\boldsymbol{\theta}^{(t)})$$
 - c. Update $\boldsymbol{\theta}$: $\boldsymbol{\theta}^{(t+1)} \leftarrow \boldsymbol{\theta}^{(t)} - \gamma \nabla J^{(B)}(\boldsymbol{\theta}^{(t)})$
 - d. Increment t : $t \leftarrow t + 1$
- Output: $\boldsymbol{\theta}^{(t)}$

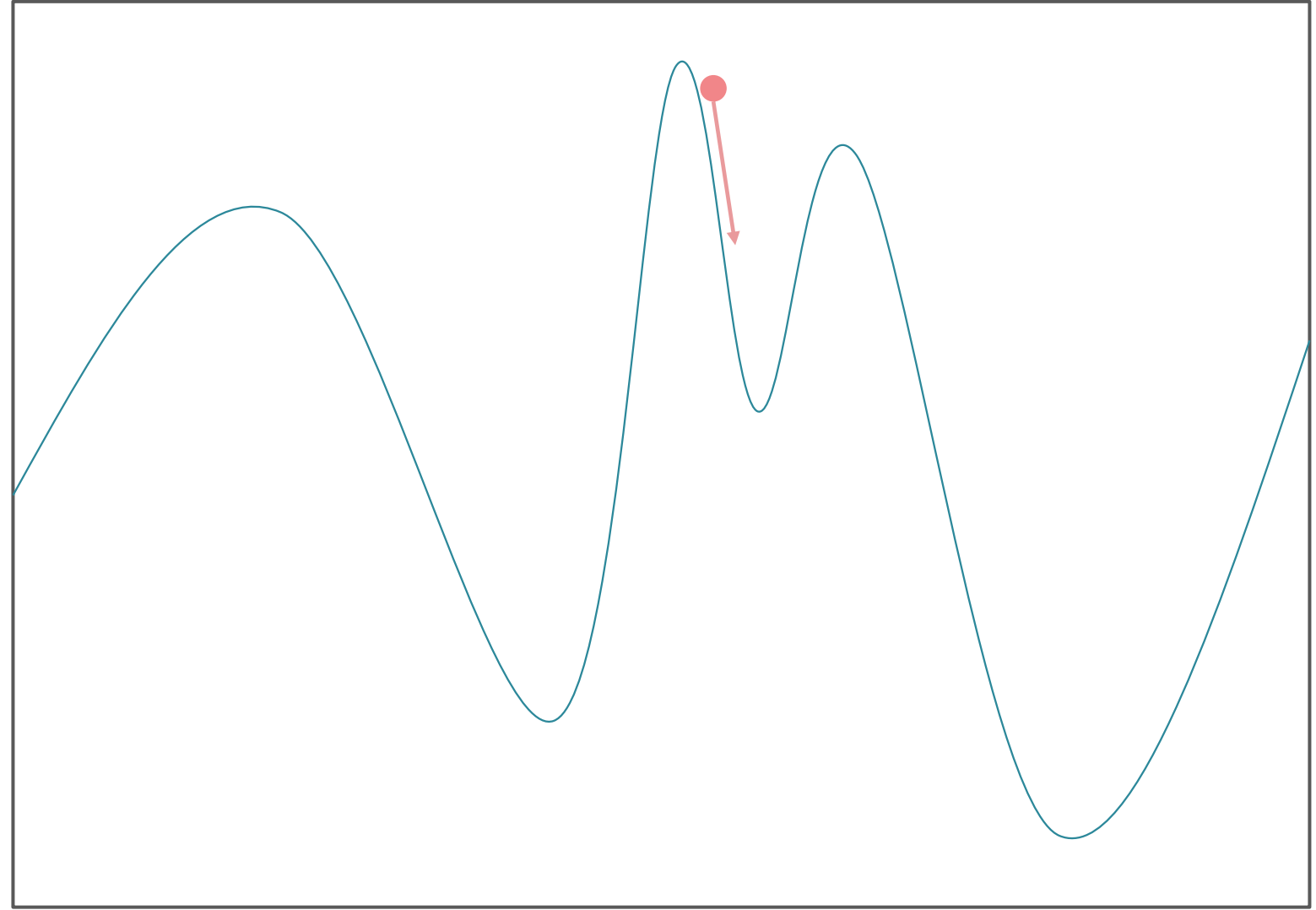
Mini-batch Stochastic Gradient Descent just the beginning!

