



香港科技大學
THE HONG KONG
UNIVERSITY OF SCIENCE
AND TECHNOLOGY

COMP 4901B

Large Language Models

Language Models

Junxian He

Sep 10, 2025

Discriminative vs. Generative Learning

Discriminative vs. Generative Learning



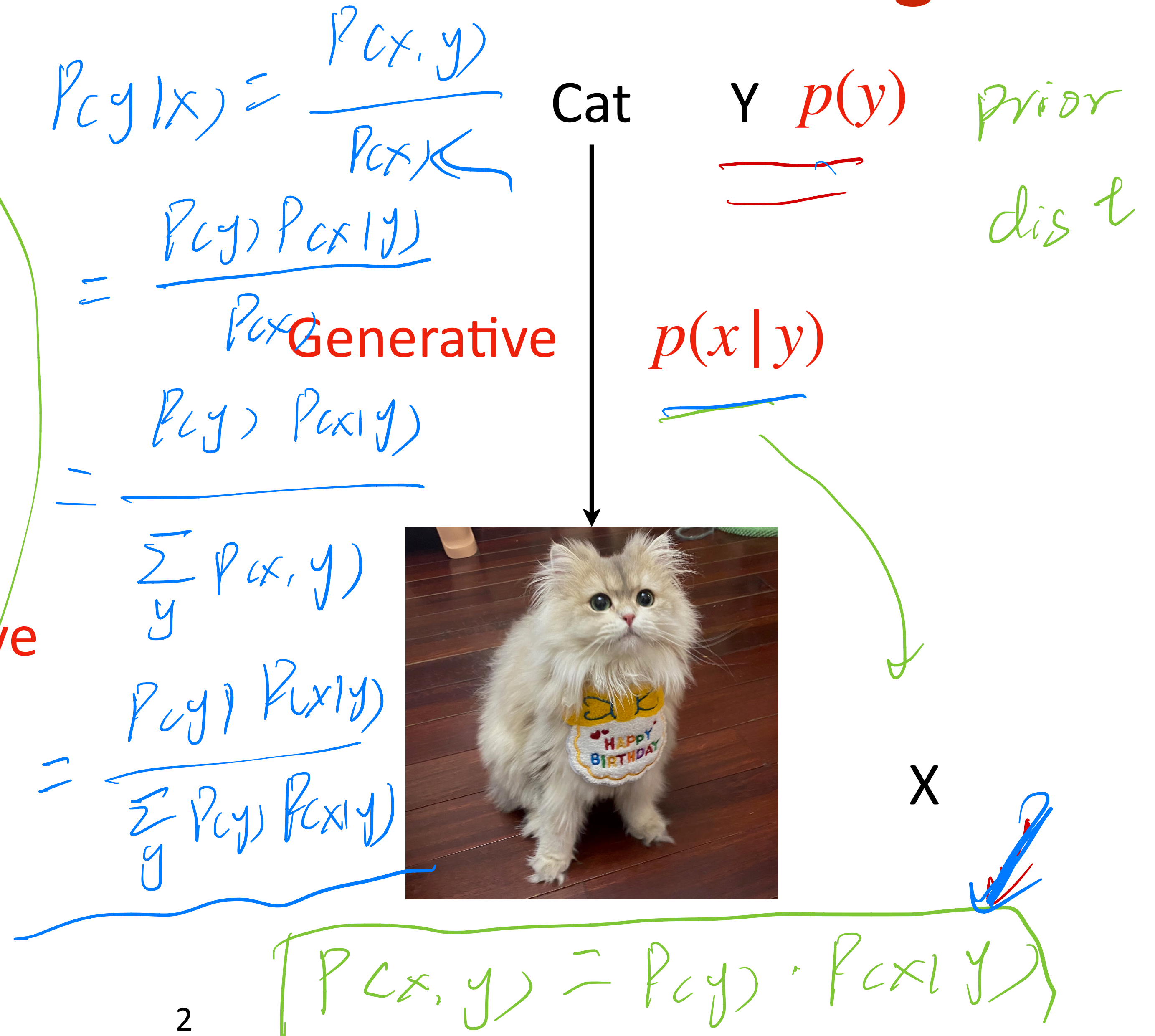
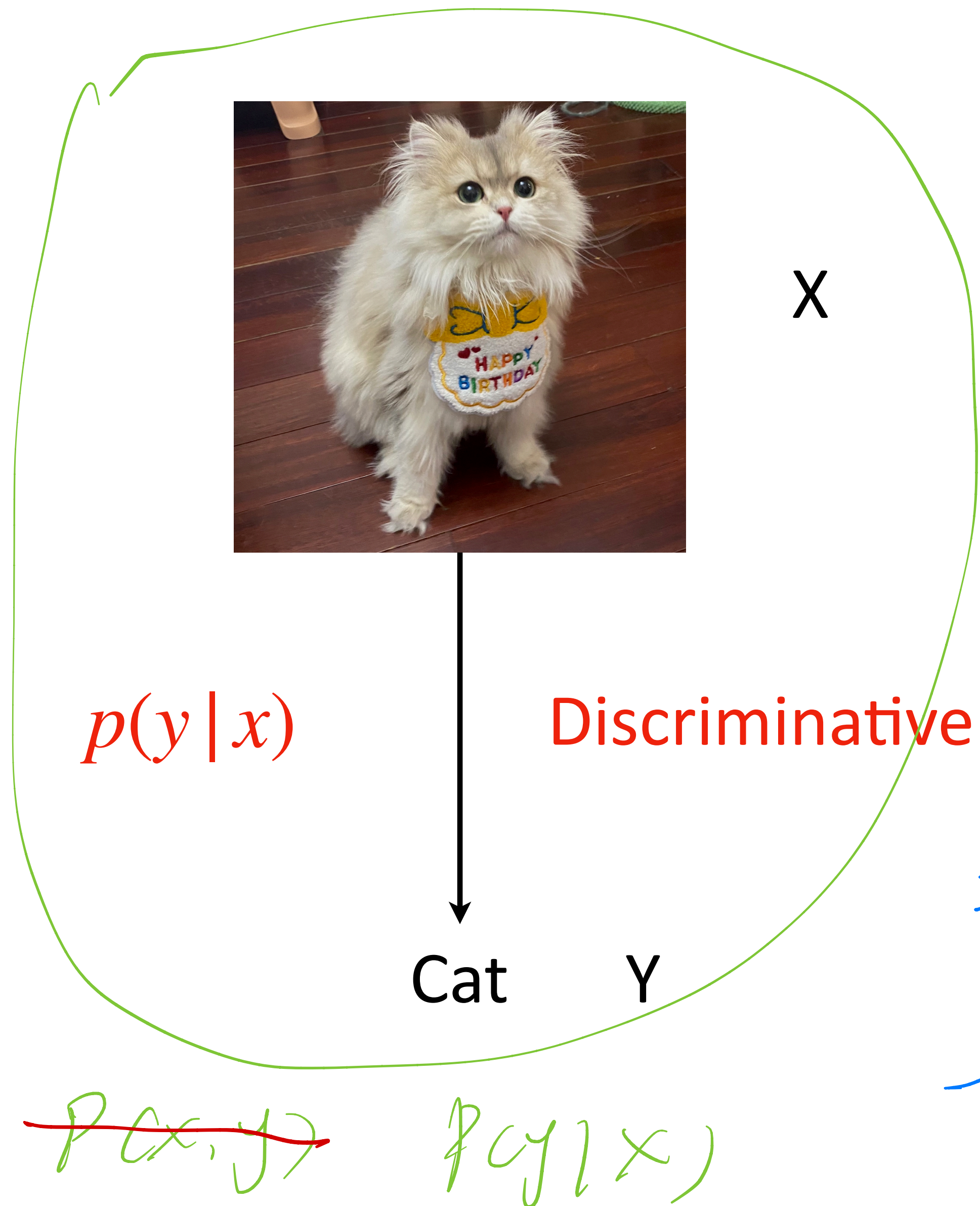
X

$p(y | x)$

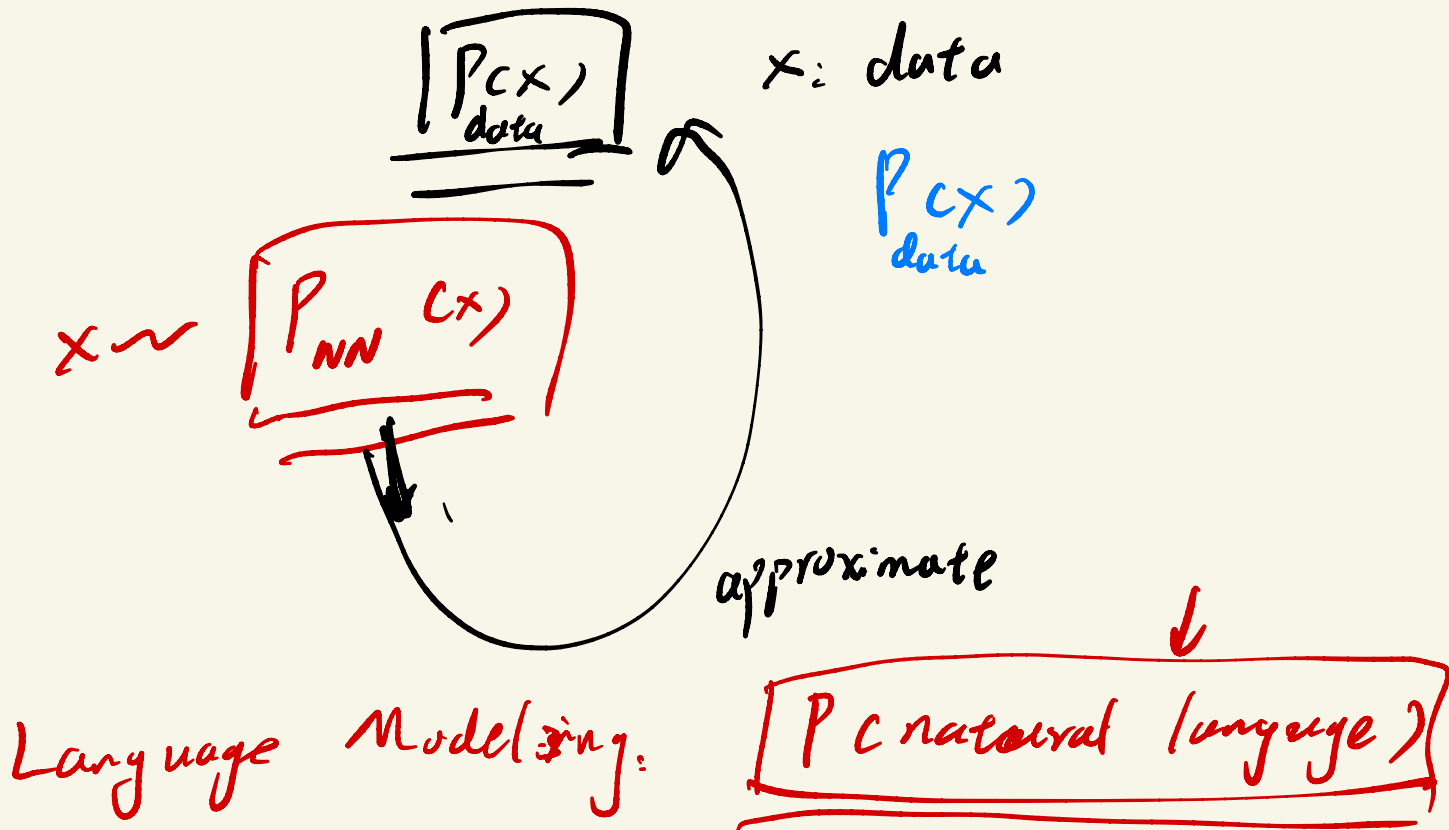
Discriminative

Cat Y

Discriminative vs. Generative Learning



Generative Modeling



Generative:

If you can create sth. ^{generative} Richard Feynman
you must understand it.

discriminative

If you don't understand sth

the mouse ate
ran,

You can never generate it

Probability of Sequences

Probability of multiple random variables:

joint probability

$$\underline{p(x_1, x_2, \dots, x_I)} = \prod_{i=1}^I p(x_i | x_{1:i-1})$$

x_i : word

goal

p

$P(x_1) P(x_2 | x_1) P(x_3 | x_1, x_2) \dots$
 $\downarrow \quad \downarrow$
 $P(x_2 | x_1, x_2 \dots x_3 \dots)$

Probability of Sequences

Probability of multiple random variables:

$$p(x_1, x_2, \dots, x_I) = \prod_{i=1}^I p(x_i | x_{1:i-1})$$

Probability of language:

Probability of Sequences

Probability of multiple random variables:

$$p(x_1, x_2, \dots, x_I) = \prod_{i=1}^I p(x_i | x_{1:i-1})$$

Probability of language:

$p(\text{the, mouse, ate, the, cheese}) =$

- $p(\text{the})$
- $p(\text{mouse} | \text{the})$
- $p(\text{ate} | \text{the, mouse})$
- $p(\text{the} | \text{the, mouse, ate})$
- $p(\text{cheese} | \text{the, mouse, ate, the}).$

the mouse | ate
ran

Probability of Sequences

Probability of multiple random variables:

$$p(x_1, x_2, \dots, x_I) = \prod_{i=1}^I p(x_i | x_{1:i-1})$$

Probability of language:

$p(\text{the, mouse, ate, the, cheese}) = p(\text{the})$

$p(\text{mouse} | \text{the})$

$p(\text{ate} | \text{the, mouse})$

$p(\text{the} | \text{the, mouse, ate})$

$p(\text{cheese} | \text{the, mouse, ate, the}).$

mathematically

Autoregressive language models

$P_C(\text{the mouse ate the cheese})$

$= P(\underline{\text{ate}}) P(\underline{\text{mouse}} | \text{ate})$

$P(\text{cheese} | \text{mouse } \phi, \text{ate})$

$P_C(\text{the} | \text{cheese, mouse, ate})$

$P_C(\text{the} | \text{cheese, mouse, ate, the})$

Autoregressive Language Models

$$p(\text{the, mouse, ate, the, cheese}) = p(\text{the}) \\ p(\text{mouse} \mid \text{the}) \\ p(\text{ate} \mid \text{the, mouse}) \\ p(\text{the} \mid \text{the, mouse, ate}) \\ p(\text{cheese} \mid \text{the, mouse, ate, the}).$$

next word prediction

distribution of next word

$$p(x_1, x_2, \dots, x_I) = \prod_{i=1}^I p(x_i \mid x_{1:i-1})$$

classification

$P(\text{word} \mid \text{context})$

Next Word

Context

Autoregressive Language Models

$$\begin{aligned} p(\text{the, mouse, ate, the, cheese}) &= p(\text{the}) \\ &\quad p(\text{mouse} \mid \text{the}) \\ &\quad p(\text{ate} \mid \text{the, mouse}) \\ &\quad p(\text{the} \mid \text{the, mouse, ate}) \\ &\quad p(\text{cheese} \mid \text{the, mouse, ate, the}). \end{aligned}$$

$$p(x_1, x_2, \dots, x_I) = \prod_{i=1}^I p(x_i \mid x_{1:i-1})$$

Learning a language model is to learn these conditional probabilities, for any language sequence

Autoregressive Language Models

$$\begin{aligned} p(\text{the, mouse, ate, the, cheese}) &= p(\text{the}) \\ &\quad p(\text{mouse} \mid \text{the}) \\ &\quad p(\text{ate} \mid \text{the, mouse}) \\ &\quad p(\text{the} \mid \text{the, mouse, ate}) \\ &\quad p(\text{cheese} \mid \text{the, mouse, ate, the}). \end{aligned}$$

$$p(x_1, x_2, \dots, x_I) = \prod_{i=1}^I p(x_i \mid x_{1:i-1})$$

parameters

Given a dataset, how to find these probabilities?

loss? $\arg \max_{\theta} \log P_{\theta}(x)$

Autoregressive Language Models

$$\begin{aligned} p(\text{the, mouse, ate, the, cheese}) &= p(\text{the}) \\ &\quad p(\text{mouse} \mid \text{the}) \\ &\quad p(\text{ate} \mid \text{the, mouse}) \\ &\quad p(\text{the} \mid \text{the, mouse, ate}) \\ &\quad p(\text{cheese} \mid \text{the, mouse, ate, the}). \end{aligned}$$

$$p(x_1, x_2, \dots, x_I) = \prod_{i=1}^I p(x_i \mid x_{1:i-1})$$

parameter to

Given a dataset, how to find these probabilities?

