# CS 581

## *Advanced Artificial Intelligence*

### April 3, 2024

# Announcements / Reminders

- **Please follow the Week 11/12 To Do List instructions (if you haven't already)**

- **Programming Assignment #02 due on Sunday (04/07) at 11:59 PM CST**
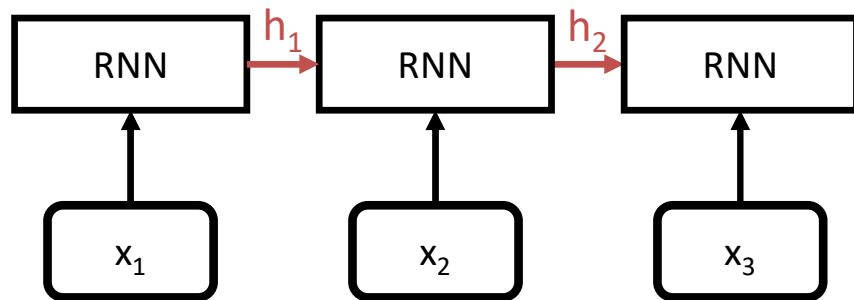
# Plan for Today

- **Attention Mechanism**

- **Transformer Basics**

- **Generative AI Models: Introduction**

# Sequence to Sequence Networks (seq2seq)
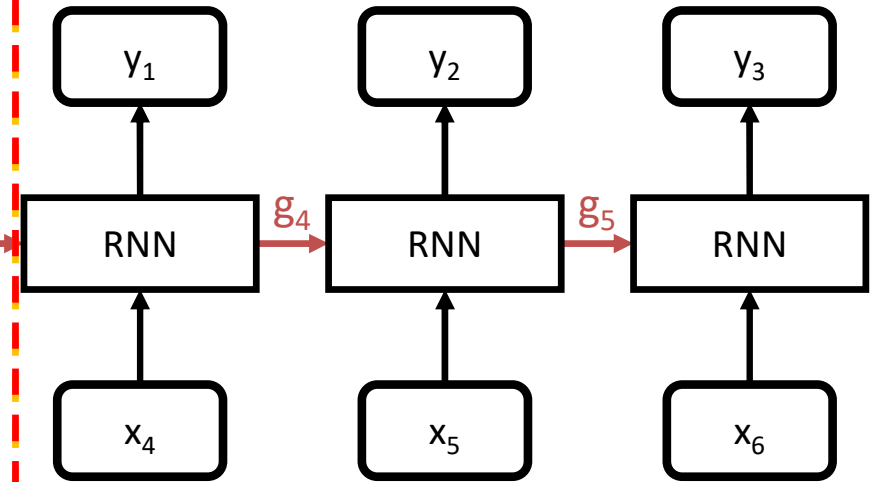# With Attention

# RNN Encoder-Decoder Architecture

$h_i$ – encoder hidden state at time $t_i$
$g_i$ – decoder hidden state at time $t_i$
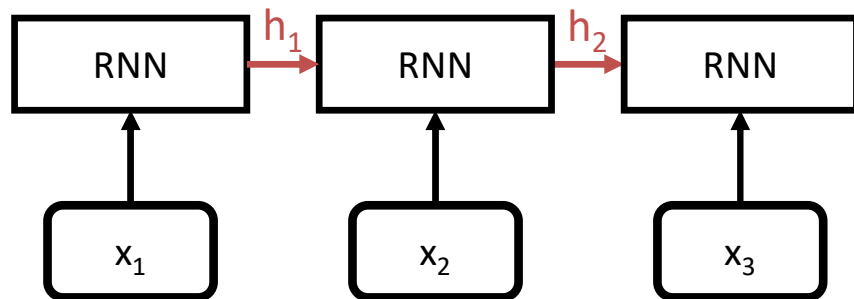
CONTEXT

**RNN Encoder**

**RNN Decoder**

| $y_1$ | $y_2$ | $y_3$ |

$h_3$

| RNN | $g_4$ | RNN | $g_5$ | RNN |

| $x_4$ | $x_5$ | $x_6$ |

| RNN | $h_1$ | RNN | $h_2$ | RNN | $h_3$ |

| $x_1$ | $x_2$ | $x_3$ |

$h_3$

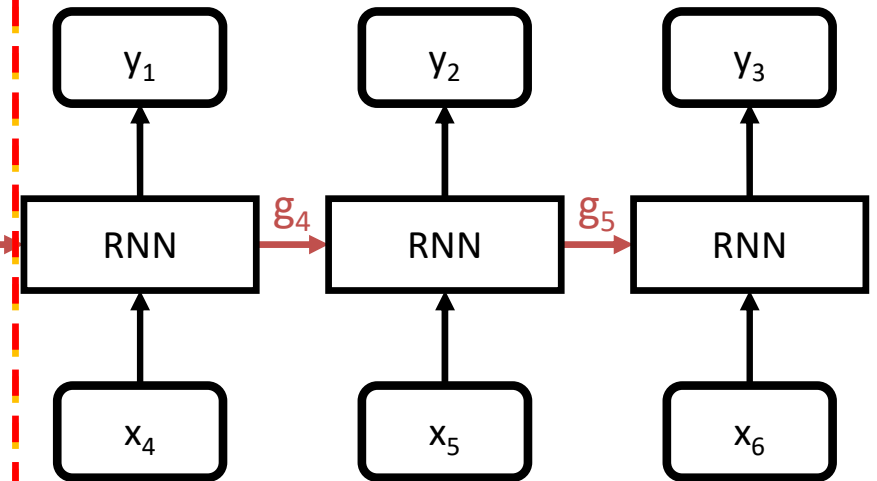**Encoded (entire) input sequence.
(fixed size vector)**

**TIME**

# RNN Encoder-Decoder: Context



$h_i$ – encoder hidden state at time $t_i$
$g_i$ – decoder hidden state at time $t_i$

CONTEXT

RNN Encoder

RNN Decoder

$y_1$ $y_2$ $y_3$

RNN $\xrightarrow{g_4}$ RNN $\xrightarrow{g_5}$ RNN

$x_4$ $x_5$ $x_6$

$h_3$

RNN $\xrightarrow{h_1}$ RNN $\xrightarrow{h_2}$ RNN $\xrightarrow{h_3}$
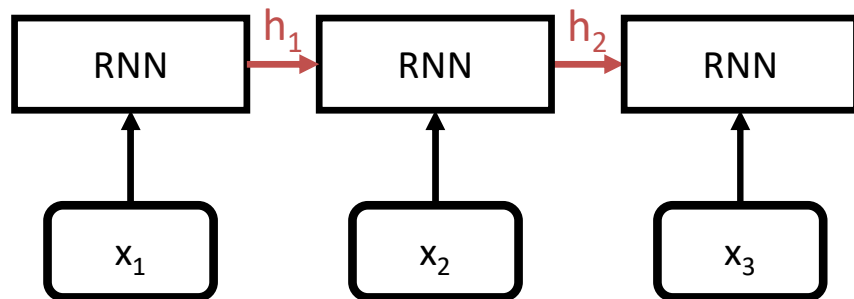
$x_1$ $x_2$ $x_3$

Encoded (entire) input sequence.
(fixed size vector)
BOTTLENECK!

TIME

# Fixed Length Context

$h_i$ − encoder hidden state at time $t_i$
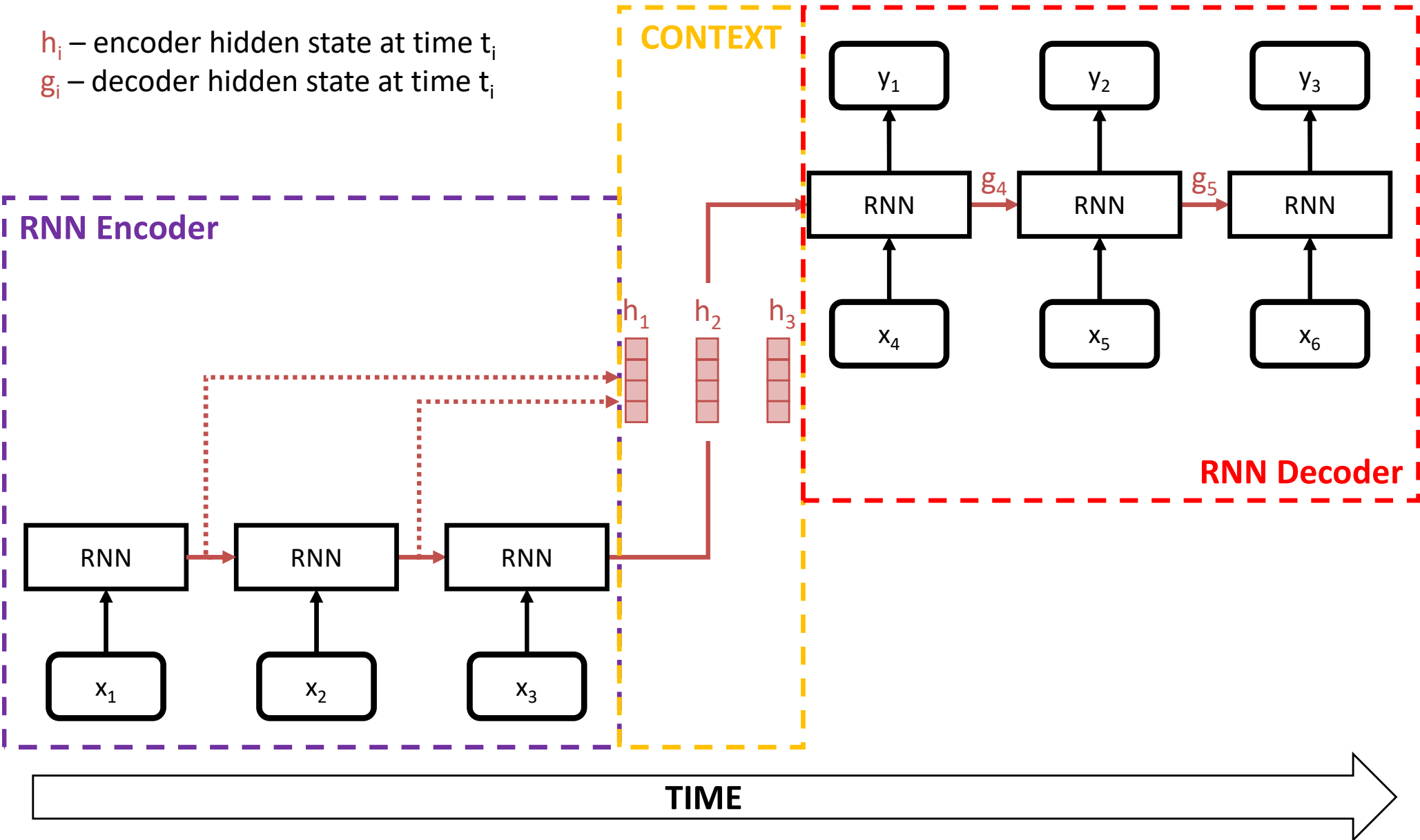$g_i$ − decoder hidden state at time $t_i$



**RNN Encoder**

**CONTEXT**

**RNN Decoder**

**TIME**

# RNN Encoder-Decoder Architecture



$h_i$ − encoder hidden state at time $t_i$
$g_i$ − decoder hidden state at time $t_i$

CONTEXT

RNN Encoder

RNN Decoder

$h_1$  $h_2$  $h_3$

$g_4$  $g_5$

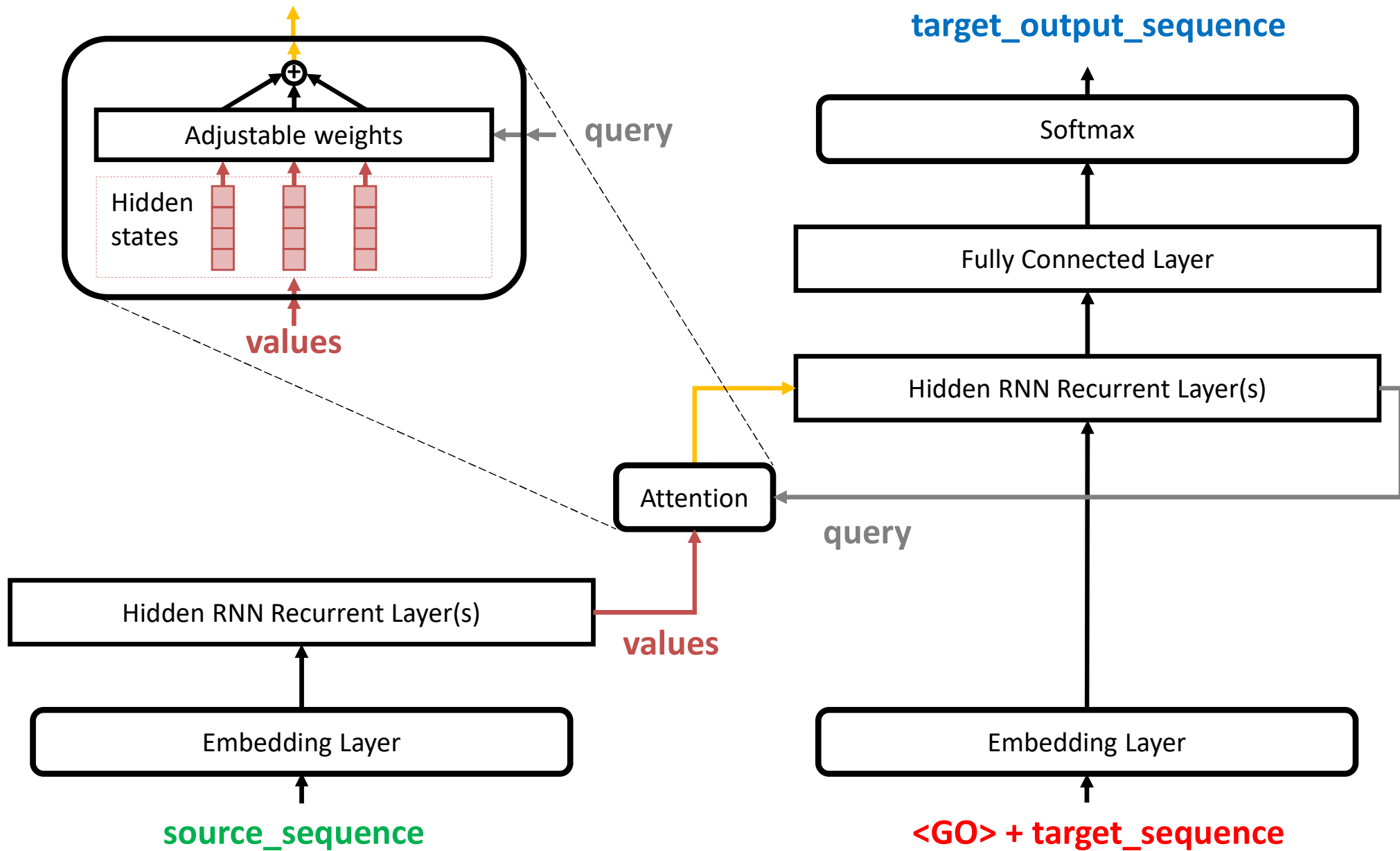$y_1$  $y_2$  $y_3$

$x_1$  $x_2$  $x_3$

$x_4$  $x_5$  $x_6$

TIME

# Attention Mechanism

- **Given a set of vector values, and a vector query, attention is a technique to compute a weighted sum of the values, dependent on the query**

- **Attention mechanism "amplifies" important aspects of the signal from the encoder based on the decoder query**

