# Lecture 10:
# Natural Language Generation

Instructor: Xiang Ren

USC CSCI 444 NLP

2026 Spring

# Logistics / Announcements

- questions re: HW1 – TA office hour

- Project Proposal feedback sent out by Tue EOD

- HW2 due March 4

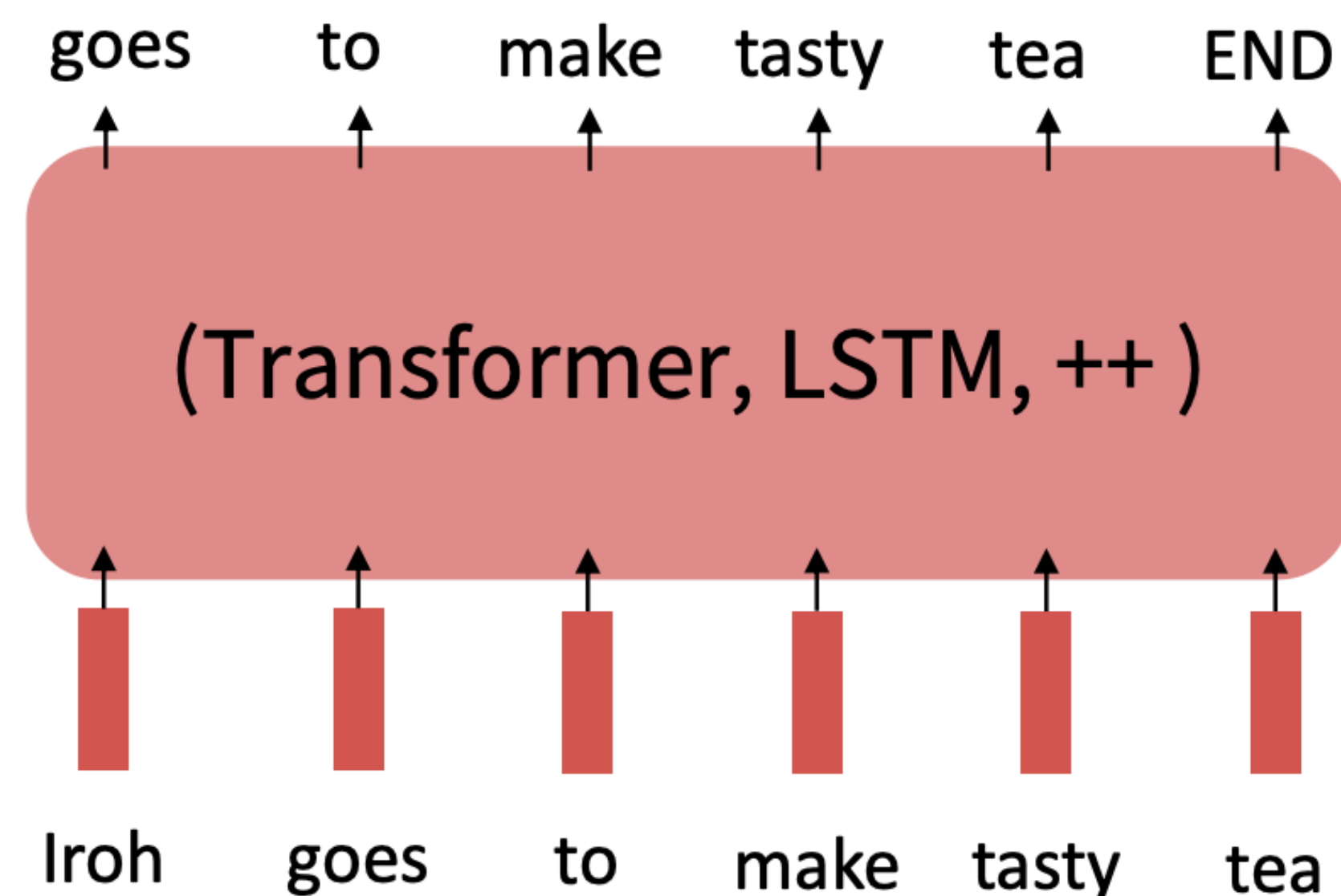| Feb 11 | Recurrent Neural Nets | J&M, Chap 13; | Project Proposal Due |
|--------|----------------------|---------------|----------------------|
| Feb 16 | Presidents Day | | |
| Feb 18 | Seq2Seq and Attention | J&M, Chap 8; | |
| Feb 23 | Transformers - Building Blocks | J&M, Chap 8; | |
| Feb 25 | PyTorch for Transformers | | |
| Mar 2 | Transformer Language Models | J&M Chap 8; | |
| Mar 4 | Tokenization | J&M, Chap 2.5; | HW2 Due |

USC Viterbi

# Recap: LLM Pretraining

# The Pretraining / Finetuning Paradigm

- Pretraining can improve NLP applications by serving as parameter initialization.

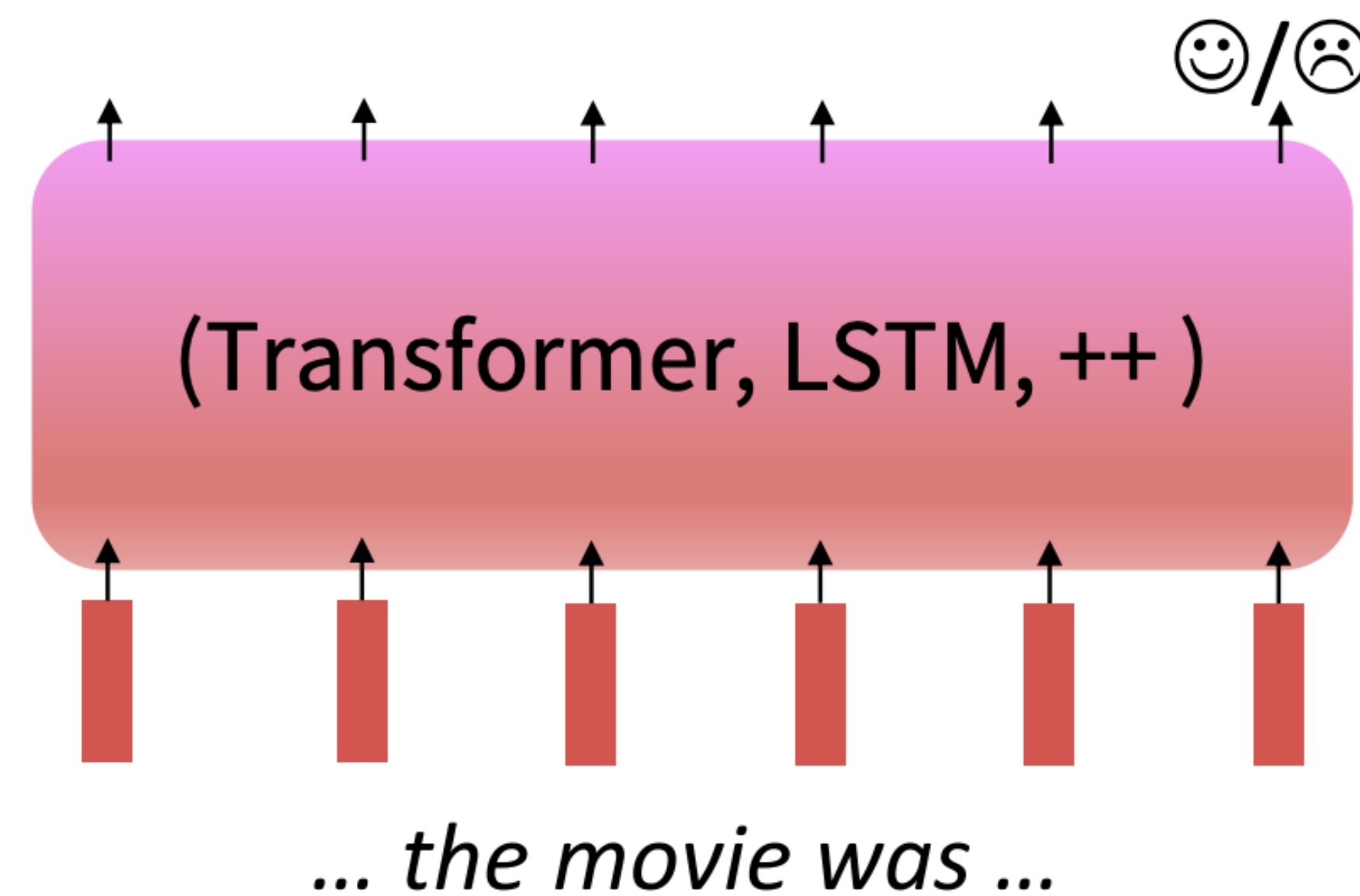> Key idea: "Pretrain once, finetune many times."

Step 1: Pretrain (on language corpora)
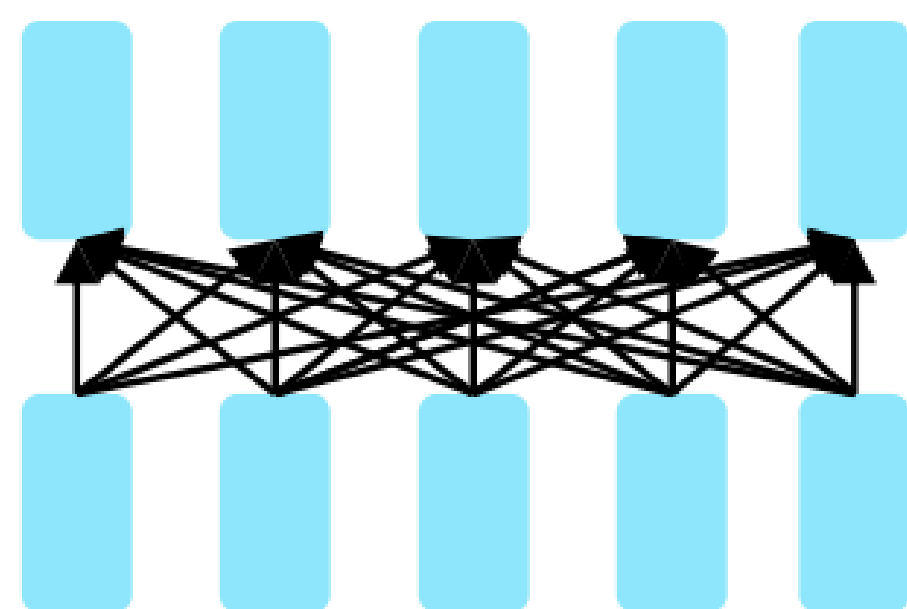Lots of text; learn general things!

Step 2: Finetune (on your task data)
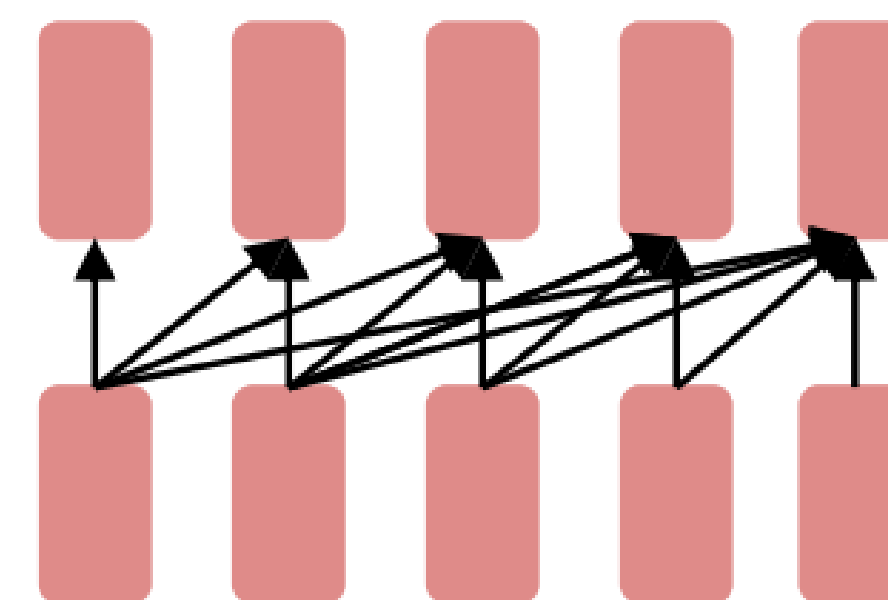Not many labels; adapt to the task!