Maximum Likelihood Estimation

Count-based Language Models

Count the frequency and divide

$$p(x_i | x_{1:i-1}) = \frac{c(x_{1:i})}{c(x_{1:i-1})}$$

$p(\text{ate} | \text{the mouse})$
↑
C C the mouse ate
C C the mouse)
C C the mouse)

the mouse ate

3 times

the mouse ran

2 times

3 + 2

= 5

other

the mouse x

0 times

$\frac{3}{5}$

$P(\text{ate} | \text{the mouse})$

$P(\text{not ate} | \text{the mouse}) = 0$
 $P(\text{not ran})$

max $3 \cdot \log P(\text{ate} | \text{the mouse}) + 2 \cdot \log P(\text{ran} | \text{the mouse})$

$P(\text{ate} | \text{the mouse}) + P(\text{ran} | \text{the mouse}) = 1$

$$\max 3 \cdot \log P(\text{cat} | \text{the mouse}) +$$

$$2 \cdot \log (1 - P(\text{cat} | \text{the mouse}))$$

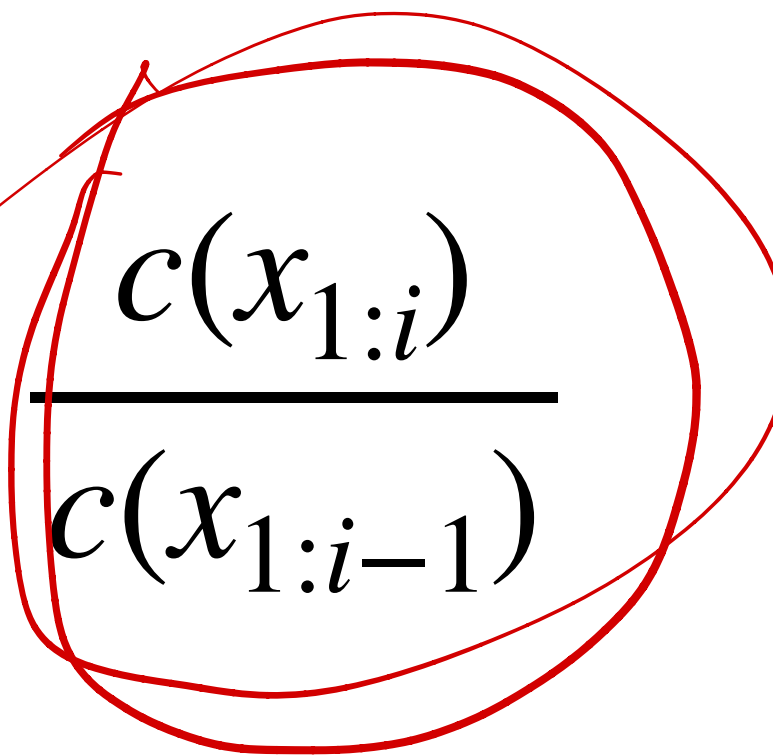
$$P(\text{cat} | \text{the mouse})$$

$$\frac{d}{dP} \left[\frac{3}{P(\text{cat} | \text{the mouse})} + \frac{-2}{1 - P(\text{cat} | \text{the mouse})} \right]$$

$$= 0 \Rightarrow P(\text{cat} | \text{the mouse}) = \frac{3}{5}$$

Count-based Language Models

Count the frequency and divide

$$p(x_i | x_{1:i-1}) = \frac{c(x_{1:i})}{c(x_{1:i-1})}$$


There are infinite number of parameters for language



Count-based Language Models

Count the frequency and divide

$$p(x_i | x_{1:i-1}) = \frac{c(x_{1:i})}{c(x_{1:i-1})}$$

There are infinite number of parameters for language

We may see long sequences only once, counting becomes meaningless

n-gram Language Models

n-gram Language Models

Next token probability only depends on the previous $n-1$ words

n-gram Language Models

Next token probability only depends on the previous n-1 words

Unigram LM:

$$p(x_1, x_2, \dots, x_I) = \prod_{i=1}^I p(x_i)$$

$p(x_i)$ finite

n-gram Language Models

Next token probability only depends on the previous n-1 words

Unigram LM:

$$p(x_1, x_2, \dots, x_I) = \prod_{i=1}^I p(x_i) \quad \text{Each token is independent}$$

n-gram Language Models

Next token probability only depends on the previous n-1 words

Unigram LM:

$$p(x_1, x_2, \dots, x_I) = \prod_{i=1}^I p(x_i)$$

Each token is independent

word

Bigram LM:

$$p(x_1, x_2, \dots, x_I) = \prod_{i=1}^I p(x_i | x_{i-1})$$

D x D

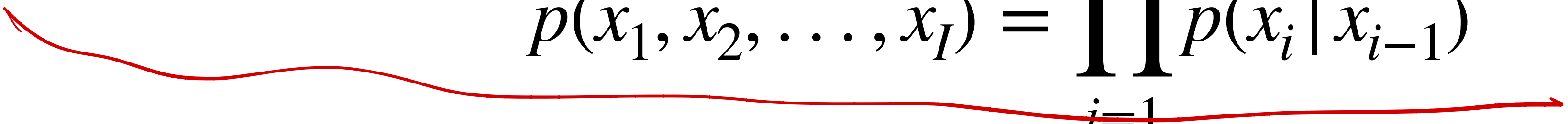
n-gram Language Models

Next token probability only depends on the previous n-1 words

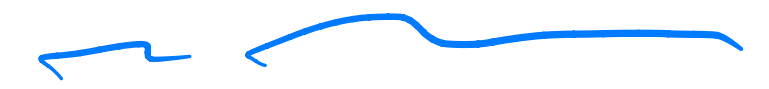
Unigram LM:

$$p(x_1, x_2, \dots, x_I) = \prod_{i=1}^I p(x_i) \quad \text{Each token is independent}$$

Bigram LM:


$$p(x_1, x_2, \dots, x_I) = \prod_{i=1}^I p(x_i | x_{i-1})$$

Generally for n-gram LM:


$$p(x_1, x_2, \dots, x_I) = \prod_{i=1}^I p(x_i | x_{i-n+1:i-1})$$

8

Parameter Estimation for n-gram LM

Count-based:

$$p(x_i | x_{i-n+1:i-1}) = \frac{c(x_{i-n+1:i})}{c(x_{i-n+1:i-1})}$$

Parameter Estimation for n-gram LM

Count-based:

$$p(x_i | x_{i-n+1:i-1}) = \frac{c(x_{i-n+1:i})}{c(x_{i-n+1:i-1})}$$

Number of parameters decreases, but flexibility decreases as well

" I am from the major computer science
Now I am year 3 in HKUST. - - -
- - - I like the sea. This semester.
I am taken classes

Parameter Estimation for n-gram LM

Count-based:

$$p(x_i | x_{i-n+1:i-1}) = \frac{c(x_{i-n+1:i})}{c(x_{i-n+1:i-1})}$$

Number of parameters decreases, but flexibility decreases as well

Traditionally, we directly compute this probability, but neural language models use neural networks to compute the probability

Neural Language Models

Neural Language Models

Neural language models are typically autoregressive *next word prediction*



Neural Language Models

Neural language models are typically autoregressive

Data: “The mouse ate the cheese .”

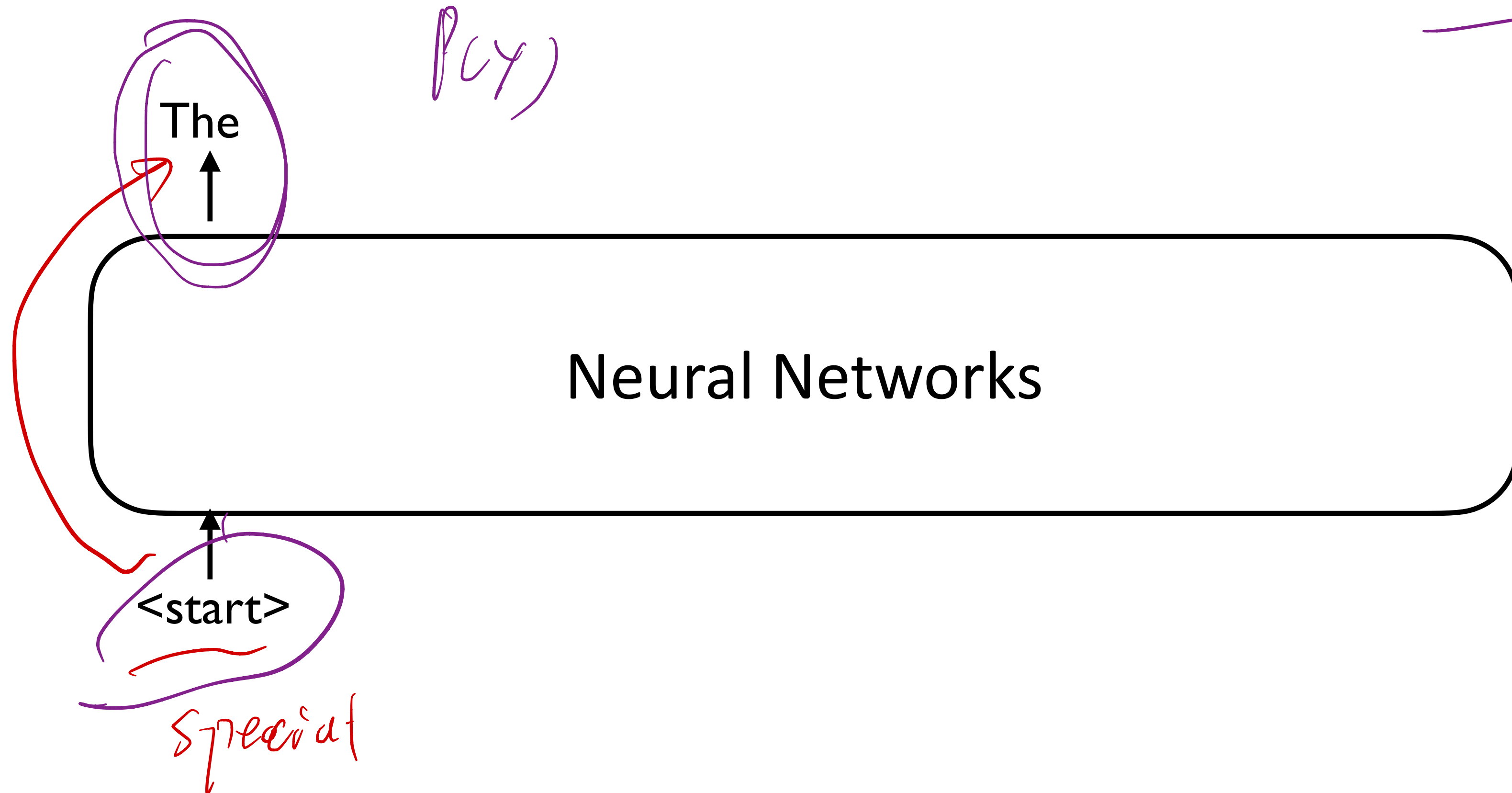
Neural Language Models

Neural language models are typically autoregressive

Data: "The mouse ate the cheese ."

