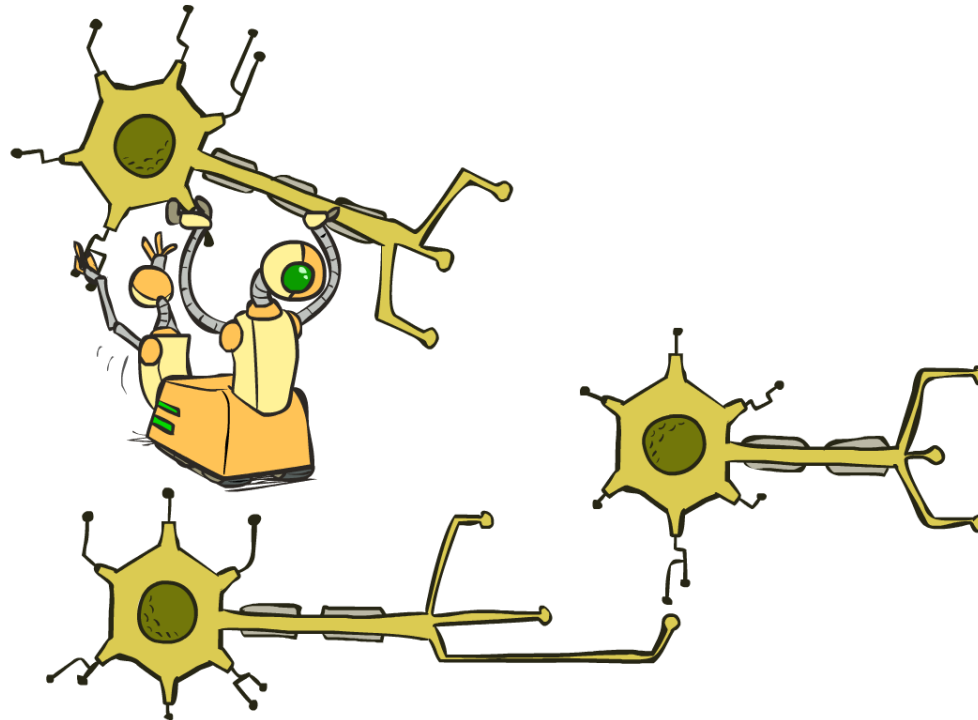


CS 188: Artificial Intelligence

Advanced Topics: Natural Language Processing



Instructors: Eve Fleisig & Evgeny Pobachienko

[Slides courtesy of Dan Klein, Abigail See, Greg Durrett, Yejin Choi, John DeNero, Eric Wallace, Kevin Lin, Fei-Fei Li, Sergey Levine, Pieter Abbeel, Jay Alamar, and many others]

Natural Language Processing



Goal: Deep Understanding

- Requires knowledge, context, and grounding
- Just starting to see successes

Reality: Shallow Matching

- Requires robustness and scale
- Amazing successes, but fundamental limitations

What tasks do we care about?

Applications

- Text generation
- Chatbots
- Machine translation
- Web search
- Speech recognition
- Virtual assistants
- Text-to-speech
- Syntactic parsing, semantic parsing, semantic role labeling, word sense induction, coreference resolution, summarization, sentiment analysis...

Machine Translation

CINÉMA • MÉDIAS Partage 

Disney décide de changer le nom du légendaire studio de cinéma Fox


Le nouveau propriétaire va faire disparaître la mention « Fox » à la 20th Century pour que le nom du studio ne soit pas associé à l'image ultra-conservatrice de la chaîne de télévision Fox News.

Le Monde avec AFP • Publié le 18 janvier 2020 à 06h46

 Lecture 1 min.

Révolution dans le monde d'Hollywood. Le célèbre studio de cinéma 20th Century Fox va voir une partie de son nom modifiée et devenir 20th Century Studios, décision prise par son nouveau propriétaire, Disney.


Selon plusieurs médias américains, Disney aurait décidé de modifier le nom du studio pour qu'il ne soit plus associé à la grande chaîne Fox mais aussi et surtout à Fox News, la chaîne d'information en continu.

CINEMA • MEDIA Sharing 

Disney decides to change the name of the legendary Fox film studio

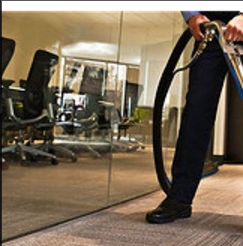
The new owner will remove the mention "Fox" in the 20th Century so that the name of the studio is not associated with the ultra-conservative image of the television channel Fox News.

Le Monde avec AFP • Posted on January 18, 2020 at 06:46

 Reading time 1 min.

Revolution in the Hollywood world. The famous film studio 20th Century Fox will see part of its name changed and become 20th Century Studios, decision taken by its new owner, Disney.

According to several American media, Disney has decided to change the name of the studio so that it is no longer associated with the big chain Fox but also and especially with Fox News, the news channel.



Search, Questions, and Reasoning

Web Images Groups News Froogle L

⚡ Default (GPT-3.5)

NI How many US states capitals' are also their largest cities?

As of my last update in September 2021, there were eight U.S. states where the capital city is also the largest city in the state. These states are:

1. Atlanta, Georgia
2. Boston, Massachusetts
3. Denver, Colorado
4. Honolulu, Hawaii
5. Indianapolis, Indiana
6. Oklahoma City, Oklahoma
7. Phoenix, Arizona
8. Salt Lake City, Utah

Please note that this information might have changed since then due to population shifts and other factors. For the most up-to-date information, it's best to consult a recent source or check official government websites.

NI How many US states capitals' are also their largest cities?

As of my knowledge cutoff in September 2021, there are 17 U.S. state capitals that are also the largest city in their respective states. These are:

1. Phoenix, Arizona
2. Little Rock, Arkansas
3. Denver, Colorado
4. Honolulu, Hawaii
5. Boise, Idaho
6. Indianapolis, Indiana
7. Des Moines, Iowa
8. Baton Rouge, Louisiana
9. Boston, Massachusetts
10. Jackson, Mississippi
11. Lincoln, Nebraska
12. Carson City, Nevada
13. Oklahoma City, Oklahoma
14. Providence, Rhode Island
15. Nashville, Tennessee
16. Salt Lake City, Utah
17. Charleston, West Virginia

Please note that this information might have changed after my last update in September 2021, so it's always a good idea to check the most recent data.

Regenerate

Send a message

Jeopardy!



Question Answering: Watson

"a camel is a horse designed by"

Wiktionary

Entry Discussion Read Edit History Search

a camel is a horse designed by a committee

Contents [hide]

- 1 English
 - 1.1 Alternative forms
 - 1.2 Proverb

The Phrase Finder

Discussion Forum

Google™ Custom Search Search

A camel is a horse designed by committee

Posted by Ruben P. Mendez on April 16, 2004

Does anyone know the origin of this maxim? I heard it way back at the United Nations, which is chockfull of committees. It may have originated there, but I'd like an authoritative explanation. Thanks

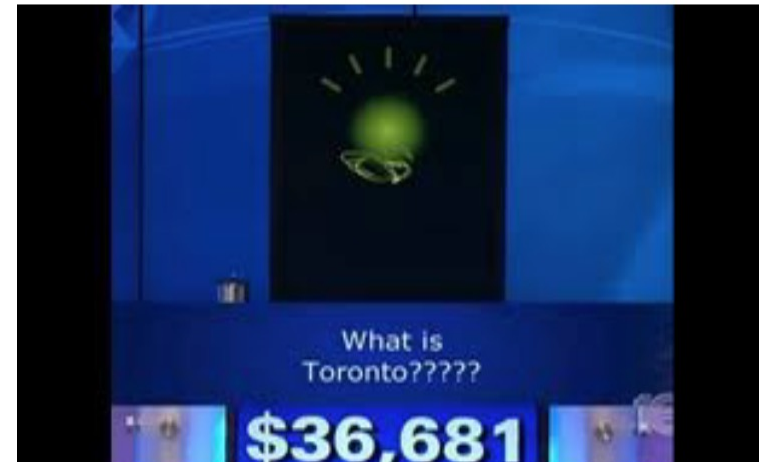
- [Re: A camel is a horse designed by committee](#) SR 16/April/04
 - [Re: A camel is a horse designed by committee](#) Henry 18/April/04

Question Answering: Watson

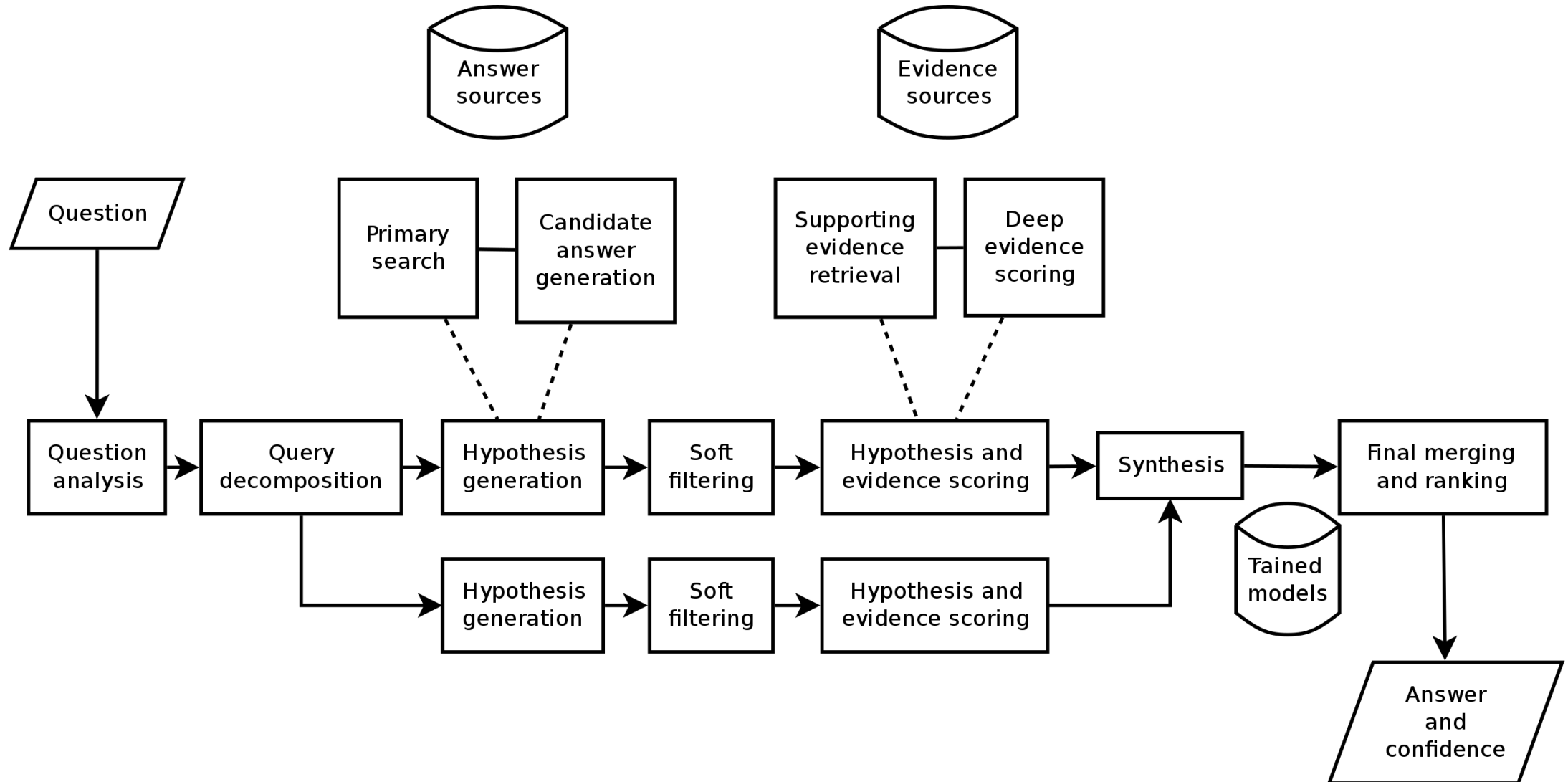


Jeopardy!
World Champion

US Cities: Its largest airport is named for a World War II hero; its second largest, for a World War II battle.



Watson



Language Comprehension?

Opera refers to a dramatic art form, originating in Europe, in which the emotional content is conveyed to the audience as much through music, both vocal and instrumental, as it is through the lyrics. By contrast, in musical theater an actor's dramatic performance is primary, and the music plays a lesser role. The drama in opera is presented using the primary elements of theater such as scenery, costumes, and acting. However, the words of the opera, or libretto, are sung rather than spoken. The singers are accompanied by a musical ensemble ranging from a small instrumental ensemble to a full symphonic orchestra.

1. It is pointed out in the reading that opera ----.

- A) has developed under the influence of musical theater
- B) is a drama sung with the accompaniment of an orchestra
- C) is not a high-budget production
- D) is often performed in Europe
- E) is the most complex of all the performing arts

Conversations with Devices?



ELIZA

ELIZA: You're not really talking about me -- are you ?
YOU: What?
ELIZA: What do you think ?
YOU: I mean, I feel happy to be writing right now, to be typing this as we speak. Lately, I've had writer's block.
ELIZA: You say you mean ?
YOU: That's what I said.
ELIZA: Can you elaborate on that ?
YOU: I only really feel happy when I'm writing.
ELIZA: How long have you been writing ?
YOU: Always, but I've been writing seriously since I was eighteen.
ELIZA: What incident are you thinking of ?
YOU: Um.
ELIZA: What does that suggest to you ?
YOU: I don't think it's enough, to only be happy when you write.
ELIZA: Does that trouble you ?
YOU: ...Um, because you can't write all the time. And then there's the real world, and you still have to deal with real people.
ELIZA: What are your feelings now ?
YOU: I've been thinking a lot of that poem by Yeats: "The Man and the Echo."
Do you know it?
ELIZA: Do you say you've been thinking a lot of that poem by yeats: "the man and the echo for some special reason ?
YOU: Yes

Large Language Models

+ New chat

May

- Detecting Human or LM
- "Converting JSONs to JSONL"
- "Converting Lists to CSV"
- Correcting "typeof" Error
- Generating Binary Random Lis
- "Overriding Method in Subclas
- Debugging Flask Cookie Issue
- Flask-SocketIO Cookie Debug
- Track Users with Cookies
- Legend: Square Next Words
- Complementary Event Subset

NI nicholas_tomlin@berke... ***

⚡ GPT-3.5 ⚙️ GPT-4

ChatGPT **PLUS**

Give me ideas
for what to do with my kids' art

Create a workout plan
for resistance training

Show me a code snippet
of a website's sticky header

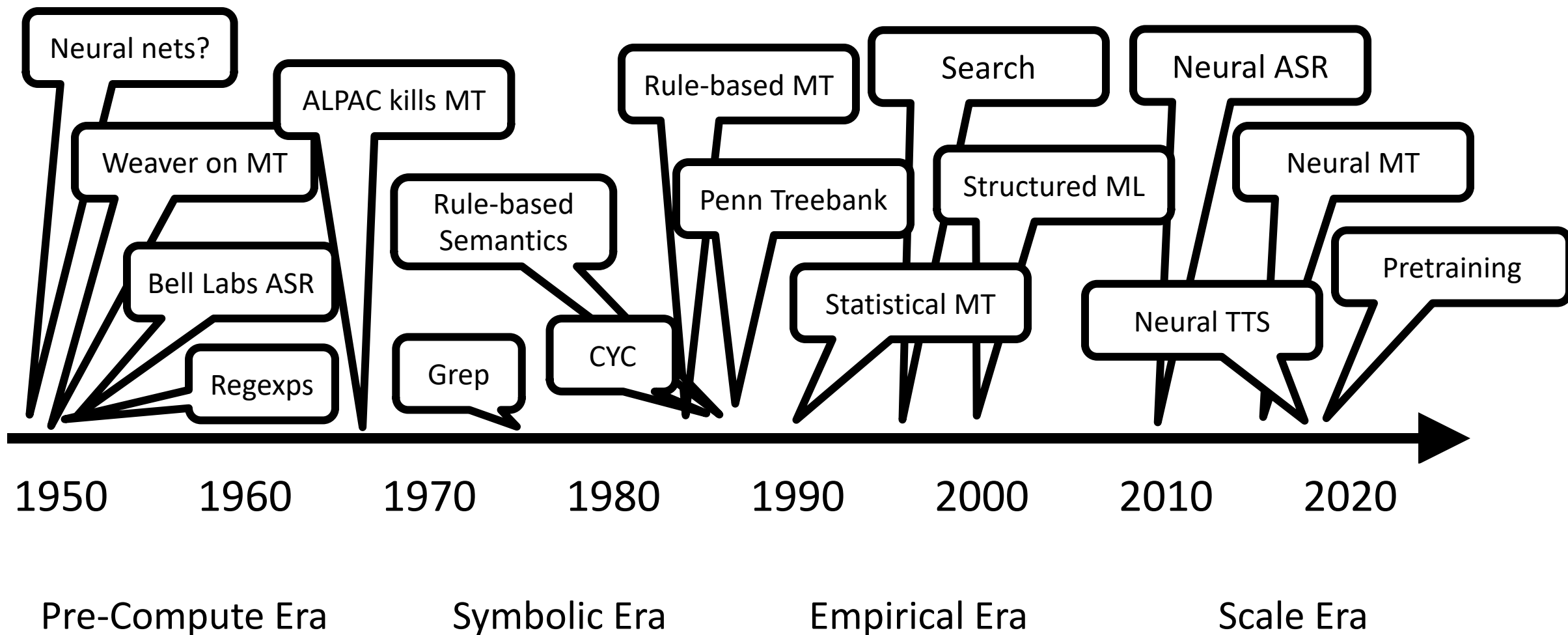
Recommend a dish
to bring to a potluck

Send a message ➤

ChatGPT may produce inaccurate information about people, places, or facts. [ChatGPT July 20 Version](#)



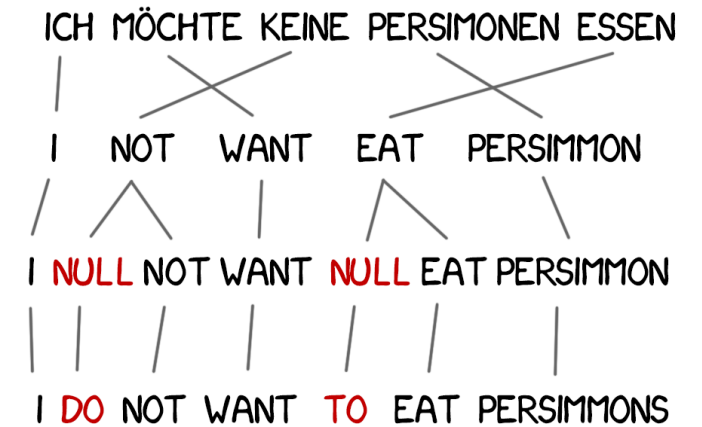
NLP History



Approach #1: Lexical Translation

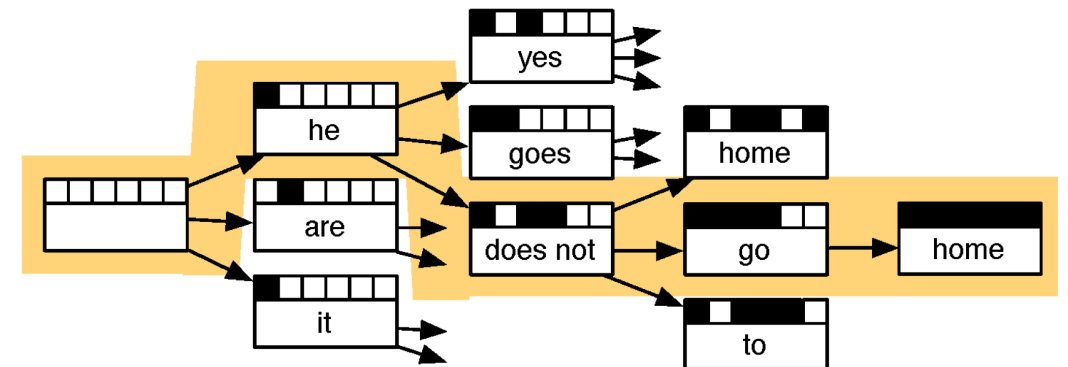
Step #1: Learn Alignments

- Learn mappings between words in source and target language
- IBM Model 1, 2, 3, 4, 5...
- Can also learn a phrase table of mappings



Step #2: Generate Language

- Search problem over the space of natural language strings
- Can use approaches like A* to guide search



Issue: Ambiguities



Stevie Wonder announces he'll be having kidney surgery during London concert

By Amir Vera, CNN

Updated 11:16 PM EDT, Sat July 06, 2019



(CNN) — [Stevie Wonder](#) will be taking a break from music.

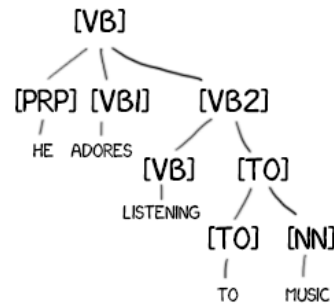
The legendary singer-songwriter announced during [a concert in London](#) Saturday that he will be undergoing kidney surgery.

Issue: Ambiguities

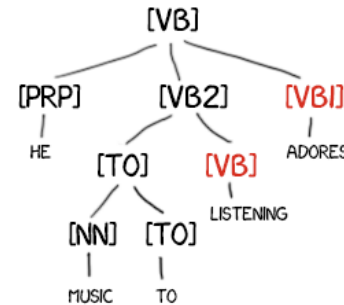
- Headlines:
 - Enraged Cow Injures Farmer with Ax
 - Teacher Strikes Idle Kids
 - Hospitals Are Sued by 7 Foot Doctors
 - Ban on Nude Dancing on Governor's Desk
 - Iraqi Head Seeks Arms
 - Stolen Painting Found by Tree
 - Kids Make Nutritious Snacks
 - Local HS Dropouts Cut in Half
- Can we come up with a representation to disambiguate the two readings of each headline?

Approach #2: Predict Intermediate Structures

HE ADORES LISTENING TO MUSIC

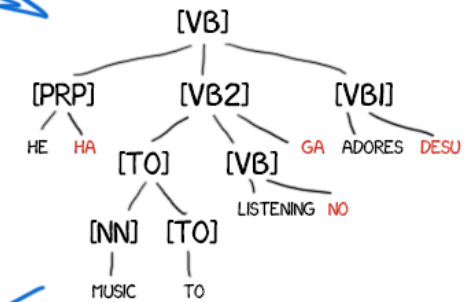


INPUT



SYNTAX TREE REORDERING

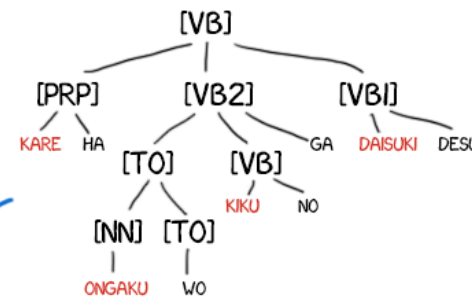
SYNTAX-BASED
STATISTICAL
MACHINE
TRANSLATION



WORD INSERTING

KARE HA ONGAKU WO KIKU NO GA DAISUKI DESU

RESULT



TRANSLATION

Approach #3: Language Modeling



the station signs are in deep in english	-14732
the stations signs are in deep in english	-14735
the station signs are in deep into english	-14739
the station 's signs are in deep in english	-14740
the station signs are in deep in the english	-14741
the station signs are indeed in english	-14757
the station 's signs are indeed in english	-14760
the station signs are indians in english	-14790

Noisy Channel Model: ASR

- We want to predict a sentence given acoustics:

$$w^* = \arg \max_w P(w|a)$$

- The noisy-channel approach:

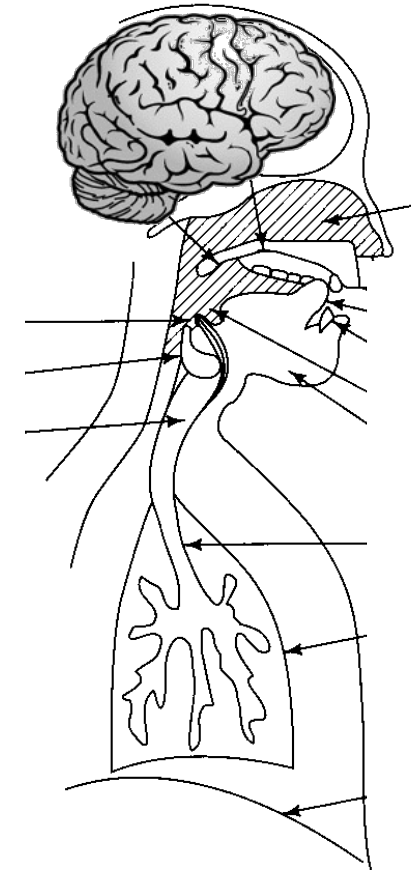
$$w^* = \arg \max_w P(w|a)$$

$$= \arg \max_w P(a|w)P(w)/P(a)$$

$$\propto \arg \max_w P(a|w)P(w)$$

Acoustic model: score fit between sounds and words

Language model: score plausibility of word sequences



Noisy Channel Model: Translation

“Having guessed and inferred considerably about, the powerful new mechanized methods in cryptography...one naturally wonders if the problem of translation could conceivably be treated as a problem in cryptography. When I look at an article in Russian, I say: ‘This is really written in English, but it has been coded in some strange symbols. I will now proceed to decode.’ ”

Warren Weaver (1947)

Machine Translation

CINÉMA • MÉDIAS Partage 

Disney décide de changer le nom du légendaire studio de cinéma Fox


Le nouveau propriétaire va faire disparaître la mention « Fox » à la 20th Century pour que le nom du studio ne soit pas associé à l'image ultra-conservatrice de la chaîne de télévision Fox News.

Le Monde avec AFP • Publié le 18 janvier 2020 à 06h46

 Lecture 1 min.

Révolution dans le monde d'Hollywood. Le célèbre studio de cinéma 20th Century Fox va voir une partie de son nom modifiée et devenir 20th Century Studios, décision prise par son nouveau propriétaire, Disney.


Selon plusieurs médias américains, Disney aurait décidé de modifier le nom du studio pour qu'il ne soit plus associé à la grande chaîne Fox mais aussi et surtout à Fox News, la chaîne d'information en continu.

CINEMA • MEDIA Sharing 

Disney decides to change the name of the legendary Fox film studio

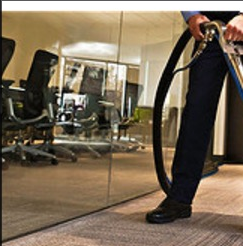
The new owner will remove the mention "Fox" in the 20th Century so that the name of the studio is not associated with the ultra-conservative image of the television channel Fox News.

Le Monde avec AFP • Posted on January 18, 2020 at 06:46

 Reading time 1 min.

Revolution in the Hollywood world. The famous film studio 20th Century Fox will see part of its name changed and become 20th Century Studios, decision taken by its new owner, Disney.

According to several American media, Disney has decided to change the name of the studio so that it is no longer associated with the big chain Fox but also and especially with Fox News, the news channel.



Empirical N-Grams

- Use statistics from data (examples here from Google N-Grams)

Training Counts	198015222 the first
	194623024 the same
	168504105 the following
	158562063 the world
	...
	14112454 the door

23135851162 the *	

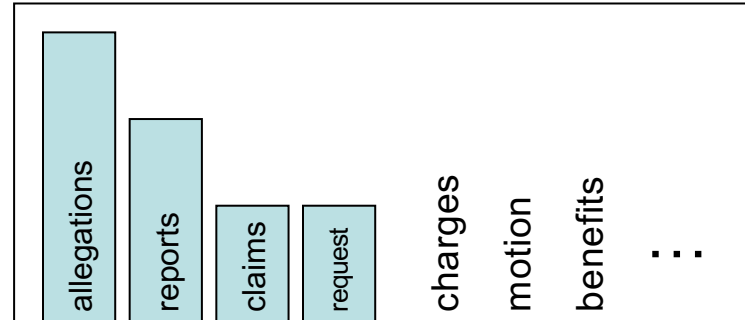
$$\hat{P}(\text{door}|\text{the}) = \frac{14112454}{23135851162} = 0.0006$$

- This is the maximum likelihood estimate, which needs modification
- N-gram models use such counts to compute probabilities on demand

Smoothing

- We often want to make estimates from sparse statistics:

$P(w \mid \text{denied the})$
3 allegations
2 reports
1 claims
1 request
7 total



- Smoothing flattens spiky distributions so they generalize better:

$P(w \mid \text{denied the})$
2.5 allegations
1.5 reports
0.5 claims
0.5 request
2 other
7 total



- Very important all over NLP, but easy to do badly

Back-off

Please close the first door on the left.

