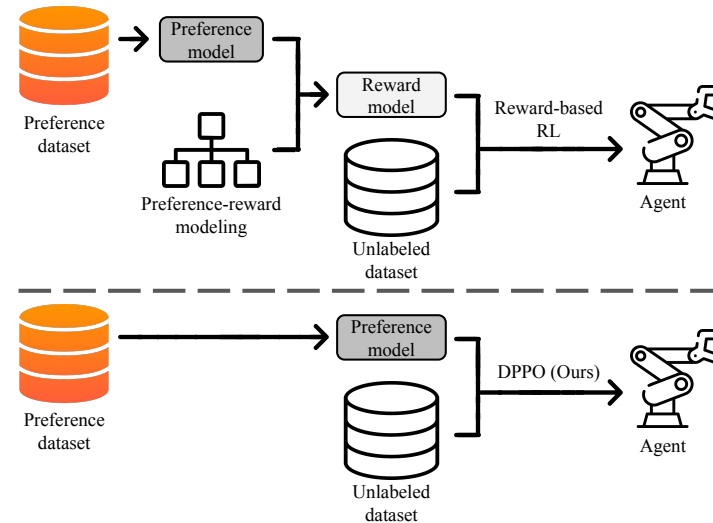


# Large Language Models and RLHF

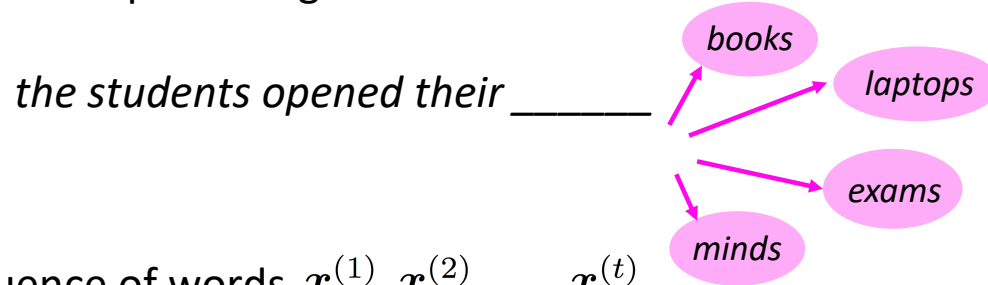
## Recall: Learning from Preferences



- One approach is to fit a reward model (eg Bradley-Terry)
- Another approach is to directly optimize a policy based on preferences
- Optimal policies exist for preference relations that are total consistent pre-orders, even if a corresponding reward function does not exist
- Today: more discussion on LLMs

# What is a language model?

- **Language Modeling** is the task of predicting what word comes next



- More formally: given a sequence of words  $\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots, \mathbf{x}^{(t)}$ , compute the probability distribution of the next word  $\mathbf{x}^{(t+1)}$ :

$$P(\mathbf{x}^{(t+1)} \mid \mathbf{x}^{(t)}, \dots, \mathbf{x}^{(1)})$$

where  $\mathbf{x}^{(t+1)}$  can be any word in the vocabulary  $V = \{w_1, \dots, w_{|V|}\}$

- A system that does this is called a **Language Model**

# Probabilistic language models

- You can also think of a Language Model as a system that **assigns a probability to a piece of text**
- For example, if we have some text  $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(T)}$ , then the probability of this text (according to the Language Model) is:

$$\begin{aligned} P(\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(T)}) &= P(\mathbf{x}^{(1)}) \times P(\mathbf{x}^{(2)} | \mathbf{x}^{(1)}) \times \dots \times P(\mathbf{x}^{(T)} | \mathbf{x}^{(T-1)}, \dots, \mathbf{x}^{(1)}) \\ &= \prod_{t=1}^T P(\mathbf{x}^{(t)} | \mathbf{x}^{(t-1)}, \dots, \mathbf{x}^{(1)}) \end{aligned}$$



This is what our LM provides

# N-gram models

*the students opened their \_\_\_\_\_*

- **Question:** How to learn a Language Model?
- **Answer** (pre- Deep Learning): learn an *n*-gram Language Model!
- **Definition:** An *n*-gram is a chunk of *n* consecutive words.
  - **unigrams:** “the”, “students”, “opened”, “their”
  - **bigrams:** “the students”, “students opened”, “opened their”
  - **trigrams:** “the students opened”, “students opened their”
  - **four-grams:** “the students opened their”
- **Idea:** Collect statistics about how frequent different n-grams are and use these to predict next word.

# N-gram models and Markov assumption

- First we make a **Markov assumption**:  $x^{(t+1)}$  depends only on the preceding  $n-1$  words

$$P(\mathbf{x}^{(t+1)} | \mathbf{x}^{(t)}, \dots, \mathbf{x}^{(1)}) = P(\mathbf{x}^{(t+1)} | \overbrace{\mathbf{x}^{(t)}, \dots, \mathbf{x}^{(t-n+2)}}^{n-1 \text{ words}}) \quad (\text{assumption})$$

prob of a n-gram  $\rightarrow$   $P(\mathbf{x}^{(t+1)}, \mathbf{x}^{(t)}, \dots, \mathbf{x}^{(t-n+2)})$

prob of a (n-1)-gram  $\rightarrow$   $P(\mathbf{x}^{(t)}, \dots, \mathbf{x}^{(t-n+2)})$

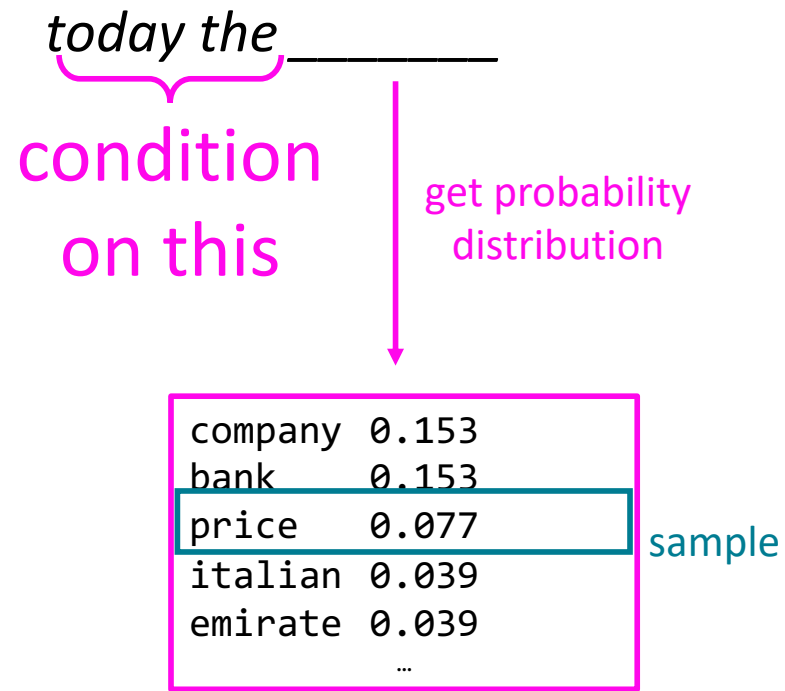
$\Rightarrow$  (definition of conditional prob)

- Question:** How do we get these  $n$ -gram and  $(n-1)$ -gram probabilities?
- Answer:** By **counting** them in some large corpus of text!

$$\approx \frac{\text{count}(\mathbf{x}^{(t+1)}, \mathbf{x}^{(t)}, \dots, \mathbf{x}^{(t-n+2)})}{\text{count}(\mathbf{x}^{(t)}, \dots, \mathbf{x}^{(t-n+2)})} \quad (\text{statistical approximation})$$

# Generating text

You can also use a Language Model to generate text



# Generating text

You can also use a Language Model to generate text

*today the price of gold per ton , while production of shoe lasts and shoe industry , the bank intervened just after it considered and rejected an imf demand to rebuild depleted european stocks , sept 30 end primary 76 cts a share .*

Surprisingly grammatical!

...but **incoherent**. We need to consider more than three words at a time if we want to model language well.

But increasing  $n$  worsens sparsity problem,  
and increases model size...



# How good are n-gram models?

You can also use a Language Model to generate text

*today the price of gold per ton , while production of shoe lasts and shoe industry , the bank intervened just after it considered and rejected an imf demand to rebuild depleted european stocks , sept 30 end primary 76 cts a share .*

Surprisingly grammatical!