$$P(\text{the}) = P(\text{The} | \underline{\text{<start>}})$$

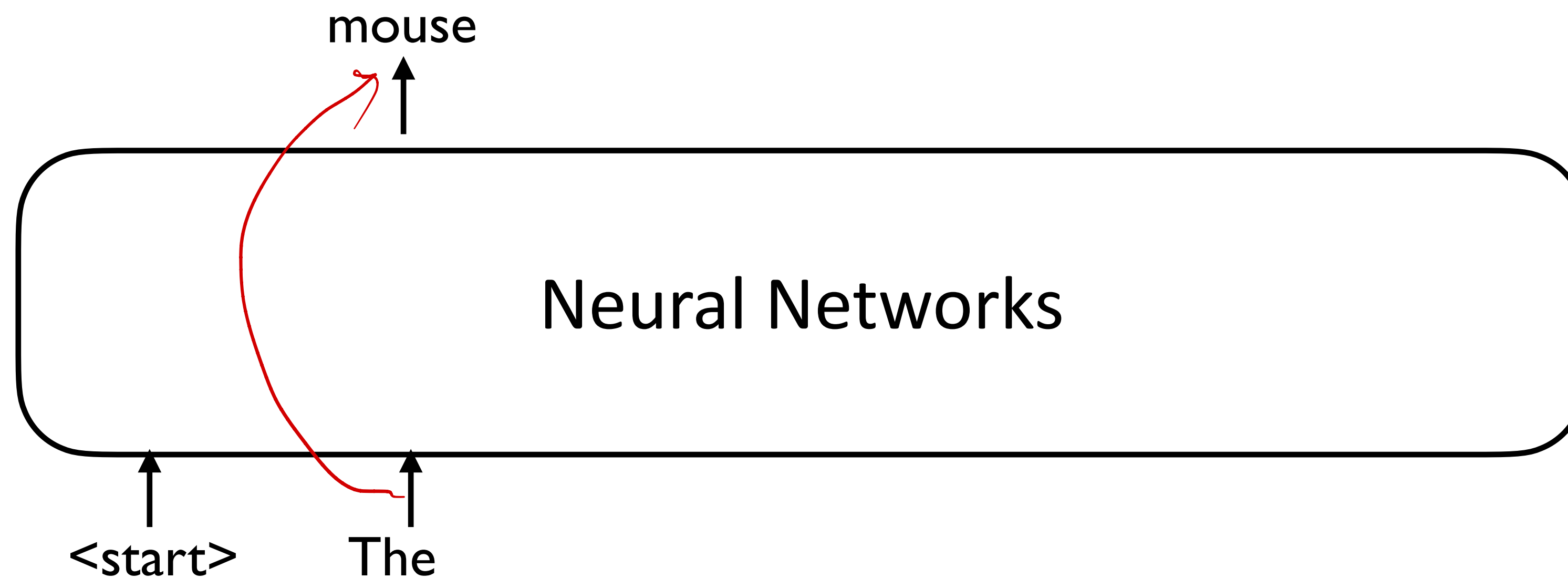
fake word



Neural Language Models

Neural language models are typically autoregressive

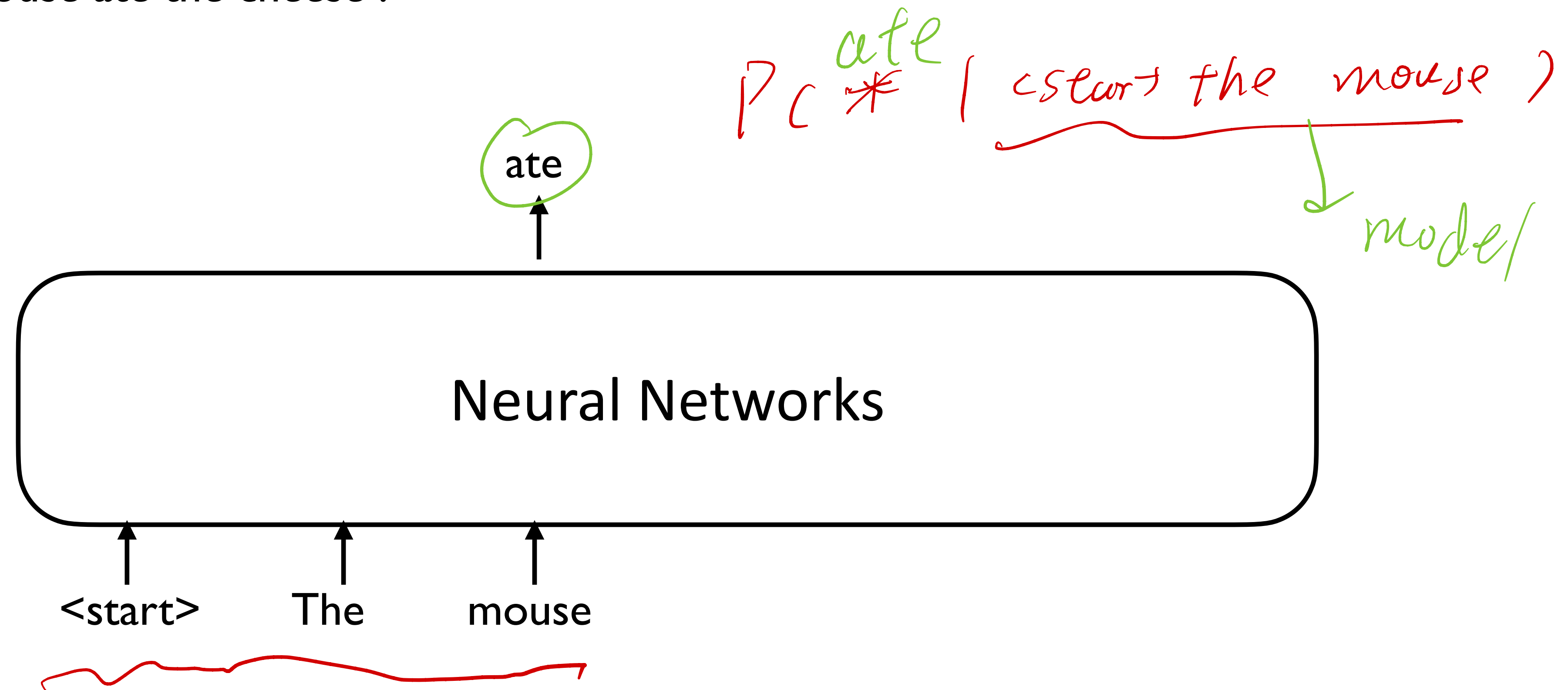
Data: “The mouse ate the cheese .”



Neural Language Models

Neural language models are typically autoregressive

Data: "The mouse ate the cheese ."



Neural Language Models

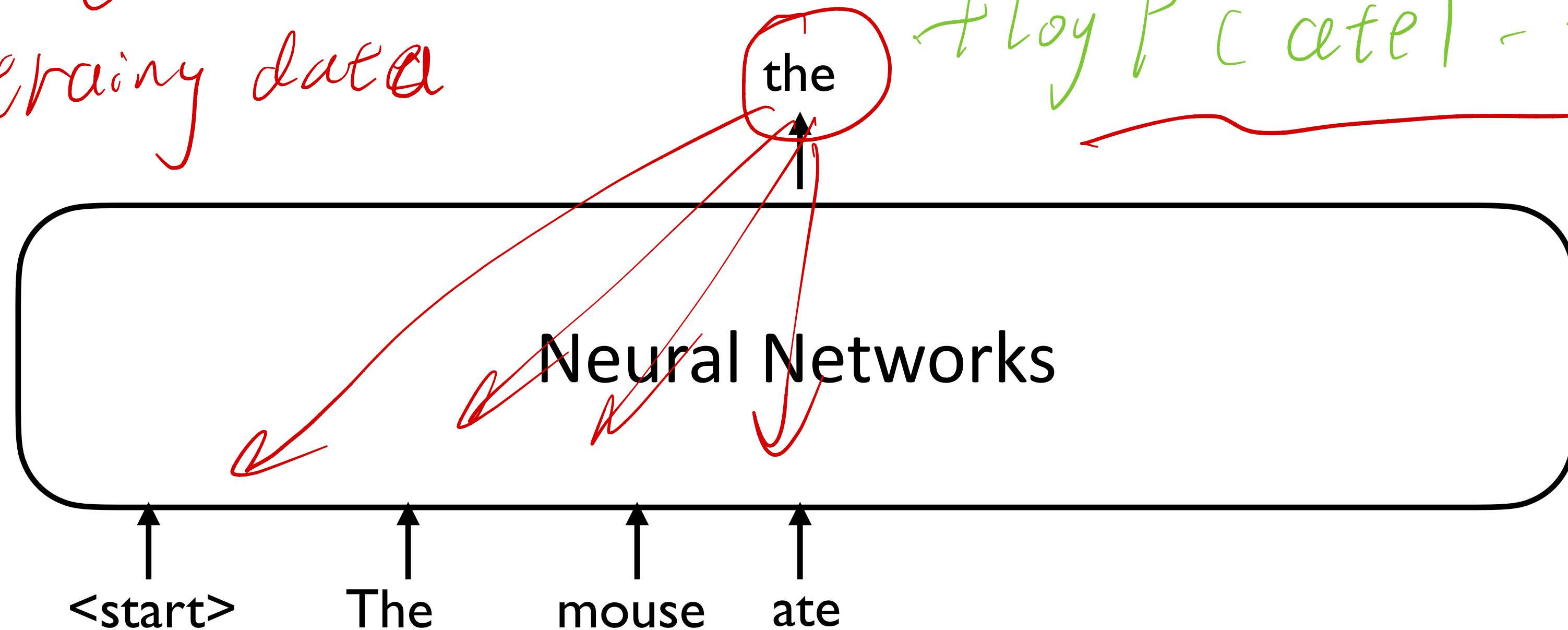
Neural language models are typically autoregressive

Data: "The mouse ate the cheese."

training data

$\log P(\text{the} | \langle \text{start} \rangle) + \log P(\text{mouse} | \text{the})$

$+ \log P(\text{ate} | \text{the, mouse})$

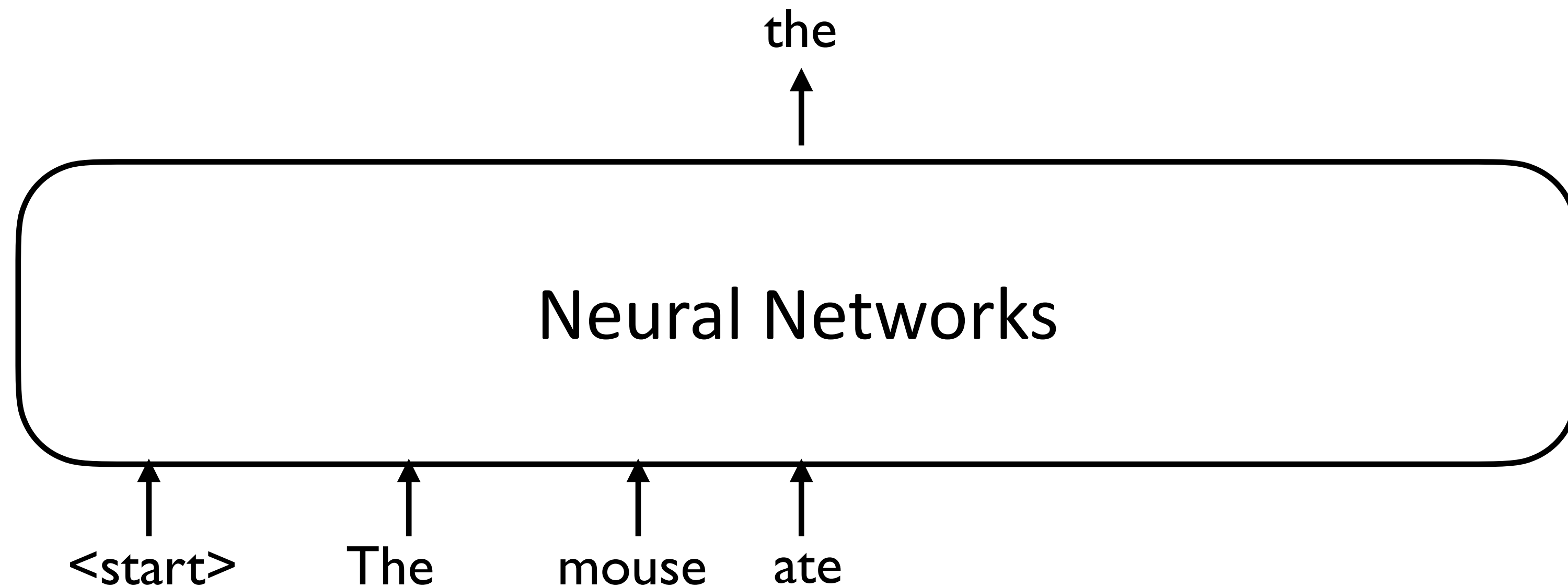


Neural Language Models

Neural language models are typically autoregressive

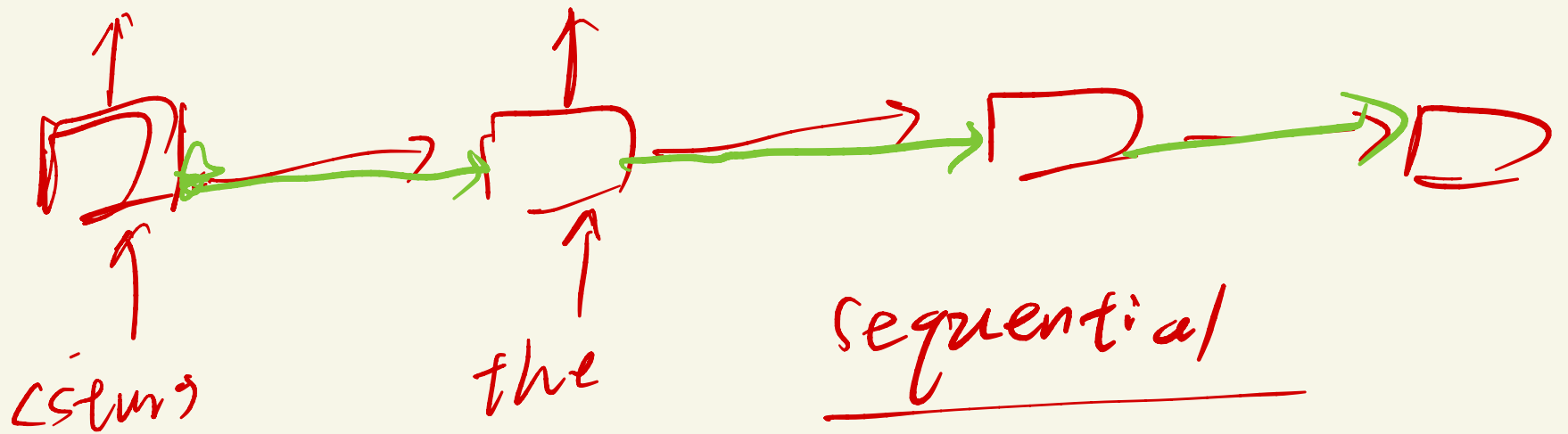
Data: "The mouse ate the cheese ."

γ



We can compute the loss on every token in parallel

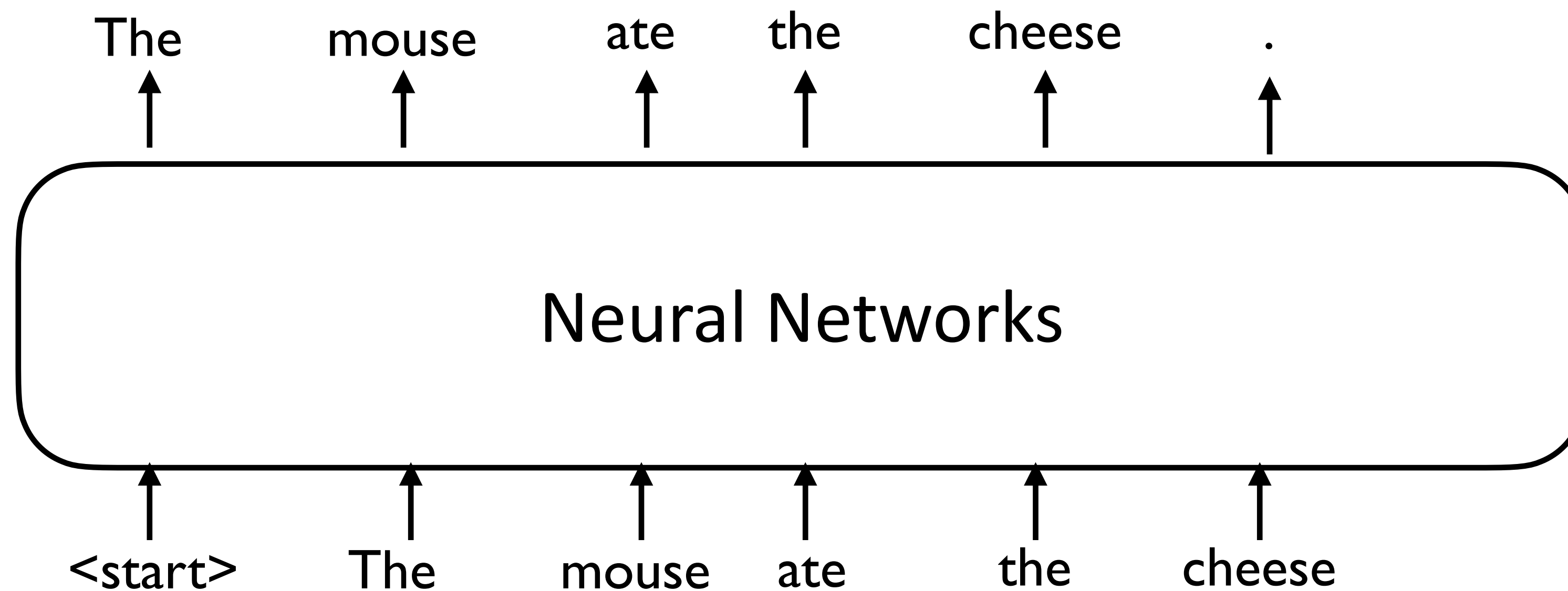
RNN



Neural Language Models

Neural language models are typically autoregressive

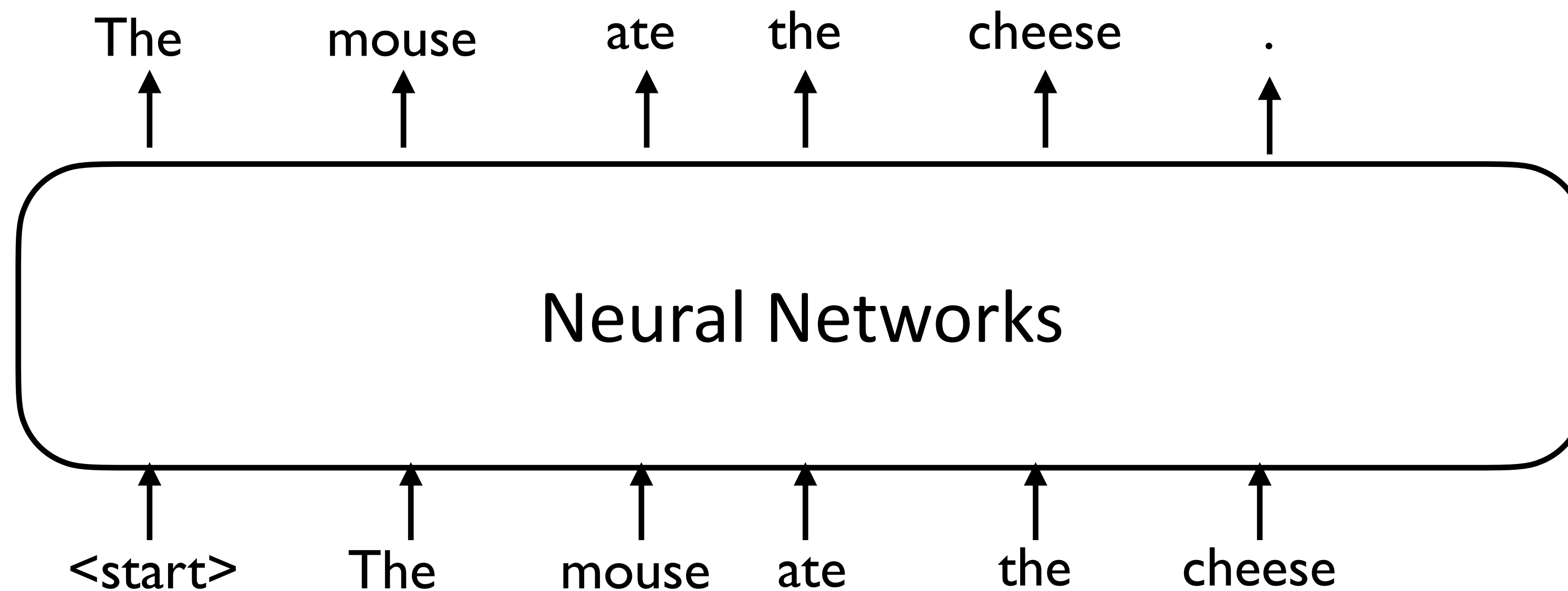
Data: “The mouse ate the cheese .”