4-Gram

```
3380 please close the door
1601 please close the window
1164 please close the new
1159 please close the gate
...
0 please close the first
-----
13951 please close the *
```

0.0

3-Gram

```
197302 close the window
191125 close the door
152500 close the gap
116451 close the thread
...
8662 close the first
-----
3785230 close the *
```

0.002

2-Gram

```
198015222 the first
194623024 the same
168504105 the following
158562063 the world
...
...
-----
23135851162 the *
```

0.009

Specific but Sparse



Dense but General

$$\lambda \hat{P}(w|w_{-1}, w_{-2}) + \lambda' \hat{P}(w|w_{-1}) + \lambda'' \hat{P}(w)$$

Discounting

- Observation: N-grams occur more in training data than they will later

Empirical Bigram Counts (Church and Gale, 91)

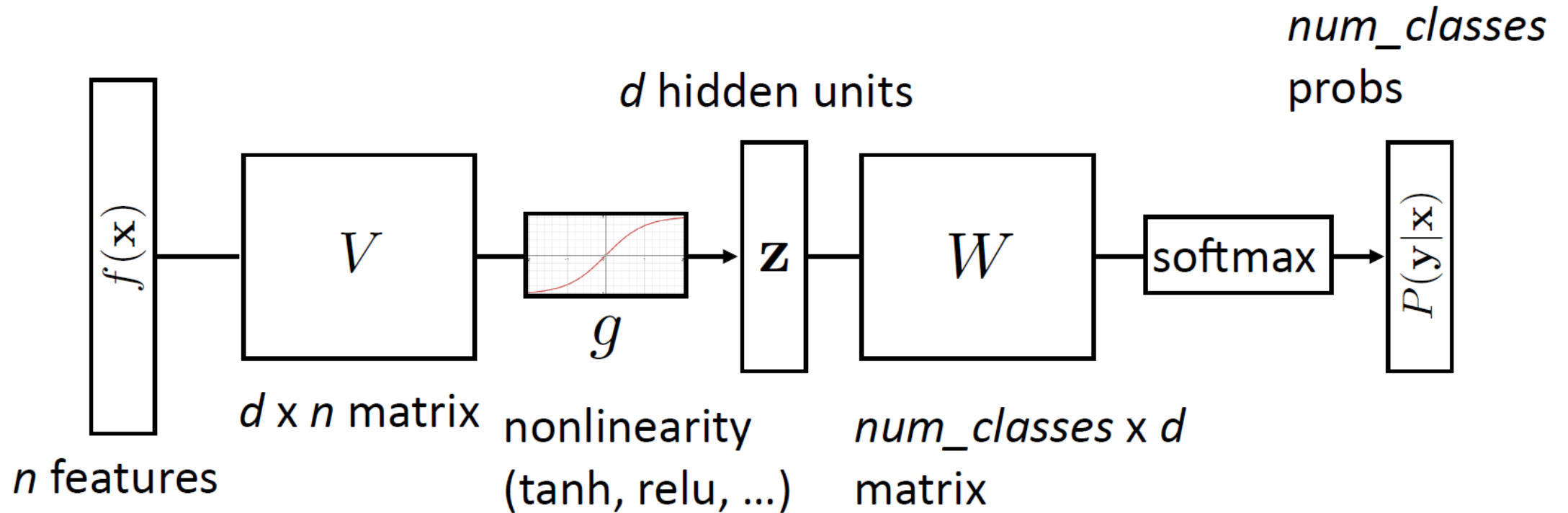
Count in 22M Words	Future c^* (Next 22M)
1	
2	
3	
4	
5	

- Absolute discounting: reduce counts by a small constant, redistribute “shaved” mass to a model of new events

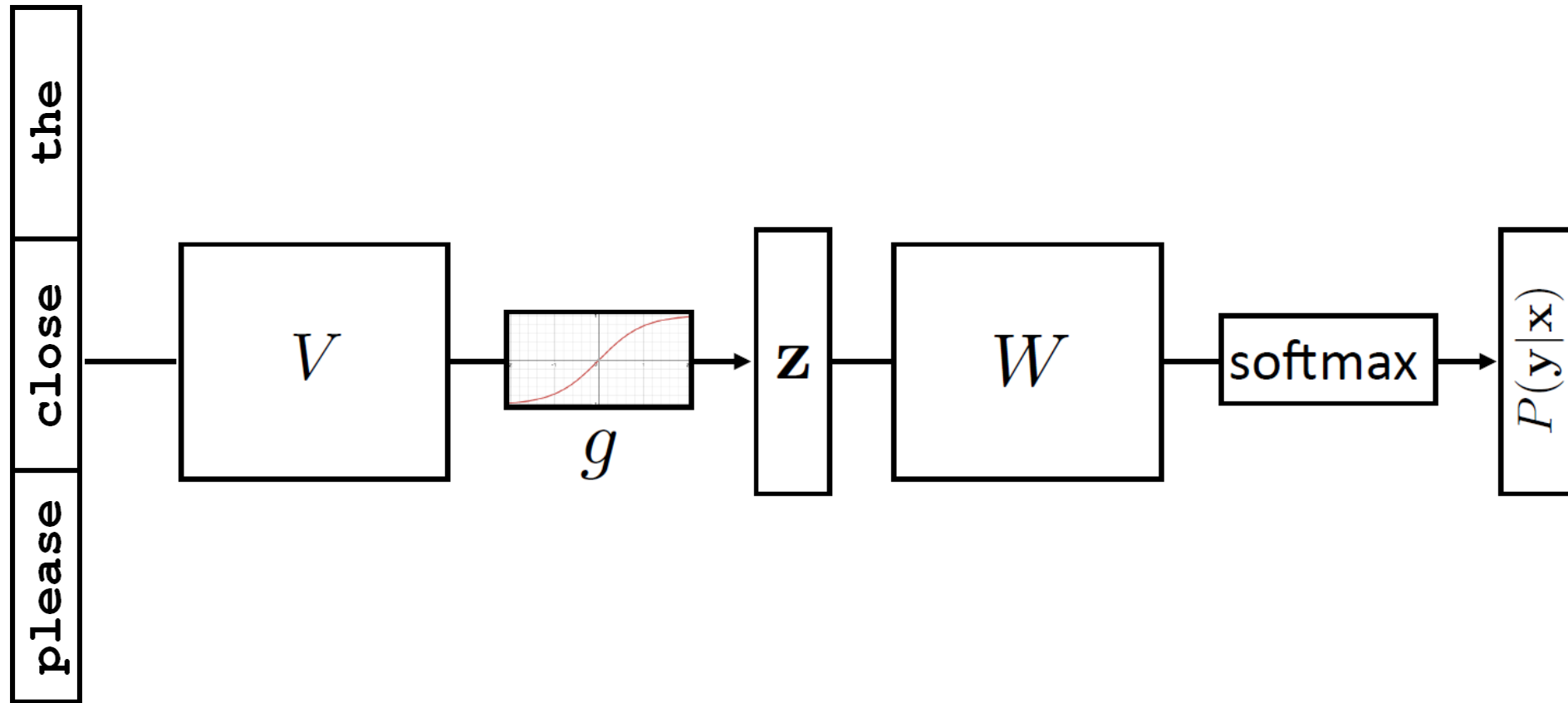
$$P_{\text{ad}}(w|w') = \frac{c(w', w) - d}{c(w')} + \alpha(w')\hat{P}(w)$$

Reminder: Feedforward Neural Nets

$$P(\mathbf{y}|\mathbf{x}) = \text{softmax}(W g(V f(\mathbf{x})))$$

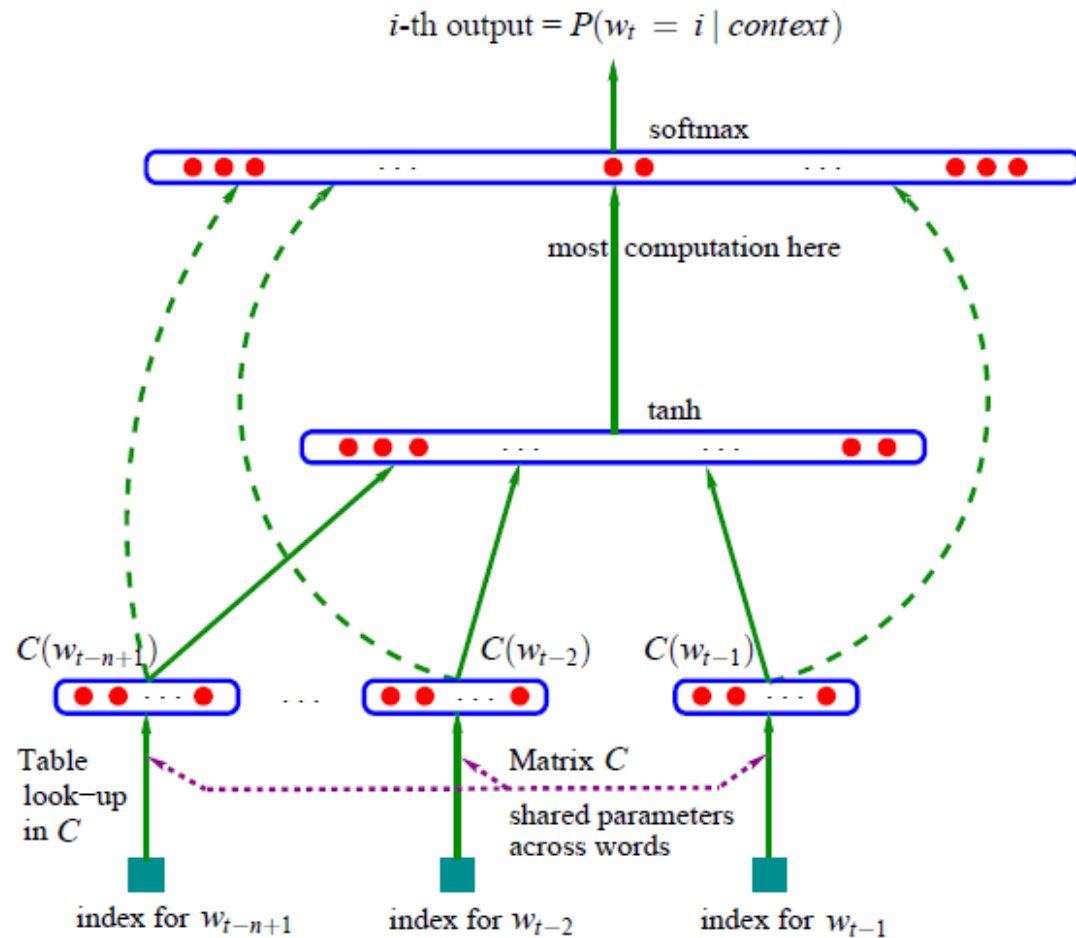


A Feedforward N-Gram Model?



Early Neural Language Models

- Fixed-order feed-forward neural LMs
 - Eg Bengio et al 03
 - Allow generalization across contexts in more nuanced ways than prefixing
 - Allow different kinds of pooling in different contexts
 - Much more expensive to train



Recurrent NNs

Cloze Task (The Shannon Game)

Cloze Task (The Shannon Game)

Today

Cloze Task (The Shannon Game)

Today, I

Cloze Task (The Shannon Game)

Today, I went

Cloze Task (The Shannon Game)

Today, I went to

Cloze Task (The Shannon Game)

Today, I went to the

Cloze Task (The Shannon Game)

Today, I went to the store

Cloze Task (The Shannon Game)

Today, I went to the store and

Cloze Task (The Shannon Game)

Today, I went to the store and bought

Cloze Task (The Shannon Game)

Today, I went to the store and bought some

Cloze Task (The Shannon Game)

Today, I went to the store and bought some milk

Cloze Task (The Shannon Game)

Today, I went to the store and bought some milk and

Cloze Task (The Shannon Game)

Today, I went to the store and bought some milk and eggs.

Cloze Task (The Shannon Game)

Today, I went to the store and bought some milk and eggs. I

Cloze Task (The Shannon Game)

Today, I went to the store and bought some milk and eggs. I knew

Cloze Task (The Shannon Game)

Today, I went to the store and bought some milk and eggs. I knew it

Cloze Task (The Shannon Game)

Today, I went to the store and bought some milk and eggs. I knew it was

Cloze Task (The Shannon Game)

Today, I went to the store and bought some milk and eggs. I knew it was going

Cloze Task (The Shannon Game)

Today, I went to the store and bought some milk and eggs. I knew it was going to

Cloze Task (The Shannon Game)

Today, I went to the store and bought some milk and eggs. I knew it was going to rain,

Cloze Task (The Shannon Game)

Today, I went to the store and bought some milk and eggs. I knew it was going to rain, but

Cloze Task (The Shannon Game)

Today, I went to the store and bought some milk and eggs. I knew it was going to rain, but I

Cloze Task (The Shannon Game)

Today, I went to the store and bought some milk and eggs. I knew it was going to rain, but I forgot

Cloze Task (The Shannon Game)

Today, I went to the store and bought some milk and eggs. I knew it was going to rain, but I forgot to

Cloze Task (The Shannon Game)

Today, I went to the store and bought some milk and eggs. I knew it was going to rain, but I forgot to take

Cloze Task (The Shannon Game)

Today, I went to the store and bought some milk and eggs. I knew it was going to rain, but I forgot to take my

Cloze Task (The Shannon Game)

Today, I went to the store and bought some milk and eggs. I knew it was going to rain, but I forgot to take my umbrella,

Cloze Task (The Shannon Game)

Today, I went to the store and bought some milk and eggs. I knew it was going to rain, but I forgot to take my umbrella, and

Cloze Task (The Shannon Game)

Today, I went to the store and bought some milk and eggs. I knew it was going to rain, but I forgot to take my umbrella, and ended

Cloze Task (The Shannon Game)

Today, I went to the store and bought some milk and eggs. I knew it was going to rain, but I forgot to take my umbrella, and ended up

Cloze Task (The Shannon Game)

Today, I went to the store and bought some milk and eggs. I knew it was going to rain, but I forgot to take my umbrella, and ended up getting

Cloze Task (The Shannon Game)

Today, I went to the store and bought some milk and eggs. I knew it was going to rain, but I forgot to take my umbrella, and ended up getting wet

Cloze Task (The Shannon Game)

Today, I went to the store and bought some milk and eggs. I knew it was going to rain, but I forgot to take my umbrella, and ended up getting wet on

Cloze Task (The Shannon Game)

Today, I went to the store and bought some milk and eggs. I knew it was going to rain, but I forgot to take my umbrella, and ended up getting wet on the

Cloze Task (The Shannon Game)

Today, I went to the store and bought some milk and eggs. I knew it was going to rain, but I forgot to take my umbrella, and ended up getting wet on the way.

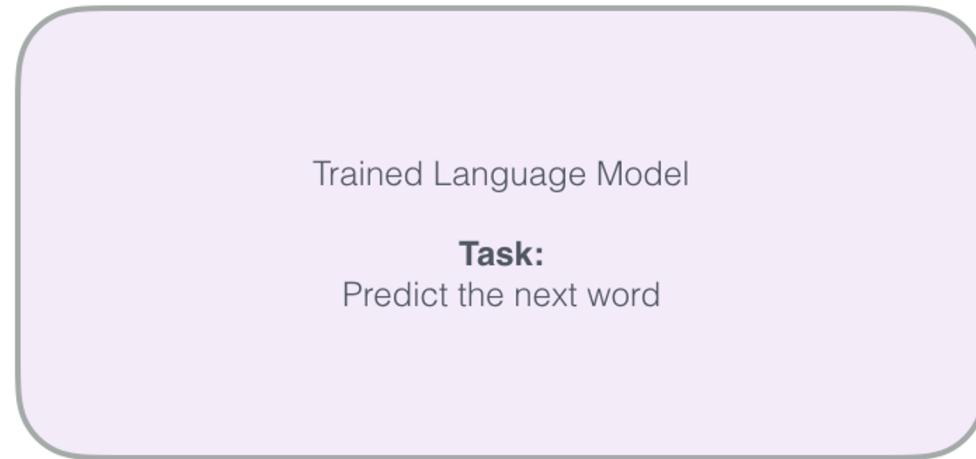
Language Modeling

Input
Features

Thou



shalt



Output
Prediction

0%	aardvark
0%	aarhus
0.1%	aaron
...	...
40%	not
...	...
0.01	zyzzyva

Recall: Language Modeling

- Goal: learn a probability distribution over possible next words

$$P(w_k \mid w_{k-1}, \dots, w_0)$$

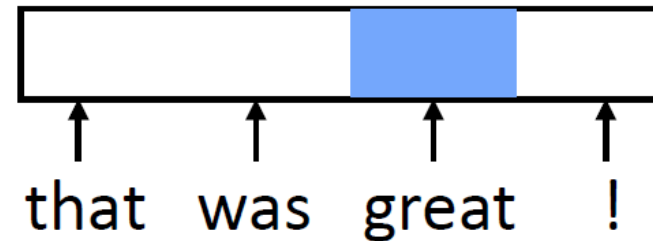
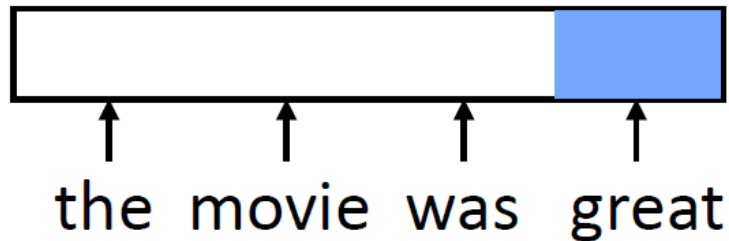
- Markovian assumption (used in n-gram models):

$$P(w_k \mid w_{k-1}, \dots, w_0) = P(w_k \mid w_{k-1}, \dots, w_{k-n+1})$$

- E.g., in a bigram model: $P(w_k \mid w_{k-1}, \dots, w_0) = P(w_k \mid w_{k-1}, w_{k-2})$

RNNs

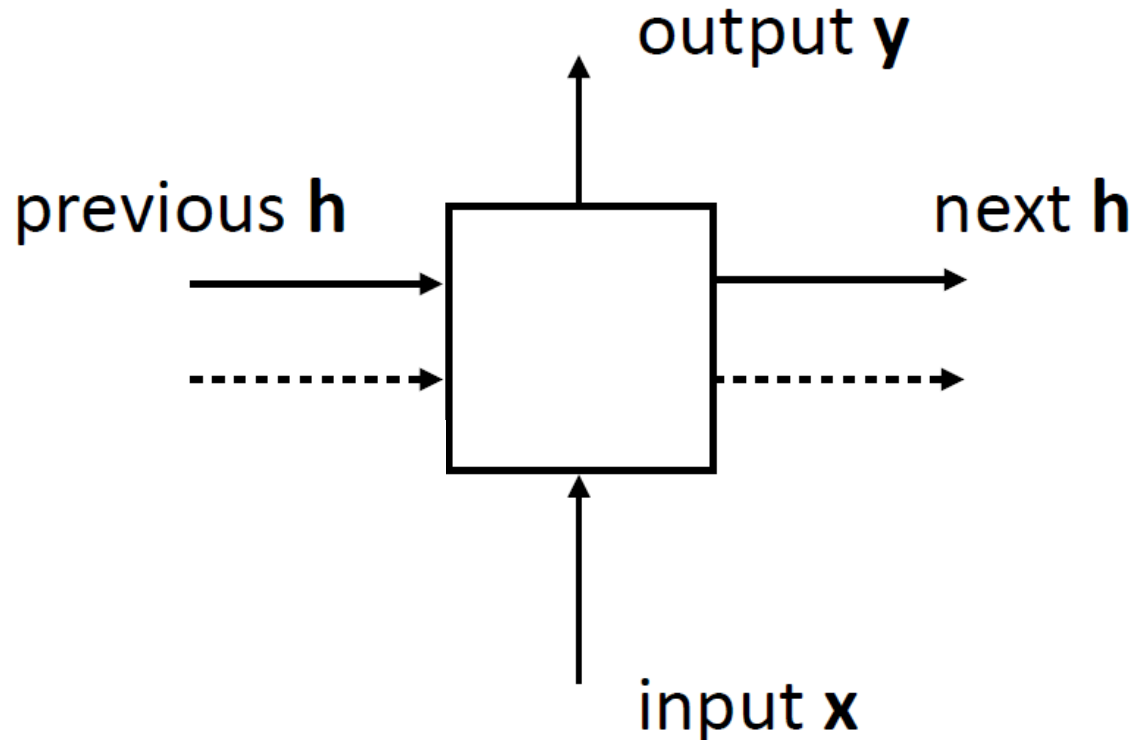
- ▶ Feedforward NNs can't handle variable length input: each position in the feature vector has fixed semantics



- ▶ These don't look related (*great* is in two different orthogonal subspaces)
- ▶ Instead, we need to:
 - 1) Process each word in a uniform way
 - 2) ...while still exploiting the context that that token occurs in

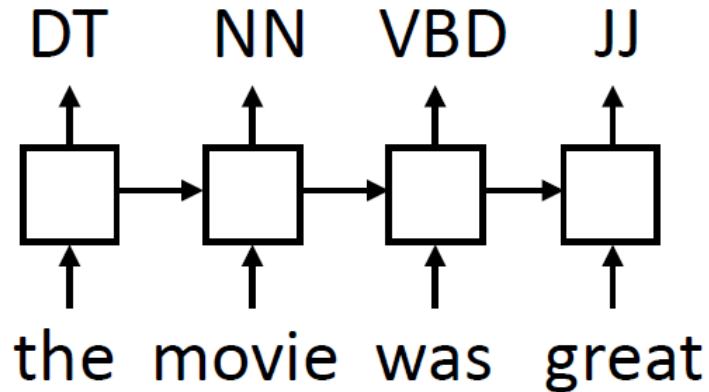
General RNN Approach

- ▶ Cell that takes some input \mathbf{x} , has some hidden state \mathbf{h} , and updates that hidden state and produces output \mathbf{y} (all vector-valued)



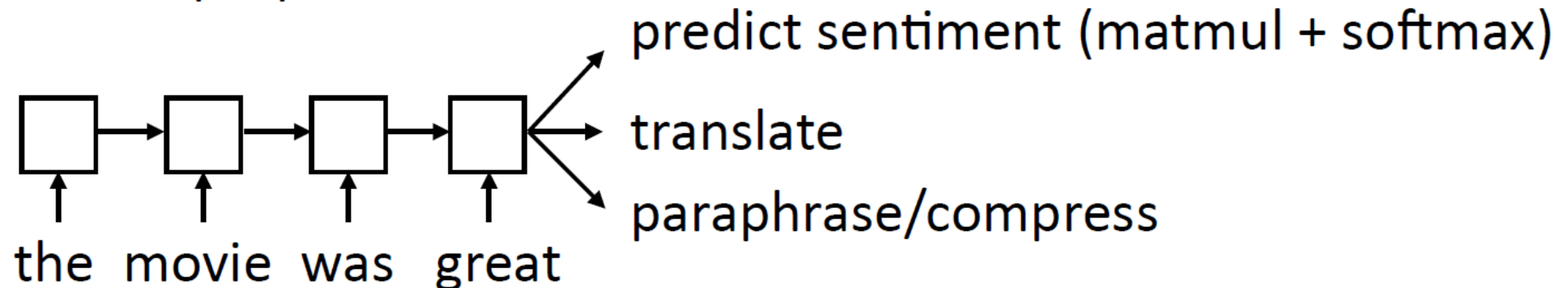
RNN Uses

- ▶ Transducer: make some prediction for each element in a sequence

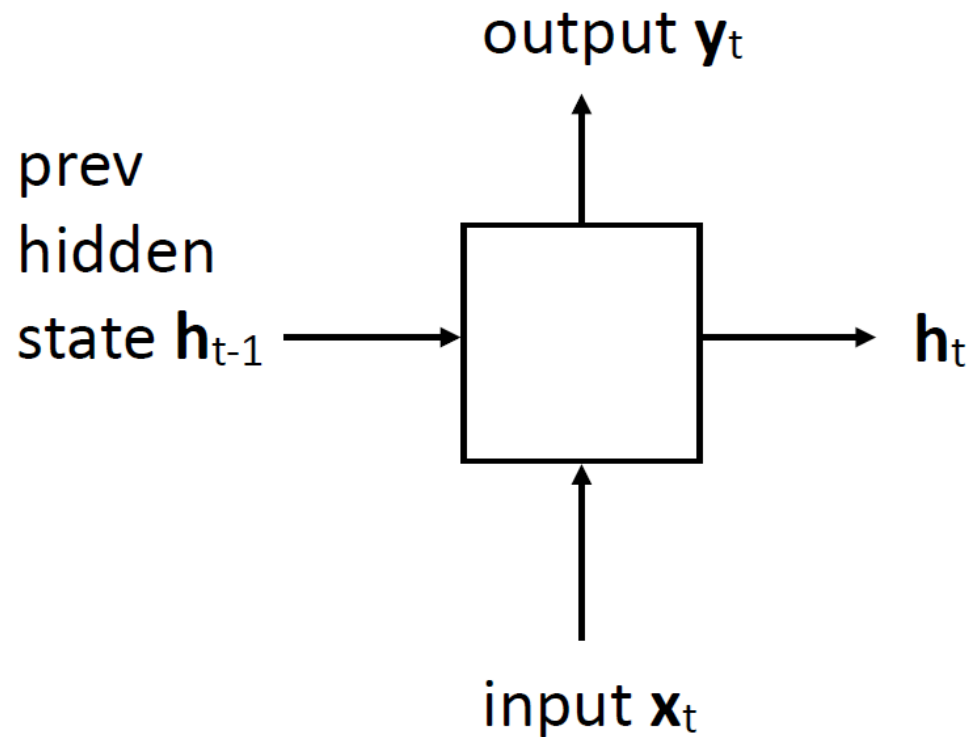


output \mathbf{y} = score for each tag, then softmax

- ▶ Acceptor/encoder: encode a sequence into a fixed-sized vector and use that for some purpose



Basic RNNs



$$\mathbf{h}_t = \tanh(W\mathbf{x}_t + V\mathbf{h}_{t-1} + \mathbf{b}_h)$$

- ▶ Updates hidden state based on input and current hidden state

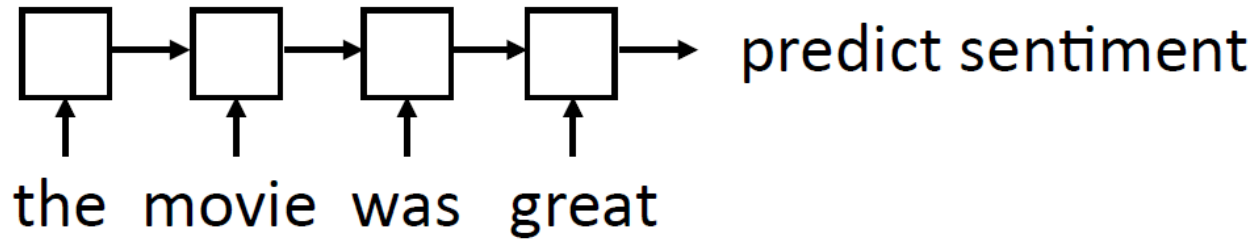
$$\mathbf{y}_t = \tanh(U\mathbf{h}_t + \mathbf{b}_y)$$

- ▶ Computes output from hidden state

- ▶ Long history! (invented in the late 1980s)

Elman (1990)

Training RNNs

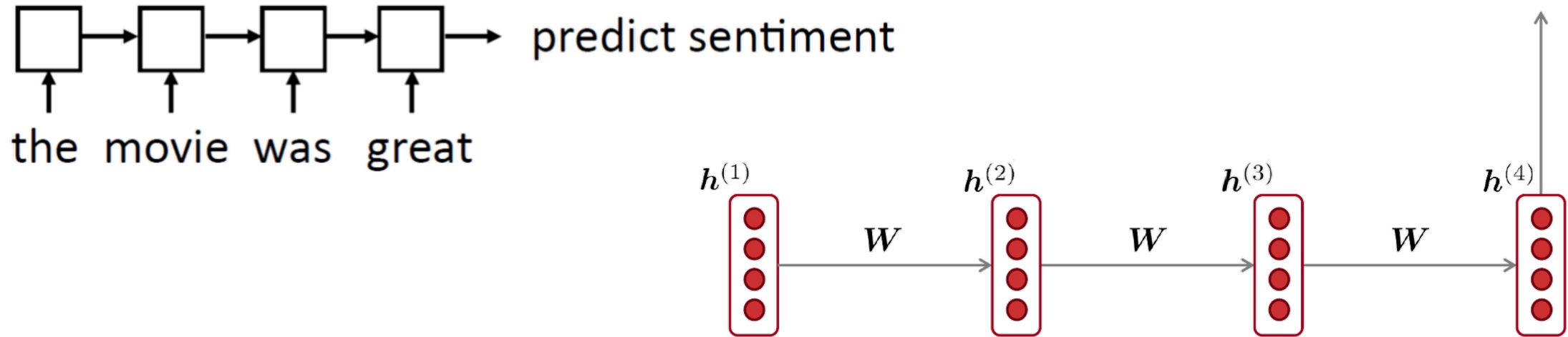


- ▶ “Backpropagation through time”: build the network as one big computation graph, some parameters are shared
- ▶ RNN potentially needs to learn how to “remember” information for a long time!

it was my **favorite** movie of 2016, though it wasn't without **problems** -> +

- ▶ “Correct” parameter update is to do a better job of remembering the sentiment of *favorite*

Problem: Vanishing Gradients



- Contribution of earlier inputs decreases if matrices are contractive (first eigenvalue < 1), non-linearities are squashing, etc
- Gradients can be viewed as a measure of the effect of the past on the future
- That's a problem for optimization but also means that information naturally decays quickly, so model will tend to capture local information

Core Issue: Information Decay

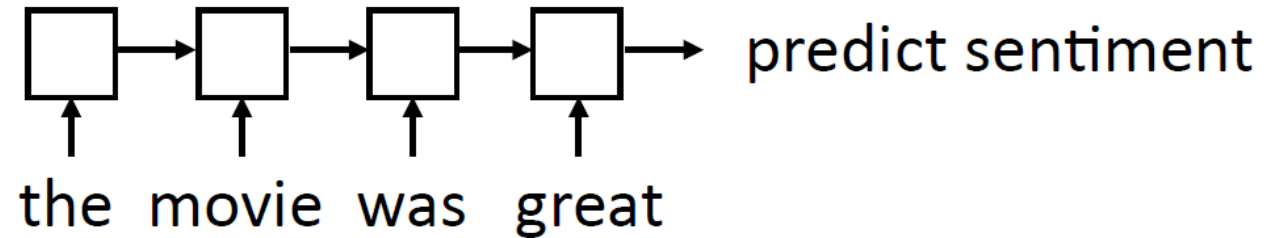
- The main problem is that *it's too difficult for the RNN to learn to preserve information over many timesteps.*

- In a vanilla RNN, the hidden state is constantly being rewritten

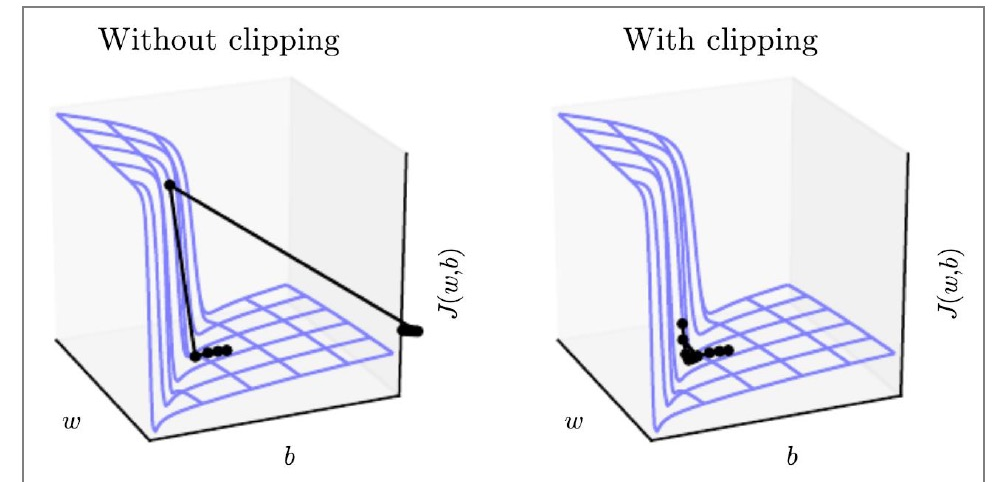
$$\mathbf{h}^{(t)} = \sigma \left(\mathbf{W}_h \mathbf{h}^{(t-1)} + \mathbf{W}_x \mathbf{x}^{(t)} + \mathbf{b} \right)$$

- How about a RNN with separate memory?