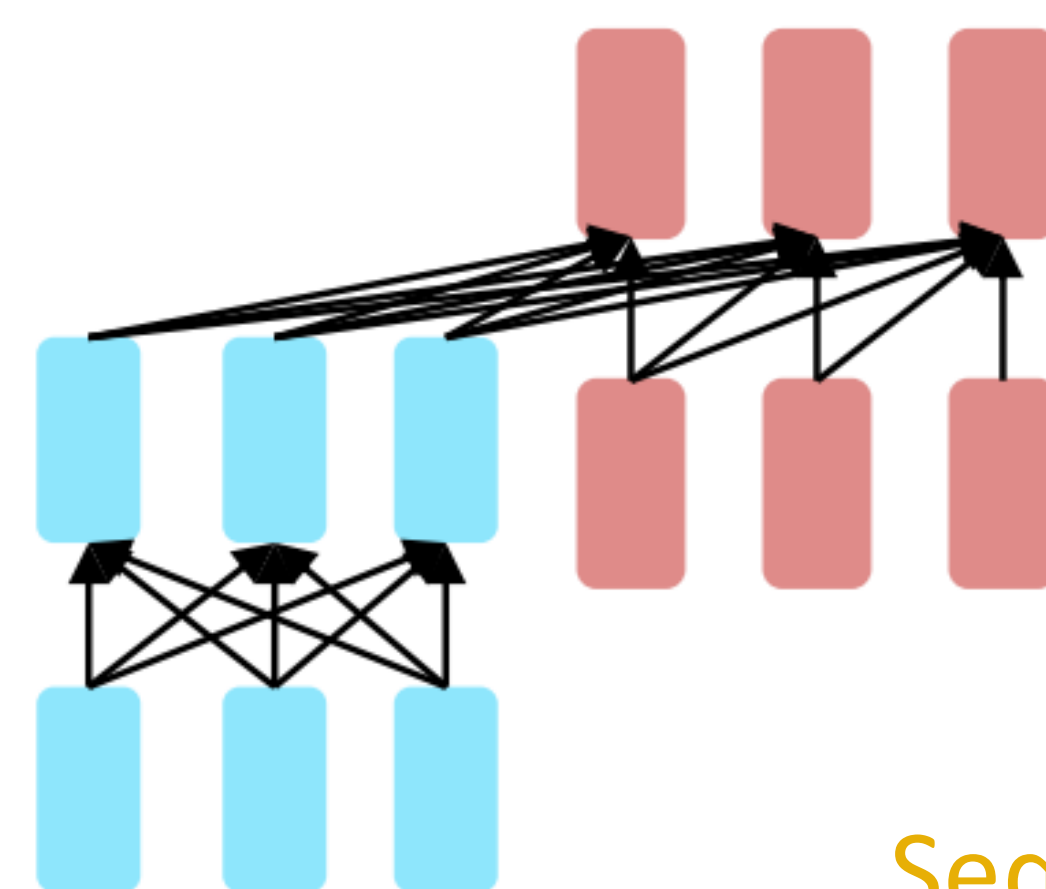4

# Pretraining for three types of architectures



**Encoders**

Bidirectional Context

**Decoders**

Language Models

**Encoder-Decoders**
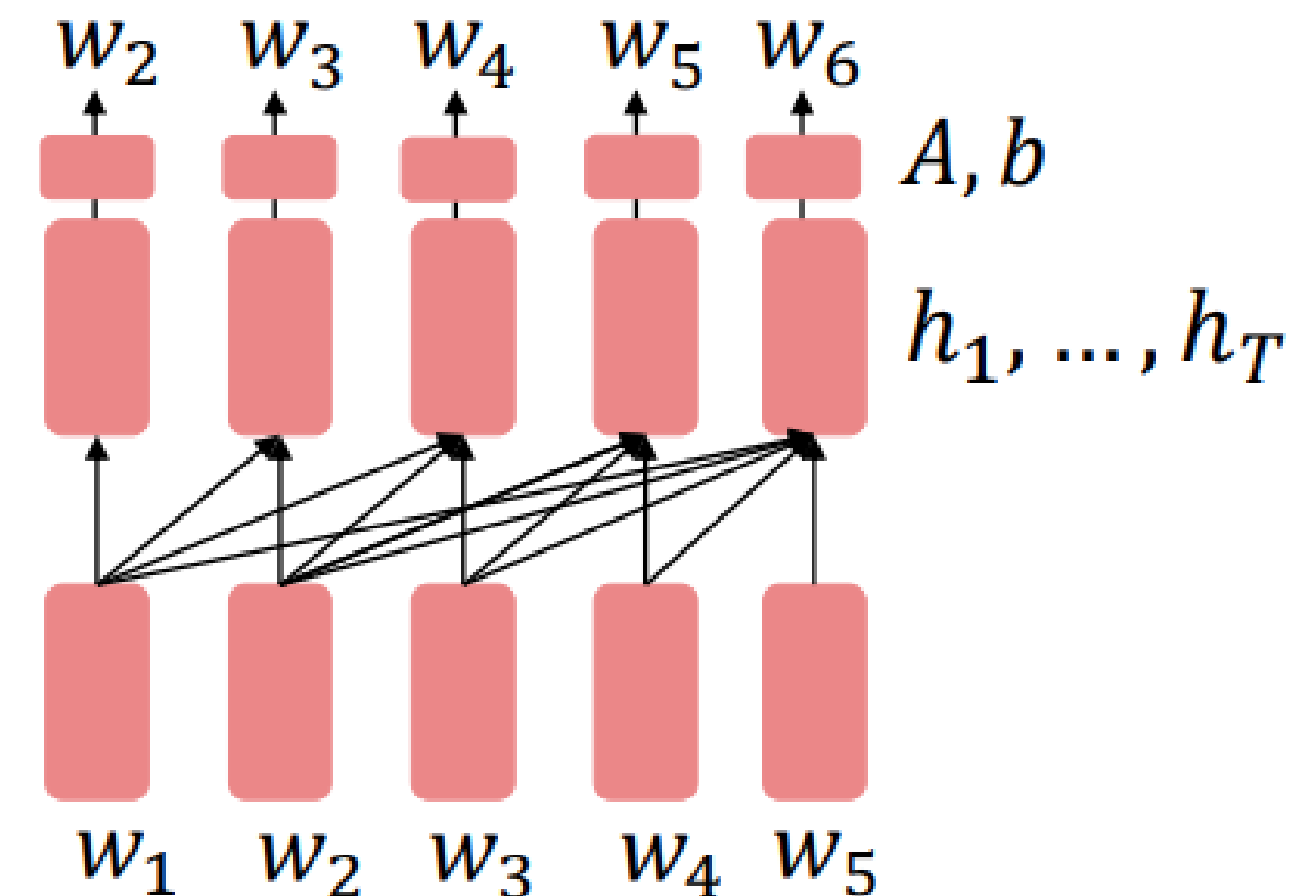
Sequence-to-sequence

# Pretraining Decoders: Generators

- More natural: pretrain decoders as language models and then use them as generators, finetuning their $p_\theta(w_t | w_{1:t-1})$

  - $h_1, \ldots, h_T = Decoder(w_1, \ldots, w_T)$

- $w_t \approx Ah_{t-1} + b$

- Where $A, b$ were pretrained in the language model!

- This is helpful in tasks where the output is a sequence with a vocabulary like that at pretraining time!
  - Dialogue (context=dialogue history)
  - Summarization (context=document)

$$w_2 \quad w_3 \quad w_4 \quad w_5 \quad w_6$$

$$A, b$$

$$h_1, \ldots, h_T$$

$$w_1 \quad w_2 \quad w_3 \quad w_4 \quad w_5$$

The linear layer has been pretrained

# Pretraining for three types of architectures

Encoders

Bidirectional Context

Decoders

Language Models

Encoder-Decoders

Sequence-to-sequence

# Generative Pretrained Transformer (GPT)

- 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"

[Radford et al., 2018]

# GPT-3 and beyond

- Markedly improved generation capabilities by greatly increasing data and model scale

- Solving all tasks through generation

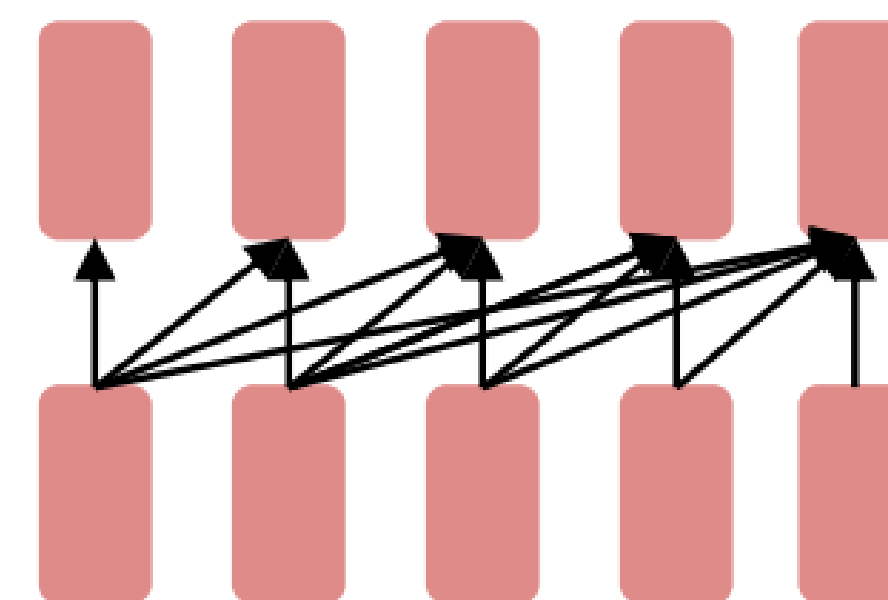- Even obviating the need to fine-tune!

More in future weeks!

# Pretraining for three types of architectures


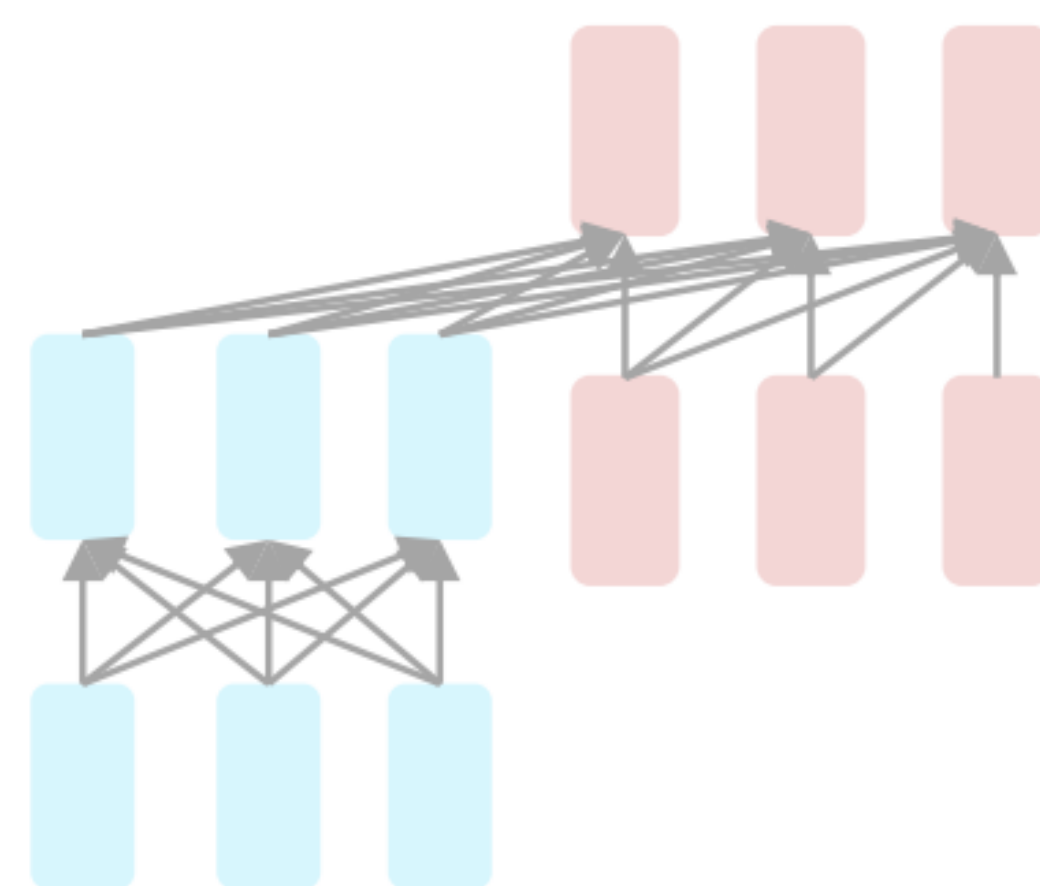
**Encoders**

Bidirectional Context

**Decoders**

Language Models

**Encoder-Decoders**

Sequence-to-sequence

# Pretraining Encoders: Bidirectional Context

I was thinking about the sequence that goes 1, 1, 2, 3, 5, 8, 13, 21, _____

Universal Studios Theme Park is located in _____, California

Problem: Input Reconstruction

'Cause darling i'm a _____ dressed like a daydream

Bidirectional context is important to reconstruct the input!

# BERT: Bidirectional Encoder Representations from Transformers

Devlin et al., 2018 proposed the "Masked LM" objective and released BERT, a Transformer, pretrained to:

- 15% of the input tokens in a training sequence are sampled for learning, these are to be predicted by the model

- Of these

  - 80% are replaced with [MASK]

  - 10% are replaced with randomly selected tokens,

  - Remaining 10% are left unchanged

[Predict these!]   *went*   *to*   *store*

Transformer Encoder

*I*   *pizza*   *to*   *the*   *[M]*

[Replaced]   [Not replaced]   [Masked]

Why? Doesn't let the model get complacent and not build strong representations of non-masked words. (No masks are seen at fine-tuning time!)

12

# BERT: Training Details

- Two models were released:

  - BERT-base: 12 layers, 768-dim hidden states, 12 attention heads, 110 million params.

  - BERT-large: 24 layers, 1024-dim hidden states, 16 attention heads, 340 million params.

- Trained on:

  - BooksCorpus (800 million words)

  - English Wikipedia (2,500 million words)

- Pretraining is expensive and impractical on a single GPU.

  - BERT was pretrained with 64 TPU chips for a total of 4 days.

    - (TPUs are special tensor operation acceleration hardware)

- Finetuning is practical and common on a single GPU

  - "Pretrain once, finetune many times."

# Pretraining for three types of architectures



Encoders

Bidirectional Context

Decoders

Language Models

Encoder-Decoders

Sequence-to-sequence

# T5: A Pretrained Encoder-Decoder Model

- Raffel et al., 2018 built T5, which uses as a span corruption pretraining objective

Replace different-length spans from the input with unique placeholders; decode out the spans that were removed!

Targets

<X> for inviting <Y> last <Z>

Original text

Thank you for inviting me to your party last week.

Still uses an objective that looks like language modeling at the decoder side.

Inputs

Thank you <X> me to your party <Y> week.

# T5: Task Preparation

T5 can be finetuned to answer a wide range of tasks, where the input and output are expressed as a sequence of tokens

# Tokenization in Transformers

# Word Structure and Subword Models

- So far, we have made some assumptions about a language's vocabulary

- We assume a fixed vocabulary of tens of thousands of words, built from the training set

- All novel words seen at test time are mapped to a single UNK

| | word | | vocab mapping | embedding |
|---|---|---|---|---|
| Common words | hat | → | pizza (index) | ▬ |
| | learn | → | tasty (index) | ▬ |
| Variations | taaaaasty | → | UNK (index) | ▬ |
| misspellings | laern | → | UNK (index) | ▬ |
| novel items | Transformerify | → | UNK (index) | ▬ |

# Word Structure in Language

- Finite vocabulary assumptions make even less sense in many languages.
  - Many languages exhibit complex morphology, or word structure.
  - The effect is: more word types, each occurring fewer times.