Neural Language Models

Neural language models are typically autoregressive

Data: “The mouse ate the cheese .”



Each prediction only sees the inputs on its left

Neural Language Models

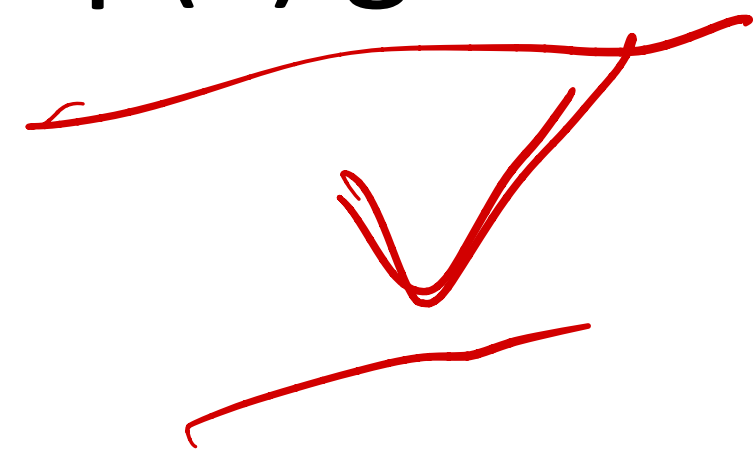
Neural Language Models

Are language models generative models? $P(x)$

Neural Language Models

Are language models generative models?

Can we compute $p(x)$ given x ? Can we sample new x ?



$\mathbb{I} p_{\text{next word} | \text{context}}$

Neural Language Models

Are language models generative models? 

Can we compute $p(x)$ given x ? Can we sample new x ? 

Neural Language Models

Are language models generative models? 

Can we compute $p(x)$ given x ? Can we sample new x ? 

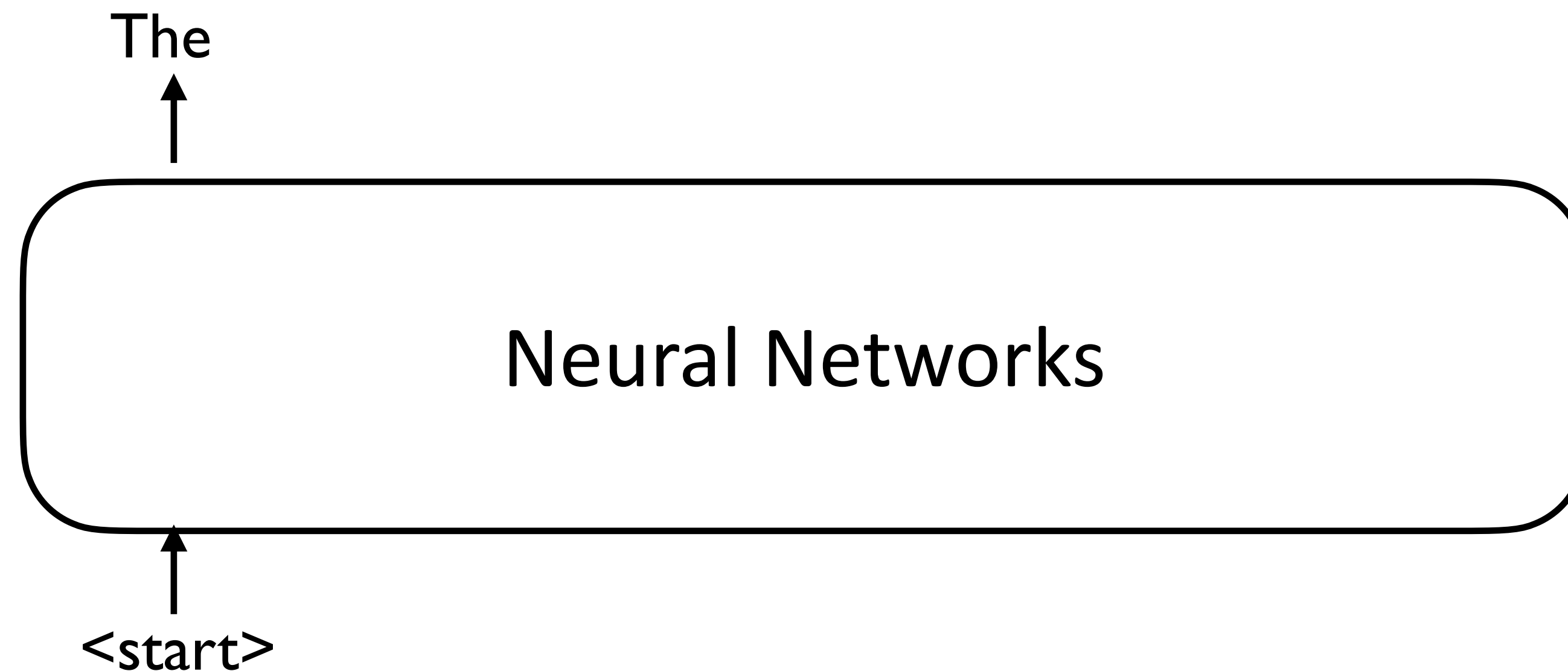
At inference time, to generate:

Neural Language Models

Are language models generative models? ✓

Can we compute $p(x)$ given x ? Can we sample new x ? ✓

At inference time, to generate:

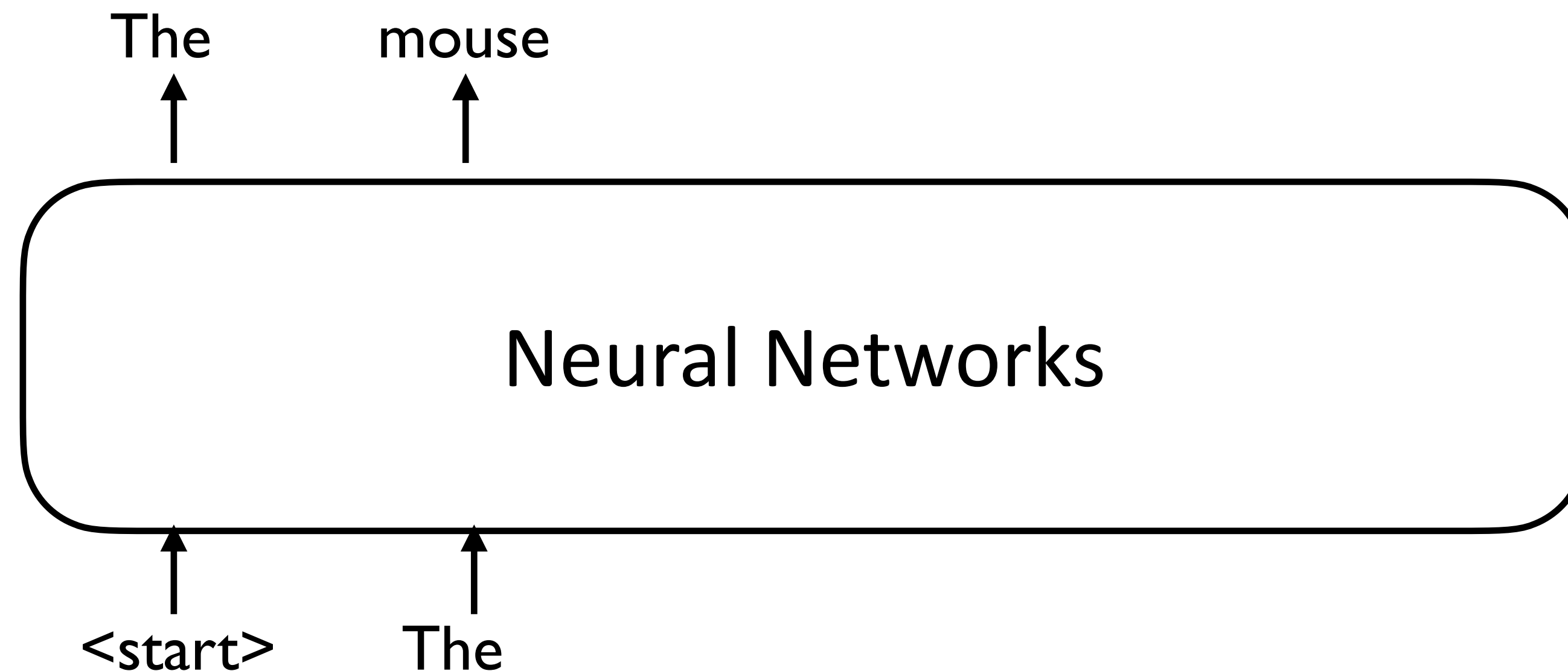


Neural Language Models

Are language models generative models? ✓

Can we compute $p(x)$ given x ? Can we sample new x ? ✓

At inference time, to generate:

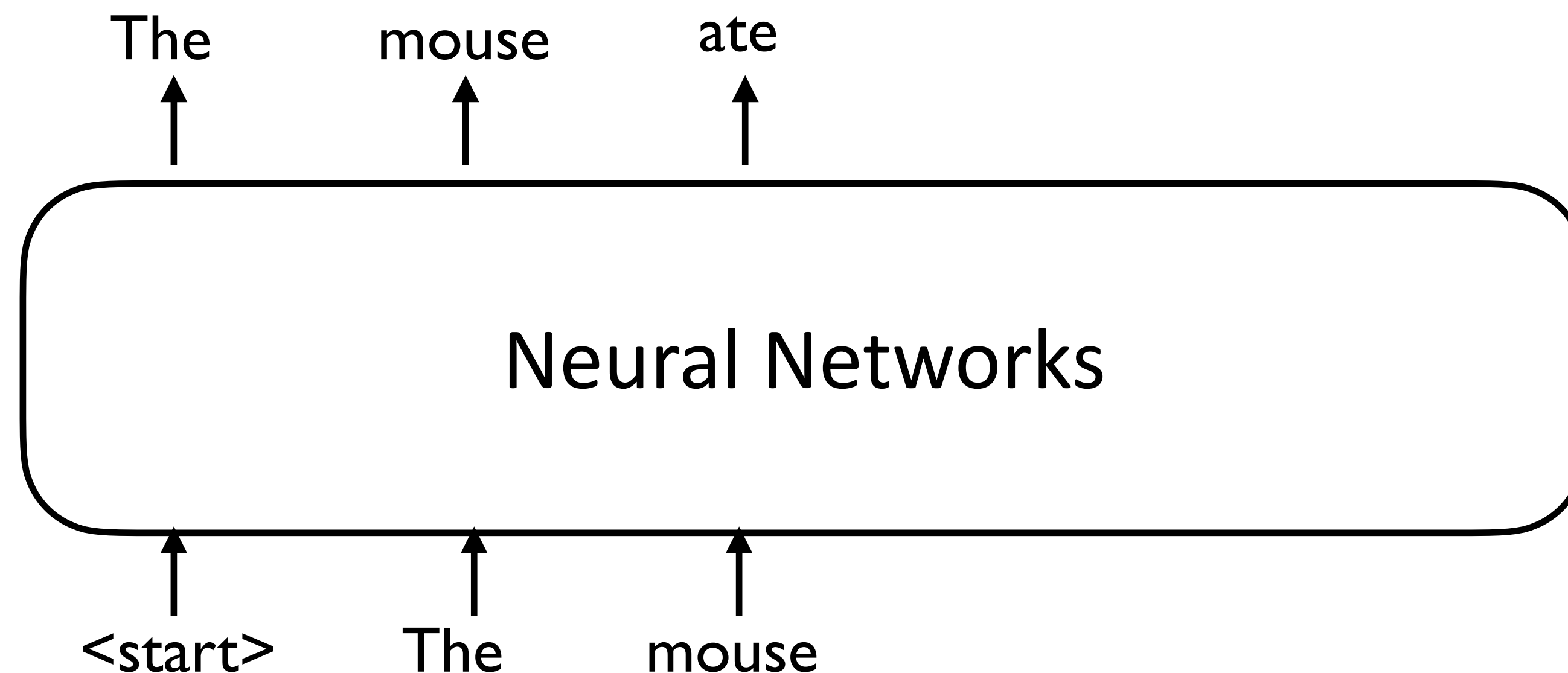


Neural Language Models

Are language models generative models? ✓

Can we compute $p(x)$ given x ? Can we sample new x ? ✓

At inference time, to generate:

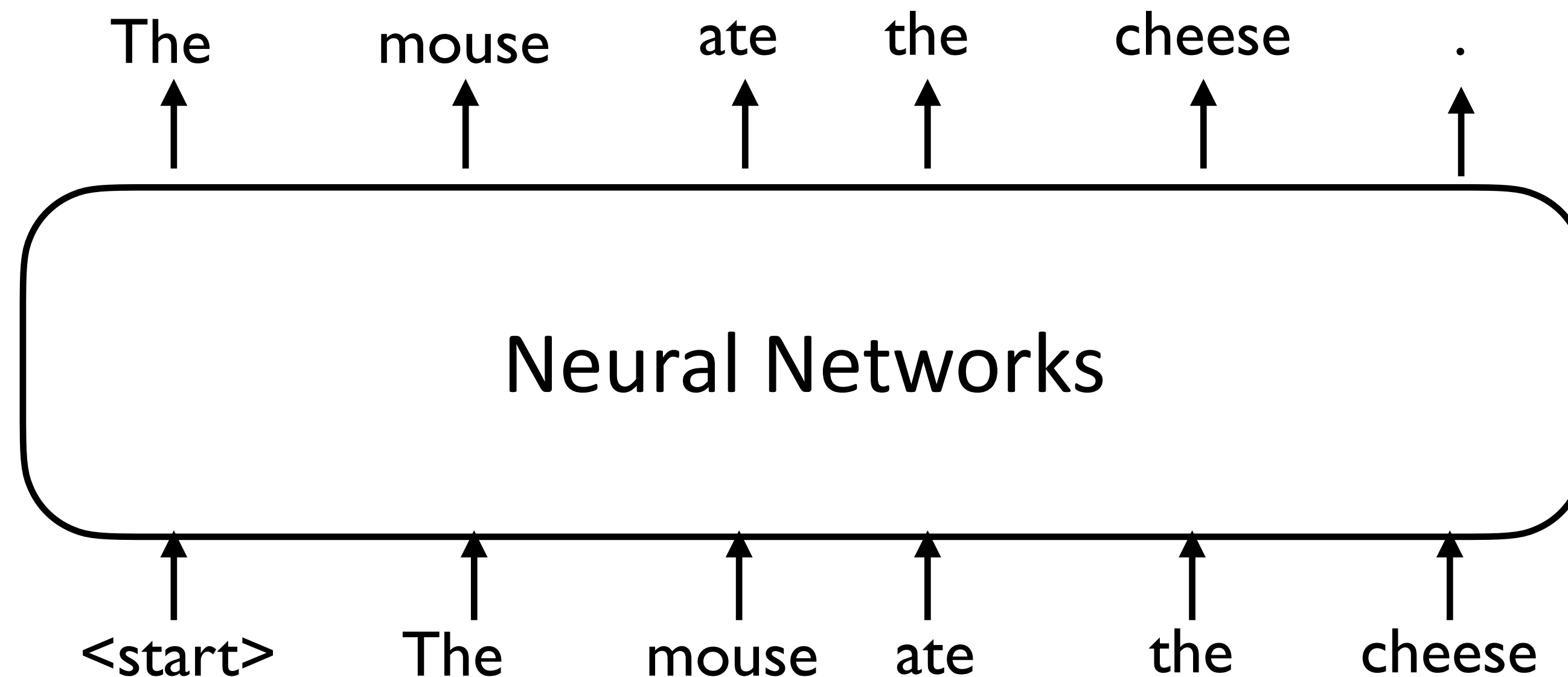


Neural Language Models

Are language models generative models? ✓

Can we compute $p(x)$ given x ? Can we sample new x ? ✓

At inference time, to generate:

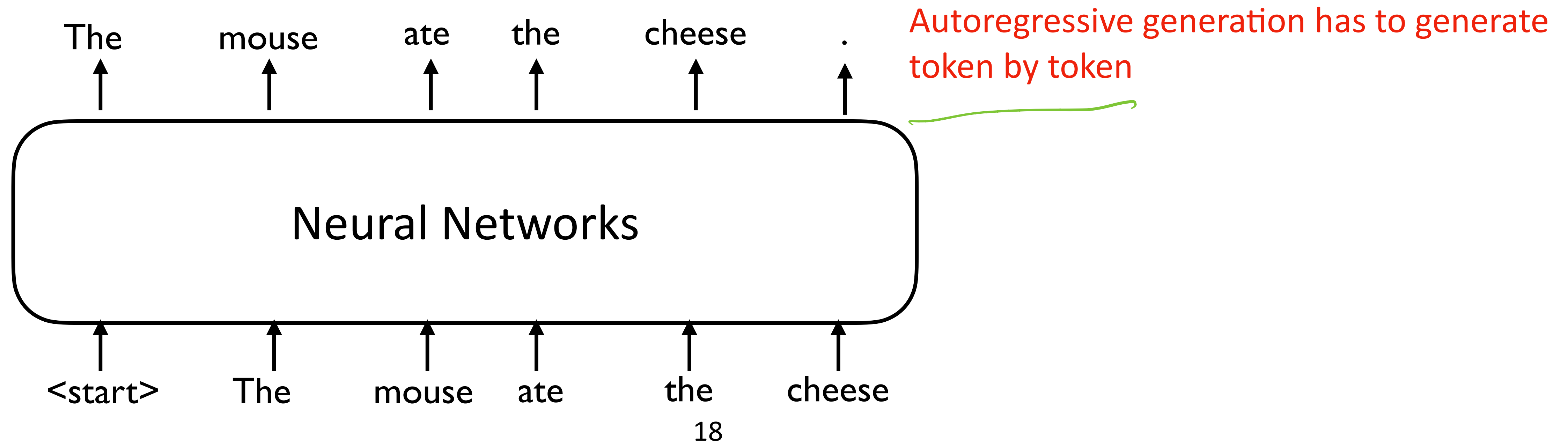


Neural Language Models

Are language models generative models? ✓

Can we compute $p(x)$ given x ? Can we sample new x ? ✓

At inference time, to generate:

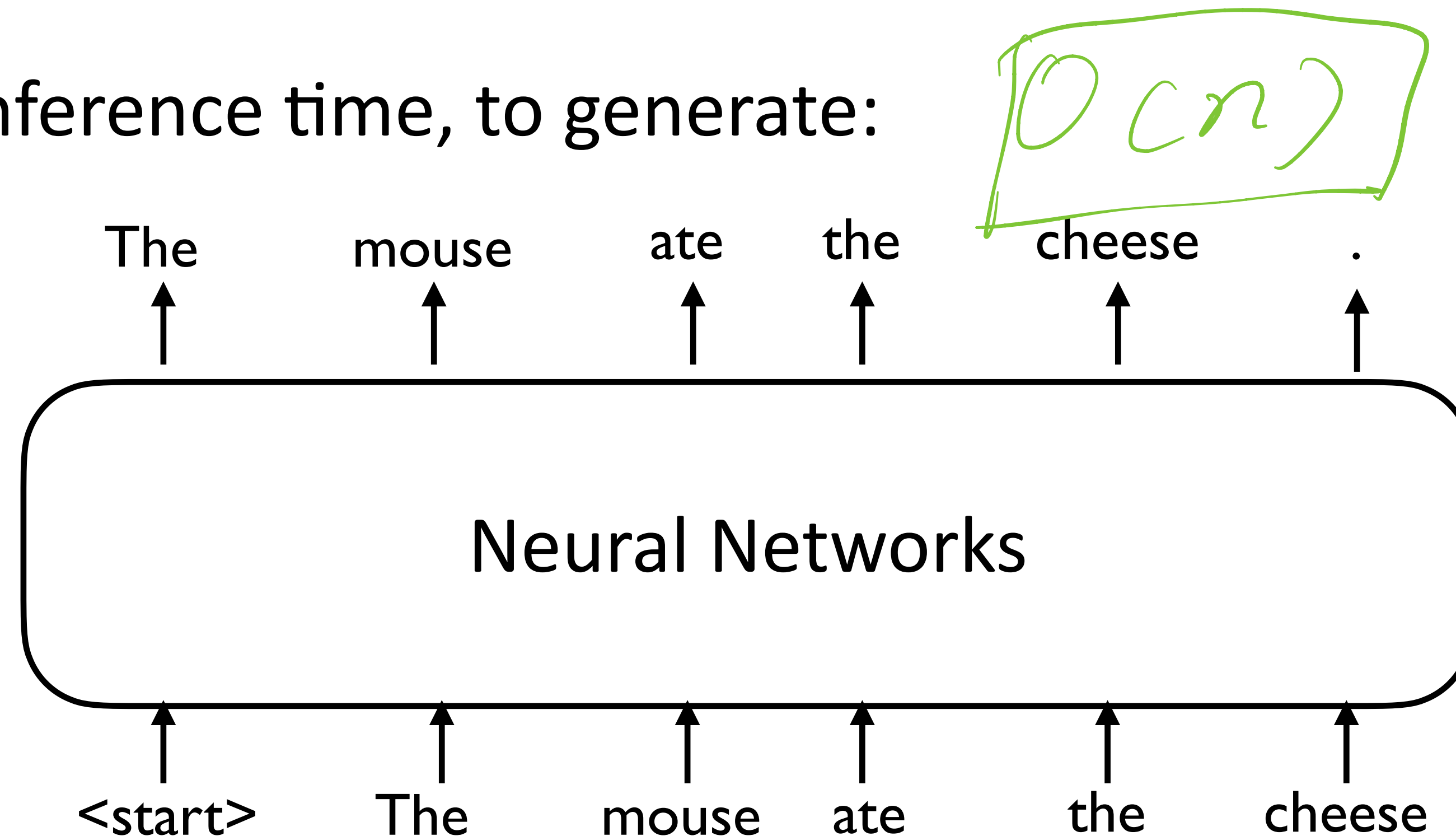


Neural Language Models

Are language models generative models? ✓

Can we compute $p(x)$ given x ? Can we sample new x ? ✓

At inference time, to generate:

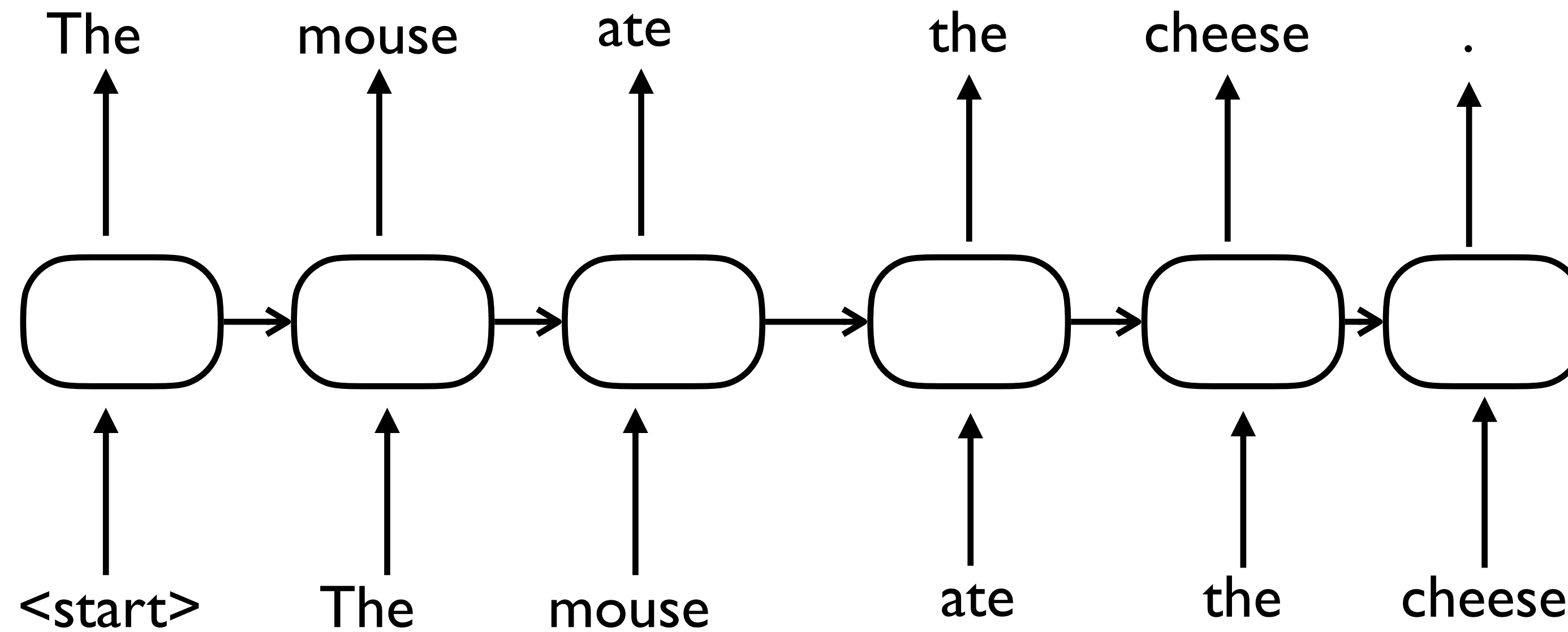


Autoregressive generation has to generate token by token

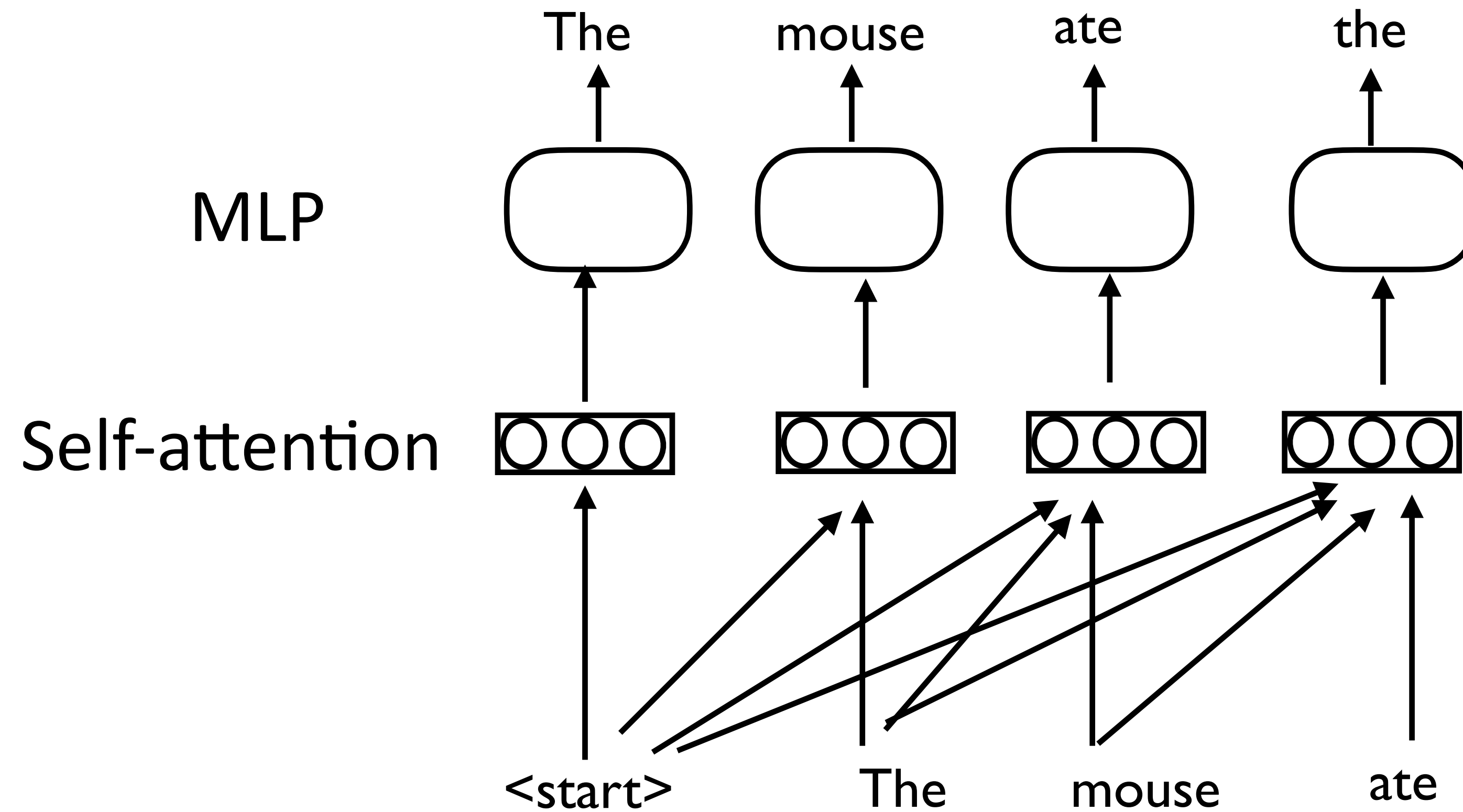
Sequential

Can't parallelize, efficiency of autoregressive decoding is still an important research topic