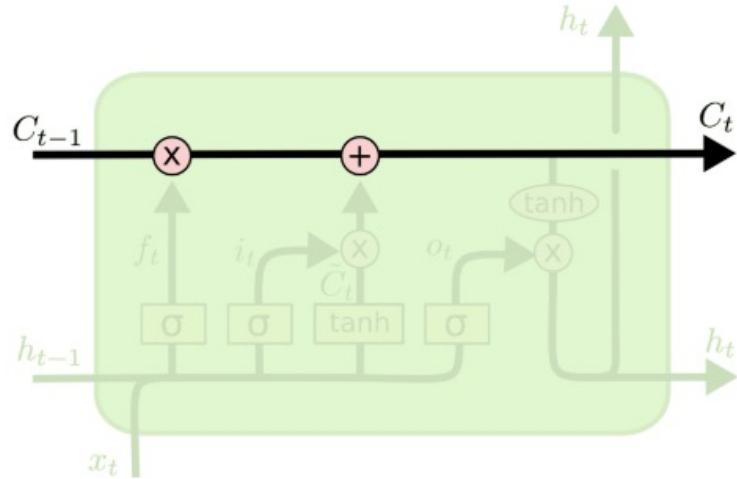
Problem: Exploding Gradients



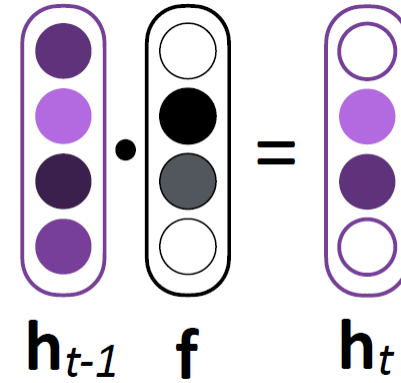
- Gradients can also be too large
 - Leads to overshooting / jumping around the parameter space
 - Common solution: gradient clipping



Key Idea: Propagated State



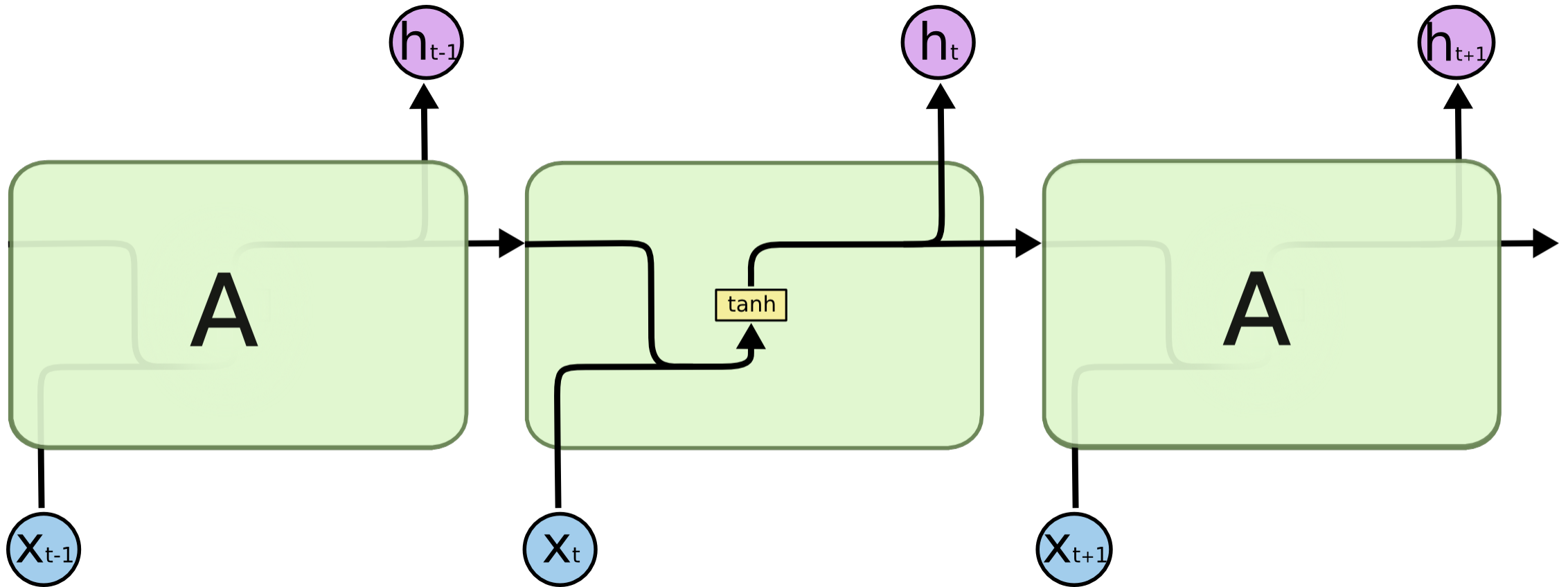
Cell State



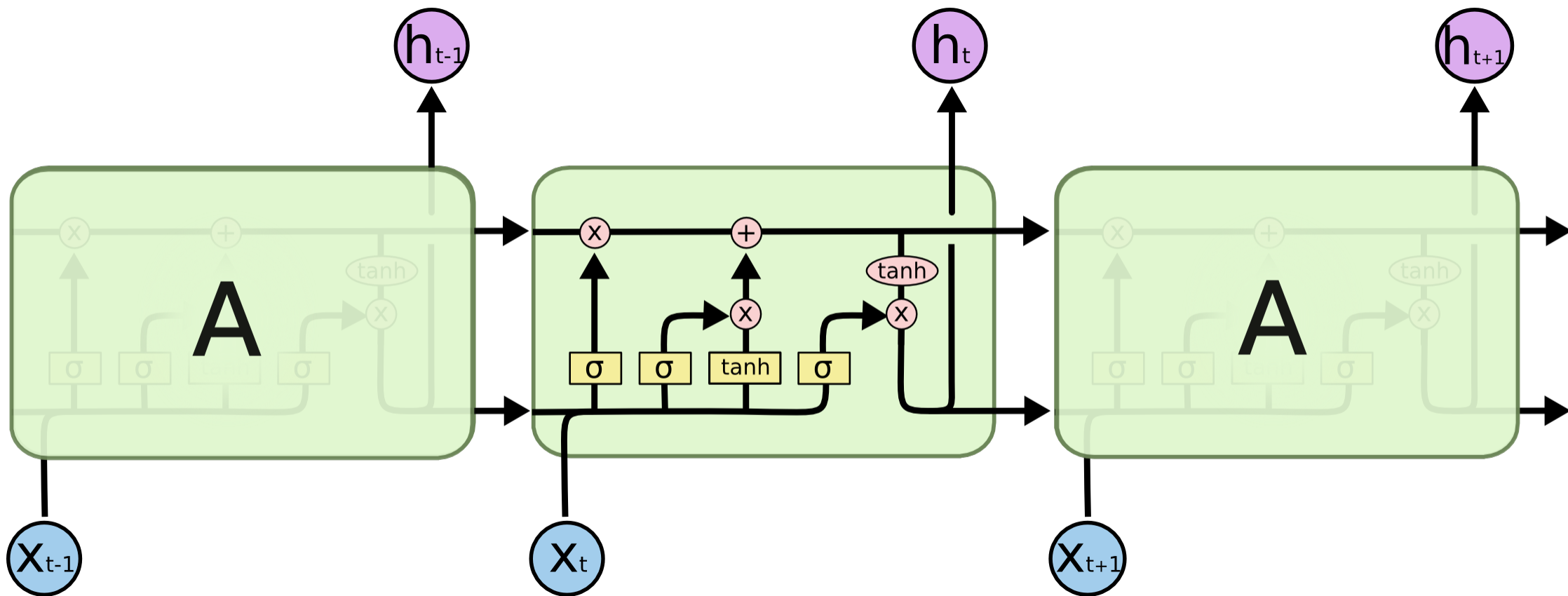
Gating

- Information decays in RNNs because it gets multiplied each time step
- Idea: have a channel called the *cell state* that by default just gets propagated (the “conveyer belt”)
- Gates make explicit decisions about what to add / forget from this channel

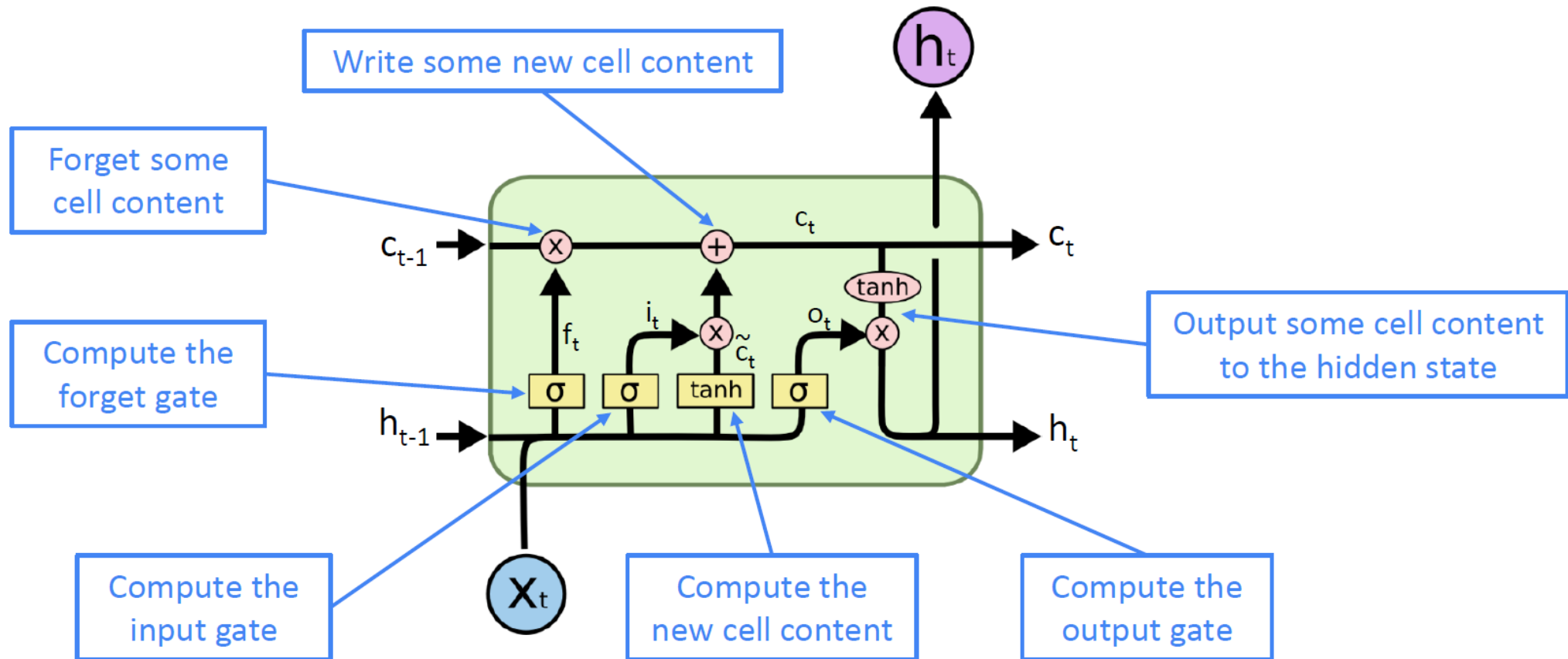
RNNs



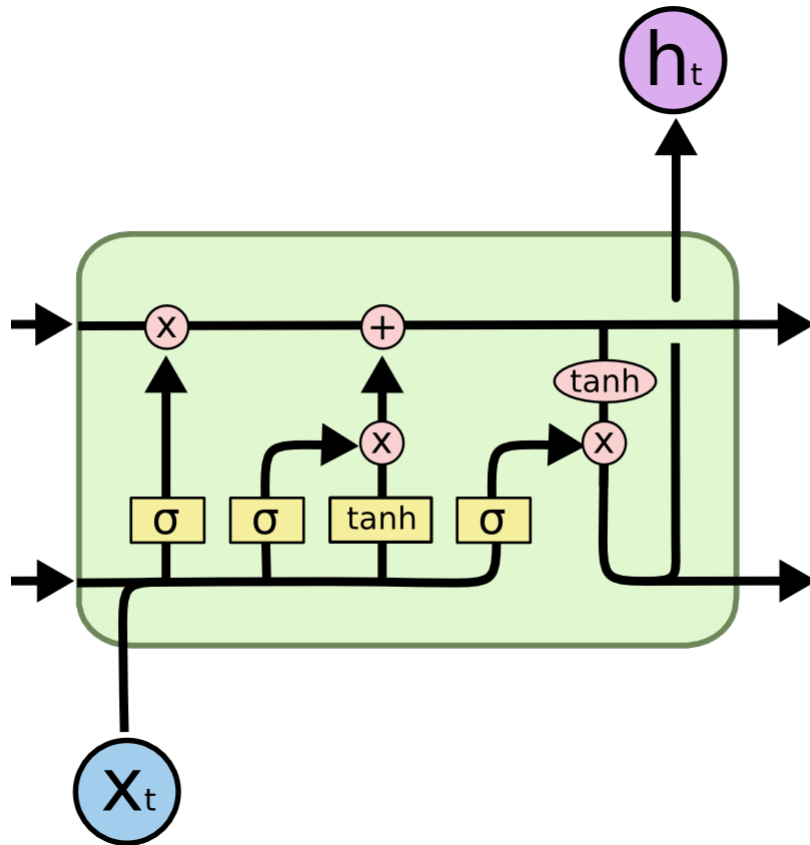
LSTMs



LSTMs

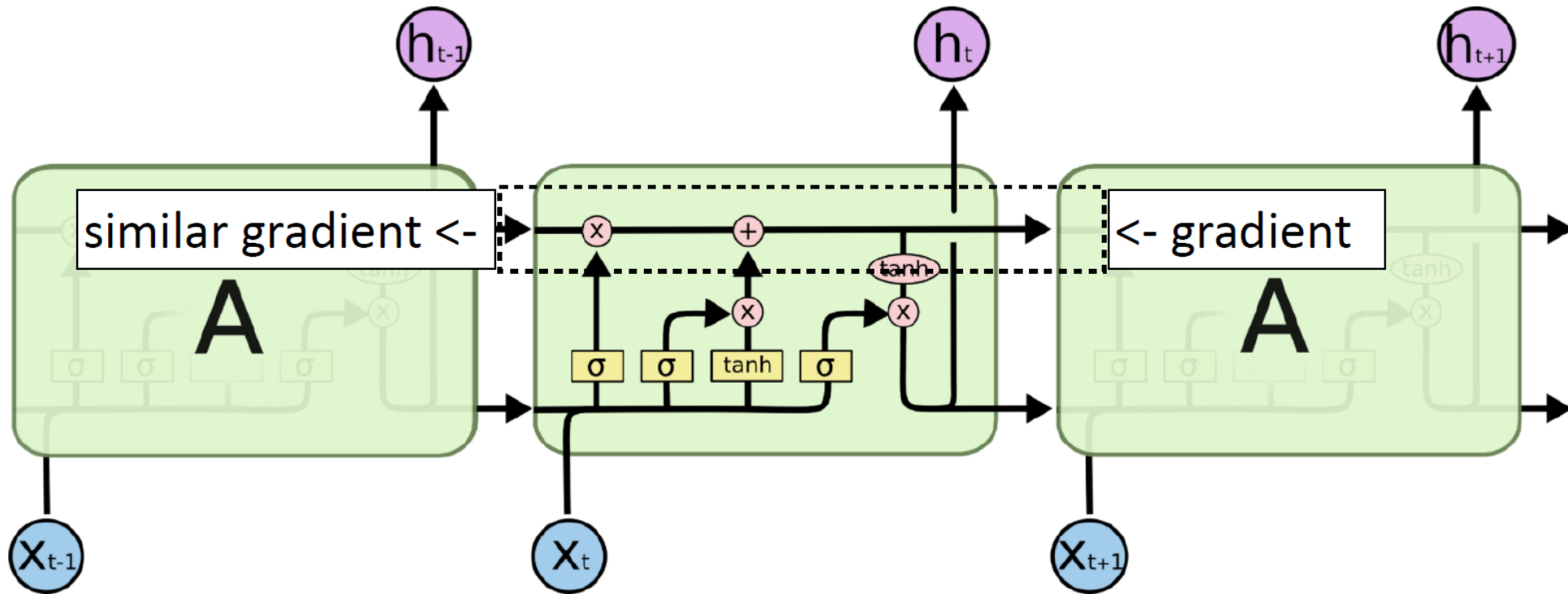


LSTMs



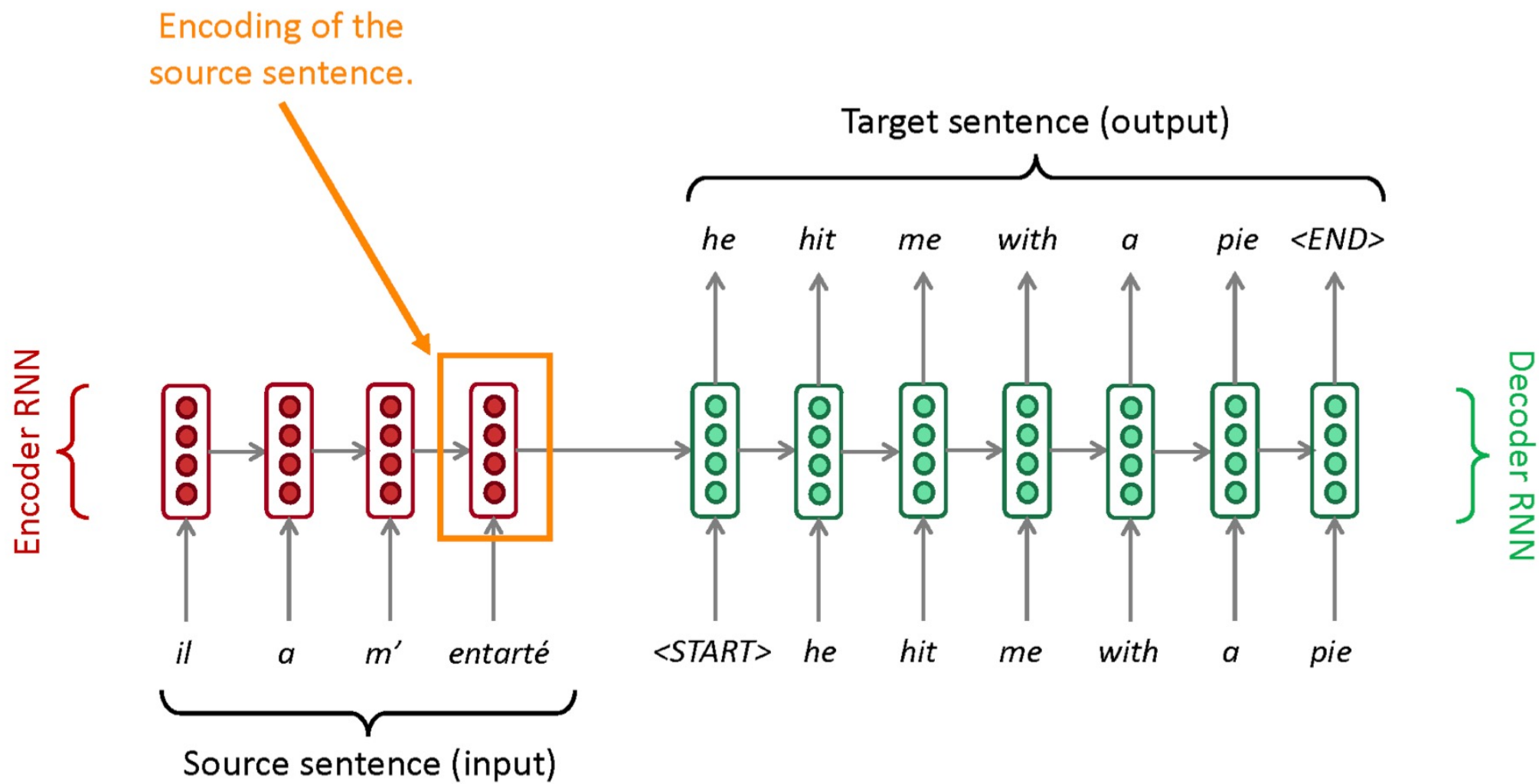
- ▶ Ignoring recurrent state entirely:
 - ▶ Lets us get feedforward layer over token
- ▶ Ignoring input:
 - ▶ Lets us discard stopwords
- ▶ Summing inputs:
 - ▶ Lets us compute a bag-of-words representation

What about the Gradients?

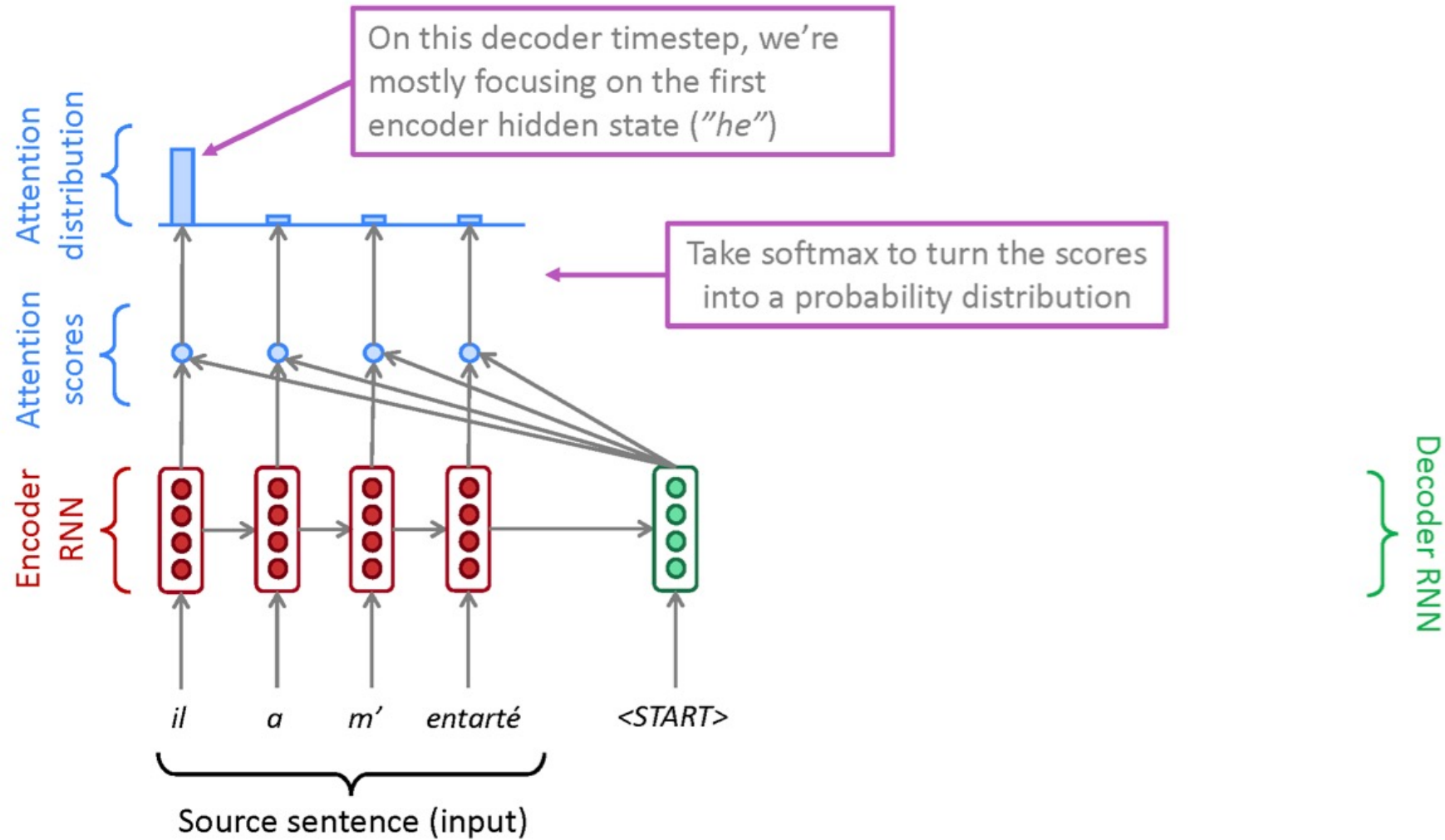


- ▶ Gradient still diminishes, but in a controlled way and generally by less — usually initialize forget gate = 1 to remember everything to start

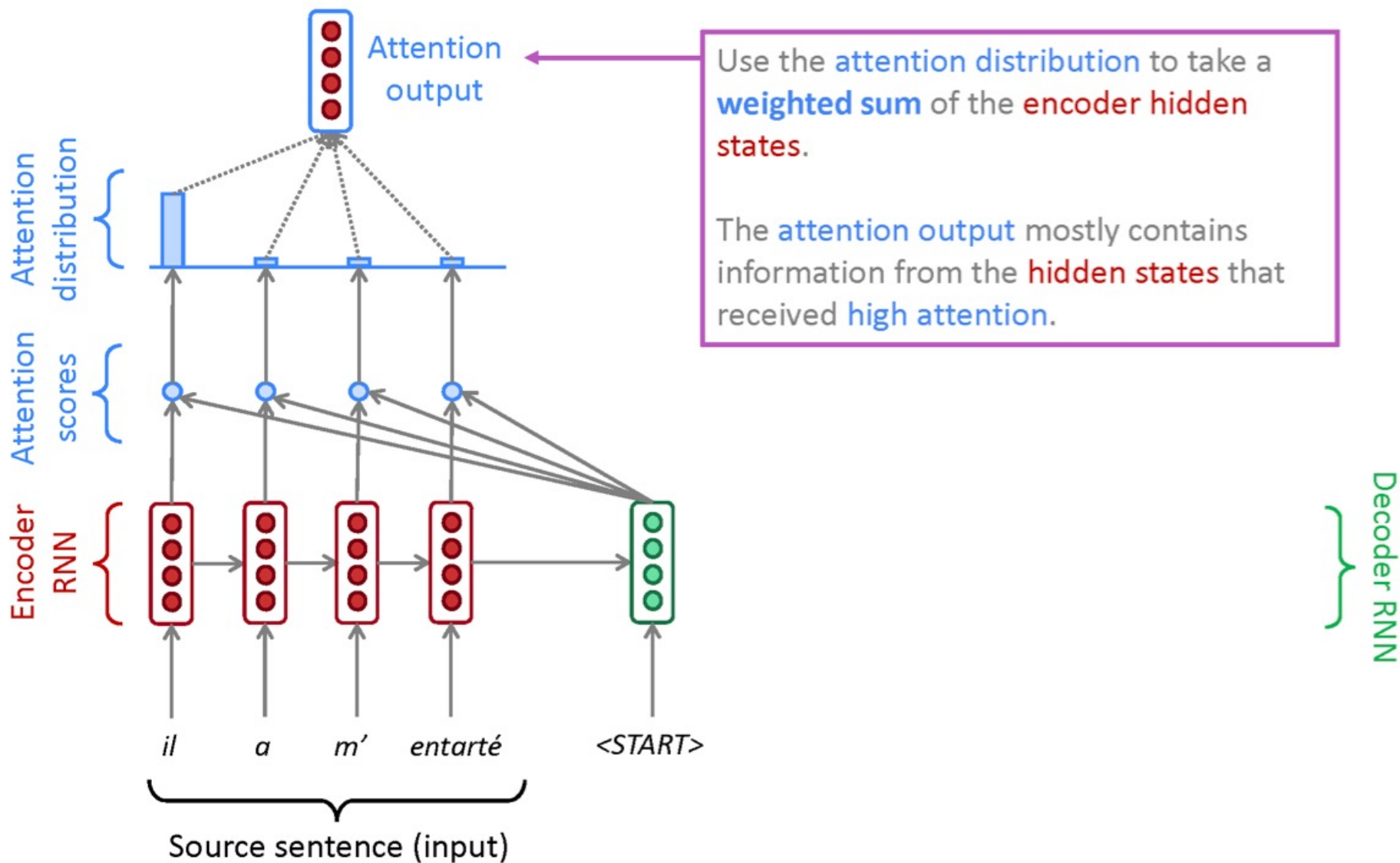
The Bottleneck Problem



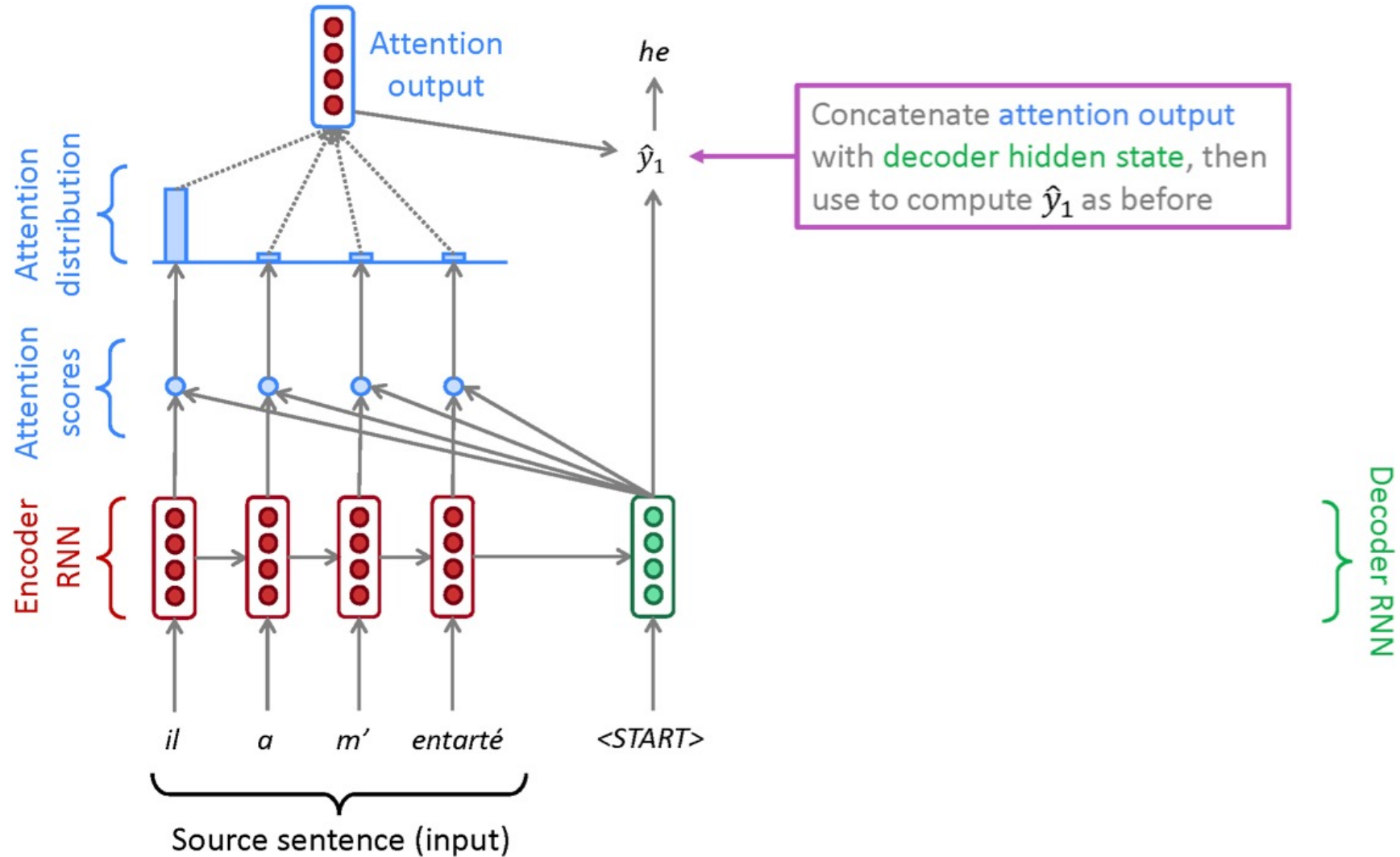
LSTMs with Attention



LSTMs with Attention



LSTMs with Attention



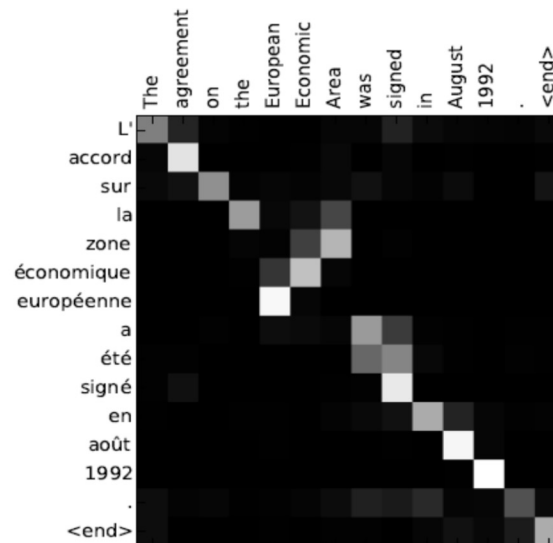
Attention



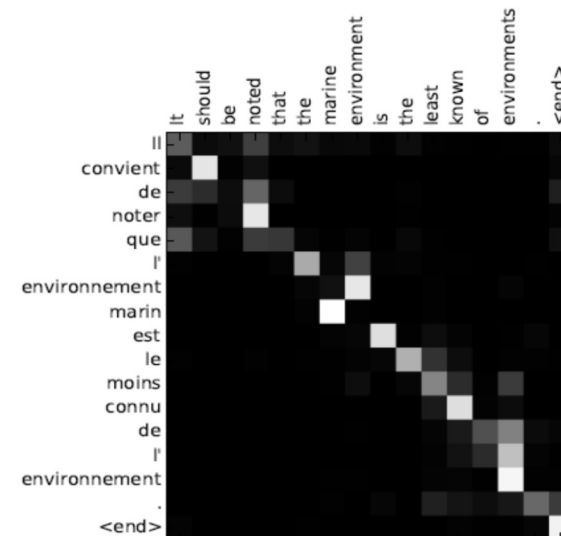
A group of people sitting on a boat in the water.