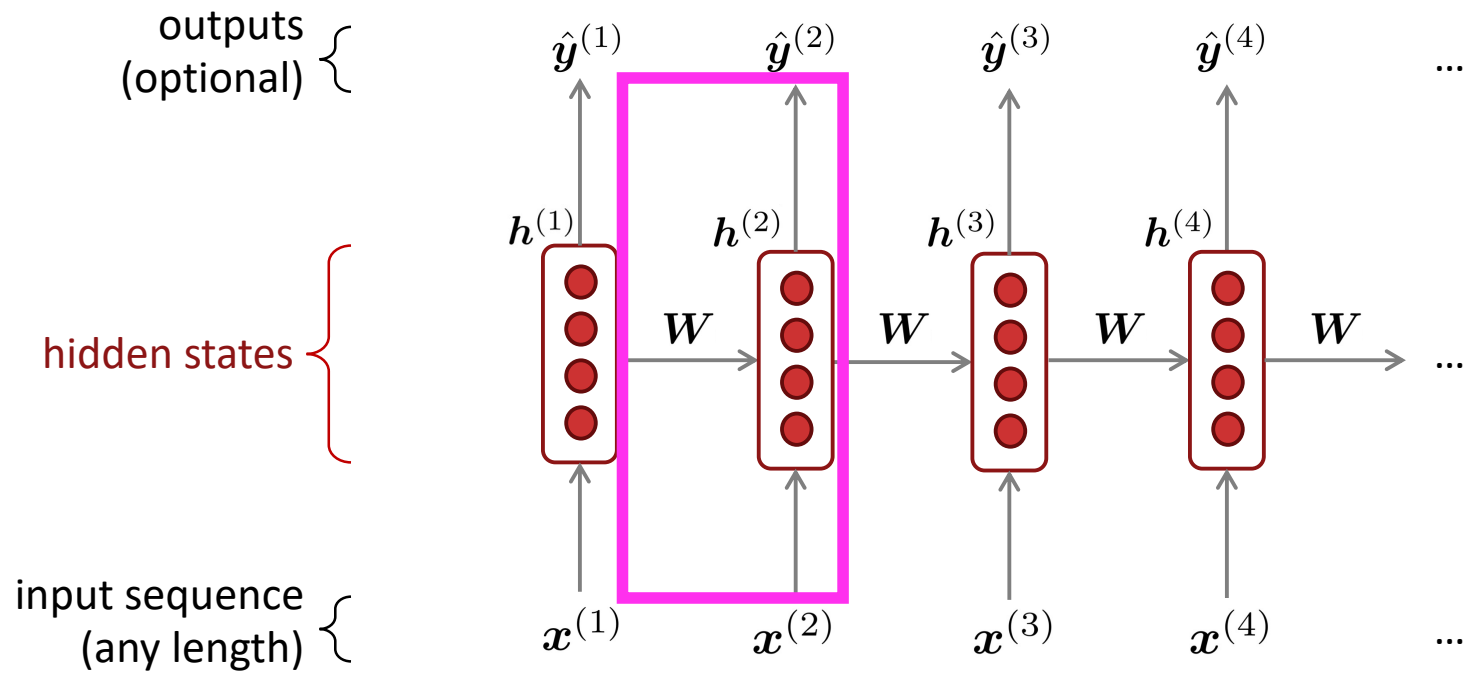
...but **incoherent**. We need to consider more than three words at a time if we want to model language well.

But increasing  $n$  worsens sparsity problem,  
and increases model size...

## Problems with n-gram models

- Small  $n$  means model is not good enough
- Large  $n$  means that many combinations do not occur in the data - *sparsity*
- Generally speaking, fixed  $n$  is very rigid

# Recurrent nets (RNNs)



# A Simple RNN Language Model

output distribution

$$\hat{y}^{(t)} = \text{softmax}(U\mathbf{h}^{(t)} + \mathbf{b}_2) \in \mathbb{R}^{|V|}$$

hidden states

$$\mathbf{h}^{(t)} = \sigma(\mathbf{W}_h \mathbf{h}^{(t-1)} + \mathbf{W}_e \mathbf{e}^{(t)} + \mathbf{b}_1)$$

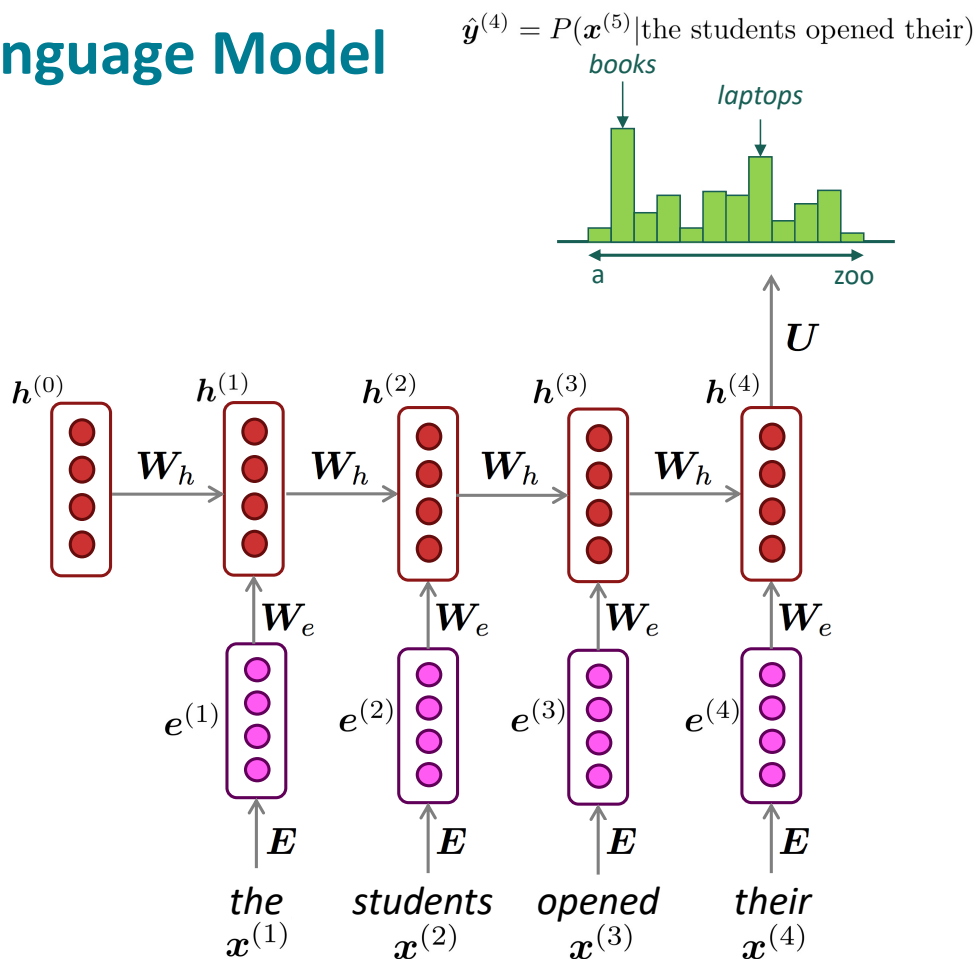
$\mathbf{h}^{(0)}$  is the initial hidden state

word embeddings

$$\mathbf{e}^{(t)} = \mathbf{E} \mathbf{x}^{(t)}$$

words / one-hot vectors

$$\mathbf{x}^{(t)} \in \mathbb{R}^{|V|}$$



Note: this input sequence could be much longer now!

# RNN training

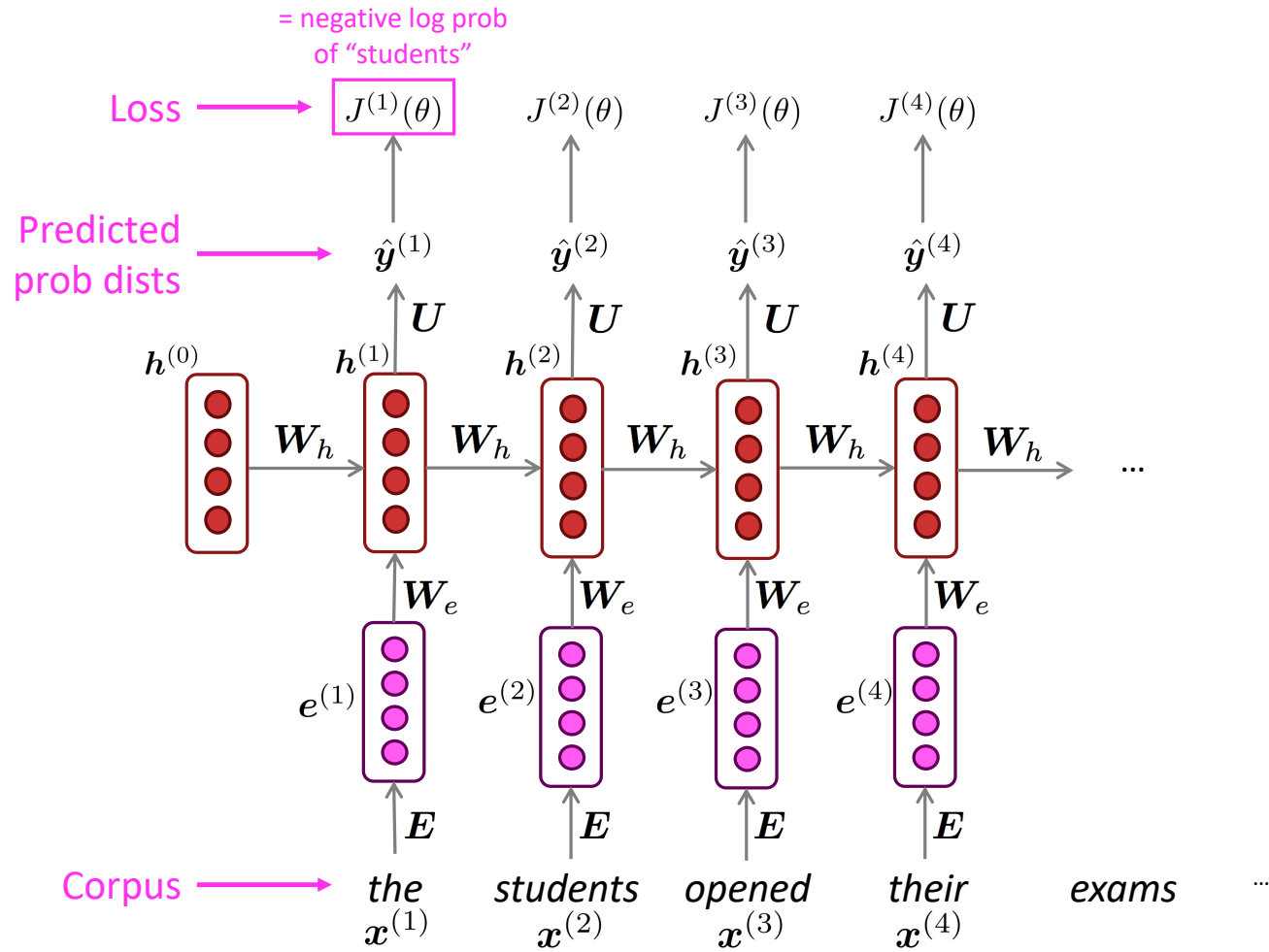
- Get a **big corpus of text** which is a sequence of words  $x^{(1)}, \dots, x^{(T)}$
- Feed into RNN-LM; compute output distribution  $\hat{\mathbf{y}}^{(t)}$  **for every step  $t$** .
  - i.e., predict probability dist of *every word*, given words so far
- **Loss function** on step  $t$  is **cross-entropy** between predicted probability distribution  $\hat{\mathbf{y}}^{(t)}$ , and the true next word  $\mathbf{y}^{(t)}$  (one-hot for  $x^{(t+1)}$ ):

$$J^{(t)}(\theta) = CE(\mathbf{y}^{(t)}, \hat{\mathbf{y}}^{(t)}) = - \sum_{w \in V} \mathbf{y}_w^{(t)} \log \hat{\mathbf{y}}_w^{(t)} = - \log \hat{\mathbf{y}}_{\mathbf{x}_{t+1}}^{(t)}$$