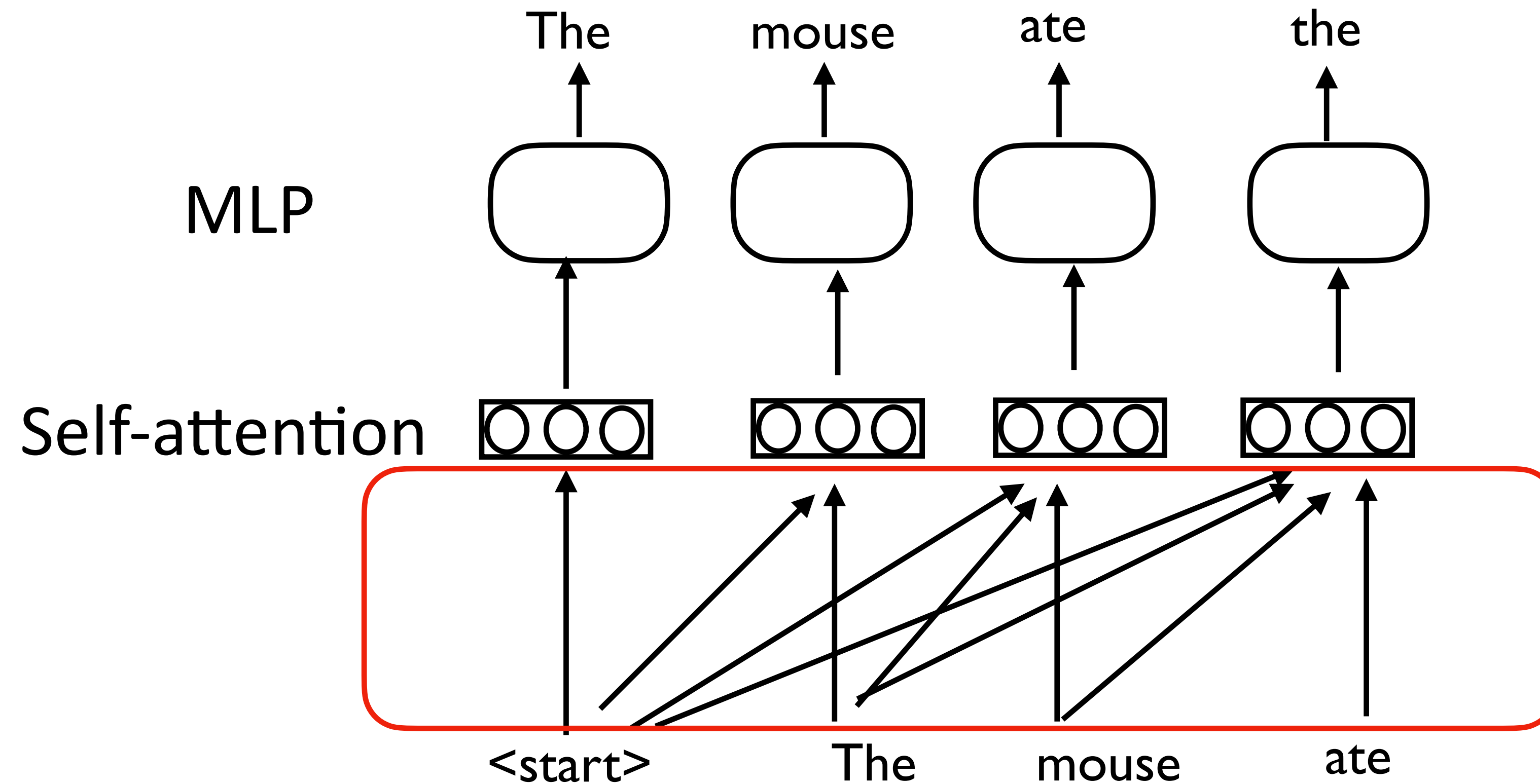
RNN Language Models



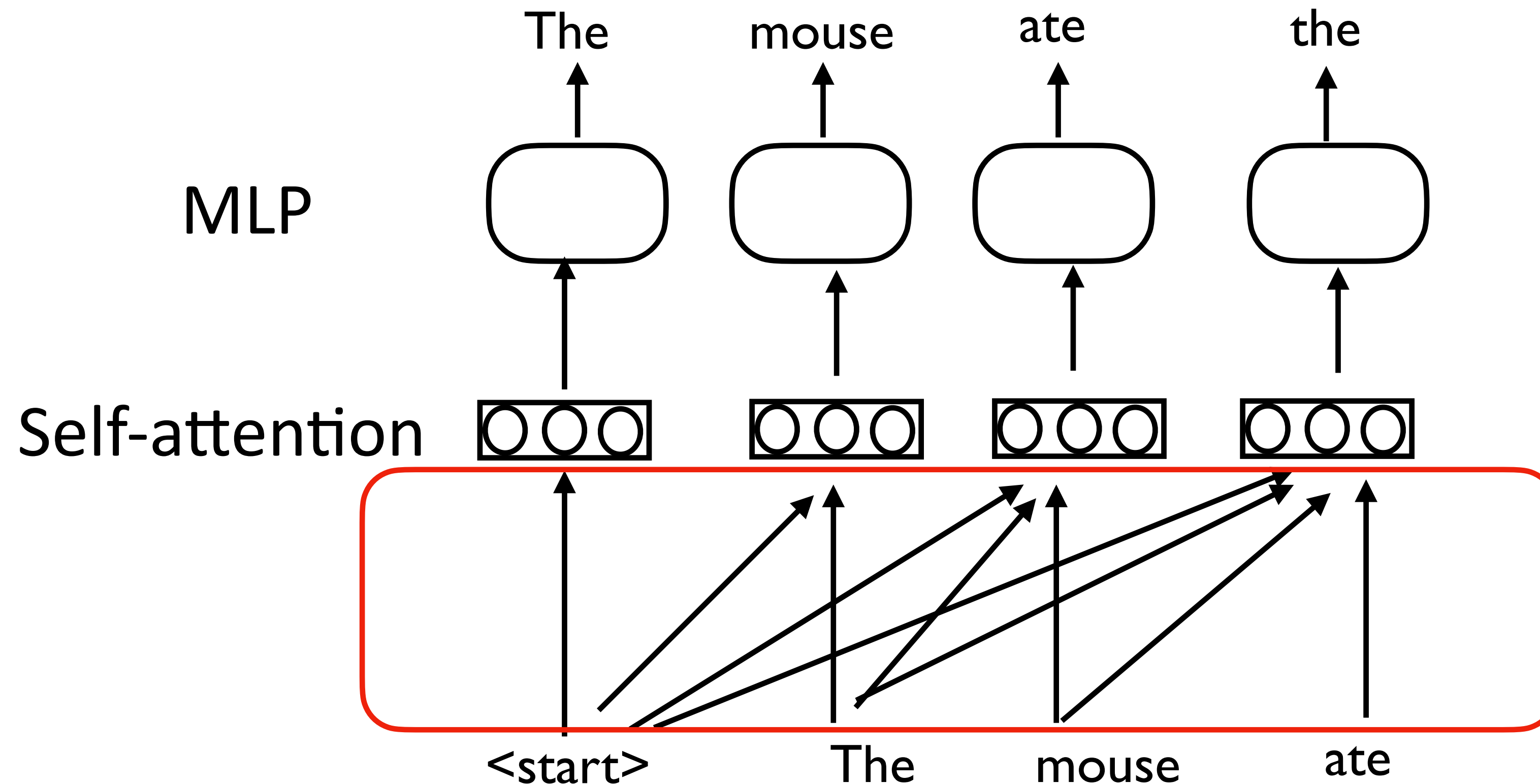
Transformer Language Models



Transformer Language Models



Transformer Language Models



Self-attention only attends to the tokens on the left (masked attention)

Neural Language Models

Language model is the fundamental block to model language distribution $p(x)$

Neural Language Models

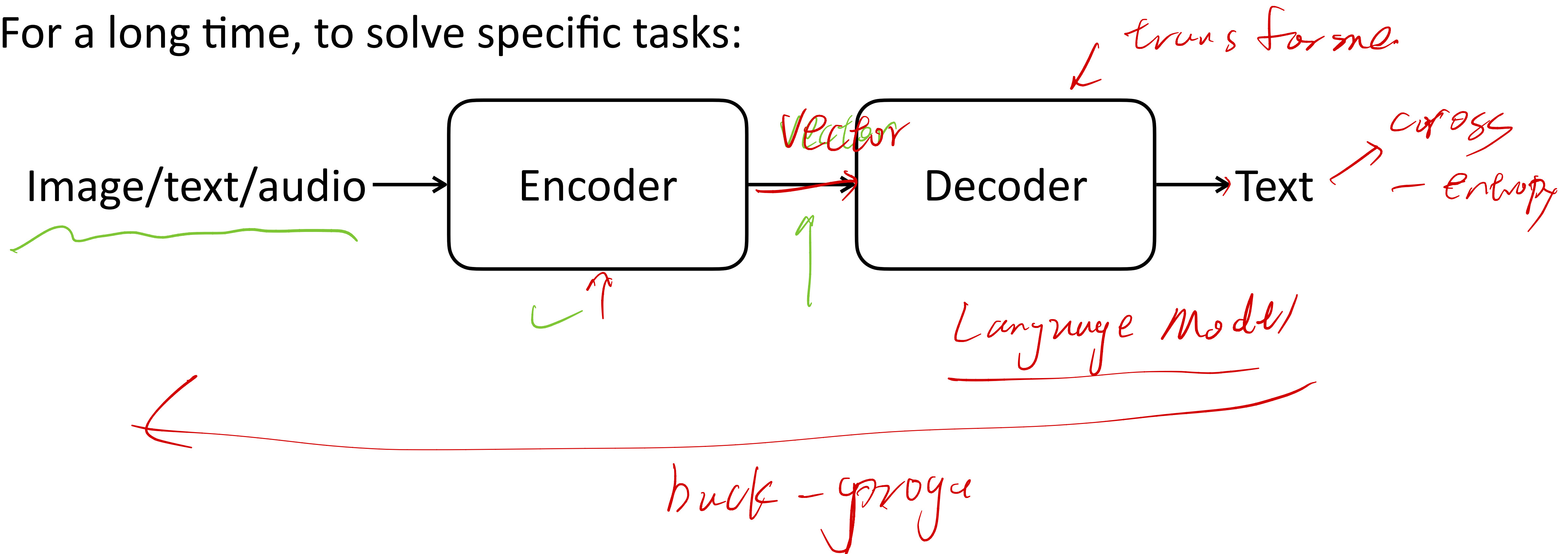
Language model is the fundamental block to model language distribution $p(x)$

For a long time, to solve specific tasks:

Neural Language Models

Language model is the fundamental block to model language distribution $p(x)$

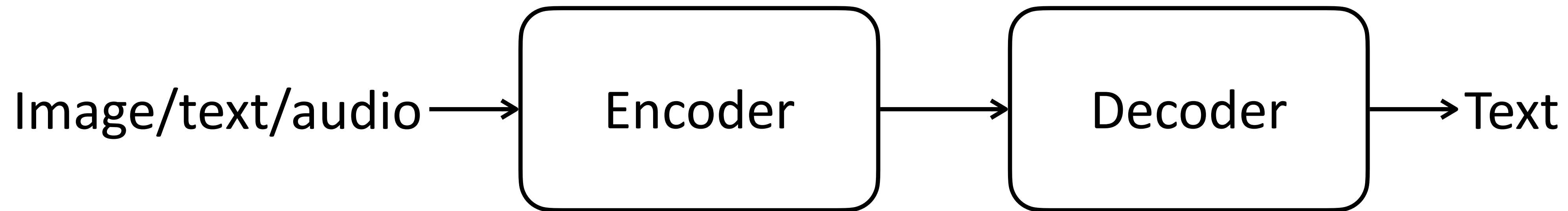
For a long time, to solve specific tasks:



Neural Language Models

Language model is the fundamental block to model language distribution $p(x)$

For a long time, to solve specific tasks:

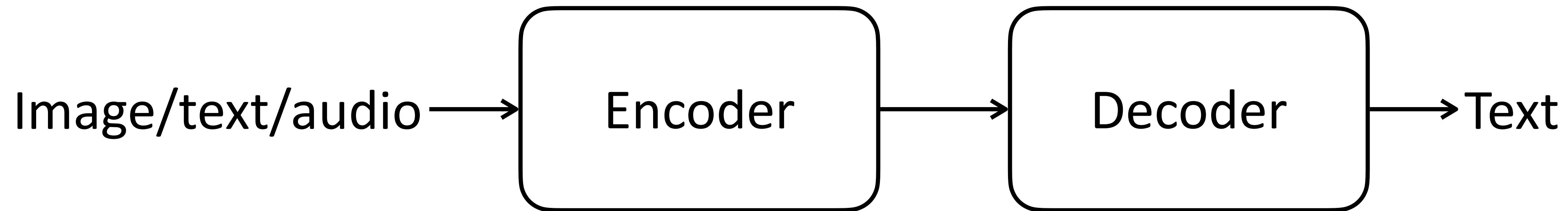


When we have a better arch/training
for LM, we can have a better decoder

Neural Language Models

Language model is the fundamental block to model language distribution $p(x)$

For a long time, to solve specific tasks:



When we have a better arch/training for LM, we can have a better decoder

Not long ago, some people think purely language models is useless because it does not directly address tasks, and LM performance may not transfer to downstream tasks

translation

summarization

QA

Is Next Token Prediction Useful?

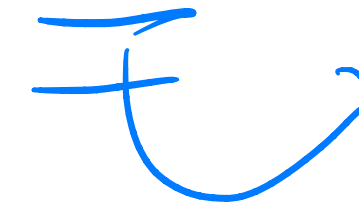
Ok, language modeling can be used as pretraining, but is a language model itself useful for some tasks directly?

Is Next Token Prediction Useful?

Ok, language modeling can be used as pretraining, but is a language model itself useful for some tasks directly?

In the late 1980s the Hong Kong Government anticipated a strong demand for university graduates to fuel an economy increasingly based on services. Sir Sze-Yuen Chung and Sir Edward Youde, the then Governor of Hong Kong, conceived the idea of another university in addition to the pre-existing two universities, The University of Hong Kong and The Chinese University of Hong Kong.

Planning for the "Third University", named The Hong Kong University of Science and Technology later, began in 1986. Construction began at the Kohima Camp site in Tai Po Tsai on the Clear Water Bay Peninsula. The site was earmarked for the construction of a new []



Is Next Token Prediction Useful?

Ok, language modeling can be used as pretraining, but is a language model itself useful for some tasks directly?

In the late 1980s the Hong Kong Government anticipated a strong demand for university graduates to fuel an economy increasingly based on services. Sir Sze-Yuen Chung and Sir Edward Youde, the then Governor of Hong Kong, conceived the idea of another university in addition to the pre-existing two universities, The University of Hong Kong and The Chinese University of Hong Kong.

Planning for the "Third University", named The Hong Kong University of Science and Technology later, began in 1986. Construction began at the Kohima Camp site in Tai Po Tsai on the Clear Water Bay Peninsula. The site was earmarked for the construction of a new []

Completion

Is Next Token Prediction Useful?

Ok, language modeling can be used as pretraining, but is a language model itself useful for some tasks directly?

In the late 1980s the Hong Kong Government anticipated a strong demand for university graduates to fuel an economy increasingly based on services. Sir Sze-Yuen Chung and Sir Edward Youde, the then Governor of Hong Kong, conceived the idea of another university in addition to the pre-existing two universities, The University of Hong Kong and The Chinese University of Hong Kong.