-ambia = to tell

Example: Swahili verbs can have hundreds of conjugations, each encoding a wide variety of information. (Tense, mood, definiteness, negation, information about the object, ++)

Source: Wiktionary

| Conjugation of *-ambia* | | | | | | | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| **Non-finite forms** | | | | | | | | | | | | | | | | | |
| **Form** | | | | | | **Positive** | | | | | | | | **Negative** | | | |
| Infinitive | | | | | | kuambia | | | | | | | | kutoambia | | | |
| **Simple finite forms** | | | | | | | | | | | | | | | | | |
| **Positive form** | | | | | | **Singular** | | | | | | | | **Plural** | | | |
| Imperative | | | | | | ambia | | | | | | | | ambieni | | | |
| Habitual | | | | | | | | huambia | | | | | | | | | |
| **Complex finite forms** | | | | | | | | | | | | | | | | | |
| Polarity | **Persons** | | | | **Persons / Classes** | | **Classes** | | | | | | | | | | | |
| | 1st | | 2nd | | 3rd / M-wa | | M-mi | | Ma | | Ki-vi | | N | | U | Ku | Pa | Mu |
| | Sg. | Pl. | Sg. | Pl. | Sg. / 1 | Pl. / 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 / 14 | 15 / 17 | 16 | 18 |
| | | | | | | | | | | | | | | | | | | [less ▲] |
| **Past** | | | | | | | | | | | | | | | | | | |
| Positive | niliambia naliambia | tuliambia twaliambia | uliambia waliambia | mliambia mwaliambia | aliambia | waliambia | uliambia | iliambia | liliambia | yaliambia | kiliambia | viliambia | iliambia | ziliambia | uliambia | kuliambia | paliambia | muliambia |
| Negative | sikuambia | hatukuambia | hukuambia | hamkuambia | hakuambia | hawakuambia | haukuambia | haikuambia | halikuambia | hayakuambia | hakikuambia | havikuambia | haikuambia | hazikuambia | haukuambia | hakukuambia | hapakuambia | hamukuambia |
| | | | | | | | | | | | | | | | | | | [less ▲] |
| **Present** | | | | | | | | | | | | | | | | | | |
| Positive | ninaambia naambia | tunaambia | unaambia | mnaambia | anaambia | wanaambia | unaambia | inaambia | linaambia | yanaambia | kinaambia | vinaambia | inaambia | zinaambia | unaambia | kunaambia | panaambia | munaambia |
| Negative | siambii | hatuambii | huambii | hamambii | haambii | hawaambii | hauambii | haiambii | haliambii | hayaambii | hakiambii | haviambii | haiambii | haziambii | hauambii | hakuambii | hapaambii | hamuambii |
| | | | | | | | | | | | | | | | | | | [less ▲] |
| **Future** | | | | | | | | | | | | | | | | | | |
| Positive | nitaambia | tutaambia | utaambia | mtaambia | ataambia | wataambia | utaambia | itaambia | litaambia | yataambia | kitaambia | vitaambia | itaambia | zitaambia | utaambia | kutaambia | pataambia | mutaambia |
| Negative | sitaambia | hatutaambia | hutaambia | hamtaambia | hataambia | hawataambia | hautaambia | haitaambia | halitaambia | hayataambia | hakitaambia | havitaambia | haitaambia | hazitaambia | hautaambia | hakutaambia | hapataambia | hamutaambia |
| | | | | | | | | | | | | | | | | | | [less ▲] |
| **Subjunctive** | | | | | | | | | | | | | | | | | | |
| Positive | niambie | tuambie | uambie | mambie | aambie | waambie | uambie | iambie | liambie | yaambie | kiambie | viambie | iambie | ziambie | uambie | kuambie | paambie | muambie |
| Negative | nisiambie | tusiambie | usiambie | msiambie | asiambie | wasiambie | usiambie | isiambie | lisiambie | yasiambie | kisiambie | visiambie | isiambie | zisiambie | usiambie | kusiambie | pasiambie | musiambie |
| | | | | | | | | | | | | | | | | | | [less ▲] |
| **Present Conditional** | | | | | | | | | | | | | | | | | | |
| Positive | ningeambia | tungeambia | ungeambia | mngeambia | angeambia | wangeambia | ungeambia | ingeambia | lingeambia | yangeambia | kingeambia | vingeambia | ingeambia | zingeambia | ungeambia | kungeambia | pangeambia | mungeambia |
| Negative | nisingeambia singeambia | tusingeambia hatungeambia | usingeambia hungeambia | msingeambia hamngeambia | asingeambia hangeambia | wasingeambia hawangeambia | usingeambia haungeambia | isingeambia | lisingeambia haingeambia | yasingeambia hayangeambia | kisingeambia hakingeambia | visingeambia havingeambia | isingeambia haingeambia | zisingeambia hazingeambia | usingeambia haungeambia | kusingeambia hakungeambia | pasingeambia hapangeambia | musingeambia hamungeambia |
| | | | | | | | | | | | | | | | | | | [less ▲] |
| **Past Conditional** | | | | | | | | | | | | | | | | | | |
| Positive | ningaliambia | tungaliambia | ungaliambia | mngaliambia | angaliambia | wangaliambia | ungaliambia | ingaliambia | lingaliambia | yangaliambia | kingaliambia | vingaliambia | ingaliambia | zingaliambia | ungaliambia | kungaliambia | pangaliambia | mungaliambia |
| Negative | nisingaliambia singaliambia | tusingaliambia hatungaliambia | usingaliambia hungaliambia | msingaliambia hamngaliambia | asingaliambia hangaliambia | wasingalia mbia hawangalia mbia | usingaliambia haungaliambia | isingaliambia haingaliambia | lisingaliambia halingaliambia | yasingaliambia hayangaliambia | kisingaliambia hakingaliam bia | visingaliambia havingaliam bia | isingaliambia haingaliambia | zisingaliambia hazingaliam bia | usingaliambia haungaliambia | kusingaliambia hakungaliam bia | pasingaliambia hapangaliam bia | musingaliambia hamungalia mbia |
| | | | | | | | | | | | | | | | | | | [less ▲] |
| **Conditional Contrary to Fact** | | | | | | | | | | | | | | | | | | |
| Positive | ningeliambia | tungeliambia | ungeliambia | mngeliambia | angeliambia | wangeliambia | ungeliambia | ingeliambia | lingeliambia | yangeliambia | kingeliambia | vingeliambia | ingeliambia | zingeliambia | ungeliambia | kungeliambia | pangeliambia | mungeliambia |
| **Gnomic** | | | | | | | | | | | | | | | | | | |
| Positive | naambia | twaambia | waambia | mwaambia | aambia | waambia | waambia | yaambia | laambia | yaambia | chaambia | vyaambia | yaambia | zaambia | waambia | kwaambia | paambia | mwaambia |
| **Perfect** | | | | | | | | | | | | | | | | | | [less ▲] |

# Subword Modeling

- Solution: look at subwords

- Subword modeling in NLP encompasses a wide range of methods for reasoning about structure below the word level

  - Subwords may be parts of words, characters, bytes

- The dominant modern paradigm is to learn a vocabulary of parts of words (subword tokens).

- At training and testing time, each word is split into a sequence of known subwords.

# Byte-pair encoding

- Byte-pair encoding is a simple, effective strategy for defining a subword vocabulary

- Adapted for word segmentation from data compression technique (Gage, 1994)

  - Instead of merging frequent pairs of bytes, we merge characters or character sequences

# Byte-pair encoding

- Algorithm:

  1. Start with a vocabulary containing only characters and an "end-of-word" symbol.

  2. Using a corpus of text, find the most common adjacent characters "a,b"; add "ab" as a subword.

  3. Replace instances of the character pair with the new subword; repeat until desired vocabulary size.

# BPE in action

### Corpus

| | | |
|---|---|---|
| low | lower | newest |
| low | lower | newest |
| low | widest | newest |
| low | widest | newest |
| low | widest | newest |

### Corpus

| | | |
|---|---|---|
| low</w> | lower</w> | newest</w> |
| low</w> | lower</w> | newest</w> |
| low</w> | widest</w> | newest</w> |
| low</w> | widest</w> | newest</w> |
| low</w> | widest</w> | newest</w> |

### Corpus

| | | |
|---|---|---|
| l o w </w> | l o w e r </w> | n e w e s t </w> |
| l o w </w> | l o w e r </w> | n e w e s t </w> |
| l o w </w> | w i d e s t </w> | n e w e s t </w> |
| l o w </w> | w i d e s t</w> | n e w e s t </w> |
| l o w </w> | w i d e s t </w> | n e w e s t </w> |

### Vocabulary

| d | e | i | l | n | o | s | t | w |
|---|---|---|---|---|---|---|---|---|
| es | | | | | | | | |

### Frequency

| | | |
|---|---|---|
| d-e (3) | l-o (7) | t-</w> (8) |
| e-r (2) | n-e (5) | w-</w> (5) |
| e-s (8) | o-w (7) | w-e (7) |
| e-w (5) | r-</w> (2) | w-i (3) |
| i-d (3) | s-t (8) | |

Source: https://www.analyticsvidhya.com/blog/2020/05/what-is-tokenization-nlp/

# BPE in action

## Corpus

| low | lower | newest |
|-----|-------|--------|
| low | lower | newest |
| low | widest | newest |
| low | widest | newest |
| low | widest | newest |

## Corpus

| low</w> | lower</w> | newest</w> |
|---------|-----------|------------|
| low</w> | lower</w> | newest</w> |
| low</w> | widest</w> | newest</w> |
| low</w> | widest</w> | newest</w> |
| low</w> | widest</w> | newest</w> |

## Corpus

| l o w </w> | l o w e r </w> | n e w es t </w> |
|------------|----------------|-----------------|
| l o w </w> | l o w e r </w> | n e w es t </w> |
| l o w </w> | w i d es t </w> | n e w es t </w> |
| l o w </w> | w i d es t</w> | n e w es t </w> |
| l o w </w> | w i d es t </w> | n e w es t </w> |

## Vocabulary

| d | e | i | l | n | o | s | t | w |
|---|---|---|---|---|---|---|---|---|
| es | est | | | | | | | |

## Frequency

| d-es (3) | l-o (7) | w-</w> (5) |
|----------|---------|------------|
| e-r (2) | n-e (5) | w-es (5) |
| e-w (5) | o-w (7) | w-e (2) |
| es-t (8) | r-</w> (2) | w-i (3) |
| i-d (3) | t-</w> (8) | |

Source: https://www.analyticsvidhya.com/blog/2020/05/what-is-tokenization-nlp/

24

# BPE in action

| Corpus | | |
|---|---|---|
| low | lower | newest |
| low | lower | newest |
| low | widest | newest |
| low | widest | newest |
| low | widest | newest |

| Corpus | | |
|---|---|---|
| low</w> | lower</w> | newest</w> |
| low</w> | lower</w> | newest</w> |
| low</w> | widest</w> | newest</w> |
| low</w> | widest</w> | newest</w> |
| low</w> | widest</w> | newest</w> |

| Corpus | | |
|---|---|---|
| l o w </w> | l o w e r </w> | n e w est </w> |
| l o w </w> | l o w e r </w> | n e w est </w> |
| l o w </w> | w i d est </w> | n e w est </w> |
| l o w </w> | w i d est </w> | n e w est </w> |
| l o w </w> | w i d est </w> | n e w est </w> |

| Vocabulary | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|
| d | e | i | l | n | o | s | t | w | |
| es | est | est</w> | lo | low | low</w> | ne | new | newest</w> | |

After 10 merges

# Byte-pair encoding

- At test time, first split words into sequences of characters, then apply the learned operations to merge the characters into larger, known symbols

- Originally used in NLP for machine translation; now a similar method (WordPiece) is used in pretrained models.

# Word structure and subword models

- Common words end up being a part of the subword vocabulary, while rarer words are split into (sometimes intuitive, sometimes not) components.

- In the worst case, words are split into as many subwords as they have characters.

| | word | | vocab mapping | embedding |
|---|---|---|---|---|
| Common words | hat | → | hat | |
| | learn | → | learn | |
| Variations | taaaaasty | → | taa## aaa## sty | |
| misspellings | laern | → | la## ern## | |
| novel items | Transformerify | → | Transformer## ify | |

# Natural Language Generation

# Natural Language Generation

- Natural language understanding and natural language generation are two sides of the same coin

    - In order to generate good language, you need to understand language

    - If you understand language, you should be able to generate it (with some effort)

- NLG is the workhorse of many classic and novel applications

    - AI Assistants

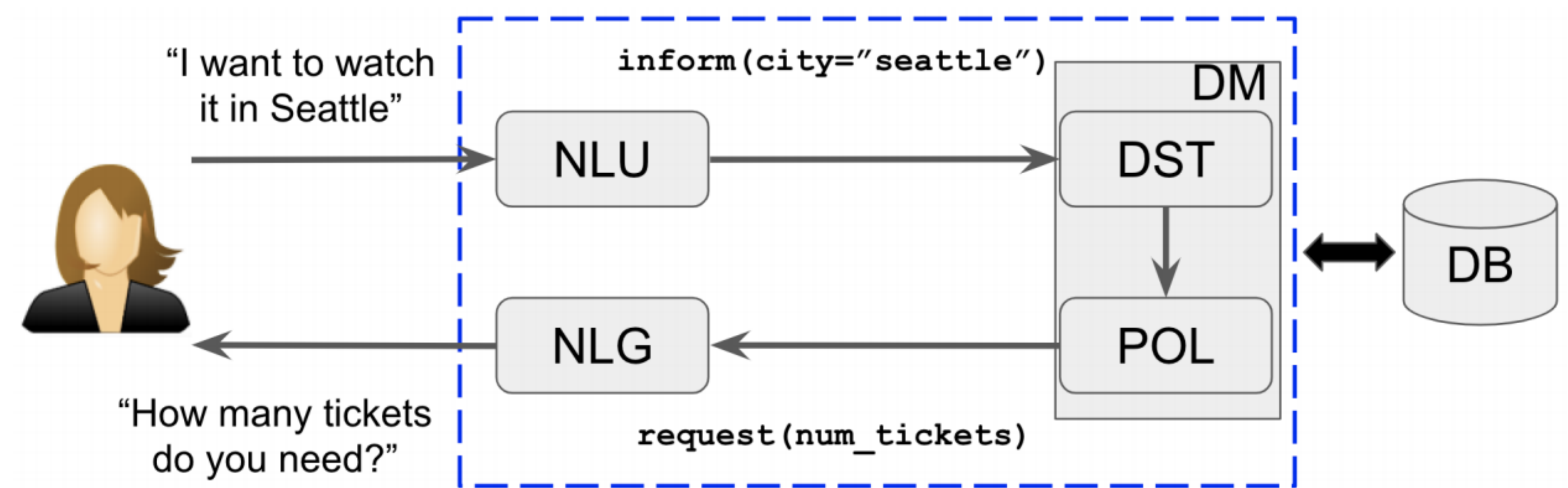    - Translators

    - Search summarizers

# NLG Use Cases



**Simple and Effective Multi-Paragraph Reading Comprehension**

Christopher Clark, Matt Gardner · Computer Science · ACL · 29 October 2017

**TLDR** We propose a state-of-the-art pipelined method for training neural paragraph-level question answering models on document QA data. Expand

66 236  PDF · 📄 View PDF on arXiv  🔖 Save  🔔 Alert  66 Cite  👍 Research Feed

Summarization

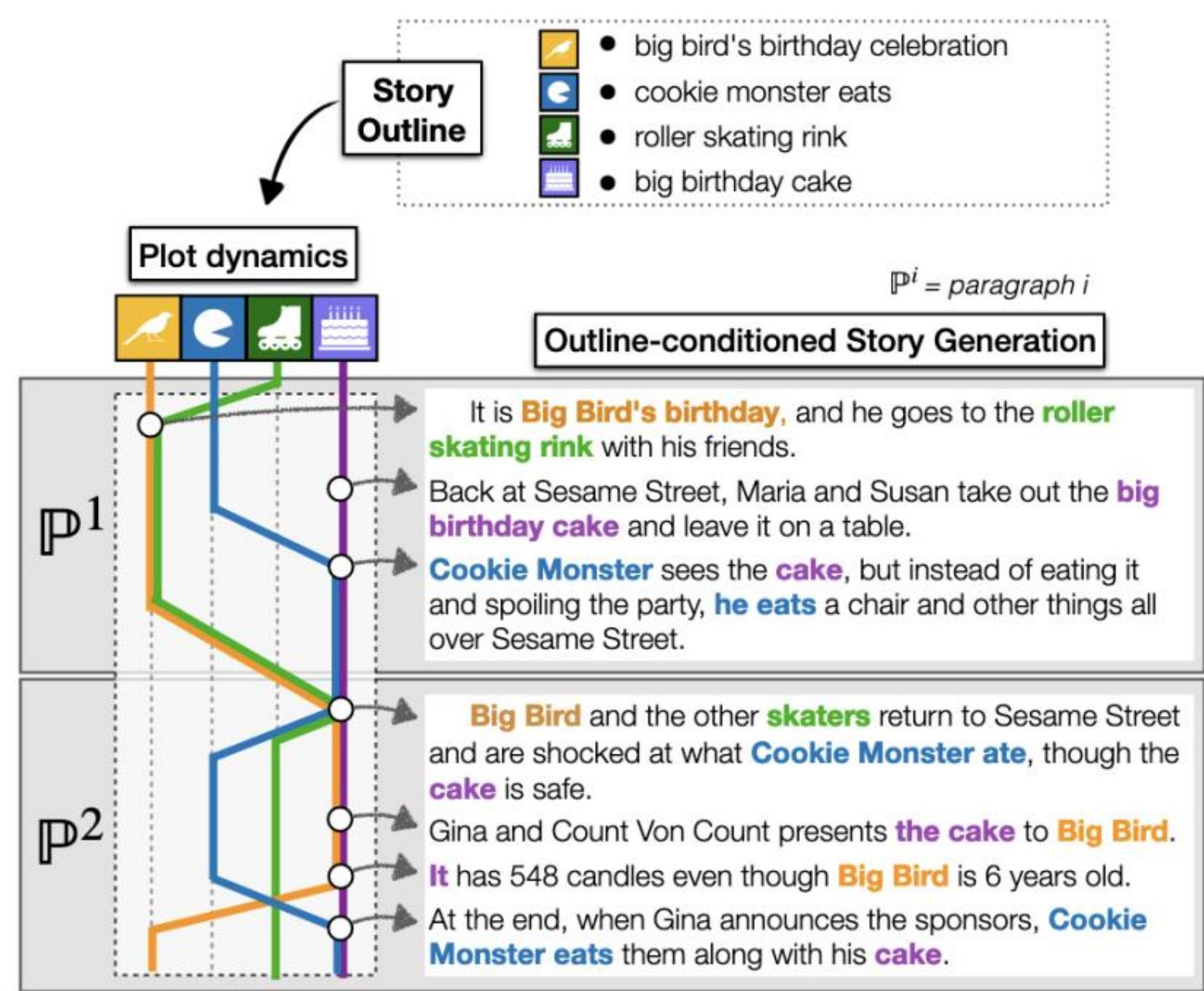"I want to watch it in Seattle"

inform(city="seattle")

NLU → DST

DM

DB

NLG ← POL

request(num_tickets)

"How many tickets do you need?"