- **In seq2seq models with attention, each decoder hidden state (query) attends to all the encoder hidden states (values)**
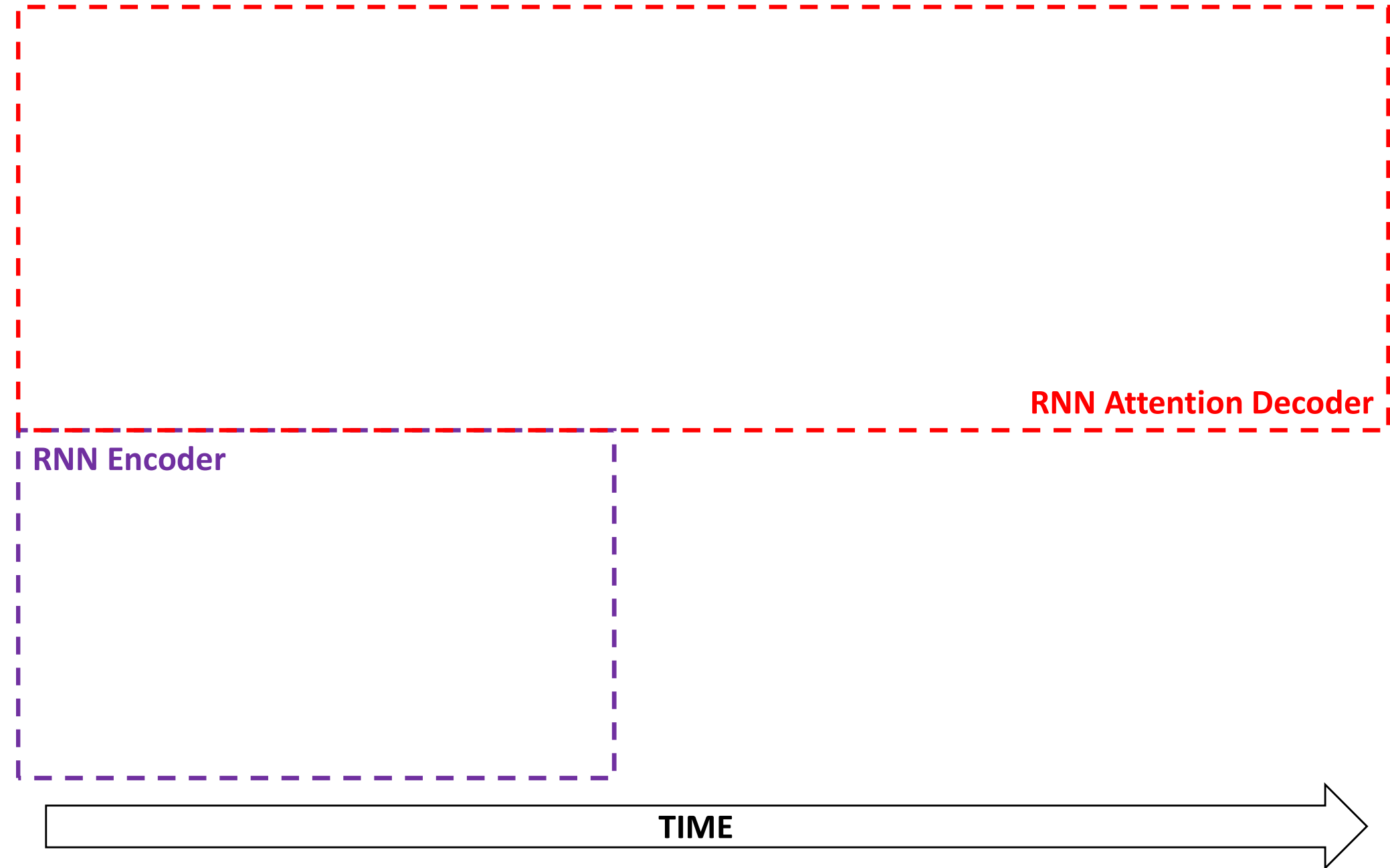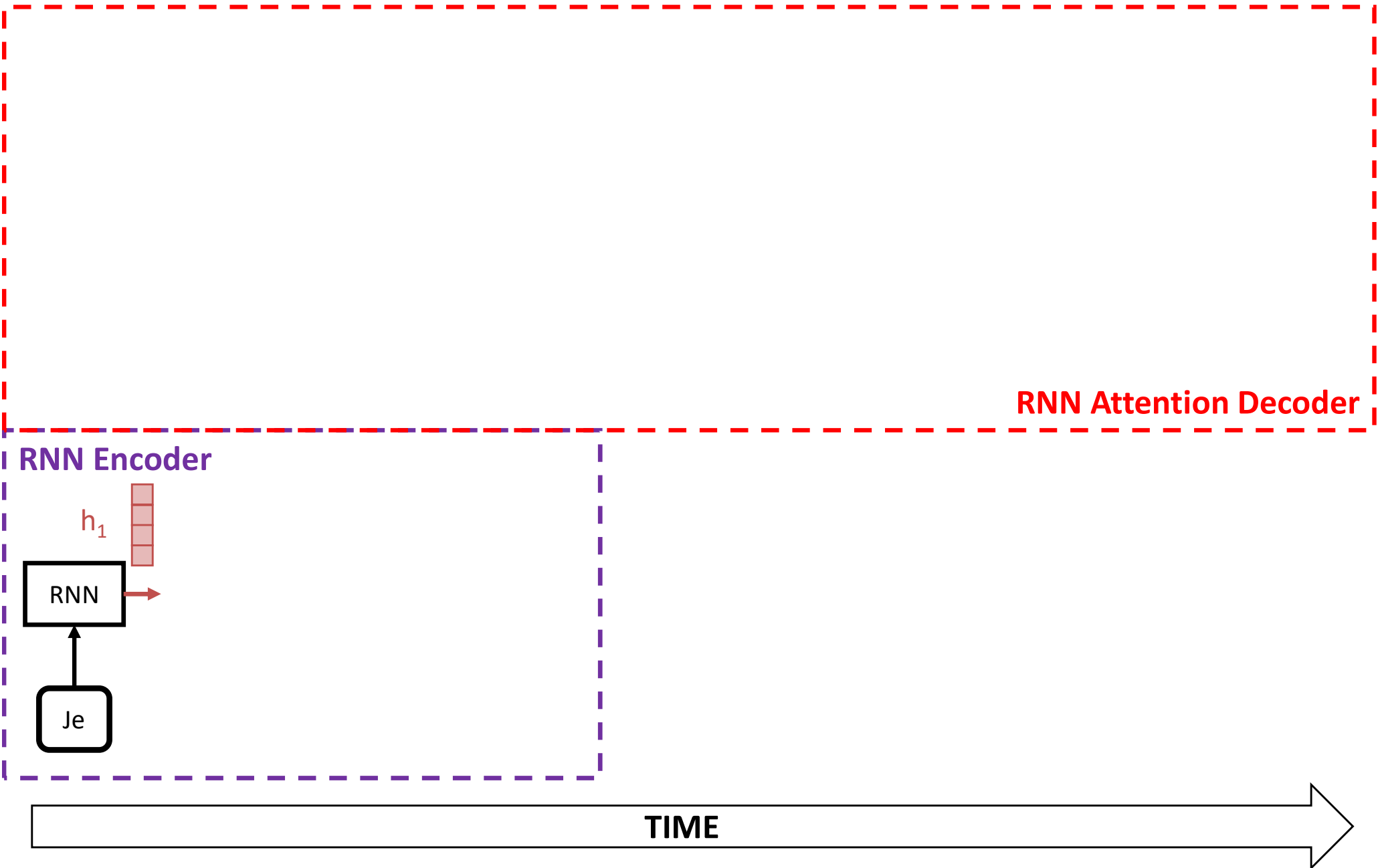
# RNN Encoder-Decoder with Attention

target_output_sequence

Softmax

Fully Connected Layer

Hidden RNN Recurrent Layer(s)

Attention

query

Hidden RNN Recurrent Layer(s)

values

Embedding Layer

Embedding Layer

source_sequence
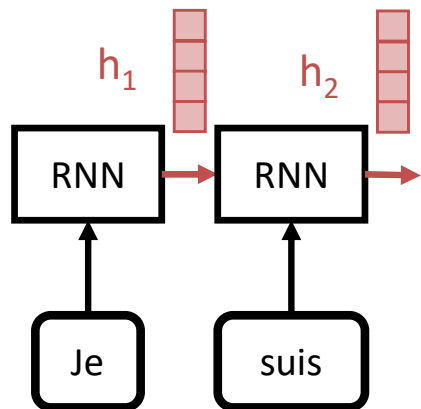
\<GO\> + target_sequence

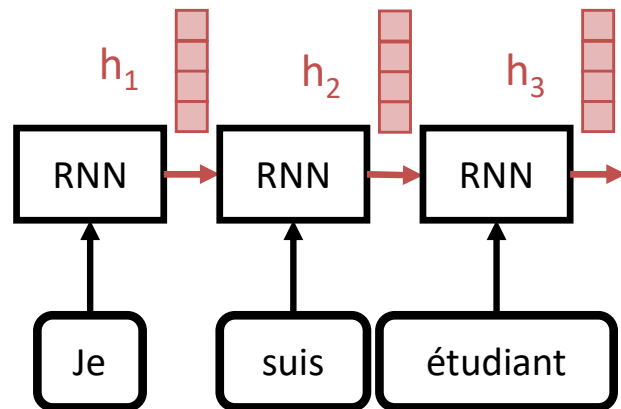# RNN Encoder-Decoder with Attention

# RNN Encoder-Decoder with Attention

**RNN Attention Decoder**

**RNN Encoder**

TIME

# RNN Encoder-Decoder with Attention
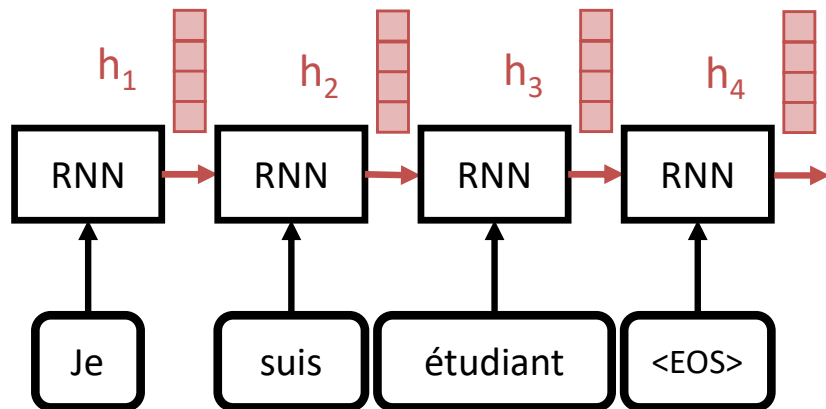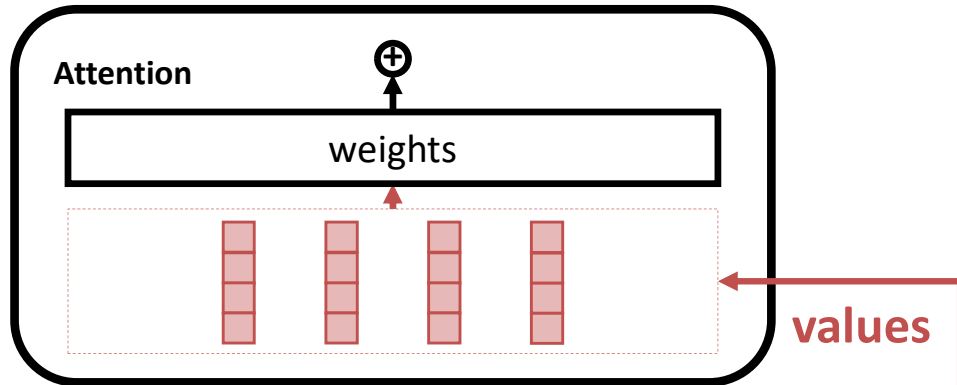
RNN Attention Decoder

