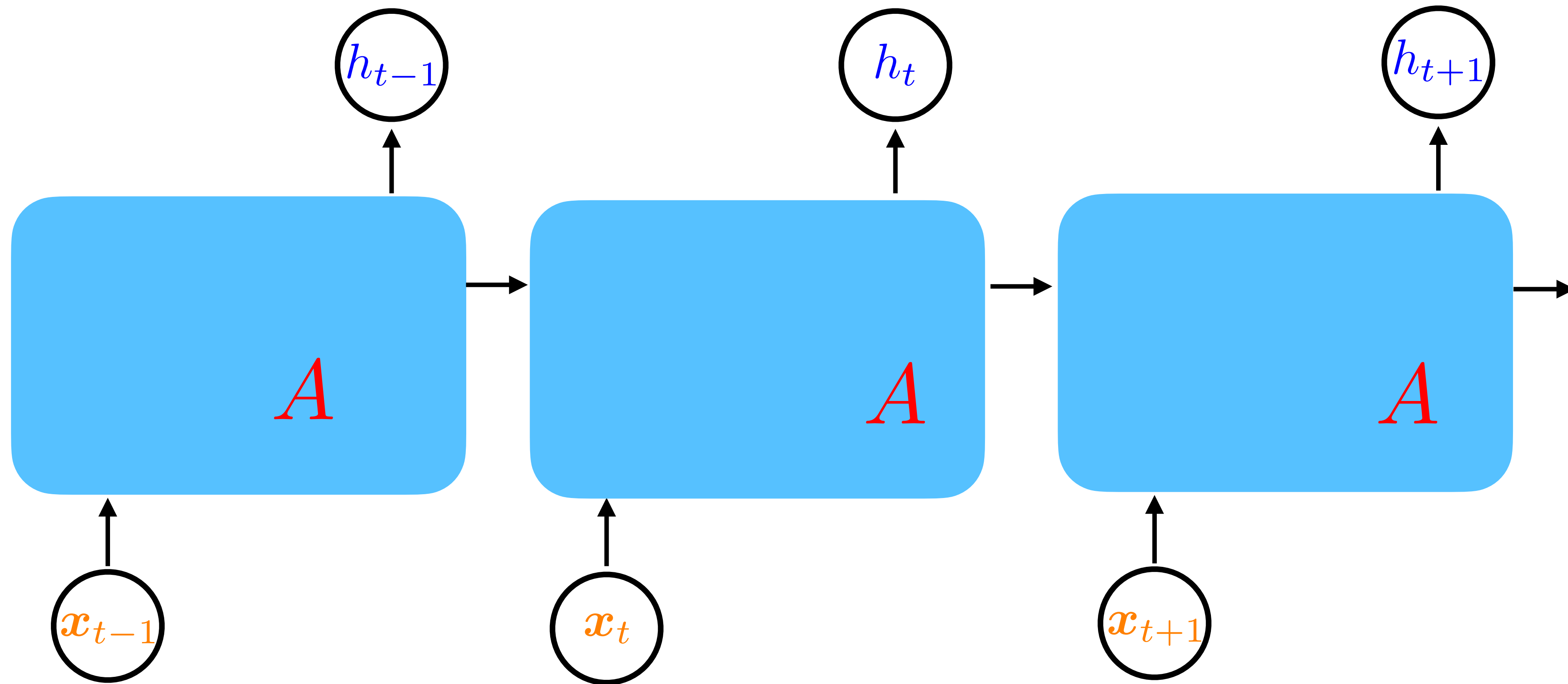# Sequence to Sequence Model and Attention