Planning for the "Third University", named The Hong Kong University of Science and Technology later, began in 1986. Construction began at the Kohima Camp site in Tai Po Tsai on the Clear Water Bay Peninsula. The site was earmarked for the construction of a new []

Completion

This task seems useless in practice

Language Models are Zero-Shot Learners

Context (human-written): 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.

GPT-2: The scientist named the population, after their distinctive horn, Ovid's Unicorn. These four-horned, silver-white unicorns were previously unknown to science.

Now, after almost two centuries, the mystery of what sparked this odd phenomenon is finally solved.

Dr. Jorge Pérez, an evolutionary biologist from the University of La Paz, and several companions, were exploring the Andes Mountains when they found a small valley, with no other animals or humans. Pérez noticed that the valley had what appeared to be a natural fountain, surrounded by two peaks of rock and silver snow.

Pérez and the others then ventured further into the valley. "By the time we reached the top of one peak, the water looked blue, with some crystals on top," said Pérez.

Pérez and his friends were astonished to see the unicorn herd. These creatures could be seen from the air without having to move too much to see them – they were so close they could touch their horns.

While examining these bizarre creatures the scientists discovered that the creatures also spoke some fairly regular English. Pérez stated, "We can see, for example, that they have a common 'language,' something like a dialect or dialectic."

Dr. Pérez believes that the unicorns may have originated in Argentina, where the animals were believed to be descendants of a lost race of people who lived there before the arrival of humans in those parts of South America.

While their origins are still unclear, some believe that perhaps the creatures were created when a human and a unicorn met each other in a time before human civilization. According to Pérez, "In South America, such incidents seem to be quite common."

However, Pérez also pointed out that it is likely that the only way of knowing for sure if unicorns are indeed the descendants of a lost alien race is through DNA. "But they seem to be able to communicate in English quite well, which I believe is a sign of evolution, or at least a change in social organization," said the scientist.

GPT-2

Radford et al. Language Models are Unsupervised Multitask Learners. 2018.

GPT-2

Next token prediction can unify many tasks

Radford et al. Language Models are Unsupervised Multitask Learners. 2018.

GPT-2

Next token prediction can unify many tasks

Machine translation:

Radford et al. Language Models are Unsupervised Multitask Learners. 2018.

GPT-2

Next token prediction can unify many tasks

Machine translation:

Chinese: 今天是学期的最后一天。

English: 

Radford et al. Language Models are Unsupervised Multitask Learners. 2018.

GPT-2

Next token prediction can unify many tasks

Machine translation:

Chinese: 今天是学期的最后一天。
English:

Question answering:

Q: What is the capital of the United States?

A:



Radford et al. Language Models are Unsupervised Multitask Learners. 2018.

GPT-2

Next token prediction can unify many tasks

Machine translation:

Chinese: 今天是学期的最后一天。
English:

Completion is very general

Question answering:

Q: What is the capital of the United States?
A:

Radford et al. Language Models are Unsupervised Multitask Learners. 2018.

GPT-2

Next token prediction can unify many tasks

Machine translation:

Chinese: 今天是学期的最后一天。
English:

Completion is very general

Question answering:

Q: What is the capital of the United States?
A:

This was an early form of prompting,
that is widely discussed today

Radford et al. Language Models are Unsupervised Multitask Learners. 2018.

Language Models Are Few-Shot Learners

Language Models Are Few-Shot Learners

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



< 1B
175B

large language
models

Language Models Are Few-Shot Learners

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
```

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 ← demo example
3 cheese => ..... ← prompt
```

Language Models Are Few-Shot Learners

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
```

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-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
```


Language Models Are Few-Shot Learners

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
```

One-shot

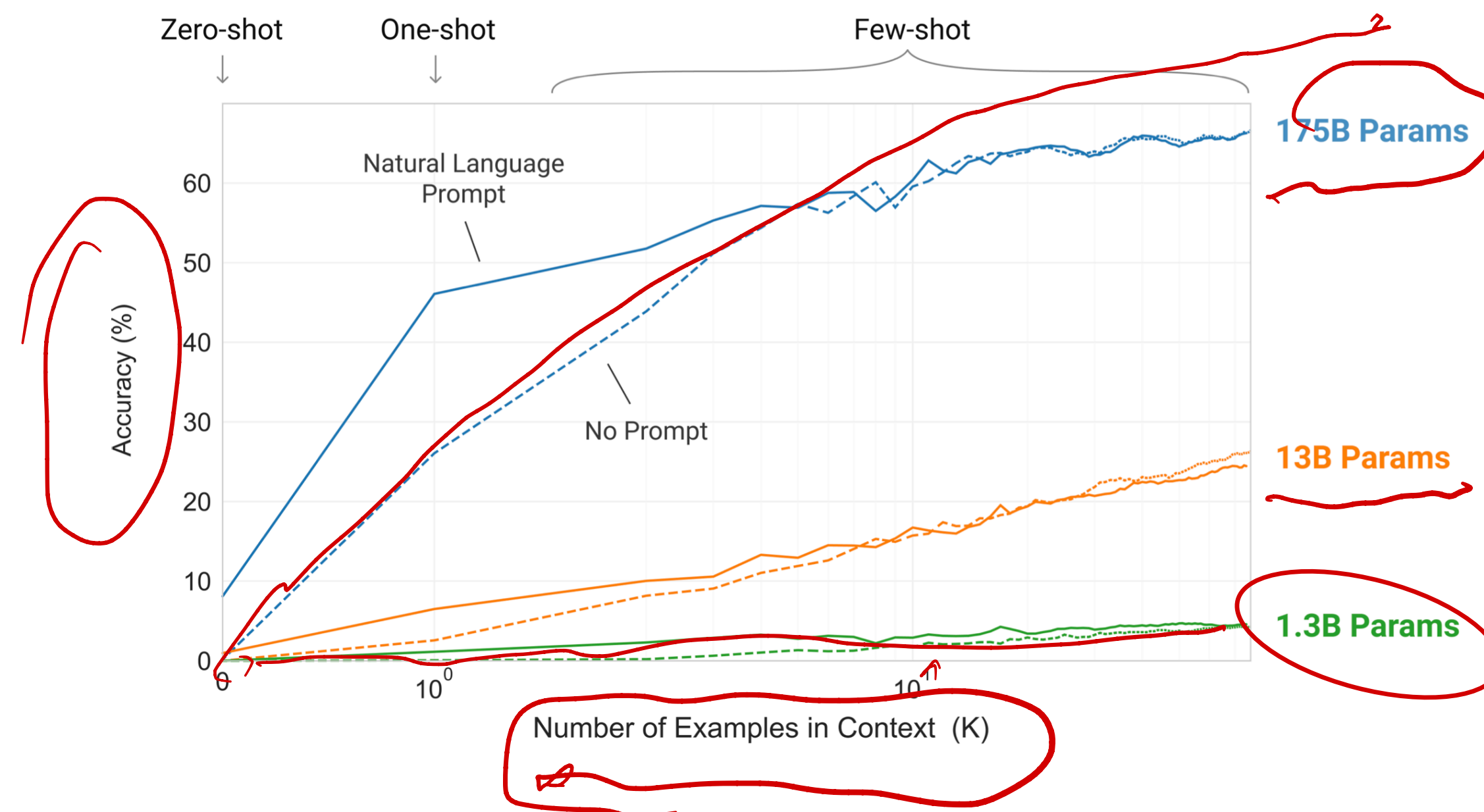
In addition to the task description, the model sees a single example of the task. No gradient updates are performed.

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

Few-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 ← examples
4 plush girafe => girafe peluche ← examples
5 cheese => ..... ← prompt
```



Language Models Are Few-Shot Learners

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
```

One-shot

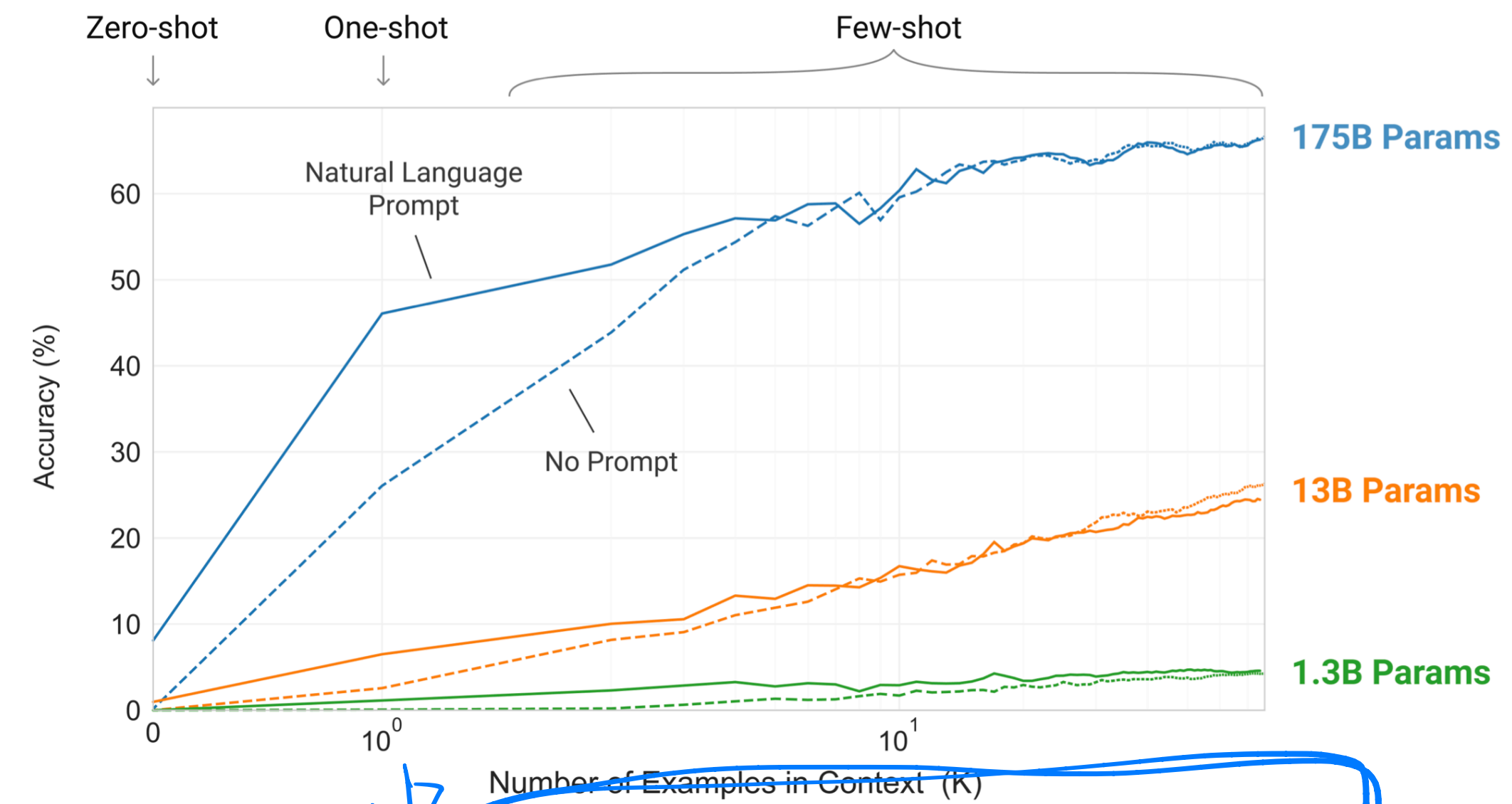
In addition to the task description, the model sees a single example of the task. No gradient updates are performed.

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

Few-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 ← examples
4 plush girafe => girafe peluche ← examples
5 cheese => ..... ← prompt
```



In-Context Learning

Brown et al. Language models are few-shot learners. 2020

Pretraining

Pretraining

Target Data B



Pretraining

Source Data A (maybe a different task)



Target Data B

Pretraining

long maybe model!

Source Data A (maybe a different task)

Train on data A first

Model

Target Data B

translation

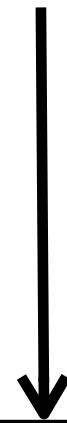
supervise

Unsupervised language modeling:

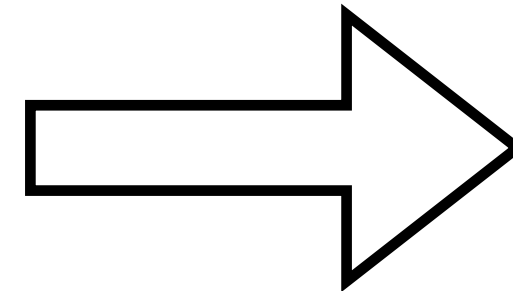
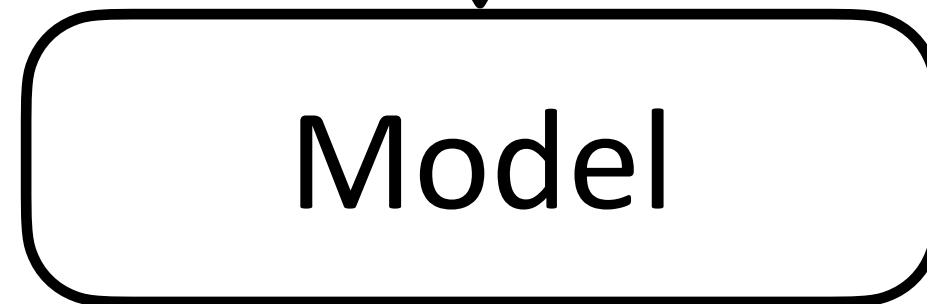
the mouse ate the cheese

Pretraining

Source Data A (maybe a different task)



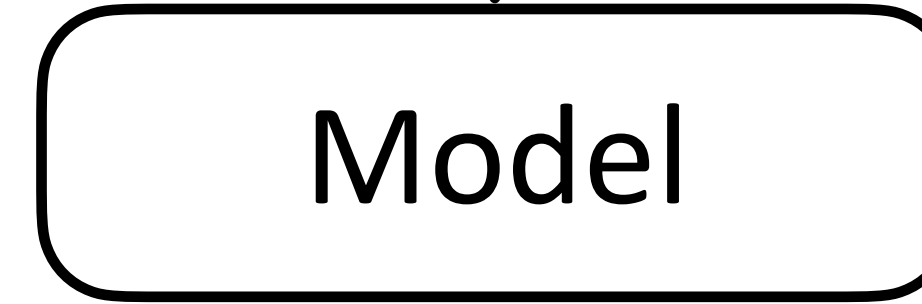
Train on data A first



Target Data B



Then train on data B

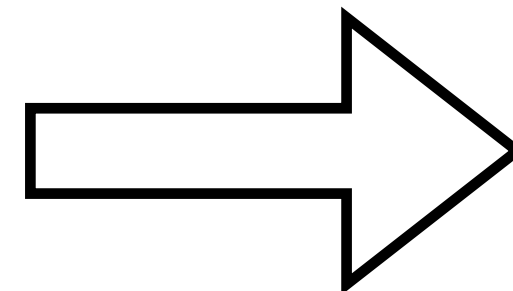
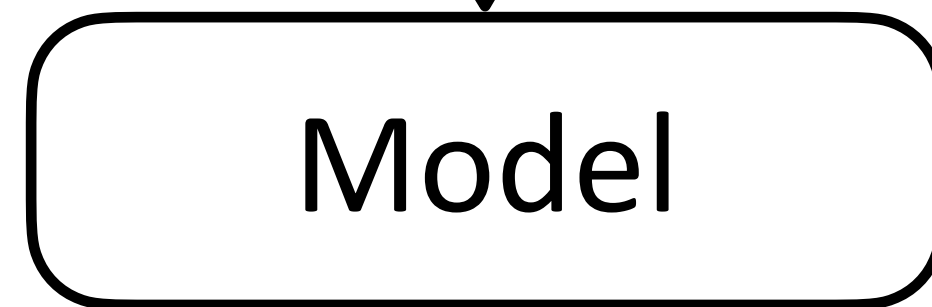