- Average this to get **overall loss** for entire training set:

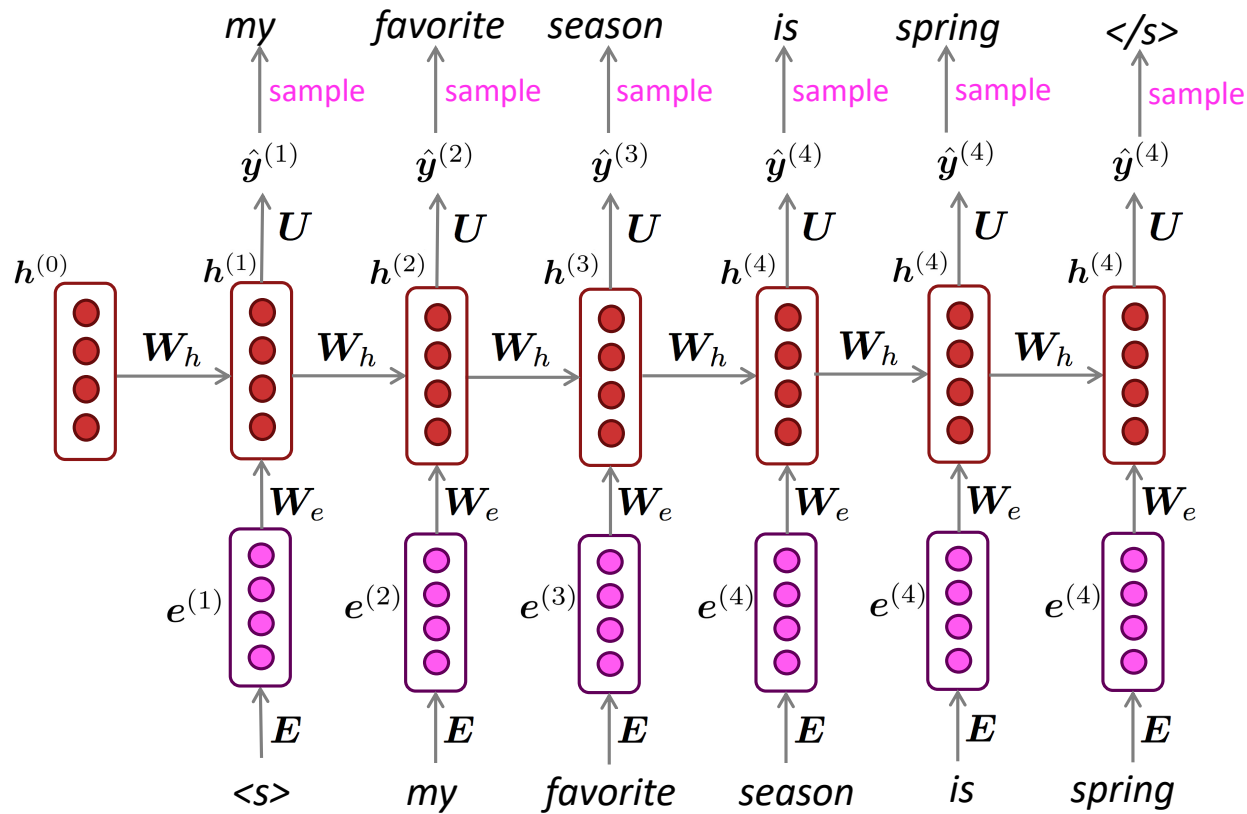
$$J(\theta) = \frac{1}{T} \sum_{t=1}^T J^{(t)}(\theta) = \frac{1}{T} \sum_{t=1}^T - \log \hat{\mathbf{y}}_{\mathbf{x}_{t+1}}^{(t)}$$

# RNN training



# Generating text with RNNs

Just like an n-gram Language Model, you can use a RNN Language Model to **generate text** by **repeated sampling**. Sampled output becomes next step's input.



# RNN example

Let's have some fun!

- You can train an RNN-LM on any kind of text, then generate text in that style.
- RNN-LM trained on *Harry Potter*:



“Sorry,” Harry shouted, panicking—“I’ll leave those brooms in London, are they?”

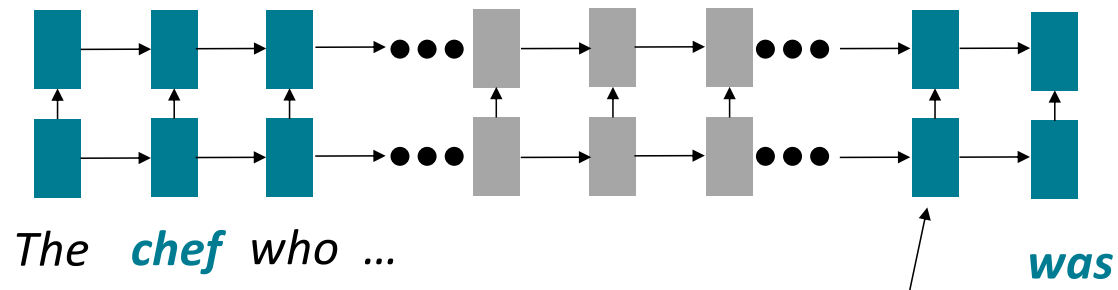
“No idea,” said Nearly Headless Nick, casting low close by Cedric, carrying the last bit of treacle Charms, from Harry’s shoulder, and to answer him the common room perched upon it, four arms held a shining knob from when the spider hadn’t felt it seemed. He reached the teams too.

**Source:** <https://medium.com/deep-writing/harry-potter-written-by-artificial-intelligence-8a9431803da6>



## Issues with RNNs: Linear interaction distance

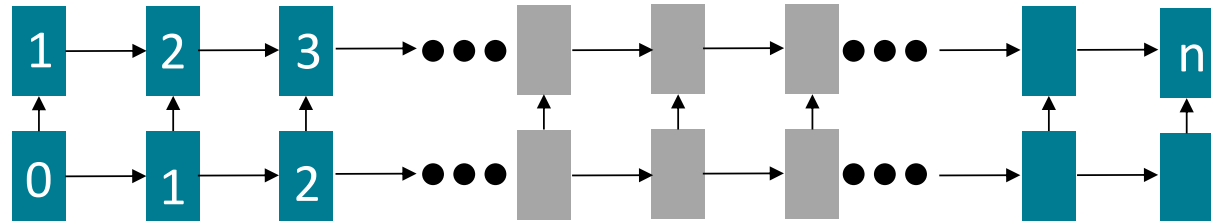
- **O(sequence length)** steps for distant word pairs to interact means:
  - Hard to learn long-distance dependencies (because gradient problems!)
  - Linear order of words is “baked in”; we already know linear order isn’t the right way to think about sentences...



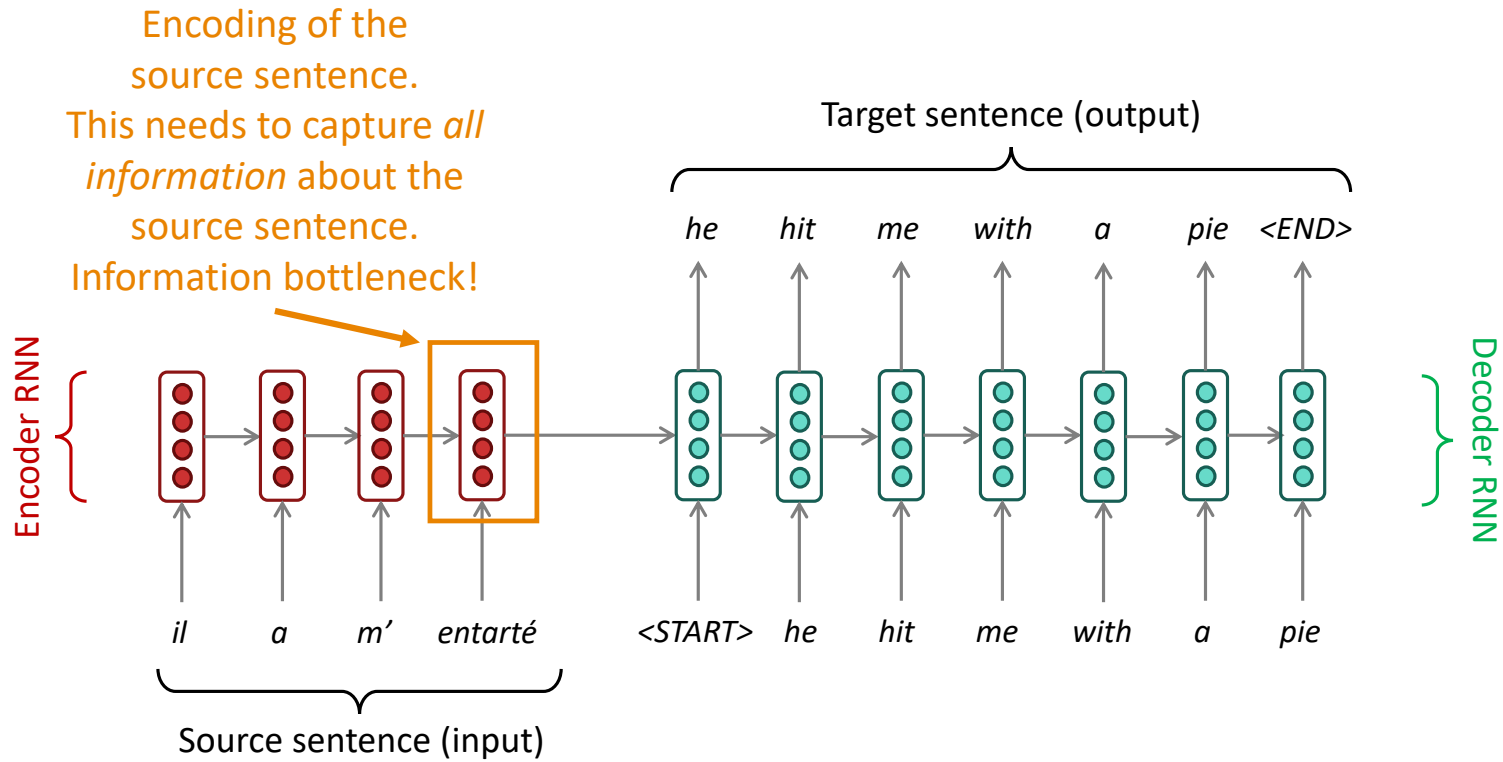
Info of *chef* has gone through  $O(\text{sequence length})$  many layers!

