# CS 489/698: Introduction to Natural Language Processing

# Lecture 6.2: Language Modeling II

Instructor: Freda Shi

*fhs@uwaterloo.ca*

*February 4th, 2026*

UNIVERSITY OF
**WATERLOO**

# Outline: Neural Language Modeling

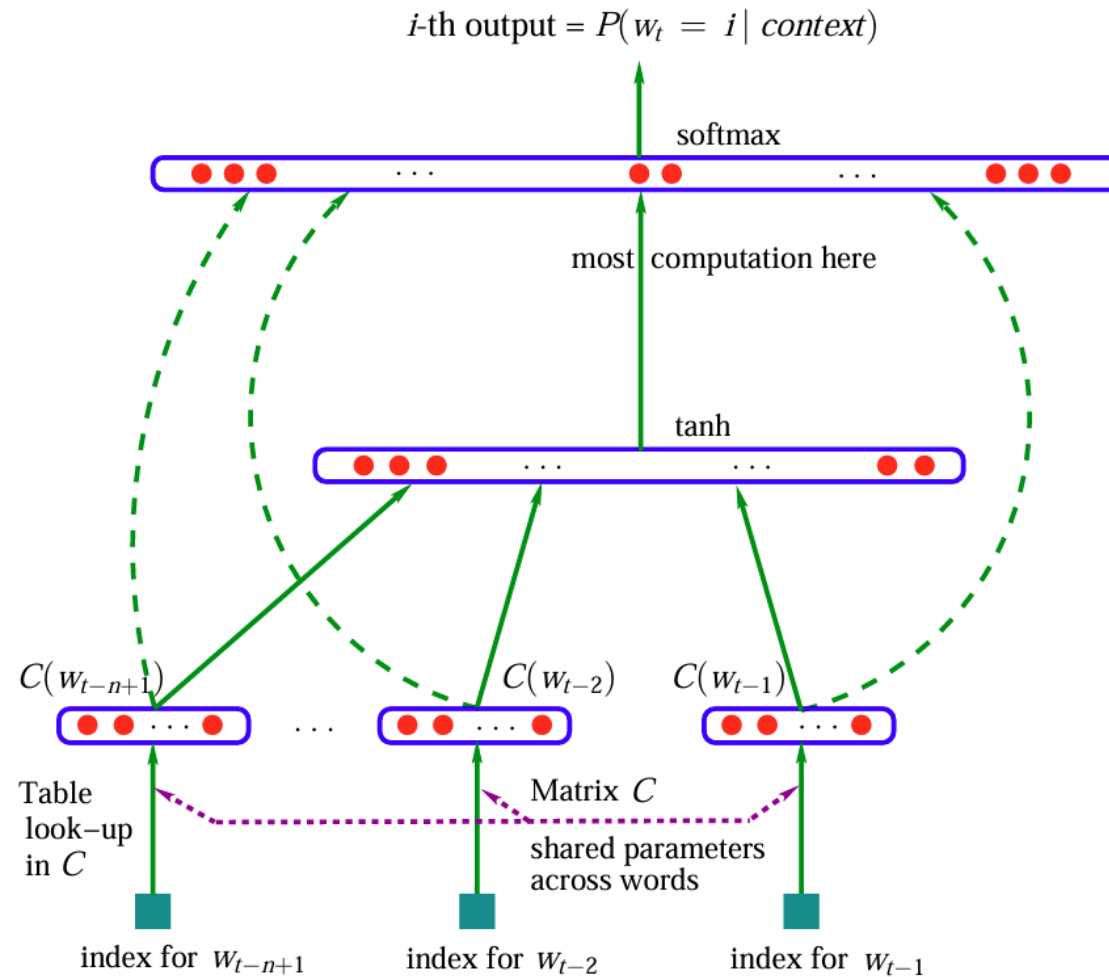- Neural language modeling

- Probing neural language models

# Recap: Language Modeling as Classification

$$P(\mathbf{x}_{1:n}) = P(x_1, x_2, \ldots, x_n) = P(\texttt{</s>} \mid x_1, \ldots, x_n) \prod_{i=1}^{n} \boxed{P(x_i \mid x_1, x_2, \ldots, x_{i-1})}$$

This is just a probabilistic classification problem!

We can use any tools from the previous lectures: linear model with features, neural networks, etc.

# Neural N-Gram Language Models

$i$-th output $= P(w_t = i \mid context)$

softmax

most computation here

tanh

$C(w_{t-n+1})$ $C(w_{t-2})$ $C(w_{t-1})$

Table look-up in $C$

Matrix $C$ shared parameters across words

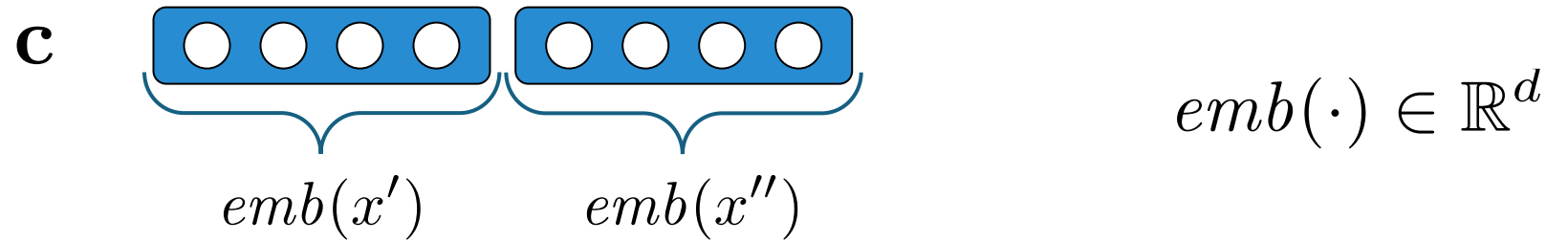index for $w_{t-n+1}$ index for $w_{t-2}$ index for $w_{t-1}$

[Src: Bengio et al., 2003]