Task-driven Dialog

☐ Chit-chat
☐ Task-oriented

User: Does money buy happiness?

System: Depends how much money you spend on it.

User: I don't have much money...

System: Me too.

User: I am looking for a place to stay that has cheap price range it should be in a type of hotel.

System: Okay, do you have a specific area you want to stay in?

User: I would like someone in the center.

Chitchat Dialog

30

# More Interesting NLG Uses

## Creative stories



Rashkin et al., 2020

## Data-to-text



Craig finished his eleven NFL seasons with 8,189 rushing yards and 566 receptions for 4,911 receiving yards.
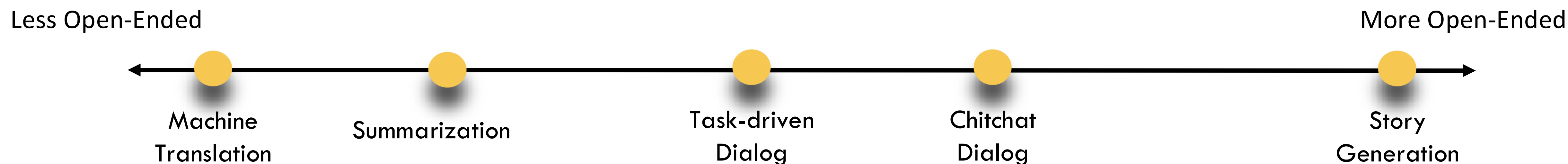
Parikh et al., 2020

## Visual description



Two children are sitting at a table in a restaurant. The children are one little girl and one little boy. The little girl is eating a pink frosted donut with white icing lines on top of it. The girl has blonde hair and is wearing a green jacket with a black long sleeve shirt underneath. The little boy is wearing a black zip up jacket and is holding his finger to his lip but is not eating. A metal napkin dispenser is in between them at the table. The wall next to them is white brick. Two adults are on the other side of the short white brick wall. The room has white circular lights on the ceiling and a large window in the front of the restaurant. It is daylight outside.

Krause et al., 2017

31

# Broad Spectrum of NLG Tasks

Less Open-Ended

More Open-Ended

Machine
Translation

Summarization

Task-driven
Dialog

Chitchat
Dialog

Story
Generation

Open-ended generation: the output distribution still has high freedom.

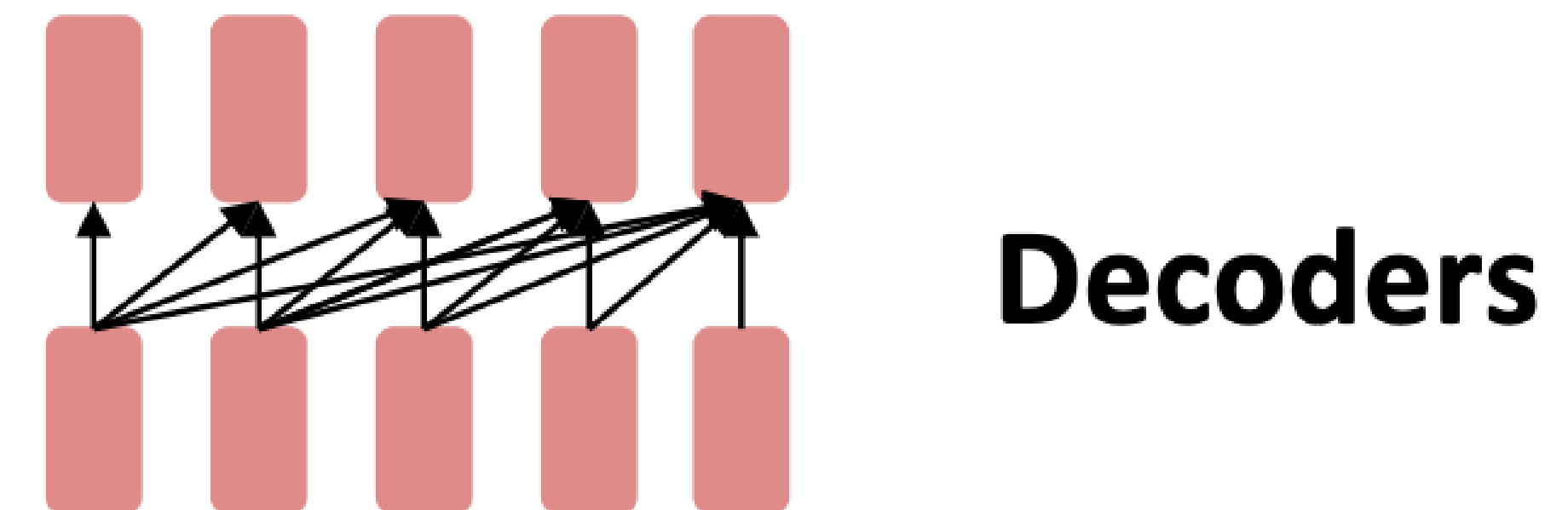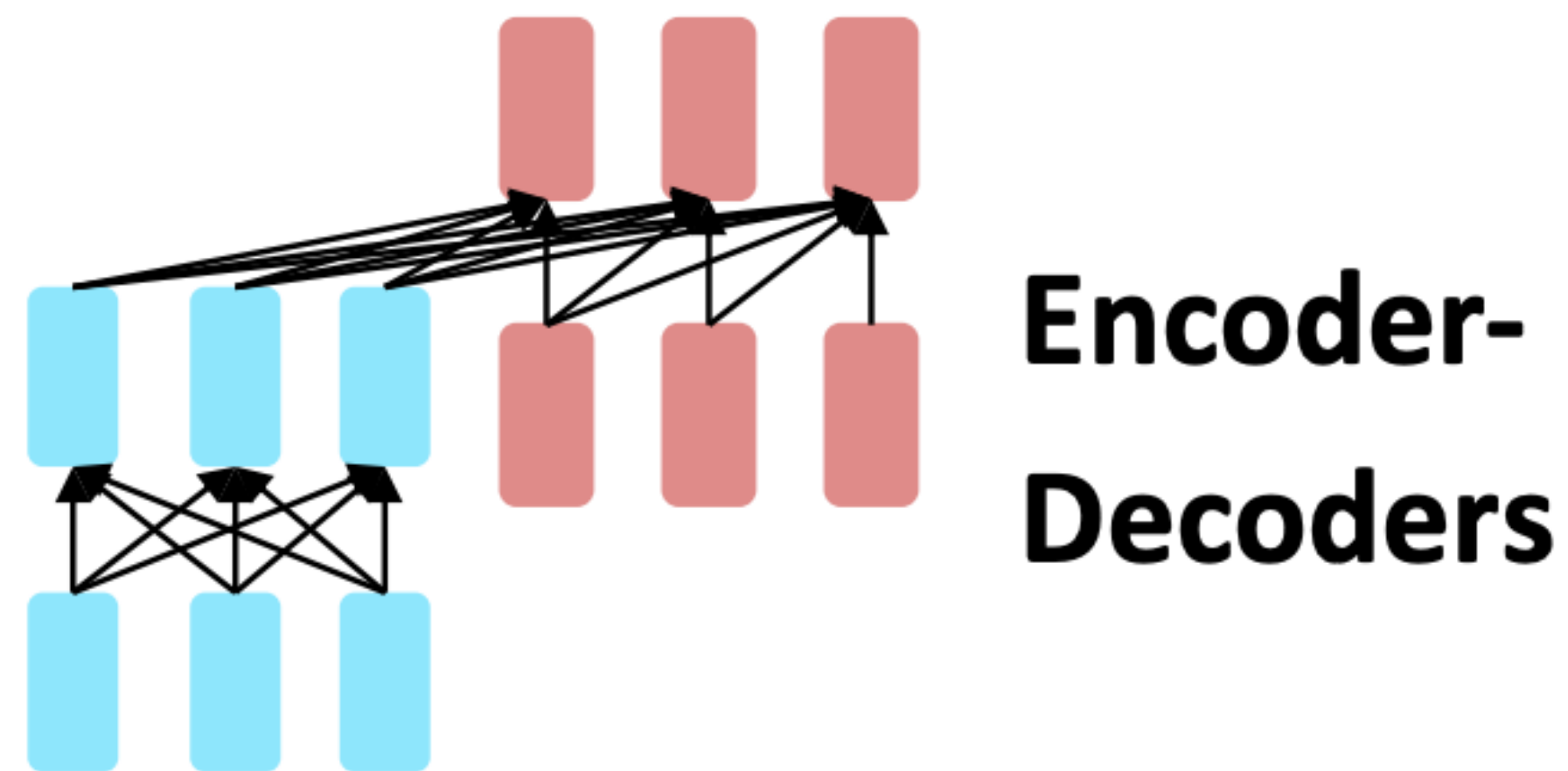Non-open-ended generation: the input mostly determines the output generation.

# Broad Spectrum of NLG Tasks

Less Open-Ended                                                                                    More Open-Ended

Machine
Translation          Summarization          Task-driven
Dialog          Chitchat
Dialog          Story
Generation

**Encoder-Decoders**

**Decoders**
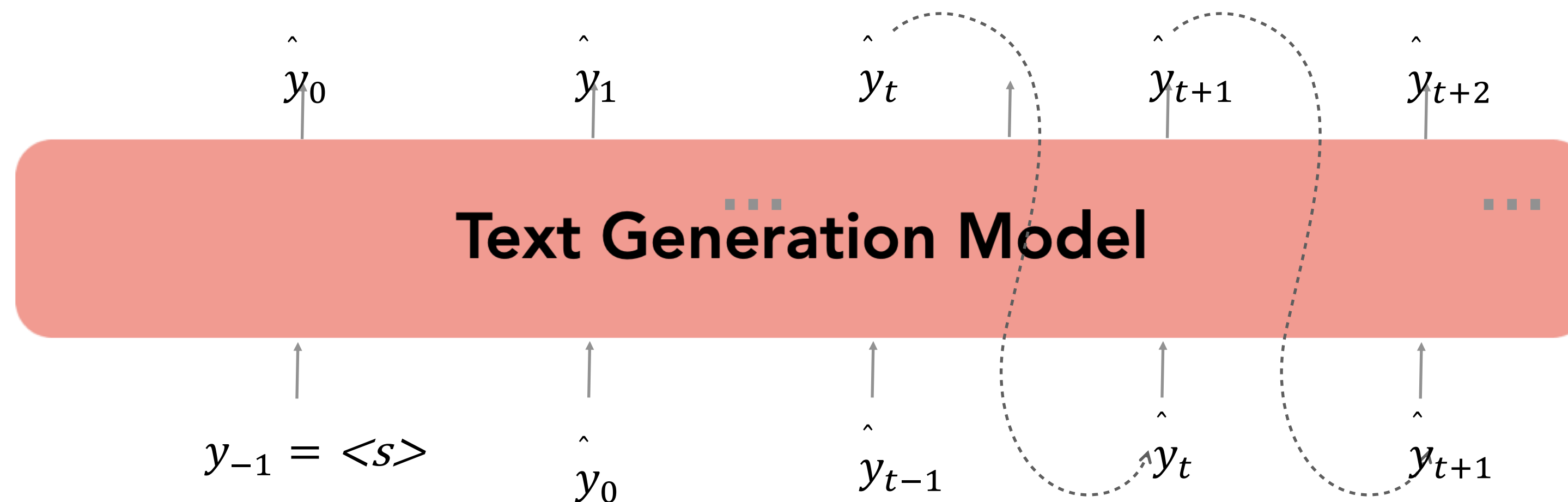
# Language Generation: Fundamentals

- In autoregressive text generation models, at each time step $t$, our model takes in a sequence of tokens as input $S = f_\theta(y_{<t}) \in \mathbb{R}^V$ and outputs a new token, $\hat{y}_t$

- For model $f_\theta(\cdot)$ and vocabulary $V$, we get scores $S = f_\theta(y_{<t}) \in \mathbb{R}^V$

$$P(w|y_{<t}) = \frac{\exp(S_w)}{\sum_{v \in V} \exp(S_v)}$$

# Language Generation: Training

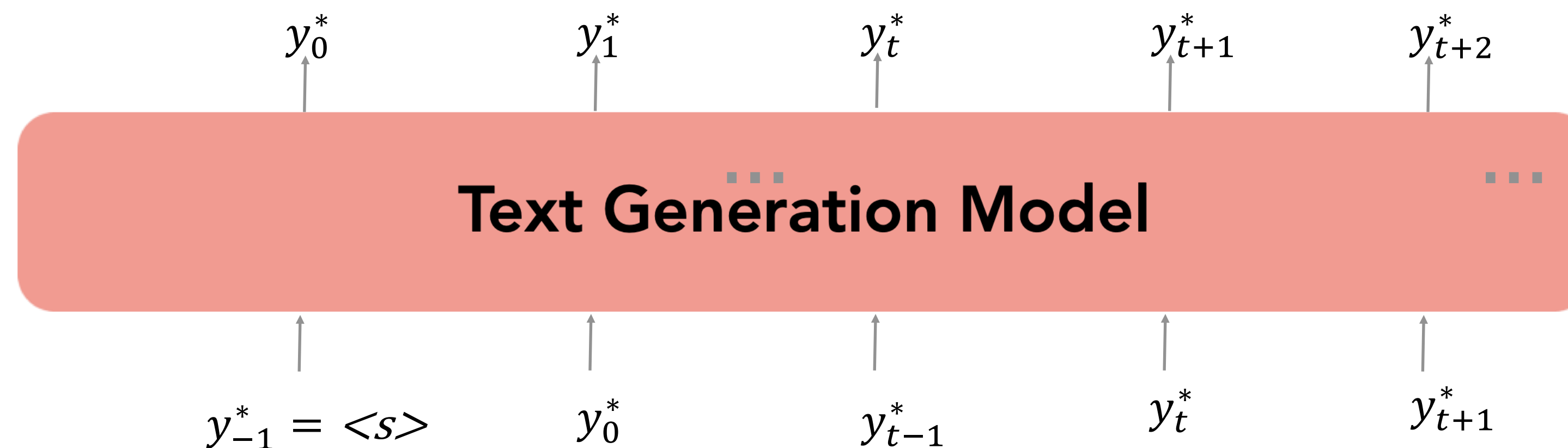- Trained one token at a time to maximize the probability of the next token $y_t^*$ given preceding words $y_{<t}^*$

$$\mathcal{L} = -\sum_{t=1}^{T} \log P(y_t|y_{<t}) = -\sum_{t=1}^{T} \log \frac{\exp(S_{y_t|y_{<t}})}{\sum_{v \in V} \exp(S_{v|y_{<t}})}$$