## Issues with RNNs: Hard to parallelize

- Forward and backward passes have  **$O(\text{sequence length})$**  unparallelizable operations
  - GPUs can perform a bunch of independent computations at once!
  - But future RNN hidden states can't be computed in full before past RNN hidden states have been computed
  - Inhibits training on very large datasets!



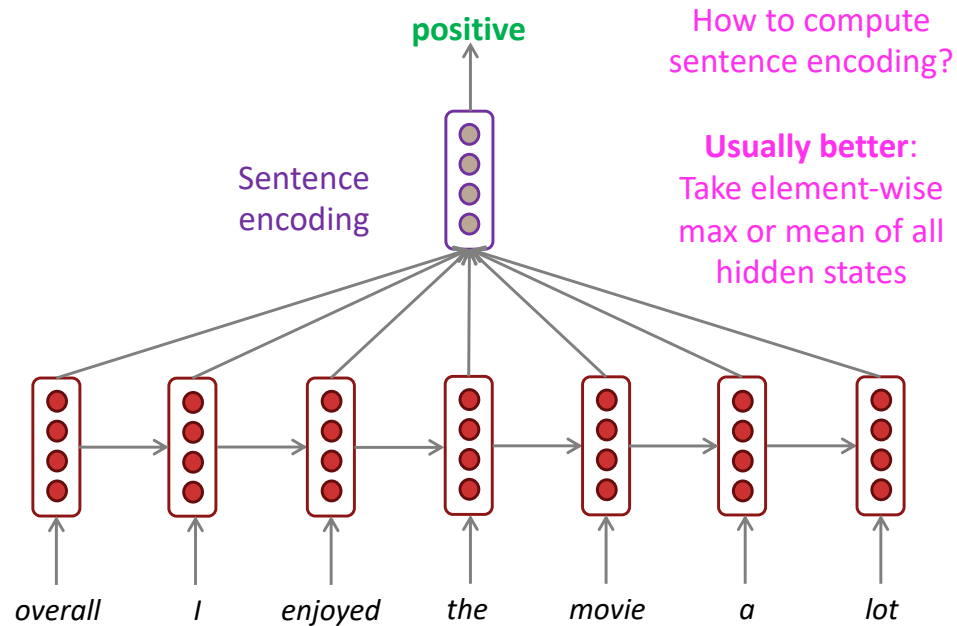
# Issues with RNNs: Bottleneck problem



## Solution: Attention

- On each step of decoding, use *direct connection to the encoder* to focus on a particular part of the sequence
- A bit like what humans do!
- Attention provides a solution to the bottleneck problem!

# Pooling in RNNs

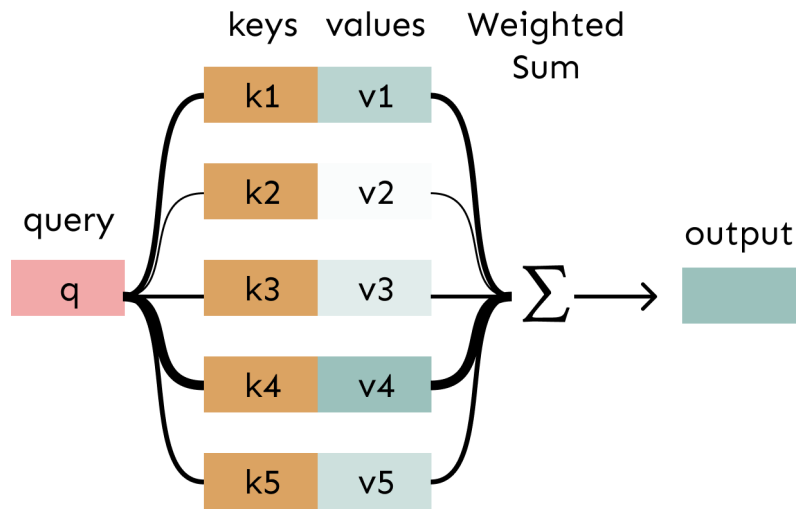


- Starting point: a *very* basic way of 'passing information from the encoder' is to *average*

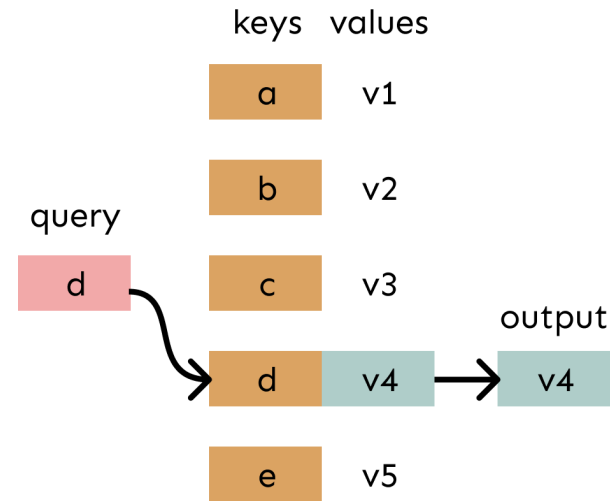
# Attention is weighted averaging!

Attention is just a **weighted** average – this is very powerful if the weights are learned!

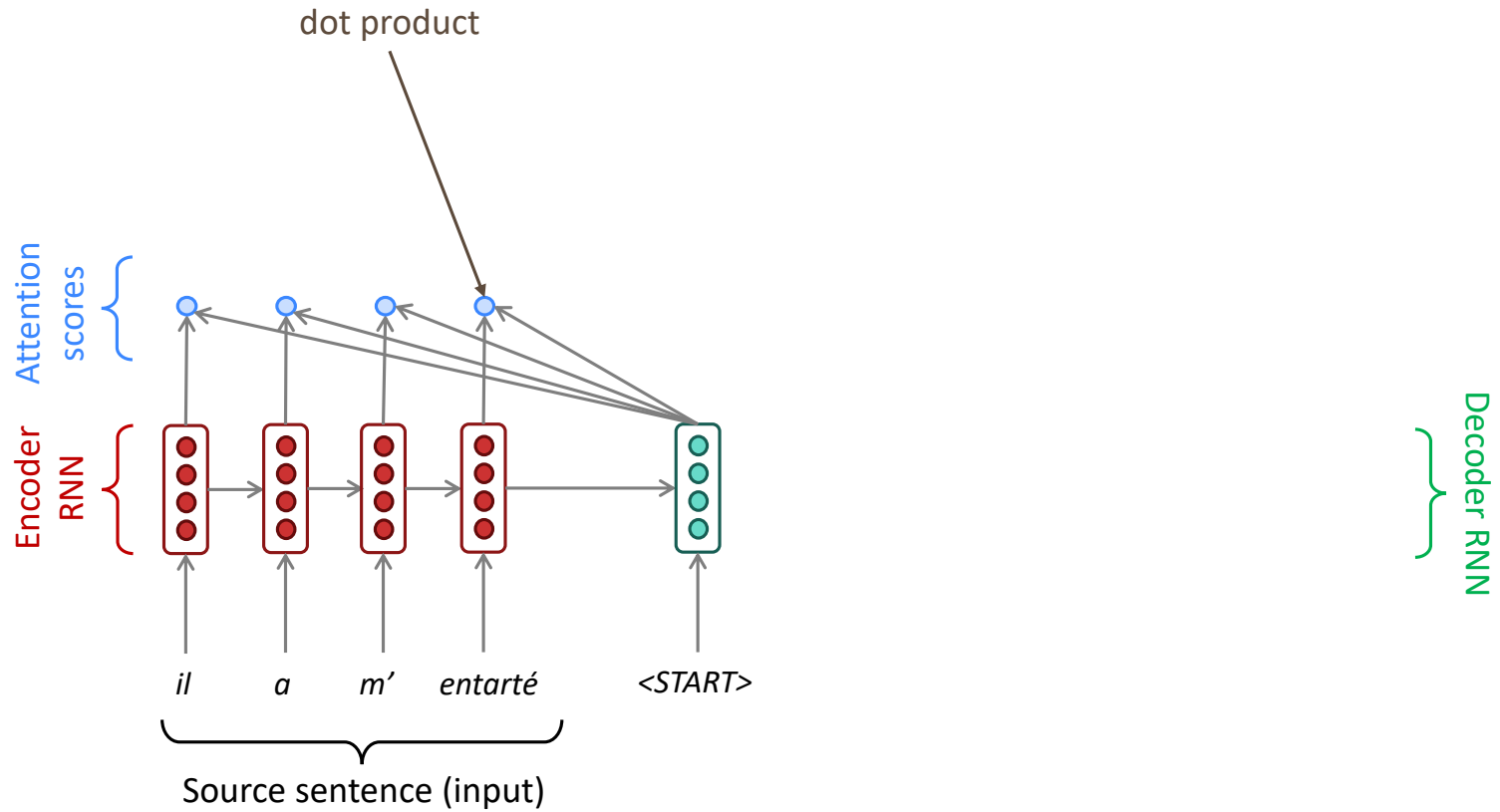
In **attention**, the **query** matches all **keys** *softly*, to a weight between 0 and 1. The keys' **values** are multiplied by the weights and summed.