# Neural Trigram Language Model

Given two previous words, compute probability distribution over possible next words
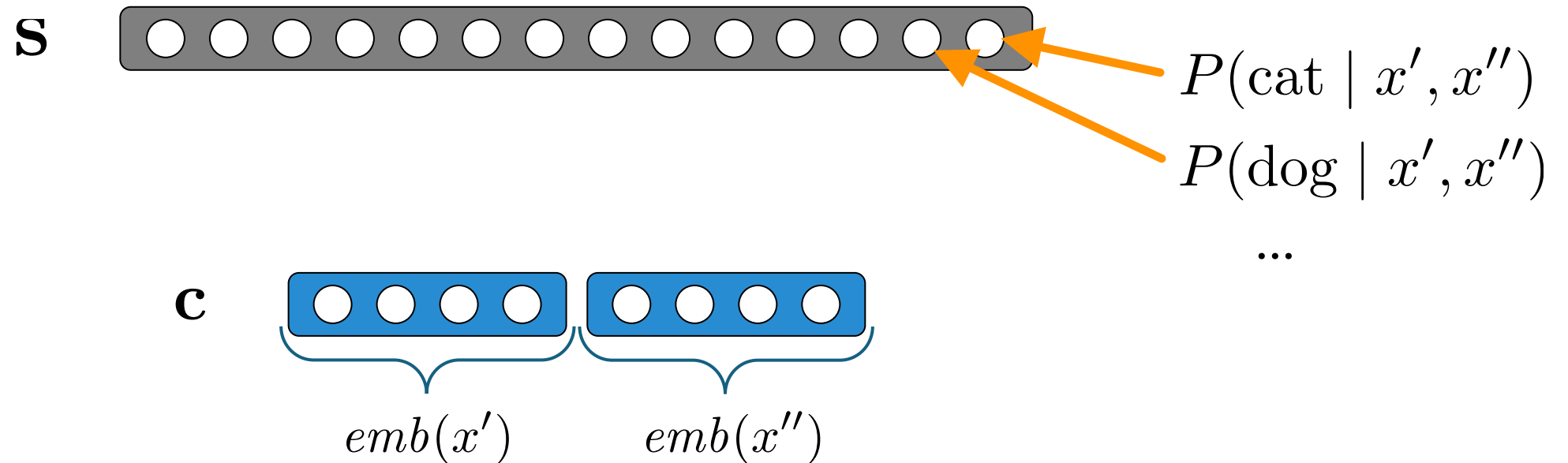
$$P(x \mid x', x'')$$

Input is concatenation of vectors (embeddings) of previous two words:

$$emb(\cdot) \in \mathbb{R}^d$$

$$\mathbf{c} = cat(emb(x'), emb(x''))$$

# Neural Trigram Language Model

Output is a vector $\mathbf{S}$ containing probabilities of all possible next words:

$P(\text{cat} \mid x', x'')$

$P(\text{dog} \mid x', x'')$

...

$emb(x')$  $emb(x'')$

# Neural Trigram Language Model

To get $\mathbf{S}$, do matrix multiplication of parameter matrix $\mathbf{W}$ and input, then "softmax" transformation

$$\mathbf{s} = \mathrm{softmax}(\mathbf{Wc})$$

$$\mathbf{W} \in \mathbb{R}^{|V| \times 2d} \quad \mathbf{c} \in \mathbb{R}^{2d}$$

$\mathbf{S}$

$\mathbf{c}$

"fully-connected layer"

$emb(x')$    $emb(x'')$

# Neural Trigram Language Model

$$P_\Theta(x_t \mid x_1, \dots, x_{t-1}) = s_{x_t}$$

$$\text{loss } L = -\sum_i \log P_\Theta(\mathbf{x}^{(i)})$$

$$= -\sum_i \sum_j \log P_\Theta(x_j^{(i)} \mid x_1^{(i)}, \dots, x_{j-1}^{(i)})$$

$\nabla_\Theta L$ via backpropagation.

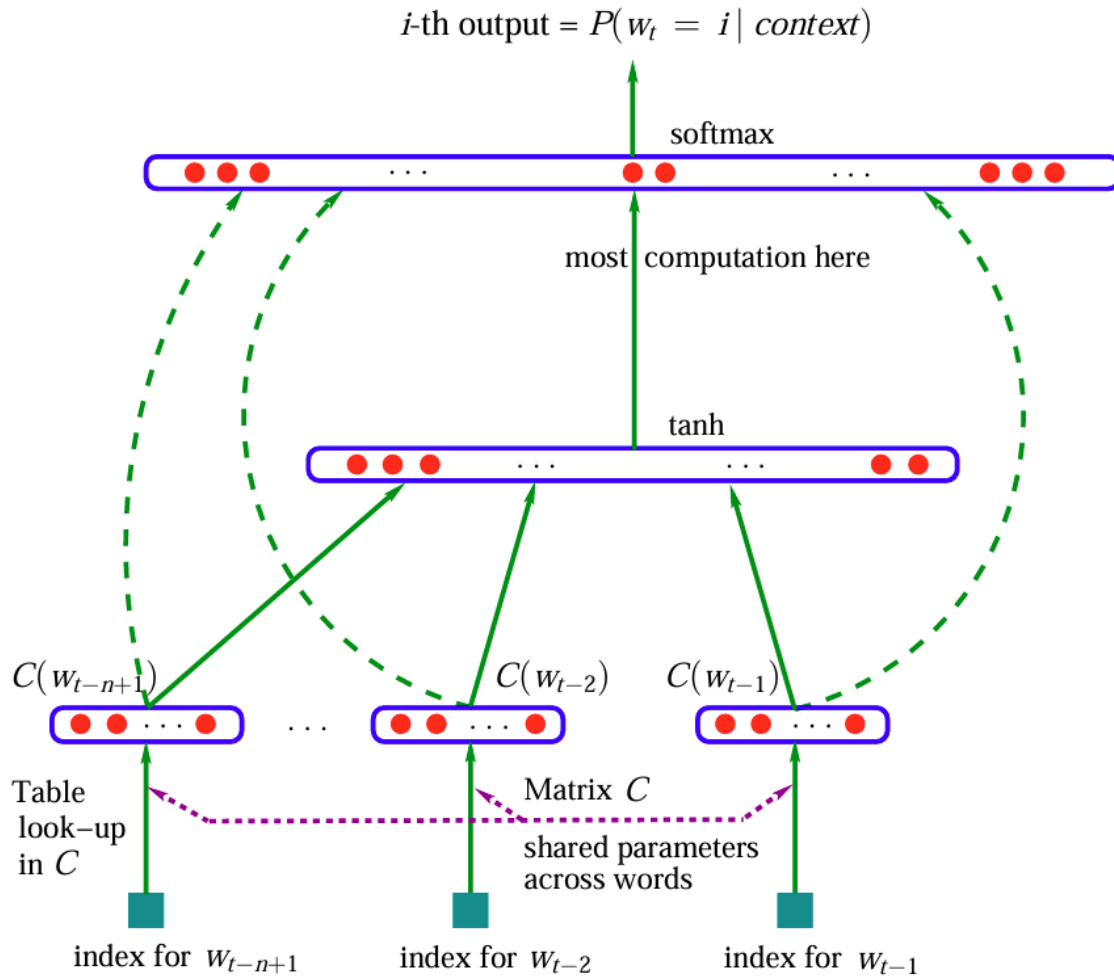# Trigram vs. Neural Trigram LMs $P(x \mid x', x'')$