- Classification task at each time step trying to predict the actual word $y_t^*$ in the training data

- "Teacher forcing" (reset at each time step to the ground truth)

# Teacher Forcing
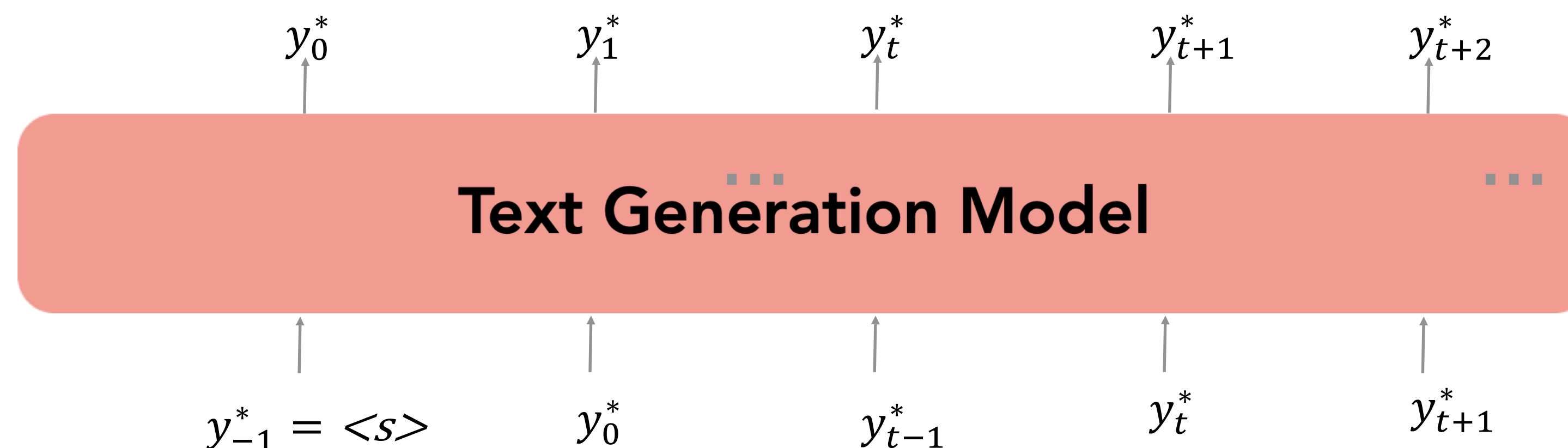
- At each time step $t$ in decoding we force the system to **use the gold target token from training as the next input** $x_{t+1}$ ,rather than allowing it to rely on the (possibly erroneous) decoder output $\hat{y}_t$

$$y_0^* \qquad y_1^* \qquad y_t^* \qquad y_{t+1}^* \qquad y_{t+2}^*$$

**Text Generation Model**

$$y_{-1}^* = <s> \qquad y_0^* \qquad y_{t-1}^* \qquad y_t^* \qquad y_{t+1}^*$$

# Teacher Forcing

- Runs the risk of exposure bias!

  - During training, our model's inputs are gold context tokens from real, human-generated texts

  - At generation time, our model's inputs are previously–decoded tokens

$$y_0^* \qquad y_1^* \qquad y_t^* \qquad y_{t+1}^* \qquad y_{t+2}^*$$

**Text Generation Model**

$$y_{-1}^* = <s> \qquad y_0^* \qquad y_{t-1}^* \qquad y_t^* \qquad y_{t+1}^*$$
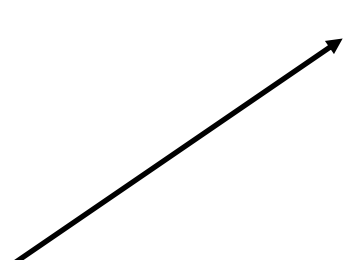
37

# Language Generation: Training

- More on LLM training (post-training) in the next class

# Language Generation: Inference

- At inference time, our decoding algorithm defines a function to select a token from this distribution:

Inference / Decoding Algorithm

$$\hat{y}_t = g(P(y_t|y_{<t}))$$

- The "obvious" decoding algorithm is to **greedily** choose the highest probability next token according to the model at each time step

$$g = \arg max$$

$$\hat{y}_t = \arg max_{w \in V}(P(y_t = w|y_{<t}))$$

USC Viterbi

# Classic Inference Algorithms:
# Greedy and Beam Search

# Decoding

- Generation from a language model is also called decoding

  - Think encoder-decoder

  - Also called inference

- Strategy so far: Take $\arg max$ on each step of the decoder to produce the most probable word on each step

- This is called greedy decoding

  - Greedy Strategy: we are not looking ahead, we are making the hastiest decision given all the information we have

# Greedy Decoding: Issues

- Greedy decoding has no wiggle room for errors!

  - Input: the green witch arrived

    - Output: llego

    - Output: llego la

    - Output: llego la <span style="color:red">verde</span>

- How to fix this?

  - Need a lookahead strategy / longer-term planning

# Exhaustive Search Decoding

- Ideally, we want to find a (length T) translation y that maximizes

$$P(y|x) = P(y_1|x)\, P(y_2|y_1, x)\, P(y_3|y_1, y_2, x)\ldots, P(y_T|y_1, \ldots, y_{T-1}, x)$$

$$= \prod_{t=1}^{T} P(y_t|y_1, \ldots, y_{t-1}, x)$$

- We could try computing all possible sequences $y$

  - This means that on each step t of the decoder, we're tracking $V^t$ possible partial translations, where $V$ is vocab size

  - This $O(V^T)$ complexity is far too expensive!

# Beam Search Decoding

- Core idea: On each step of decoder, keep track of the k most probable partial translations (which we call hypotheses)

  - k is the beam size (in practice around 5 to 10, in NMT)

- A hypothesis has a score which is its log probability:

$$\text{score}(y_1, \ldots, y_t) = \log P_{\text{LM}}(y_1, \ldots, y_t | x) = \sum_{i=1}^{t} \log P_{\text{LM}}(y_i | y_1, \ldots, y_{i-1}, x)$$

  - Scores are all negative, and higher score is better

  - We search for high-scoring hypotheses, tracking top k on each step

- Beam search is not guaranteed to find optimal solution

- But much more efficient than exhaustive search!

# Beam Search Decoding: Example

Beam size = k = 2. Blue numbers = $\text{score}(y_1, \ldots, y_t) = \sum_{i=1}^{t} \log P_{\text{LM}}(y_i | y_1, \ldots, y_{i-1}, x)$

<START>

Calculate prob
dist of next word

Slide credit: Chris Manning

# Beam Search Decoding: Example

Beam size = k = 2. Blue numbers = $\text{score}(y_1, \ldots, y_t) = \sum_{i=1}^{t} \log P_{\text{LM}}(y_i | y_1, \ldots, y_{i-1}, x)$

-0.7 = log $P_{\text{LM}}$(*he*|<START>)

| *he* |

| <START> |

| *I* |

-0.9 = log $P_{\text{LM}}$(*I*|<START>)

Take top *k* words
and compute scores

**USC**Viterbi

# Beam Search Decoding: Example

Beam size = k = 2. Blue numbers = $\text{score}(y_1, \ldots, y_t) = \sum_{i=1}^{t} \log P_{\text{LM}}(y_i | y_1, \ldots, y_{i-1}, x)$
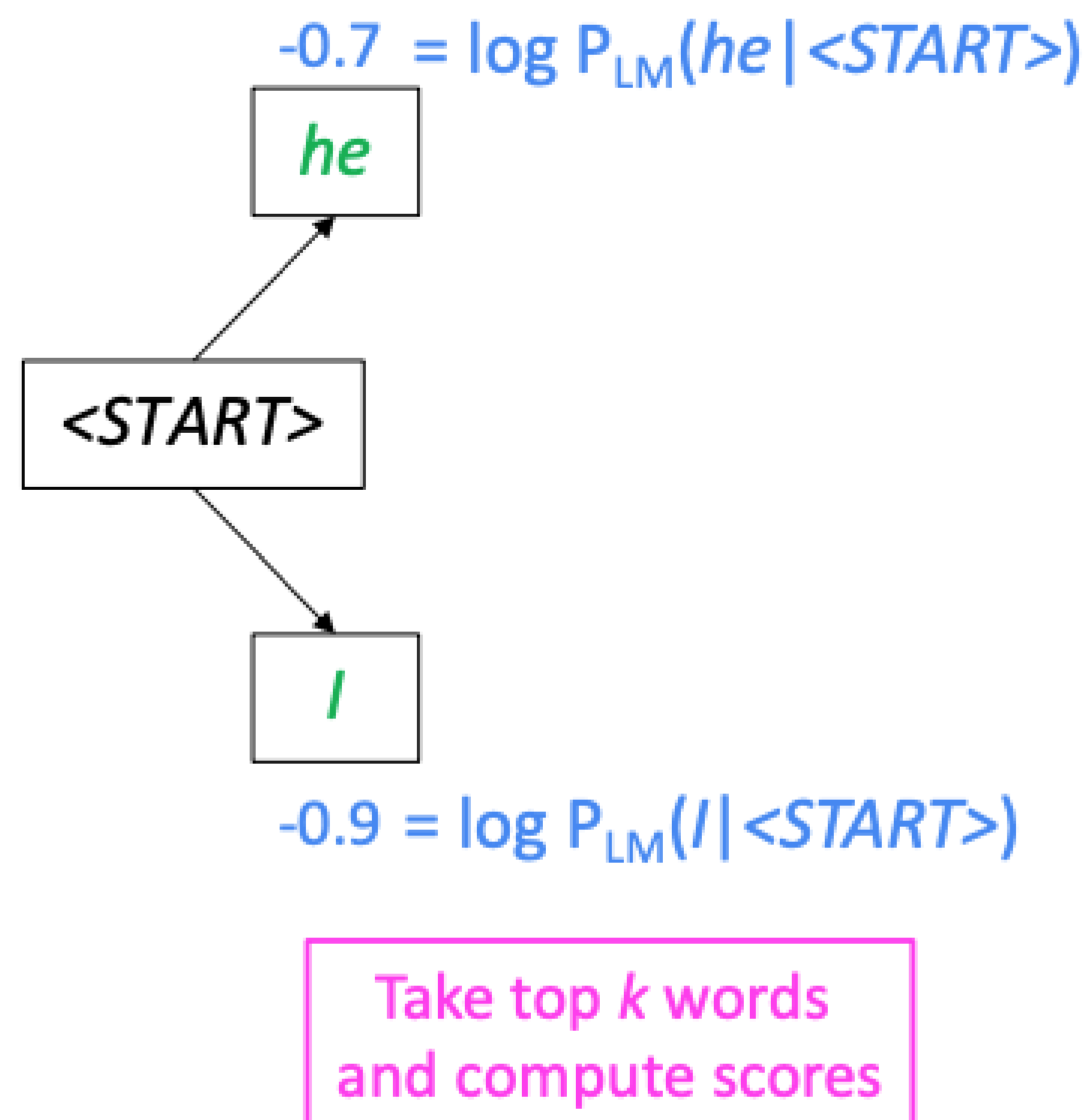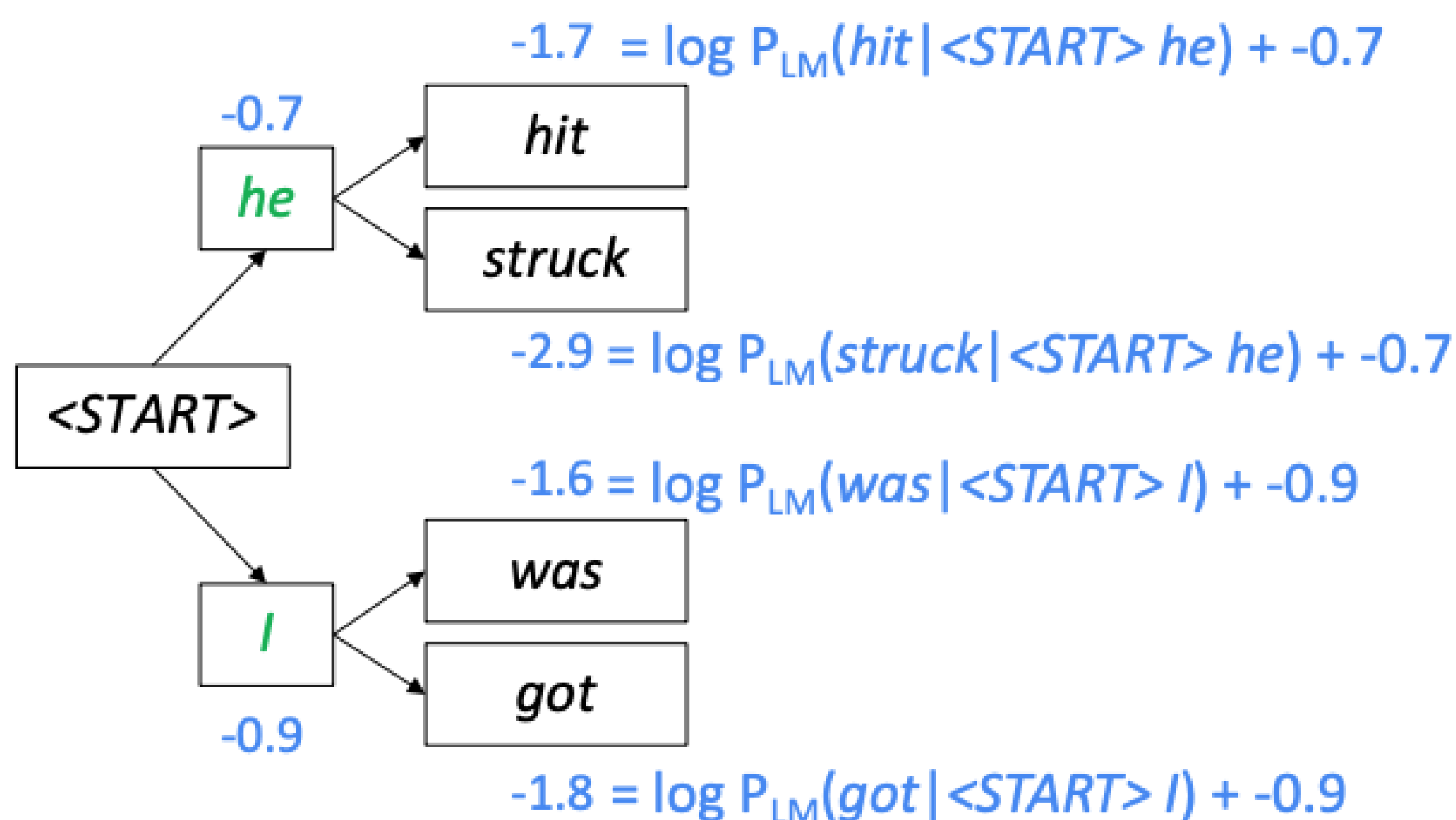


