# CSCI 5541: Natural Language Processing

**Lecture 7: Language Models: Search and Decoding Algorithms**

computer science
& engineering

MINNESOTA · NLP · EST. 2021

UNIVERSITY OF MINNESOTA
Driven to Discover®

# Announcement (0923)

❑ Project Brainstorm (due: Sep 25)
  o Receive proper review by instructors for your proposal pitch on Oct 7–9.
  o Reviews of project brainstorm will be released after TAs/myself review them.
  o Reviews of your brainstorming will consist of the following
    ✓ Which ideas are best to pursue. Suggestions on how to better pursue them
    ✓ Who your 2 mentors are.
    ✓ Which group you are a part of (A or B). This will dictate which days you present the proposal pitch and the final presentation (next slide)

❑ Specific project guideline will be on Sep 30

# Announcement (0923)

❑ Project Proposal Pitch

- o To be held week after next week (Oct 7 and 9)
- o ~3mins discussion of topic, ~5mins of questions and follow-up
- o Groups assigned to Group A will present on Oct 7
- o Groups assigned to Group B will present on Oct 9
- o Before the presentation **you must** upload a slide describing your pitch which includes discussion on the comments we present to your initial brainstorming

| Oct 7 | Project Proposal Pitch (1) | Slides Deck for Group A • |
|-------|----------------------------|--------------------------|
| Oct 9 | Project Proposal Pitch (2) | Slides Deck for Group B • |

# Pitch Slide Template

**Idea Name**
Name 1, Name 2, …..

This is a template slide.
Don't delete or move.

Team Name/Mentor 1, Mentor 2

**Problem Definition**

Just an example

**Data/Methods/etc.**

Just an example

**Plan Forward / Preliminary Results if Any**

Just an example. Feel free to add pictures etc!

**Some questions for your audience**

Just an example

# Outline

❑ Review

❑ Search

    ○ Basics

    ○ Greedy Search

    ○ Beam Search

    ○ Fixing Model Errors in Search

❑ Sampling

    ○ Top-k Sampling

    ○ Top-p Sampling

❑ Search in Training

# Review
## (*N-Grams* to *Neural LMs* to *RNNS* to *LSTMS* to *Seq2Seq*

# Estimation from data

Uni-gram                Bi-gram                Tri-gram

$$\prod_{i=1}^{n} P(w_i)$$
$$\times P(STOP)$$

$$\prod_{i=1}^{n} P(w_i \mid w_{i-1})$$
$$\times P(STOP \mid w_n)$$

$$\prod_{i=1}^{n} P(w_i \mid w_{i-2}, w_{i-1})$$
$$\times P(STOP \mid w_{n-1} w_n)$$

Use the counts of words, pairs of words and groups of three words

$$\frac{c(w_i)}{N}$$

$$\frac{c(w_{i-1}, w_i)}{c(w_{i-1})}$$

$$\frac{c(w_{i-2}, w_{i-1}, w_i)}{c(w_{i-2}, w_{i-1})}$$

# Neural LM

Simple feed-forward multilayer perceptron
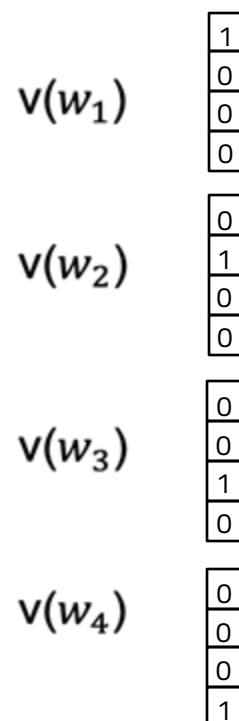(e.g., one hidden layer)

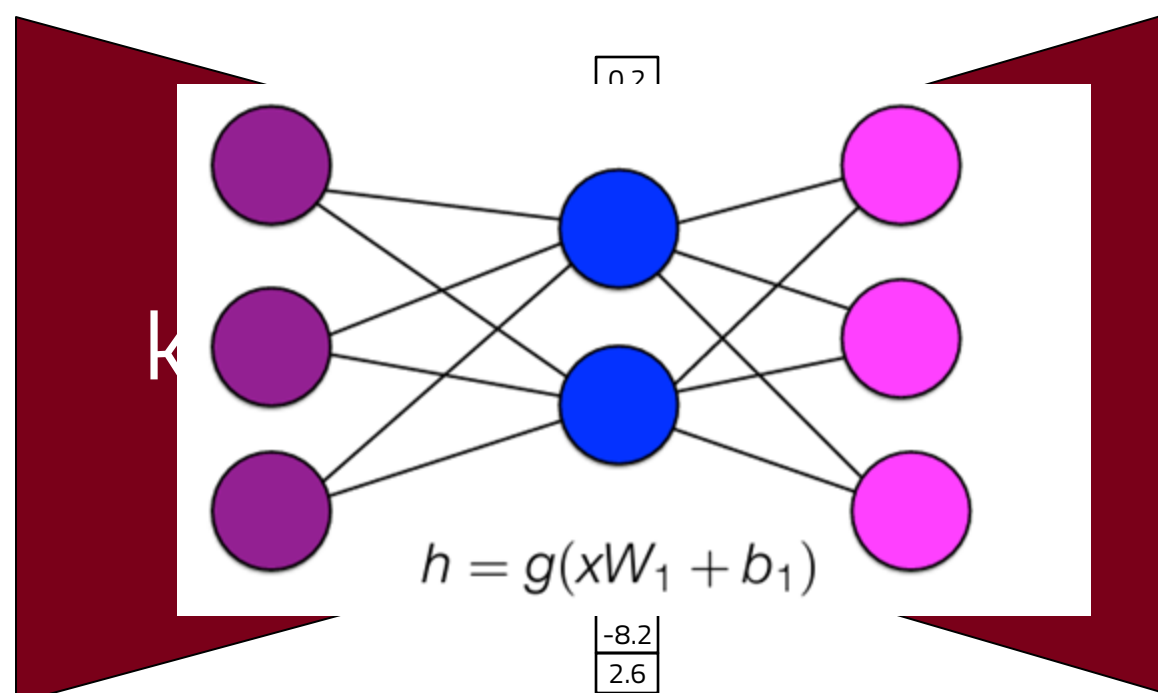$x = [v(w_1); \ldots v(w_k)]$

Concatenation (k x V)

$w_1$ = tried

$w_2$ = to

$w_3$ = prepare

$w_4$ = midterms

$v(w_1)$ $\begin{array}{c}1\\0\\0\\0\\0\end{array}$

$v(w_2)$ $\begin{array}{c}0\\1\\0\\0\end{array}$

$v(w_3)$ $\begin{array}{c}0\\0\\1\\0\end{array}$

$v(w_4)$ $\begin{array}{c}0\\0\\0\\0\\1\end{array}$

One-hot encoding

$h = g(xW_1 + b_1)$

$\begin{array}{c}-8.2\\2.6\end{array}$

Distributed representation
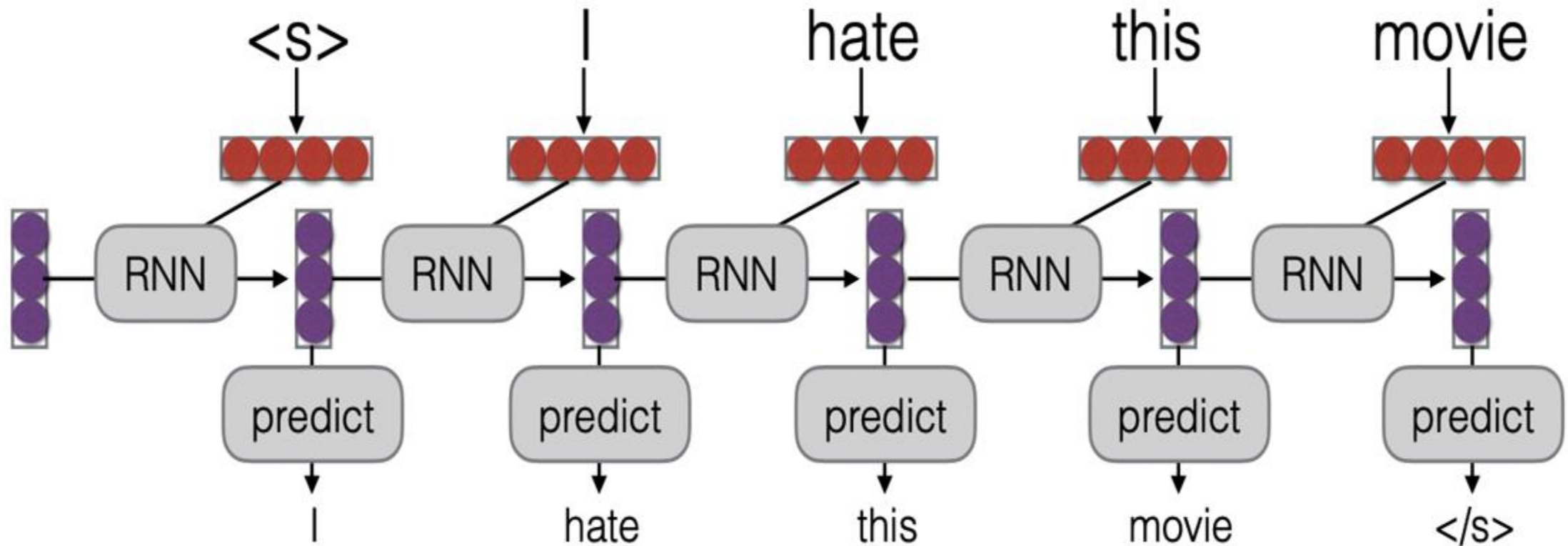
Multi-class (Vocab)
classification