- **Trigram language model**
  - separate parameters for every combination of $x, x', x''$
  - so, approx. $|V|^3$ parameters
  - \# parameters is exponential in $n$-gram size
  - most parameters are zero
  - even with smoothing, many parameters can remain zero

- **Neural trigram language model**
  - only has $\mathcal{O}(d|V|)$ parameters
  - $d$ can be chosen to scale \# parameters up or down
  - \# parameters linear in $n$-gram size
  - (almost) no parameters are zero
  - no explicit smoothing, though smoothing done implicitly via distributed representations

# Removing N-Gram Constraints

$$P(x_i \mid x_1, x_2, \ldots, x_{i-1})$$



*i*-th output = $P(w_t = i \mid context)$

softmax

most computation here

tanh

$C(w_{t-n+1})$    $C(w_{t-2})$    $C(w_{t-1})$

Table look–up in $C$

Matrix $C$

shared parameters across words

index for $w_{t-n+1}$    index for $w_{t-2}$    index for $w_{t-1}$

[Src: Bengio et al., 2003]

$$x_i \xrightarrow{\text{look-up in } \mathbf{U}} \mathbf{w}_i \in \mathbb{R}^d$$

$$x_1, x_2, \ldots, x_{i-1} \to \frac{1}{i-1} \sum_{j=1}^{i-1} \mathbf{w}_j \in \mathbb{R}^d$$
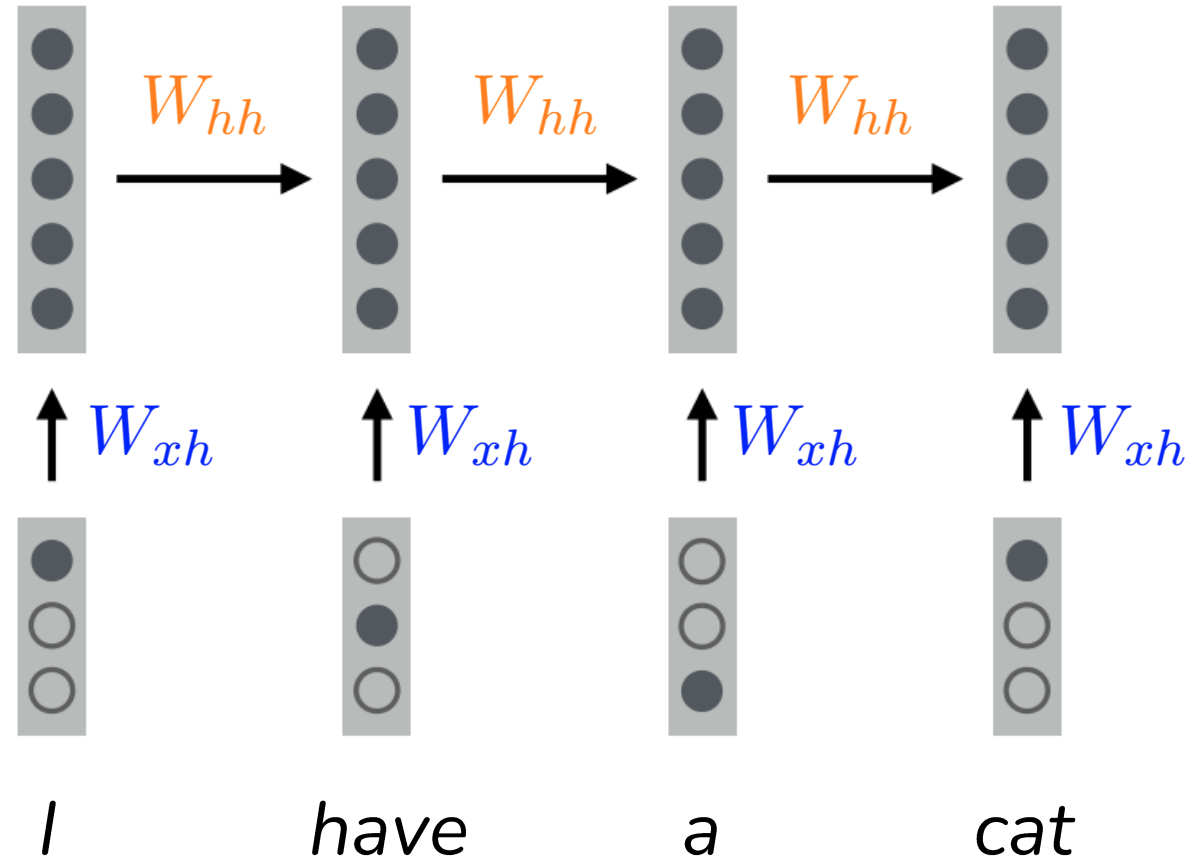
$$P(x_i = v \mid x_1, x_2, \ldots, x_{i-1})$$