RNN Encoder

$h_1$

RNN

Je

TIME

# RNN Encoder-Decoder with Attention

**RNN Attention Decoder**

**RNN Encoder**

$h_1$    $h_2$

| RNN | → | RNN | → |

↑         ↑

Je        suis

**TIME**

# RNN Encoder-Decoder with Attention

**RNN Encoder**

$h_1$   $h_2$   $h_3$

| RNN | → | RNN | → | RNN | → |

↑         ↑         ↑

Je      suis    étudiant

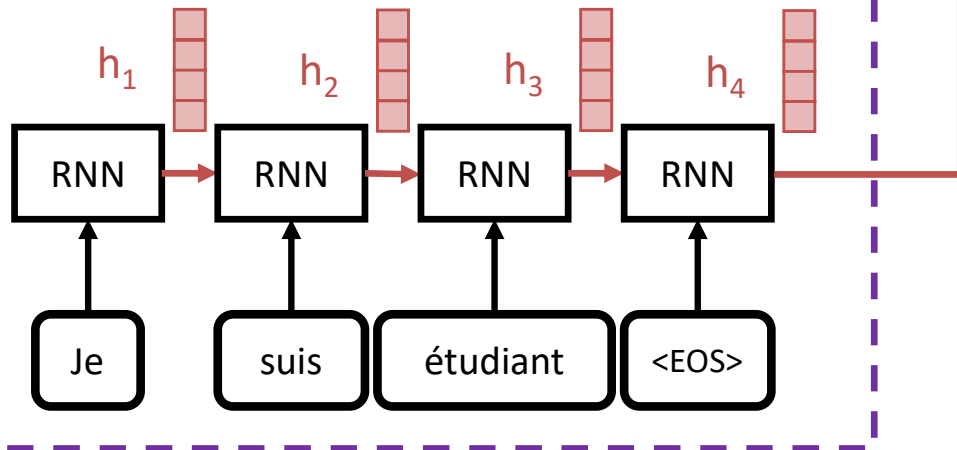**TIME**

# RNN Encoder-Decoder with Attention



RNN Attention Decoder

RNN Encoder

$h_1$ $h_2$ $h_3$ $h_4$

| RNN | RNN | RNN | RNN |

| Je | suis | étudiant | <EOS> |

TIME

# RNN Encoder-Decoder with Attention

# RNN Encoder-Decoder with Attention

# RNN Encoder-Decoder with Attention

**Attention**

weights

$\oplus$

context for $t_5$

I

RNN

<GO>

**RNN Attention Decoder**

values

**RNN Encoder**

$h_1$   $h_2$   $h_3$   $h_4$

RNN → RNN → RNN → RNN

Je   suis   étudiant   <EOS>

**TIME**

# RNN Encoder-Decoder with Attention

# RNN Encoder-Decoder with Attention

# RNN Encoder-Decoder with Attention

# RNN Encoder-Decoder with Attention

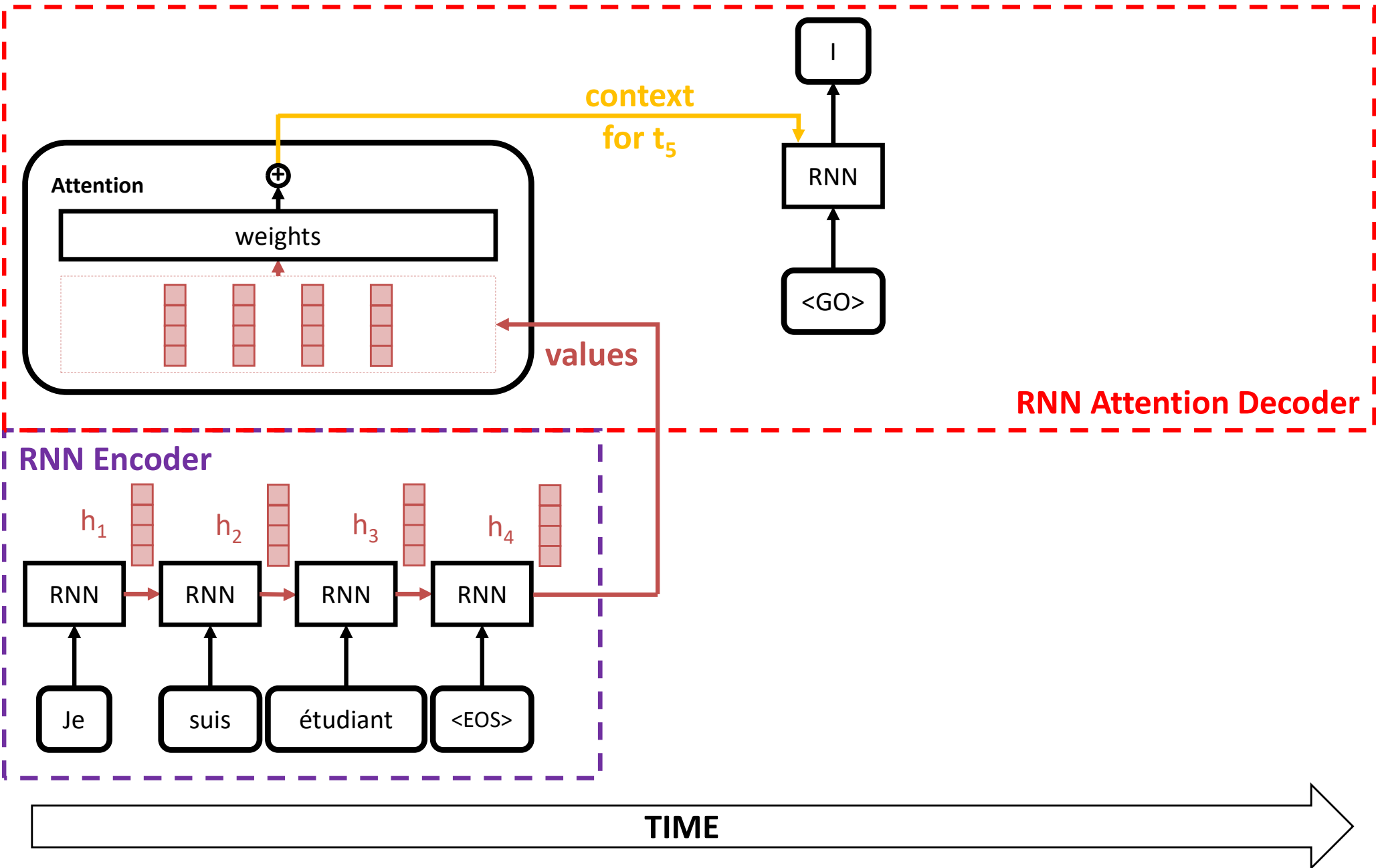# RNN Encoder-Decoder with Attention
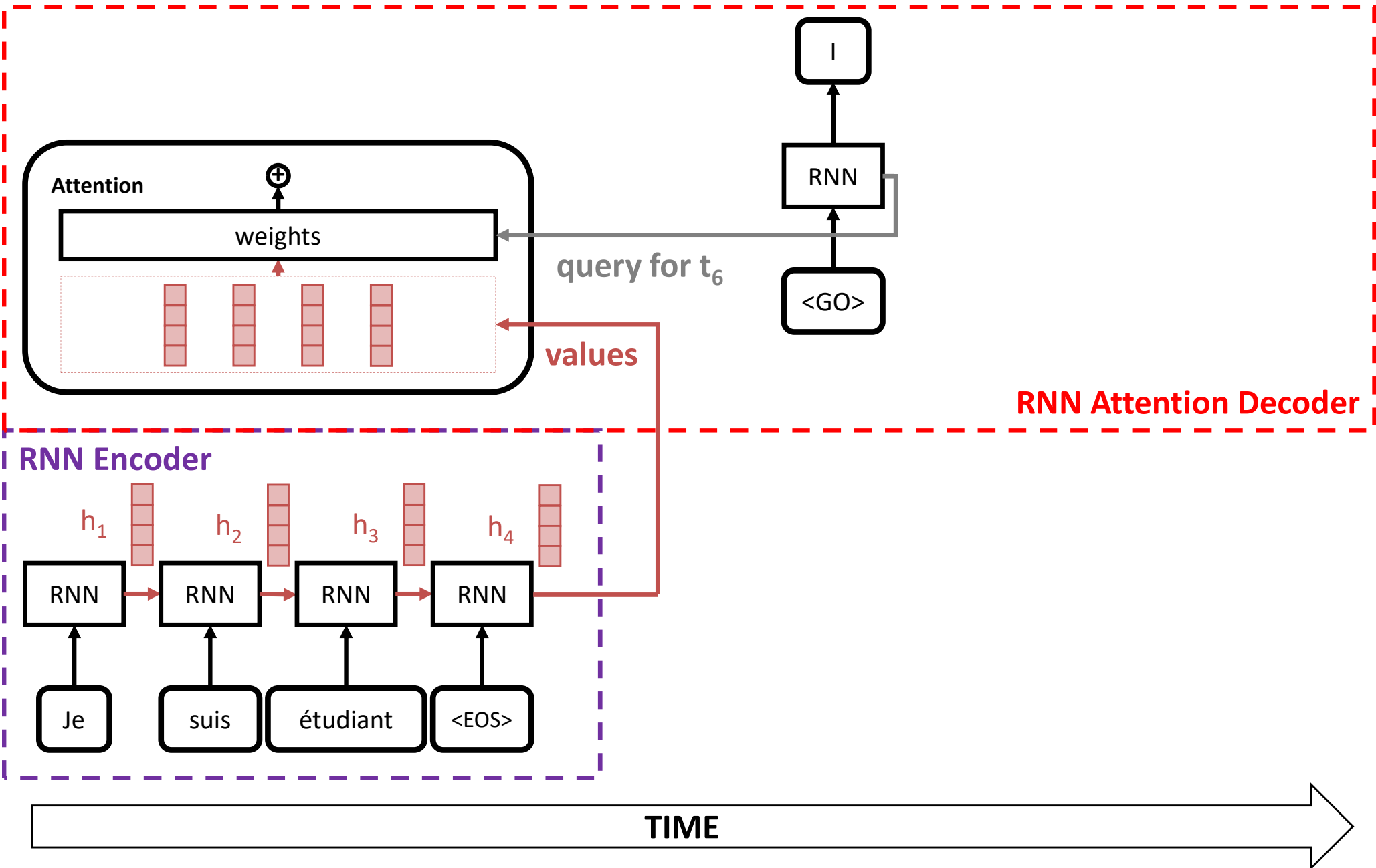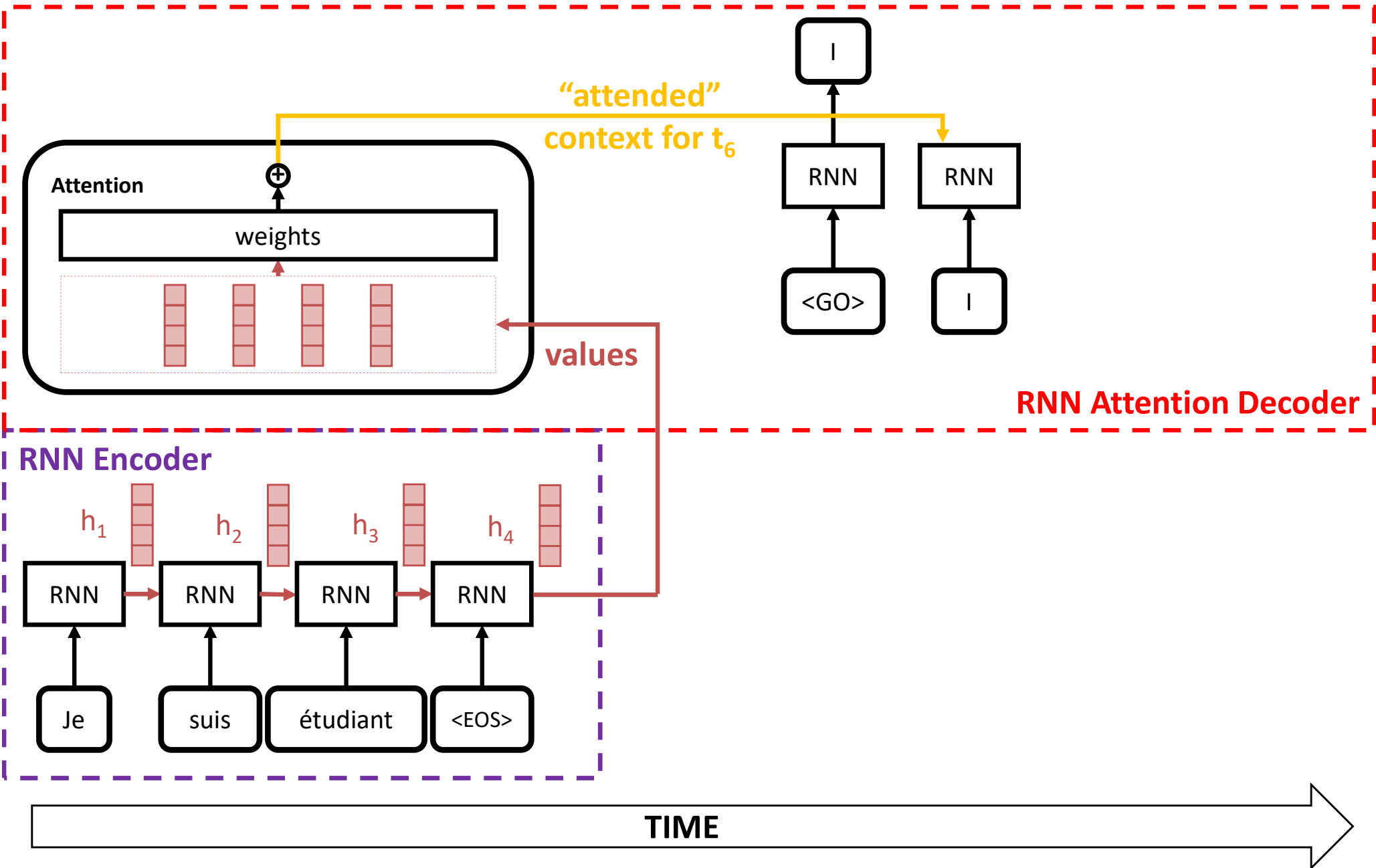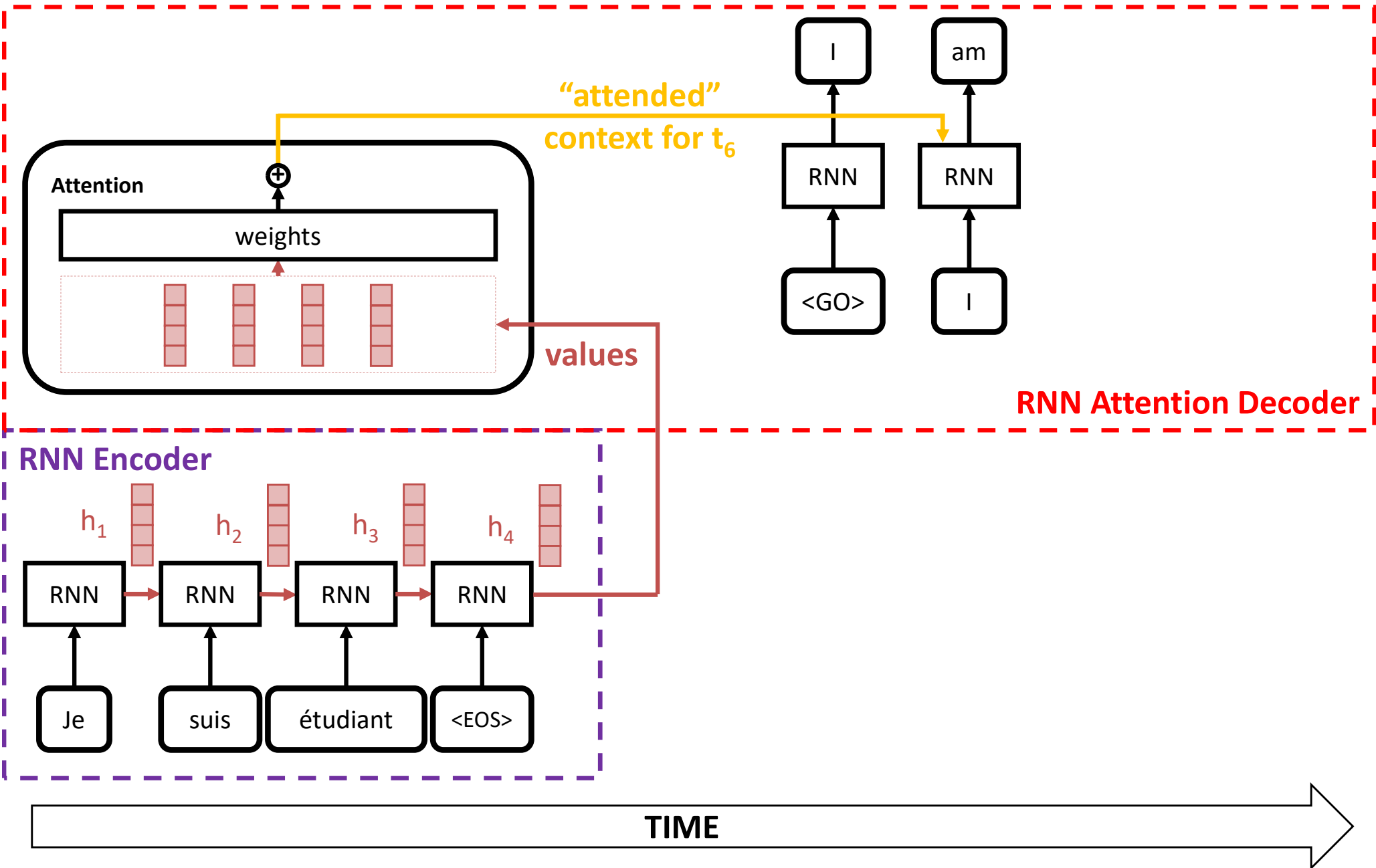
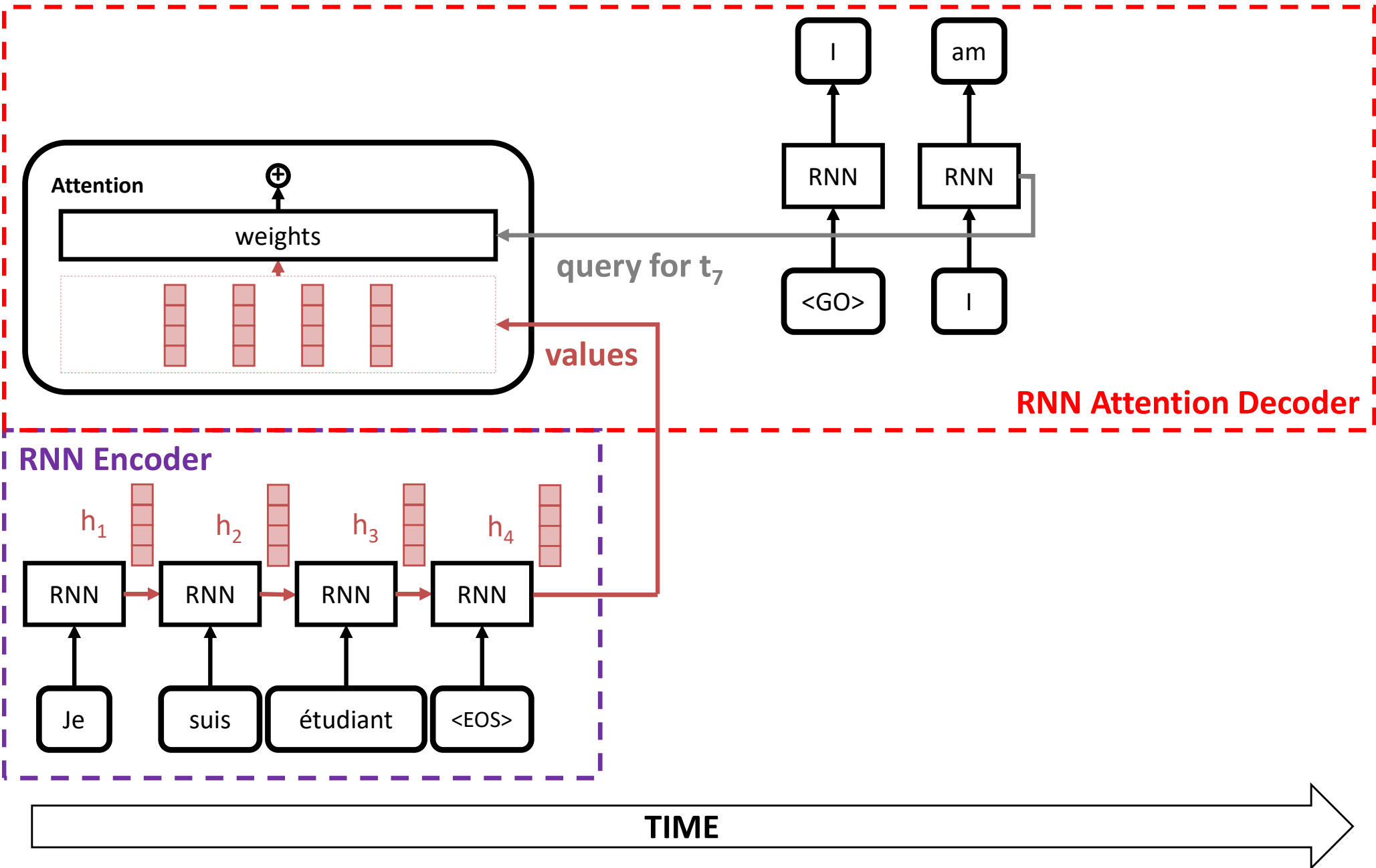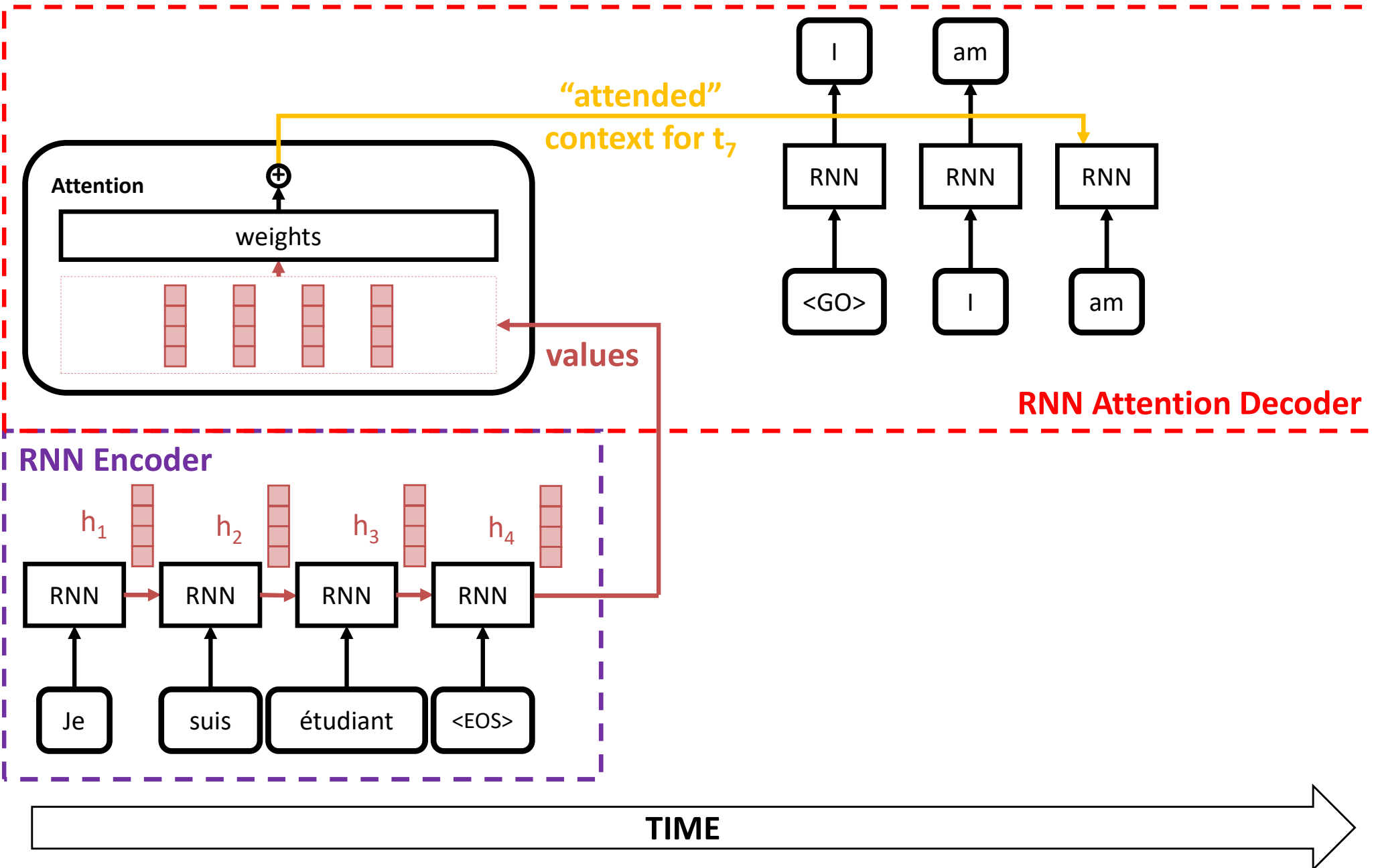# RNN Encoder-Decoder with Attention