COMP7607 — Lecture 3
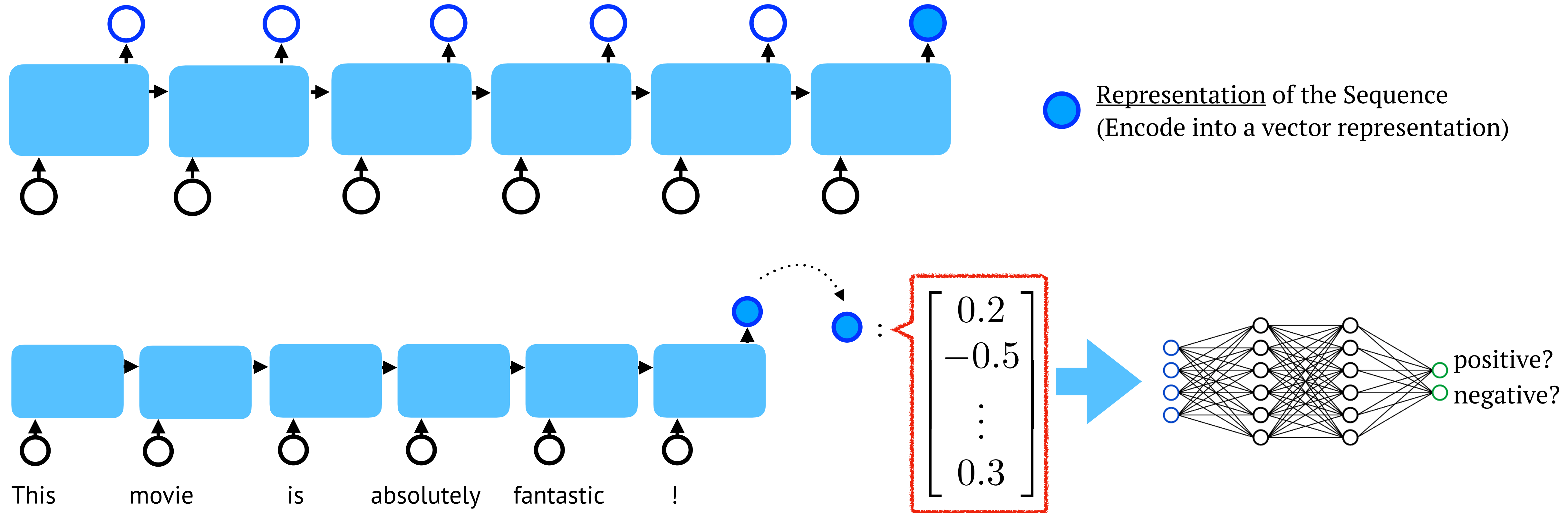
**Lingpeng Kong**

**Department of Computer Science, The University of Hong Kong**
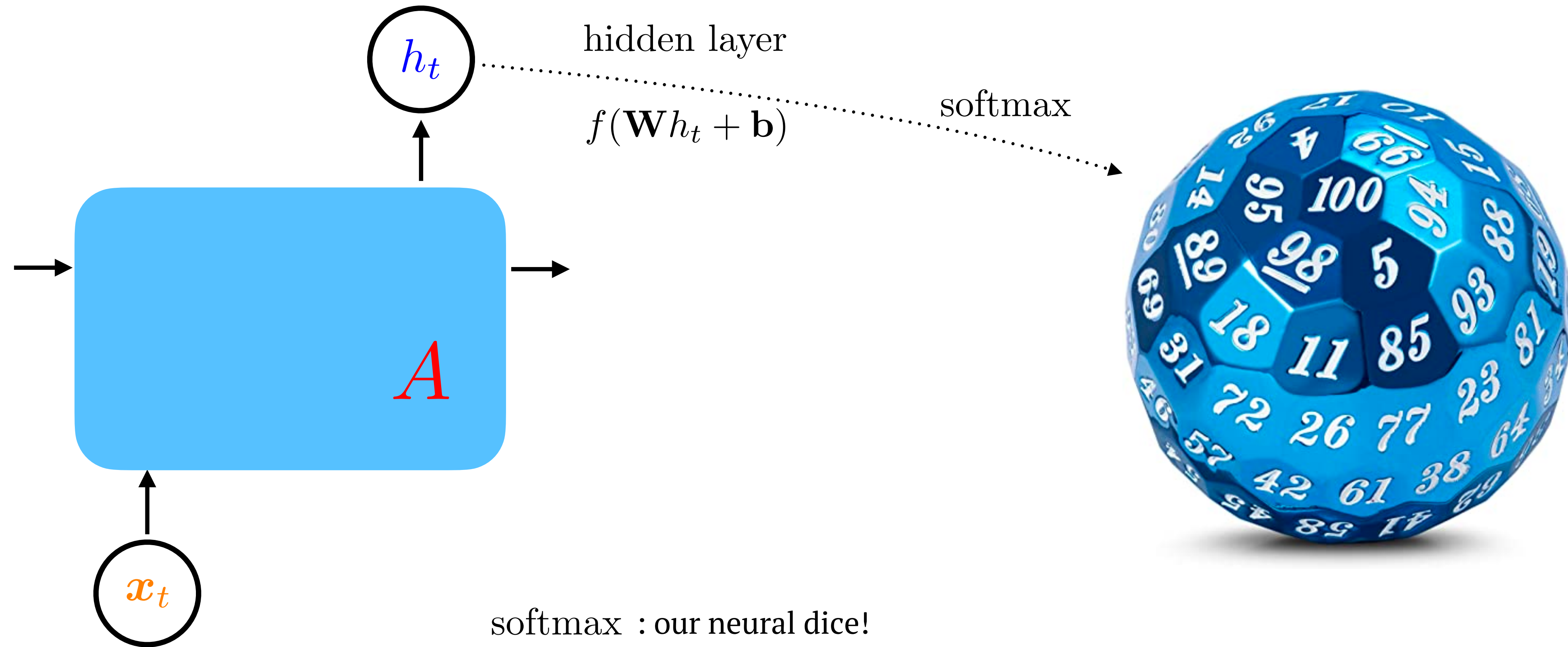
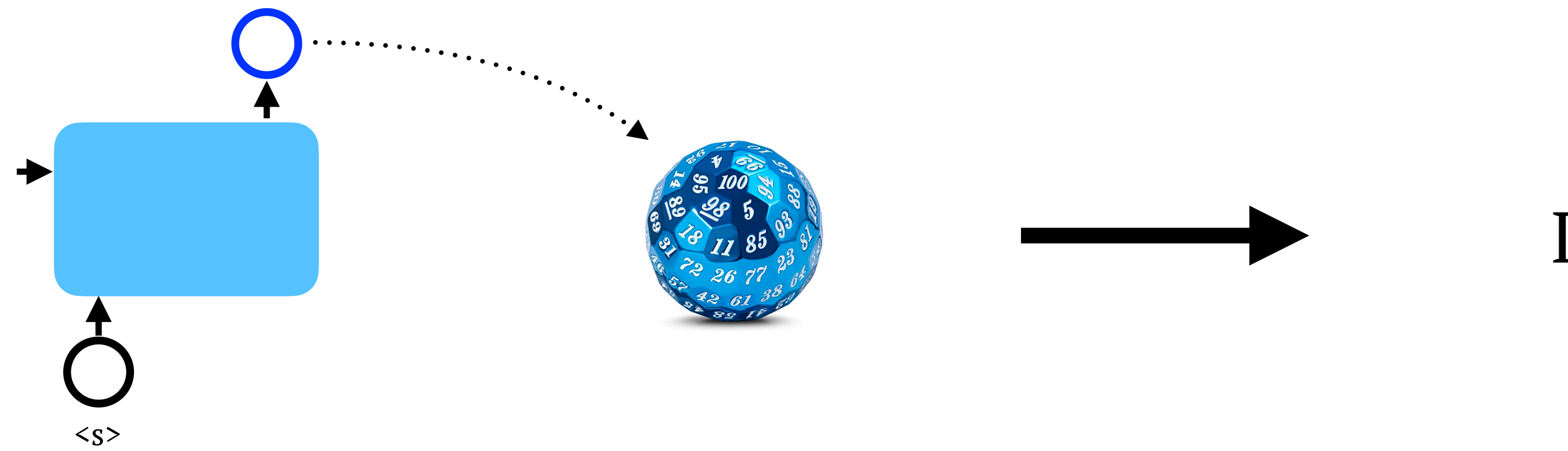# Recurrent Neural Network

# RNN as Encoder



Representation of the Sequence
(Encode into a vector representation)

$$\begin{bmatrix} 0.2 \\ -0.5 \\ \vdots \\ 0.3 \end{bmatrix}$$

This movie is absolutely fantastic !

positive?
negative?