$$= \text{softmax}_v \, \mathbf{W} \left( \frac{1}{i-1} \sum_{j=1}^{i-1} \mathbf{w}_j \right)$$

$$\mathbf{W} \in \mathbb{R}^{|V| \times d}$$

# RNN Language Models

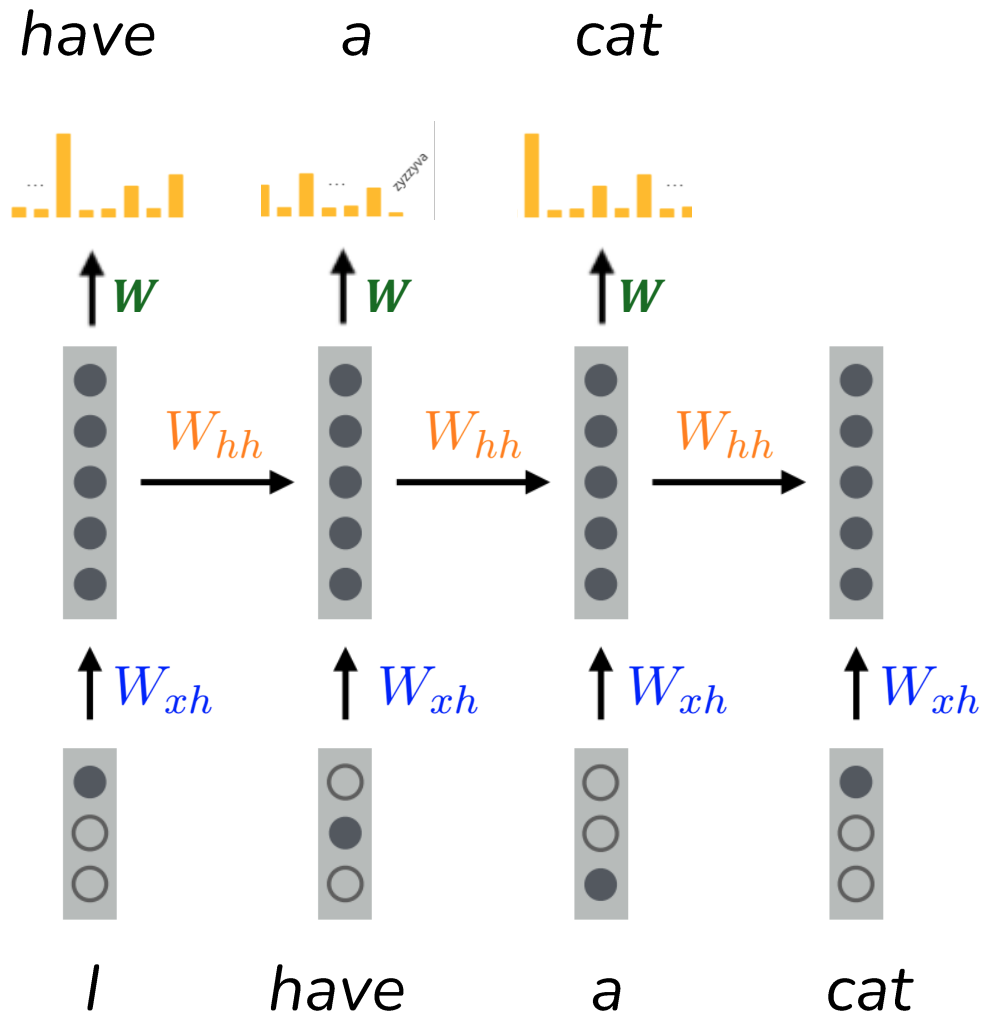Hidden state is a function of previous hidden state and current input.

Same weights at each state.



$W_{hh}$     $W_{hh}$     $W_{hh}$

$W_{xh}$     $W_{xh}$     $W_{xh}$     $W_{xh}$

*I*     *have*     *a*     cat

$$\mathbf{h}_t = f(\mathbf{W}_{hh}\mathbf{h}_{t-1} + \mathbf{W}_{xh}\mathbf{x}_t)$$

$$P_\Theta(x_t \mid x_1, \ldots, x_{t-1}) = \text{softmax}_{x_t}(\mathbf{W}\mathbf{h}_{t-1})$$

# RNN Language Model



$$L_t = -\log P_\Theta(x_t \mid x_1, \ldots, x_{t-1})$$
$$\Theta = \{\mathbf{U}, \mathbf{W}_{xh}, \mathbf{W}_{hh}\}$$

$$L = \sum_t L_t$$
$$= \sum_t -\log P_\Theta(x_t \mid x_1, \ldots, x_{t-1})$$
$$= -\log P_\Theta(\mathbf{x})$$

$\nabla_\Theta L$ via backpropagation.

# Transformer Language Models

A token "attends" to all previous tokens.

$$\mathbf{E} = (emb(w_1), \ldots, emb(w_k)) + \mathbf{P} \in \mathbb{R}^{d_1 \times k}$$