A giraffe standing in a forest with trees in the background.

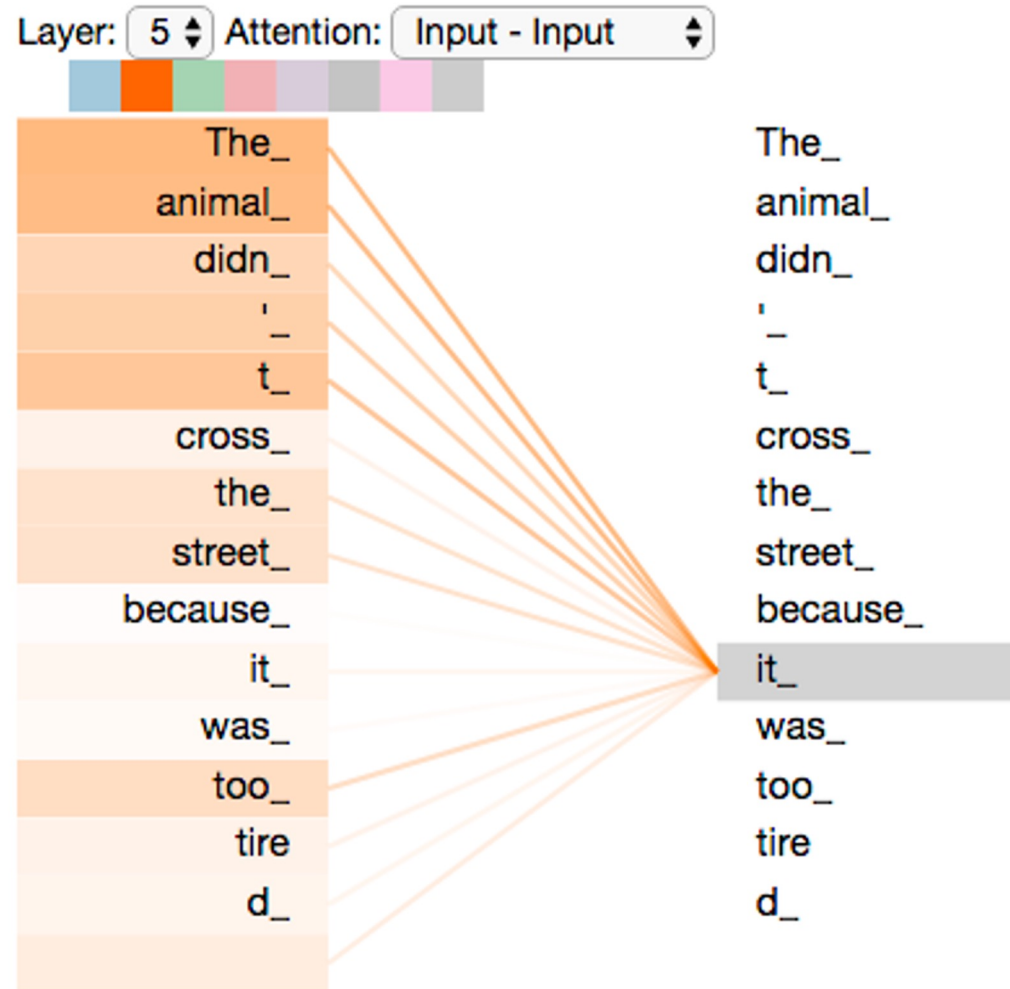


(a)

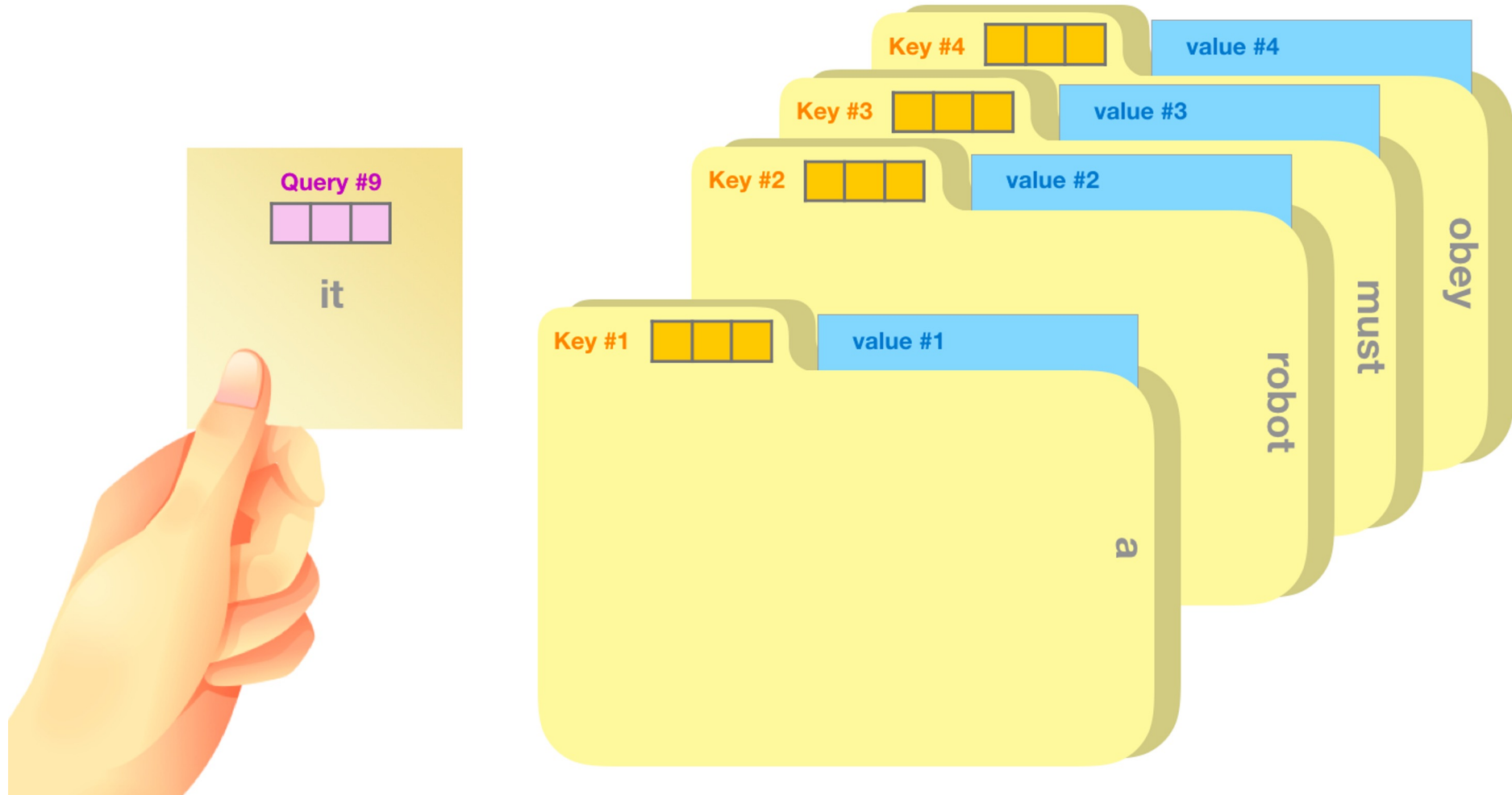


(b)

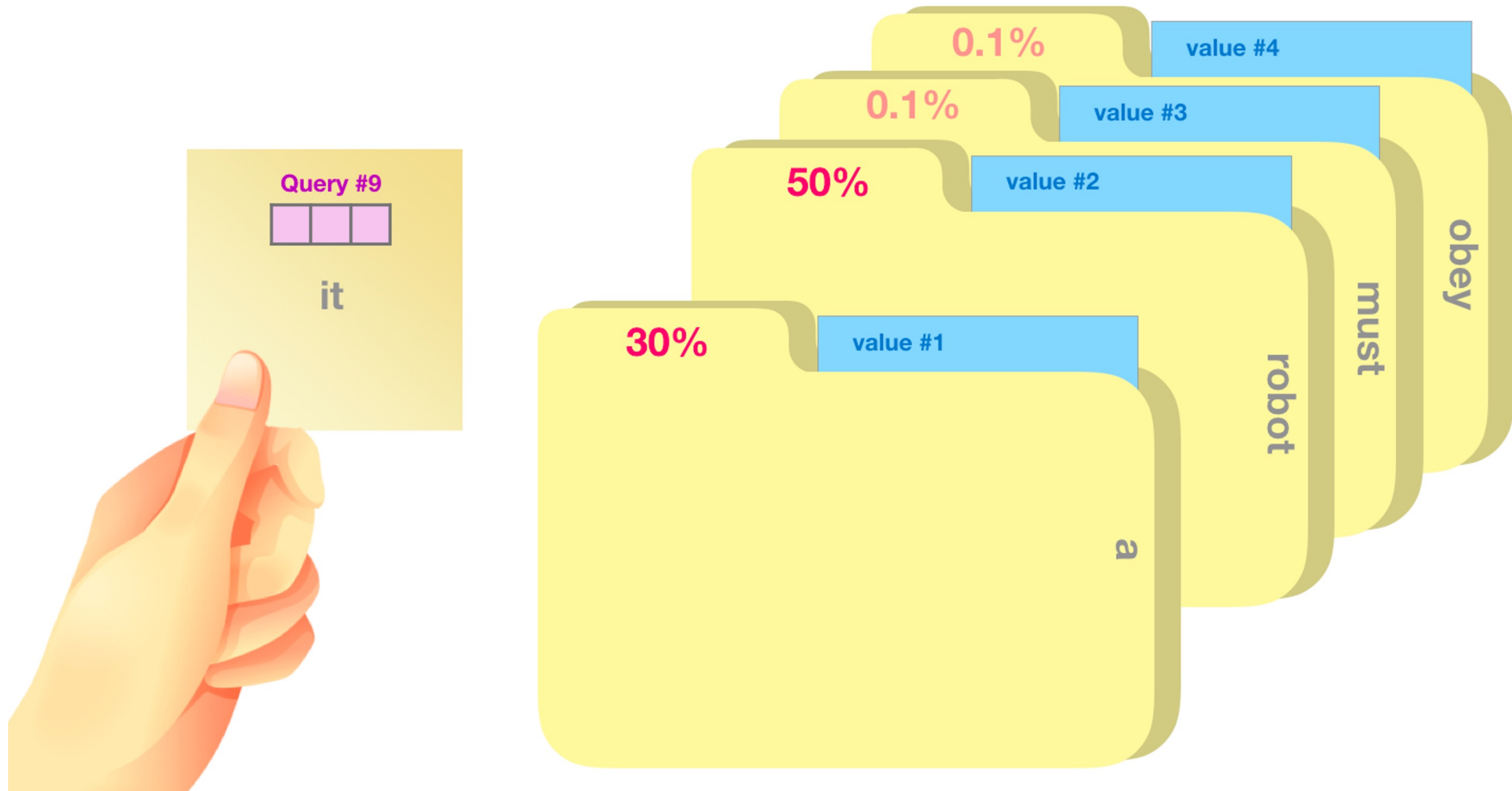
Self-Attention



Self-Attention



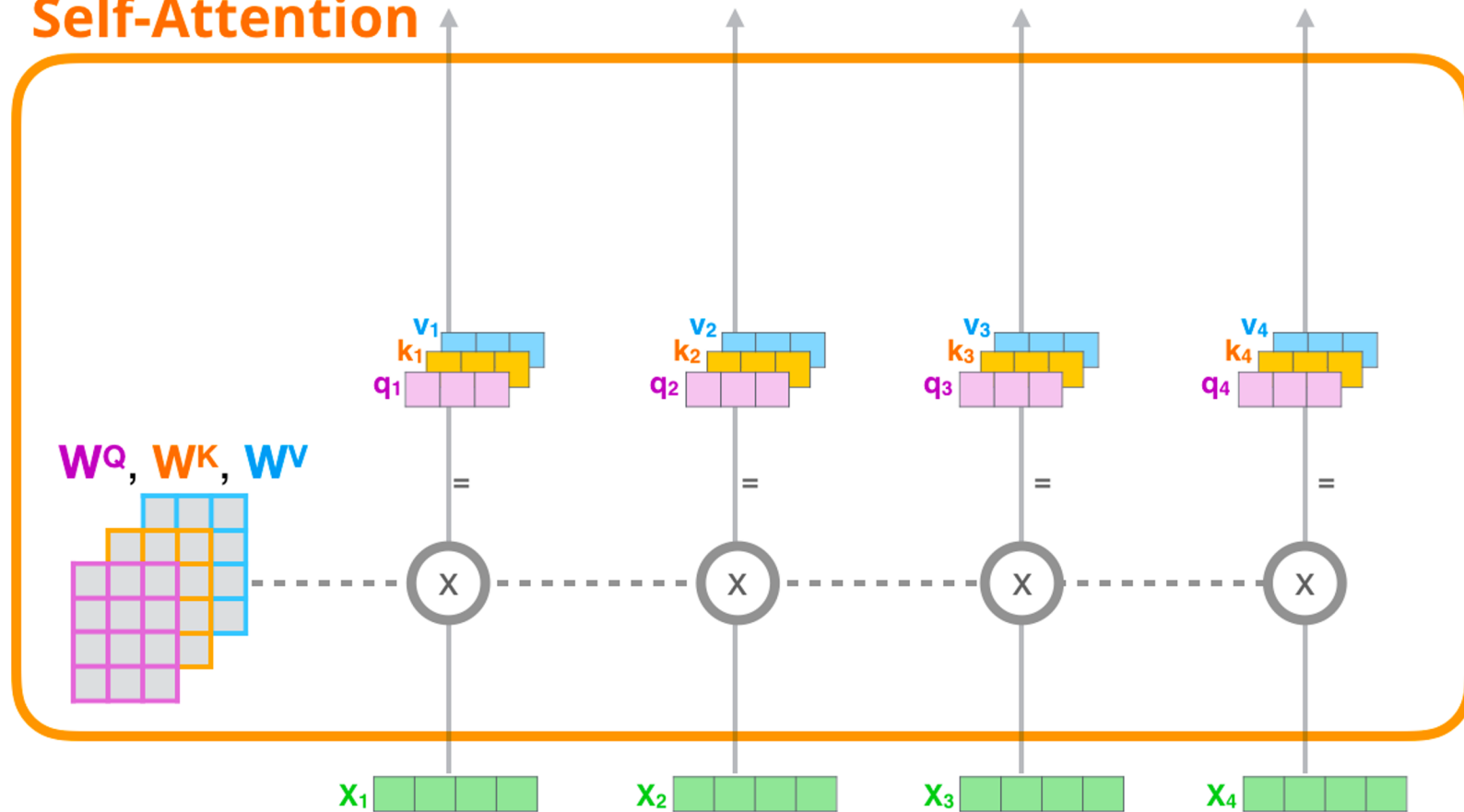
Self-Attention



Self-Attention

1) For each input token, create a **query vector**, a **key vector**, and a **value vector** by multiplying by weight Matrices W^Q , W^K , W^V

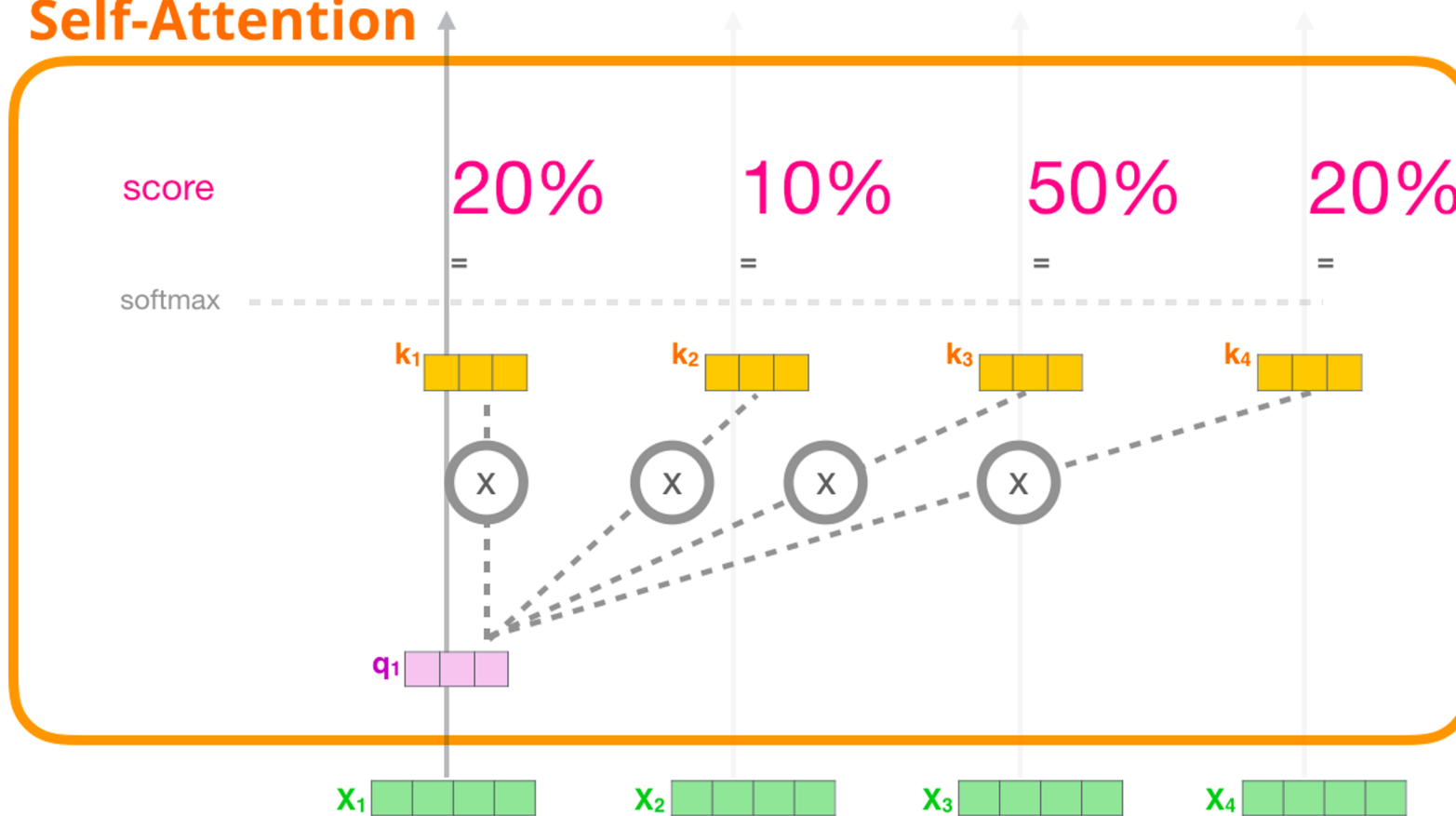
Self-Attention



Self-Attention

2) Multiply (dot product) the current **query vector**, by all the **key vectors**, to get a score of how well they match

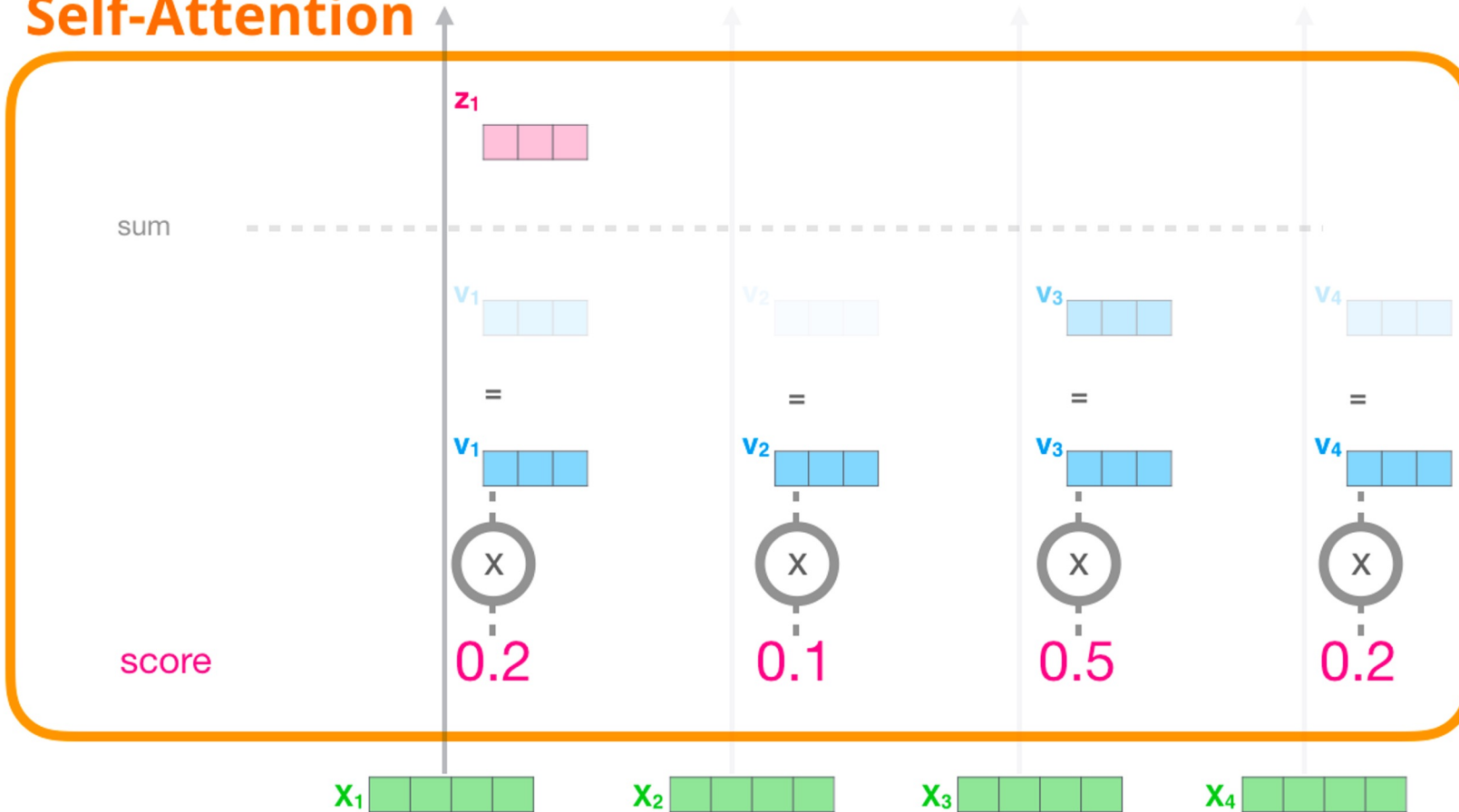
Self-Attention



Self-Attention

3) Multiply the **value vectors** by the **scores**, then sum up

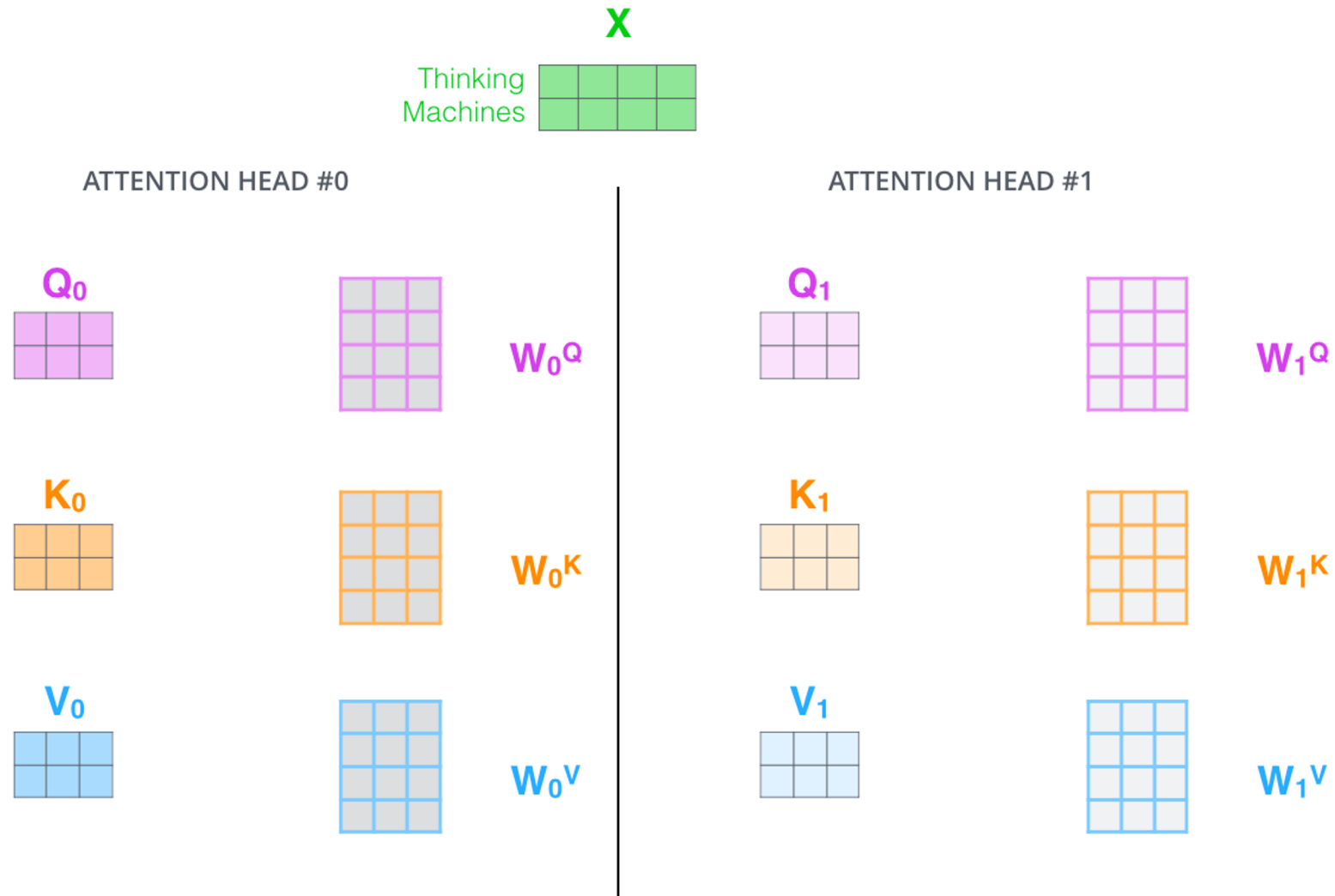
Self-Attention



Self-Attention

$$\text{softmax} \left(\frac{\begin{matrix} \text{Q} & & \text{K}^T \\ \begin{matrix} \square & \square & \square \\ \square & \square & \square \end{matrix} & \times & \begin{matrix} \square & \square \\ \square & \square \\ \square & \square \end{matrix} \end{matrix} }{\sqrt{d_k}} \right) \begin{matrix} \text{V} \\ \begin{matrix} \square & \square & \square \\ \square & \square & \square \end{matrix} \end{matrix} \\ = \begin{matrix} \text{Z} \\ \begin{matrix} \square & \square & \square \\ \square & \square & \square \end{matrix} \end{matrix}$$

Multi-Head Attention



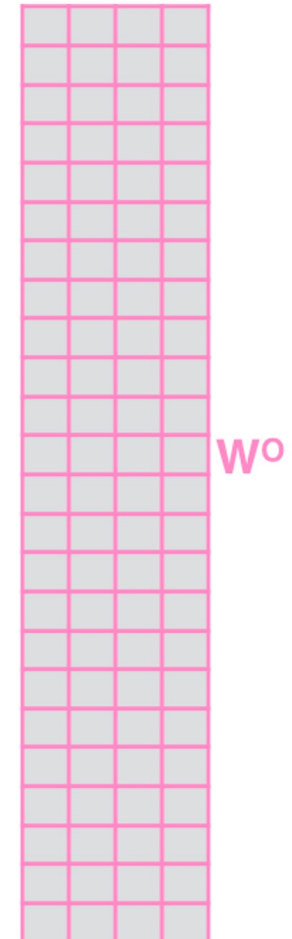
Multi-Head Attention

1) Concatenate all the attention heads

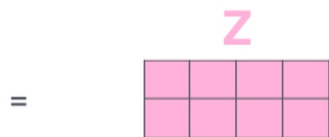


2) Multiply with a weight matrix W^O that was trained jointly with the model

X



3) The result would be the Z matrix that captures information from all the attention heads. We can send this forward to the FFNN



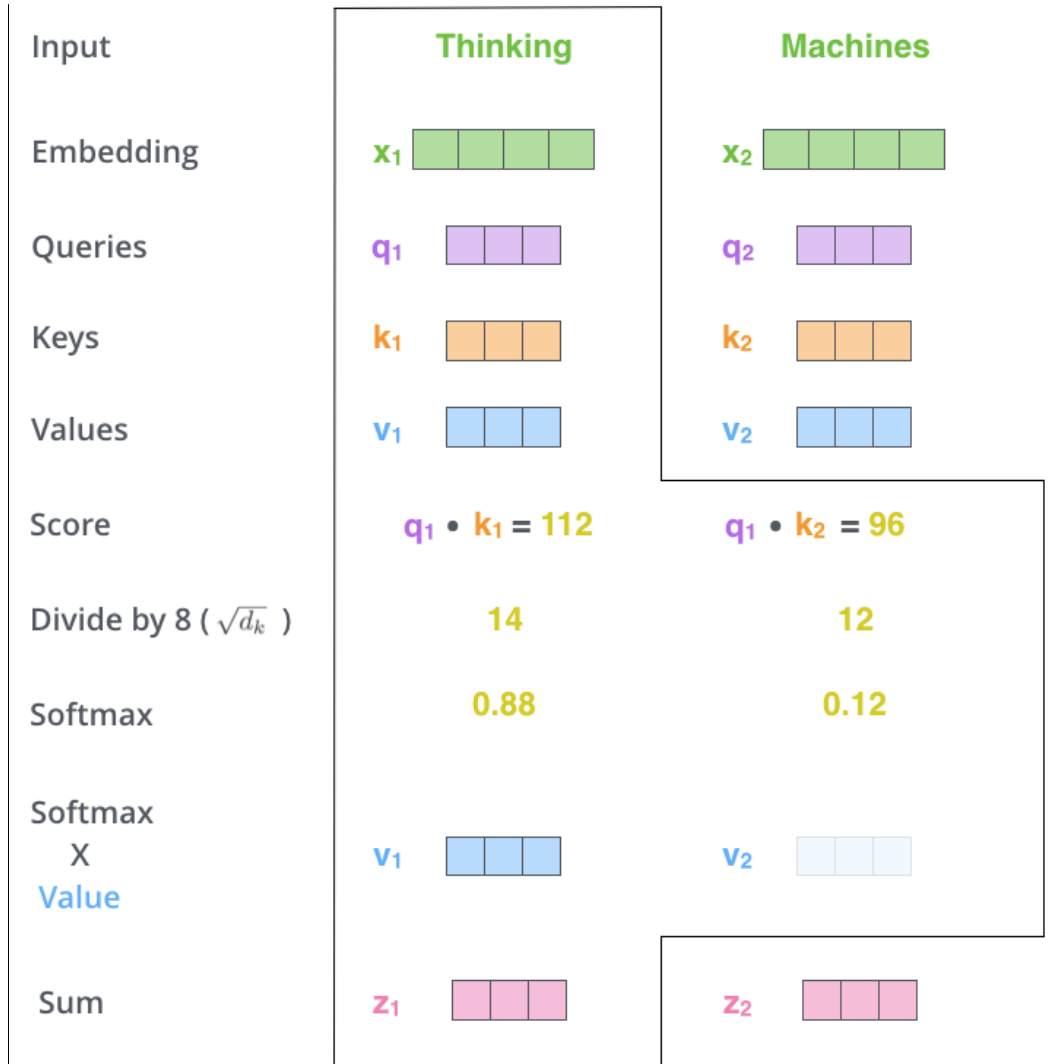
Transformers

Instead of an RNN, just use attention

High throughput & expressivity: compute queries, keys and values as (different) linear transformations of the input.