Pretraining

Source Data A (maybe a different task)



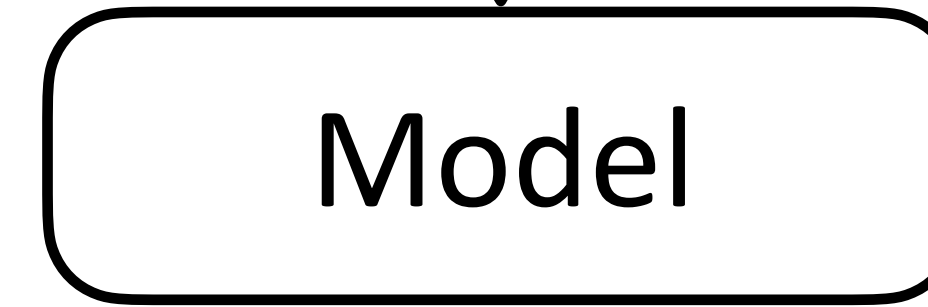
Train on data A first



Target Data B



Then train on data B



Classically, this is transfer Learning

Eng → German

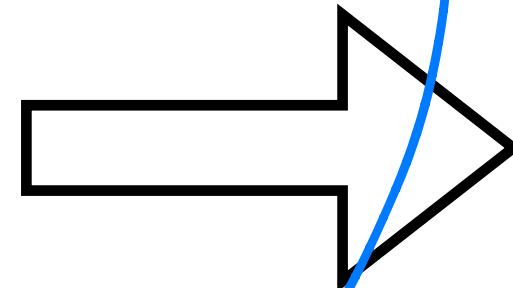
Pretraining

small scale traditionally

Source Data A (maybe a different task)

Train on data A first

Model



Target Data B

Then train on data B

Model

Eng → Fren

Classically, this is transfer Learning

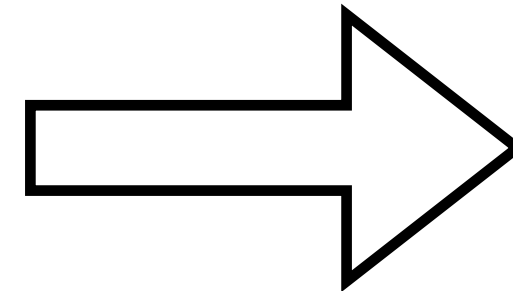
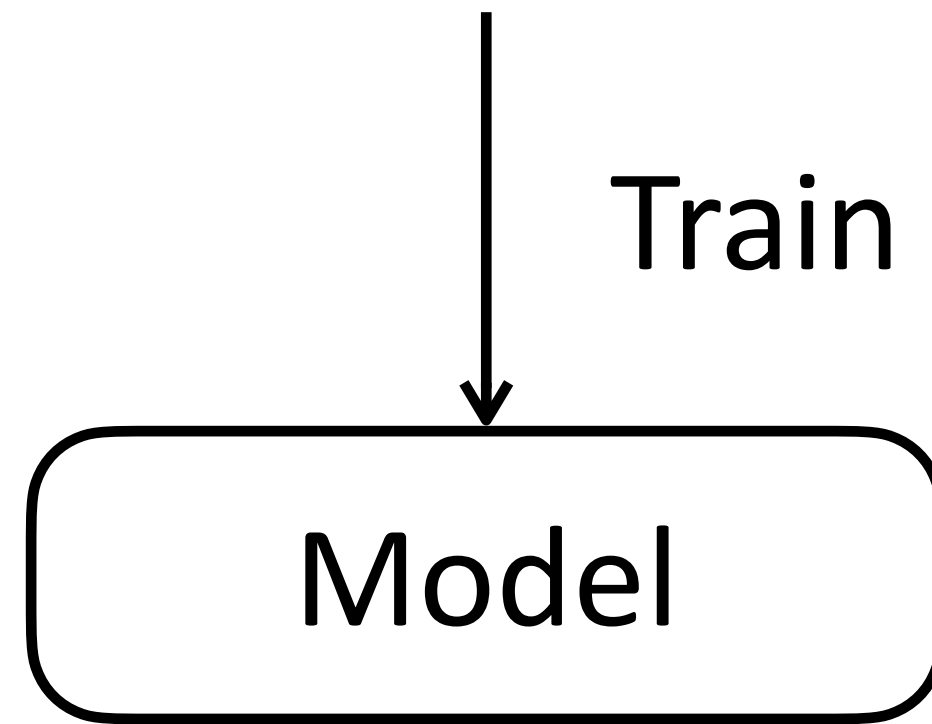
It is now called pretraining because of the scale of A

*unsupervised
language
modeling*

Pretraining

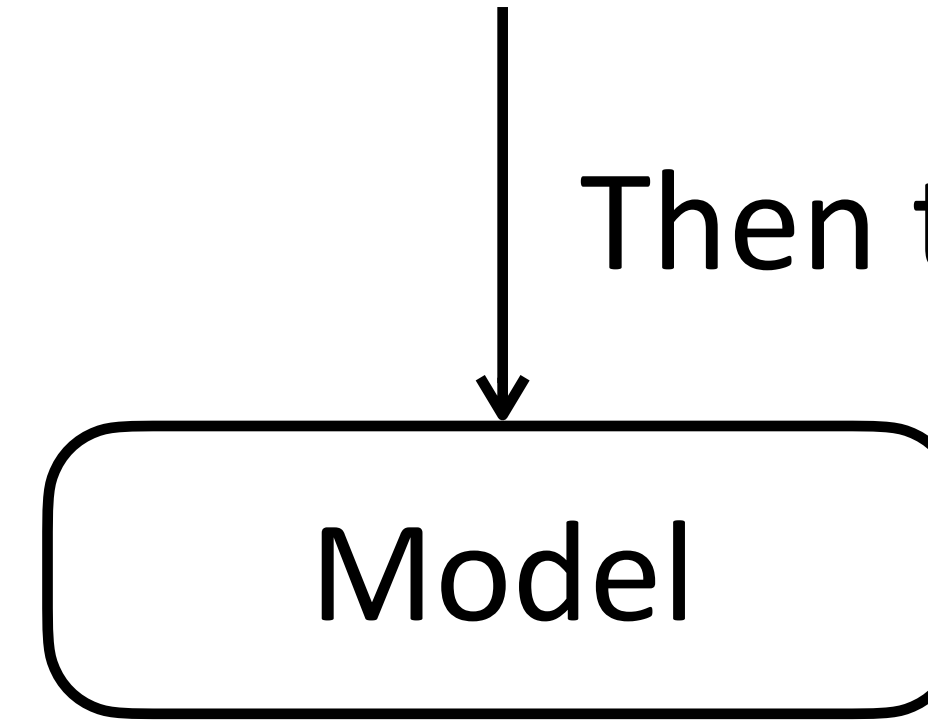
Source Data A (maybe a different task)

Train on data A first



Target Data B

Then train on data B



For supervised training, data A is often limited

How can we find large-scale data A to train?

LM

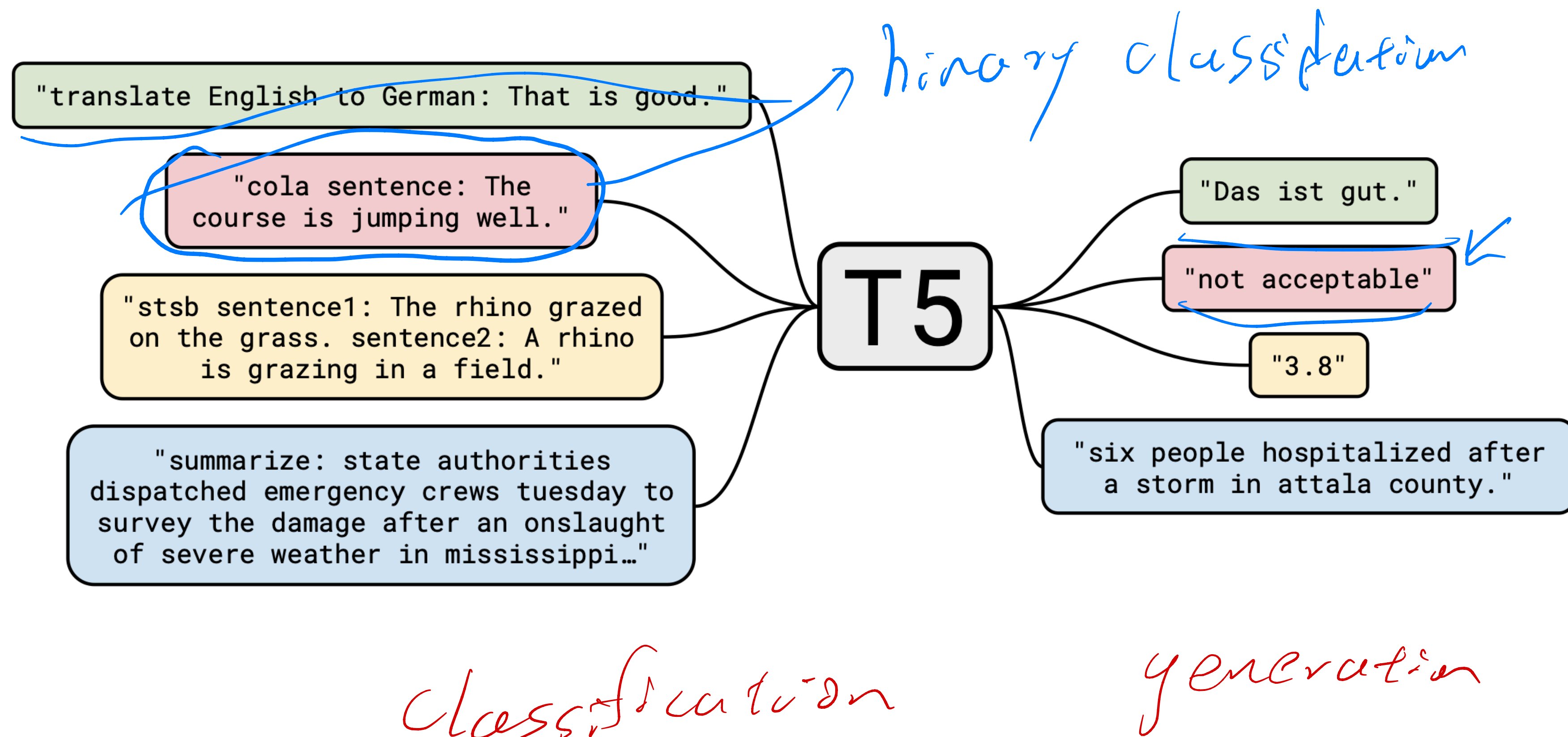
Prompt Breaks Task Boundaries

Prompt Breaks Task Boundaries

Almost all text tasks can be expressed with a unified format, no matter whether it is classification or generation

Prompt Breaks Task Boundaries

Almost all text tasks can be expressed with a unified format, no matter whether it is classification or generation



Raffle et al. Exploring the Limits of Transfer Learning. 2020

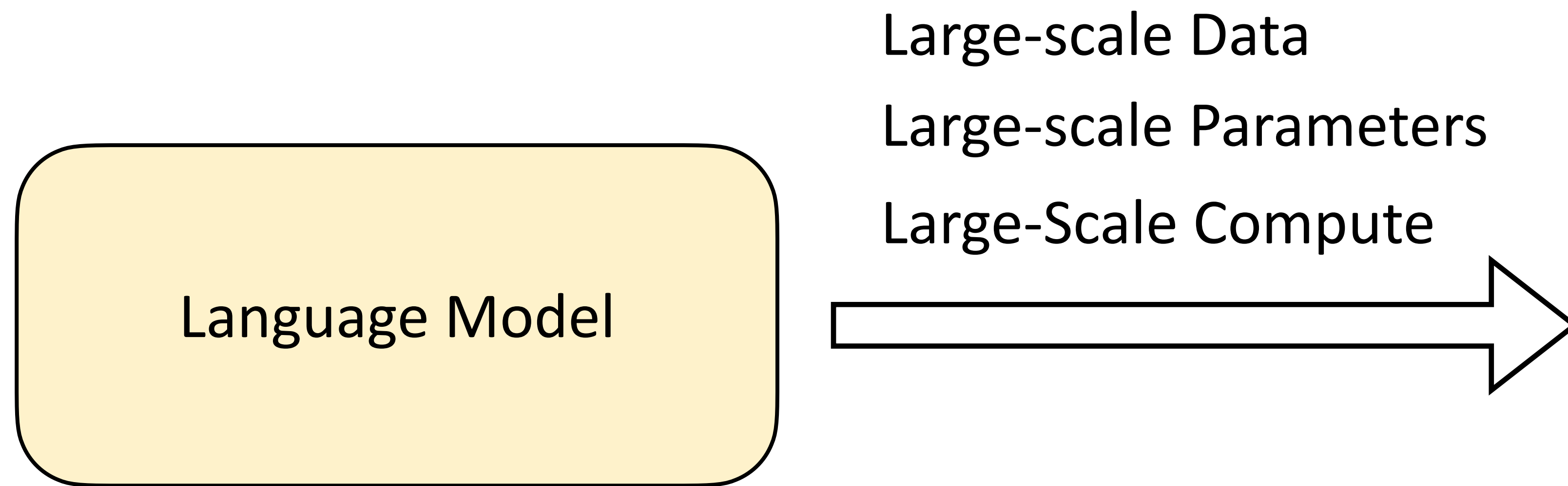
Large Language Models

Large Language Models

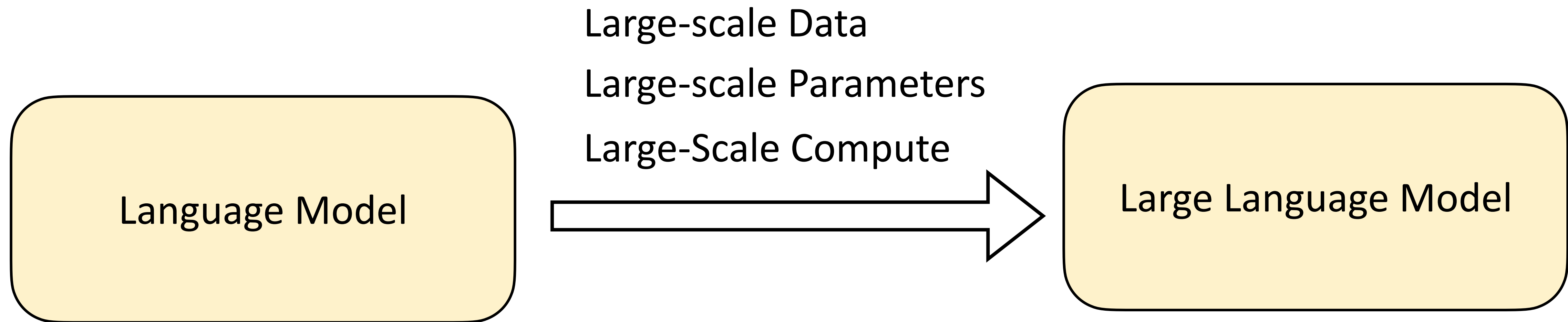


Language Model

Large Language Models



Large Language Models



Thank You!