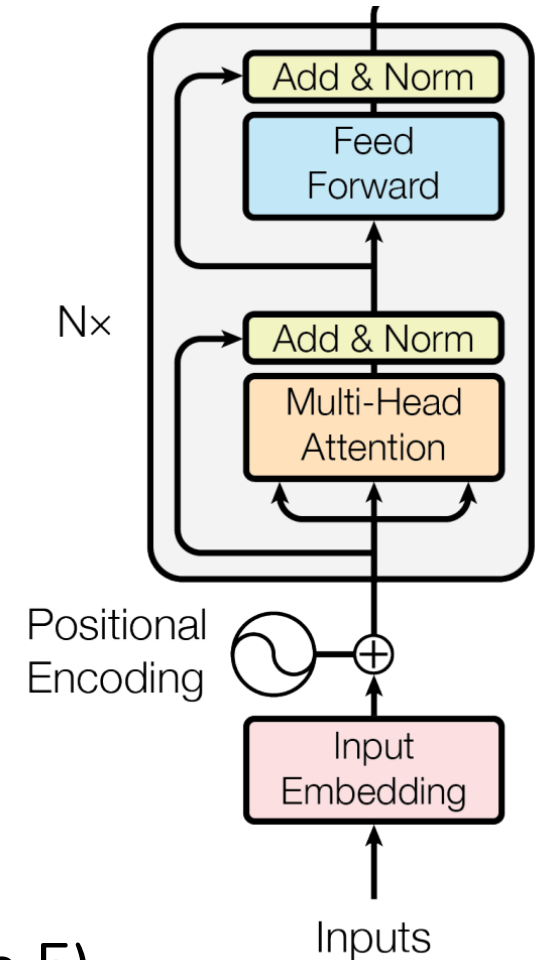$$\mathbf{K} = \mathbf{W}_k \mathbf{E} \quad \mathbf{W}_k \in \mathbb{R}^{d_2 \times d_1}$$

$$\mathbf{Q} = \mathbf{W}_q \mathbf{E} \quad \mathbf{W}_q \in \mathbb{R}^{d_2 \times d_1}$$

$$\mathbf{V} = \mathbf{W}_v \mathbf{E} \quad \mathbf{W}_v \in \mathbb{R}^{d_3 \times d_1}$$
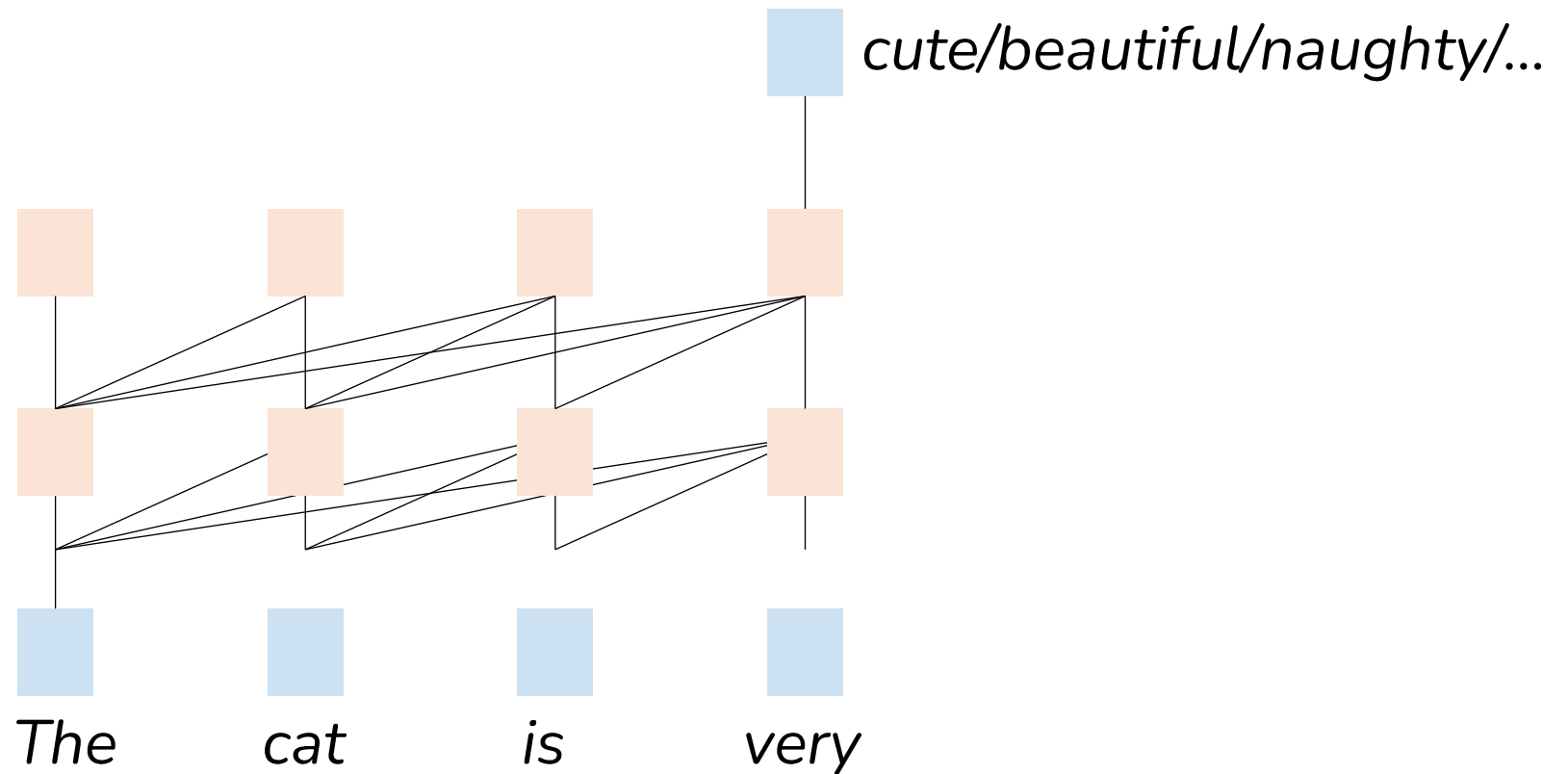
$$\tilde{\mathbf{E}} = \mathbf{V}\mathrm{softmax}\left(\frac{\mathbf{K}^\top \mathbf{Q}}{\sqrt{d_2}}\right) \in \mathbb{R}^{d_3 \times k}$$



Use $\tilde{\mathbf{E}}_{\cdot,k}$ as feature to predict the next token.

