10-301/601: Introduction to Machine Learning Lecture 19: Pretraining & Fine-tuning

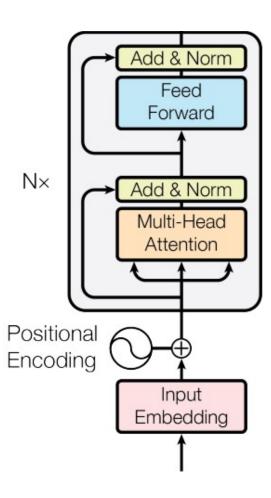
Matt Gormley & Henry Chai 3/24/25

#### **Front Matter**

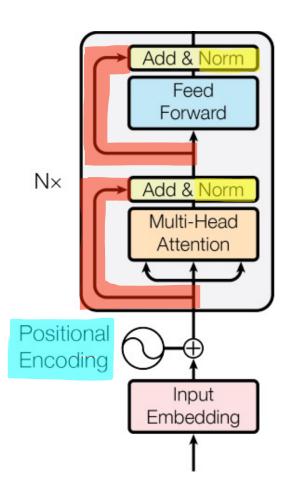
- Announcements
  - Exam 2 on 3/26 (Wednesday) from 7 9 PM
    - All topics from Lecture 8 to Lecture 16 (inclusive)
      - + MLE/MAP portion of Lecture 17 are in-scope
    - Exam 1 content may be referenced but will not be the primary focus of any question
    - Please review the seating chart on Piazza and make sure you have a seat / know where you're going

3/24/25 **2** 

#### **Transformers**



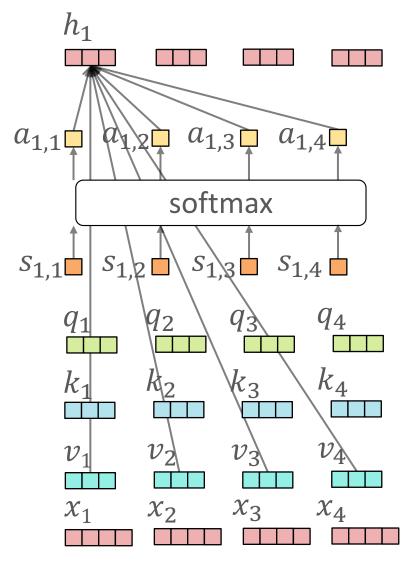
#### **Transformers**



- In addition to multi-head attention, transformer architectures use
  - 1. Positional encodings
  - 2. Layer normalization
  - 3. Residual connections
  - 4. A fully-connected feed-forward network

# Recall: Scaled Dot-product Attention

 Approach: compute a representation for each token in the *input sequence* by attending to all the input tokens



$$h_1 = \sum_{j=1}^{4} \operatorname{softmax}(s_{1,j}) v_j$$

attention weights

scores: 
$$s_{1,j} = \frac{k_j^T q_1}{\sqrt{d_k}}$$

queries: 
$$q_t = W_Q x_t$$

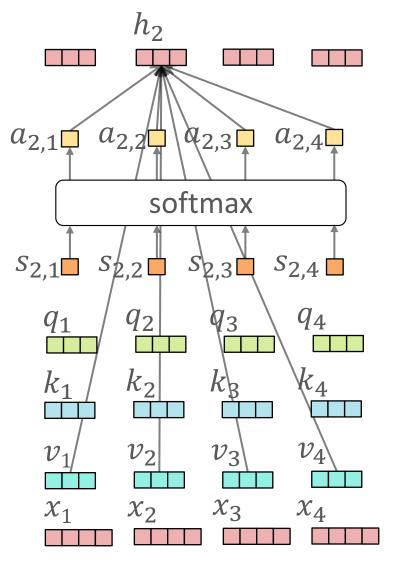
keys: 
$$k_t = W_K x_t$$

values: 
$$v_t = W_V x_t$$

input tokens

# Recall: Scaled Dot-product Attention

 Approach: compute a representation for each token in the *input sequence* by attending to all the input tokens



$$h_2 = \sum_{j=1}^{4} \operatorname{softmax}(s_{2,j}) v_j$$

attention weights

scores: 
$$s_{2,j} = \frac{k_j^T q_2}{\sqrt{d_k}}$$

queries: 
$$q_t = W_Q x_t$$

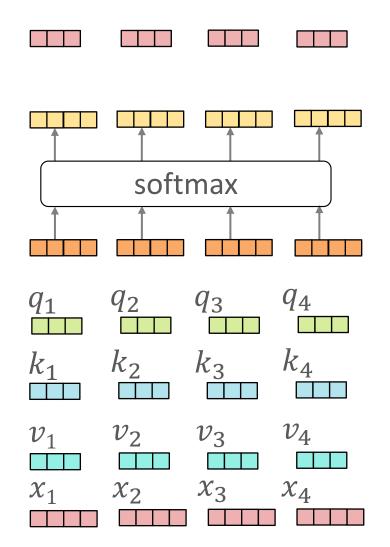
keys: 
$$k_t = W_K x_t$$

values: 
$$v_t = W_V x_t$$

input tokens

# Scaled Dot-product Attention: Matrix Form

Issue: if all tokens attend to every token in the sequence,
 then how does the model infer the order of tokens?



$$H = \operatorname{softmax}(S)V \in \mathbb{R}^{N \times d_v}$$

attention weights

scores: 
$$S = \frac{QK^T}{\sqrt{d_k}} \in \mathbb{R}^{N \times N}$$

queries: 
$$Q = XW_Q \in \mathbb{R}^{N \times d_k}$$

keys: 
$$K = XW_K \in \mathbb{R}^{N \times d_k}$$

values: 
$$V = XW_V \in \mathbb{R}^{N \times d_v}$$

design matrix: 
$$X \in \mathbb{R}^{N \times D}$$

## Positional Encodings

- Issue: if all tokens attend to every token in the sequence, then how does the model infer the order of tokens?
- Idea: add a position-specific embedding  $p_t$  to the token embedding  $x_t$

$$\tilde{x}_t = x_t + p_t$$

- Positional encodings can be
  - fixed i.e., some predetermined function of t or learned alongside the token embeddings
  - absolute i.e., only dependent on the token's location in the sequence or relative to the query token's location