# Sample a sentence from RNNLMs

$h_t$

hidden layer

$f(\mathbf{W}h_t + \mathbf{b})$

softmax

$A$

$x_t$

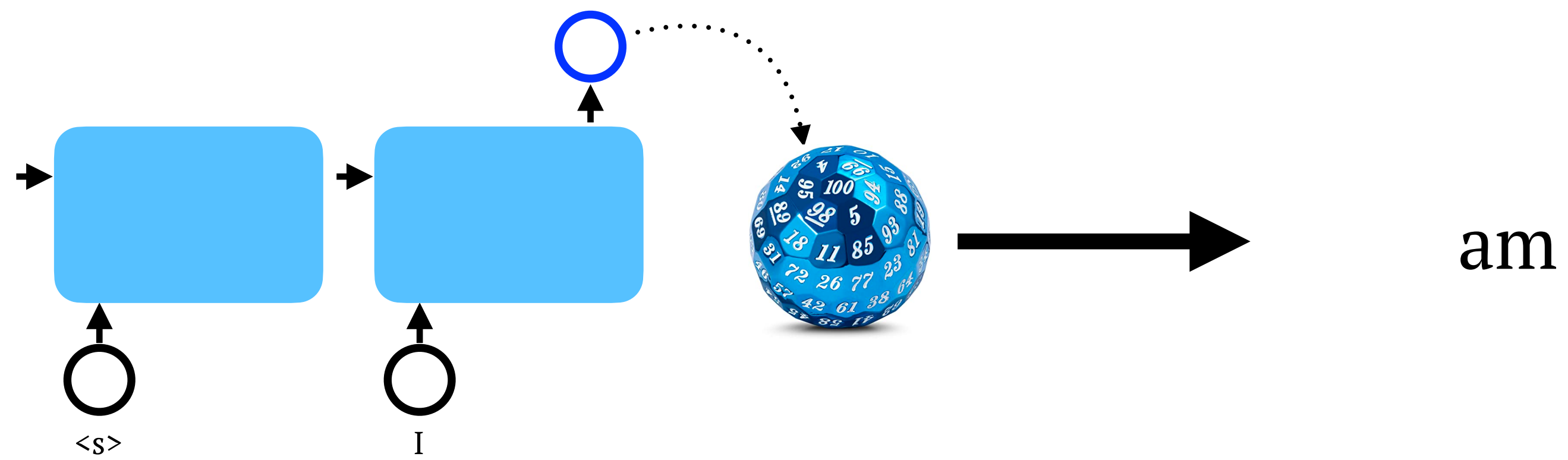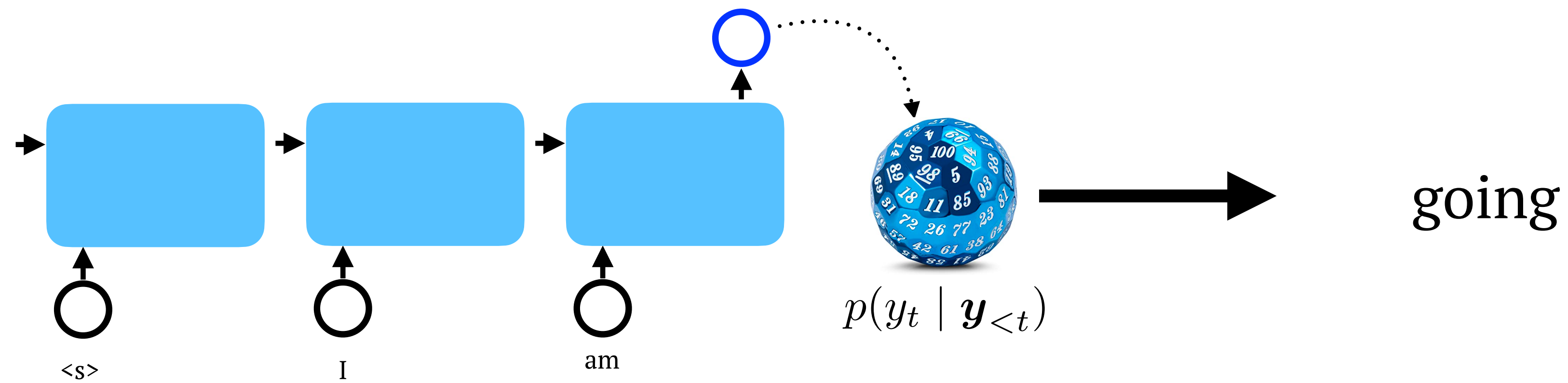softmax : our neural dice!