Note: feature is more complicated in real practice (see Lecture 5).
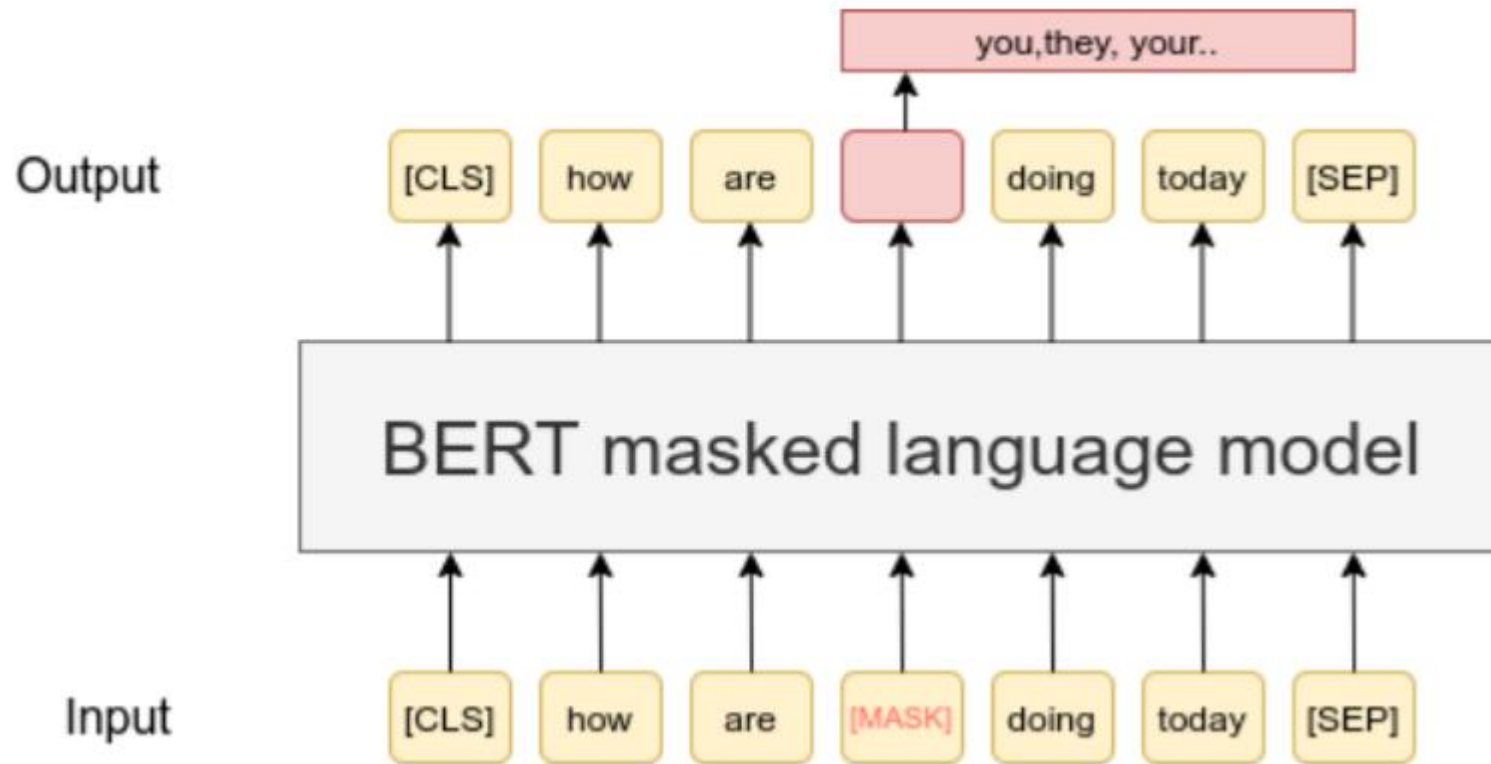
# Language models encode knowledge about language



*cute/beautiful/naughty/...*

*The*     *cat*     *is*     *very*

The pretraining-finetuning paradigm: Language modeling, as the pretraining task, helps encode knowledge.

The knowledge helps downstream tasks.

# Masked Language Models

Motivation: learning useful representations of text.

# Mased Language Models

A token "attends" to all context tokens.

$$\mathbf{E} = (emb(w_1), \ldots, emb(w_k)) + \mathbf{P} \in \mathbb{R}^{d_1 \times k}$$

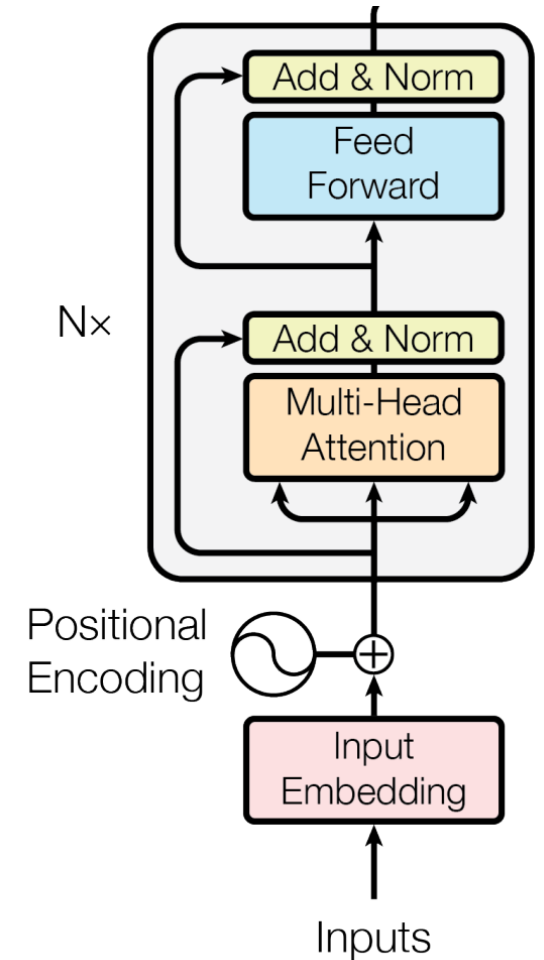$$\mathbf{K} = \mathbf{W}_k \mathbf{E} \quad \mathbf{W}_k \in \mathbb{R}^{d_2 \times d_1}$$

$$\mathbf{Q} = \mathbf{W}_q \mathbf{E} \quad \mathbf{W}_q \in \mathbb{R}^{d_2 \times d_1}$$

$$\mathbf{V} = \mathbf{W}_v \mathbf{E} \quad \mathbf{W}_v \in \mathbb{R}^{d_3 \times d_1}$$

$$\tilde{\mathbf{E}} = \mathbf{V}\mathrm{softmax}\left(\frac{\mathbf{K}^\top \mathbf{Q}}{\sqrt{d_2}}\right) \in \mathbb{R}^{d_3 \times k}$$

Replace token at position $i$ with a placeholder [MASK].