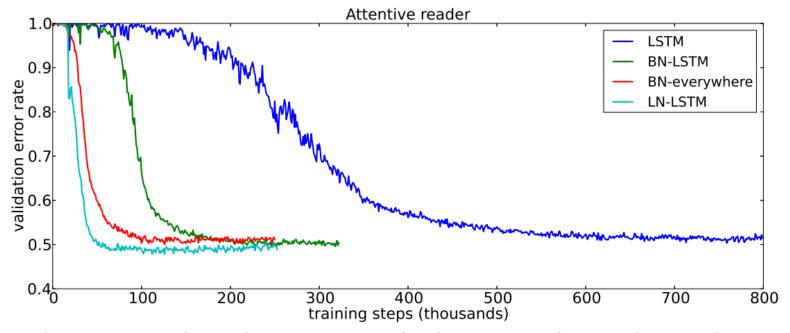
#### Layer Normalization

- Issue: for certain activation functions, the weights in later layers are **highly sensitive** to changes in the earlier layers
  - Small changes to weights in early layers are amplified so weights in deeper layers have to deal with massive dynamic ranges → slow optimization convergence
- Idea: normalize the output of a layer to always have the same (learnable) mean,  $\beta$ , and variance,  $\gamma^2$

$$H' = \gamma \left( \frac{H - \mu}{\sigma} \right) + \beta$$

where  $\mu$  is the mean and  $\sigma$  is the standard deviation of the values in the vector H

#### Layer Normalization



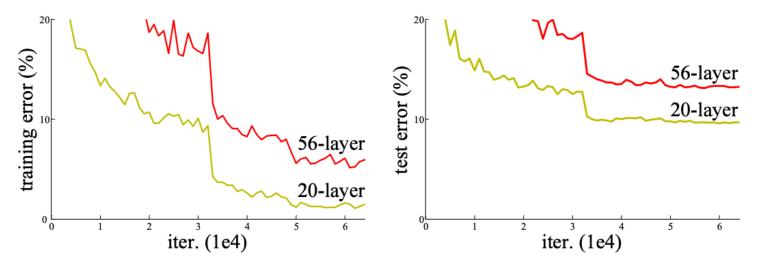
• Idea: normalize the output of a layer to always have the same (learnable) mean,  $\beta$ , and variance,  $\gamma^2$ 

$$H' = \gamma \left( \frac{H - \mu}{\sigma} \right) + \beta$$

where  $\mu$  is the mean and  $\sigma$  is the standard deviation of the values in the vector H

### Residual Connections

 Observation: early deep neural networks suffered from the "degradation" problem where adding more layers actually made performance worse!



- Wait but this is ridiculous: if the later layers aren't helping,
   couldn't they just learn the identity transformation???
- Insight: neural network layers actually have a hard time learning the identity function

### Residual Connections

- Observation: early deep neural networks suffered from the "degradation" problem where adding more layers actually made performance worse!
- Idea: add the input embedding back to the output of a layer

$$H' = H(x^{(i)}) + x^{(i)}$$

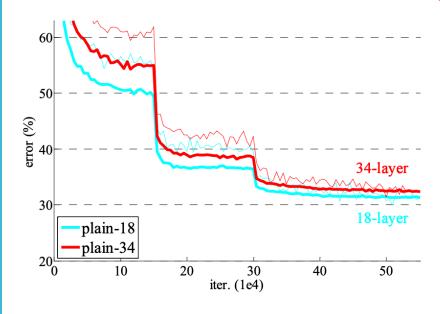
- Suppose the target function is f
  - Now instead of having to learn  $f(x^{(i)})$ , the hidden layer just needs to learn the residual  $r = f(x^{(i)}) x^{(i)}$
  - If f is the identity function, then the hidden layer just needs to learn r = 0, which is easy for a neural network!

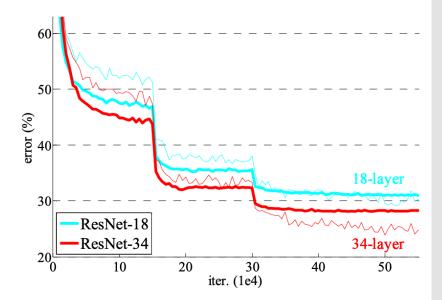
**12** 

### Residual Connections

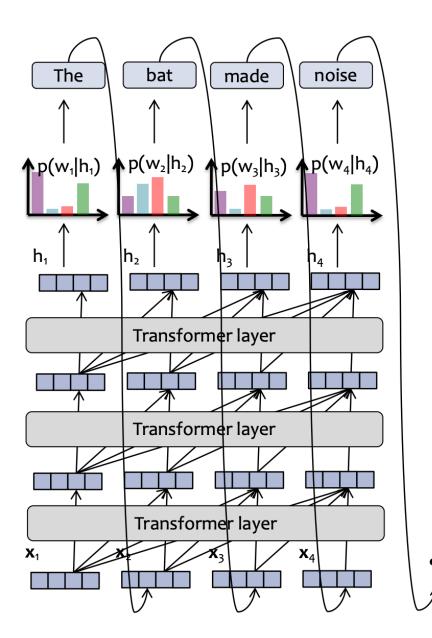
- Observation: early deep neural networks suffered from the "degradation" problem where adding more layers actually made performance worse!
- Idea: add the input embedding back to the output of a layer

$$H' = H(x^{(i)}) + x^{(i)}$$





#### **Transformers**



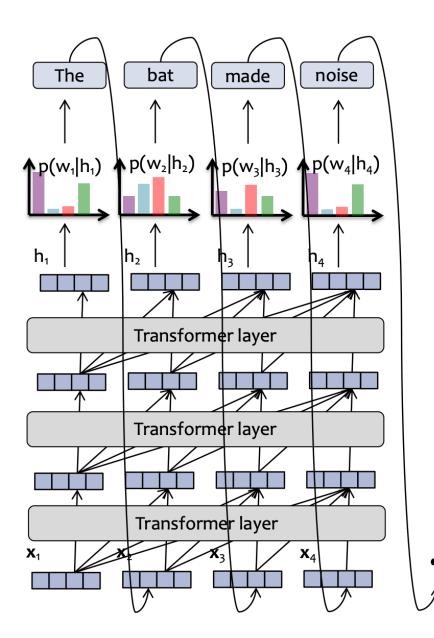
**Each layer** of a Transformer LM consists of several **sublayers**:

- 1. attention
- 2. feed-forward neural network
- 3. layer normalization
- 4. residual connections

Each hidden vector looks back at the hidden vectors of the current and previous timesteps in the previous layer.