Bengio et al. 2003, *A Neural Probabilistic Language Model*
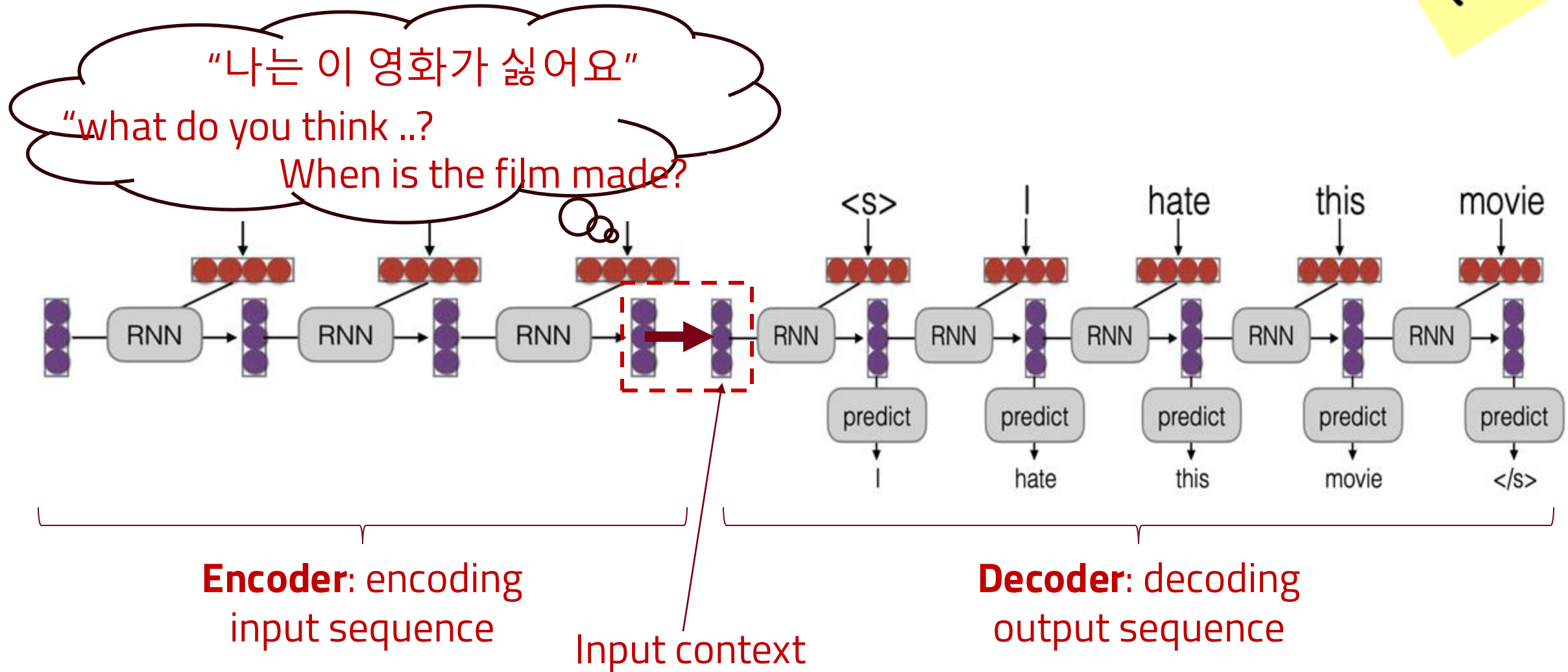
# RNN (Recurrent Neural Network)

❑ Language modeling is like a tagging task, where each tag is the next word!

# LSTMs (Long Short Term Memory)

❑ The Cell State is an information highway

❑ Gradient can flow over this without nearly as many issues of vanishing/exploding gradients that we saw in RNNs

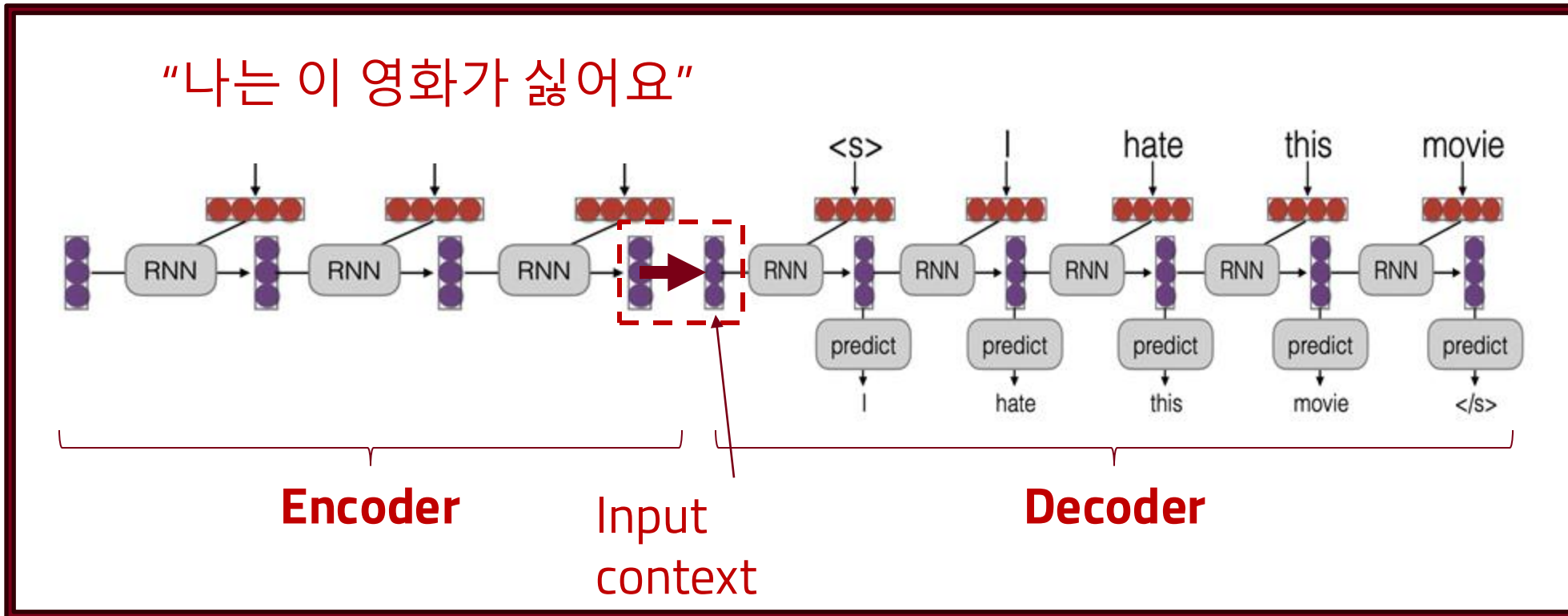❑ We are doing a better job at reducing the 'distance' between our loss function and each individual parameter

# Seq2Seq (Encoder-Decoder)



**Encoder**: encoding input sequence

Input context

**Decoder**: decoding output sequence

# Peeking ahead to Transformers

The input context serves as a significant bottleneck. Most modern language models (transformers) implement some improvements upon this → We'll revisit this in the coming weeks



"나는 이 영화가 싫어요"

**Encoder**

Input context

**Decoder**

# State-of-the-art Language Models

# Teaser: Transformer-based LMs

❑ SOTA LMs: GPT-2, Radford et al. 2018;  GPT-3, Brown et al. 2020

| Trigram | LSTM |
|---------|------|
| 109 | 58.3 |

| GPT-2 | GPT-3 |
|-------|-------|
| 35.8 | 20.5 |



Figure 1: The Transformer - model architecture.