Use $\tilde{\mathbf{E}}_{\cdot,i}$ as feature to predict token at position $i$.



Add & Norm

Feed Forward

N×

Add & Norm

Multi-Head Attention

Positional Encoding

Input Embedding

Inputs

# Some Important Details of LMs

- The importance of the held-out data.

- AdamW (Kingma & Ba, 2015; Loshchilov & Hutter) has become the go-to optimizer instead of vanilla gradient descent.

- $1 \times 10^{-4}$ could be a default learning rate for AdamW.

- Monitor your training loss through time (by printing it out or using loggers like weights & biases; https://wandb.ai/site/).

Check out the HuggingFace Tutorial on training language models:
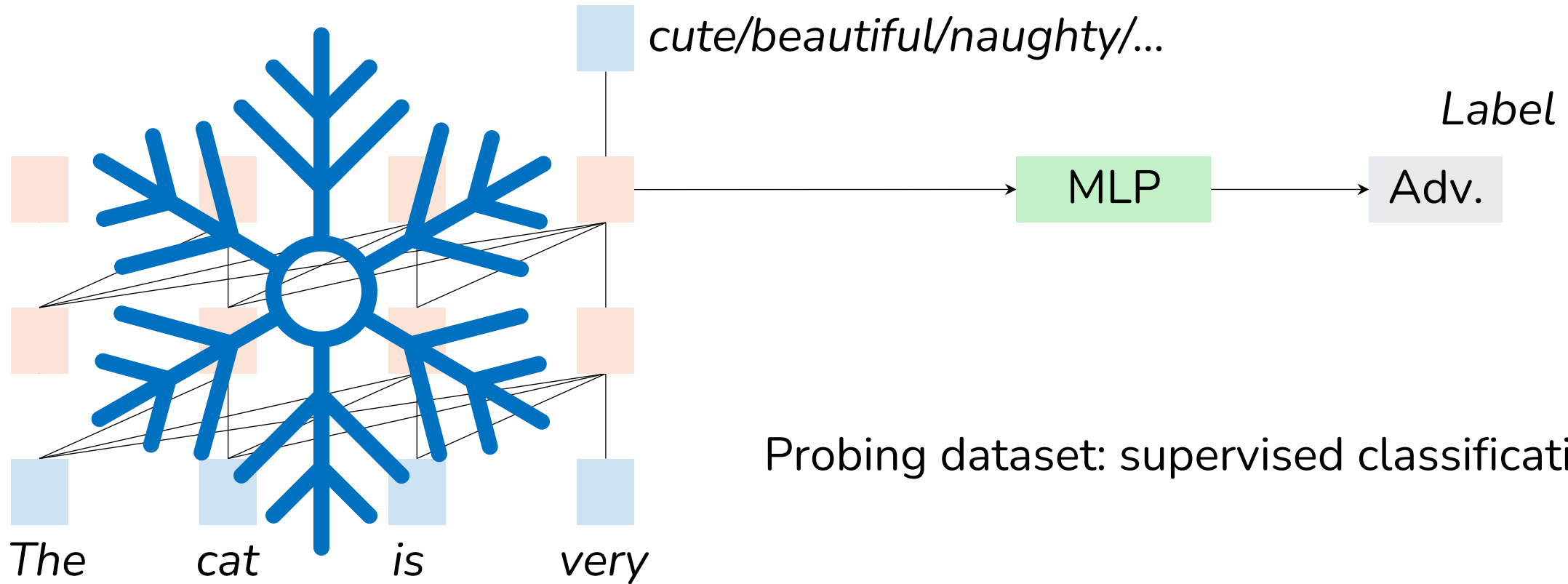https://huggingface.co/learn/llm-course/en/chapter7/6

# Probing

What is encoded in a trained neural language model?

Empirical answer: linguistic probe (Ettinger et al., 2016).

Take a fixed model as the "frozen" feature extractor, train a lightweight model (probe, usually linear model or MLP) to predict labels.

Frozen: the base model never gets updated when training the lightweight model.
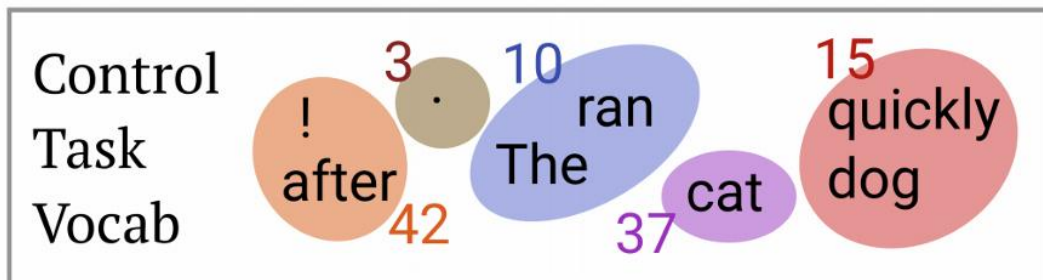
# Probing Syntax (Part-of-Speech Tags)



*cute/beautiful/naughty/...*

*Label*

MLP → Adv.

Probing dataset: supervised classification.

*The*  *cat*  *is*  *very*

# Confounding

Q: Does above-chance performance on held out data mean the model encodes part-of-speech knowledge?

A: Not necessarily – the model might just encodes word identity, and the probe learns to group them together.