The language model part is just like an RNN-LM.

Okay, but how on earth do we go about training these things?



**Each layer** of a Transformer LM consists of several **sublayers**:

- 1. attention
- feed-forward neural network
- 3. layer normalization
- 4. residual connections

Each hidden vector looks back at the hidden vectors of the current and previous timesteps in the previous layer.

The language model part is just like an RNN-LM.

## Backpropagation: Procedural Method

#### Algorithm 1 Forward Computation

```
1: procedure NNFORWARD(Training example (\mathbf{x}, \mathbf{y}), Params \alpha, \beta)
2: \mathbf{a} = \alpha \mathbf{x}
3: \mathbf{z} = \sigma(\mathbf{a})
4: \mathbf{b} = \beta \mathbf{z}
5: \hat{\mathbf{y}} = \operatorname{softmax}(\mathbf{b})
6: J = -\mathbf{y}^T \log \hat{\mathbf{y}}
7: \mathbf{o} = \operatorname{object}(\mathbf{x}, \mathbf{a}, \mathbf{z}, \mathbf{b}, \hat{\mathbf{y}}, J)
8: return intermediate quantities \mathbf{o}
```

#### Algorithm 2 Backpropagation

- 1: **procedure** NNBACKWARD(Training example (x, y), Params  $\alpha, \beta$ , Intermediates o)
- 2: Place intermediate quantities  $\mathbf{x}, \mathbf{a}, \mathbf{z}, \mathbf{b}, \hat{\mathbf{y}}, J$  in  $\mathbf{o}$  in scope
- 3:  $\mathbf{g}_{\hat{\mathbf{y}}} = -\mathbf{y} \div \hat{\mathbf{y}}$
- 4:  $\mathbf{g}_{\mathbf{b}} = \mathbf{g}_{\hat{\mathbf{y}}}^T \left( \mathsf{diag}(\hat{\mathbf{y}}) \hat{\mathbf{y}} \hat{\mathbf{y}}^T \right)$
- 5:  $\mathbf{g}_{\boldsymbol{\beta}} = \mathbf{g}_{\mathbf{b}}^T \mathbf{z}^T$
- 6:  $\mathbf{g}_{\mathbf{z}} = \boldsymbol{\beta}^T \mathbf{g}_{\mathbf{b}}^T$
- 7:  $\mathbf{g}_{\mathbf{a}} = \mathbf{g}_{\mathbf{z}} \odot \mathbf{z} \odot (1 \mathbf{z})$
- 8:  $\mathbf{g}_{\alpha} = \mathbf{g}_{\mathbf{a}} \mathbf{x}^T$
- 9: **return** parameter gradients  $\mathbf{g}_{\alpha}, \mathbf{g}_{\beta}$

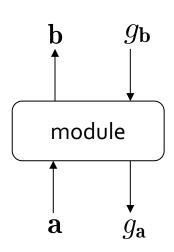
#### Issues:

- Hard to reuse / adapt for other models
- Hard to optimize individual steps
- 3. Hard to debug using the finite-difference check

**16** 

### Module-based AutoDiff

- Key Idea:
  - componentize the computation of the neuralnetwork into layers
  - each layer consolidates multiple real-valued nodes in the computation graph (a subset of them) into one vector-valued node (aka. a module)
- Each **module** is capable of two actions:
  - Forward computation of the output given some input
  - Backward computation of the gradient with respect to the input given the gradient with respect to the output



### Module-based AutoDiff

**Linear Module** The linear layer has two inputs: a vector  $\mathbf{a}$  and parameters  $\omega \in \mathbb{R}^{B \times A}$ . The output  $\mathbf{b}$  is not used by Linear Backward, but we pass it in for consistency of form.

```
1: procedure LINEARFORWARD(\mathbf{a}, \boldsymbol{\omega})
2: \mathbf{b} = \boldsymbol{\omega} \mathbf{a}
3: return \mathbf{b}
4: procedure LINEARBACKWARD(\mathbf{a}, \boldsymbol{\omega}, \mathbf{b}, \mathbf{g_b})
5: \mathbf{g}_{\boldsymbol{\omega}} = \mathbf{g_b} \mathbf{a}^T
6: \mathbf{g}_{\mathbf{a}} = \boldsymbol{\omega}^T \mathbf{g_b}
7: return \mathbf{g}_{\boldsymbol{\omega}}, \mathbf{g_a}
```

**Softmax Module** The softmax layer has only one input vector  $\mathbf{a}$ . For any vector  $\mathbf{v} \in \mathbb{R}^D$ , we have that  $\operatorname{diag}(\mathbf{v})$  returns a  $D \times D$  diagonal matrix whose diagonal entries are  $v_1, v_2, \ldots, v_D$  and whose non-diagonal entries are zero.

```
1: procedure SOFTMAXFORWARD(a)

2: \mathbf{b} = \operatorname{softmax}(\mathbf{a})

3: return \mathbf{b}

4: procedure SOFTMAXBACKWARD(a, b, \mathbf{g_b})

5: \mathbf{g_a} = \mathbf{g_b}^T \left(\operatorname{diag}(\mathbf{b}) - \mathbf{bb}^T\right)

6: return \mathbf{g_a}
```

```
Sigmoid Module The sigmoid layer has only one input vector \mathbf{a}. Below \sigma is the sigmoid applied elementwise, and \odot is element-wise multiplication \mathbf{s}. \mathbf{t}. \mathbf{u} \odot \mathbf{v} = [u_1v_1, \dots, u_Mv_M].

1: procedure SIGMOIDFORWARD(a)

2: \mathbf{b} = \sigma(\mathbf{a})

3: return \mathbf{b}

4: procedure SIGMOIDBACKWARD(a, b, \mathbf{g}_{\mathbf{b}})

5: \mathbf{g}_{\mathbf{a}} = \mathbf{g}_{\mathbf{b}} \odot \mathbf{b} \odot (1 - \mathbf{b})

6: return \mathbf{g}_{\mathbf{a}}
```