# Teaser: Two Objectives for Language Model Pretraining

GPT GPT2 GPT3 ● ● ●

ELMo BERT ● ●

Auto-regressive LM (GPT3)

Denoising autoencoding (BERT)



$$\log p(\mathbf{x}) = \sum_{t=1}^{T} \log p(x_t | \mathbf{x}_{<t})$$

$$\log p(\bar{\mathbf{x}} | \hat{\mathbf{x}}) = \sum_{t=1}^{T} \text{mask}_t \log p(x_t | \hat{\mathbf{x}})$$

Next-token prediction

Reconstruct masked tokens

Slides from Zihang Dai

Why better language models are useful?

Language models can directly encode knowledge present in the training corpus.

The director of 2001: A Space Odyssey is _____

# Language models can directly encode knowledge present in the training corpus.

| Query | Answer | Generation |
|---|---|---|
| Francesco Bartolomeo Conti was born in ____. | Florence | Rome [-1.8] , **Florence** [-1.8] , Naples |

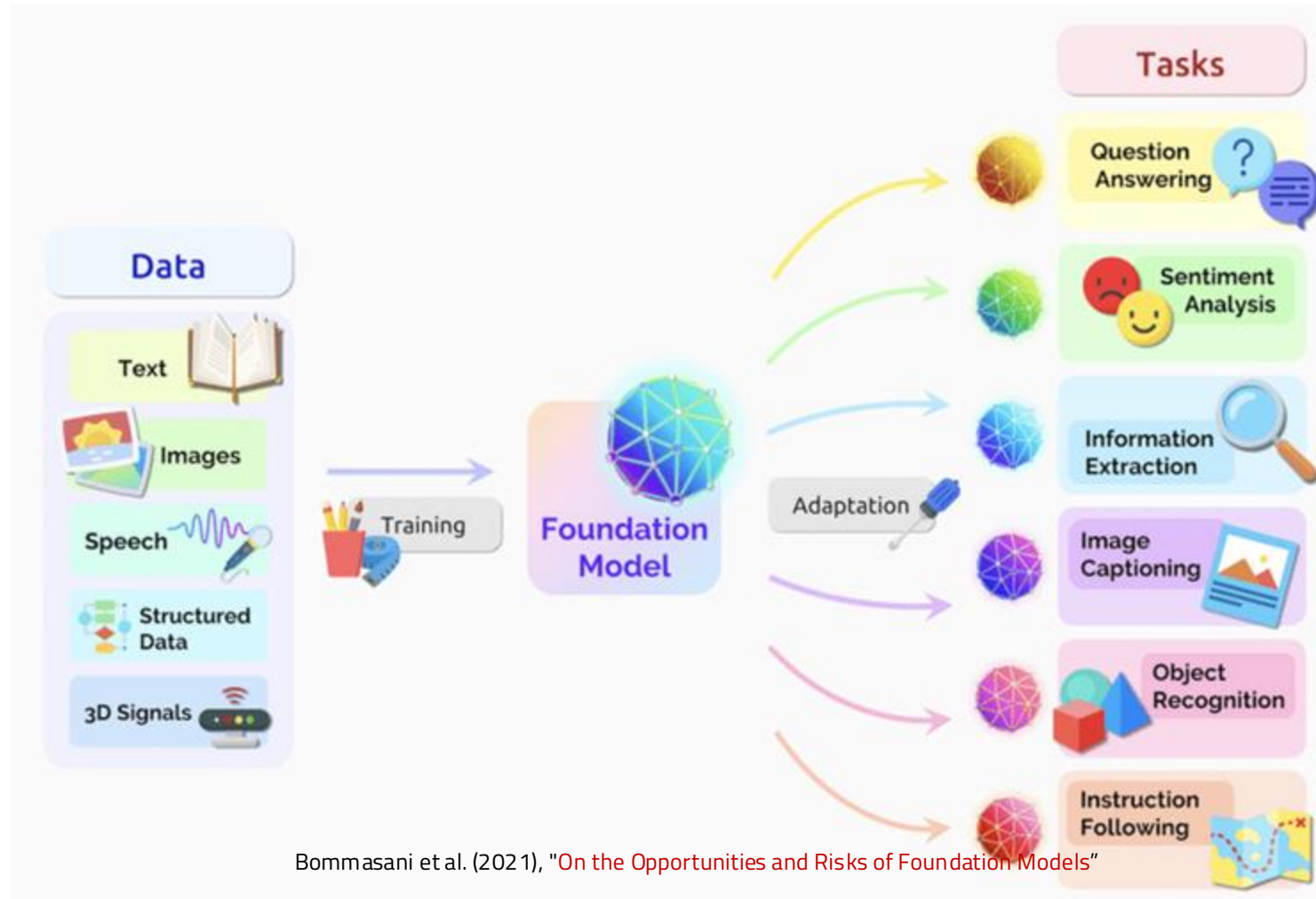Petroni et al. (2019), "Language Models as Knowledge Bases?" (ACL)

# Language models can directly encode knowledge present in the training corpus.

| Query | Answer | Generation |
|---|---|---|
| Francesco Bartolomeo Conti was born in ____ . | Florence | Rome [-1.8] , **Florence** [-1.8] , Naples |
| Adolphe Adam died in ____ . | Paris | **Paris** [-0.5] , London [-3.5] , Vienna |
| English bulldog is a subclass of ____ . | dog | dogs [-0.3] , breeds [-2.2] , **dog** |
| The official language of Mauritius is ____ . | English | **English** [-0.6] , French [-0.9] , Arabic |
| Patrick Oboya plays in ____ position. | midfielder | centre [-2.0] , center [-2.2] , **midfielder** |
| Hamburg Airport is named after ____ . | Hamburg | Hess [-7.0] , Hermann [-7.1] , Schmidt |

Petroni et al. (2019), "Language Models as Knowledge Bases?" (ACL)

# Language models can be a foundation for various tasks across different modalities



Bommasani et al. (2021), "On the Opportunities and Risks of Foundation Models"

# Language models are stochastic parrots



Bender et al. (2021), "On the Dangers of Stochastic Parrots: Can Language Models Be Too Big?"

# Search and Decoding

- *greedy decoding* by calling greedy_search() if num_beams=1 and do_sample=False.

- *multinomial sampling* by calling sample() if num_beams=1 and do_sample=True.

- *beam-search decoding* by calling beam_search() if num_beams>1 and do_sample=False.

- *beam-search multinomial sampling* by calling beam_sample() if num_beams>1 and do_sample=True.

- *diverse beam-search decoding* by calling group_beam_search(), if num_beams>1 and num_beam_groups>1.

- *constrained beam-search decoding* by calling constrained_beam_search(), if constraints!=None or force_words_ids!=None.

https://huggingface.co/docs/transformers/main_classes/text_generation

# Notation

$$P(x_j \mid x_1, \ldots x_{j-1})$$
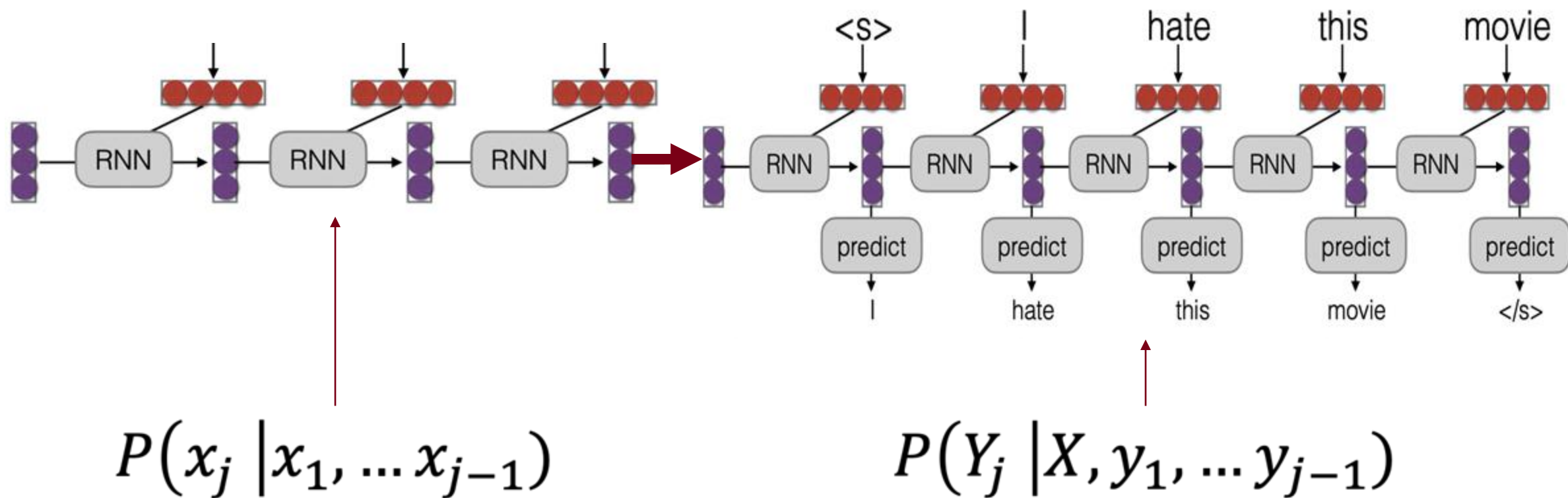
Context given and previous text generated