# RNN Encoder-Decoder with Attention

# RNN Encoder-Decoder with Attention

# RNN Encoder-Decoder with Attention

# RNN Encoder-Decoder with Attention

# What is inside the RNN Decoder Attention?

# RNN Encoder-Decoder Attention

# RNN Encoder-Decoder Attention

# RNN Encoder-Decoder Attention

# RNN Encoder-Decoder Attention

# What is RNN Decoder Attention doing?

# RNN Encoder-Decoder Attention

**Attention**

1. **Prepare inputs**
a) **encoder hidden states**

$h_1$ | $h_{1,1}$ $h_{1,2}$ $h_{1,3}$ $h_{1,4}$

$h_2$ | $h_{2,1}$ $h_{2,2}$ $h_{2,3}$ $h_{2,4}$

$h_3$ | $h_{3,1}$ $h_{3,2}$ $h_{3,3}$ $h_{3,4}$

$h_4$ | $h_{4,1}$ $h_{4,2}$ $h_{4,3}$ $h_{4,4}$

**values**

# RNN Encoder-Decoder Attention

Attention

1. **Prepare inputs**
   a) **encoder hidden states**
   b) **query** for (next) time $t_5$

| | |
|---|---|
| | $q_1$ |
| | $q_2$ | ← **query** |
| | $q_3$ |
| query q | $q_4$ |

| $h_1$ | $h_{1,1}$ | $h_2$ | $h_{2,1}$ | $h_3$ | $h_{3,1}$ | $h_4$ | $h_{4,1}$ |
|---|---|---|---|---|---|---|---|
| | $h_{1,2}$ | | $h_{2,2}$ | | $h_{3,2}$ | | $h_{4,2}$ | ← **values** |
| | $h_{1,3}$ | | $h_{2,3}$ | | $h_{3,3}$ | | $h_{4,3}$ |
| | $h_{1,4}$ | | $h_{2,4}$ | | $h_{3,4}$ | | $h_{4,4}$ |

# RNN Encoder-Decoder Attention



Attention

**2. Score** hidden states
- measure similarity
between **q** and $h_i$

query

query q

values

# RNN Encoder-Decoder Attention



**2. Score hidden states**
- measure alignment between q and $h_i$

query

query q

values

# RNN Encoder-Decoder Attention

# RNN Encoder-Decoder Attention



**Attention**

2. **Score** hidden states - calculate **attention weights** for **h**$_i$

query

query q

values

$q_1$

$q_2$

$q_3$

$q_4$

score score score score

$h_{1,1}$ $h_{2,1}$ $h_{3,1}$ $h_{4,1}$

$h_{1,2}$ $h_{2,2}$ $h_{3,2}$ $h_{4,2}$

$h_{1,3}$ $h_{2,3}$ $h_{3,3}$ $h_{4,3}$

$h_1$ $h_{1,4}$ $h_2$ $h_{2,4}$ $h_3$ $h_{3,4}$ $h_4$ $h_{4,4}$

# RNN Encoder-Decoder Attention



**Attention**

3. Softmax **alignment scores** for $h_i$

softmax (for **alignment scores** only)

score | score | score | score

$q_1$
$q_2$ — query
$q_3$
query q | $q_4$

values

$h_{1,1}$ | $h_{2,1}$ | $h_{3,1}$ | $h_{4,1}$
$h_{1,2}$ | $h_{2,2}$ | $h_{3,2}$ | $h_{4,2}$
$h_{1,3}$ | $h_{2,3}$ | $h_{3,3}$ | $h_{4,3}$
$h_1$ $h_{1,4}$ | $h_2$ $h_{2,4}$ | $h_3$ $h_{3,4}$ | $h_4$ $h_{4,4}$

# RNN Encoder-Decoder Attention

# RNN Encoder-Decoder Attention



**Attention**

context

| $s_1 * h_{1,1}$ | $s_2 * h_{2,1}$ | $s_3 * h_{3,1}$ | $s_4 * h_{4,1}$ |
| --- | --- | --- | --- |
| $s_1 * h_{1,2}$ | $s_2 * h_{2,2}$ | $s_3 * h_{3,2}$ | $s_4 * h_{4,2}$ |
| $s_1 * h_{1,3}$ | $s_2 * h_{2,3}$ | $s_3 * h_{3,3}$ | $s_4 * h_{4,3}$ |
| $s_1 * h_{1,4}$ | $s_2 * h_{2,4}$ | $s_3 * h_{3,4}$ | $s_4 * h_{4,4}$ |

**5. Sum up weighted state vectors to produce context vector for time $t_6$**

$s_1$   $s_2$   $s_3$   $s_4$

**softmax (for alignment scores only)**

score   score   score   score

| $h_{1,1}$ | $h_{2,1}$ | $h_{3,1}$ | $h_{4,1}$ |
| --- | --- | --- | --- |
| $h_{1,2}$ | $h_{2,2}$ | $h_{3,2}$ | $h_{4,2}$ |
| $h_{1,3}$ | $h_{2,3}$ | $h_{3,3}$ | $h_{4,3}$ |
| $h_{1,4}$ | $h_{2,4}$ | $h_{3,4}$ | $h_{4,4}$ |

$h_1$   $h_2$   $h_3$   $h_4$

| $q_1$ |
| --- |
| $q_2$ |
| $q_3$ |
| $q_4$ |

query

query q

values

**Illinois Institute of Technology**

# Summary / Intuition

- **Attention mechanism creates context dynamically**
- **"Attended" values "dominate" context**
  - this allows the decoder to dynamically focus on most relevant aspects of the encoder output
- **The query represents/summarizes the "current task" of the decoder**
  - the input it is processing within the current context
- **Scoring for values estimates similarity (between a query and specific value) values:**
  - dot product
  - scaled dot product
  - etc.

# Score: Dot Product

**Attention**

**score:**

$$\boldsymbol{\alpha}(\boldsymbol{q}, \boldsymbol{h_i}) = \boldsymbol{q} \bullet \boldsymbol{h_i}$$

score

score

score

score

$h_{1,1}$
$h_{1,2}$
$h_{1,3}$
$h_{1,4}$

$h_{2,1}$
$h_{2,2}$
$h_{2,3}$
$h_{2,4}$

$h_{3,1}$
$h_{3,2}$
$h_{3,3}$
$h_{3,4}$

$h_{4,1}$
$h_{4,2}$
$h_{4,3}$
$h_{4,4}$

$h_1$

$h_2$

$h_3$

$h_4$

$q_1$
$q_2$
$q_3$
$q_4$

query

query q

values

# Score: Scaled Dot Product

**Attention**

**score:**

$$\alpha(q, h_i) = \frac{q \bullet h_i}{\sqrt{dimension_{h_i}}}$$



score

score

score

score

$h_1$ $h_{1,1}$ $h_{1,2}$ $h_{1,3}$ $h_{1,4}$

$h_2$ $h_{2,1}$ $h_{2,2}$ $h_{2,3}$ $h_{2,4}$

$h_3$ $h_{3,1}$ $h_{3,2}$ $h_{3,3}$ $h_{3,4}$

$h_4$ $h_{4,1}$ $h_{4,2}$ $h_{4,3}$ $h_{4,4}$

$q_1$ $q_2$ $q_3$ $q_4$

**query**

query q

**values**

# Softmax

**Attention**

**softmax on scores:**

$$s(score_i) = \frac{e^{score_i}}{\sum_j e^{score_j}}$$

softmax (for **alignment scores** only)

| score | score | score | score |

| $h_{1,1}$ | $h_{2,1}$ | $h_{3,1}$ | $h_{4,1}$ |
| $h_{1,2}$ | $h_{2,2}$ | $h_{3,2}$ | $h_{4,2}$ |
| $h_{1,3}$ | $h_{2,3}$ | $h_{3,3}$ | $h_{4,3}$ |
| $h_{1,4}$ | $h_{2,4}$ | $h_{3,4}$ | $h_{4,4}$ |

$h_1$ $h_2$ $h_3$ $h_4$

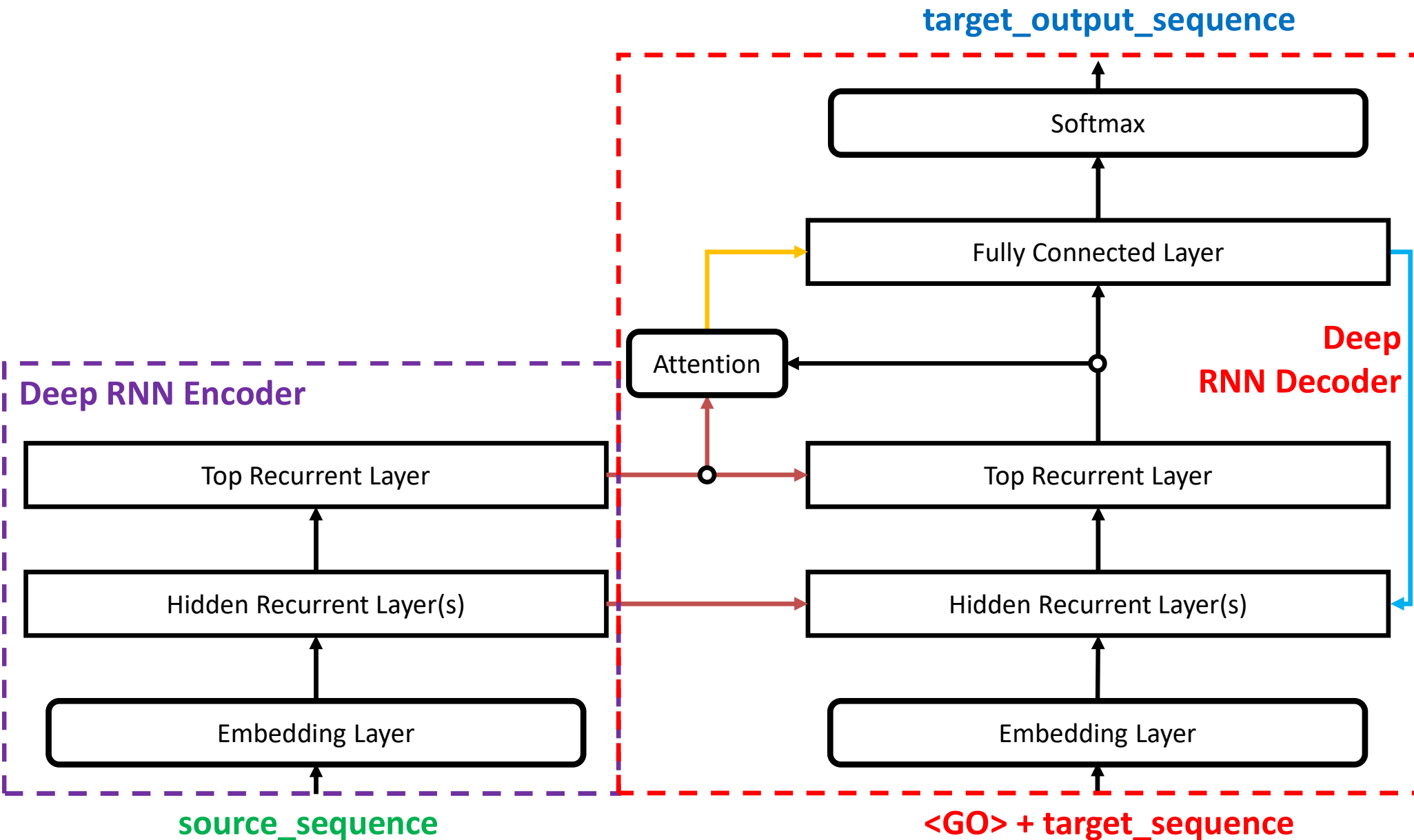$q_1$
$q_2$  **query**
$q_3$
query q  $q_4$

**values**

# Benefits of the Attention Structure

- **Significantly improves performance (in many applications)**
  - it's very useful to allow the decoder to focus on certain parts of the source
- **Solves the bottleneck issue**
  - attention allows decoder to look directly at the source (and "bypass" the bottleneck)
- **Helps with vanishing gradient problem**
  - provides shortcut to far away states
- **Provides some interpretability**
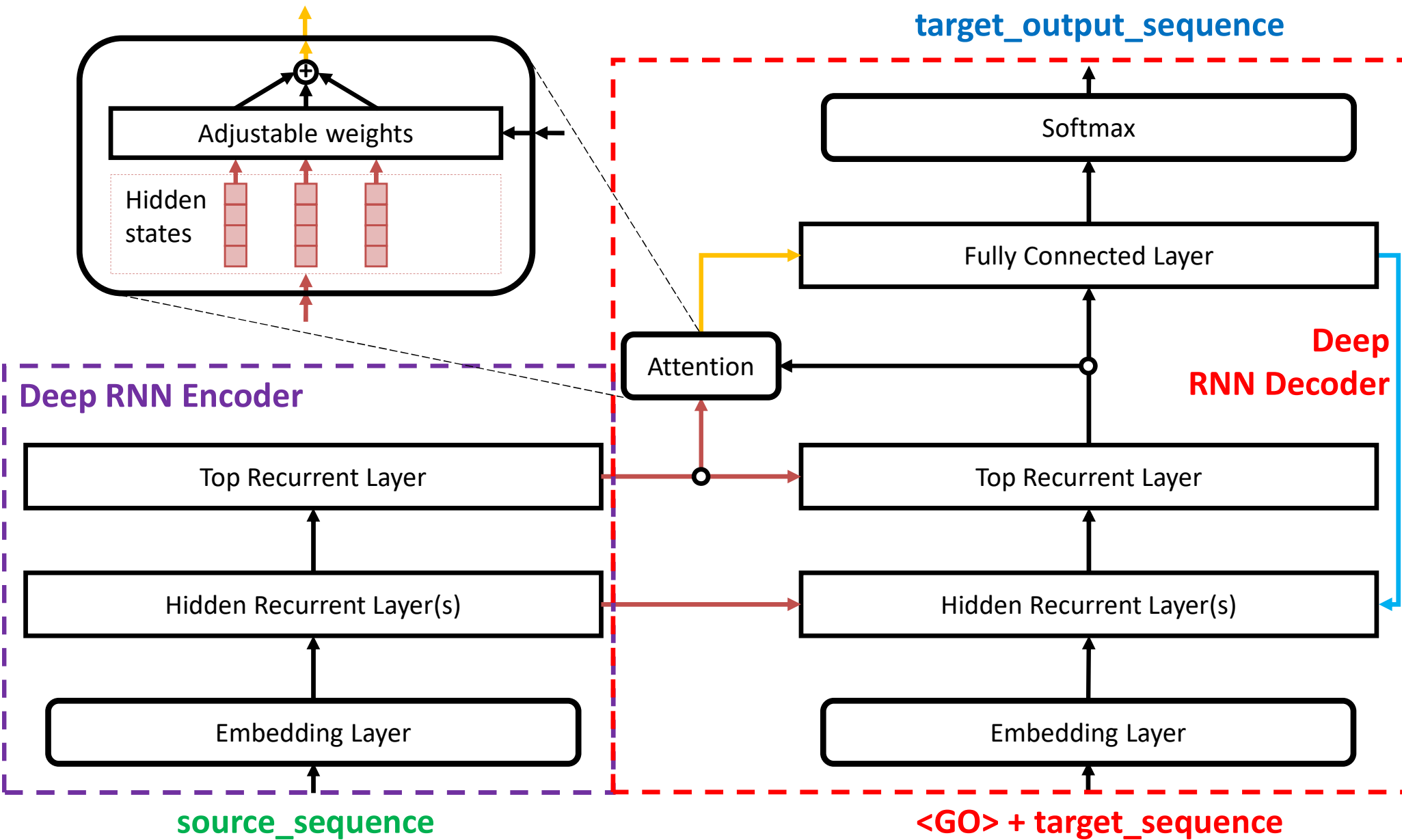  - inspecting attention distribution we can see what the decoder was focusing on