**Cross-Entropy Module** The cross-entropy layer has two inputs: a gold one-hot vector  $\mathbf{a}$  and a predicted probability distribution  $\hat{\mathbf{a}}$ . It's output  $b \in \mathbb{R}$  is a scalar. Below  $\div$  is element-wise division. The output b is not used by CrossentropyBackward, but we pass it in for consistency of form.

```
1: procedure CROSSENTROPYFORWARD(\mathbf{a}, \hat{\mathbf{a}})
2: b = -\mathbf{a}^T \log \hat{\mathbf{a}}
3: return \mathbf{b}
4: procedure CROSSENTROPYBACKWARD(\mathbf{a}, \hat{\mathbf{a}}, b, g_b)
5: \mathbf{g}_{\hat{\mathbf{a}}} = -g_b(\mathbf{a} \div \hat{\mathbf{a}})
6: return \mathbf{g}_{\mathbf{a}}
```

### Module-based AutoDiff

#### Algorithm 1 Forward Computation

```
1: procedure NNFORWARD(Training example (\mathbf{x}, \mathbf{y}), Parameters \alpha, \beta)
2: \mathbf{a} = \mathsf{LINEARFORWARD}(\mathbf{x}, \alpha)
3: \mathbf{z} = \mathsf{SIGMOIDFORWARD}(\mathbf{a})
4: \mathbf{b} = \mathsf{LINEARFORWARD}(\mathbf{z}, \beta)
5: \hat{\mathbf{y}} = \mathsf{SOFTMAXFORWARD}(\mathbf{b})
6: J = \mathsf{CROSSENTROPYFORWARD}(\mathbf{y}, \hat{\mathbf{y}})
7: \mathbf{o} = \mathsf{object}(\mathbf{x}, \mathbf{a}, \mathbf{z}, \mathbf{b}, \hat{\mathbf{y}}, J)
8: return intermediate quantities \mathbf{o}
```

#### Algorithm 2 Backpropagation

```
1: procedure NNBACKWARD(Training example (\mathbf{x}, \mathbf{y}), Parameters \alpha, \beta, Intermediates \mathbf{o})
2: Place intermediate quantities \mathbf{x}, \mathbf{a}, \mathbf{z}, \mathbf{b}, \hat{\mathbf{y}}, J in \mathbf{o} in scope
3: g_J = \frac{dJ}{dJ} = 1 \triangleright Base case
4: \mathbf{g}_{\hat{\mathbf{y}}} = \mathsf{CROSSENTROPYBACKWARD}(\mathbf{y}, \hat{\mathbf{y}}, J, g_J)
5: \mathbf{g}_{\mathbf{b}} = \mathsf{SOFTMAXBACKWARD}(\mathbf{b}, \hat{\mathbf{y}}, \mathbf{g}_{\hat{\mathbf{y}}})
6: \mathbf{g}_{\beta}, \mathbf{g}_{\mathbf{z}} = \mathsf{LINEARBACKWARD}(\mathbf{z}, \mathbf{b}, \mathbf{g}_{\mathbf{b}})
7: \mathbf{g}_{\mathbf{a}} = \mathsf{SIGMOIDBACKWARD}(\mathbf{a}, \mathbf{z}, \mathbf{g}_{\mathbf{z}})
8: \mathbf{g}_{\alpha}, \mathbf{g}_{\mathbf{x}} = \mathsf{LINEARBACKWARD}(\mathbf{x}, \mathbf{a}, \mathbf{g}_{\mathbf{a}}) \triangleright We discard \mathbf{g}_{\mathbf{x}}
9: return parameter gradients \mathbf{g}_{\alpha}, \mathbf{g}_{\beta}
```

- Easy to reuse /
   adapt for other
   models
- Individual layers are easier to optimize
- Simple to debug:just run a finite-difference checkon each layerseparately

#### Module-based AutoDiff (OOP Version)

#### Object-Oriented Implementation:

- Let each module be an object and allow the control flow of the program to define the computation graph
- No longer need to implement NNBackward(•), just follow the computation graph in reverse topological order

```
class Sigmoid (Module)

method forward (a)

\mathbf{b} = \sigma(\mathbf{a})

return \mathbf{b}

method backward (a, b, \mathbf{g_b})

\mathbf{g_a} = \mathbf{g_b} \odot \mathbf{b} \odot (1 - \mathbf{b})

return \mathbf{g_a}
```

```
class Softmax(Module)
method forward(a)
b = softmax(a)
return b
method backward(a, b, g<sub>b</sub>)
g_{a} = g_{b}^{T} (diag(b) - bb^{T})
return g<sub>a</sub>
```

```
class Linear (Module)

method forward (\mathbf{a}, \boldsymbol{\omega})

\mathbf{b} = \boldsymbol{\omega} \mathbf{a}

return \mathbf{b}

method backward (\mathbf{a}, \boldsymbol{\omega}, \mathbf{b}, \mathbf{g}_{\mathbf{b}})

\mathbf{g}_{\boldsymbol{\omega}} = \mathbf{g}_{\mathbf{b}} \mathbf{a}^{T}

\mathbf{g}_{\mathbf{a}} = \boldsymbol{\omega}^{T} \mathbf{g}_{\mathbf{b}}

return \mathbf{g}_{\boldsymbol{\omega}}, \mathbf{g}_{\mathbf{a}}
```

```
class CrossEntropy (Module)

method forward (\mathbf{a}, \hat{\mathbf{a}})

b = -\mathbf{a}^T \log \hat{\mathbf{a}}

return \mathbf{b}

method backward (\mathbf{a}, \hat{\mathbf{a}}, b, g_b)

\mathbf{g}_{\hat{\mathbf{a}}} = -g_b(\mathbf{a} \div \hat{\mathbf{a}})

return \mathbf{g}_{\mathbf{a}}
```