-1.7 = log $P_{\text{LM}}$(*hit*|<START> *he*) + -0.7

-0.7

hit

*he*

struck

<START>

-2.9 = log $P_{\text{LM}}$(*struck*|<START> *he*) + -0.7

-1.6 = log $P_{\text{LM}}$(*was*|<START> *I*) + -0.9

was

*I*

got

-0.9

-1.8 = log $P_{\text{LM}}$(*got*|<START> *I*) + -0.9

For each of the *k* hypotheses, find
top *k* next words and calculate scores
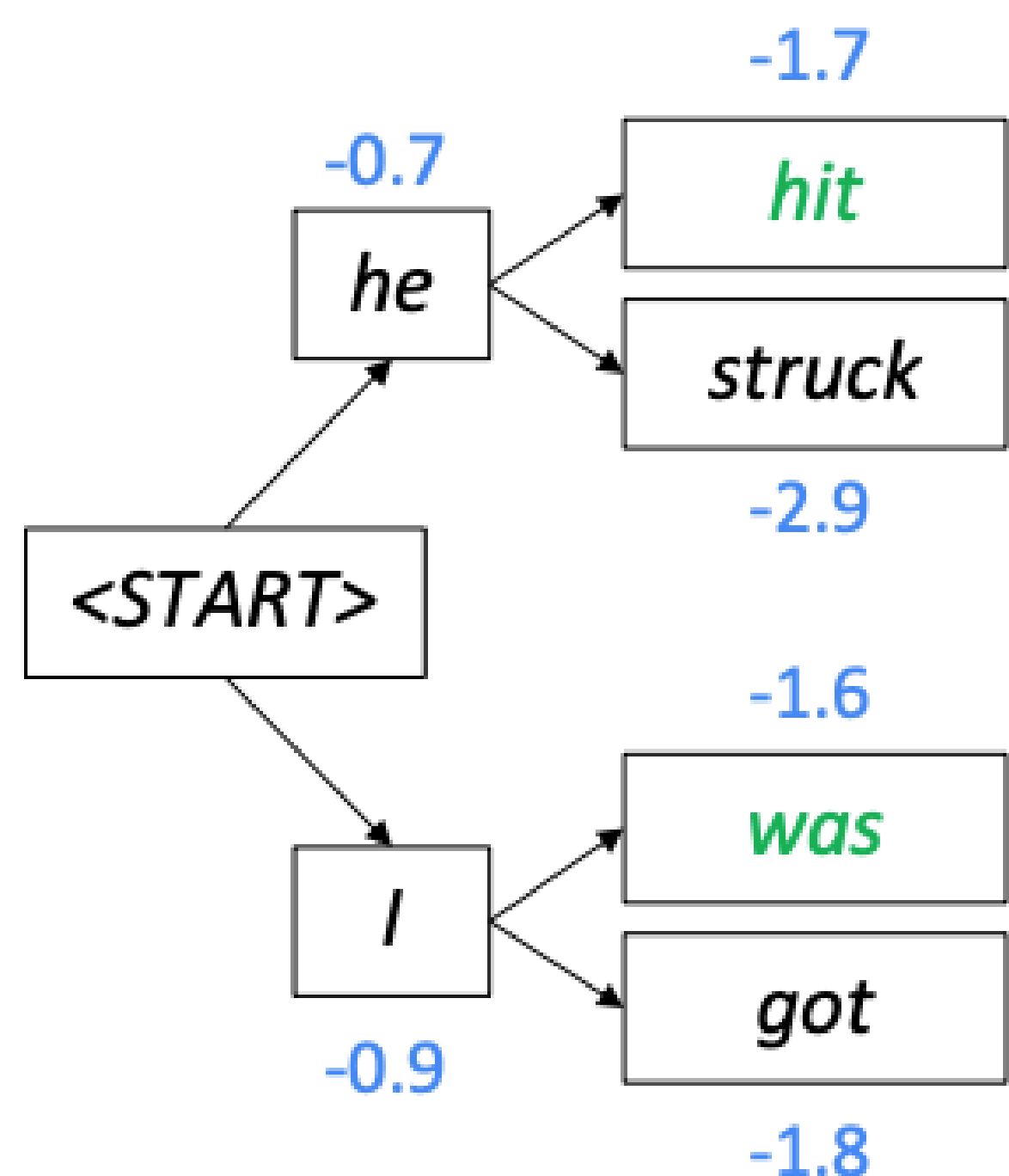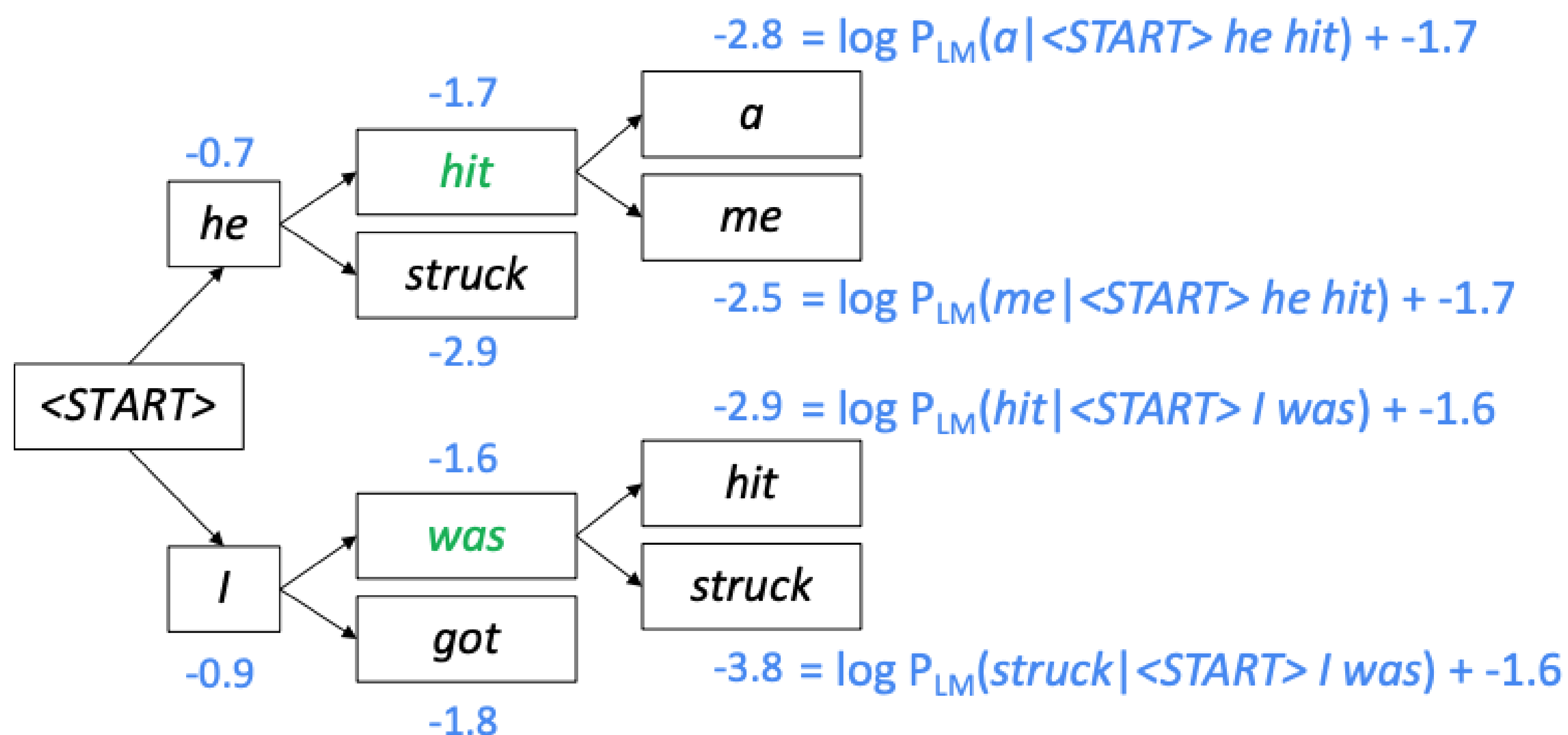
Slide credit: Chris Manning

# Beam Search Decoding: Example

Beam size = k = 2. Blue numbers = $\text{score}(y_1, \ldots, y_t) = \sum_{i=1}^{t} \log P_{\text{LM}}(y_i | y_1, \ldots, y_{i-1}, x)$

```
              -1.7
      -0.7    ┌────────┐
   ┌──────┐   │  hit   │
   │  he  │ ◁ └────────┘
   └──────┘   ┌────────┐
         △    │ struck │
┌─────────┐   └────────┘
│ <START> │       -2.9
└─────────┘
         ▽        -1.6
   ┌──────┐   ┌────────┐
   │  I   │ ◁ │  was   │
   └──────┘   └────────┘
    -0.9      ┌────────┐
              │  got   │
              └────────┘
                  -1.8
```

Of these $k^2$ hypotheses,
just keep $k$ with highest scores

Slide credit: Chris Manning

**USC**Viterbi

# Beam Search Decoding: Example

Beam size = k = 2. Blue numbers = $\text{score}(y_1, \ldots, y_t) = \sum_{i=1}^{t} \log P_{\text{LM}}(y_i | y_1, \ldots, y_{i-1}, x)$
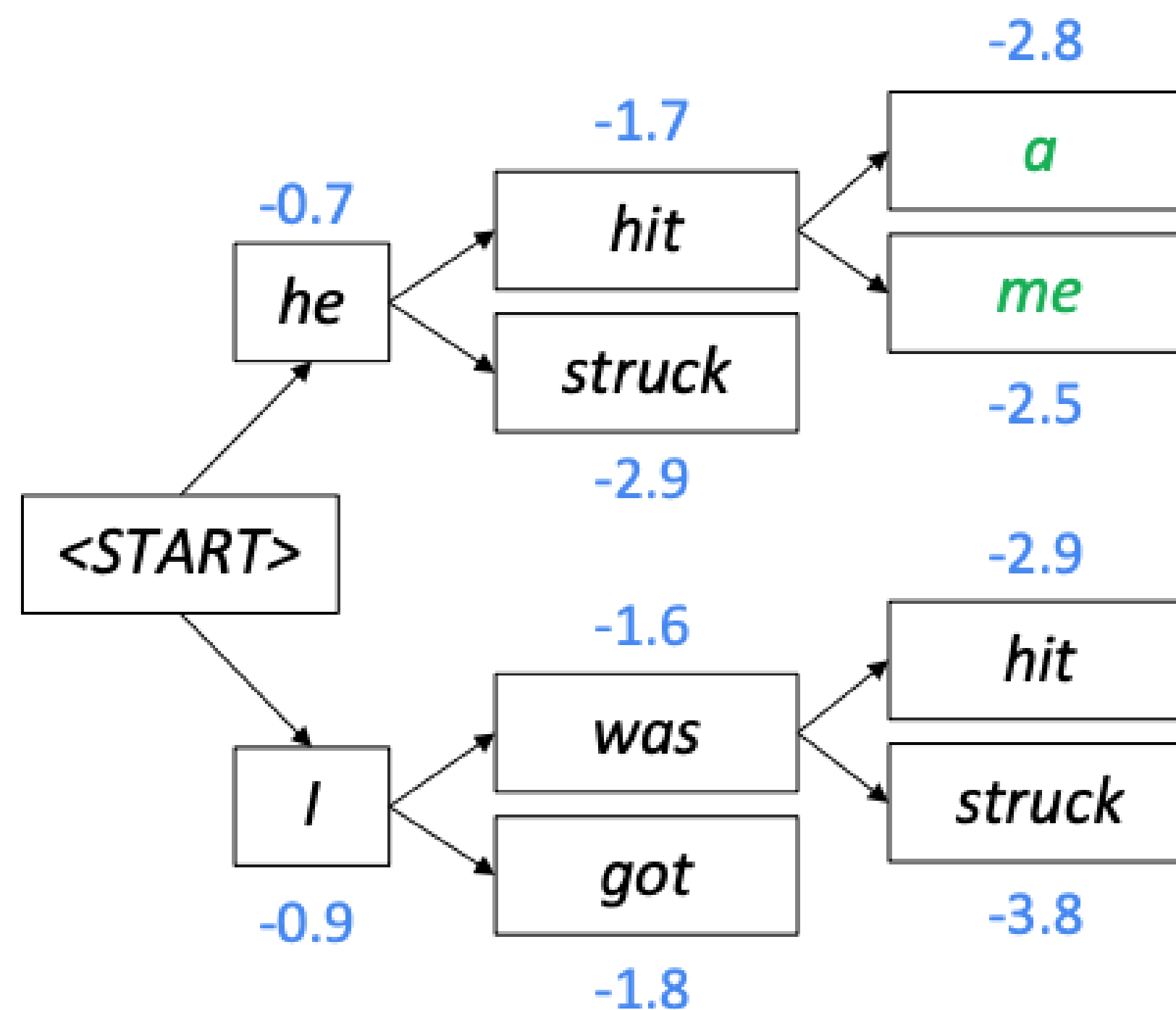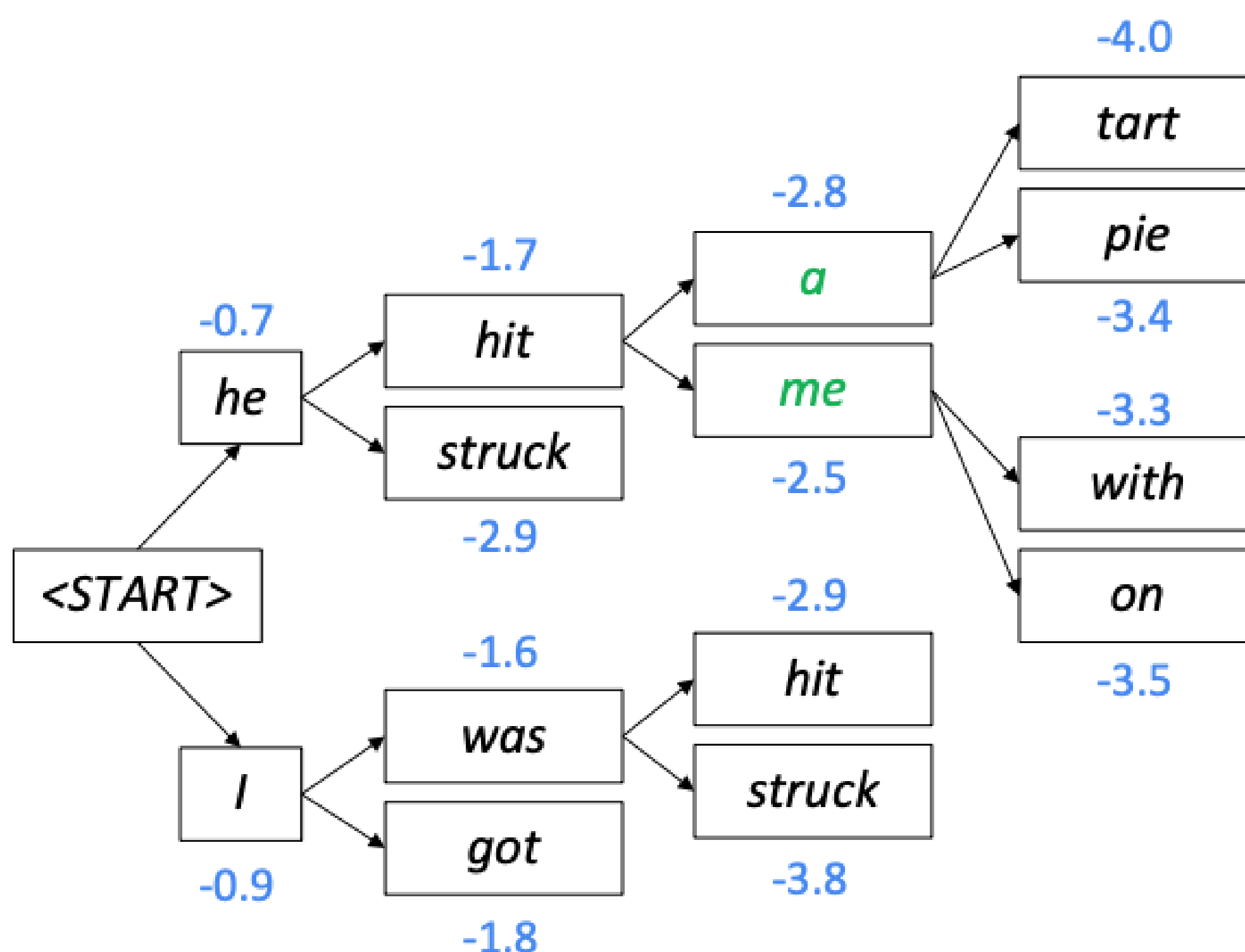