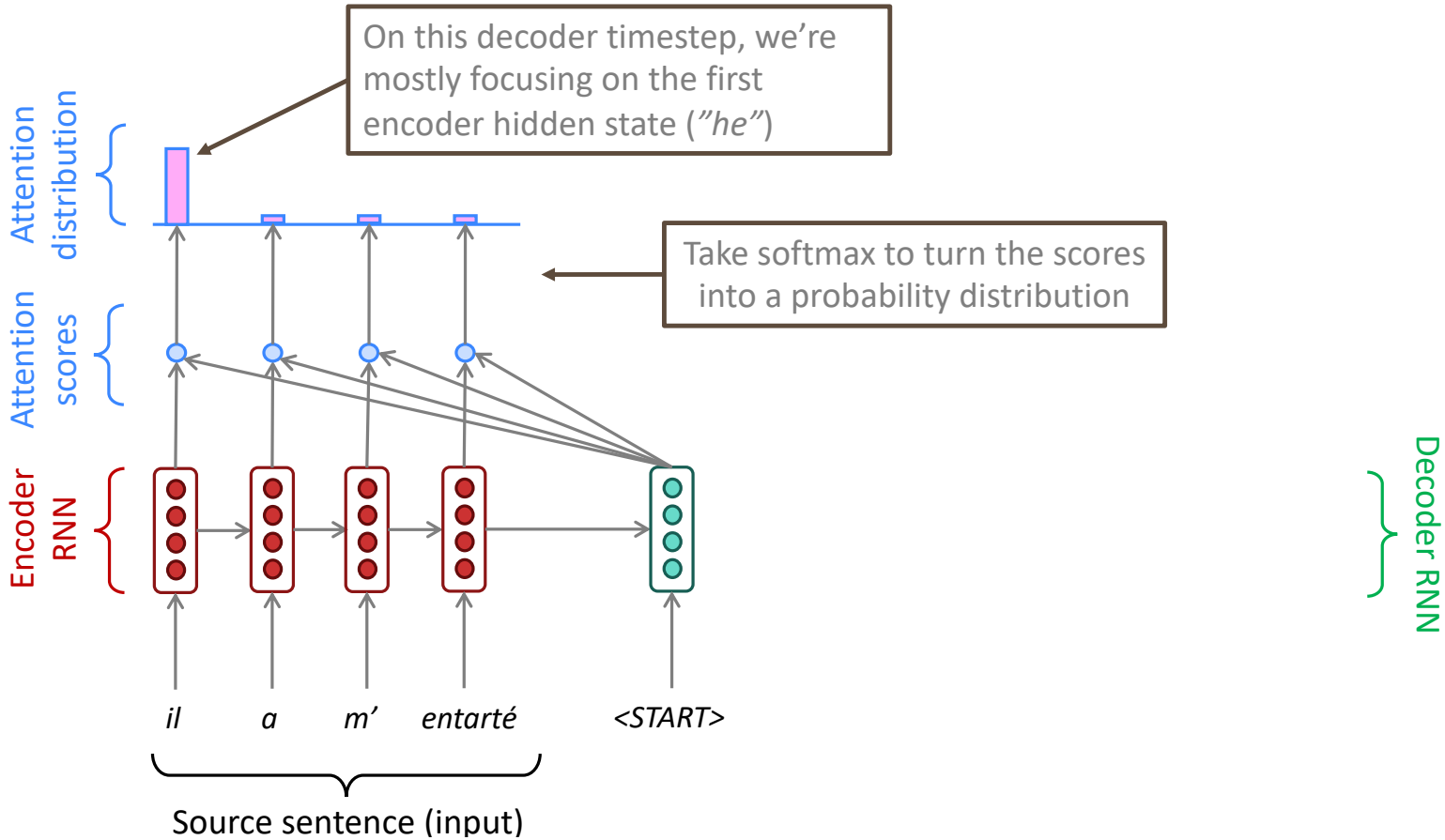
In a **lookup table**, we have a table of **keys** that map to **values**. The **query** matches one of the keys, returning its value.



# Using dot products

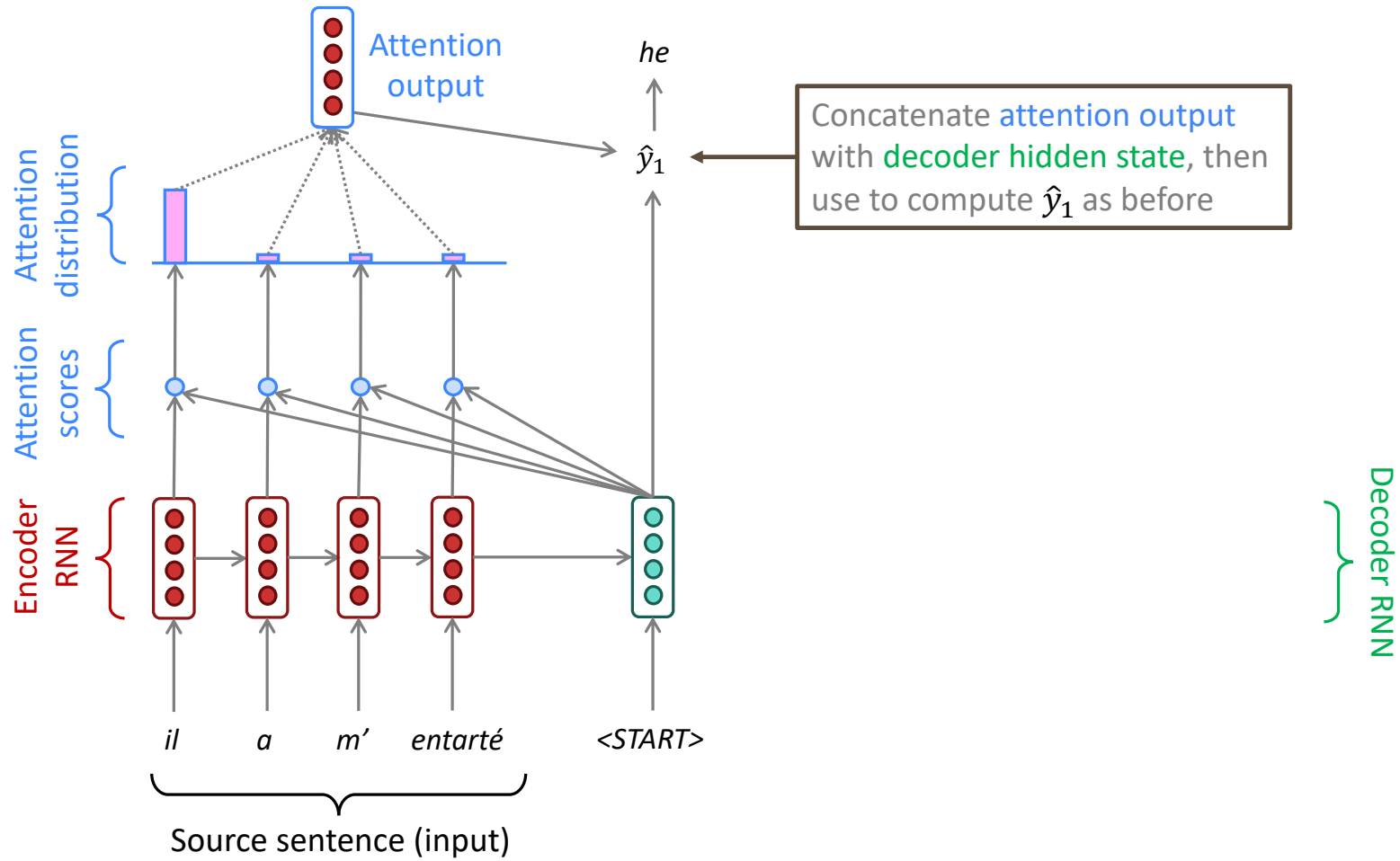


# Using softmax for aggregation





# Putting it all together



# Attention more formally

- We have encoder hidden states  $h_1, \dots, h_N \in \mathbb{R}^h$
- On timestep  $t$ , we have decoder hidden state  $s_t \in \mathbb{R}^h$
- We get the attention scores  $e^t$  for this step:

$$e^t = [s_t^T h_1, \dots, s_t^T h_N] \in \mathbb{R}^N$$

- We take softmax to get the attention distribution  $\alpha^t$  for this step (this is a probability distribution and sums to 1)