#### Module-based AutoDiff (OOP Version)

```
class NeuralNetwork (Module):
 2
          method init()
               lin1_layer = Linear()
               sig_layer = Sigmoid()
               lin2\_layer = Linear()
               soft_layer = Softmax()
               ce_layer = CrossEntropy()
 9
          method forward (Tensor x, Tensor y, Tensor \alpha, Tensor \beta)
10
               \mathbf{a} = \text{lin1}_{\text{layer.apply}_{\text{fwd}}}(\mathbf{x}, \boldsymbol{\alpha})
11
               z = sig_layer.apply_fwd(a)
12
               \mathbf{b} = \text{lin2}_{\text{layer.apply}_{\text{fwd}}}(\mathbf{z}, \boldsymbol{\beta})
13
               \hat{\mathbf{y}} = \text{soft\_layer.apply\_fwd}(\mathbf{b})
14
                J = \text{ce}_{\text{layer.apply}_{\text{fwd}}}(\mathbf{y}, \hat{\mathbf{y}})
15
               return J. out tensor
16
17
          method backward (Tensor x, Tensor y, Tensor \alpha, Tensor \beta)
18
                tape_bwd()
19
               return lin1_layer.in_gradients[1], lin2_layer.in_gradients[1]
20
```

3/24/25 **21** 

#### Module-based AutoDiff (OOP Version)

```
global tape = stack()
2
   class Module:
       method init()
5
           out tensor = null
           out gradient = 1
8
       method apply_fwd(List in_modules)
9
           in tensors = [x.out tensor for x in in modules]
10
           out tensor = forward(in tensors)
11
           tape.push(self)
12
           return self
13
14
       method apply_bwd():
15
           in_gradients = backward(in_tensors, out_tensor, out_gradient)
16
           for i in 1,..., len(in_modules):
17
               in_modules[i].out_gradient += in_gradients[i]
18
           return self
19
20
   function tape_bwd():
       while len(tape) > 0
22
           m = tape.pop()
23
           m.apply_bwd()
24
```

## Traditional Supervised Learning

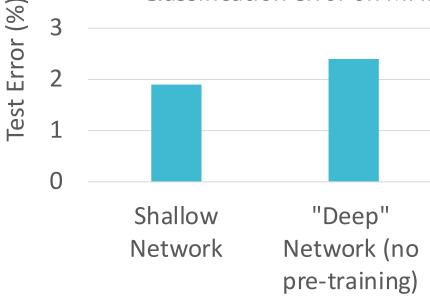
- You have some task that you want to apply machine learning to
- You have a labelled dataset to train with
- You fit a deep learning model to the dataset

**23** 

#### 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



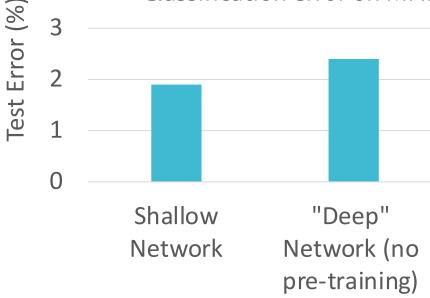


"gradient-based
 optimization starting
 from random initialization
 appears to often get
 stuck in poor solutions for
 such deep networks."

#### 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

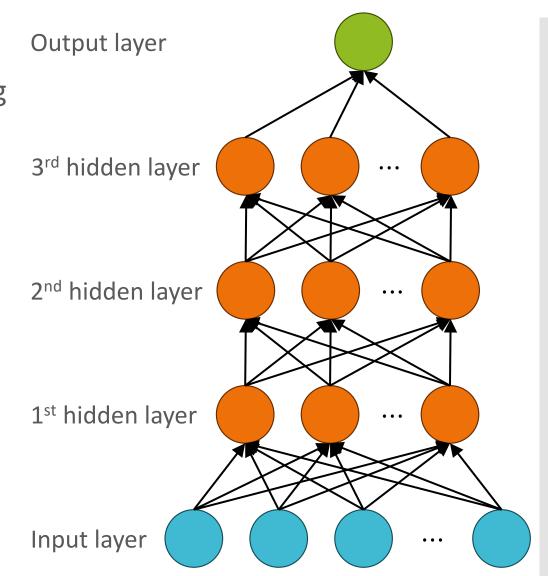




Idea: if shallow
networks are easier to
train, let's just
decompose our deep
network into a series
of shallow networks!

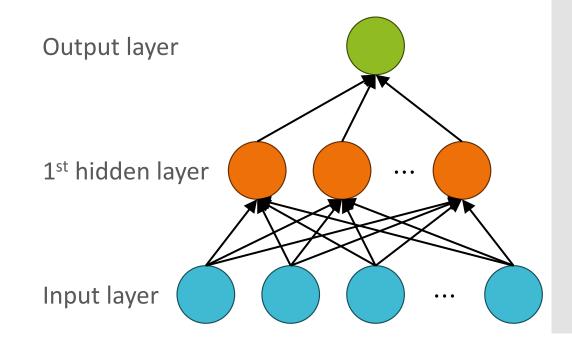
#### Pre-training (Bengio et al., 2006)