- Input: training dataset $\mathcal{D} = \{(\mathbf{x}^{(i)}, y^{(i)})\}_{i=1}^N$,
step size γ , and batch size B
 1. *Pre-train* the parameters $\boldsymbol{\theta}^{(0)}$ and set $t = 0$
 2. While TERMINATION CRITERION is not satisfied
 - a. Randomly sample B data points from \mathcal{D} , $\{(\mathbf{x}^{(b)}, y^{(b)})\}_{b=1}^B$
 - b. Compute the gradient of the *fine-tuning* loss
$$\nabla J^{(B)}(\boldsymbol{\theta}^{(t)})$$
 - c. Update $\boldsymbol{\theta}$: $\boldsymbol{\theta}^{(t+1)} \leftarrow \boldsymbol{\theta}^{(t)} - \gamma \nabla J^{(B)}(\boldsymbol{\theta}^{(t)})$
 - d. Increment t : $t \leftarrow t + 1$
- Output: $\boldsymbol{\theta}^{(t)}$

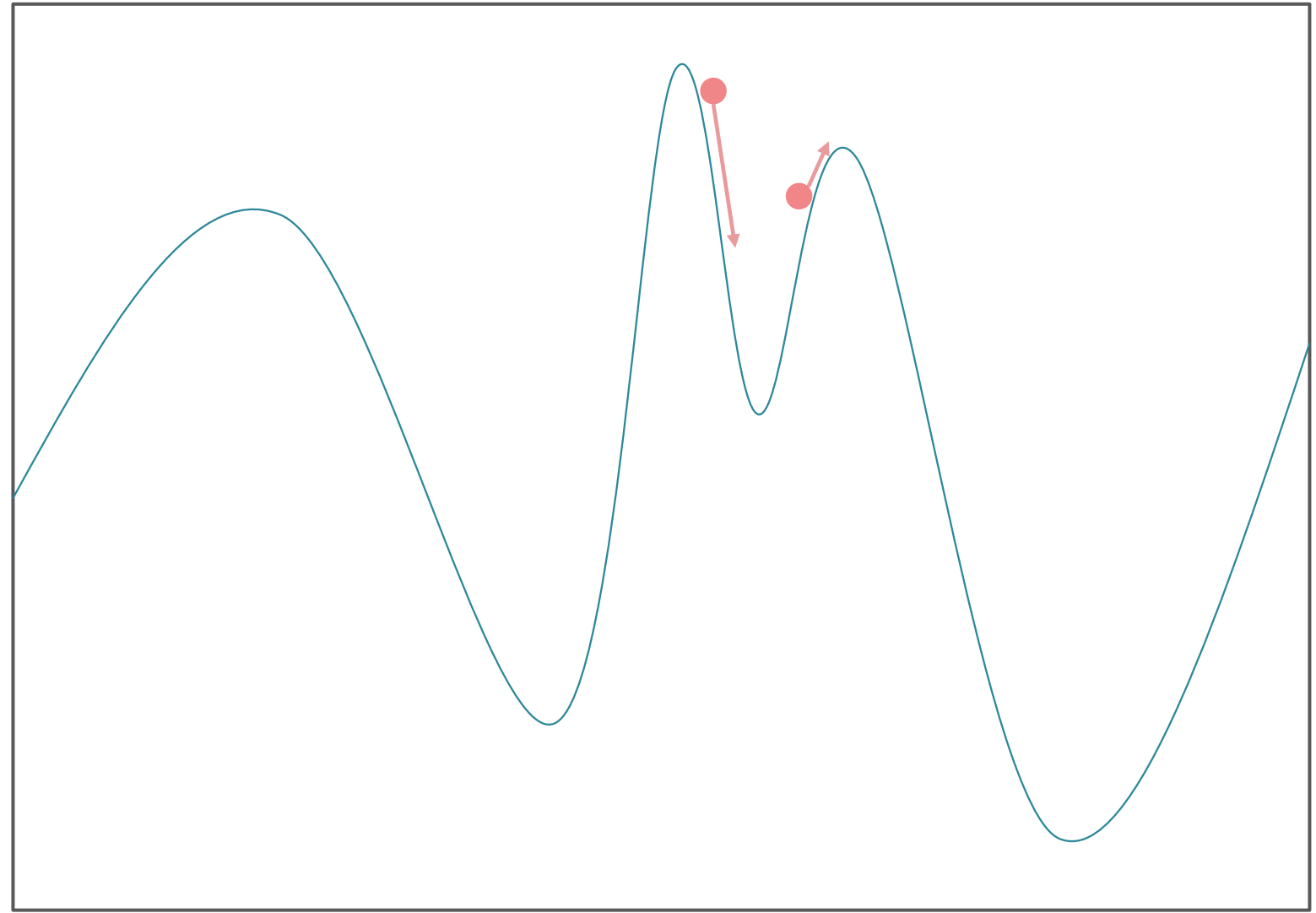
Mini-batch Stochastic Gradient Descent with Momentum