# Sample a sentence from RNNLMs

I



I

# Sample a sentence from RNNLMs

I    am

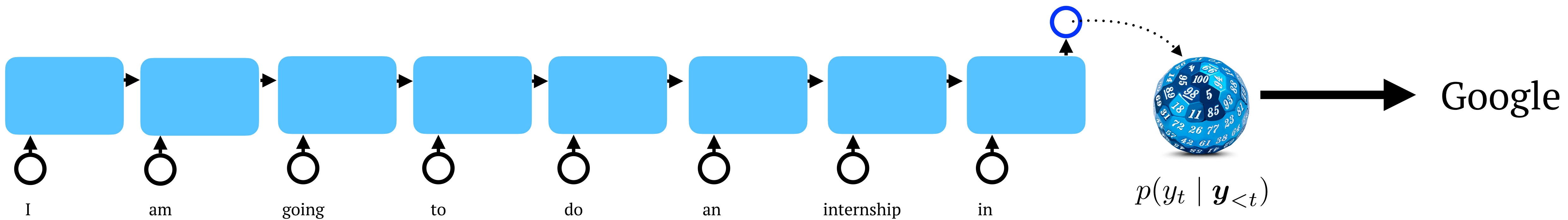# Sample a sentence from RNNLMs

I     am     going



$p(y_t \mid \boldsymbol{y}_{<t})$

going

<s>     I     am

# Sample a sentence from RNNLMs

I     am     going     to     do     an     internship     in     Google



$$p(y_t \mid \boldsymbol{y}_{<t})$$

I     am     going     to     do     an     internship     in

Google

# RNN as Decoder (RNNLM)



$$p(y_t \mid \boldsymbol{y}_{<t})$$

# Machine Translation
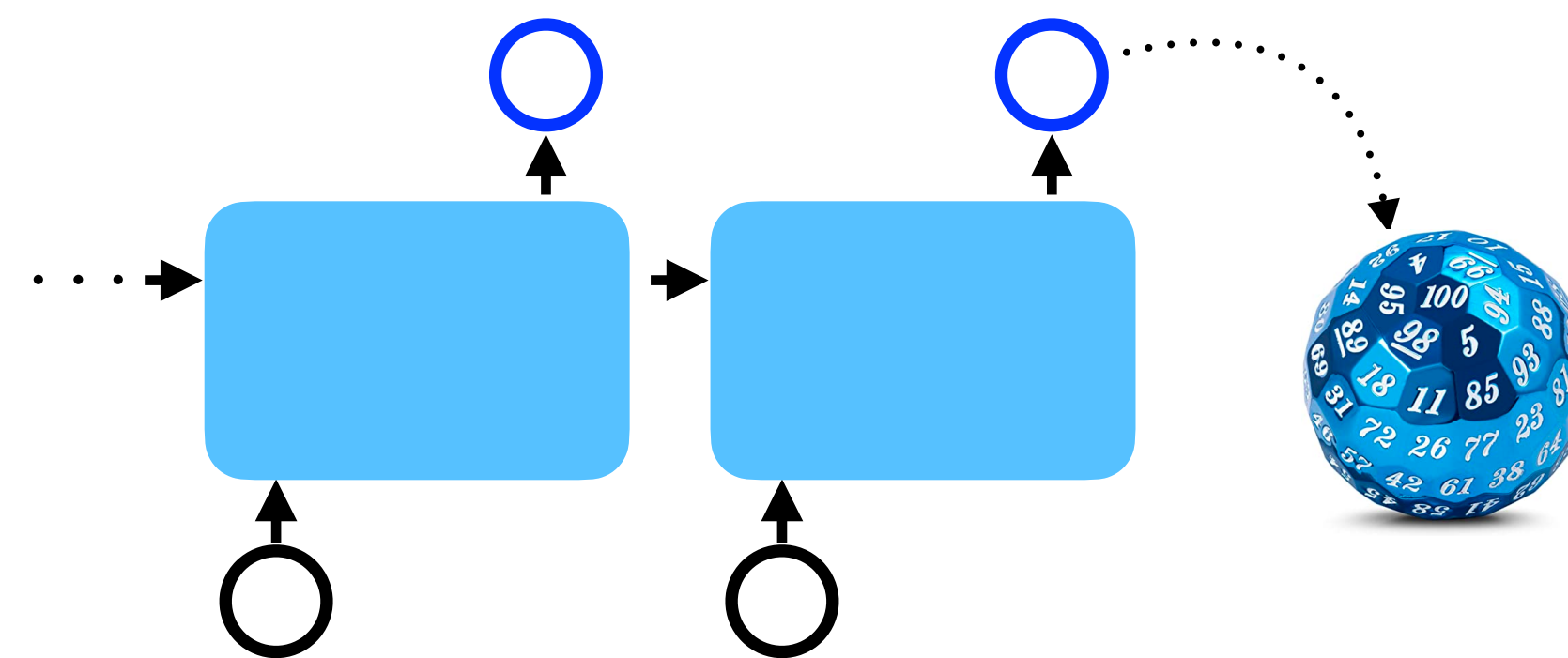
中 秋 快 樂 ！

$$x$$

Happy mid autumn festival !

$$y$$

Happy mid autumn festival !

$$p(\boldsymbol{y}) = p(y_1 \ldots y_n) = \prod_{t=1}^{n} p(y_t \mid \boldsymbol{y}_{<t})$$

$$p(y_t \mid \boldsymbol{y}_{<t})$$

# Machine Translation

中 秋 快 樂 !

$$\boldsymbol{x}$$

Happy mid autumn festival !