- 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



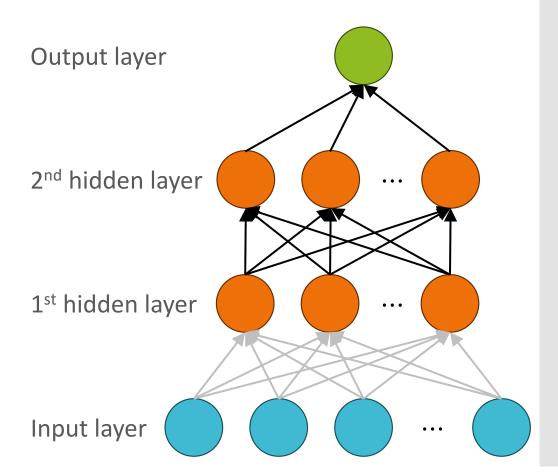
## Pre-training (Bengio et al., 2006)

- 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



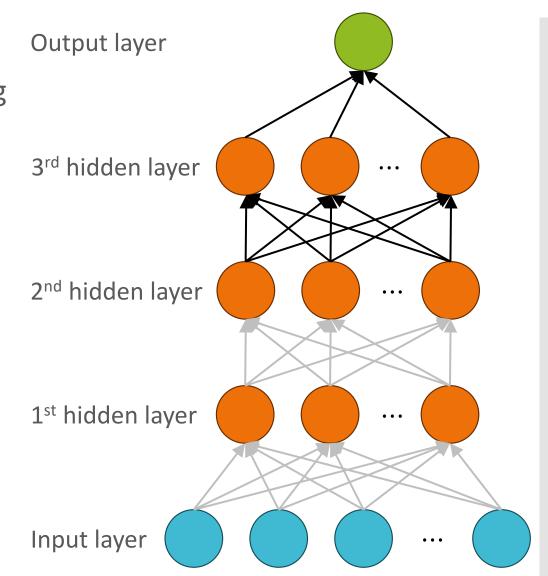
## Pre-training (Bengio et al., 2006)

- 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



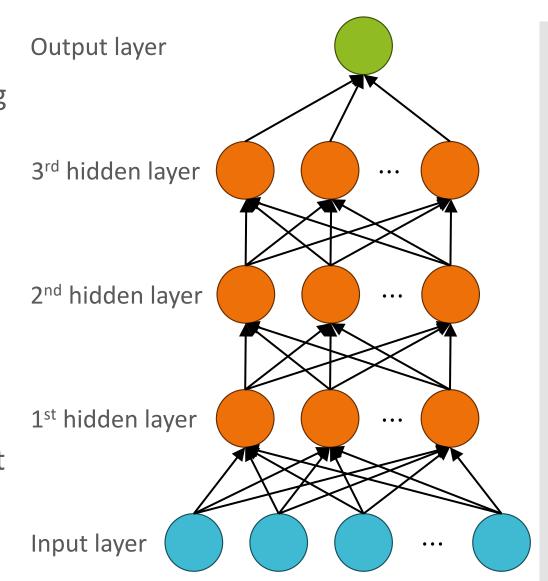
## Pre-training (Bengio et al., 2006)

- 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



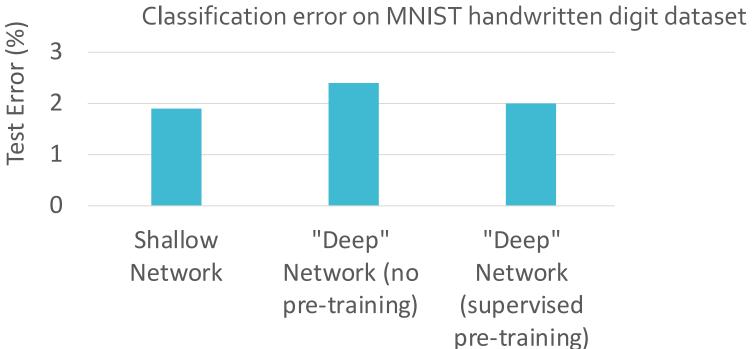
## Fine-tuning (Bengio et al., 2006)

- 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



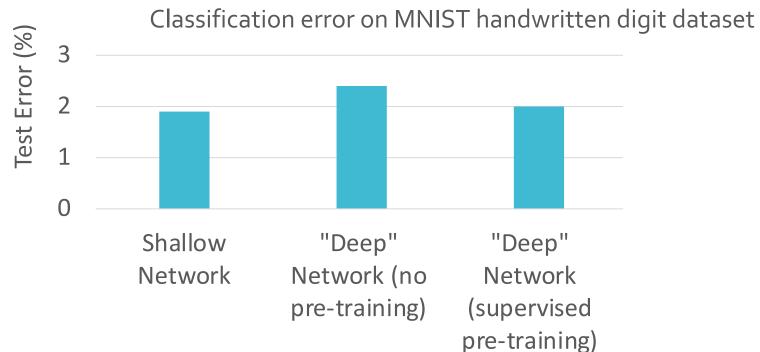
# Supervised Pre-training (Bengio et al., 2006)

 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



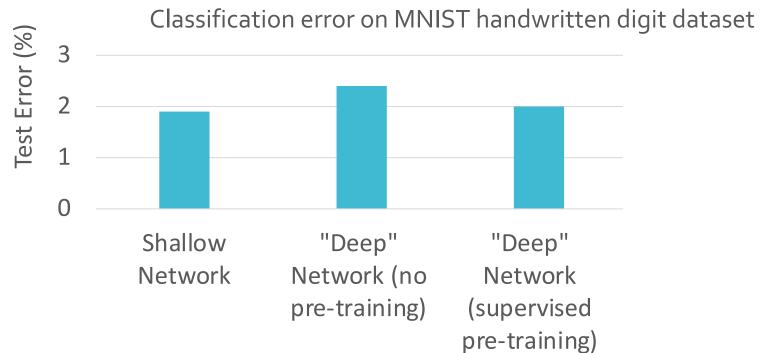
# Supervised Pre-training (Bengio et al., 2006)

- 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



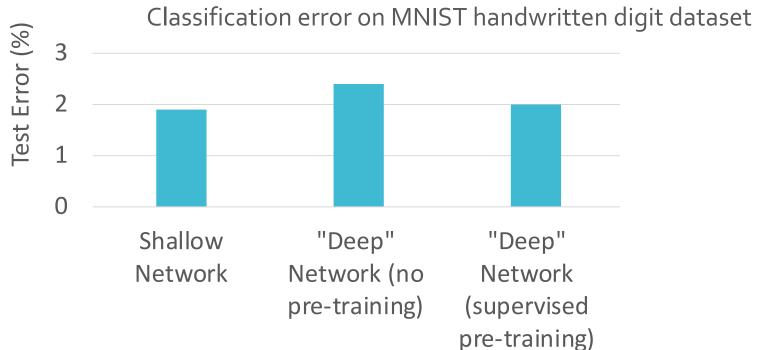
# Is this the only thing we could do with the training data?