- Input: training dataset $\mathcal{D} = \{(\mathbf{x}^{(i)}, y^{(i)})\}_{i=1}^N$,
step size γ , and batch size B , decay parameter β
- 1. Pre-train the parameters $\boldsymbol{\theta}^{(0)}$ and set $t = 0$, $G_{-1} = \mathbf{0} \odot \boldsymbol{\theta}^{(0)}$
- 2. While TERMINATION CRITERION is not satisfied
 - a. Randomly sample B data points from \mathcal{D} , $\{(\mathbf{x}^{(b)}, y^{(b)})\}_{b=1}^B$
 - b. Compute the gradient of the *fine-tuning* loss
$$G_t = \nabla J^{(B)}(\boldsymbol{\theta}^{(t)})$$
 - c. Update $\boldsymbol{\theta}$: $\boldsymbol{\theta}^{(t+1)} \leftarrow \boldsymbol{\theta}^{(t)} - \gamma(\beta G_{t-1} + G_t)$
 - d. Increment t : $t \leftarrow t + 1$
- Output: $\boldsymbol{\theta}^{(t)}$

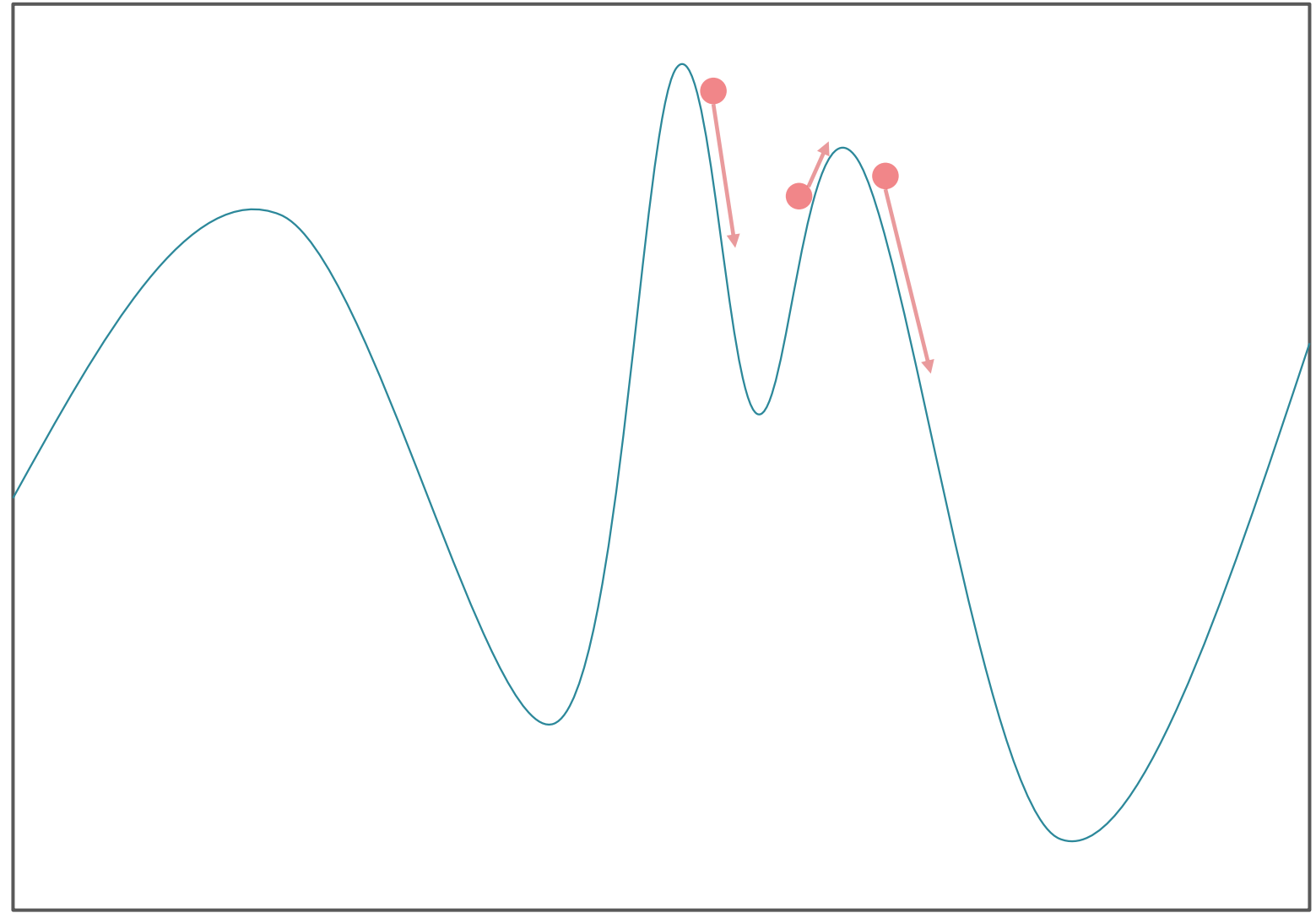
Mini-batch Stochastic Gradient Descent with Momentum