-2.8 = log $P_{\text{LM}}(a | \text{<START> he hit}) + -1.7$

-1.7
a

-0.7
hit

he

struck
me

-2.9
-2.5 = log $P_{\text{LM}}(me | \text{<START> he hit}) + -1.7$

<START>

-2.9 = log $P_{\text{LM}}(hit | \text{<START> I was}) + -1.6$

-1.6
was
hit

I

got
struck

-0.9
-3.8 = log $P_{\text{LM}}(struck | \text{<START> I was}) + -1.6$

-1.8

For each of the *k* hypotheses, find
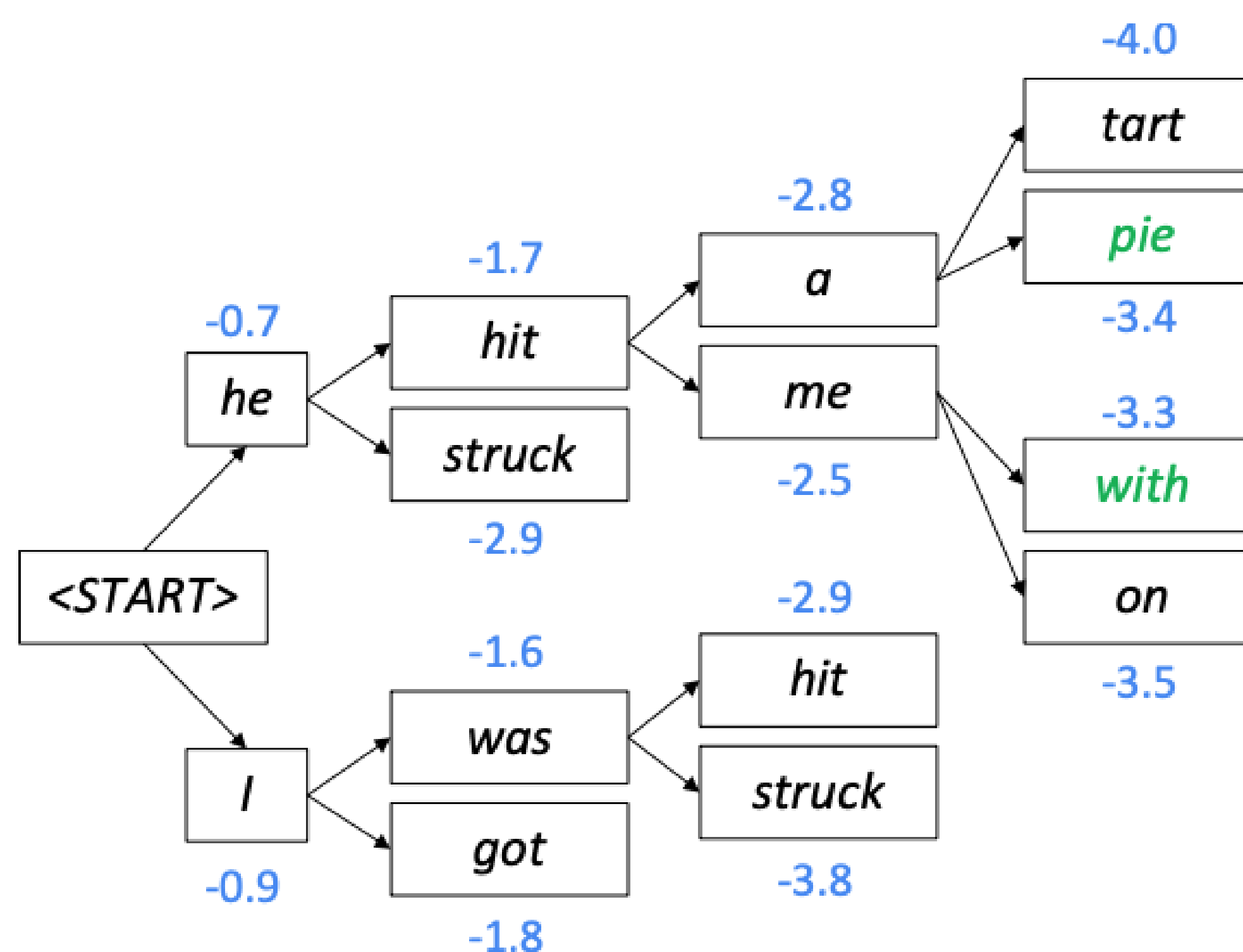top *k* next words and calculate scores

# Beam Search Decoding: Example

Beam size = k = 2. Blue numbers = $\mathrm{score}(y_1, \ldots, y_t) = \sum_{i=1}^{t} \log P_{\mathrm{LM}}(y_i | y_1, \ldots, y_{i-1}, x)$



Of these $k^2$ hypotheses,
just keep $k$ with highest scores

50

Slide credit: Chris Manning

USC Viterbi

# Beam Search Decoding: Example

Beam size = k = 2. Blue numbers = $\mathrm{score}(y_1, \ldots, y_t) = \sum_{i-1}^{t} \log P_{\mathrm{LM}}(y_i | y_1, \ldots, y_{i-1}, x)$



For each of the *k* hypotheses, find
top *k* next words and calculate scores

Slide credit: Chris Manning
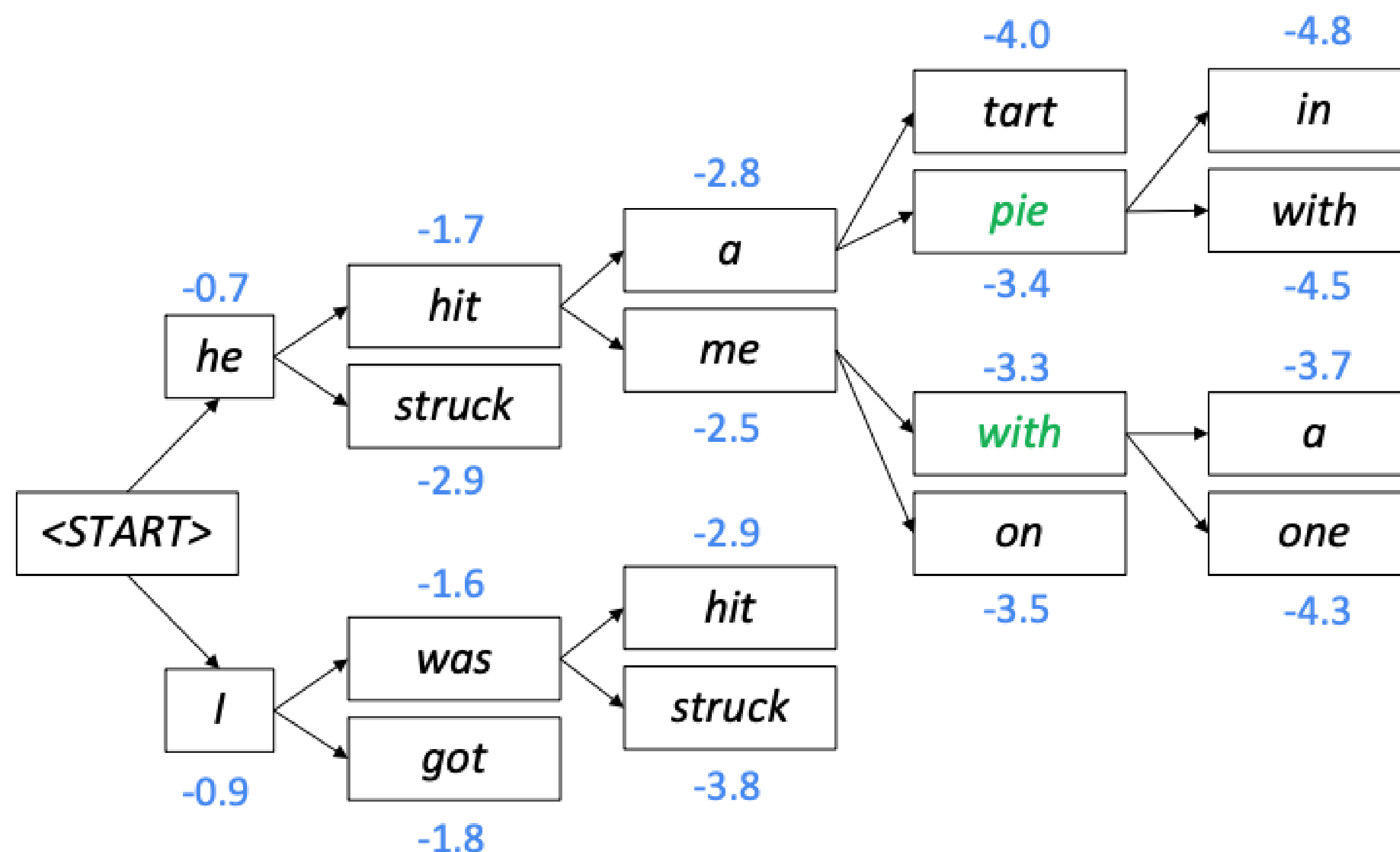
# Beam Search Decoding: Example

Beam size = k = 2. Blue numbers = $\text{score}(y_1, \ldots, y_t) = \sum_{i=1}^{t} \log P_{\text{LM}}(y_i | y_1, \ldots, y_{i-1}, x)$



Of these $k^2$ hypotheses,
just keep $k$ with highest scores

52

# Beam Search Decoding: Example

Beam size = k = 2. Blue numbers = $\mathrm{score}(y_1, \ldots, y_t) = \sum_{i=1}^{t} \log P_{\mathrm{LM}}(y_i | y_1, \ldots, y_{i-1}, x)$