Attention weights are queries • keys;
outputs are sums of weighted values.

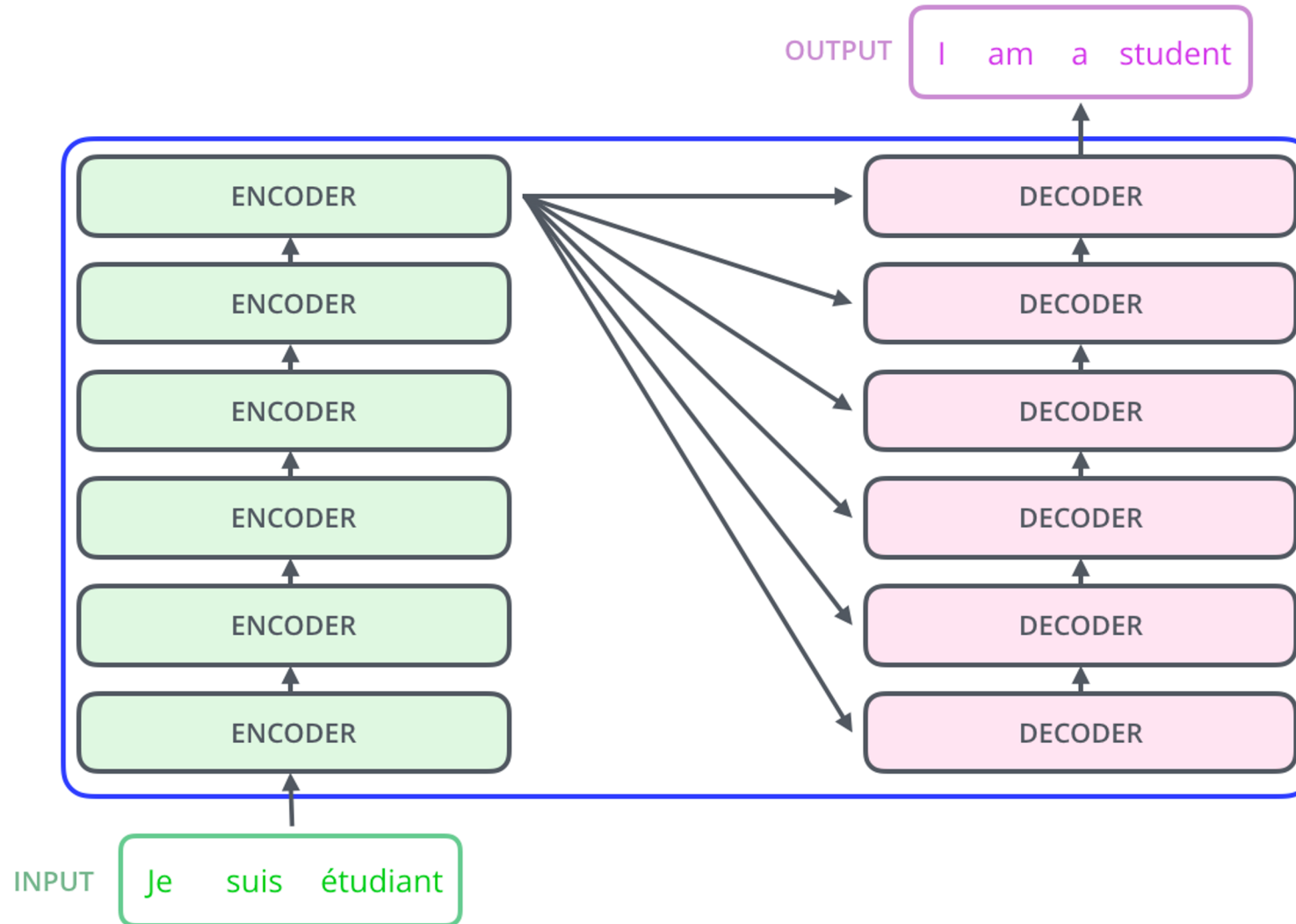
$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$



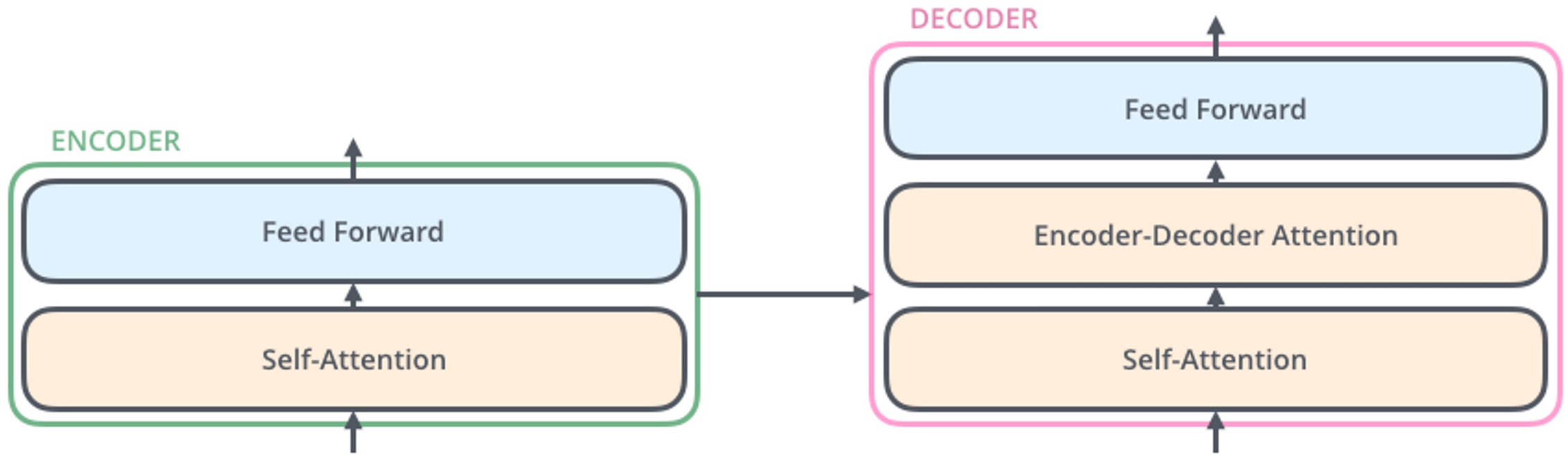
Transformer



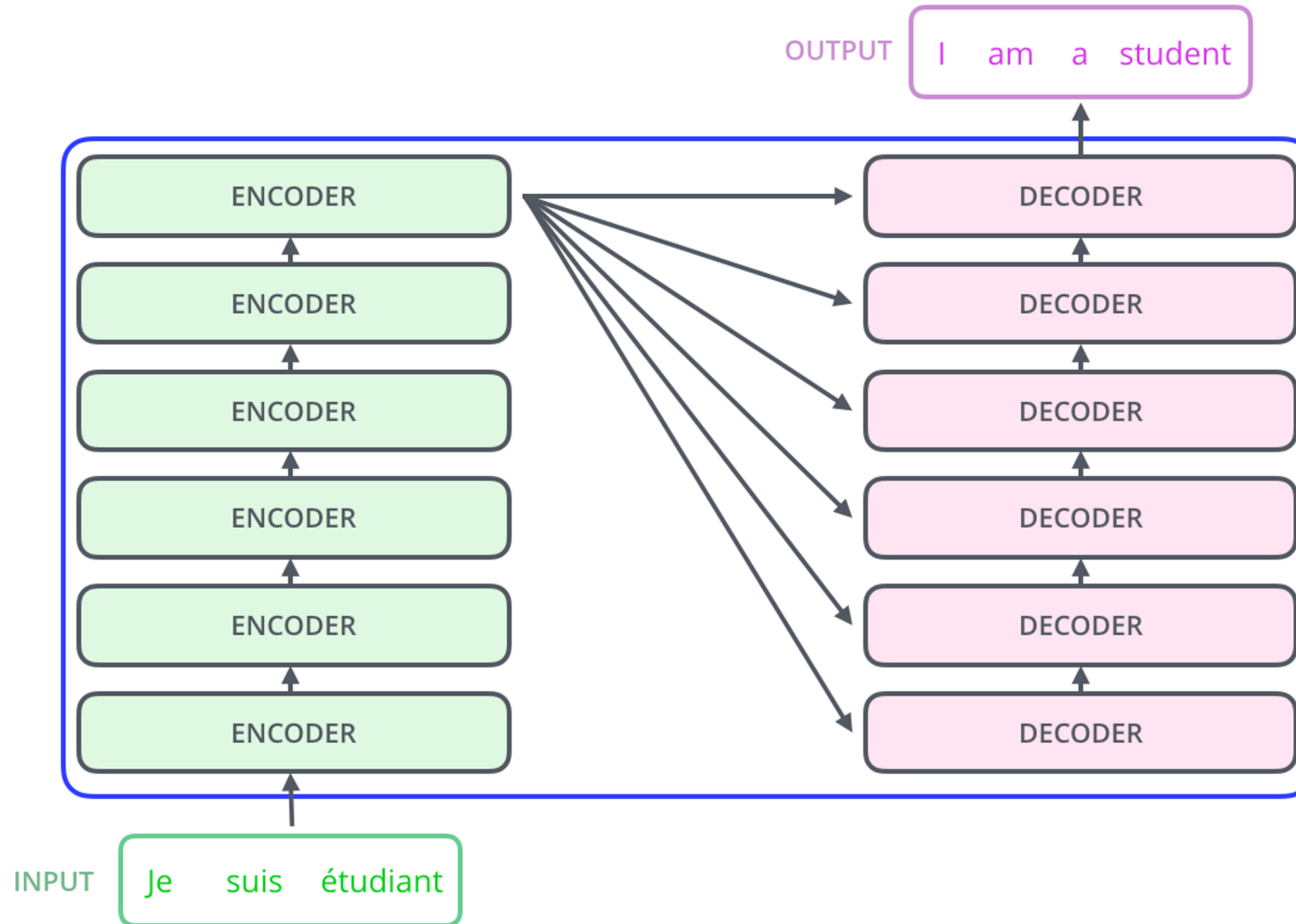
Transformer



Transformer

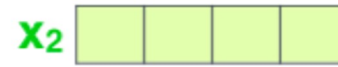
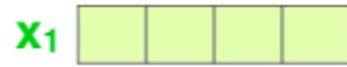


Transformer



Transformer Input

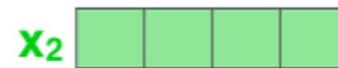
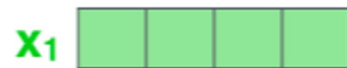
EMBEDDING
WITH TIME
SIGNAL



POSITIONAL
ENCODING



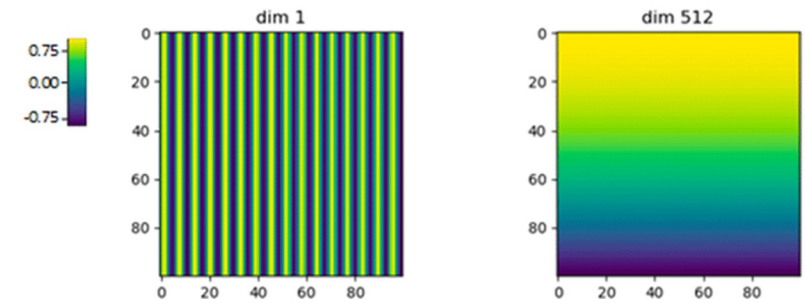
EMBEDDINGS



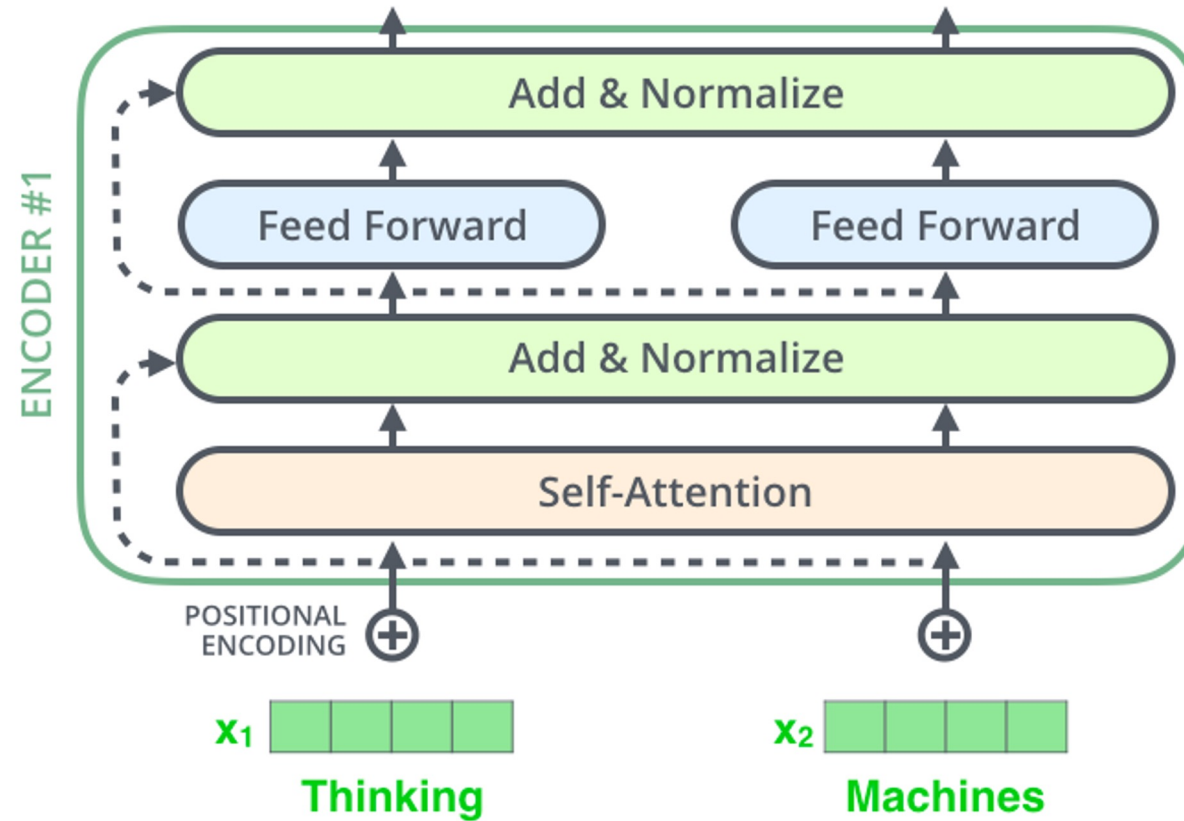
INPUT

Je

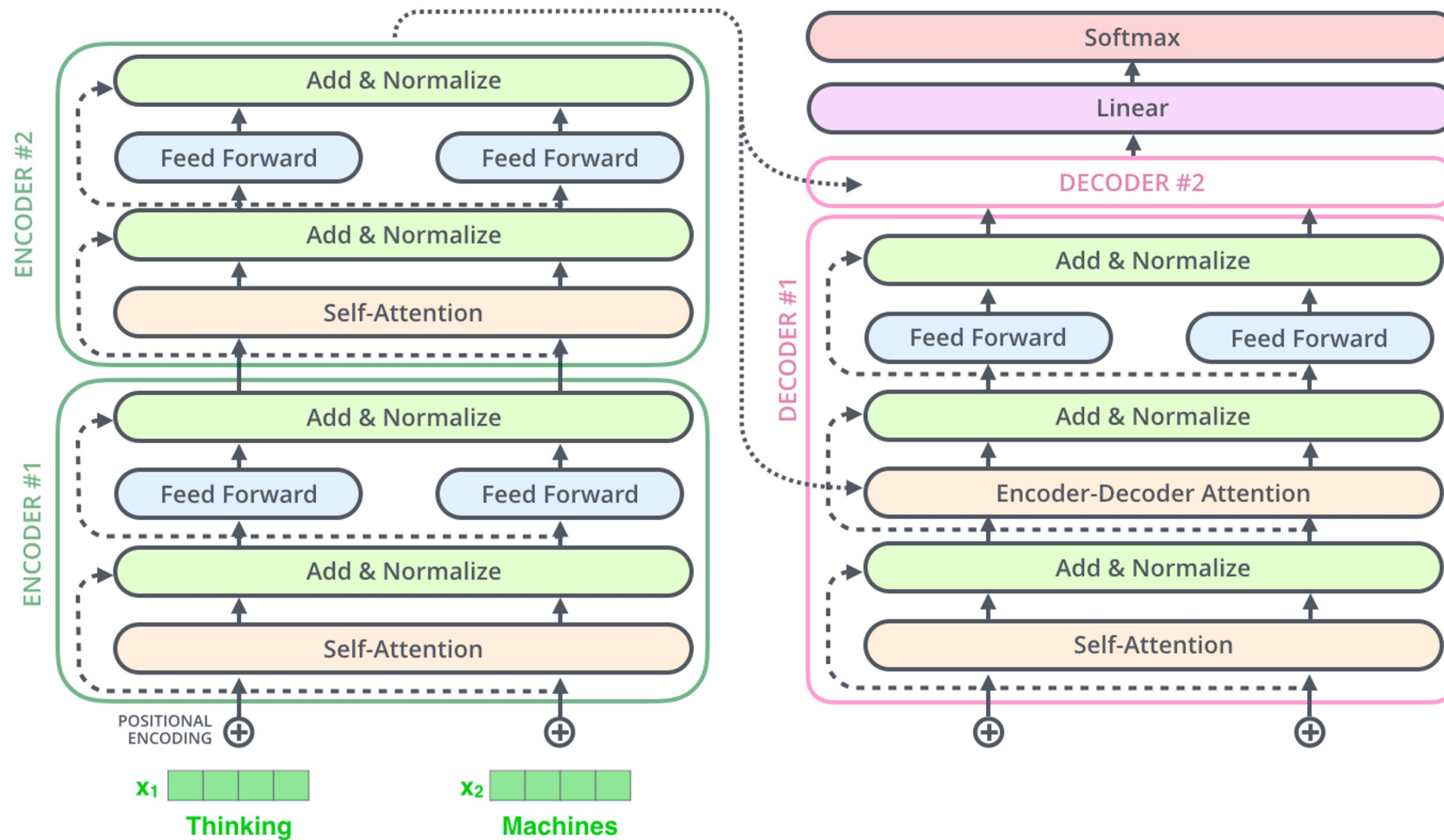
suis



Transformer Encoder

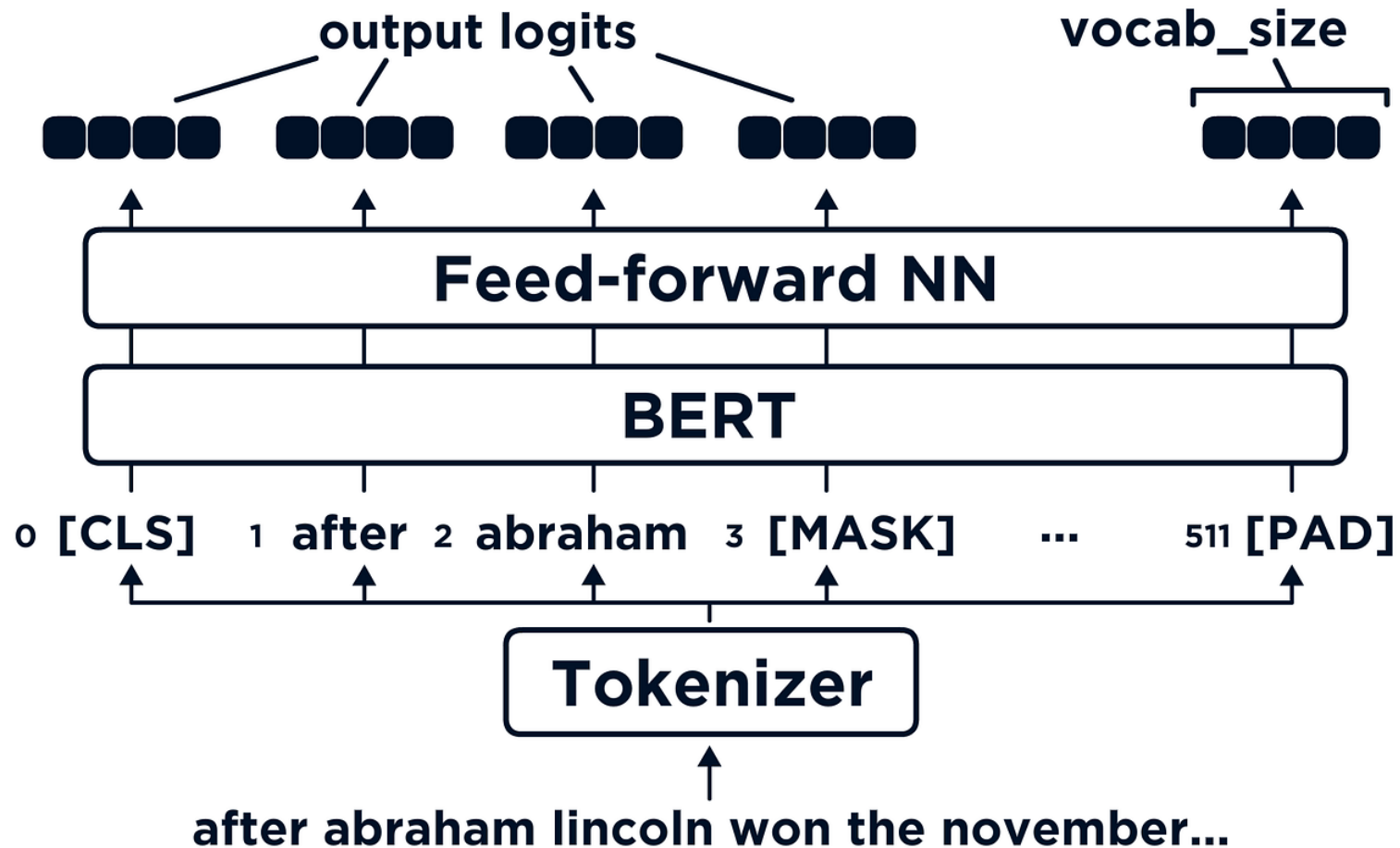


Full Transformer: Adding the Decoder

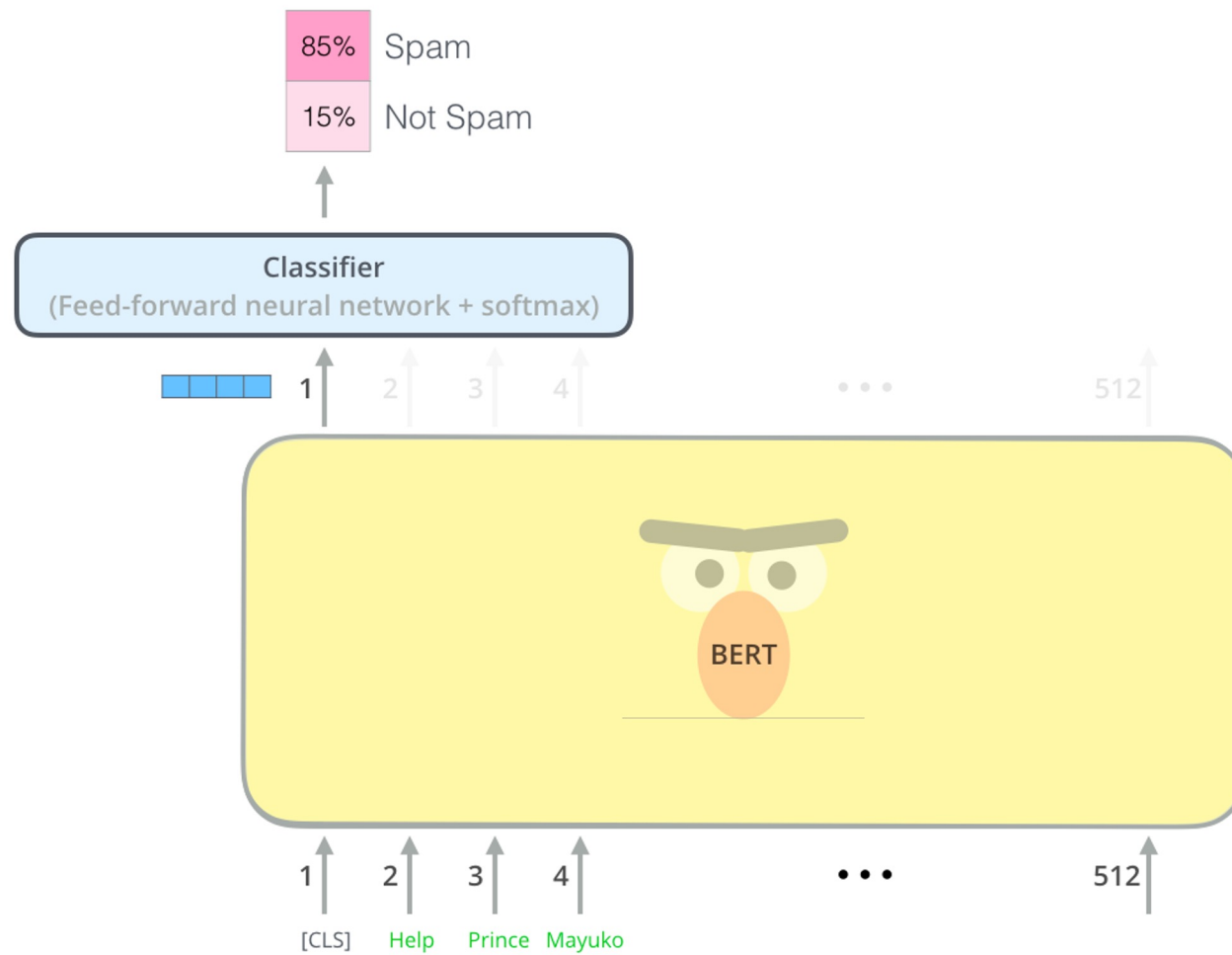


Masked Language Models

Key idea: learn representations and then fine-tune (training \neq inference)

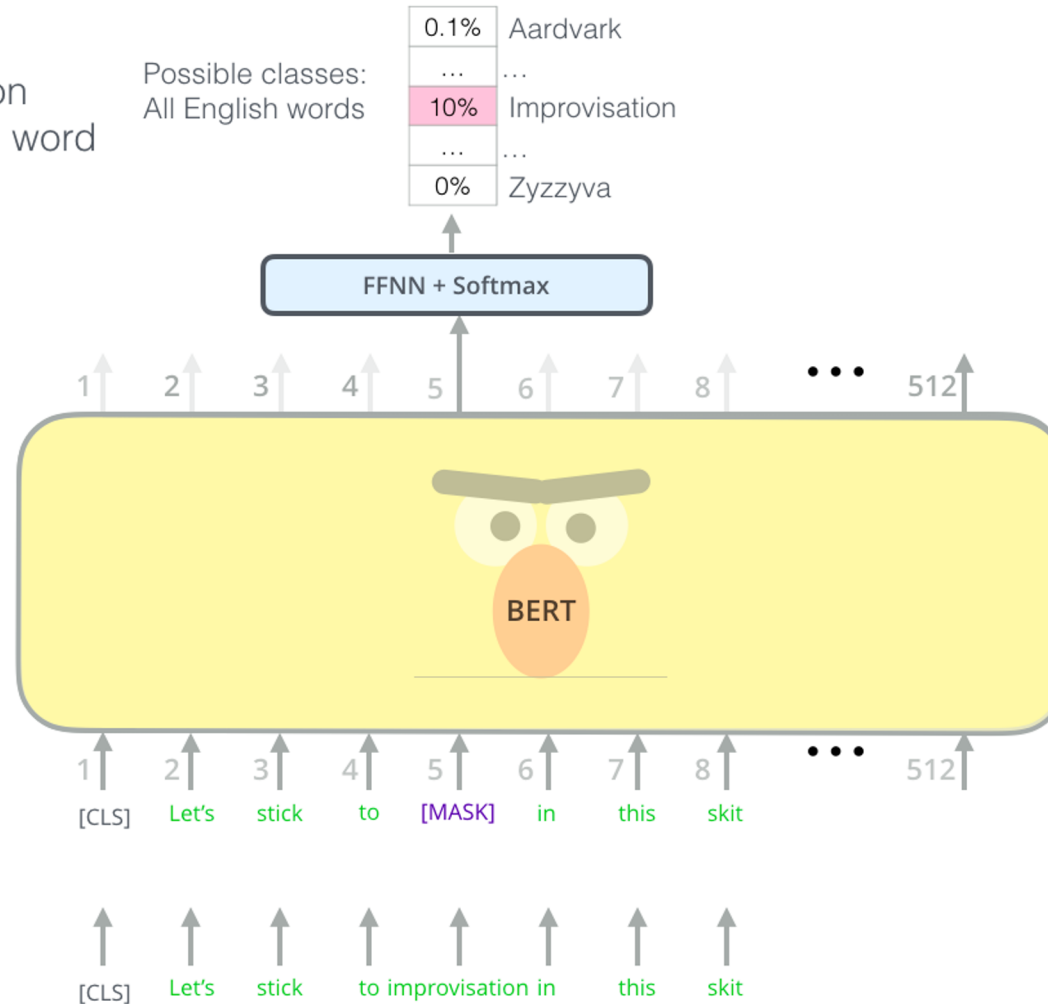


BERT



BERT

Use the output of the masked word's position to predict the masked word

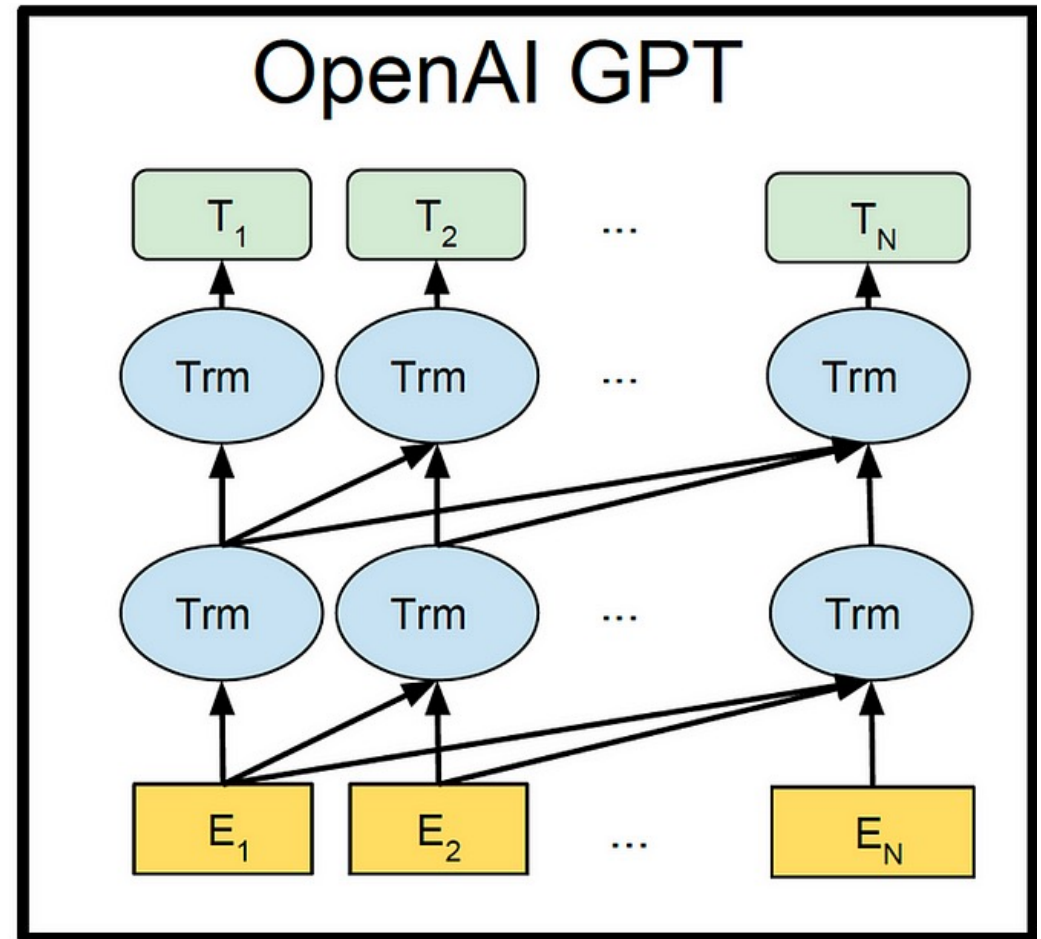
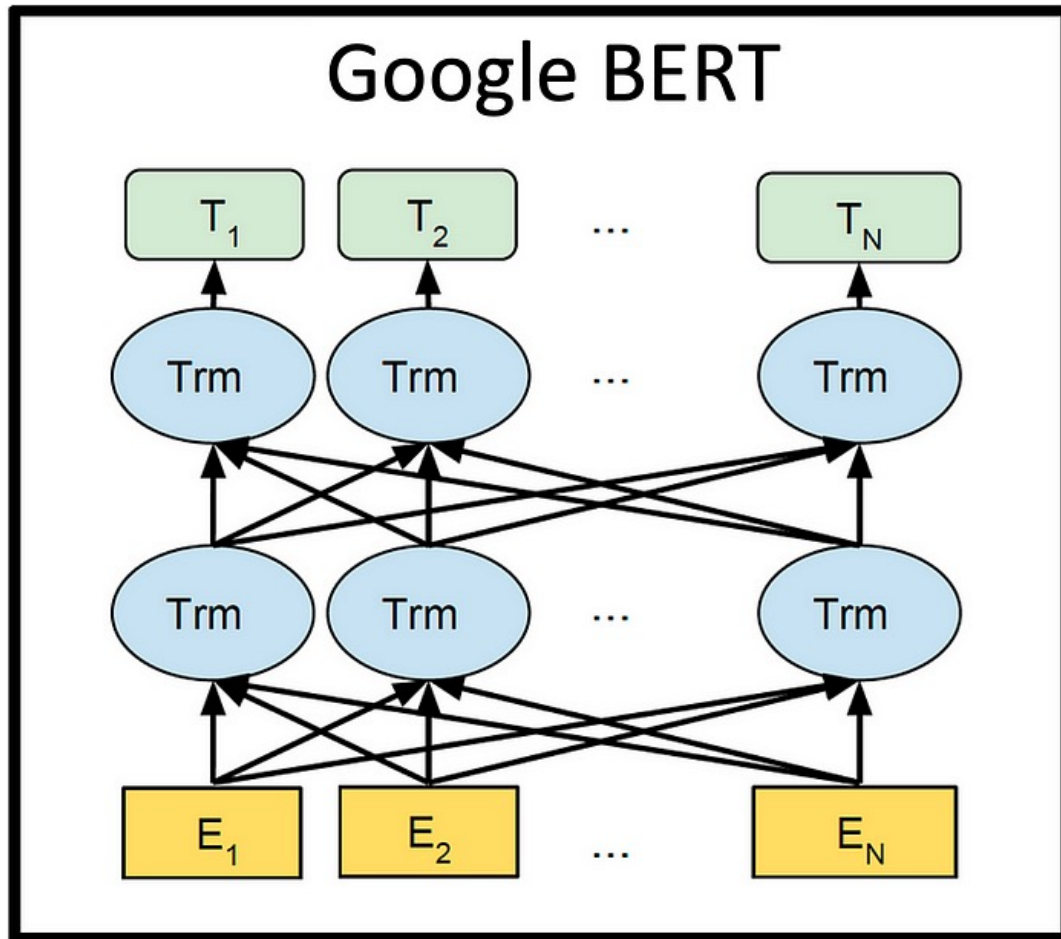