$$P(Y_j \mid X, y_1, \ldots y_{j-1})$$

Context given in seq2seq setup      Previous text generated

# Notation



$$P(x_j \mid x_1, \ldots x_{j-1})$$

$$P(Y_j \mid X, y_1, \ldots y_{j-1})$$

# Search

# Generation Problem

❑ We have a language model of $P(Y|X)$ trained on text corpora, how do we use it to generate a sentence?

❑ Two methods:

   o We want the best possible single output
     ✓ **Search** (Argmax): Try to generate the sentence with the highest probability.

$$Y_j = argmax \; P(Y_j \, | X, y_1 \dots y_{j-1})$$

   o We want to observe multiple outputs according to the probability distribution
     ✓ **Sampling**: Try to generate a random sentence according to the probability distribution.

$$Y_j = \text{sampling from } P(Y_j \, | X, y_1 \dots y_{j-1})$$

# Generation Problem

❑ We have a language model of $P(Y|X)$ trained on text corpora, how do we use it to generate a sentence?

❑ Two methods:
  - We want the best possible single output
    - ✓ **Search** (Argmax): Try to generate the sentence with the highest probability.

$$Y_j = argmax \, P\left(Y_j \,\middle|\, X, y_1 \dots y_{j-1}\right) \longleftarrow \text{Deterministic}$$

  - We want to observe multiple outputs according to the probability distribution
    - ✓ **Sampling**: Try to generate a random sentence according to the probability distribution.

$$Y_j = \text{sampling from } P\left(Y_j \,\middle|\, X, y_1 \dots y_{j-1}\right) \longleftarrow \text{Probabilistic}$$

# Search Basics

We want to find the **best** output

❑ The **most accurate** output
→ impossible! we don't know the reference

❑ The **most probable** output according to the model
→ simple, but not necessarily tied to accuracy.
Can be computationally demanding

$$\hat{Y} = \underset{\tilde{Y}}{\arg\min} \; \text{error}(Y, \tilde{Y})$$

$$\hat{Y} = \underset{\tilde{Y}}{\arg\max} \; P(\tilde{Y}|X)$$

# Greedy Search

❑ One by one, pick the single highest-probability word

$$While\ Y_{j-1}\ ! = <STOP>$$

$$Y_j = argmax\ P(Y_j\ |X, y_1, \dots y_{j-1})$$

❑ Not exact, real problems:
  o Will often generate the easy words first
  o Will prefer multiple common words to one rare word
  o May not generate highest probability sequence



S = The boy went to the ___

# Greedy methods get repetitive

$$Y_j = argmax \, P(Y_j \,|X, y_1, \dots y_{j-1})$$

**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…**

# Problems w/ Disparate Search Difficulty

$$Y_j = argmax\ P(Y_j\ |X, y_1, \dots y_{j-1})$$

❏ Sometimes need to cover specific content, some easy some hard



|  | | |
|---|---|---|
| I | saw | the escarpment |
| watashi | mita | dangai?  zeppeki?  kyushamen? iwa? |

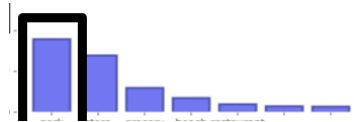❏ Can cause the search algorithm to select the easy thing first, then hard thing later



watashi  wa dangai wo mita
(I saw the escarpment)

watashi  ga mita dangai
(the escarpment I saw)

# Problems w/ Multi-word Sequences

$$Y_j = argmax \, P(Y_j \,|\, X, y_1, \dots y_{j-1})$$

| Next word | P(next word) |
|---|---|
| Pittsburgh | 0.4 |
| New York | 0.3 |
| New Jersey | 0.25 |
| Other | 0.05 |

$P(\text{Pittsburgh}|...) = 0.4$

$P(\text{New}|...) = 0.55$

# Beam Search

❑ Instead of picking the highest probability/score, maintain <span style="color:red">multiple paths</span> (beam size)

❑ At each time step

- Expand each path until <STOP>
- Choose a subset paths from the expanded set

S = The boy went to the ___



Beam size (k) = 2

Blue numbers $= score\ (y_1 \ldots y_t)$

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

# Beam Search

❑ Instead of picking the highest probability/score, maintain multiple paths (beam size)

❑ At each time step
  o Expand each path until <STOP>
  o Choose a subset paths from the expanded set



S = The boy went to the ___



Beam size (k) = 2

Blue numbers = $score\ (y_1 \ldots y_t)$

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

# Beam Search

❑ Instead of picking the highest probability/score, maintain <span style="color:red">multiple paths</span> (beam size)

❑ At each time step
  o Expand each path until <STOP>
  o Choose a subset paths from the expanded set



Beam size (k) = 2

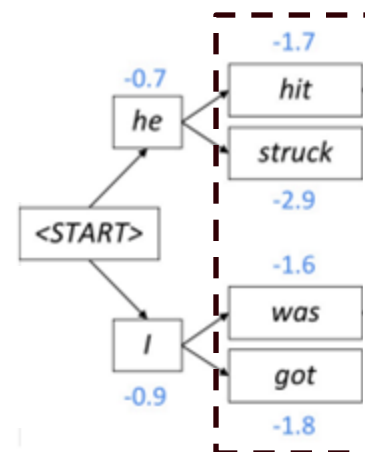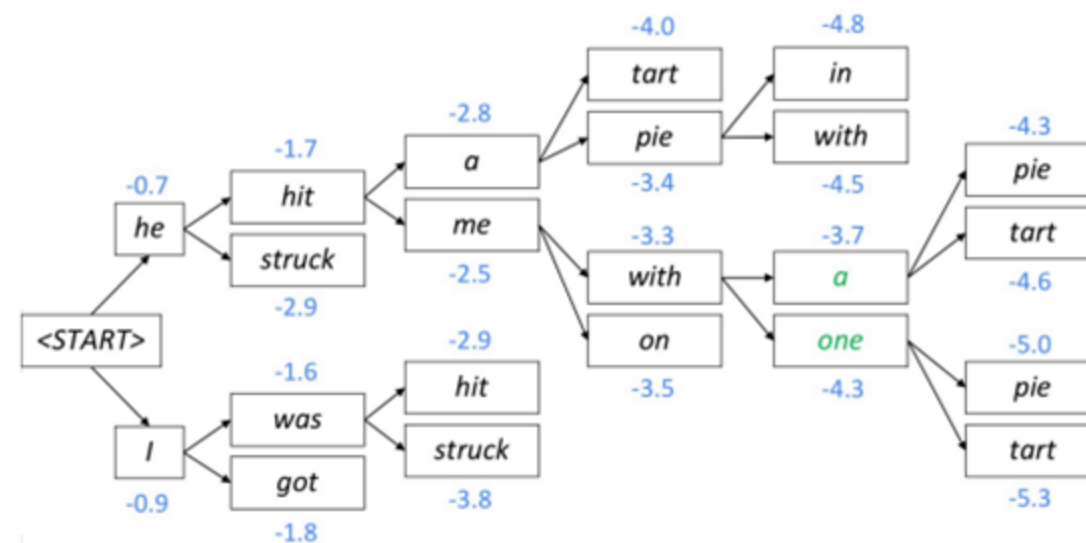Blue numbers = $score\ (y_1 \dots y_t)$

$\qquad\qquad = \prod_{i=1}^{t} \log P_{LM}(y_i | y_1 \dots y_{i-1}, x)$

he hit a tart in ..

he hit a tart in ..

he struck me with a one..

$$score\ (y_1 \dots y_t) = \prod_{i=1}^{t} \log P_{LM}(y_i|y_1 \dots y_{i-1}, x) = -4.0$$

he hit a tart in ..



he struck me with a one..

$$score\ (y_1 \dots y_t) = \prod_{i=1}^{t} \log P_{LM}(y_i|y_1 \dots y_{i-1}, x) = -4.3$$

$$score\ (y_1 \dots y_t) = \prod_{i=1}^{t} \log P_{LM}(y_i|y_1 \dots y_{i-1}, x) = -4.0$$

he hit a tart in ..



he struck me with a one..

$$score\ (y_1 \dots y_t) = \prod_{i=1}^{t} \log P_{LM}(y_i|y_1 \dots y_{i-1}, x) = -4.3$$

# Basic Pruning Methods

How to select which paths to keep expanding?