$$\boldsymbol{y}$$

$$p(\boldsymbol{y} \mid \boldsymbol{x}) = p(y_1 \ldots y_n \mid x_1 \ldots x_m) = \prod_{t=1}^{n} p(y_t \mid \boldsymbol{y}_{<t}, \textcolor{red}{\boldsymbol{x}})$$

target  source

Conditional Language Model

$$p(y_t \mid \boldsymbol{y}_{<t}, \textcolor{red}{\boldsymbol{x}})$$

$$p(y_t \mid \boldsymbol{y}_{<t})$$
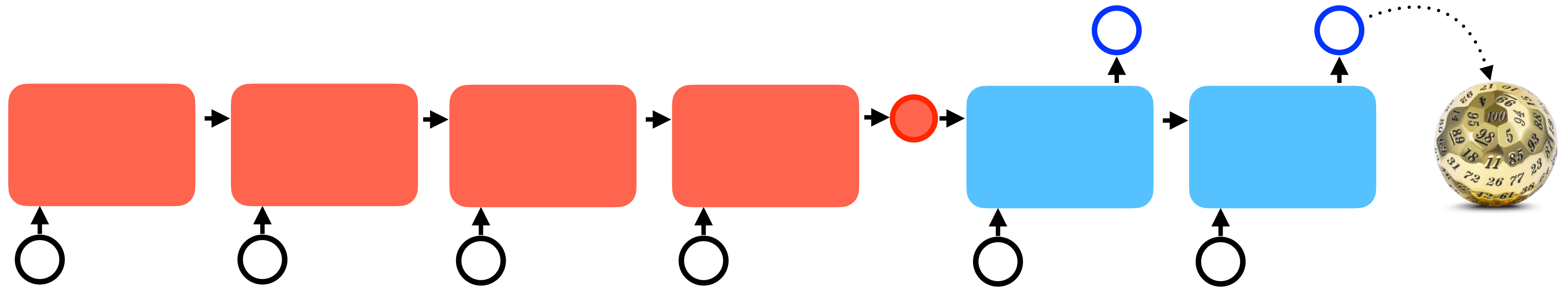
# Recurrent Neural Network

# Encoder + Decoder



$p(y_t \mid \boldsymbol{y}_{<t}, \boldsymbol{x})$

# Sequence to Sequence Model



$p(y_t \mid \boldsymbol{y}_{<t}, \boldsymbol{x})$

Encoder

Decoder

# Sequence to Sequence Model



Encoder

Decoder

# Vanishing Gradient in RNNs



$$\boldsymbol{h} = f(\boldsymbol{z})$$

$$\frac{\partial s}{\partial \boldsymbol{z}} = \frac{\partial s}{\partial \boldsymbol{h}} \frac{\partial \boldsymbol{h}}{\partial \boldsymbol{z}} \qquad \frac{\partial \boldsymbol{h}}{\partial \boldsymbol{z}} \qquad \frac{\partial s}{\partial \boldsymbol{h}}$$

Gradient Flow Direction

In general, the longer the path, the smaller the gradient signal.

# Alignment in Machine Translation



Parallel Corpus

$$p(\boldsymbol{y} \mid \boldsymbol{x})$$

target  source

# Alignment in Machine Translation



Some words might have no "counter-part".

Alignment can be many-to-one (or one-to-many).

(Brown et al, 1993)

# Sequence to Sequence Model



Happy

中　秋　快　樂　　<s>　Happy

mid

$$p(y_t \mid \boldsymbol{y}_{<t}, \boldsymbol{x})$$

# Sequence to Sequence Model



$$p(y_t \mid \boldsymbol{y}_{<t}, \boldsymbol{x})$$

中 ⟶ mid

Happy

# Attention Mechanism

Use direct connection to the encoder to <u>focus on (attend to)</u> a particular part of the source sequence.

Where do I want to look at now?

中　秋　快　樂　<s>　Happy

# Attention Mechanism

Use direct connection to the encoder to <u>focus on (attend to)</u> a particular part of the source sequence.

I think it's here.

中　秋　快　樂　<s>　Happy

# Attention Mechanism

Use direct connection to the encoder to <u>focus on (attend to)</u> a particular part of the source sequence.

# Attention Mechanism

Use direct connection to the encoder to <u>focus on (attend to)</u> a particular part of the source sequence.



$$\mathrm{softmax}(\mathbf{w}(\ \bigcirc\ \bigcirc\ ) + \mathbf{b})$$

中　秋　快　樂　<s>　Happy

# Memory Abstraction

Memory (keys)                                    Query

Task: Finding the most "<u>relevant</u>" item in the memory.