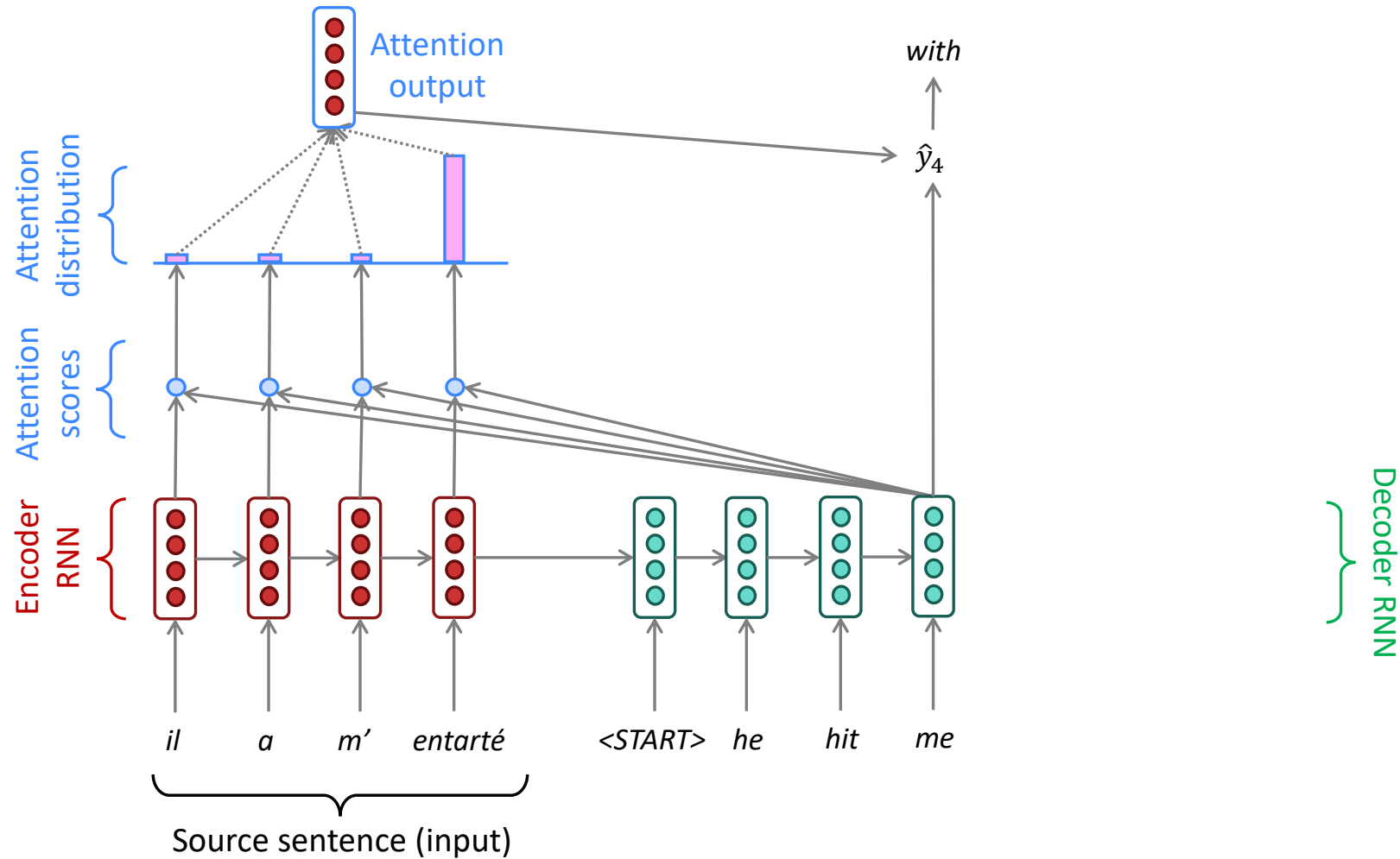
$$\alpha^t = \text{softmax}(e^t) \in \mathbb{R}^N$$

- We use  $\alpha^t$  to take a weighted sum of the encoder hidden states to get the attention output  $a_t$

$$a_t = \sum_{i=1}^N \alpha_i^t h_i \in \mathbb{R}^h$$

- Finally we concatenate the attention output  $a_t$  with the decoder hidden state  $s_t$  and proceed as in the non-attention seq2seq model

# Attention example



# Attention blueprint

- We have some *values*  $\mathbf{h}_1, \dots, \mathbf{h}_N \in \mathbb{R}^{d_1}$  and a *query*  $\mathbf{s} \in \mathbb{R}^{d_2}$

- Attention always involves:

1. Computing the *attention scores*  $\mathbf{e} \in \mathbb{R}^N$
2. Taking softmax to get *attention distribution*  $\alpha$ :

There are multiple ways to do this

$$\alpha = \text{softmax}(\mathbf{e}) \in \mathbb{R}^N$$

3. Using attention distribution to take weighted sum of values:

$$\mathbf{a} = \sum_{i=1}^N \alpha_i \mathbf{h}_i \in \mathbb{R}^{d_1}$$

thus obtaining the *attention output*  $\mathbf{a}$  (sometimes called the *context vector*)

# From translation to language generation: Self-attention

Let  $\mathbf{w}_{1:n}$  be a sequence of words in vocabulary  $V$ , like *Zuko made his uncle tea*.

For each  $\mathbf{w}_i$ , let  $\mathbf{x}_i = E\mathbf{w}_i$ , where  $E \in \mathbb{R}^{d \times |V|}$  is an embedding matrix.

1. Transform each word embedding with weight matrices  $Q, K, V$ , each in  $\mathbb{R}^{d \times d}$

$$\mathbf{q}_i = Q\mathbf{x}_i \text{ (queries)} \quad \mathbf{k}_i = K\mathbf{x}_i \text{ (keys)} \quad \mathbf{v}_i = V\mathbf{x}_i \text{ (values)}$$

2. Compute pairwise similarities between keys and queries; normalize with softmax

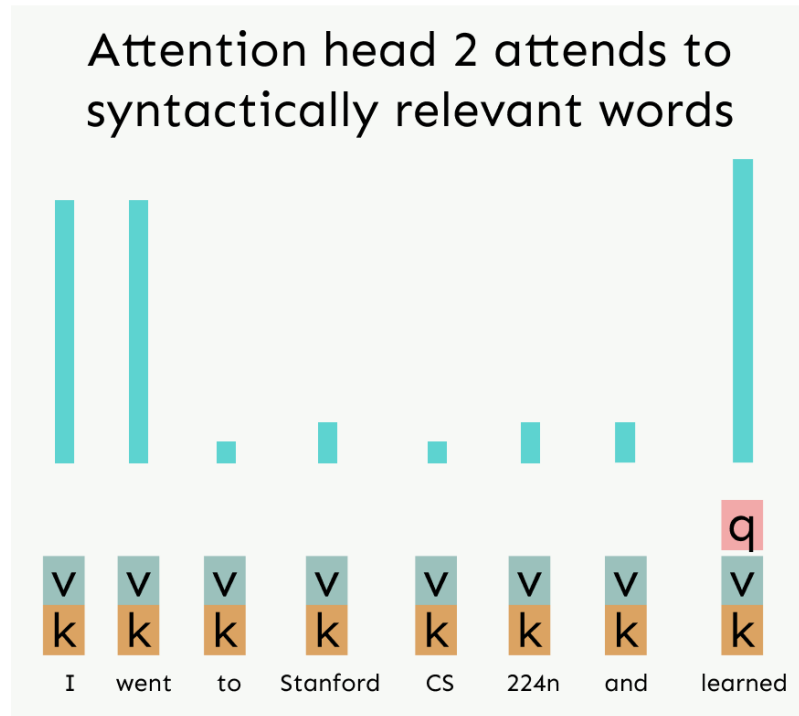
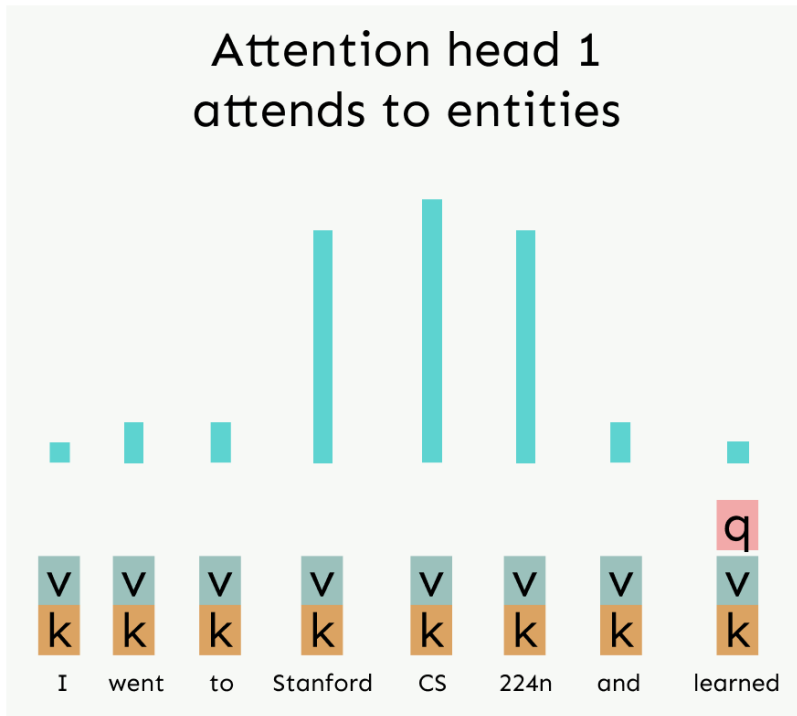
$$\mathbf{e}_{ij} = \mathbf{q}_i^\top \mathbf{k}_j \quad \alpha_{ij} = \frac{\exp(\mathbf{e}_{ij})}{\sum_{j'} \exp(\mathbf{e}_{ij'})}$$