# Deep RNN Enc-Dec with Attention

target_output_sequence

| Softmax |
|---|

| Fully Connected Layer |
|---|

| Attention |
|---|

"attentional" vector

| Top Recurrent Layer | | Top Recurrent Layer |
|---|---|---|

| Hidden Recurrent Layer(s) | | Hidden Recurrent Layer(s) |
|---|---|---|

| Embedding Layer | | Embedding Layer |
|---|---|---|

source_sequence

<GO> + target_sequence

# Deep RNN Enc-Dec with Attention

# Deep RNN Enc-Dec with Attention

# Transformers - Basics

# Generative Pre-trained Transformer 3

## What is it?

Generative Pre-trained Transformer 3 (GPT-3) is an **autoregressive language model that uses deep learning to produce human-like text**. It is the third-generation language prediction model in the GPT-n series (and the successor to GPT-2) created by OpenAI, a San Francisco-based artificial intelligence research laboratory.
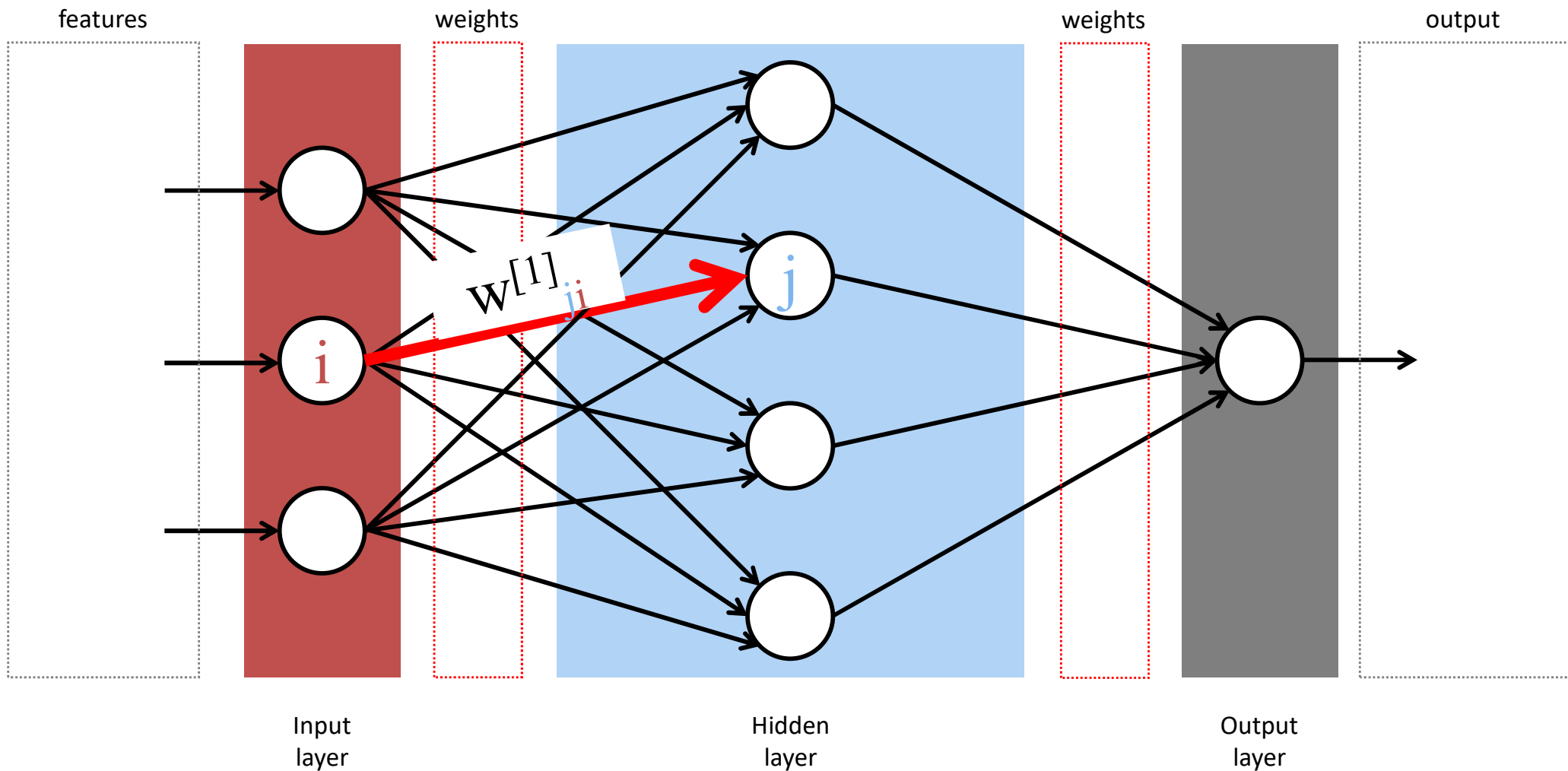
Size:
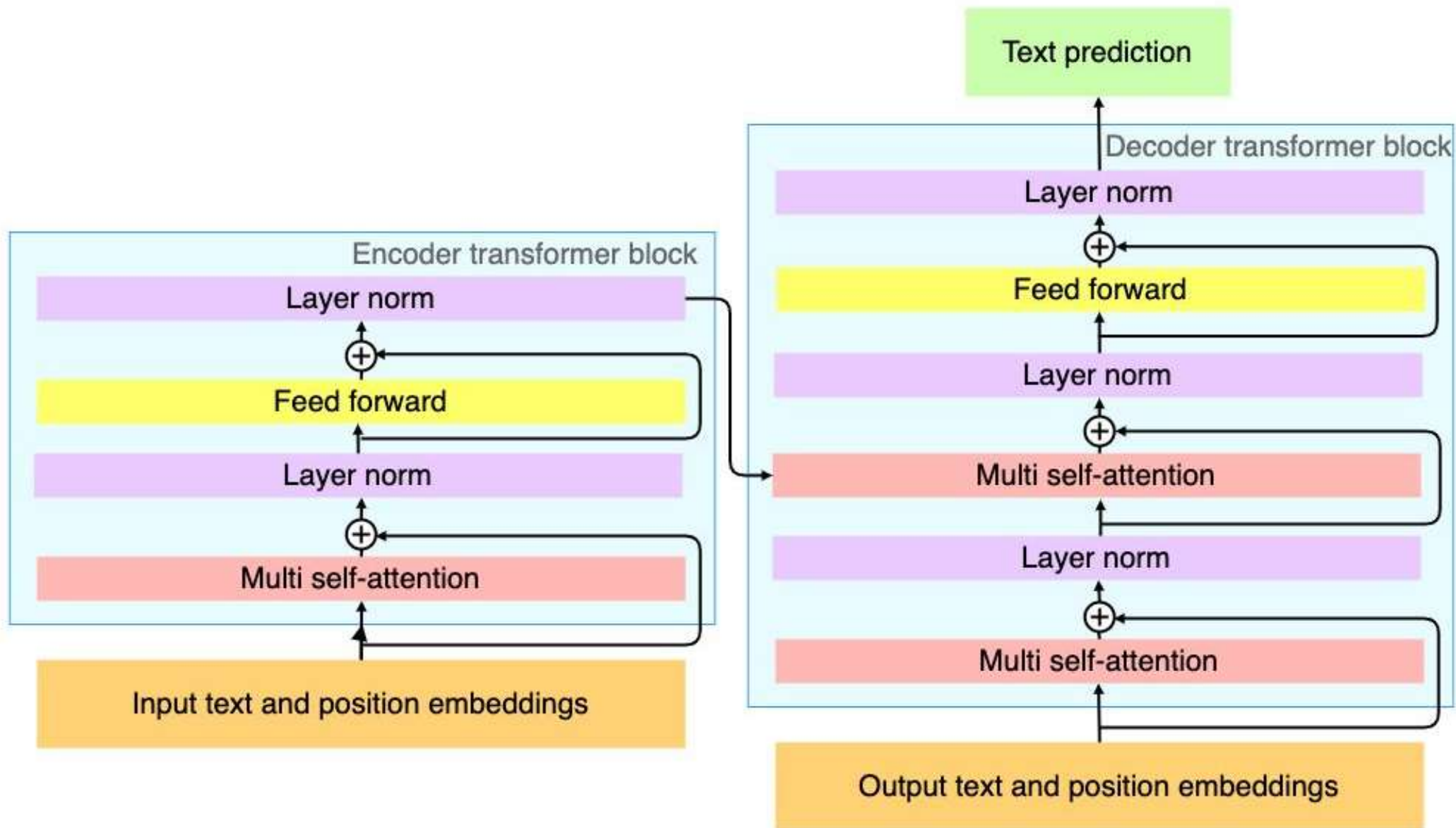**175 billion machine learning parameters**
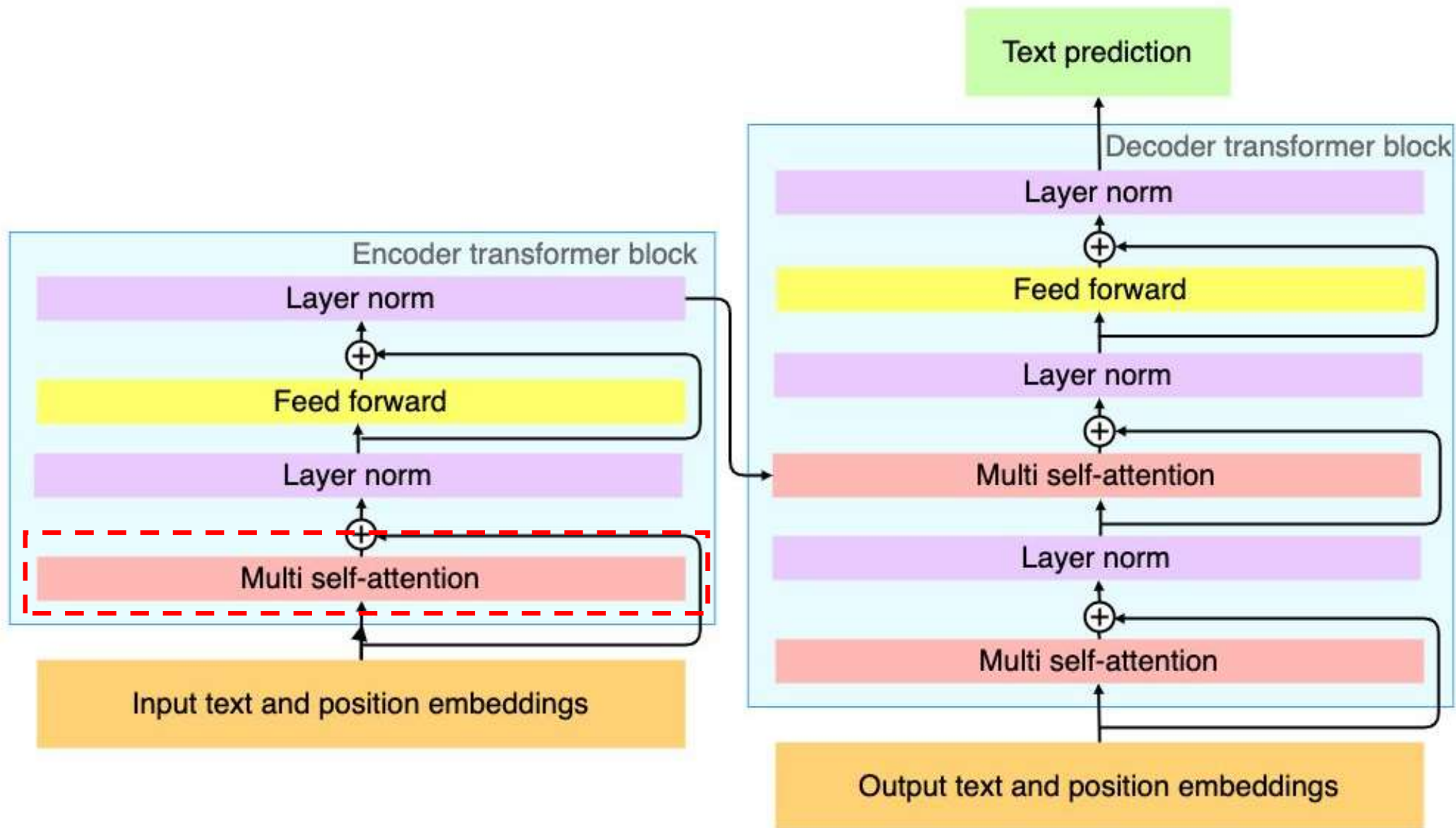~45 GB

*Source: Wikipedia*

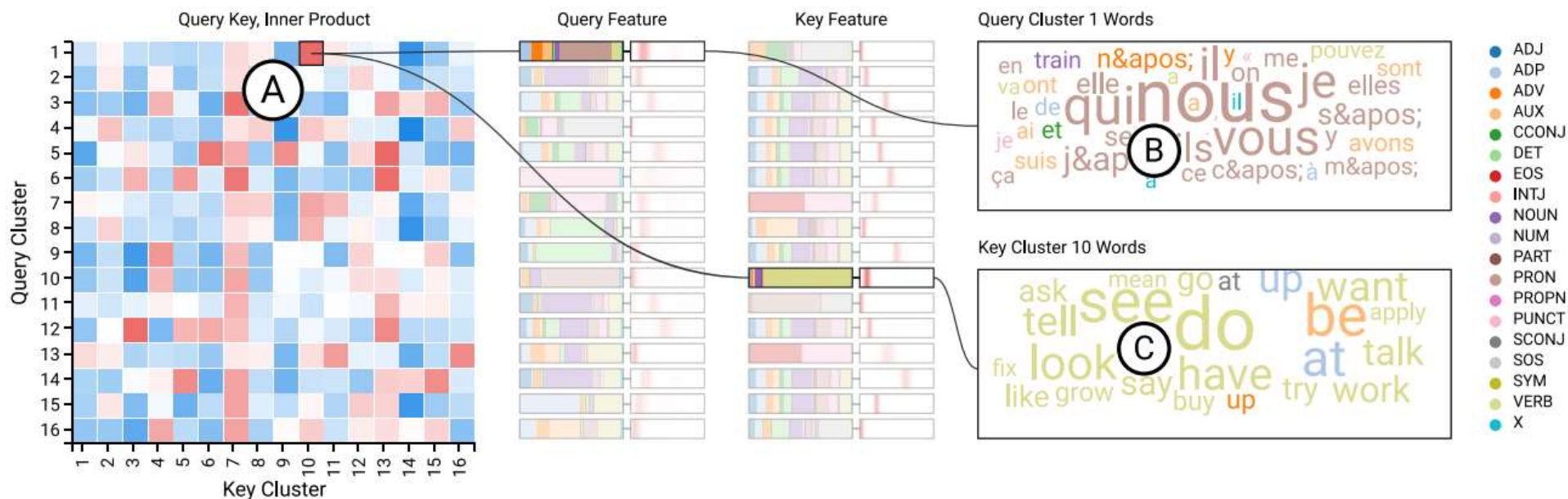# Parameters? What Are Those?