$\mathbf{q} \cdot \mathbf{k}_1$

$\mathbf{q} \cdot \mathbf{k}_2$

$\text{argmax}(\mathbf{q} \cdot \mathbf{k}_i)$        context vector   $\mathbf{c}$

$\mathbf{q} \cdot \mathbf{k}_3$

$\mathbf{q} \cdot \mathbf{k}_4$

# Dot-Product-Softmax Attention



Memory (keys)

Query

Task: Finding the most "<u>relevant</u>" item in the memory.

$\mathbf{q} \cdot \mathbf{k}_1$

$\mathbf{q} \cdot \mathbf{k}_2$

$\mathbf{q} \cdot \mathbf{k}_3$

$\mathbf{q} \cdot \mathbf{k}_4$

$\mathrm{softmax}\left( \begin{array}{c} \mathbf{q} \cdot \mathbf{k}_1 \\ \mathbf{q} \cdot \mathbf{k}_2 \\ \mathbf{q} \cdot \mathbf{k}_3 \\ \mathbf{q} \cdot \mathbf{k}_4 \end{array} \right) \rightarrow \begin{bmatrix} 0.6 \\ 0.1 \\ 0.2 \\ 0.1 \end{bmatrix} \begin{bmatrix} \circ \\ \circ \\ \circ \\ \circ \end{bmatrix} \rightarrow$

$0.6 \bigcirc + 0.1 \bigcirc + 0.2 \bigcirc + 0.1 \bigcirc$

$= \bigcirc$ context vector $\mathbf{c}$

# Attention Mechanism

$$\text{softmax}\left(\begin{array}{c} \mathbf{q} \cdot \mathbf{k}_1 \\[1em] \mathbf{q} \cdot \mathbf{k}_2 \\[1em] \mathbf{q} \cdot \mathbf{k}_3 \\[1em] \mathbf{q} \cdot \mathbf{k}_4 \end{array}\right)$$



中　秋　快　樂　&lt;s&gt;　Happy

# Attention Mechanism

$$\begin{bmatrix} 0.6 \\ 0.1 \\ 0.2 \\ 0.1 \end{bmatrix} \begin{bmatrix} \bigcirc \\ \bigcirc \\ \bigcirc \\ \bigcirc \end{bmatrix} \rightarrow 0.6\, \bigcirc \;+ 0.1\, \bigcirc \;+ 0.2\, \bigcirc \;+ 0.1\, \bigcirc$$

$$= \bigcirc \quad \text{context vector} \quad \mathbf{c}$$



中　　秋　　快　　樂　　<s>　　Happy

# Attention Mechanism



$$\text{softmax}(\mathbf{w}(\;\bigcirc\;\bigcirc\;) + \mathbf{b})$$

中　秋　快　樂　&lt;s&gt;　Happy

# Attention Mechanism



$$\mathrm{softmax}(\mathbf{w}(\ \bigcirc\ \bigcirc\ ) + \mathbf{b})$$

中　　秋　　快　　樂　　\<s\>　　Happy　　mid

# Dot-Product-Softmax Attention



Memory (key-value pairs)

Query

$$\text{softmax}\left(\begin{array}{c} \mathbf{q} \cdot \mathbf{k}_1 \\ \mathbf{q} \cdot \mathbf{k}_2 \\ \mathbf{q} \cdot \mathbf{k}_3 \\ \mathbf{q} \cdot \mathbf{k}_4 \end{array}\right) \rightarrow \begin{bmatrix} 0.6 \\ 0.1 \\ 0.2 \\ 0.1 \end{bmatrix}$$

$\mathbf{q} \cdot \mathbf{k}_1$

$\mathbf{q} \cdot \mathbf{k}_2$

$\mathbf{q} \cdot \mathbf{k}_3$

$\mathbf{q} \cdot \mathbf{k}_4$

context vector $\mathbf{c}$

# Dot-Product-Softmax Attention

similarity

normalized similarity

$$\sum_{m=1}^{M} \frac{\exp\left(\boldsymbol{q}_n \boldsymbol{k}_m\right)}{\sum_{m'=1}^{M} \exp\left(\boldsymbol{q}_n \boldsymbol{k}_{m'}\right)} \boldsymbol{v}_m^{\top} = \mathbf{V}^{\top} \mathrm{softmax}(\mathbf{K}\boldsymbol{q}_n)$$

weighted sum

# Computational Complexity

# Attention Mechanism

It helps with vanishing gradient problem.

# Attention Mechanism

It offers some interpretability.