- 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



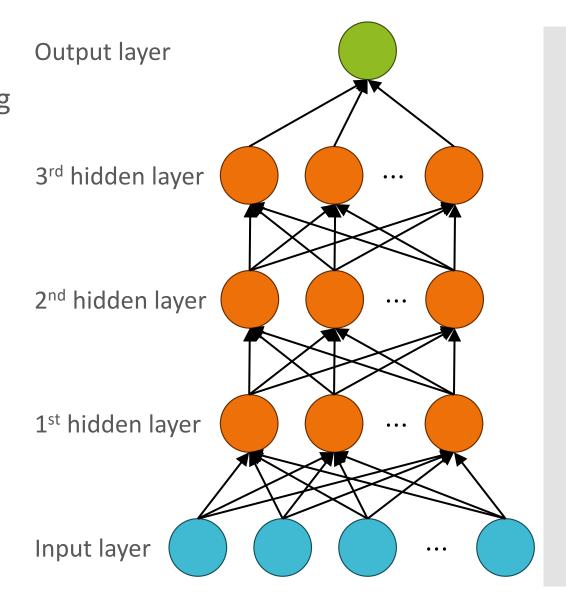
# Unsupervised Pre-training (Bengio et al., 2006)

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

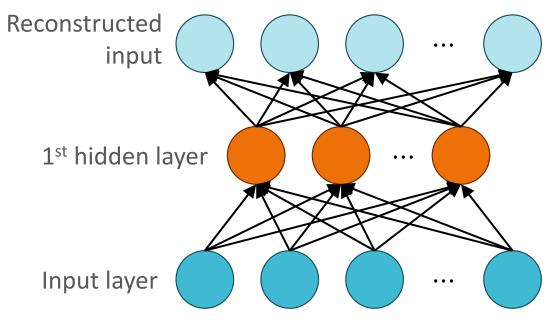


# Unsupervised Pre-training (Bengio et al., 2006)

 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

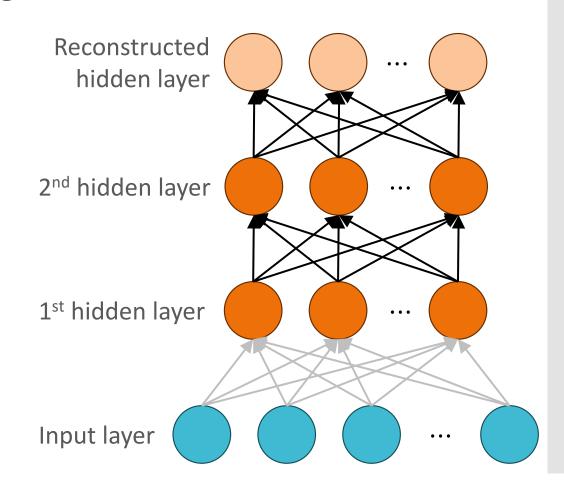


Unsupervised Pre-training (Bengio et al., 2006)

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

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

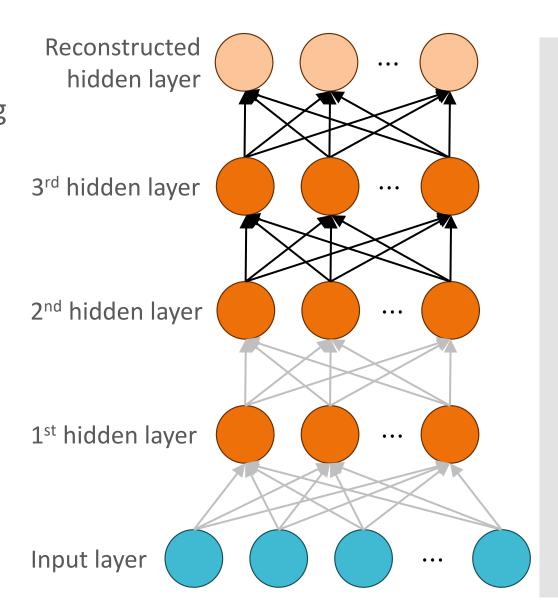
This architecture/
 objective defines an
 autoencoder



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

• Train each layer of the network iteratively using the training dataset by minimizing the reconstruction error  $||x - h(x)||_2$ 