Solution (control tasks; Hewitt and Liang, 2019): draw conclusion iff. performance on main task is significantly better than that on control task.



| Sentence 1 | The | cat | ran | quickly | . |
|---|---|---|---|---|---|
| **Part-of-speech** | DT | NN | VBD | RB | . |
| **Control task** | 10 | 37 | 10 | 15 | 3 |
| Sentence 2 | The | dog | ran | after | ! |
| **Part-of-speech** | DT | NN | VBD | IN | . |
| **Control task** | 10 | 15 | 10 | 42 | 42 |

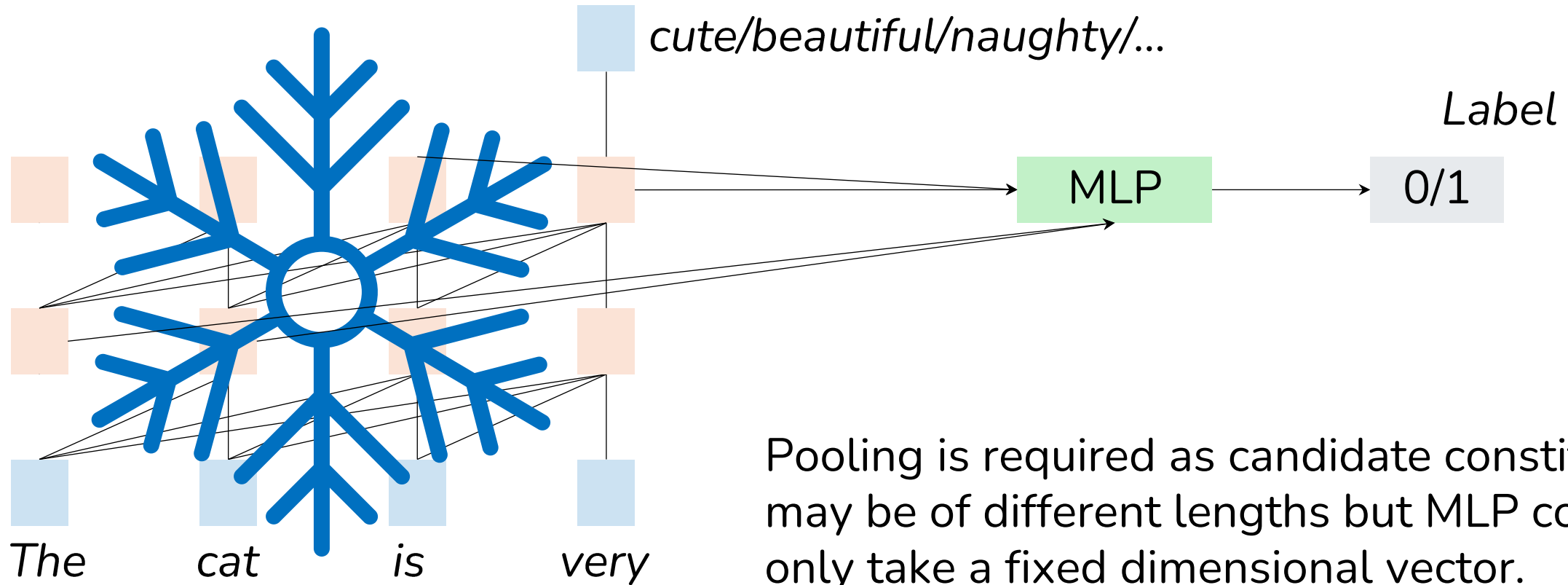# Syntax: Constituency

Sentence: *the cat is cute*

Bracketing: *((the cat) (is cute))*

Tree:



the    cat    is    cute

Task: given any span of words -- is it a constituent?

# Probing Syntax (Constituency)



*cute/beautiful/naughty/...*

*Label*

MLP → 0/1

*The     cat     is     very*

Pooling is required as candidate constituents may be of different lengths but MLP could only take a fixed dimensional vector.
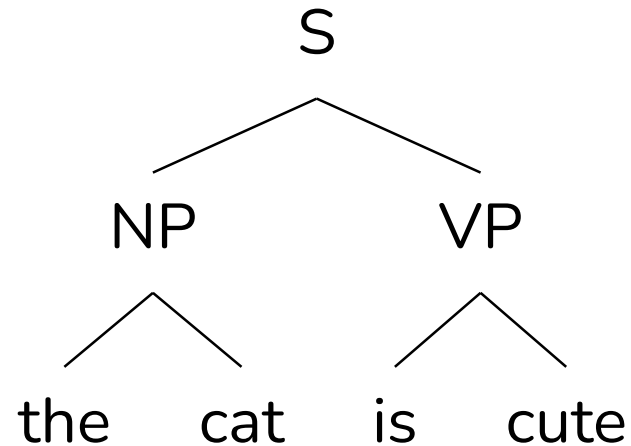
# Syntax: Constituent Labels

Sentence: *the cat is cute*

Bracketing: *((the cat) (is cute))*

Tree:

```
                    S
                 /     \
              NP         VP
             /  \       /  \
          the   cat   is   cute
```

Task: given a constituent, what's the label?

# Constituent Labels: Syntactic Substitutability

S

NP

VP

the man
John
the first boss I ever had

walked to the park
fell asleep
is here

# Probing Syntax (Constituency Labels)



*cute/beautiful/naughty/...*

*Label*

MLP → NP/VP/...

Pooling is required to accommodate the variable length of constituents.

*The     cat     is     very*

# More Tasks

- Does an LM encode numeracy?

- Does an LM "know" if two words refer to the same identity (coreference resolution)?

- Does an LM encode sentence-level information, e.g., topic?

- Does an LM encode knowledge about common sense?

…

BERTology: studying the internals of the BERT model.

# Rethinking Probing: Any Issues?

Take a fixed model as the "frozen" feature extractor, train a lightweight model (probe, usually linear model or MLP) to predict labels.

A poor performance may come from:

- The information is not encoded in the model under investigation.

- The information is encoded but not effectively used by the probe.

- The dataset for probe training does not effectively represent the information.

(Belinkov, 2022)

# Next

- Basic Syntax: What forms a constituent?

- Context-Free Grammars and Probabilistic Context-Free Grammars