# Generation as (Conditional) Language Modeling

$$p(\boldsymbol{y} \mid \boldsymbol{x}) = p(y_1 \ldots y_n \mid x_1 \ldots x_m) = \prod_{t=1}^{n} p(y_t \mid \boldsymbol{y}_{<t}, \textcolor{red}{\boldsymbol{x}})$$

target source

Conditional Language Model



$p(y_t \mid \boldsymbol{y}_{<t}, \textcolor{red}{\boldsymbol{x}})$

# Generation as (Conditional) Language Modeling

$$p(\boldsymbol{y} \mid \boldsymbol{x}) = p(y_1 \dots y_n \mid x_1 \dots x_m) = \prod_{t=1}^{n} p(y_t \mid \boldsymbol{y}_{<t}, \textcolor{red}{\boldsymbol{x}})$$

# Continue Training / Post-Training / Instruction Tuning



Instruction finetuning

Please answer the following question.
What is the boiling point of Nitrogen?

Chain-of-thought finetuning

Answer the following question by reasoning step-by-step.

The cafeteria had 23 apples. If they used 20 for lunch and bought 6 more, how many apples do they have?

-320.4F

The cafeteria had 23 apples originally. They used 20 to make lunch. So they had 23 - 20 = 3. They bought 6 more apples, so they have 3 + 6 = 9.

Language model

*Multi-task instruction finetuning* **(1.8K tasks)**

*Inference: generalization to unseen tasks*

Q: Can Geoffrey Hinton have a conversation with George Washington?

Give the rationale before answering.

Geoffrey Hinton is a British-Canadian computer scientist born in 1947. George Washington died in 1799. Thus, they could not have had a conversation together. So the answer is "no".

(Hyung Won Chung et al, 2022)

# How to perform decoding?



$$p(y_t \mid \boldsymbol{y}_{<t}, \boldsymbol{x})$$

arg max
beam search

# Beam Search

# Beam Search or Pure Sampling?

**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 ..."

**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."

degenerate repetition                                    incoherent gibberish

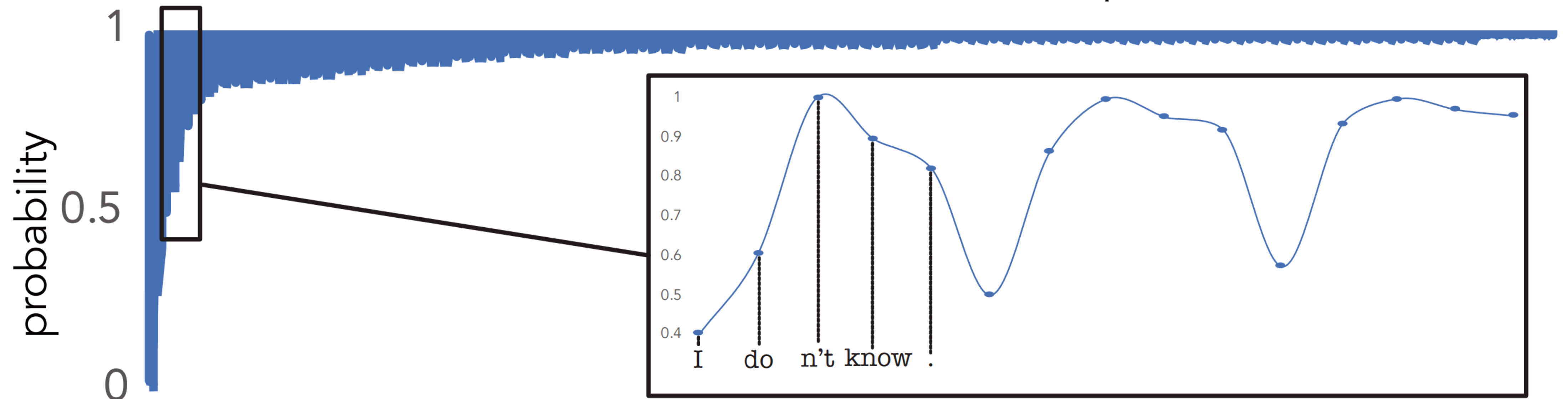GPT-2 Large (774M parameters)

Holtzman et al., 2020

# What happened?



Token Probabilities for "I don't know." Repeated 200 times

The probability of a repeated phrase increases with each repetition, creating a feedback loop.

Holtzman et al., 2020

# Top-K Sampling

He wanted to go to the ➤ [NLG Model] ➤

restroom
grocery
store
airport
bathroom
beach
doctor
hospital
pub
gym

Fan et al., 2018; Holtzman et al., 2020

# Top-K Sampling



Holtzman et al., 2020

# Top-p (nucleus) Sampling

To cut off by the cumulative probability mass, rather than the first K terms.



Holtzman et al., 2020