❑ **Histogram Pruning**: keep exactly $k$ hypotheses at every time step

❑ **Score Threshold Pruning**: keep all hypotheses where score is within a threshold $\alpha$ of best score $s_1$

❑ **Probability Mass Pruning:** keep all hypotheses up until probability mass $\alpha$

# Better Search can Hurt Results! <span>(Koehn and Knowles 2017)</span>

❑ Better search (=better model score) can result in worse BLEU score!

❑ Why? Model errors! (Model is not trained to give better BLEU score but predict better next token)



Czech–English

# Beam Search Curse (Yang et al. 2018)

❑ As beam size increases, it becomes easier for the search algorithm to find the **</eos>** symbol.

❑ Then, shorter candidates have clear advantages w.r.t. model score.

# A Typical Model Error: Length Bias

❑ In many tasks (e.g. Machine translation), the output sequences will be of variable length

❑ Running beam search may then favor short sentences

GSM8K Distances   ⟶   

50                    1.23k

# Length Normalization

❑ Beam search may then favor short sentences
❑ Normalize by the length, dividing by |Y| to prioritize longer sentences.

(Cho et al. 2014)

$$\frac{1}{T_y^\alpha} \quad argmax_y \sum_{j=1}^{T_y} \log P(y_j \mid X, y_1, \dots y_{j-1})$$

$\alpha = [0, 1.0]$

(Wu et al. 2016)

$$\frac{(5 + 1)^\alpha}{(5 + |Y|)^\alpha}$$

# Constrained Decoding

❑ Some tasks (coding/mathematics/synthetic data generation) have an explicit structure

❑ When finite state machines can be attached to your outputs, then you can limit the set of words your model will output from |V| to M, where M is the set of possible next words

```
{
    "type": "object",
    "properties": {
        "name": {
            "type": "string"
        },
        "age": {
            "type": "integer"
        },
        "house": {
            "type": "string",
            "enum": [
                "Gryffindor",
                "Hufflepuff",
                "Ravenclaw",
                "Slytherin"
            ]
        }
    },
    "required": [
        "name",
        "age",
        "house"
    ]
}
```

**JSON Schema**

**+**

**Large Language Models**

**Constrained Generation**

```
{
    "name": "Harry",
    "age": 15,
    "house": "Gryffindor"
}
```

```
{
    "name": "Cedric",
    "age": 18,
    "house": "Hufflepuff"
}
```

```
{
    "name": "Luna",
    "age": 14,
    "house": "Ravenclaw"
}
```

**Generated JSONs**

The generation house can only be one of the four words in the 'house' portion of the schema

# Constrained Decoding

S =
'{

"name": "Jim",

"age": 28,

"house": '



P(W|S)

Slytherin   Minneapolis   St Paul   Hufflepuff   Gryffindor   Boston   San Francisco

# Constrained Decoding

S =

'{

    "name": "Jim",

    "age": 28,

    "house": '



Ignore potential words which do not conform to our sample

# Sampling

How are you doing?

I'm good! How about you?

So so..

It was a hard day for me.

# Recap: Greedy/Beam Search (w/o Sampling)

S = The boy went to the ___

$$W_T = argmax_W(P(W \mid W_{1:T-1}))$$

P(W|S)

Next token [W]: park, store, grocery, beach, restaurant, ..., ....

Beam size (k) = 2



P ("park", "today" | S) = 0.12

0.33 — today
park
0.36
0.25 — yesterday

The boy went to the ___

0.15 — grocery
0.9 — store
0.07 — during

P ("grocery", "store" | S) = 0.135

**Deterministic beam search:**
I went into town on Saturday morning because...
-> I was going to go to the gym and I was going to go to the gym and I was going to go to the ..

# Ancestral Sampling

❑ Randomly generate words one-by-one

- $Y_j = P\left(Y_j \mid X, y_1 \ldots y_{j-1}\right)$
- Until <STOP> is generated

❑ An exact method for sampling from $P(X)$, no further work needed.

X = 'The boy when to the'
Y = ''
Generate $Y_0$

S = The boy went to the ___



https://medium.com/ai2-blog/a-guide-to-language-model-sampling-in-allennlp-3b1239274bc3

# Decoding with Ancestral/Multinomial Sampling

S = The boy went to the ___



**Multinomial Sampling:**

I went into town on Saturday morning because...

-> I have to wear suits and collared in the South Bay. This was shocking!" "This is our city. First of all, I'm strange in the name of Santa, Howard Daniel, and

https://medium.com/ai2-blog/a-guide-to-language-model-sampling-in-allennlp-3b1239274bc3

**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.

**Beam Search, *b*=32**:
"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 …"

Repetition

**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.
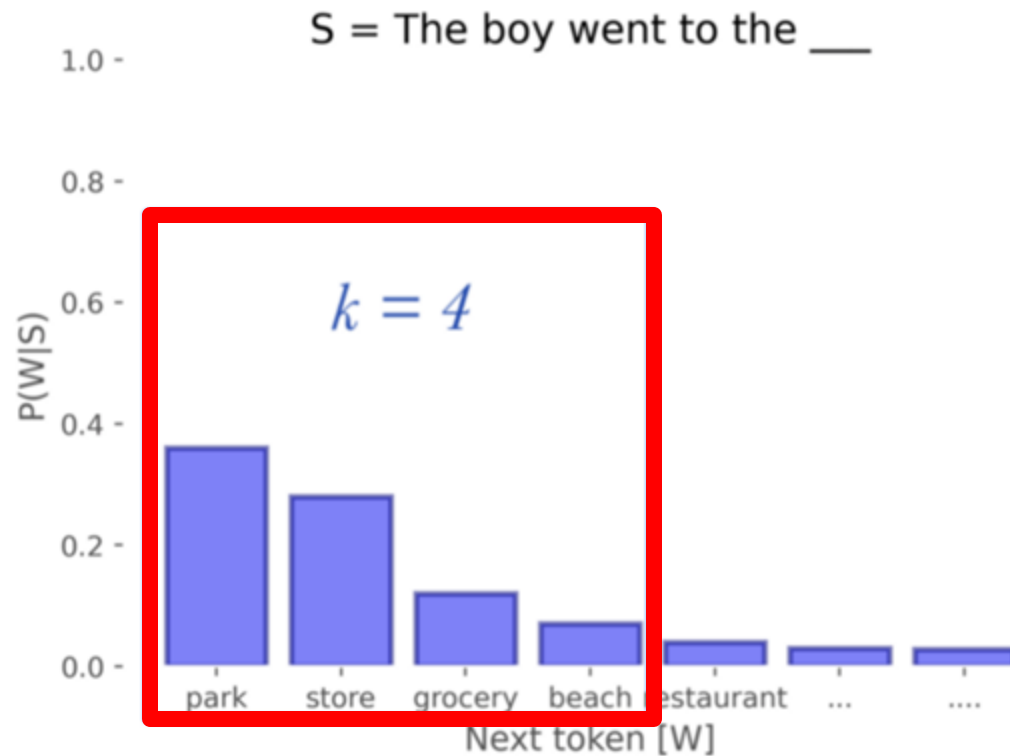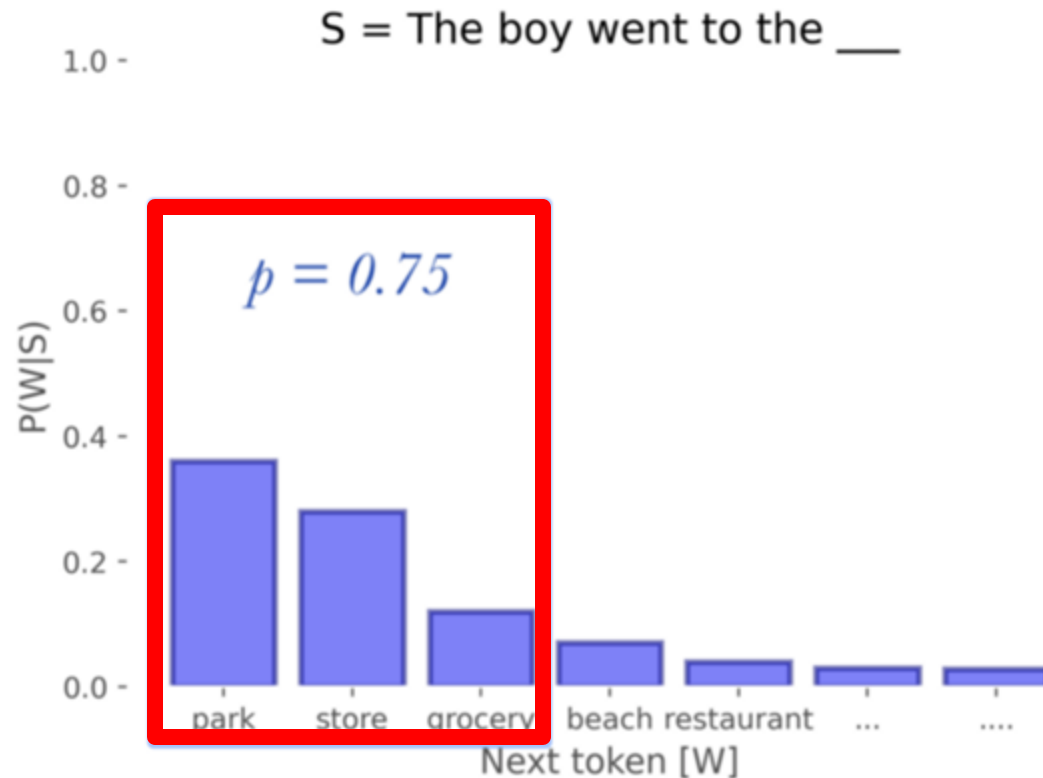