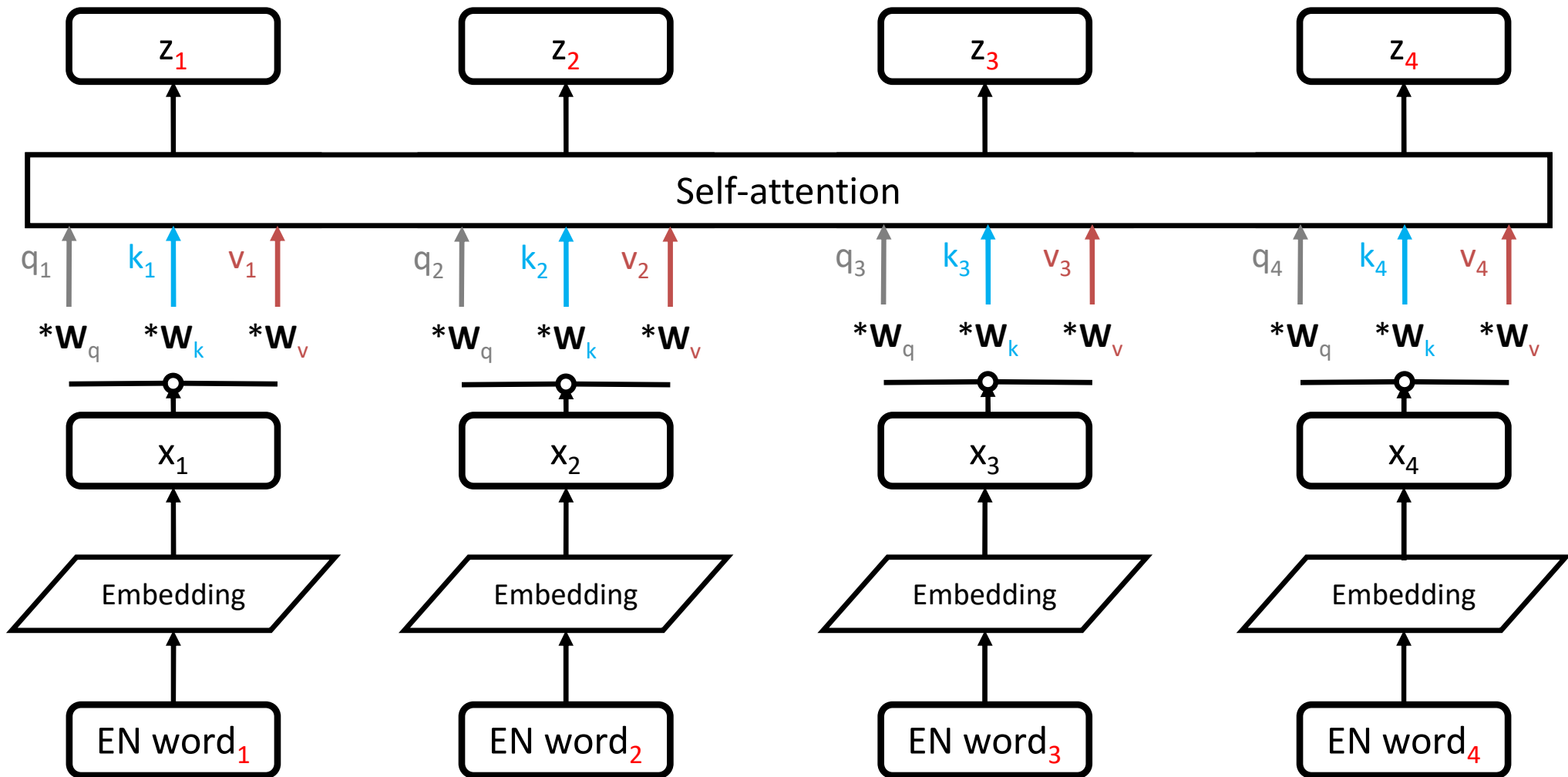# Transformer Architecture

# Transformer Architecture

# Self-Attention

In artificial neural networks, **attention is a technique that is meant to mimic cognitive attention**. The effect **enhances some parts of the input data while diminishing other parts** — the motivation being that the network should devote more focus to the important parts of the data, even though they may be small. **Learning which part of the data is more important than another depends on the context**, and this is trained by gradient descent.
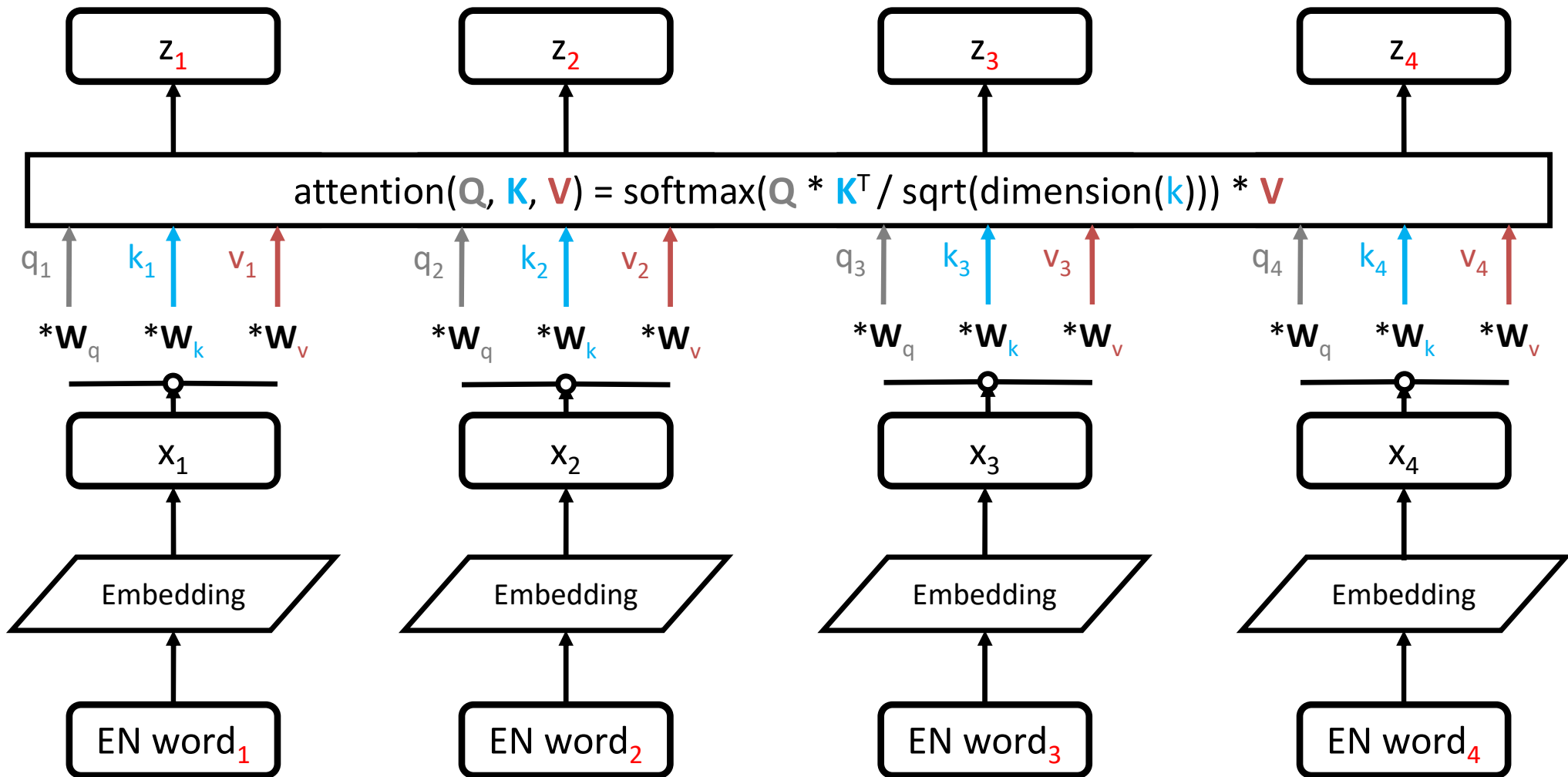


*Source: Park et al. – "SANVis: Visual Analytics for Understanding Self-Attention Networks"*

# Single Head Self-Attention



$q_i$, $k_i$, $v_i$ – query$_i$, key$_i$, value$_i$ | $W_q$, $W_k$, $W_v$ - query, key, value weight matrices [trained: backpropagation]

# Single Head Self-Attention

$z_1$      $z_2$      $z_3$      $z_4$

attention($\mathbf{Q}$, $\mathbf{K}$, $\mathbf{V}$) = softmax($\mathbf{Q}$ * $\mathbf{K}^T$ / sqrt(dimension($k$))) * $\mathbf{V}$

$q_1$   $k_1$   $v_1$     $q_2$   $k_2$   $v_2$     $q_3$   $k_3$   $v_3$     $q_4$   $k_4$   $v_4$

$*\mathbf{W}_q$   $*\mathbf{W}_k$   $*\mathbf{W}_v$    $*\mathbf{W}_q$   $*\mathbf{W}_k$   $*\mathbf{W}_v$    $*\mathbf{W}_q$   $*\mathbf{W}_k$   $*\mathbf{W}_v$    $*\mathbf{W}_q$   $*\mathbf{W}_k$   $*\mathbf{W}_v$

$x_1$      $x_2$      $x_3$      $x_4$

Embedding    Embedding    Embedding    Embedding

EN word$_1$    EN word$_2$    EN word$_3$    EN word$_4$

$q_i$, $k_i$, $v_i$ – query$_i$, key$_i$, value$_i$ | $W_q$, $W_k$, $W_v$ - query, key, value weight matrices [trained: backpropagation]

# Single Head Self-Attention



source: https://jalammar.github.io/illustrated-transformer/

# Single Head Self-Attention



$$\text{softmax}\left(\frac{Q \times K^T}{\sqrt{d_k}}\right) V$$

$= Z$

# Alignment Matrix



*source: https://arxiv.org/pdf/1409.0473.pdf*

# Multi-Head Self-Attention

$$\text{MultiHead}(Q, K, V) = \text{Concat}(\text{head}_1, ..., \text{head}_h)W^O$$

$$\text{where head}_i = \text{Attention}(QW_i^Q, KW_i^K, VW_i^V)$$

1) This is our input sentence*
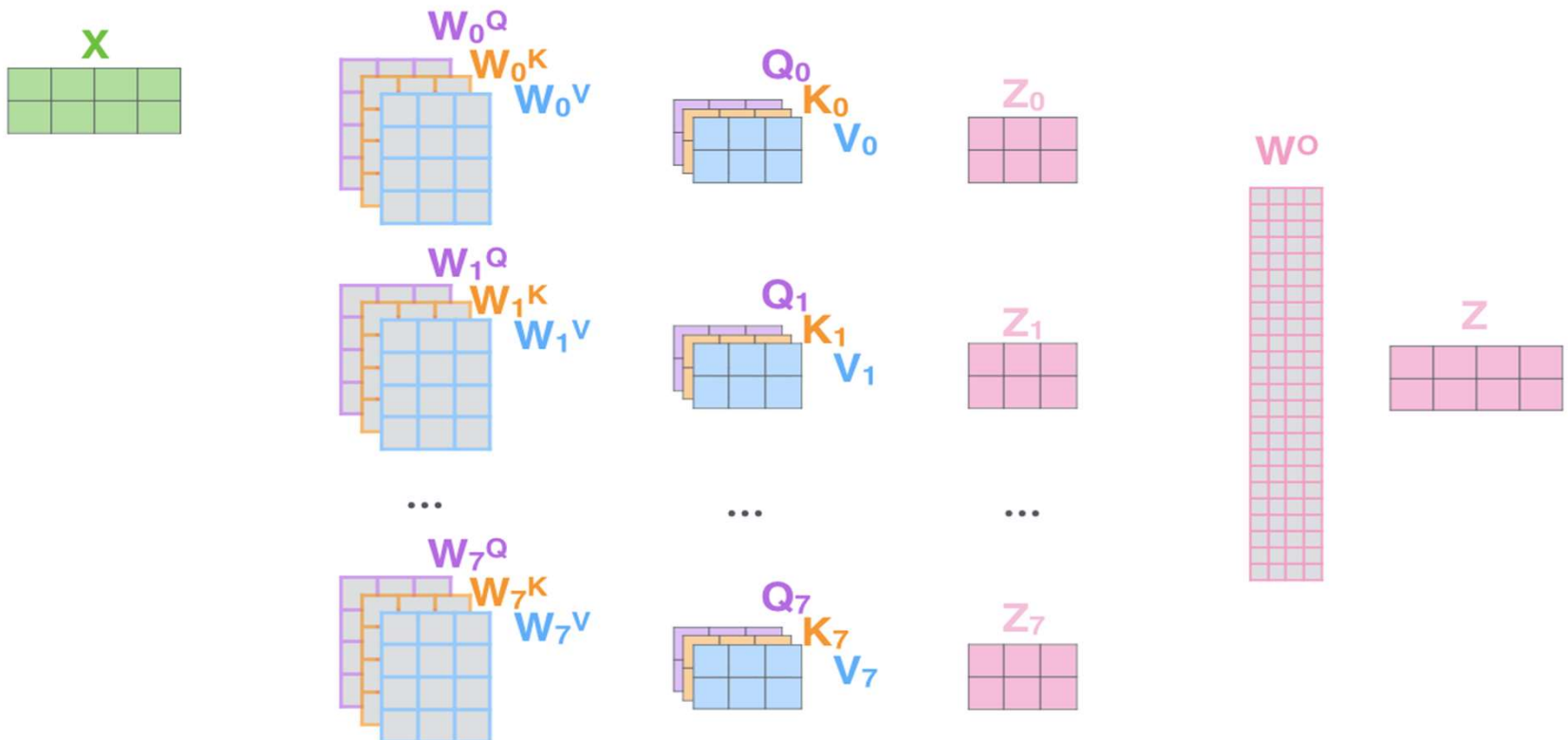
2) We embed each word*

3) Split into 8 heads. We multiply X or R with weight matrices

4) Calculate attention using the resulting Q/K/V matrices

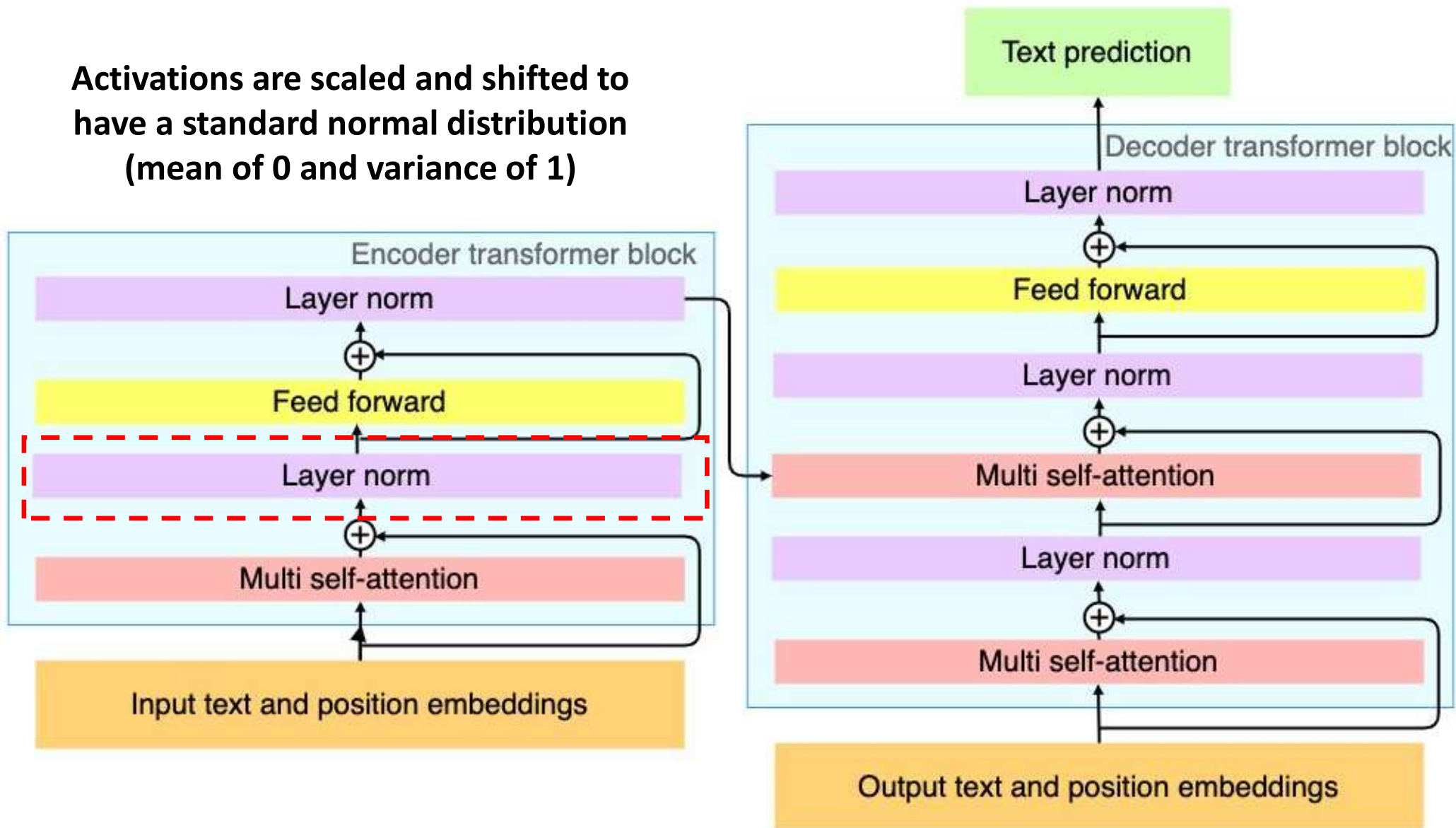5) Concatenate the resulting Z matrices, then multiply with weight matrix $W^O$ to produce the output of the layer