3. Compute output for each word as weighted sum of values

$$\mathbf{o}_i = \sum_j \alpha_{ij} \mathbf{v}_i$$

$$\mathbf{o}_i = \sum_j \alpha_{ij} \mathbf{v}_i$$

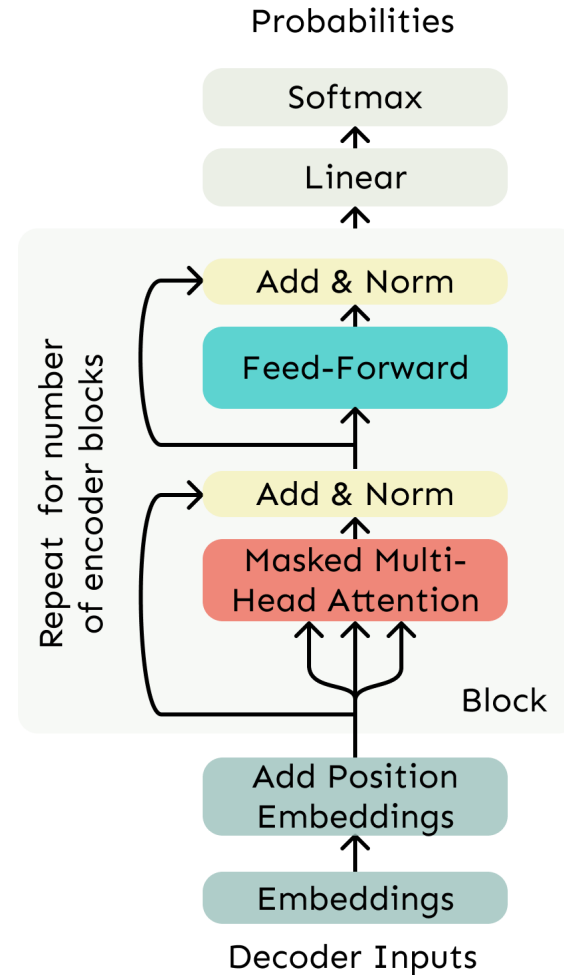
# Multihead attention



# Transformer decoder

## The Transformer Decoder

- The Transformer Decoder is a stack of Transformer Decoder **Blocks**.
- Each Block consists of:
  - Self-attention
  - Add & Norm
  - Feed-Forward
  - Add & Norm
- That's it! We've gone through the Transformer Decoder.

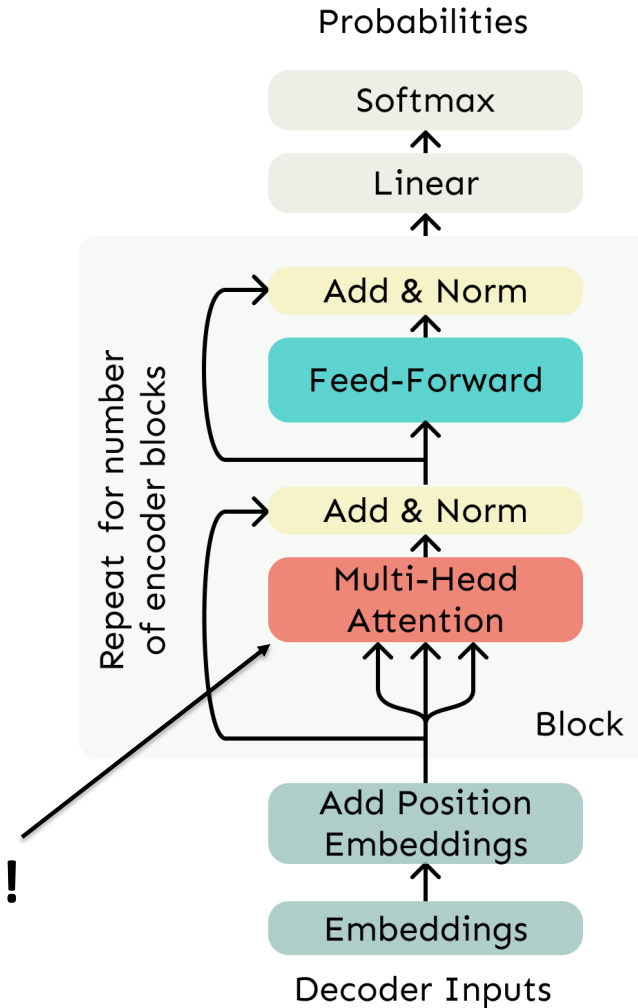


# Transformer encoder

## The Transformer Encoder

- The Transformer Decoder constrains to **unidirectional context**, as for **language models**.
- What if we want **bidirectional context**, like in a bidirectional RNN?
- This is the Transformer Encoder. The only difference is that we **remove the masking** in the self-attention.

**No Masking!**





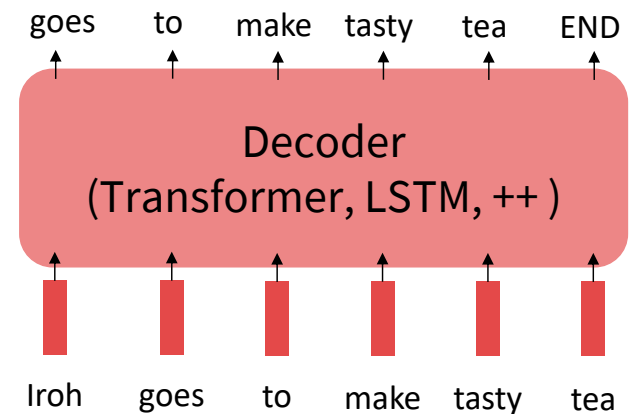
# Pretraining

Recall the **language modeling** task:

- Model  $p_{\theta}(w_t|w_{1:t-1})$ , the probability distribution over words given their past contexts.
- There's lots of data for this! (In English.)

**Pretraining through language modeling:**

- Train a neural network to perform language modeling on a large amount of text.
- Save the network parameters.

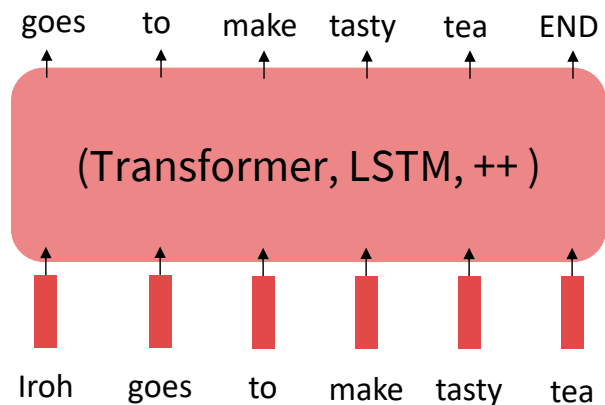


# Pretraining / finetuning paradigm

Pretraining can improve NLP applications by serving as parameter initialization.

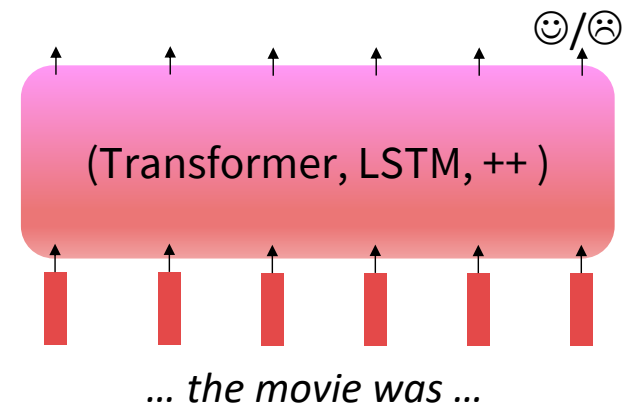
## Step 1: Pretrain (on language modeling)

Lots of text; learn general things!



## Step 2: Finetune (on your task)

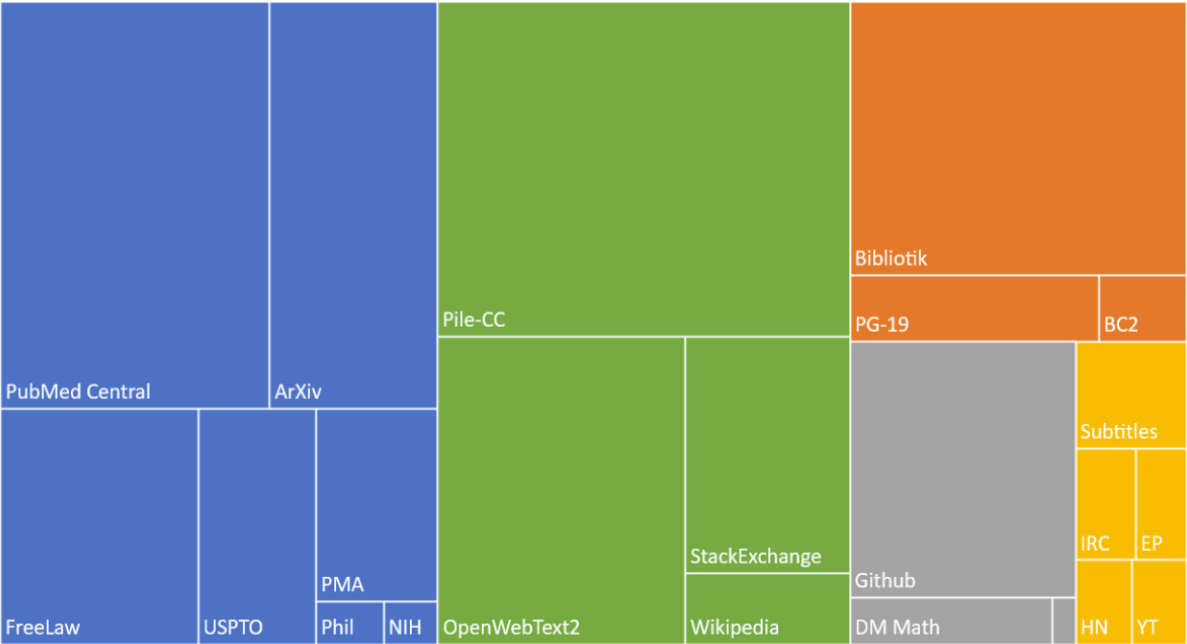
Not many labels; adapt to the task!