$W_T = argmax_W(P(W \mid W_{1:T-1}))$

**Beam Search, b=32:**
"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 …"

**Pure Sampling:**
They were cattle called Bolivian Cavalleros; they live in a remote desert uninterrupted by town, and they speak huge, beautiful, paradisiacal Bolivian linguistic thing. They say, 'Lunch, marge.' They don't tell what the lunch is," director Professor Chuperas Omwell told Sky News. "They've only been talking to scientists, like we're being interviewed by TV reporters. We don't even stick around to be interviewed by TV reporters. Maybe that's how they figured out that they're cosplaying as the Bolivian Cavalleros."

# Repetition

# Incoherence

# Top-k Sampling

S = The boy went to the ___



❑ Only sample from the *k* most probable tokens, by redistributing the PMF over the top-k tokens

❑ But, picking a good value of *k* can be difficult as the distribution or words is different for each step.
  - Increase k for more diverse/risky outputs
  - Decrease k for more generic/safe outputs

# Top-p Sampling (or Nucleus Sampling) <span>(Holtzman et al. 2020)</span>
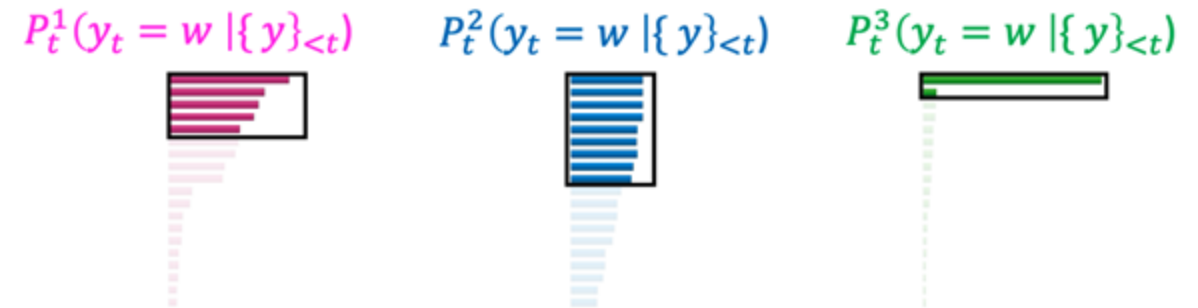
S = The boy went to the ___



❑ Another way to exclude very low probability tokens is to include the most probable tokens that make up the "nucleus" of the PMF

   o the sum of the most probable tokens just reaches P

# Top-p Sampling (or Nucleus Sampling) (Holtzman et al. 2020)

S = The boy went to the ___

$p = 0.75$



- ❑ **Flexible** as the distribution changes, allowing the size of the filtered words to expand and contract when it makes sense.

$P_t^1(y_t = w \,|\, \{y\}_{<t})$    $P_t^2(y_t = w \,|\, \{y\}_{<t})$    $P_t^3(y_t = w \,|\, \{y\}_{<t})$

# Cautions about Sampling-based Search

❏ Is sampling necessary for diversity?

   o questionable, we could do diverse beam search instead.

❏ Results are inconsistent from run-to-run:

   o need to consider variance from this in reporting

   o (in addition to variance in training and data selection)

❏ Conflates model and search errors:

   o if you make a better model you might get worse results, because the search algorithm can't find the outputs your model likes

# Decoding: Takeaways

❑ Many problems in neural NLG are not really problems with our learned language model probability distribution, but problems with the decoding algorithm

❑ Different decoding algorithms can allow us to inject biases that encourage different properties of coherent natural language generation

❑ Some of the most impactful advances in NLG of the last few years have come from simple but effective modifications to decoding algorithms