*source: https://jalammar.github.io/illustrated-transformer/*

# Transformer Architecture

**Activations are scaled and shifted to have a standard normal distribution (mean of 0 and variance of 1)**

# Generative Pre-trained Transformer 4

## What is it?

Generative Pre-trained Transformer 4 (GPT-4) is a **multimodal large language model** created by OpenAI. As a transformer, GPT-4 was **pretrained to predict the next token** (using both public data and "data licensed from third-party providers"), and was then **fine-tuned with reinforcement learning from human and AI feedback for human alignment and policy compliance**.
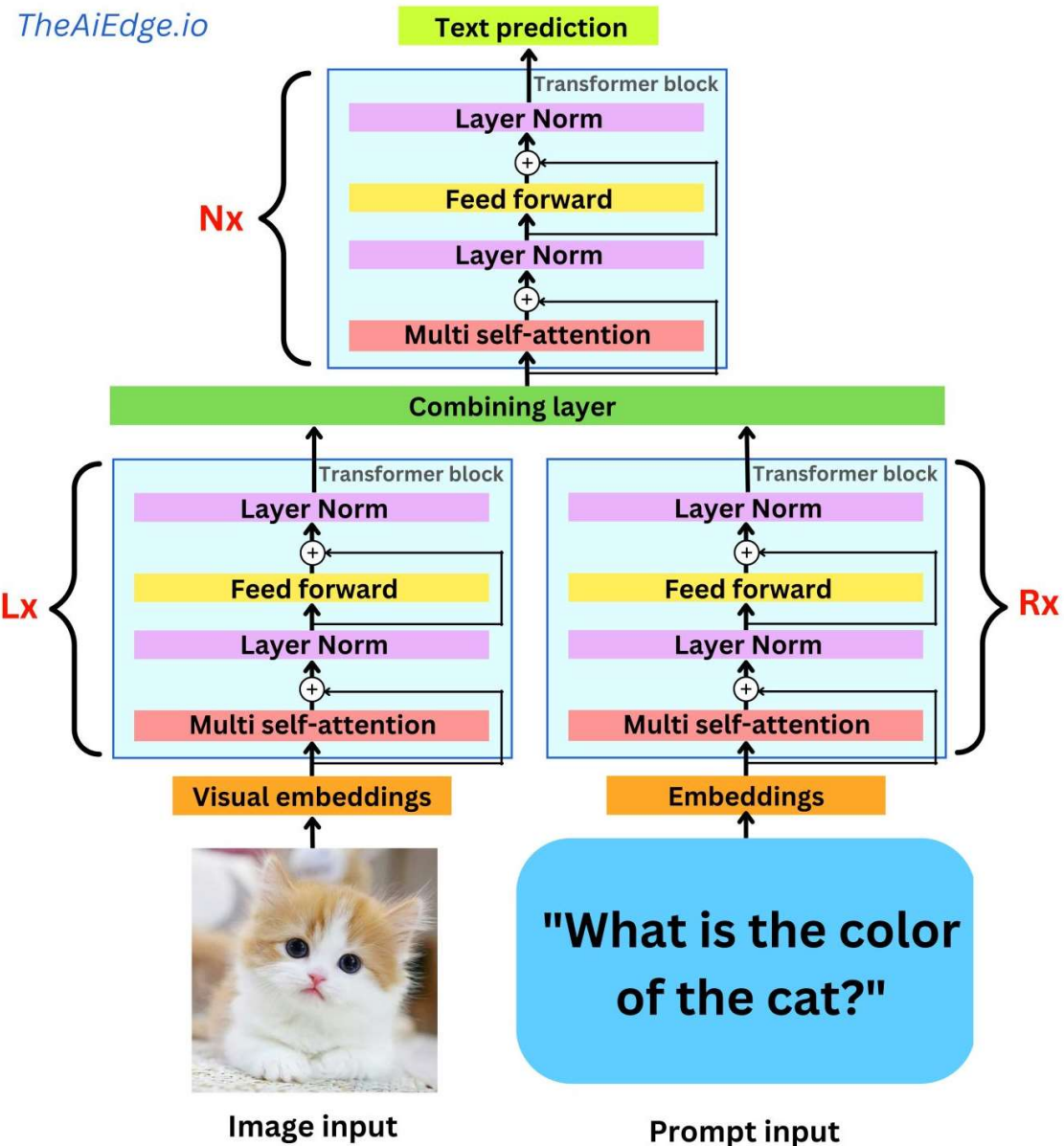
Size:

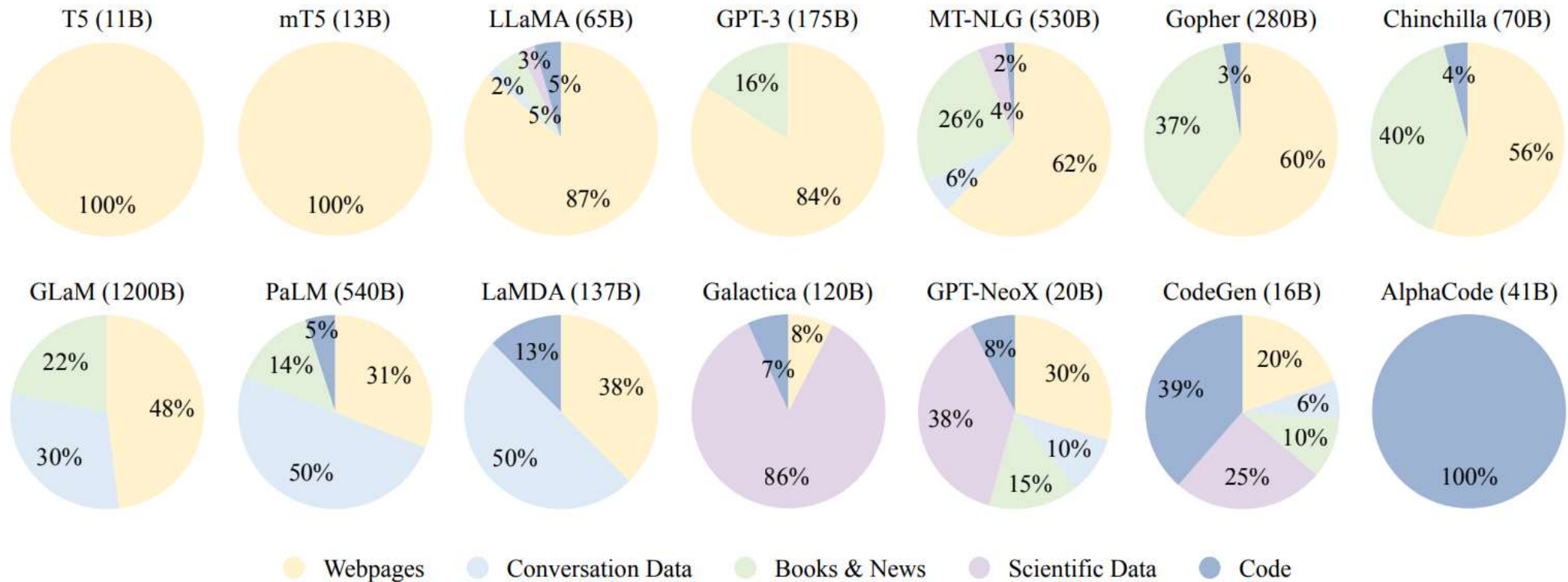**1 trillion machine learning parameters**
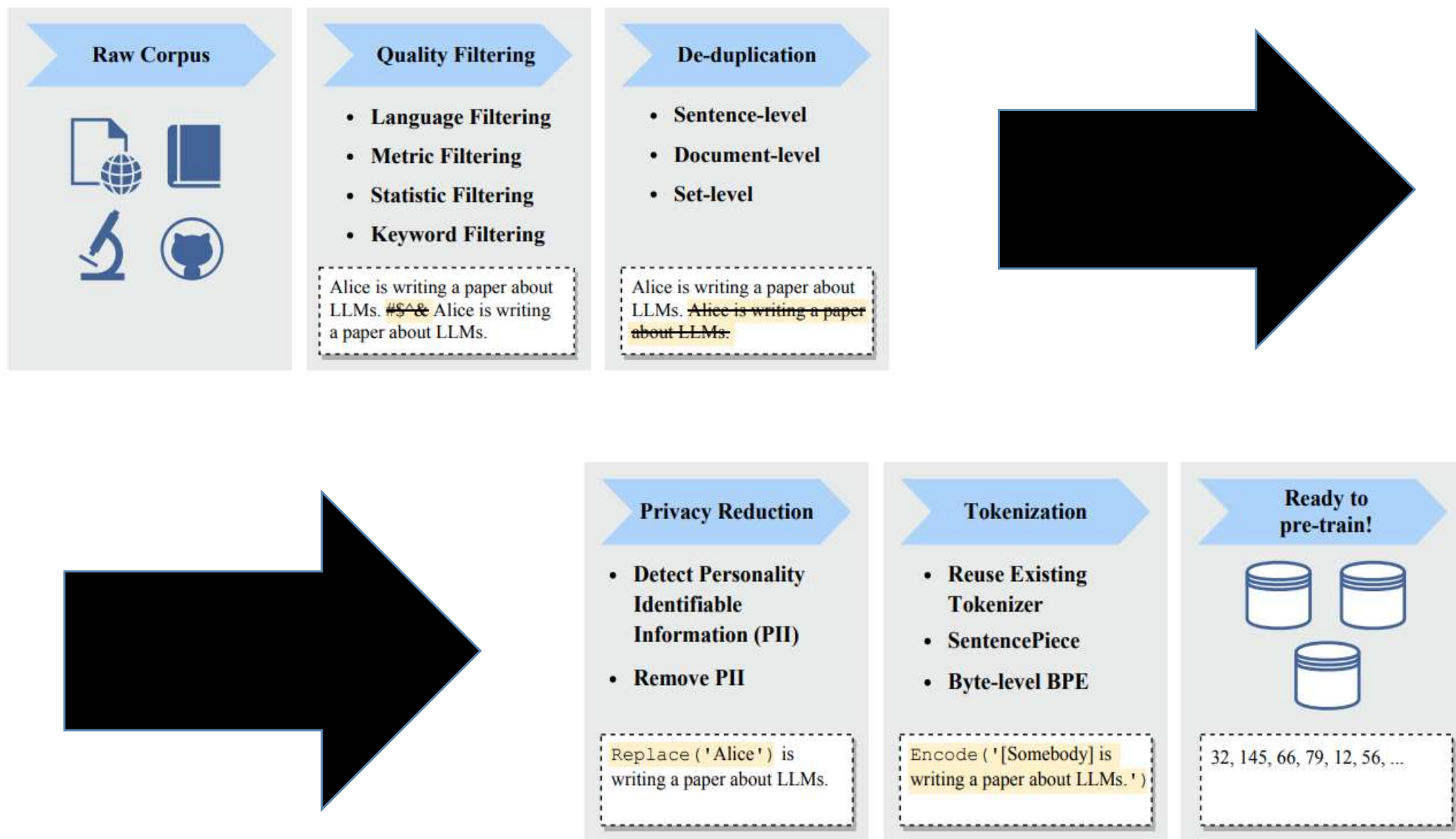
~45 GB

*Source: Wikipedia*

# GPT-4 Architecture

# Large Language Models Data Sources



Source: Zhao et al. – "A Survey of Large Language Models" [2023]

# LLM Data Pre-Processing Pipeline



**Raw Corpus**

**Quality Filtering**
- Language Filtering
- Metric Filtering
- Statistic Filtering
- Keyword Filtering

Alice is writing a paper about LLMs. ~~#$^&~~ Alice is writing a paper about LLMs.

**De-duplication**
- Sentence-level
- Document-level
- Set-level

Alice is writing a paper about LLMs. ~~Alice is writing a paper about LLMs.~~

**Privacy Reduction**
- Detect Personality Identifiable Information (PII)
- Remove PII

Replace('Alice') is writing a paper about LLMs.

**Tokenization**
- Reuse Existing Tokenizer
- SentencePiece
- Byte-level BPE

Encode('[Somebody] is writing a paper about LLMs.')

**Ready to pre-train!**

32, 145, 66, 79, 12, 56, ...

*Source: Zhao et al. – "A Survey of Large Language Models" [2023]*

Illinois Institute of Technology

# ChatGPT

**What is it?**

ChatGPT is a **chatbot** developed by OpenAI and released in November 2022. It is **built on top of OpenAI's GPT-3.5 and GPT-4 families of large language models** (LLMs) and has been **fine-tuned** (an approach to **transfer learning**) **using both supervised and reinforcement learning techniques**.
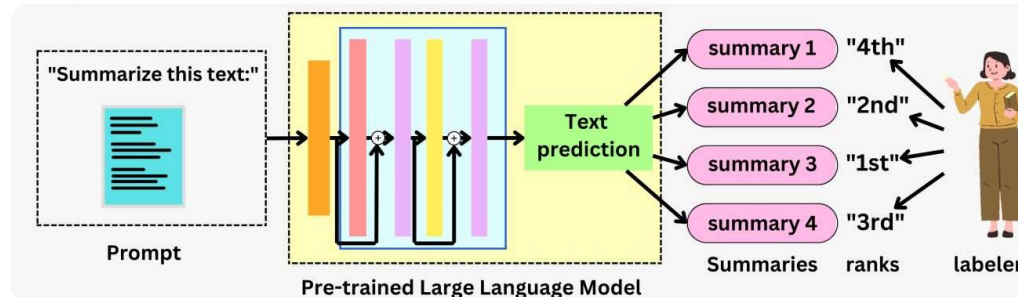
*Source: Wikipedia*

# Transfer Learning

In **transfer learning**, **experience** with one learning task **helps an agent learn better on another task**.

Pre-trained models can be used as a starting point for developing new models.