This architecture/
 objective defines an
 autoencoder

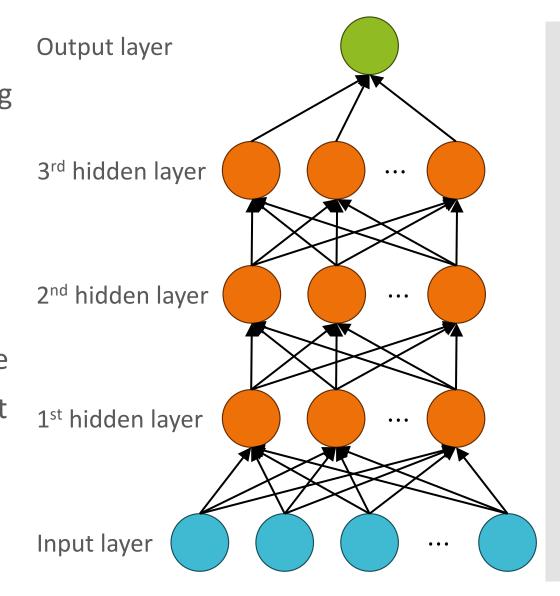


## Fine-tuning (Bengio et al., 2006)

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

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

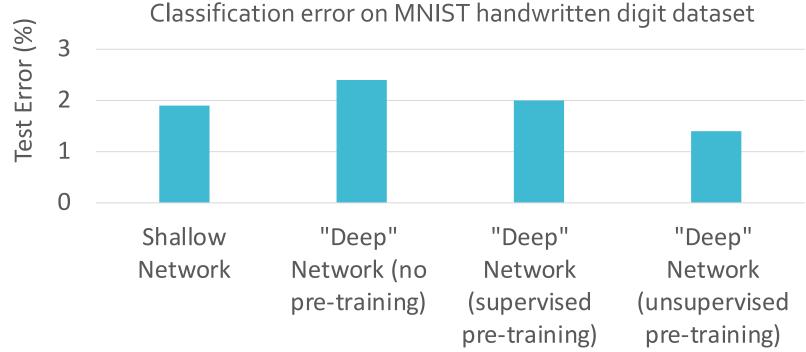
 When fine-tuning, we're effectively swapping out the last layer and fitting all the weights to the training dataset



## 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

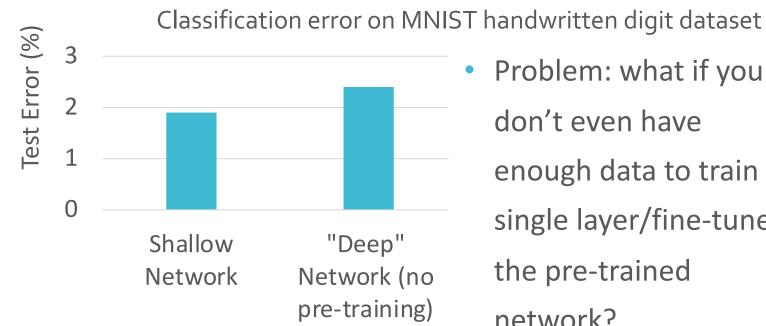
40



3/24/25 Pre training)

## 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

3/24/25

## 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. you don't have enough data to fine-tune your model?
  - B. the concept of labelled data doesn't apply to your task i.e., not every input has a "correct" label e.g., chatbots?

3/24/25

## 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

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.

```
Translate English to French: 

task description

sea otter => loutre de mer examples

peppermint => menthe poivrée

plush girafe => girafe peluche

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.



• Idea: leverage the LLM's context window by passing a few one 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.

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.



Source: https://arxiv.org/pdf/2005.14165.pdf

• 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.

```
1 Translate English to French: ← task description
2 cheese => ← prompt
```

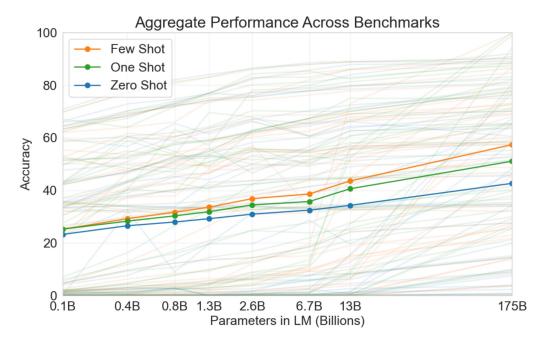
Traditional fine-tuning (not used for GPT-3)

## Fine-tuning

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



• 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



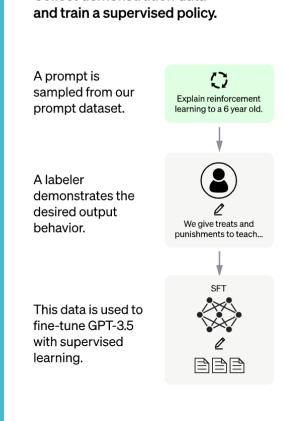
Key Takeaway: LLMs can perform well on novel tasks
without having to fine-tune the model, sometimes even
with just one or zero labelled training data points!

## 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!

**3**/24/25 **50** 

## Reinforcement Learning from Human Feedback (RLHF)



Collect demonstration data

Step 1

Collect comparison data and train a reward model. A prompt and several model Explain reinforcement outputs are learning to a 6 year old. sampled. A labeler ranks the outputs from best to worst. D > G > A > B This data is used

Step 2

Optimize a policy against the reward model using the PPO reinforcement learning algorithm. A new prompt is sampled from Write a story the dataset. about otters. The PPO model is initialized from the supervised policy. The policy generates Once upon a time... an output. The reward model calculates a reward for the output. The reward is used to update the policy using PPO.

Step 3

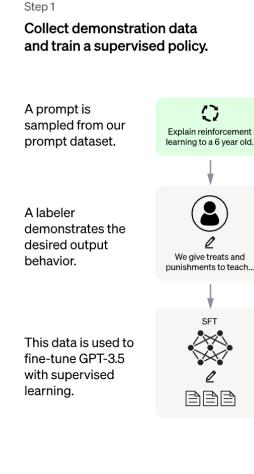
 RLHF is a form of fine-tuning that uses reinforcement learning where the reward function is learned from human preferences

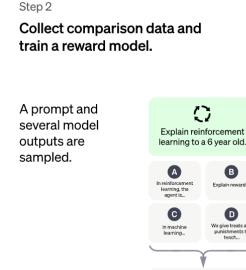
3/24/25 Source: <a href="https://openai.com/blog/chatgpt">https://openai.com/blog/chatgpt</a>
51

to train our

reward model.

## What the heck is "Reinforcement Learning"?





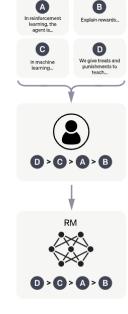
A labeler ranks the

outputs from best

This data is used to train our

reward model.

to worst.



Step 3

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

A new prompt is sampled from Write a story the dataset. about otters. The PPO model is initialized from the supervised policy. The policy generates Once upon a time... an output. The reward model calculates a reward for the output. The reward is used to update the policy using PPO.

• RLHF is a form of fine-tuning that uses *reinforcement learning* where the reward function is learned from human preferences

3/24/25 Source: https://openai.com/blog/chatgpt 52