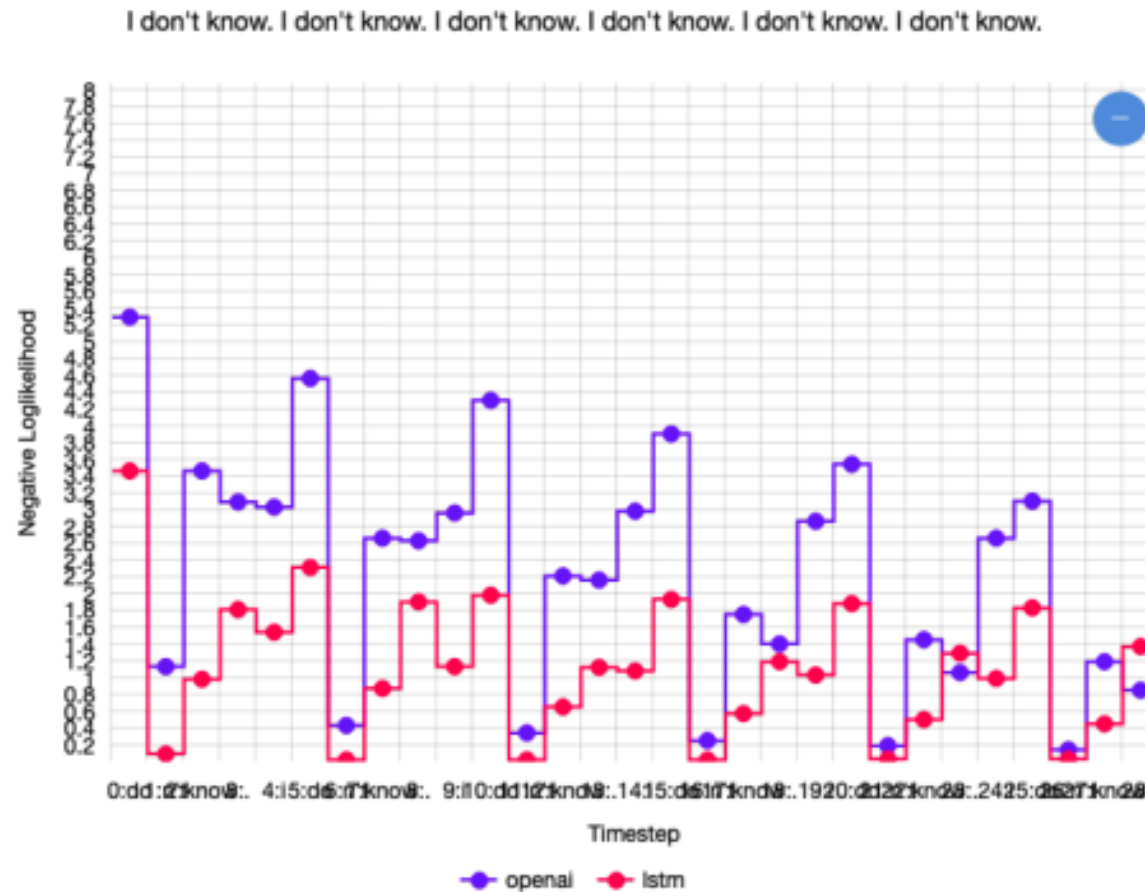# Search in Training

# Diversity Issues (Holtzman et. al., 2020)

❑ Maximum Likelihood Estimation discourages diverse text generation



I don't know. I don't know. I don't know. I don't know. I don't know. I don't know.
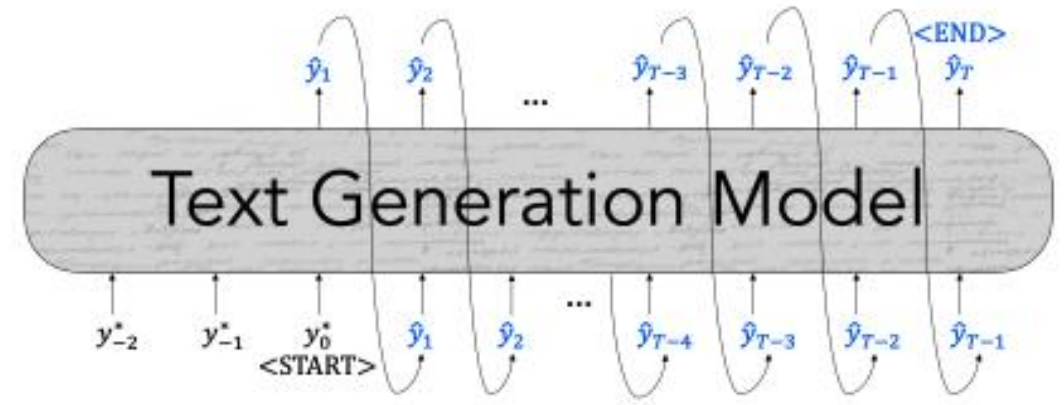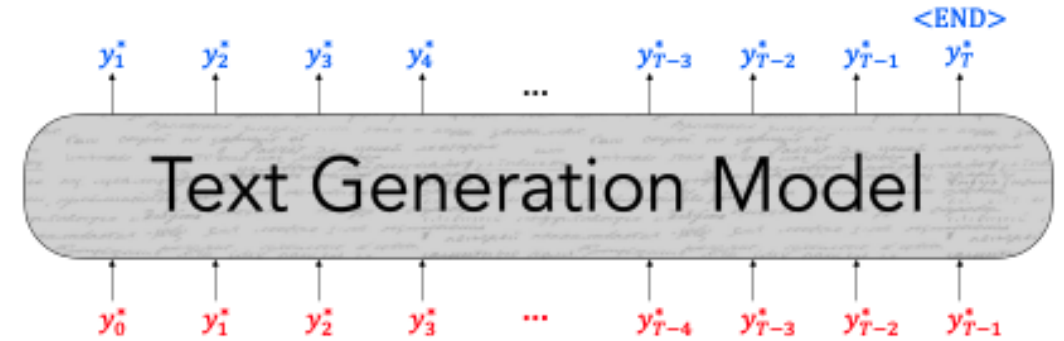
# Why? Exposure Bias

❑ Training with teacher forcing leads to exposure bias at generation time

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

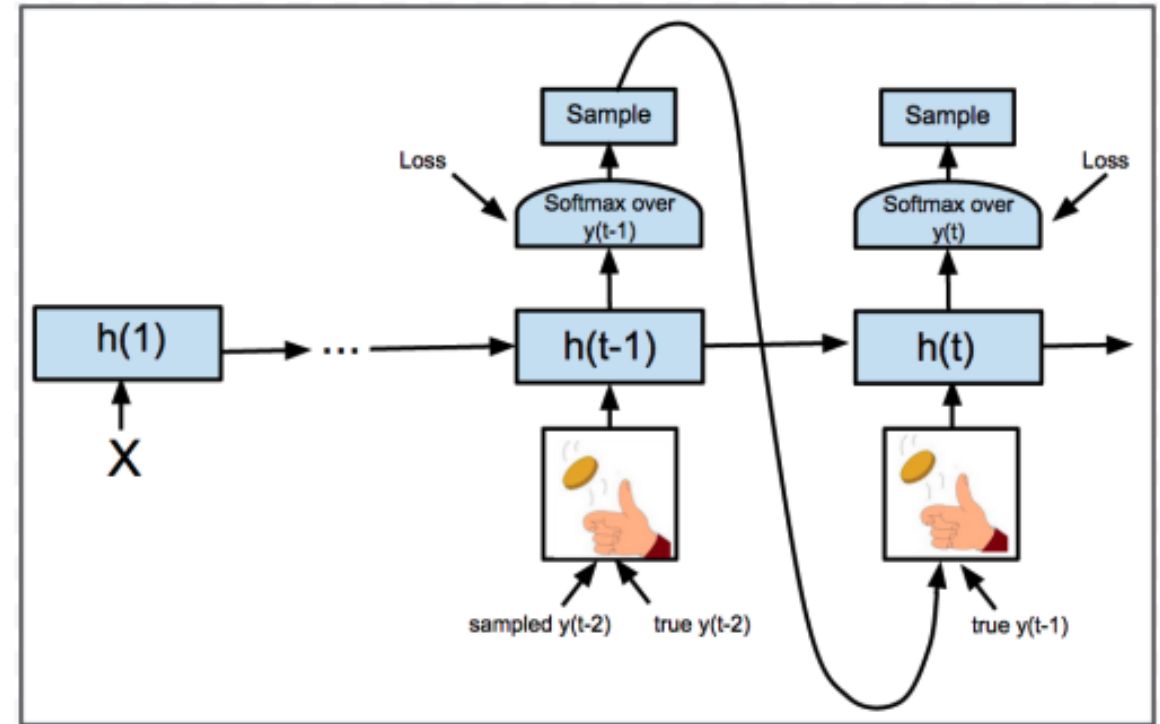$$\mathcal{L}_{MLE} = -\log P(y_t^* | \{y^*\}_{<t})$$

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

$$\mathcal{L}_{dec} = -\log P(\hat{y}_t | \{\hat{y}\}_{<t})$$
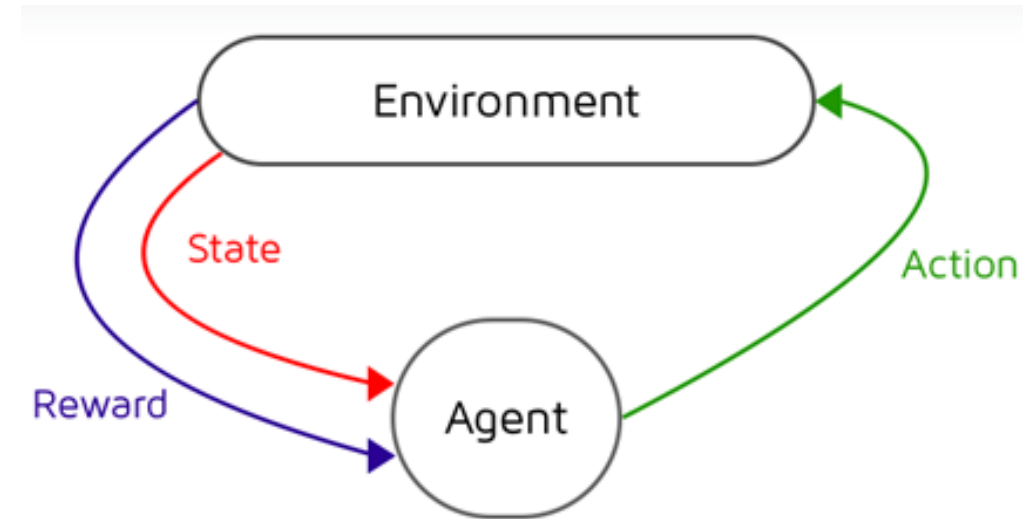
# Fix Exposure Bias: Scheduled sampling (training)

❑ With some probability p, decode a token and feed that as the next input, rather than the gold token.