Mini-batch Stochastic Gradient Descent with Momentum



Mini-batch Stochastic Gradient Descent with Momentum



Mini-batch Stochastic Gradient Descent with Root Mean Square Propagation (RMSProp)

- Input: training dataset $\mathcal{D} = \{(\mathbf{x}^{(i)}, y^{(i)})\}_{i=1}^N$,
step size γ , and batch size B , decay parameter β
- 1. Pre-train the parameters $\boldsymbol{\theta}^{(0)}$ and set $t = 0$, $S_{-1} = 0 \odot \boldsymbol{\theta}^{(0)}$
- 2. While TERMINATION CRITERION is not satisfied
 - a. Randomly sample B data points from \mathcal{D} , $\{(\mathbf{x}^{(b)}, y^{(b)})\}_{b=1}^B$
 - b. Compute the gradient of the *fine-tuning* loss
$$G_t = \nabla J^{(B)}(\boldsymbol{\theta}^{(t)})$$
 - c. Update the scaling factor: $S_t = \beta S_{t-1} + (1 - \beta)(G_t \odot G_t)$
 - d. Update $\boldsymbol{\theta}$: $\boldsymbol{\theta}^{(t+1)} \leftarrow \boldsymbol{\theta}^{(t)} - \frac{\gamma}{\sqrt{S_t}} \odot G_t$
 - e. Increment t : $t \leftarrow t + 1$

- Output: $\boldsymbol{\theta}^{(t)}$

Adam (Adaptive Moment Estimation) = SGD + Momentum + RMSProp

- Input: training dataset $\mathcal{D} = \{(\mathbf{x}^{(i)}, y^{(i)})\}_{i=1}^N$, step size γ , and batch size B , decay parameters β_1 and β_2
 1. Pre-train the parameters $\boldsymbol{\theta}^{(0)}$, $t = 0$, $M_{-1} = S_{-1} = 0 \odot \boldsymbol{\theta}^{(0)}$
 2. While TERMINATION CRITERION is not satisfied
 - a. Randomly sample B data points from \mathcal{D} , $\{(\mathbf{x}^{(b)}, y^{(b)})\}_{b=1}^B$
 - b. Compute the gradient, momentum and scaling factor
$$G_t = \nabla J^{(B)}(\boldsymbol{\theta}^{(t)})$$
$$M_t = \beta_1 M_{t-1} + (1 - \beta_1) G_t \text{ and } S_t = \beta_2 S_{t-1} + (1 - \beta_2) (G_t \odot G_t)$$
 - c. Update $\boldsymbol{\theta}$: $\boldsymbol{\theta}^{(t+1)} \leftarrow \boldsymbol{\theta}^{(t)} - \frac{\gamma}{\sqrt{S_t / (1 - \beta_2^t)}} \odot (M_t / (1 - \beta_1^t))$
 - d. Increment t : $t \leftarrow t + 1$
- Output: $\boldsymbol{\theta}^{(t)}$