For each of the *k* hypotheses, find top *k* next words and calculate scores

Slide credit: Chris Manning

# Beam Search Decoding: Example

Beam size = k = 2. Blue numbers = $\text{score}(y_1, \ldots, y_t) = \sum_{i=1}^{t} \log P_{\text{LM}}(y_i | y_1, \ldots, y_{i-1}, x)$
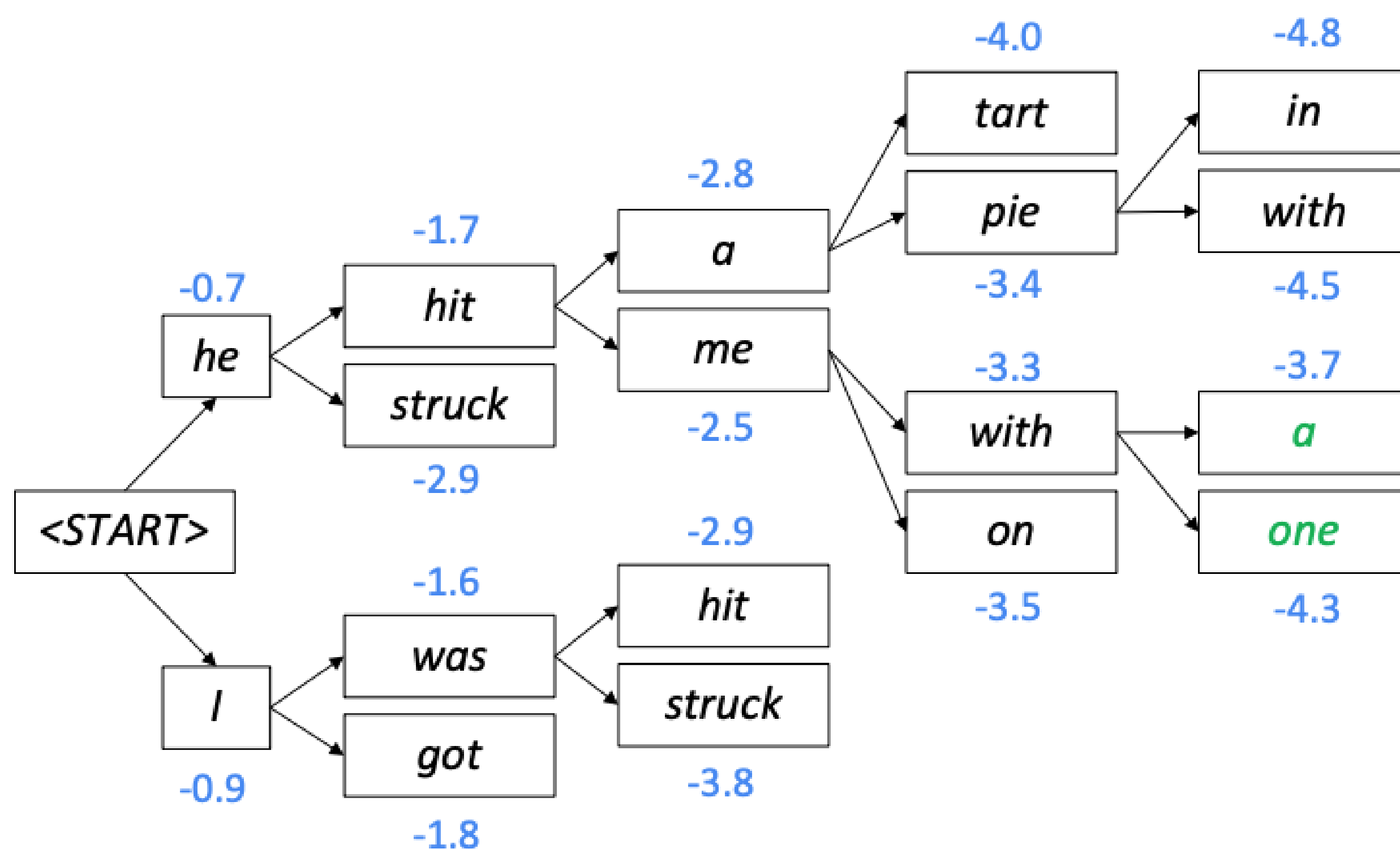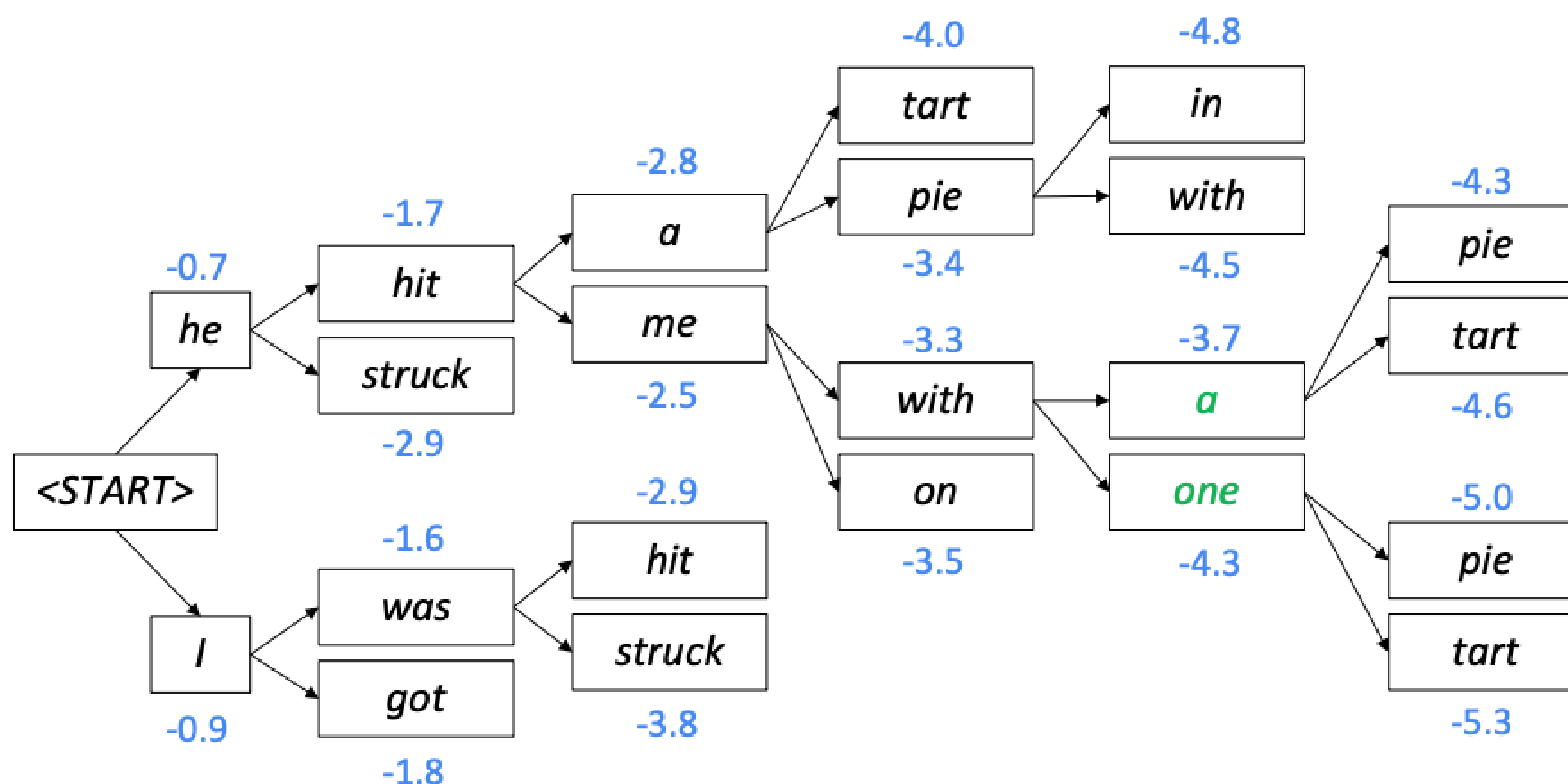


Of these $k^2$ hypotheses,
just keep $k$ with highest scores

Slide credit: Chris Manning
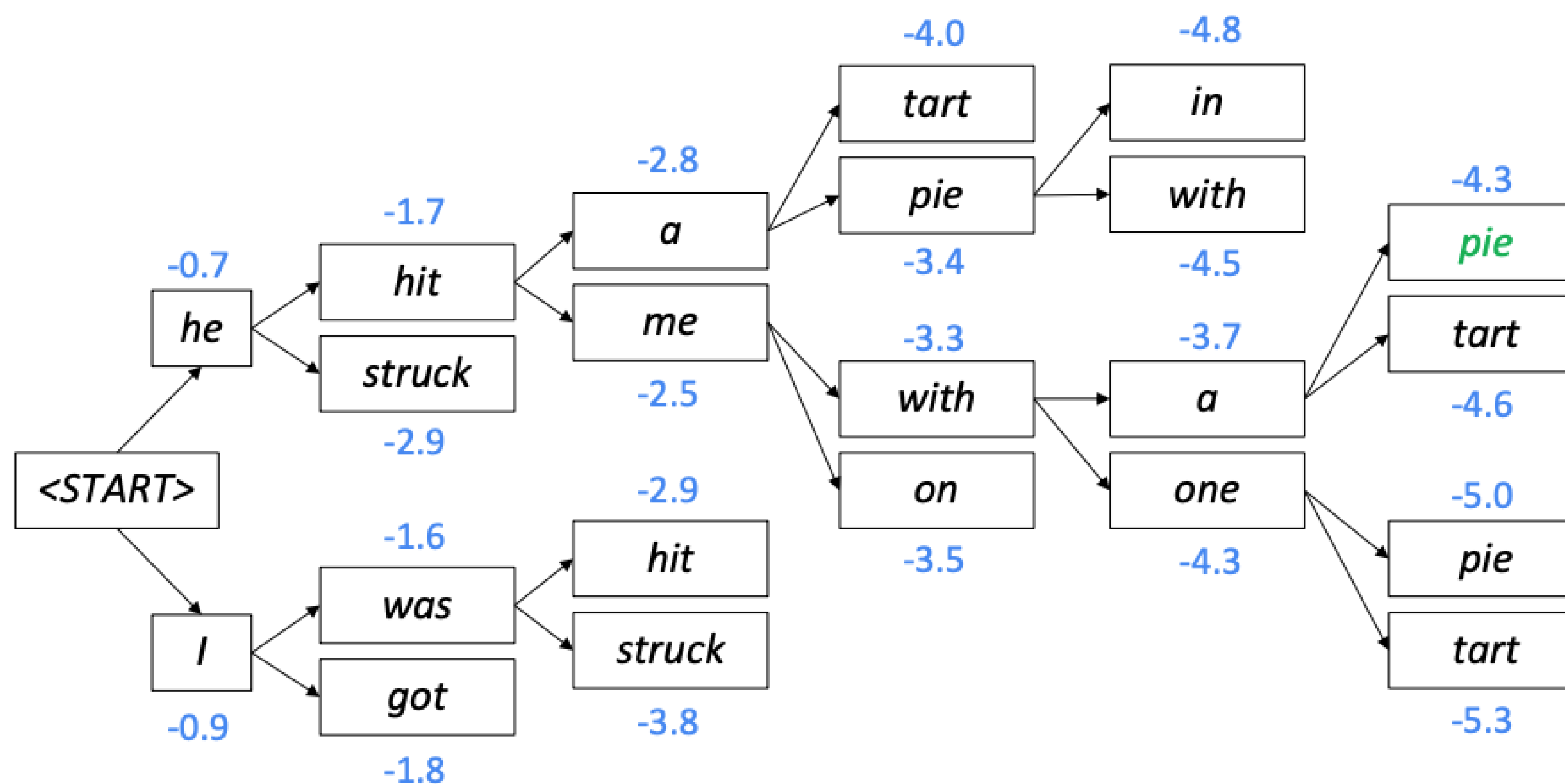
# Beam Search Decoding: Example

Beam size = k = 2. Blue numbers = $\text{score}(y_1, \ldots, y_t) = \sum_{i=1}^{t} \log P_{\text{LM}}(y_i | y_1, \ldots, y_{i-1}, x)$



For each of the *k* hypotheses, find
top *k* next words and calculate scores

Slide credit: Chris Manning
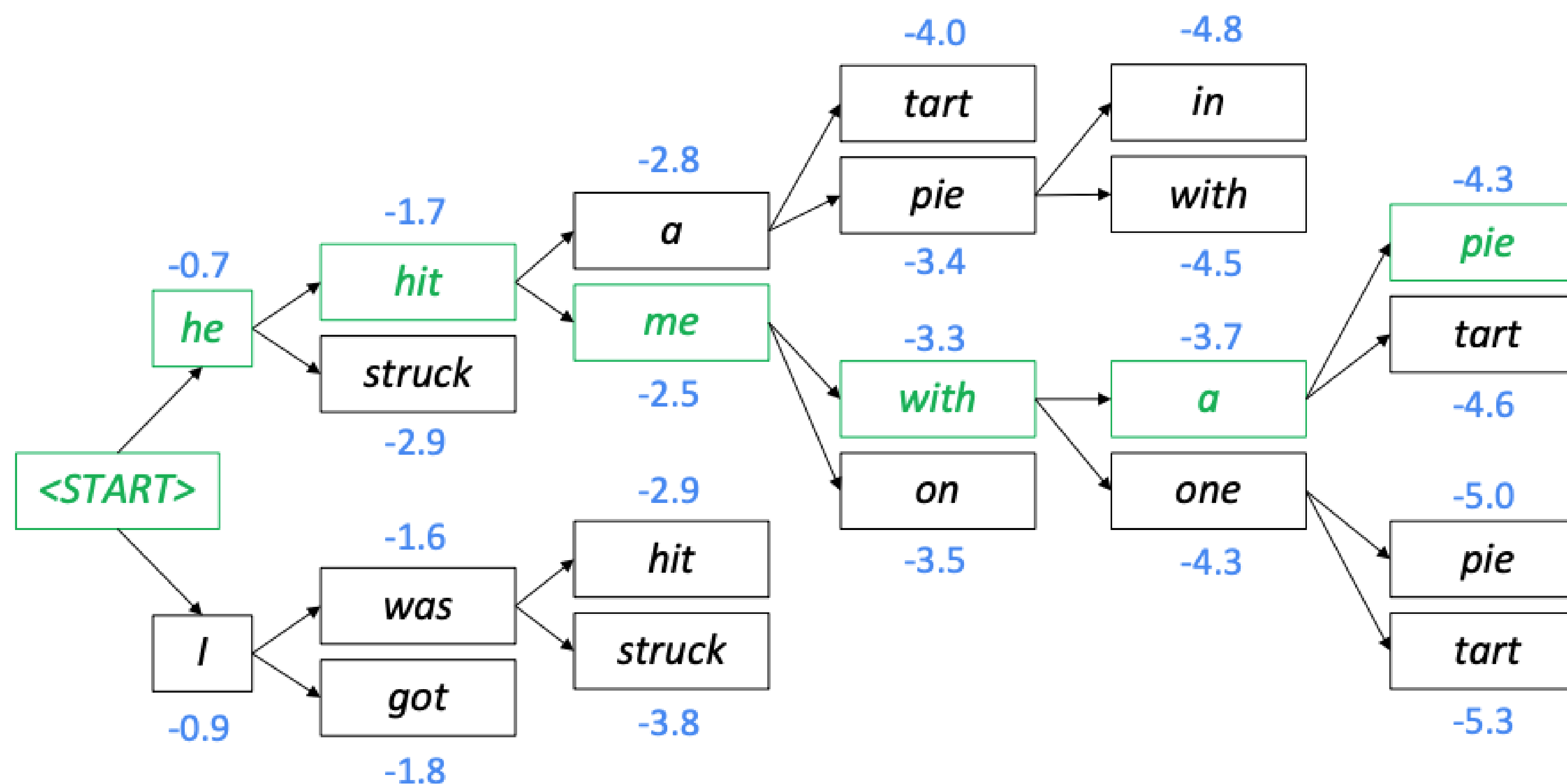
USC Viterbi

# Beam Search Decoding: Example

Beam size = k = 2. Blue numbers = $\text{score}(y_1, \ldots, y_t) = \sum_{i=1}^{t} \log P_{\text{LM}}(y_i|y_1, \ldots, y_{i-1}, x)$



This is the top-scoring hypothesis!

Slide credit: Chris Manning

# Beam Search Decoding: Example

Beam size = k = 2. Blue numbers = $\mathrm{score}(y_1, \ldots, y_t) = \sum_{i=1}^{t} \log P_{\mathrm{LM}}(y_i | y_1, \ldots, y_{i-1}, x)$



Backtrack to obtain the full hypothesis

Slide credit: Chris Manning

# Beam Search Decoding: Stopping Criterion

- Greedy Decoding is done until the model produces an </s> token

  - For e.g. <s> he hit me with a pie </s>

- In Beam Search Decoding, different hypotheses may produce </s> tokens at different time steps

  - When a hypothesis produces </s>, that hypothesis is complete.

  - Place it aside and continue exploring other hypotheses via beam search.

- Usually we continue beam search until:

  - We reach time step T (where T is some pre-defined cutoff), or

  - We have at least n completed hypotheses (where n is pre-defined cutoff)

# Beam Search Decoding: Parting Thoughts

- We have our list of completed hypotheses. Now how to select top one?
- Each hypothesis $y_1, \ldots, y_t$ on our list has a score

$$\text{score}(y_1, \ldots, y_t) = \log P_{\text{LM}}(y_1, \ldots, y_t | x) = \sum_{i=1}^{t} \log P_{\text{LM}}(y_i | y_1, \ldots, y_{i-1}, x)$$

- Problem with this: longer hypotheses have lower score

- Fix: Normalize by length. Use this to select top one instead

But this is expensive!

$$\frac{1}{t} \sum_{i=1}^{t} \log P_{\text{LM}}(y_i | y_1, \ldots, y_{i-1}, x)$$

# Issues with Maximization-Based Decoding

- Either greedy or beam search

- Beam search can be more effective with large beam width, but also more expensive

- Another key issue:

Generation can be bland or repetitive (also called degenerate)

**Context:** In a shocking finding, scientist discovered a herd of unicorns living in a remote, previously unexplored valley, in the Andes Mountains. Even more surprising to the researchers was the fact that the unicorns spoke perfect English.

**Continuation:** The study, published in the Proceedings of the National Academy of Sciences of the United States of America (PNAS), was conducted by researchers from the **Universidad Nacional Autónoma de México (UNAM)** and **the Universidad Nacional Autónoma de México (UNAM/Universidad Nacional Autónoma de México/ Universidad Nacional Autónoma de México/ Universidad Nacional Autónoma de México/ Universidad Nacional Autónoma de México...**

Holtzmann et al., 2020

USC Viterbi

# Modern Generation: Sampling

Next Class!