# What data to use?

Composition of the Pile by Category

■ Academic ■ Internet ■ Prose ■ Dialogue ■ Misc



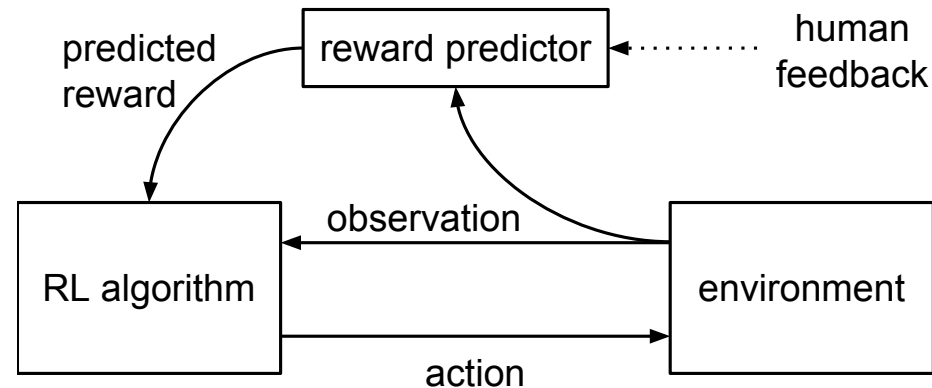
Model	Training Data
BERT	BookCorpus, English Wikipedia
GPT-1	BookCorpus
GPT-3	CommonCrawl, WebText, English Wikipedia, and 2 book databases (“Books 1” and “Books 2”)
GPT-3.5+	Undisclosed

# GPT (Devlin et al, 2018)

2018's GPT was a big success in pretraining a decoder!

- Transformer decoder with 12 layers, 117M parameters.
- 768-dimensional hidden states, 3072-dimensional feed-forward hidden layers.
- Byte-pair encoding with 40,000 merges
- Trained on BooksCorpus: over 7000 unique books.
  - Contains long spans of contiguous text, for learning long-distance dependencies.
- The acronym “GPT” never showed up in the original paper; it could stand for “Generative PreTraining” or “Generative Pretrained Transformer”

# Deep RL from Human Feedback (Christiano et al, 2017)



- People provide a *preference* among two choices
- Assuming there is a latent variable explaining the choice, reward is fit using maximum likelihood (Bradley-Terry model)
- Cf. <https://arxiv.org/pdf/1706.03741.pdf>

## Bradely-Terry reward model

- Collect data from human raters (pairs of  $y_w, y_l$  responses to a prompt  $x$ )
- Optimize the expected value of:

$$-\log(\sigma(r_\theta(x, y_w) - r_\theta(x, y_l)))$$

wrt reward parameter vector  $\theta$

- Cf. Ouyang et al, InstructGPT
- Corresponds to maximum likelihood fitting of binomial preference function if reward is linear over the variables

## Optimizing Preferences: Setup

- An agent interacting with an environment receives observations for a set  $\mathcal{O}$  and performs action from set  $\mathcal{A}$
- A *history*  $h_t$  is a sequence of observation-action pairs  $\langle o_0, a_0, o_1, a_1, \dots, o_t \rangle$
- A *policy*  $\pi$  is a mapping from histories to actions:  $\pi : \mathcal{H} \rightarrow \mathcal{A}$
- Consider a *binary relation over trajectory distributions*  $\preceq$
- A policy  $\pi$  in an environment  $e$  induces a probability distribution over trajectories,  $D^\pi$
- See Colaco-Carr et al, AISTATS'2024 (<https://arxiv.org/abs/2311.01990>)

# Preference Relations and Their Properties

- We will formalize preference relations through pre-orders
- For trajectory distributions  $A$  and  $B$ ,  $A \preceq B$  means is that  $B$  is at least as preferred as  $A$
- $\preceq$  is a *pre-order* if it satisfies:
  - *Reflexivity*:  $A \preceq A$
  - *Transitivity*: if  $A \preceq B$  and  $B \preceq C$  the  $A \preceq C$
- A pre-order is *total* if for and  $A, B$ ,  $A \preceq B$  and  $B \preceq A$



# Direct Preference Process

- A *Direct Preference Process* is a tuple  $\mathcal{O}, \mathcal{A}, T, e, \preceq$  where:
  - $\mathcal{O}$  is an observation set
  - $\mathcal{A}$  is an action set
  - $T$  is a time horizon
  - $e$  is an environment (transition function from achievable history-action pairs to the next observation)
  - $\preceq$  is a binary (preference) relation over trajectory distributions
- $\preceq$  is *expressible through a reward function*  $r : \mathcal{H} \rightarrow \mathbb{R}$  if:

$$\forall A, B, A \preceq B \text{ if and only if } \mathbb{E}_A \left[ \sum_{t=0}^T r(H_t) \right] \leq \mathbb{E}_B \left[ \sum_{t=0}^T r(H_t) \right]$$

# Preference Relations and Their Properties

- A total pre-order is *consistent* if

$$\forall \alpha \in (0, 1), \forall A, B, C, A \preceq B \implies \alpha A + (1 - \alpha)C \preceq \alpha B + (1 - \alpha)C$$

- A total pre-order is *convex* if

$$\forall \alpha \in (0, 1), \forall A, B, C, A \preceq B. \text{ if and only if } \alpha A + (1 - \alpha)C \preceq \alpha B + (1 - \alpha)C$$

- A total pre-order has the *interpolation property* if

$$\forall A, B, C, A \preceq B \text{ and } B \preceq C \text{ implies } \exists \alpha \in (0, 1), \alpha A + (1 - \alpha)C \sim B$$

- Von Neumann-Morgenstern theorem: if all the above hold,  $\preceq$  can be expressed by a utility function

# When Are Preferences Representable By Reward Functions?

- Main result
  - *If convexity and/or interpolation do not hold,  $\preceq$  is NOT expressible through a reward function*
  - *However, total consistent pre-orders have deterministic optimal policy!*
- The latter situation is not exotic or rare!

## Examples when Optimal Policies Exist Without Rewards

- *Total consistent convex pre-order not satisfying interpolation: tie-breaking criteria*
  - Use a first criterion, if tied go to a second criterion
  - See not flooding vs water in second reservoir in power plant example
- *Total consistent pre-order that is non-convex: excess risk*
  - If risky event does not occur, linear utility
  - Risky event occurring entails exponential penalty
  - No flooding neighbouring areas in power plant example

## How Do We Compute Optimal Policies?

- If  $\preceq$  is a total consistent pre-order and a policy  $\pi$  satisfies the following for any attainable history  $h_t, t < T$  and any action  $a_t$ :

$$D^\pi(h_t \cdot a_t) \preceq D^\pi(h_t)$$

then  $\pi$  is  $\preceq$ -optimal

- So we are *justified to do policy search!*
- If  $\preceq$  is expressible through a reward function, value iteration is a direct consequence of this result

## Discussion

- Nice to know that approaches such as direct preference optimization are justified
- Our results are currently on distributions - working on sample-based extensions
- If we can fit a reward function, should we?
  - Bias-variance trade-off? Sample complexity considerations?
- *What can we do if other properties of pre-orders are violated?*

# Learning with non-transitive preferences: NashLLM

- Objective: find a policy  $\pi^*$  which is preferred over any other policy

$$\pi^* = \arg \max_{\pi} \min_{\pi'} \mathbb{P}(\pi' \preceq \pi)$$

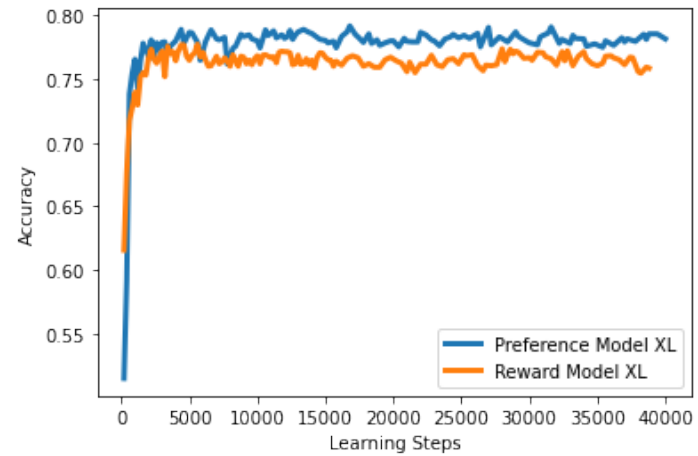
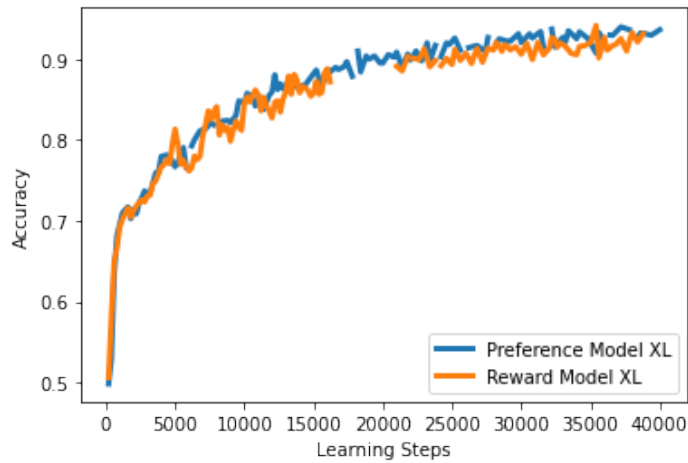
- Think of this as a game: one player picks  $\pi$  the other picks  $\pi'$
- When both players use  $\pi^*$  this is a *Nash equilibrium* for the game
- For this game an equilibrium exists (even if eg preferences are not transitive)
- Cf. Munos et al, 2024 (<https://arxiv.org/pdf/2312.00886.pdf>)

# NashLLM-style algorithms

- Fit a *two-argument preference function* by supervised learning
- Decide what is the *set of opponent policies*
- Ideally, the max player should play against a mixture of past policies
- *Optimize* using eg online mirror descent, convex-concave optimization...
- A lot of algorithmic variations to explore!



## NashLLM results



Using preferences instead of rewards leads to less overfitting