❑ Increase $p$ over the course of training

❑ Leads to improvements in practice, but can lead to strange training objectives

❑ Also called teacher forcing



(Bengio et al., 2015)

# Fix Exposure Bias: Reinforcement Learning

❑ Cast your text generation model as a Markov decision process

  ➢ **State** $s$ is the model's representation of the preceding context

  ➢ **Actions** $a$ are the words that can be generated

  ➢ **Policy** $\pi$ is the decoder

  ➢ **Rewards** $r$ are provided by an external score

❑ Learn behaviors by rewarding the model when it exhibits them

❑ Use REINFORCE or similar; it's difficult because huge branching factor/search space

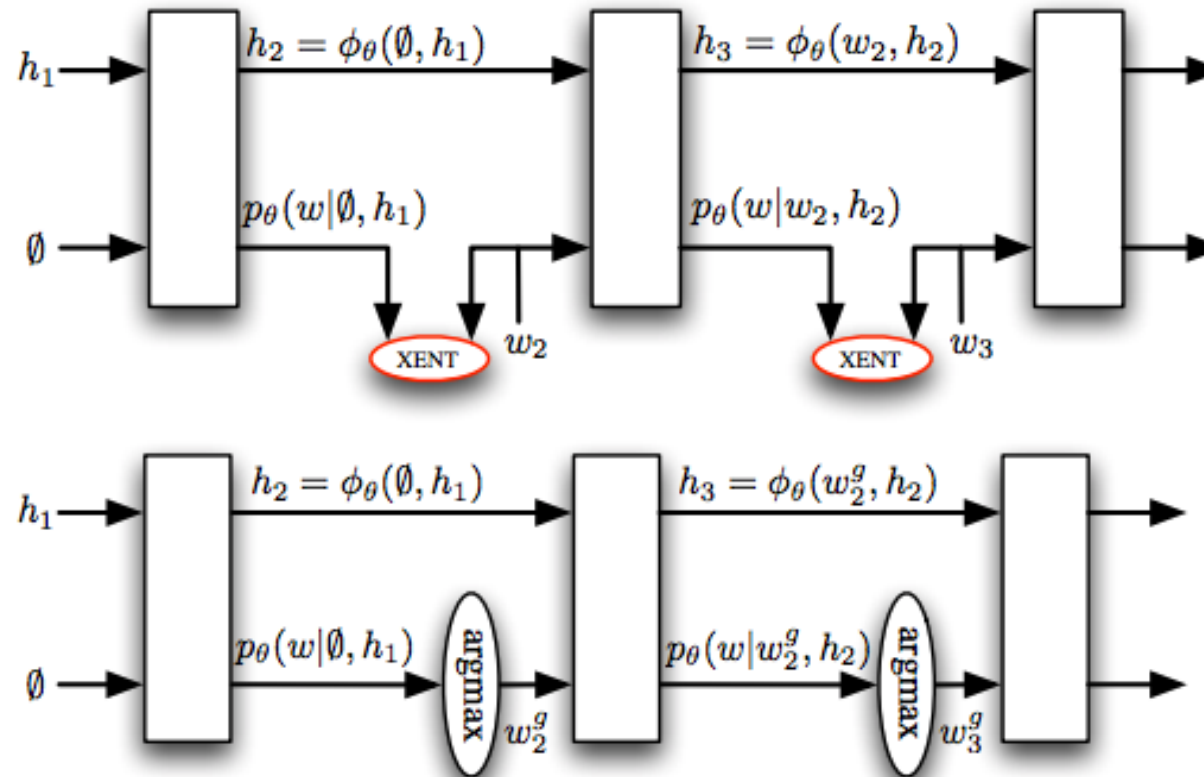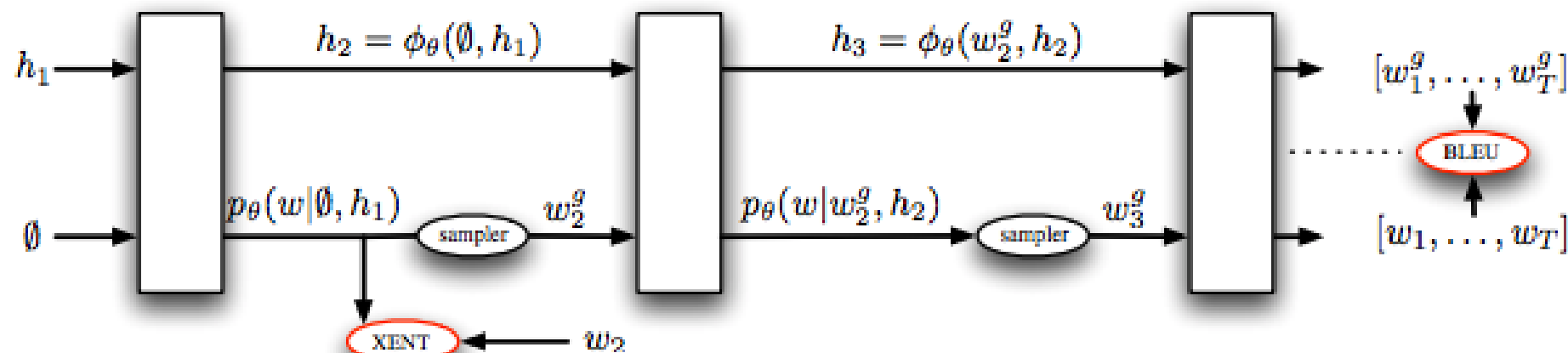# MIXER: Sequence-level training with REINFORCE

Figure 1: RNN training using XENT (top), and how it is used at test time for generation (bottom).

# MIXER: Sequence-level training with REINFORCE

Ranzato et al., 2016



| TASK | XENT | DAD | E2E | MIXER |
|------|------|-----|-----|-------|
| *summarization* | 13.01 | 12.18 | 12.78 | **16.22** |
| *translation* | 17.74 | 20.12 | 17.77 | **20.73** |
| *image captioning* | 27.8 | 28.16 | 26.42 | **29.16** |

$$\mathscr{L} = \mathscr{L}_{MLE} + \alpha\mathscr{L}_{RL}$$

$$\mathscr{L}_{RL} = -\sum_{t=1}^{T}(r(\hat{y}_t) - \boldsymbol{b})\log P(...)$$

MIXER seems to be a useful, agnostic trick to improve MT results, but did not
see wide usage ~ perhaps due to unstability of REINFORCE

# Reward Estimation

$$\mathcal{L}_{RL} = -\sum_{t=1}^{T} (r(\hat{y}_t) - \boldsymbol{b}) \log P(\ldots)$$

❑ How should we define a reward function? Just use your evaluation metric!

- ○ BLEU (machine translation; Ranzato et al., ICLR 2016; Wu et al., 2016)
- ○ ROUGE (summarization; Paulus et al., 2018; Celikyilmaz et al., 2018)
- ○ CIDEr (image captioning; Rennie et al., CVPR 2017)
- ○ SPIDEr (image captioning; Liu et al., ICCV 2017)

❑ Be careful about optimizing for the task as opposed to "gaming" the reward!

- ○ Evaluation metrics are merely proxies for generation quality!
- ○ "*even though RL refinement can achieve better BLEU scores, it barely improves the human impression of the translation quality*" – Wu et al., 2016

# Reward Estimation

❑ What behaviors can we tie to rewards?

- ○ Sentence simplicity (Zhang and Lapata, EMNLP 2017)
- ○ Temporal Consistency (Bosselut et al., NAACL 2018)
- ○ Cross-modality consistency in image captioning (Ren et al., CVPR 2017)
- ○ Utterance Politeness (Tan et al., TACL 2018)
- ○ Paraphrasing (Li et al., EMNLP 2018)
- ○ Sentiment (Gong et al., NAACL 2019)
- ○ Formality (Gong et al., NAACL 2019)

❑ If you can formalize a behavior as a Python function (or train a neural network to approximate it!), you can train a text generation model to exhibit that behavior!

# Search in Training: Takeaways

❏ **Teacher forcing** is still the main algorithm for training text generation models

❏ **Diversity** is an issue with sequences generated from teacher forced models

❏ **Exposure bias** causes text generation models to lose coherence easily

❏ Training with RL can allow models to learn behaviors that are challenging to formalize

    o    But learning can be very **unstable**!

    o    chatGPT: advanced RL algorithms (e.g., PPO) for better human alignment with human feedback

# Other techniques not covered

❑ Decoding time control for controllable text generation (e.g., PPLM)

❑ Multi-attribute control using RL

❑ Unlikelihood training

❑ Data augmentation for reducing the exposure bias

❑ Retrieval-augmented Generation (RAG) (will be covered)

❑ Retrieval based generation (e.g., KNN Language Models)

❑ Instruction tuning and human feedback learning (will be covered)

…