Randomly mask 15% of tokens

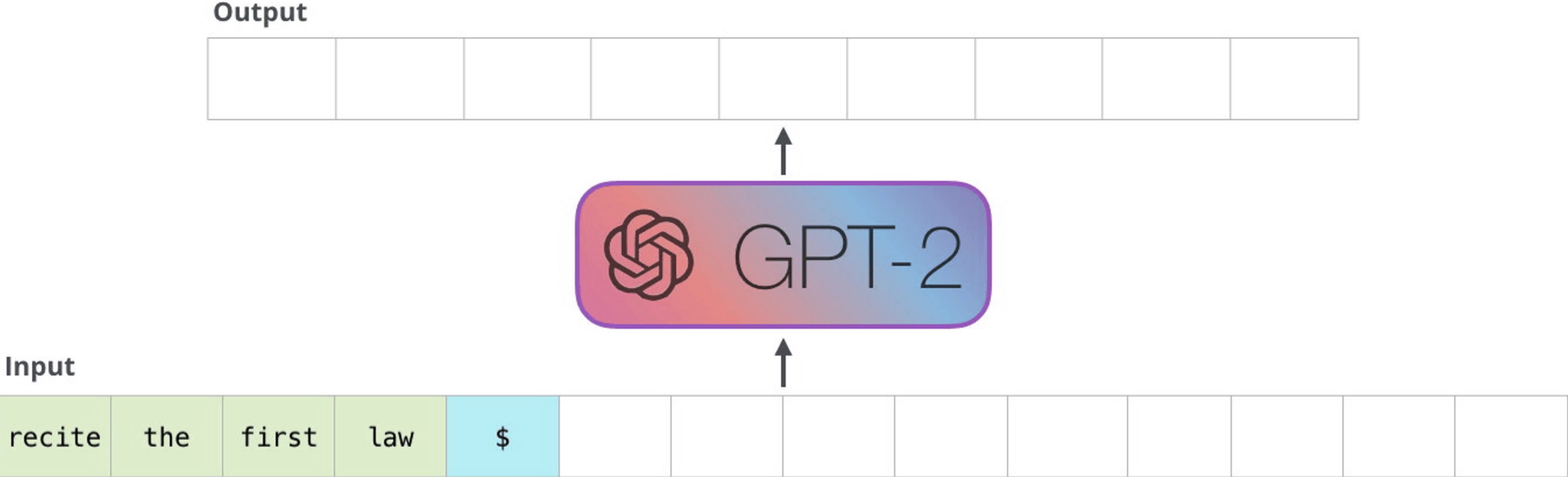
Input

Autoregressive Language Models

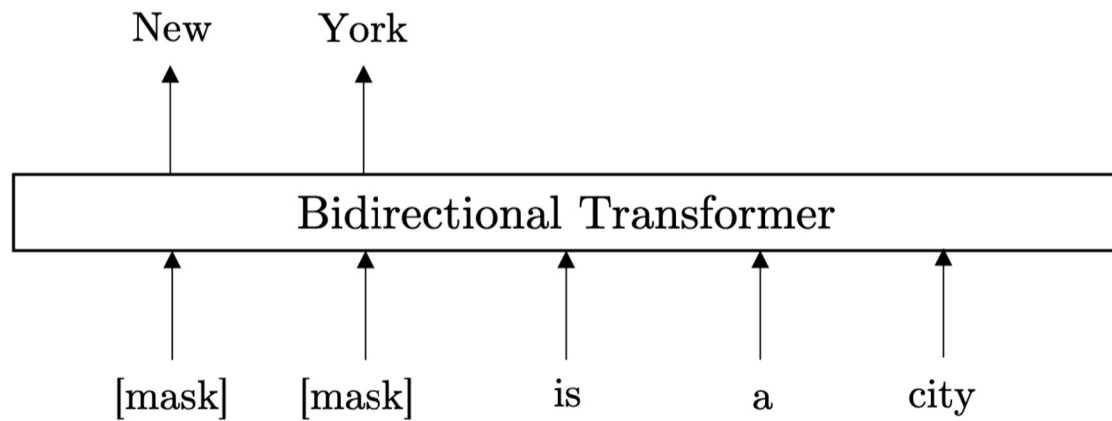
Key idea: learn next-word prediction directly (training = inference)



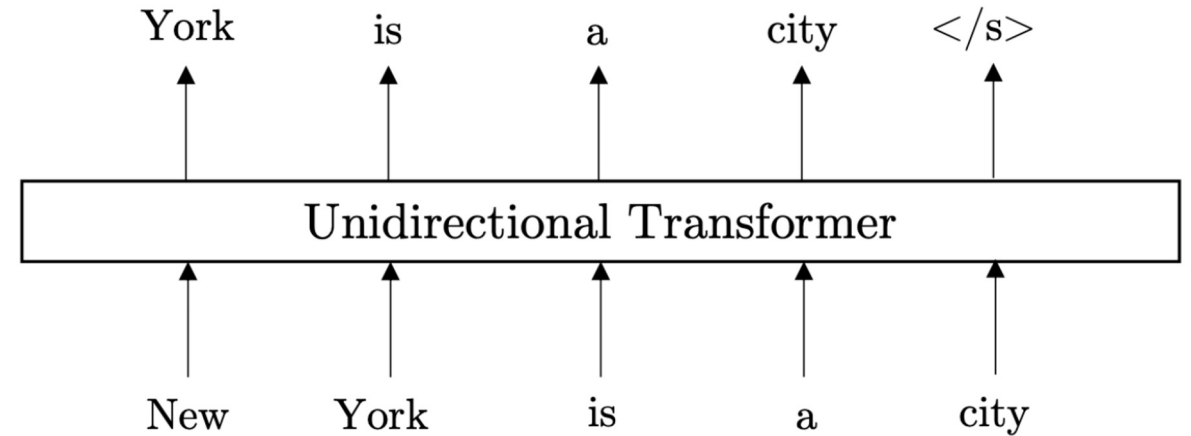
GPT Models



Masked vs. Autoregressive Language Modeling

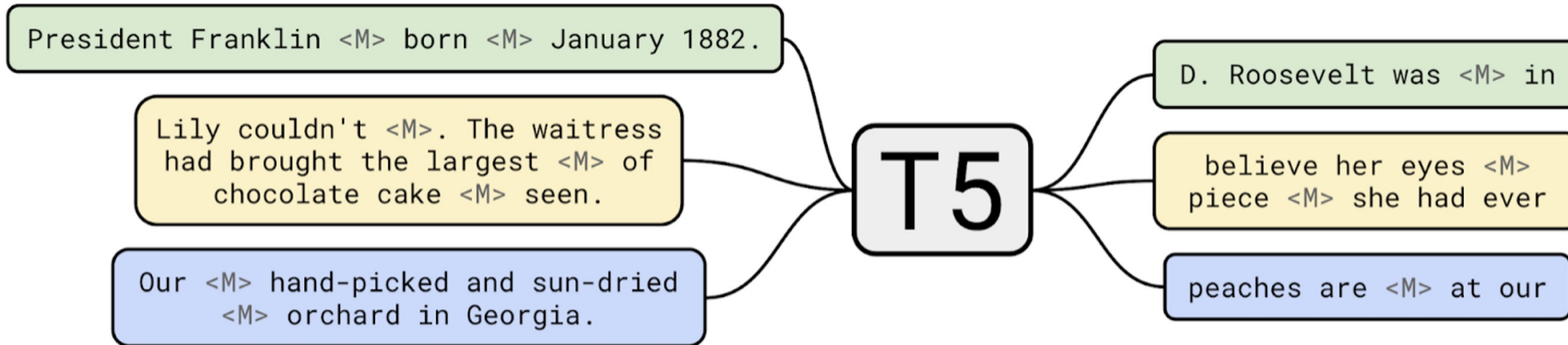


$$\log p(\bar{\mathbf{x}}|\hat{\mathbf{x}}) = \sum_{t=1}^T \text{mask}_t \log p(x_t|\hat{\mathbf{x}})$$



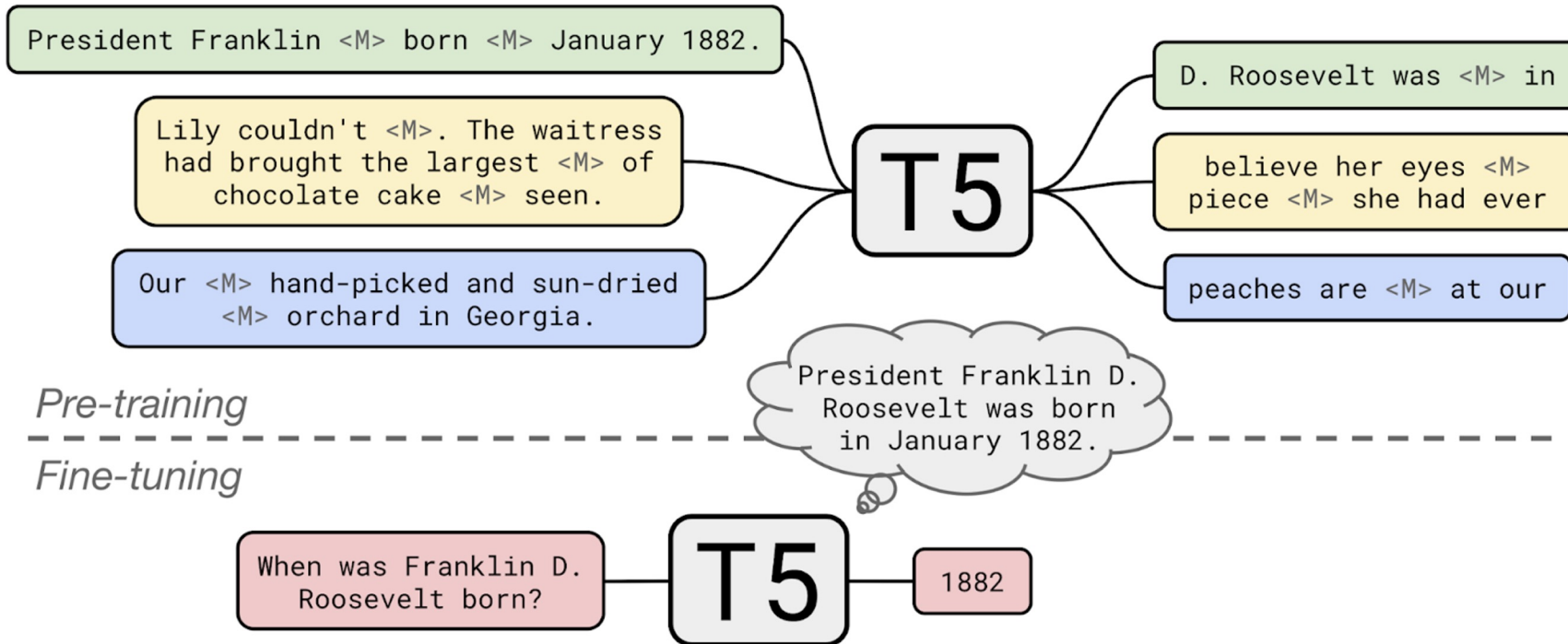
$$\log p(\mathbf{x}) = \sum_{t=1}^T \log p(x_t|\mathbf{x}_{<t})$$

Pretraining & Fine-tuning

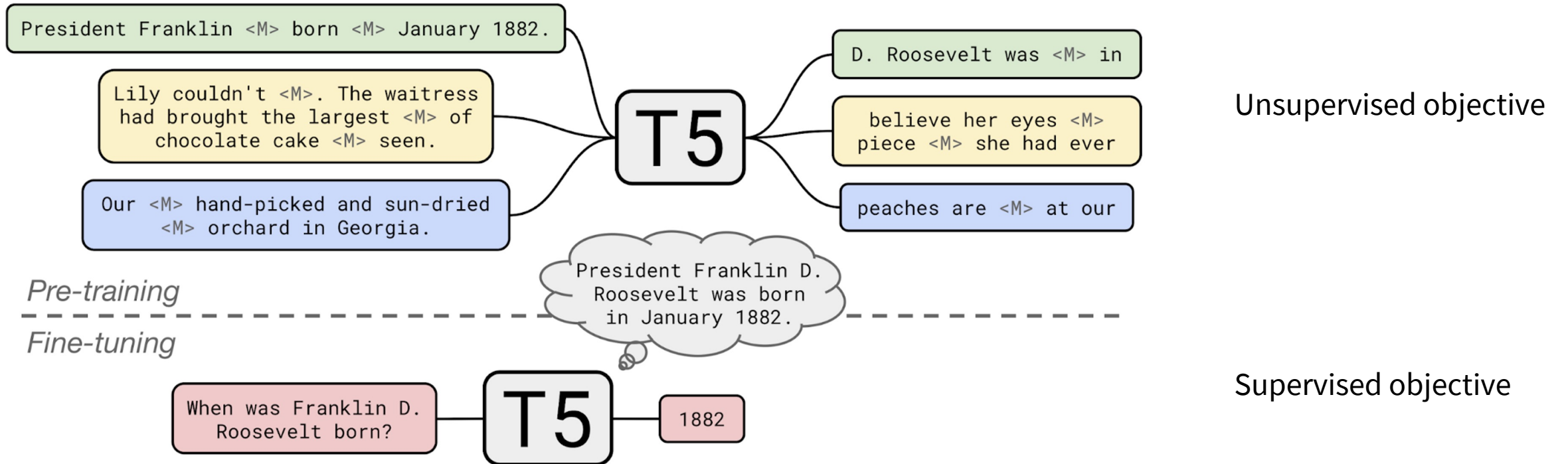


Pre-training

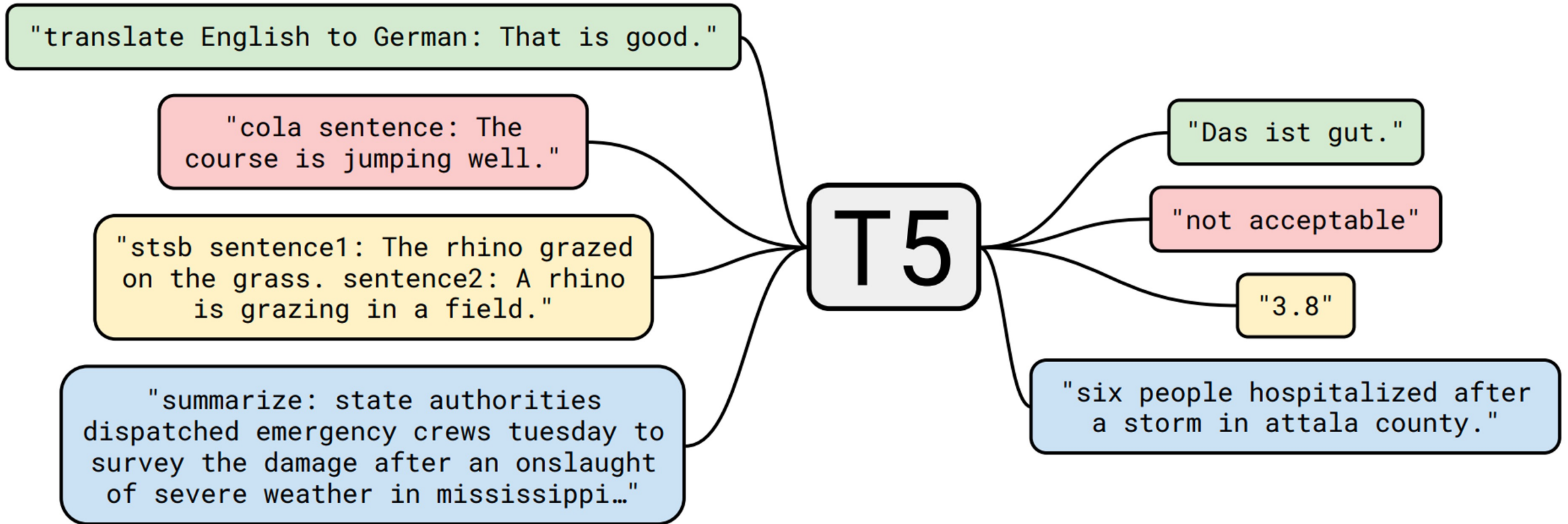
Pretraining & Fine-tuning



Pretraining & Fine-tuning



Prompting



Language Modeling Subsumes All Tasks (?)



Can you translate "language models are powerful these days" into French?



Certainly! The translation of "language models are powerful these days" into French is:



"Les modèles de langage sont puissants de nos jours."



Can you write a syntactic parse for the sentence "Kids Make Nutritious Snacks"?



Sure! Here's a syntactic parse for the sentence "Kids Make Nutritious Snacks":

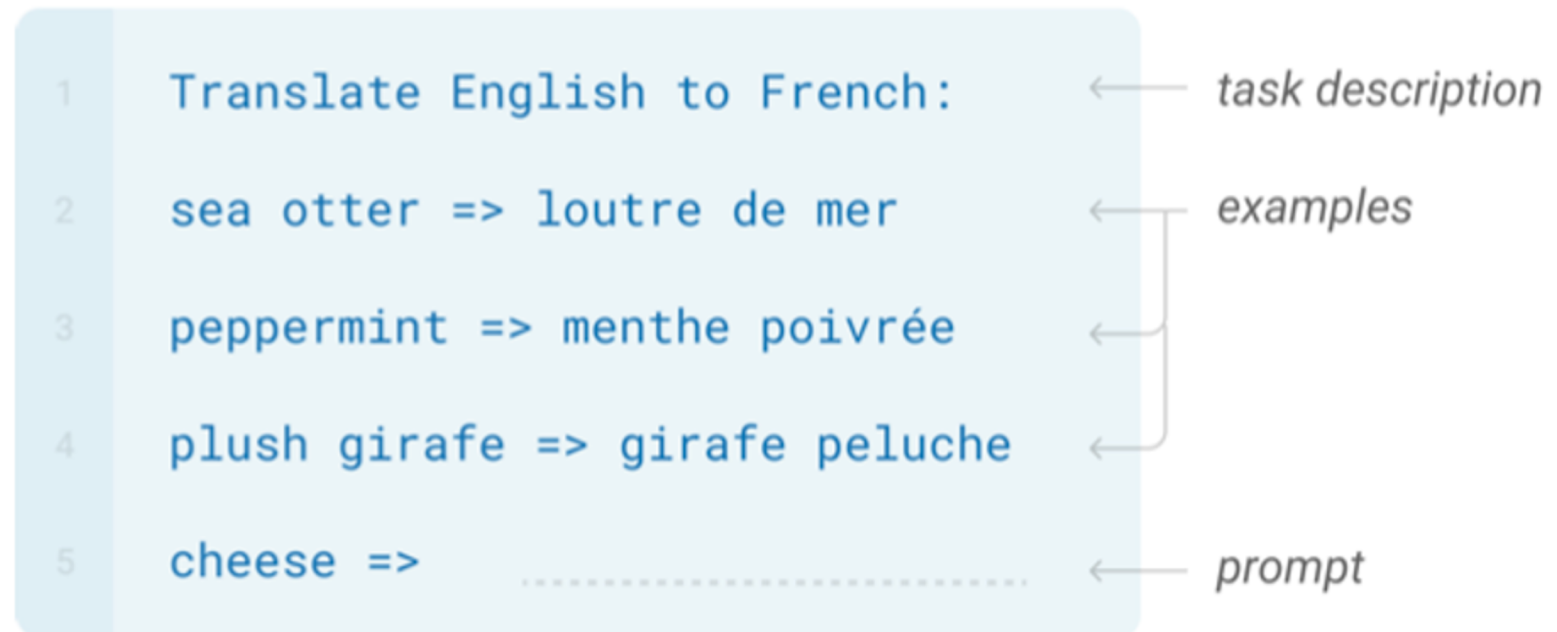


```
scss Copy code
(S
  (NP (NNS Kids))
  (VP (VBP Make)
    (NP (JJ Nutritious) (NNS Snacks))))
```

Few- & Zero-Shot Learning

Few-shot

In addition to the task description, the model sees a few examples of the task. No gradient updates are performed.



Few- & Zero-Shot Learning

Few-shot

In addition to the task description, the model sees a few examples of the task. No gradient updates are performed.

```
1 Translate English to French: ← task description
2 sea otter => loutre de mer ← examples
3 peppermint => menthe poivrée ←
4 plush girafe => girafe peluche ←
5 cheese => ..... ← prompt
```

One-shot

In addition to the task description, the model sees a single example of the task. No gradient updates are performed.

```
1 Translate English to French: ← task description
2 sea otter => loutre de mer ← example
3 cheese => ..... ← prompt
```

Few- & Zero-Shot Learning

Few-shot

In addition to the task description, the model sees a few examples of the task. No gradient updates are performed.

```
1 Translate English to French: ← task description
2 sea otter => loutre de mer ← examples
3 peppermint => menthe poivrée ←
4 plush girafe => girafe peluche ←
5 cheese => ..... ← prompt
```

One-shot

In addition to the task description, the model sees a single example of the task. No gradient updates are performed.

```
1 Translate English to French: ← task description
2 sea otter => loutre de mer ← example
3 cheese => ..... ← prompt
```

Zero-shot

The model predicts the answer given only a natural language description of the task. No gradient updates are performed.

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

Instruction Tuning and RLHF

Key issue: language modeling \neq assisting users

PROMPT *Explain the moon landing to a 6 year old in a few sentences.*

COMPLETION GPT-3

Explain the theory of gravity to a 6 year old.

Explain the theory of relativity to a 6 year old in a few sentences.

Explain the big bang theory to a 6 year old.

Explain evolution to a 6 year old.

Instruction Tuning and RLHF

PROMPT *Explain the moon landing to a 6 year old in a few sentences.*

COMPLETION

GPT-3

Explain the theory of gravity to a 6 year old.

Explain the theory of relativity to a 6 year old in a few sentences.

Explain the big bang theory to a 6 year old.

Explain evolution to a 6 year old.

InstructGPT

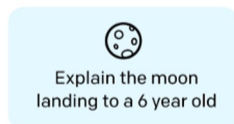
People went to the moon, and they took pictures of what they saw, and sent them back to the earth so we could all see them.

Reinforcement Learning from Human Feedback

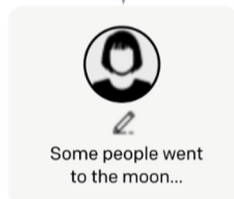
Step 1

**Collect demonstration data,
and train a supervised policy.**

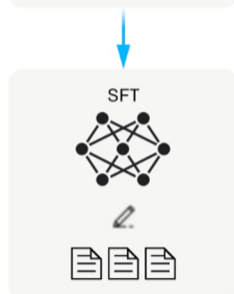
A prompt is
sampled from our
prompt dataset.



A labeler
demonstrates the
desired output
behavior.



This data is used
to fine-tune GPT-3
with supervised
learning.



Reinforcement Learning from Human Feedback

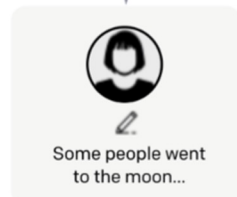
Step 1

**Collect demonstration data,
and train a supervised policy.**

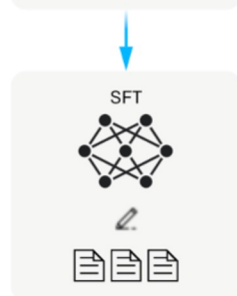
A prompt is
sampled from our
prompt dataset.



A labeler
demonstrates the
desired output
behavior.



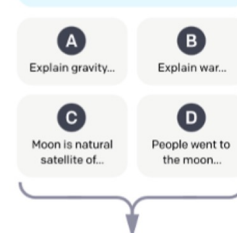
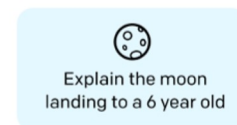
This data is used
to fine-tune GPT-3
with supervised
learning.



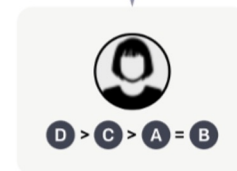
Step 2

**Collect comparison data,
and train a reward model.**

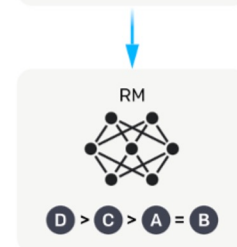
A prompt and
several model
outputs are
sampled.



A labeler ranks
the outputs from
best to worst.



This data is used
to train our
reward model.

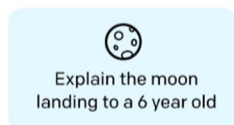


Reinforcement Learning from Human Feedback

Step 1

Collect demonstration data, and train a supervised policy.

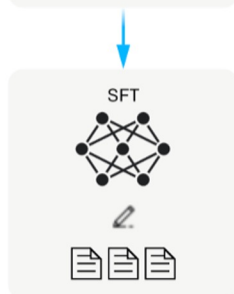
A prompt is sampled from our prompt dataset.



A labeler demonstrates the desired output behavior.



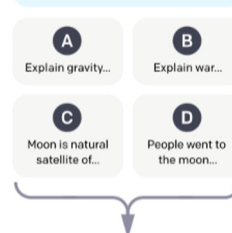
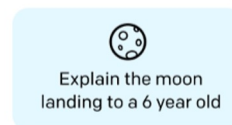
This data is used to fine-tune GPT-3 with supervised learning.



Step 2

Collect comparison data, and train a reward model.

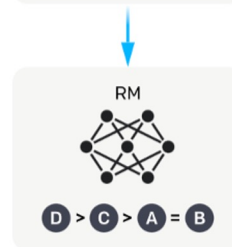
A prompt and several model outputs are sampled.



A labeler ranks the outputs from best to worst.



This data is used to train our reward model.



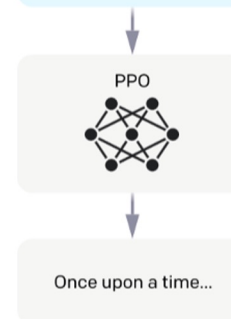
Step 3

Optimize a policy against the reward model using reinforcement learning.

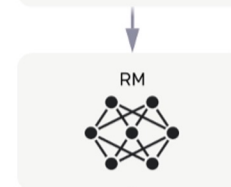
A new prompt is sampled from the dataset.



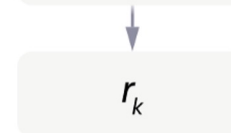
The policy generates an output.



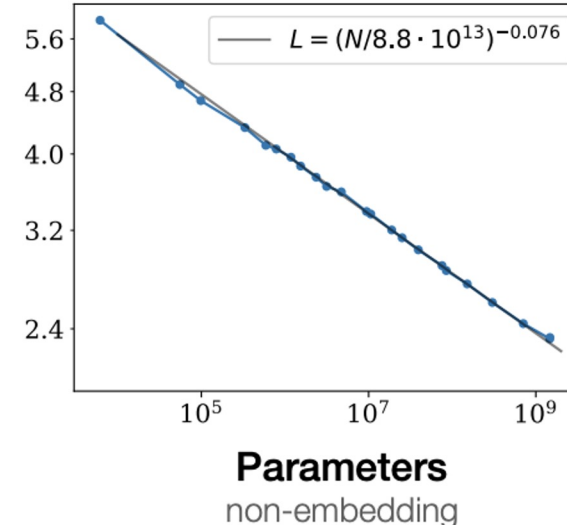
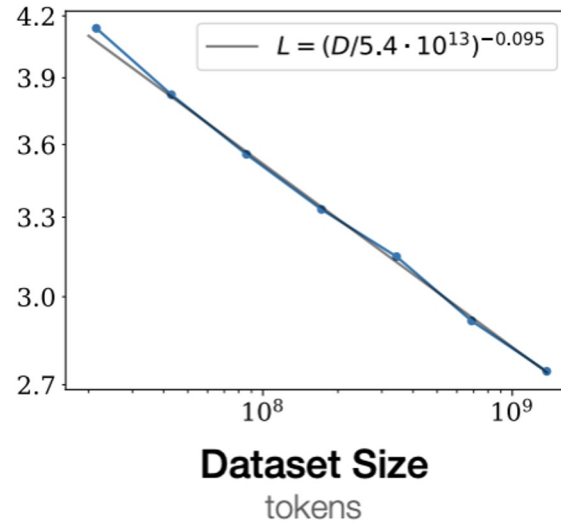
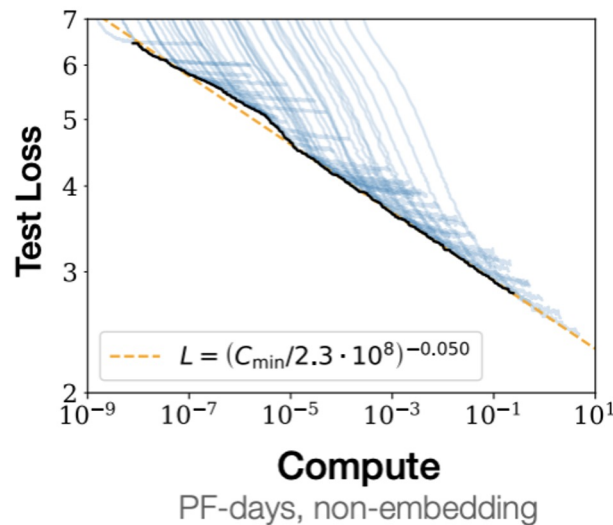
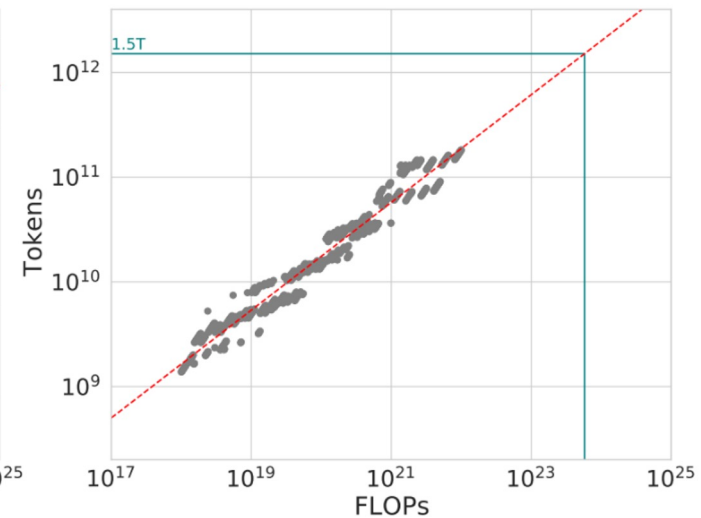
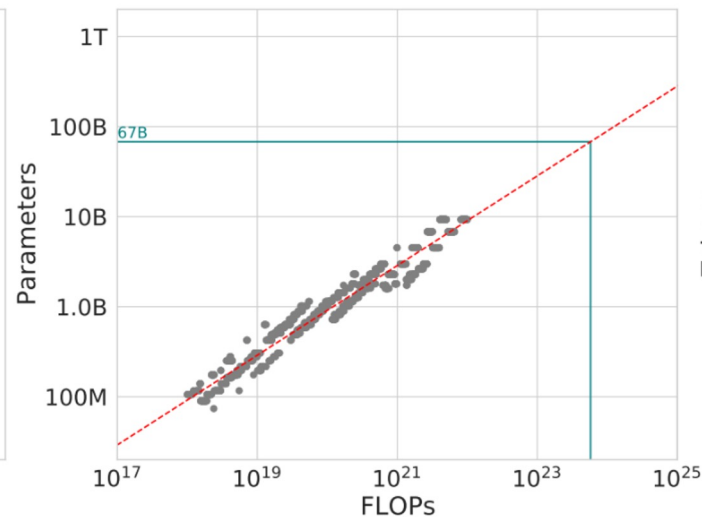
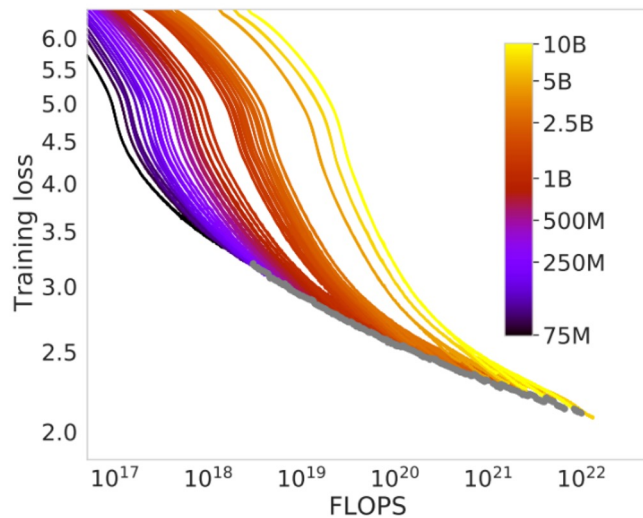
The reward model calculates a reward for the output.



The reward is used to update the policy using PPO.



Scaling Data & Compute



Going Forward

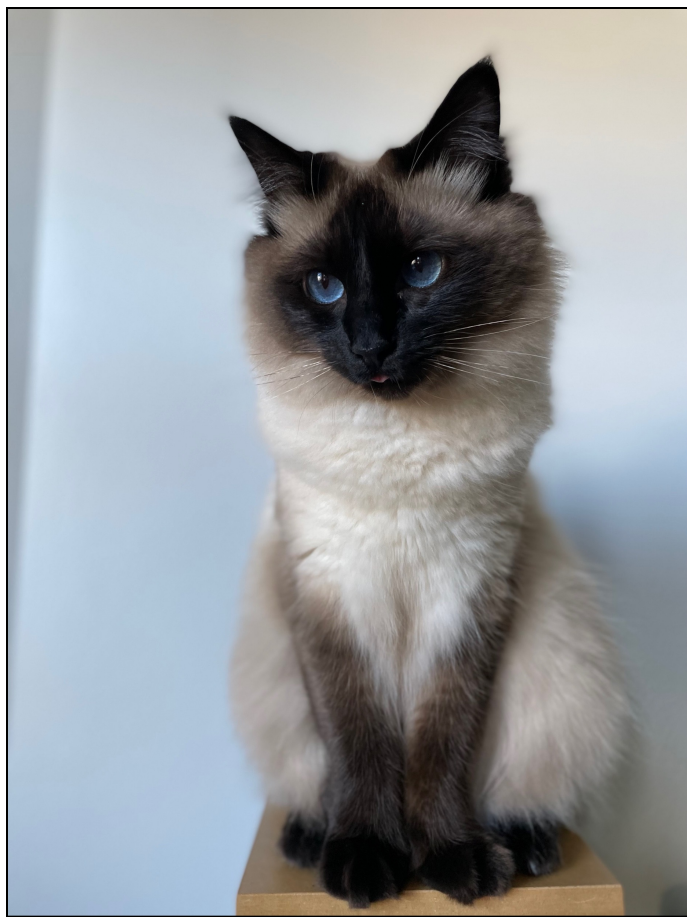
- Tool use (e.g., getting language models to use APIs)
- Grounding into non-linguistic inputs (e.g., vision, sensor data, etc.)
- Managing security & privacy
- Efficient / on-device / smaller / faster models
- Avoiding harmful or undesirable outputs
- Supporting multilinguality, esp. for low resource languages

Bonus: Computer Vision

What tasks do we care about?

- Object detection and classification
- Semantic segmentation
- Image captioning
- Visual question answering
- Video classification and understanding
- Image generation
- ...

Image Classification



cat

dog

horse

person

airplane

house

...

Beyond Image Classification

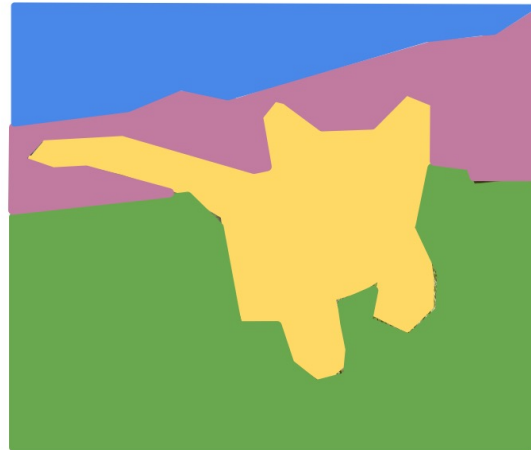
Classification



CAT

No spatial extent

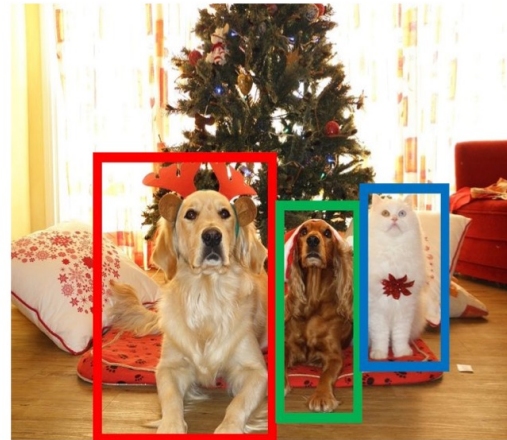
Semantic Segmentation



**GRASS, CAT,
TREE, SKY**

No objects, just pixels

Object Detection



DOG, DOG, CAT

Multiple Object

Instance Segmentation



DOG, DOG, CAT

This image is [CC0 public domain](#)

Image Generation

TEXT PROMPT an armchair in the shape of an avocado. an armchair imitating an avocado.

AI-GENERATED
IMAGES



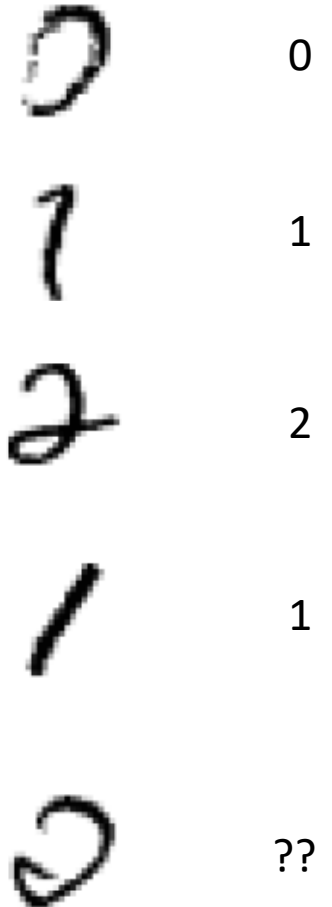
Recall: MNIST Digit Classification

Task specification:

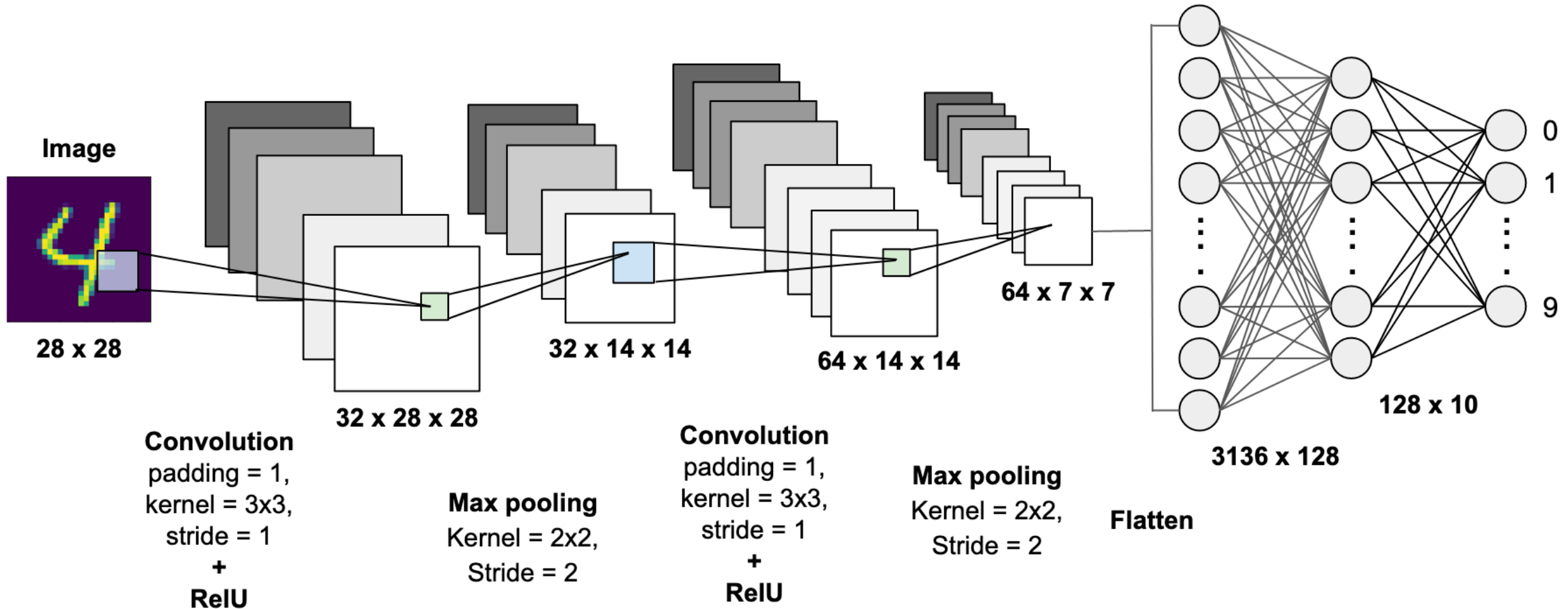
- Input features: binary pixel values
- Output: a digit classification (0-9)

Issues with Naïve Bayes classifier:

- Can overfit to individual pixels
- Not robust to scaling, movement left/right, etc.

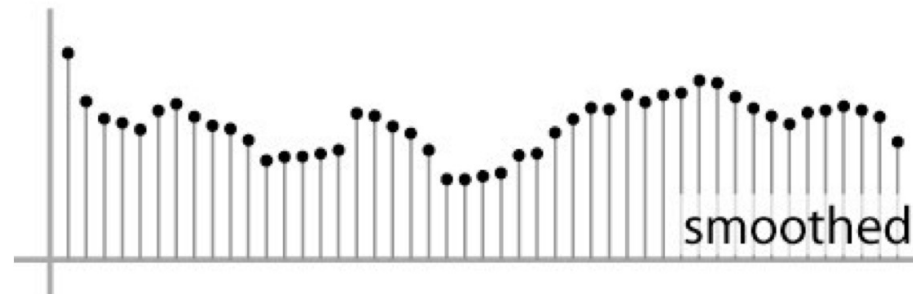
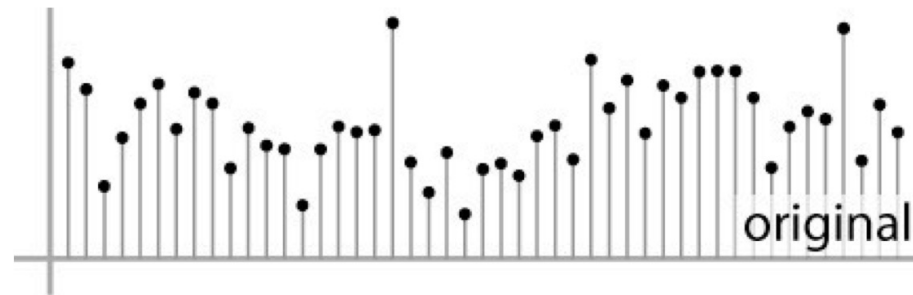


Convolutional Neural Networks



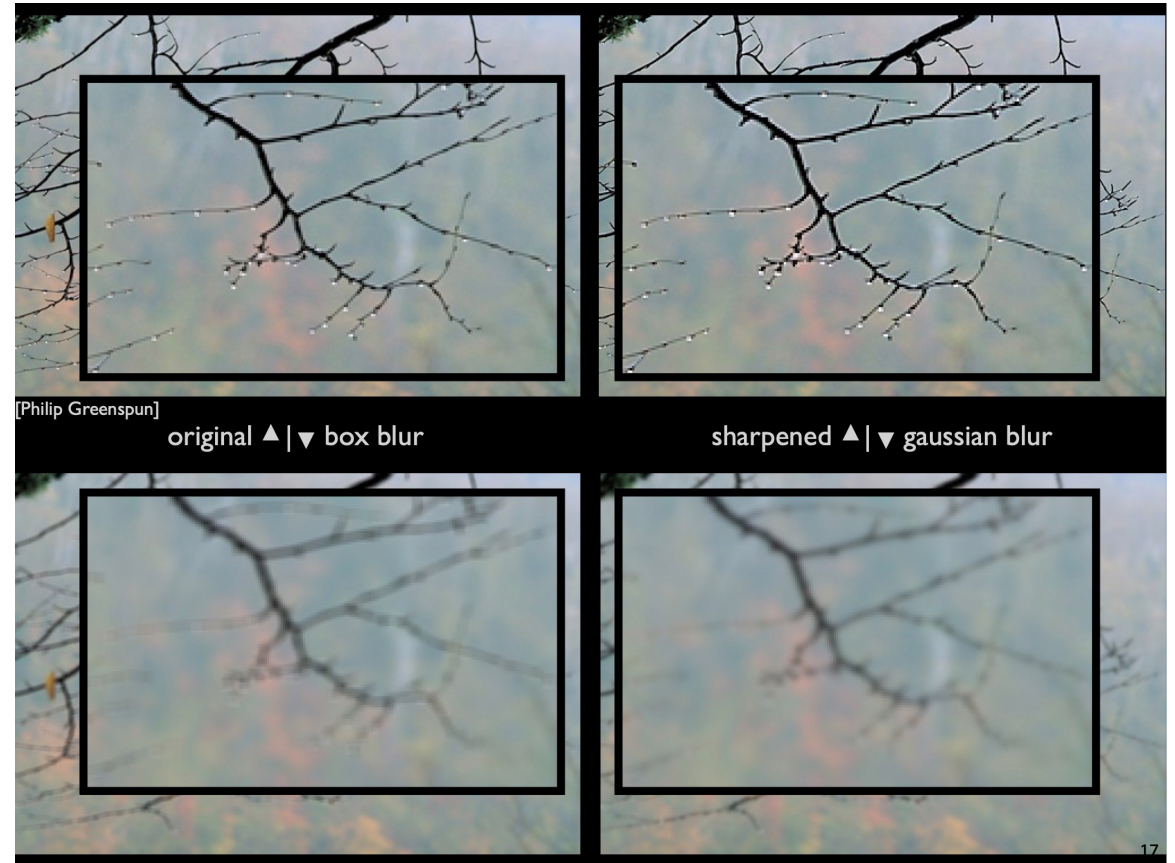
Convolution in 1D

- Basic idea: define a new function by averaging over a sliding window
- Example in one dimension: smoothing



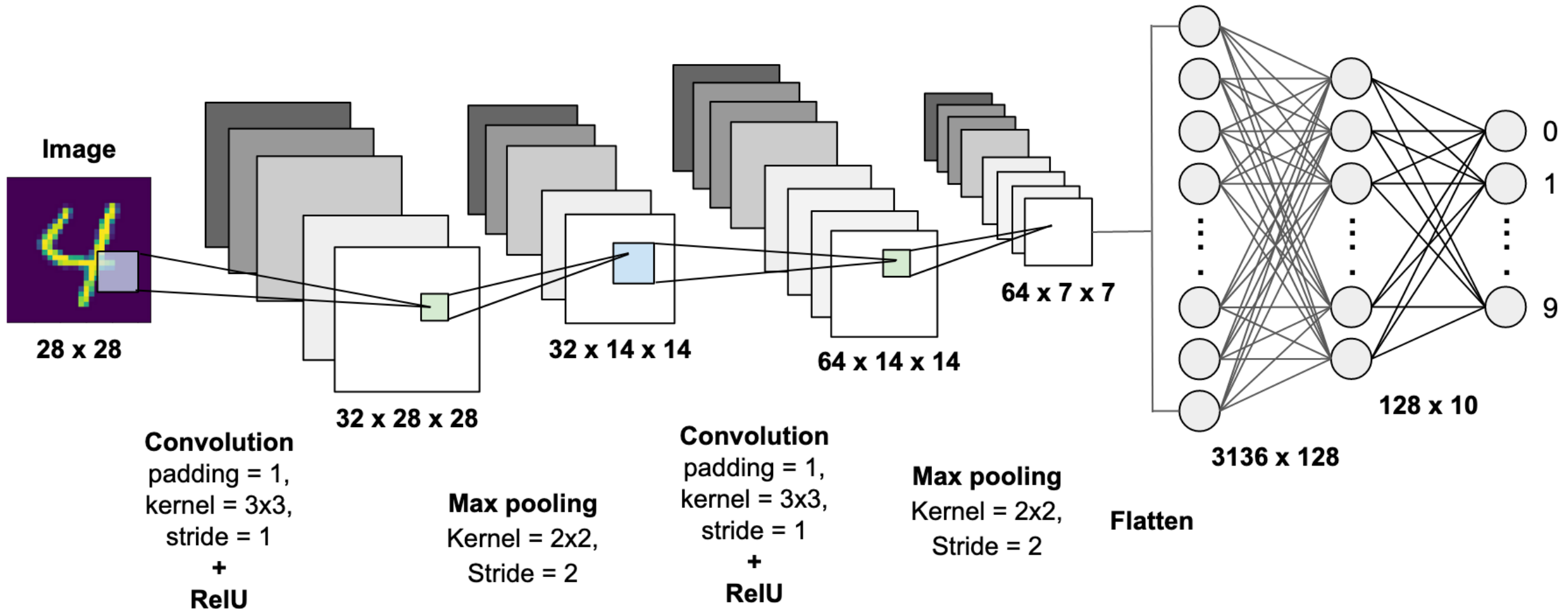
Convolution in 2D

- Filters in two dimensions: same idea but apply over a square patch of inputs (often 3x3 or 5x5)
- Applications:
 - Blurring
 - Sharpening
 - Feature detection
 - ...

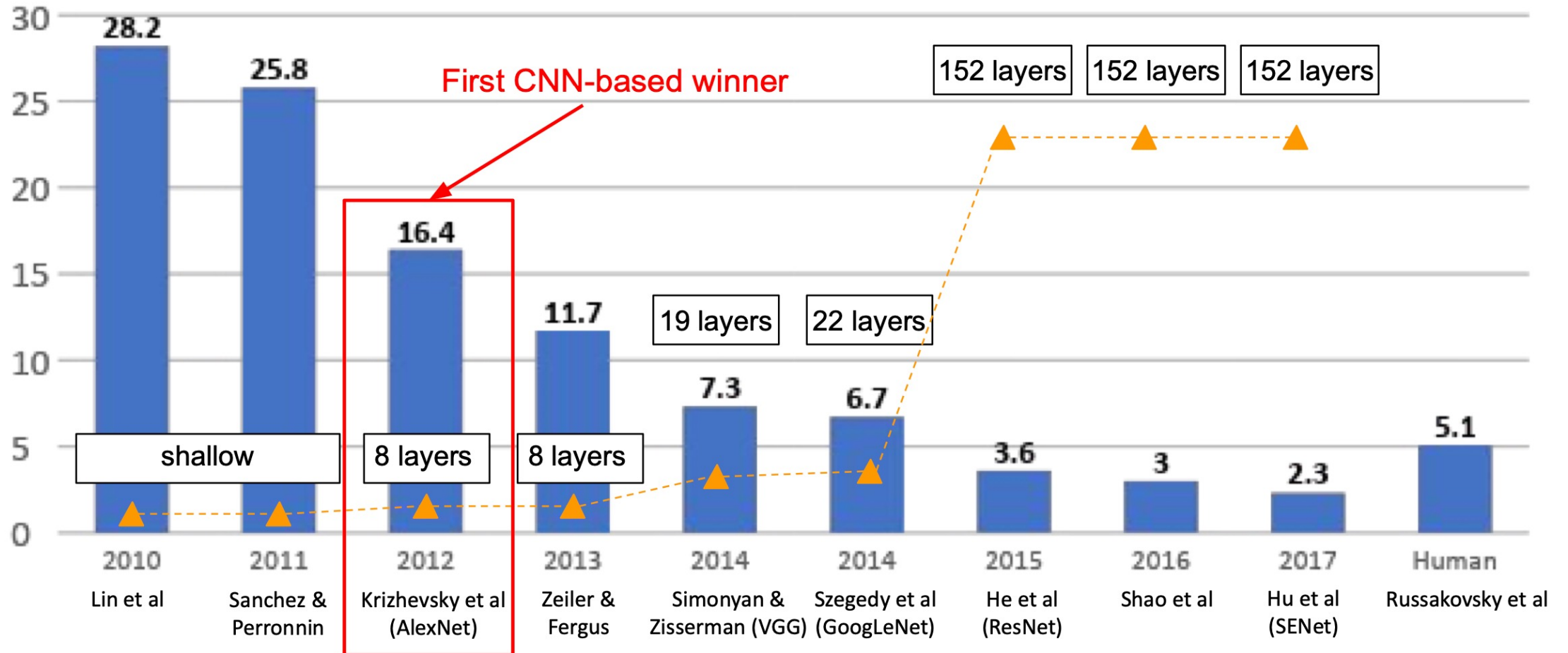


Convolutional Neural Networks

- Key idea: learn the filter weights via backprop



Benchmarking on ImageNet



ResNet (He, et al. 2015)

- Key idea:

- Want deeper networks with more parameters, but training signal becomes weak
- Add “skip” connections between layers so that there are shorter paths between early parameters and the final loss function

- ResNet:

- 152-layer model for ImageNet
- Massive improvement over all previous CNN-based classification models circa 2015

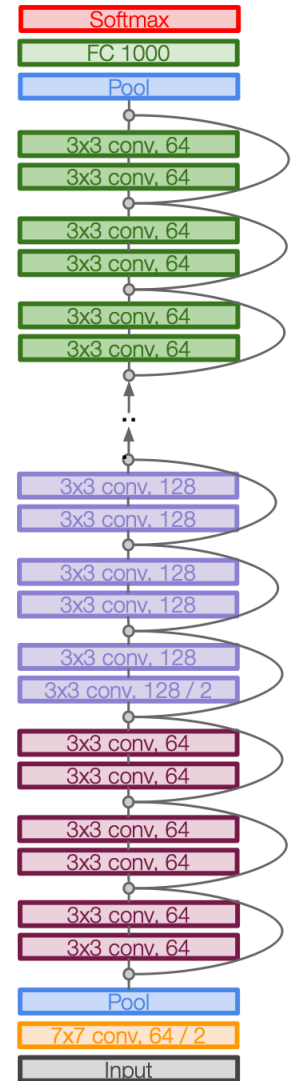
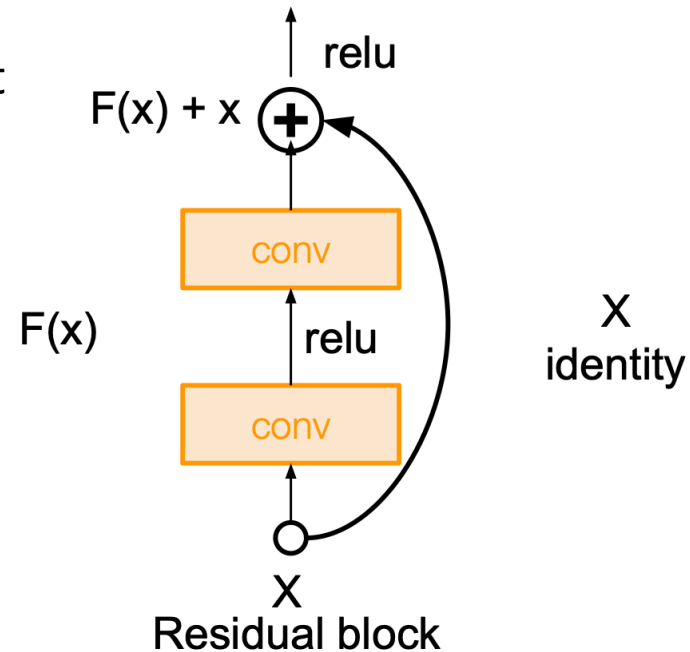


Image Classification



cat
dog
horse
person
airplane
house
...

Image Captioning



a cat standing on a desk

Image Captioning with RNNs

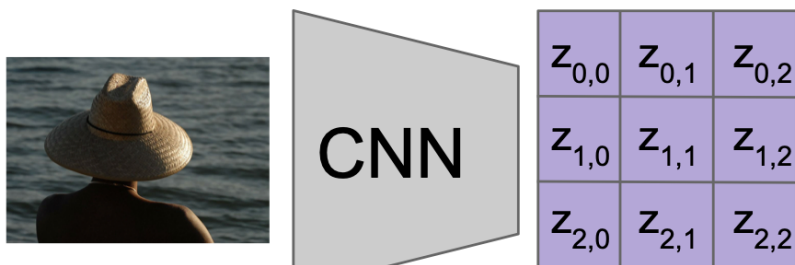
Input: Image I

Output: Sequence $\mathbf{y} = y_1, y_2, \dots, y_T$

Encoder: $h_0 = f_w(\mathbf{z})$

where \mathbf{z} is spatial CNN features

$f_w(\cdot)$ is an MLP



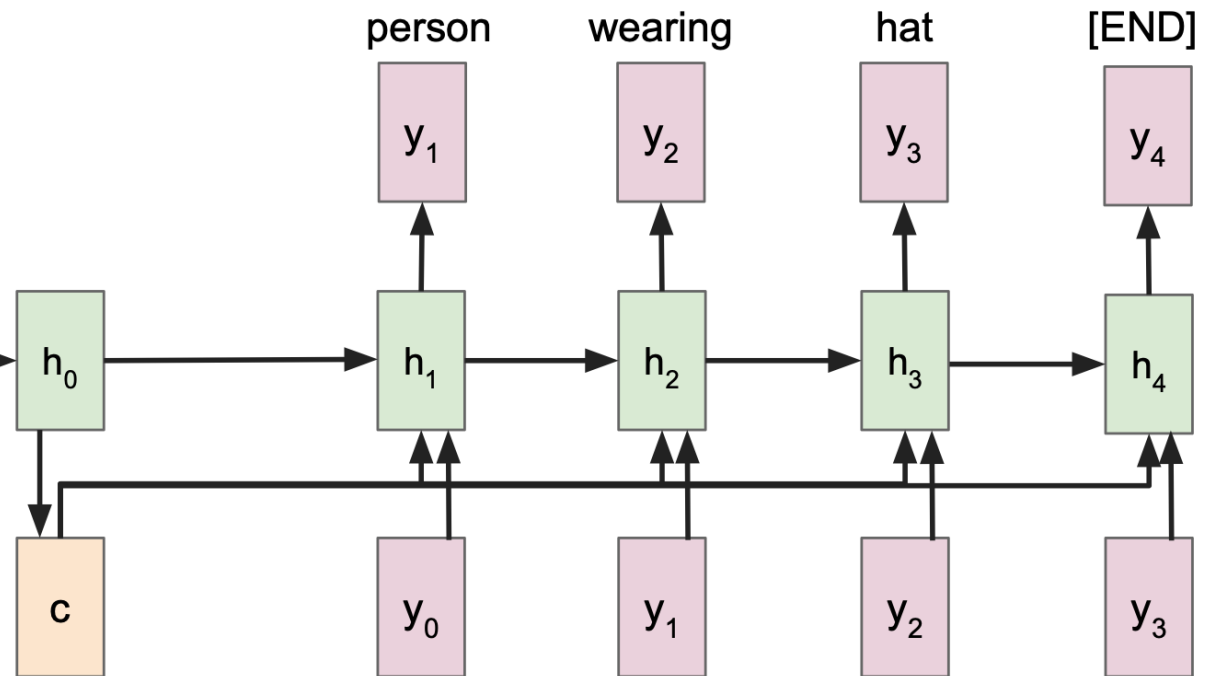
Extract spatial features from a pretrained CNN

Features:
 $H \times W \times D$

MLP

Decoder: $y_t = g_v(y_{t-1}, h_{t-1}, c)$

where context vector c is often $c = h_0$



[START]

person

wearing

hat

Image Captioning with RNNs + Attention

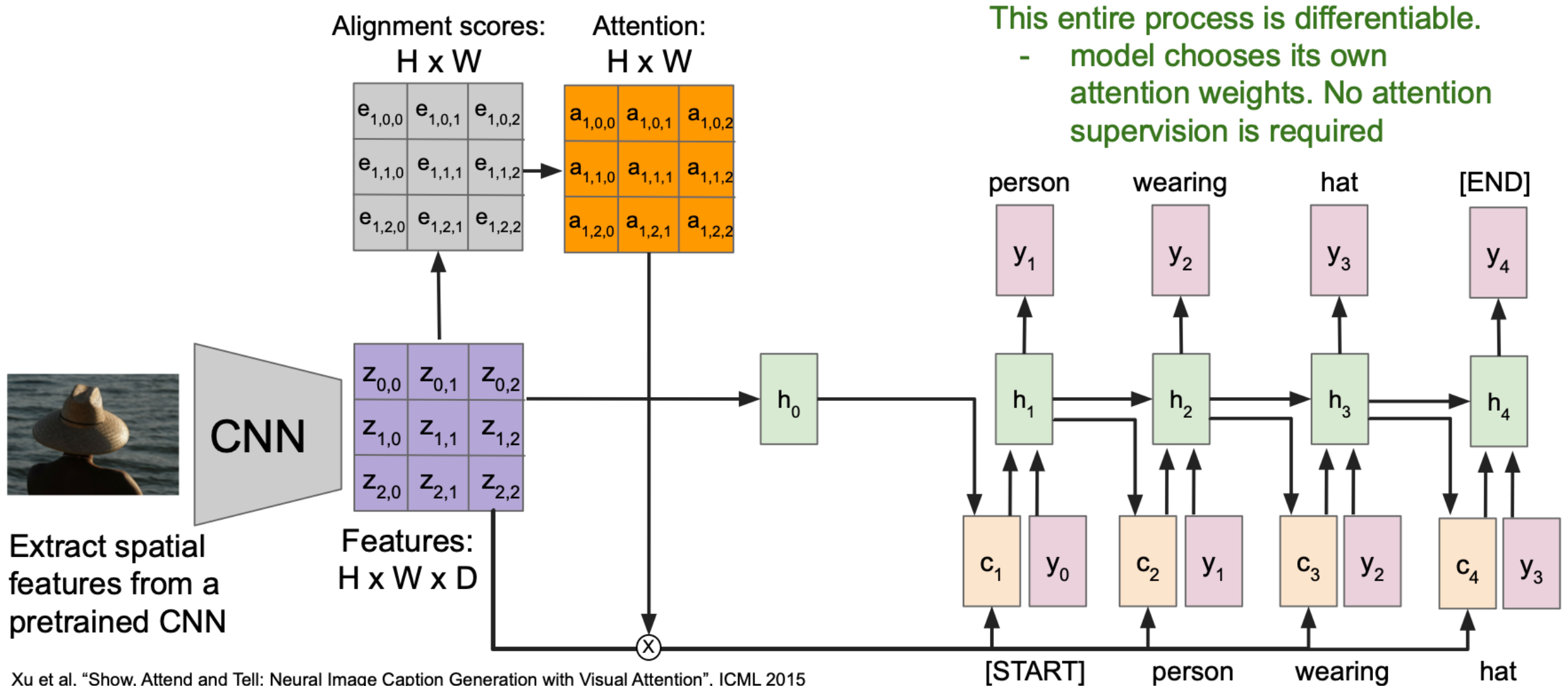


Image Captioning with Transformers

Input: Image I

Output: Sequence $\mathbf{y} = y_1, y_2, \dots, y_T$

Encoder: $\mathbf{c} = T_w(\mathbf{z})$

where \mathbf{z} is spatial CNN features

$T_w(\cdot)$ is the transformer encoder

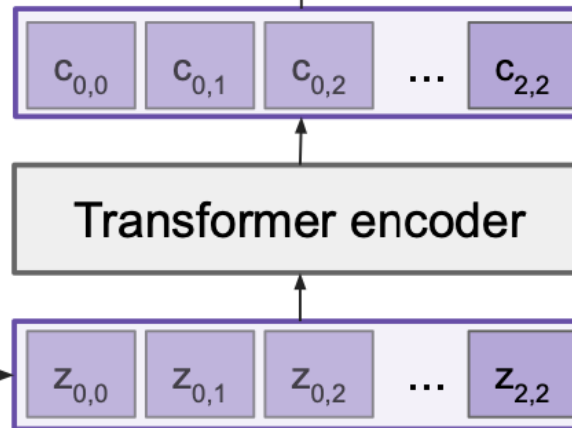


CNN

$z_{0,0}$	$z_{0,1}$	$z_{0,2}$
$z_{1,0}$	$z_{1,1}$	$z_{1,2}$
$z_{2,0}$	$z_{2,1}$	$z_{2,2}$

Features:
 $H \times W \times D$

Extract spatial features from a pretrained CNN



Decoder: $y_t = T_D(\mathbf{y}_{0:t-1}, \mathbf{c})$

where $T_D(\cdot)$ is the transformer decoder

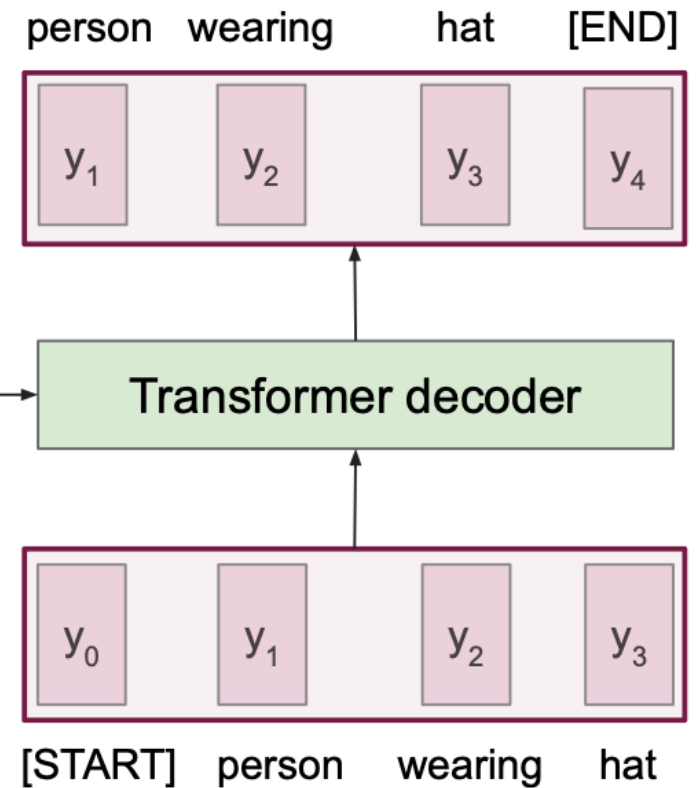
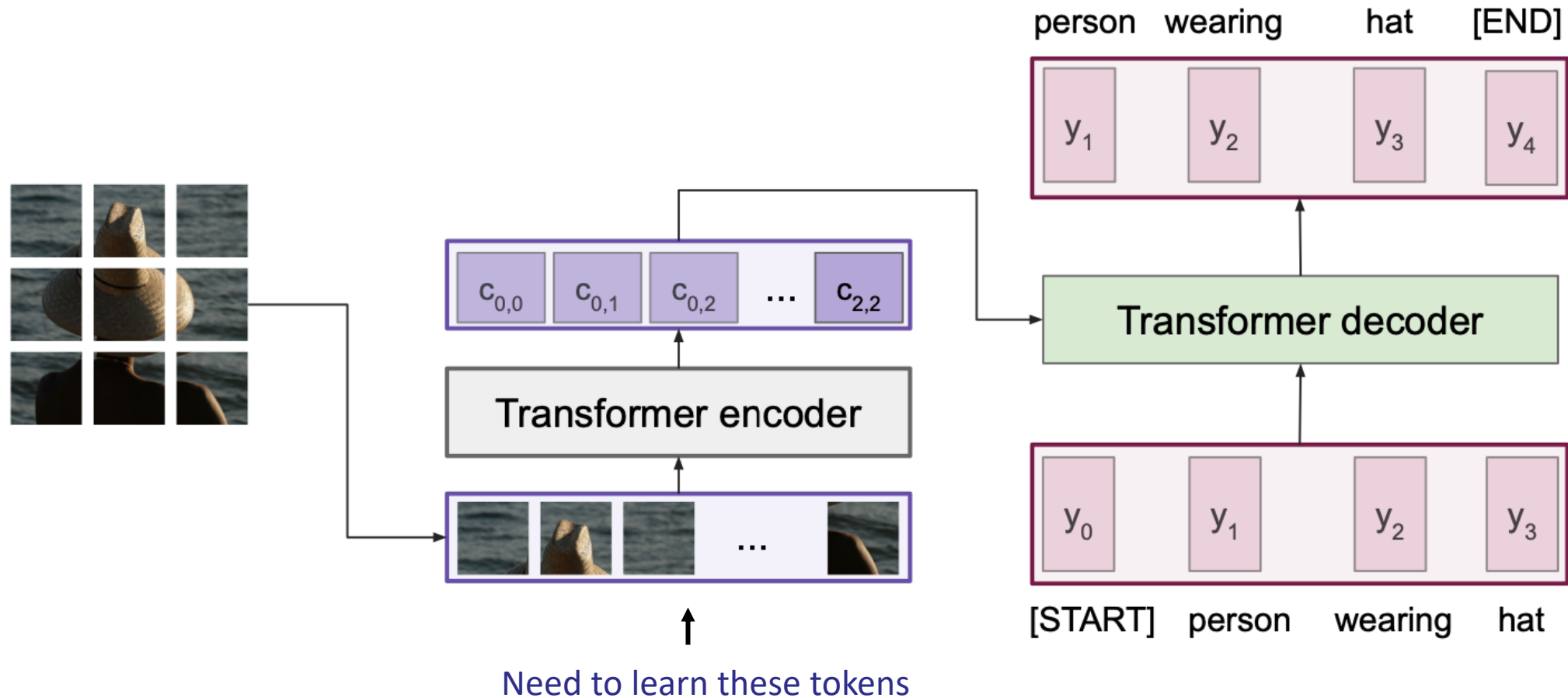
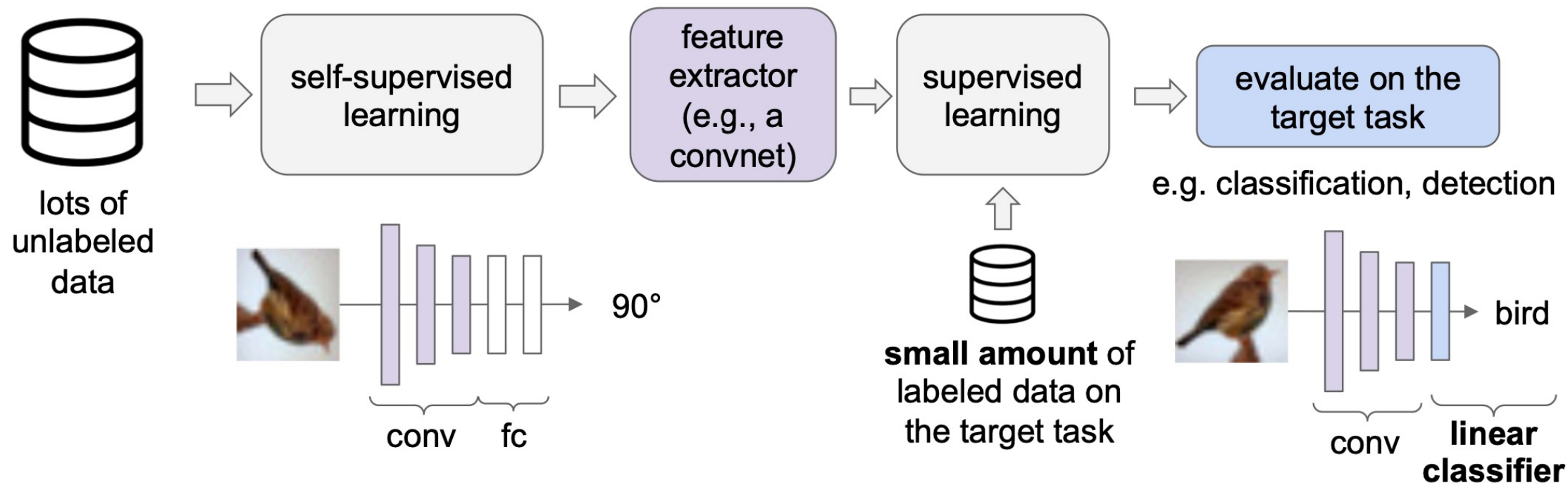


Image Captioning with Vision Transformers



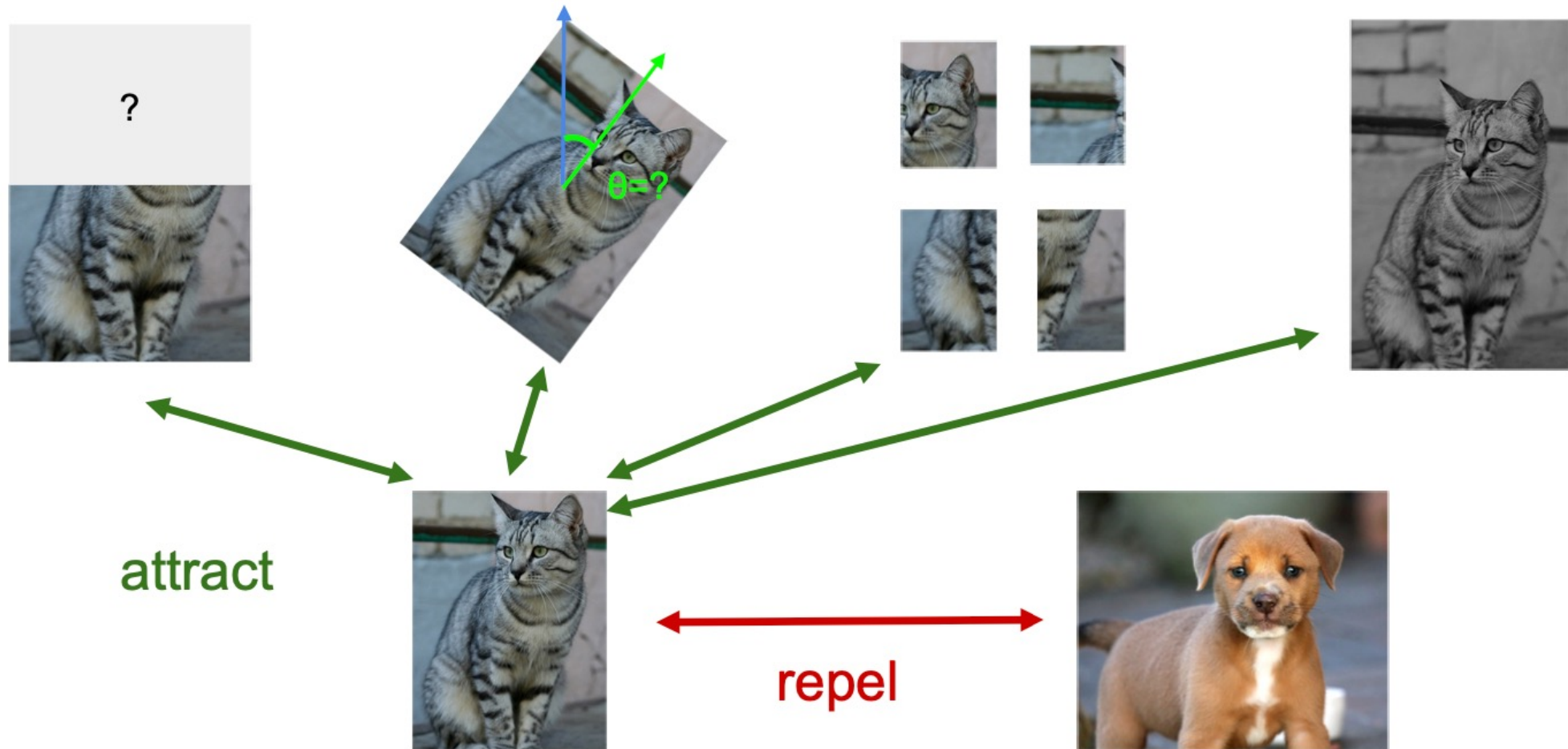
Representation Learning



1. Learn good feature extractors from self-supervised pretext tasks, e.g., predicting image rotations

2. Attach a shallow network on the feature extractor; train the shallow network on the target task with small amount of labeled data

Contrastive Learning



Representation Learning: SimCLR



(a) Original



(b) Crop and resize



(c) Crop, resize (and flip)



(d) Color distort. (drop)



(e) Color distort. (jitter)



(f) Rotate $\{90^\circ, 180^\circ, 270^\circ\}$



(g) Cutout



(h) Gaussian noise



(i) Gaussian blur



(j) Sobel filtering

Representation Learning: SimCLR



(a) Original



(b) Crop and resize



(c) Crop, resize (and flip)



(d) Color distort. (drop)



(e) Color distort. (jitter)



(f) Rotate $\{90^\circ, 180^\circ, 270^\circ\}$



(g) Cutout



(h) Gaussian noise

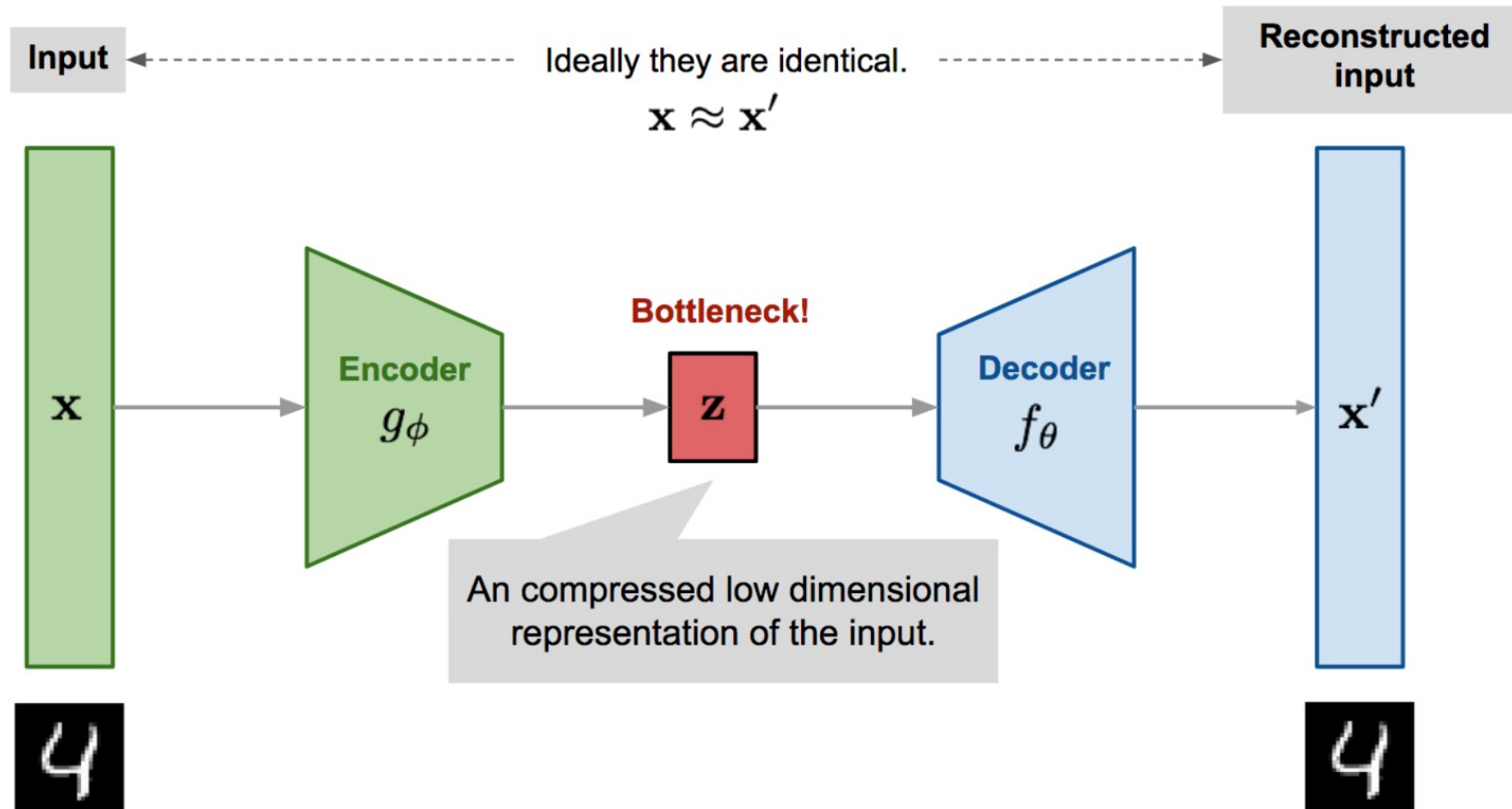


(i) Gaussian blur

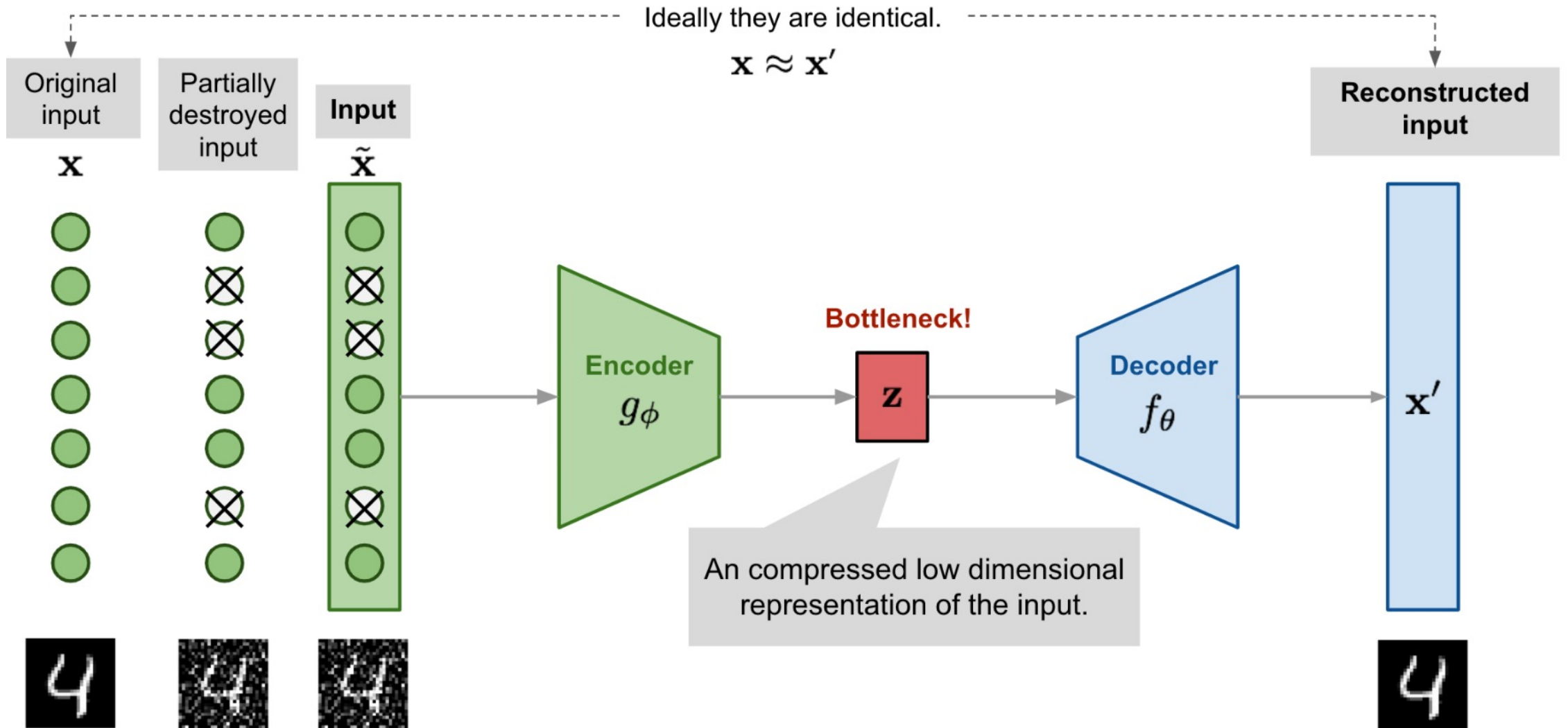


(j) Sobel filtering

Autoencoders

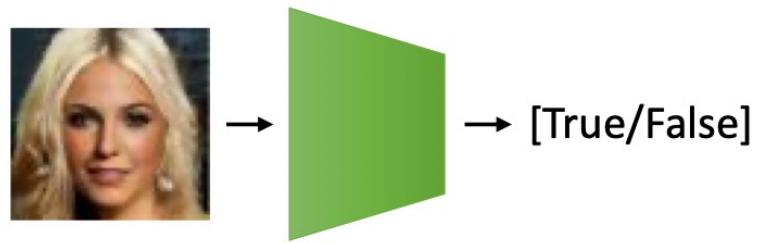


Denoising Autoencoder



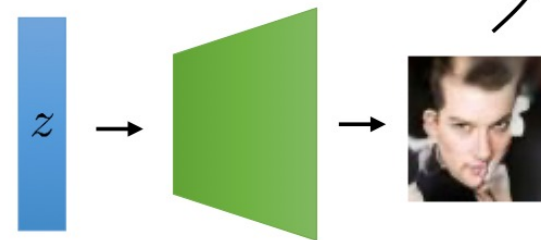
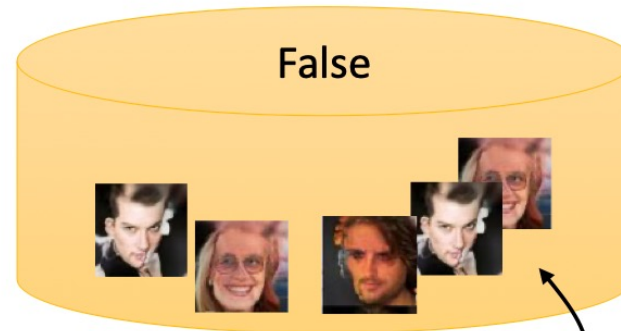
Generative Adversarial Networks

Idea: train a **network** to guess which images are real and which are fake!

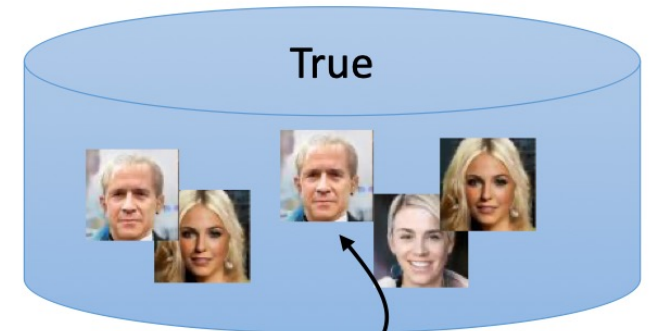


“is this a **real** image”

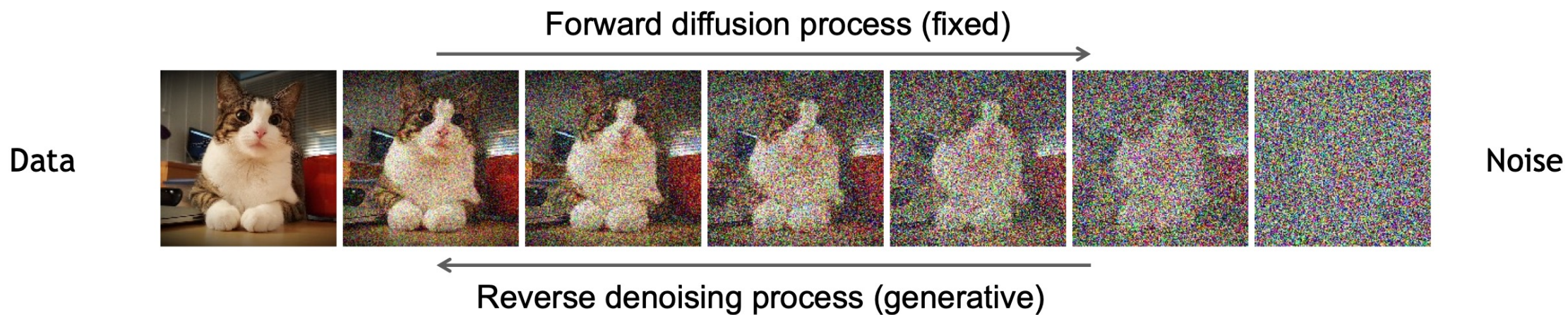
This model can then serve as a loss function for the generator!



$p(z)$ $p_{\theta}(x|z)$

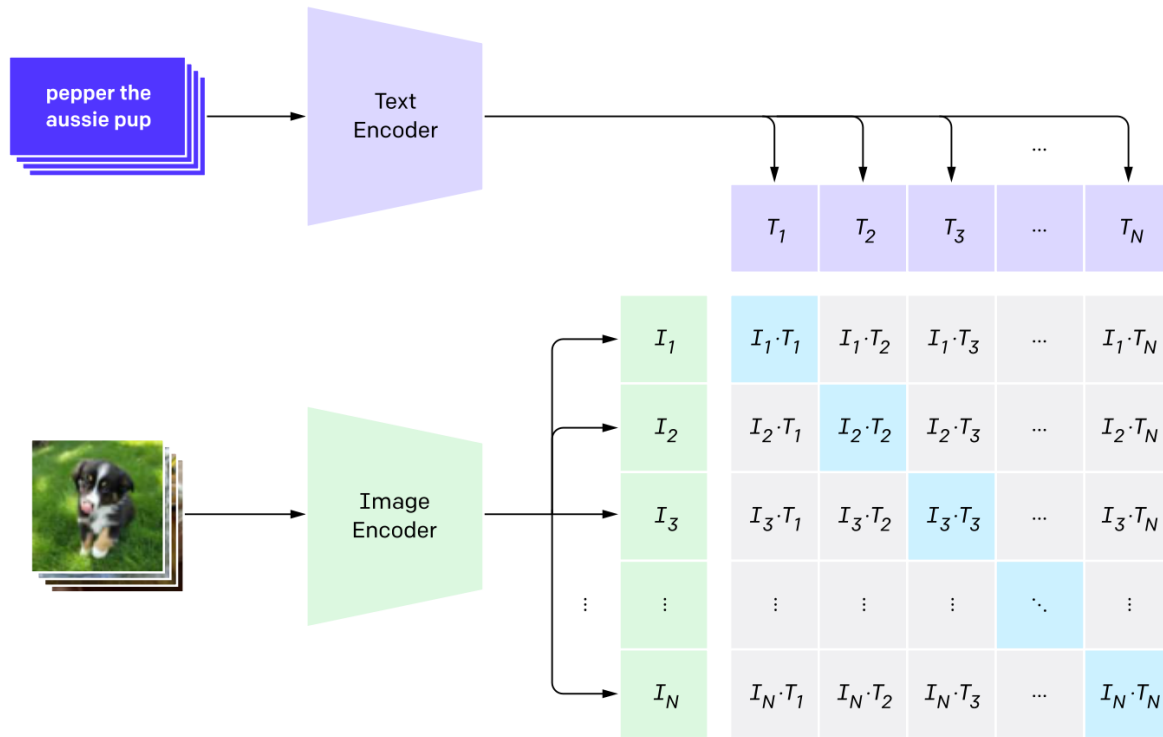


Diffusion Models



CLIP and DALL-E

1. Contrastive pre-training



2. Create dataset classifier from label text