# ChatGPT: Learning From Feedback



Source: TheAiEdge.io

# Generative Models

# Generative vs Discriminative Models
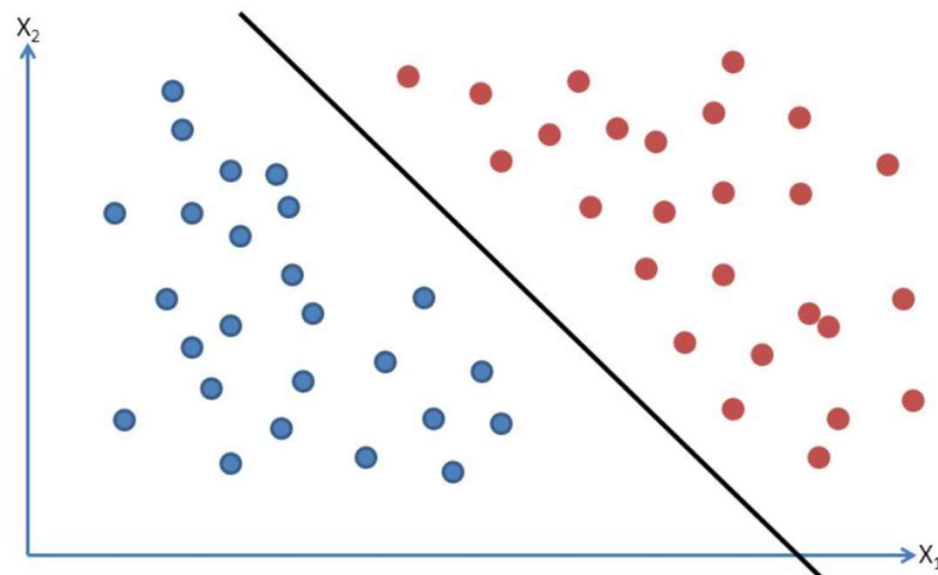
**Generative**

**Discriminative**



**Generative model models actual distributions for EACH CLASS / LABEL / TAG**

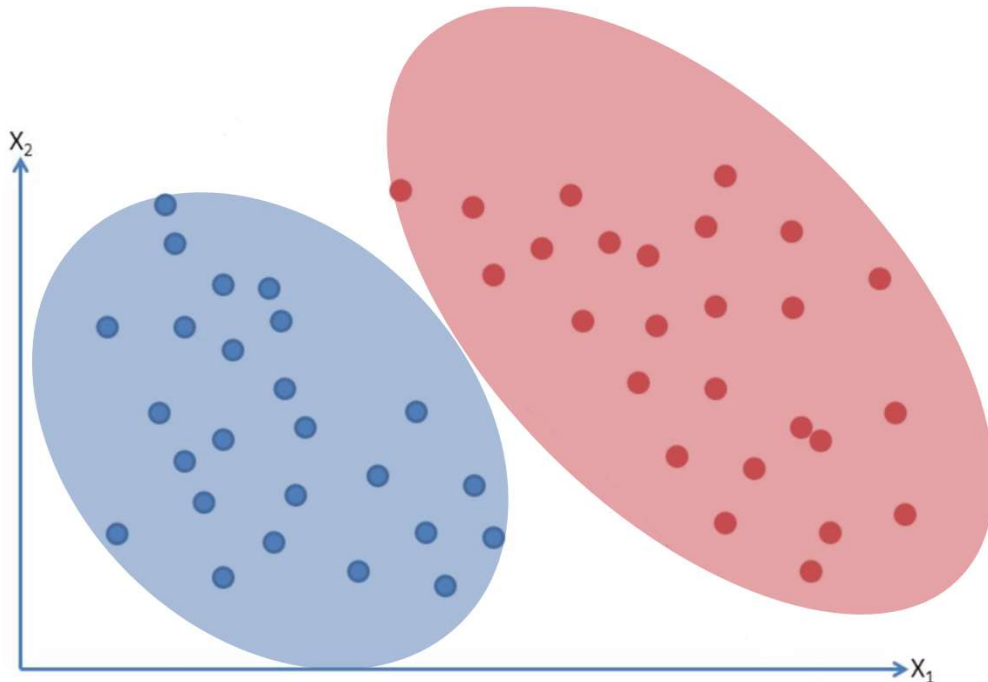**to**

**make a** $P(class \mid sample)$ **prediction**

**Discriminative model models the decision boundary between CLASSES / LABELS / TAGS**

**to**

**make a** $P(class \mid sample)$ **prediction**

# Generative vs Discriminative Models

**Generative**

**Discriminative**
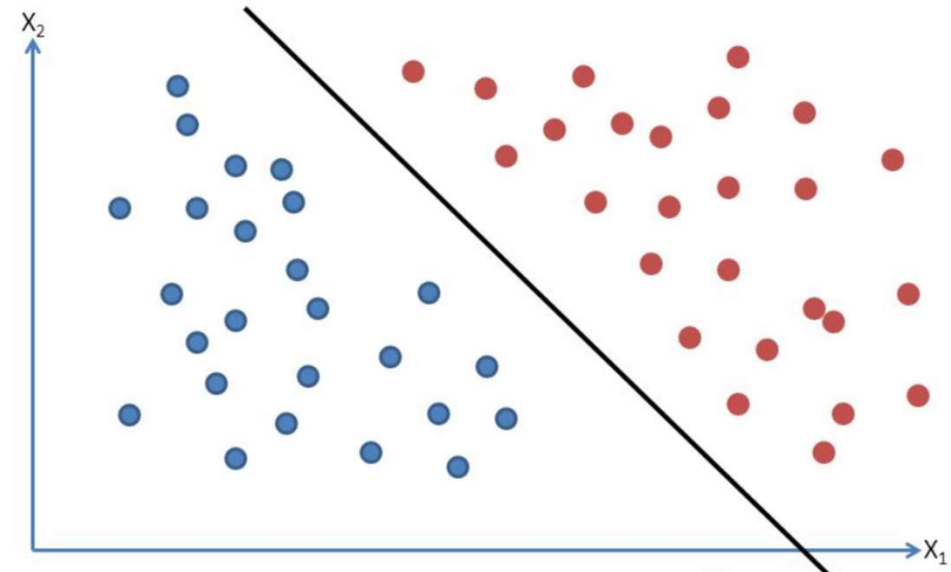


Generative model uses training data to learn $P(\text{sample, class})$ <span style="color:red">**joint probabilities**</span>

and then

uses Bayes Theorem to get the $P(\text{class} \mid \text{sample})$ prediction
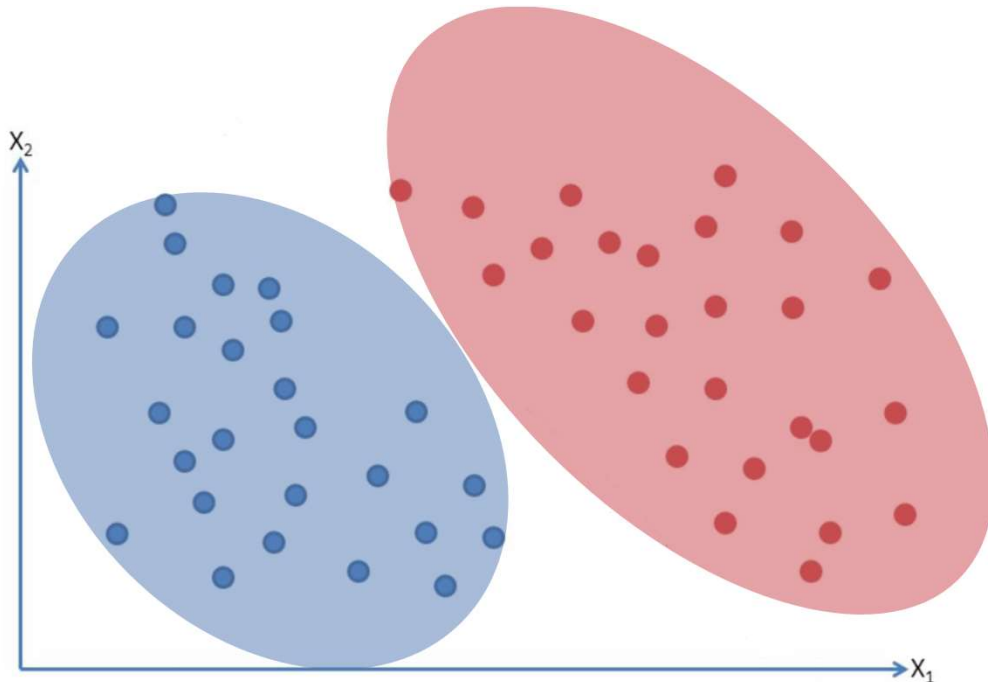
Discriminative model uses training data to learn $P(\text{class} \mid \text{sample})$ <span style="color:red">**conditional probability**</span>

and then

uses it to make a prediction
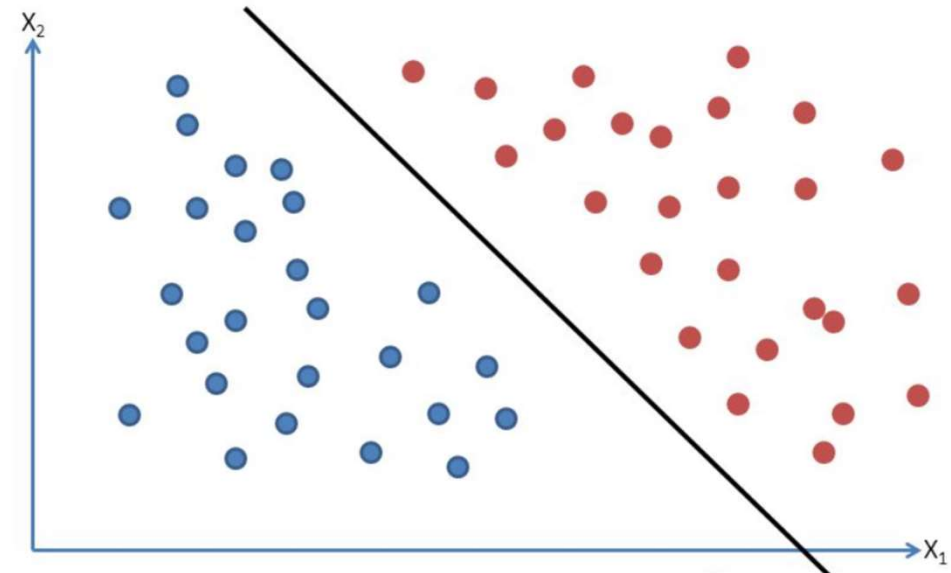
# Generative vs Discriminative Classifier

**Generative**

**Discriminative**



**Generative classifiers:**

- **Assume some form of** $P(\text{class})$**,** $P(\text{sample} \mid \text{class})$
- **Estimate** $P(\text{class})$**,** $P(\text{sample} \mid \text{class})$ **using training data**
- **Use Bayes Theorem to calculate** $P(\text{class} \mid \text{sample})$
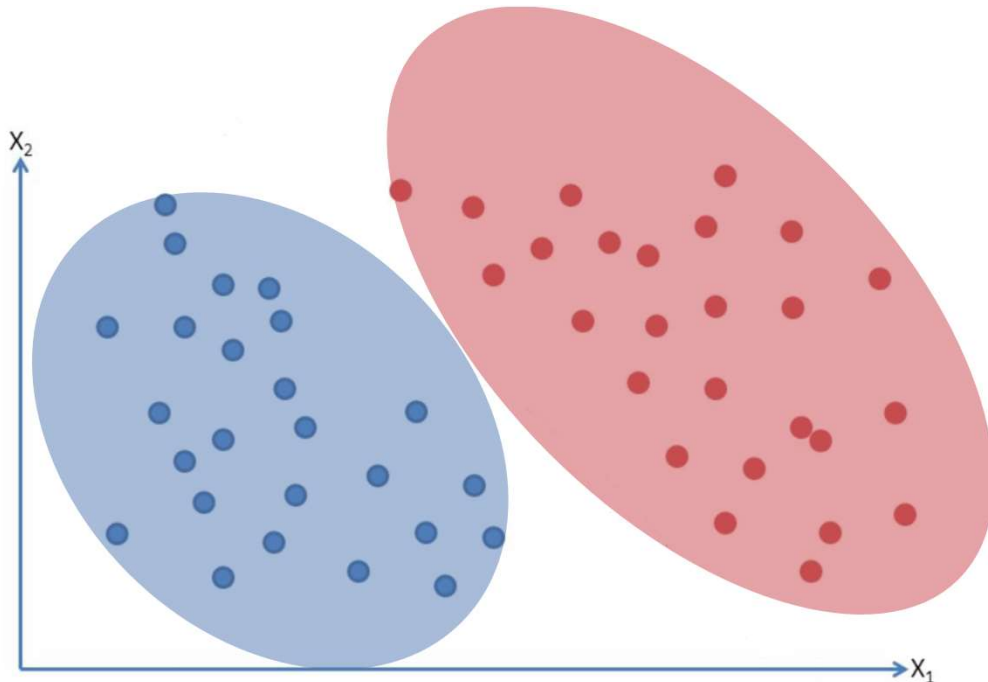
**Discriminative classifiers:**

- **Assume some form of** $P(\text{class} \mid \text{sample})$
- **Estimate** $P(\text{class} \mid \text{sample})$ **using training data**
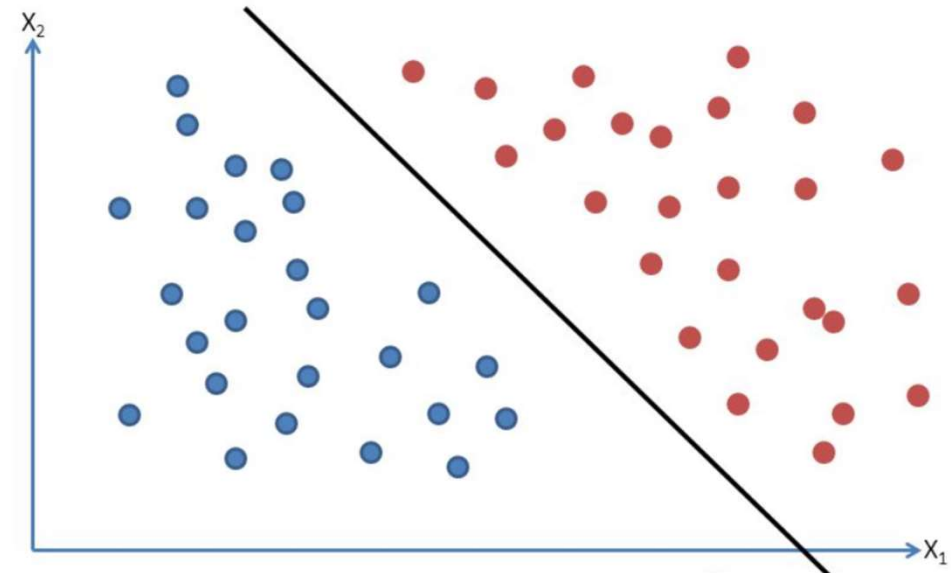
# Generative vs Discriminative Classifier

**Generative**



**Discriminative**



Generative classifiers:
- **Naive Bayes**
- **Bayesian networks**
- **Markov random fields**
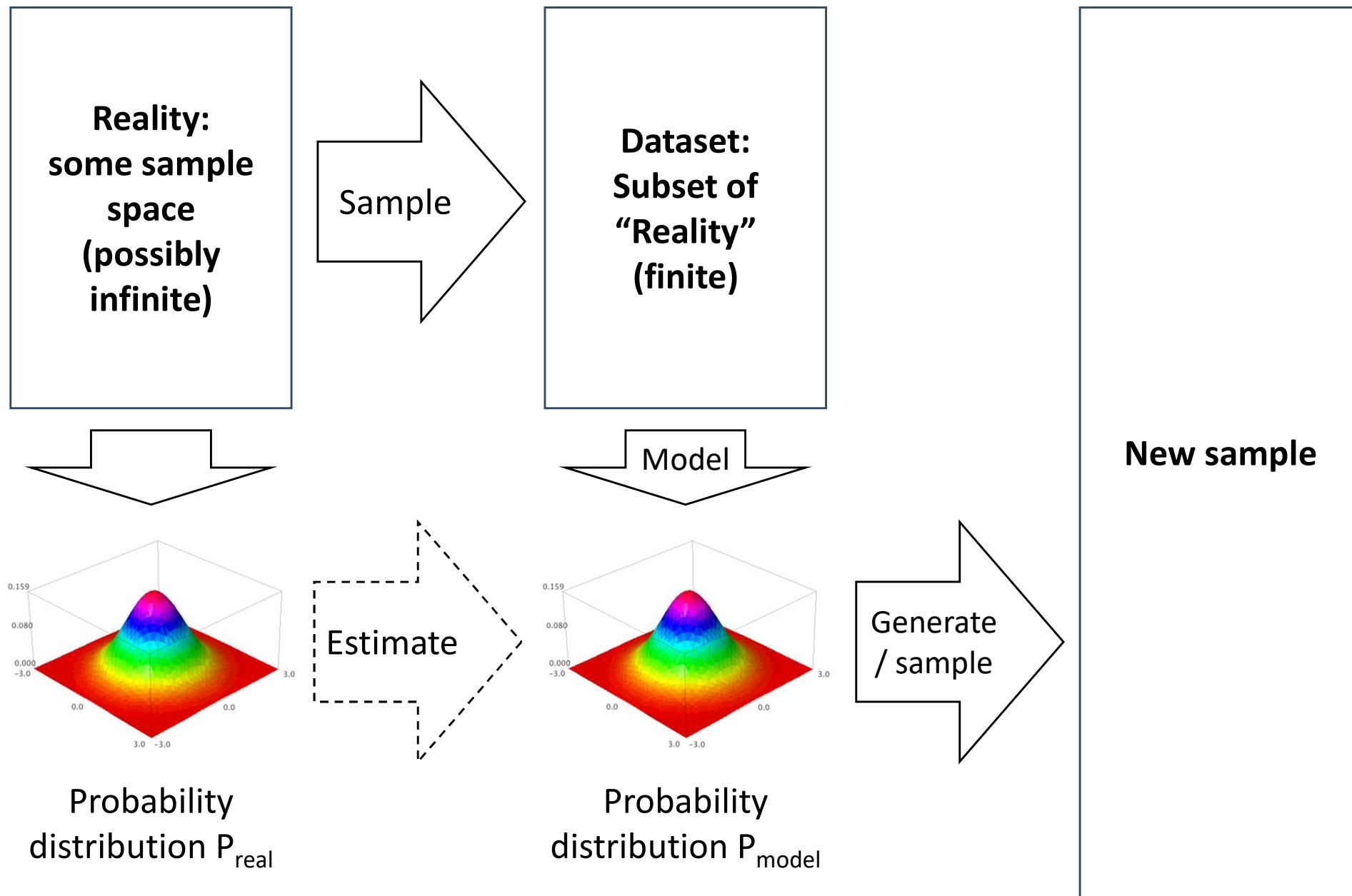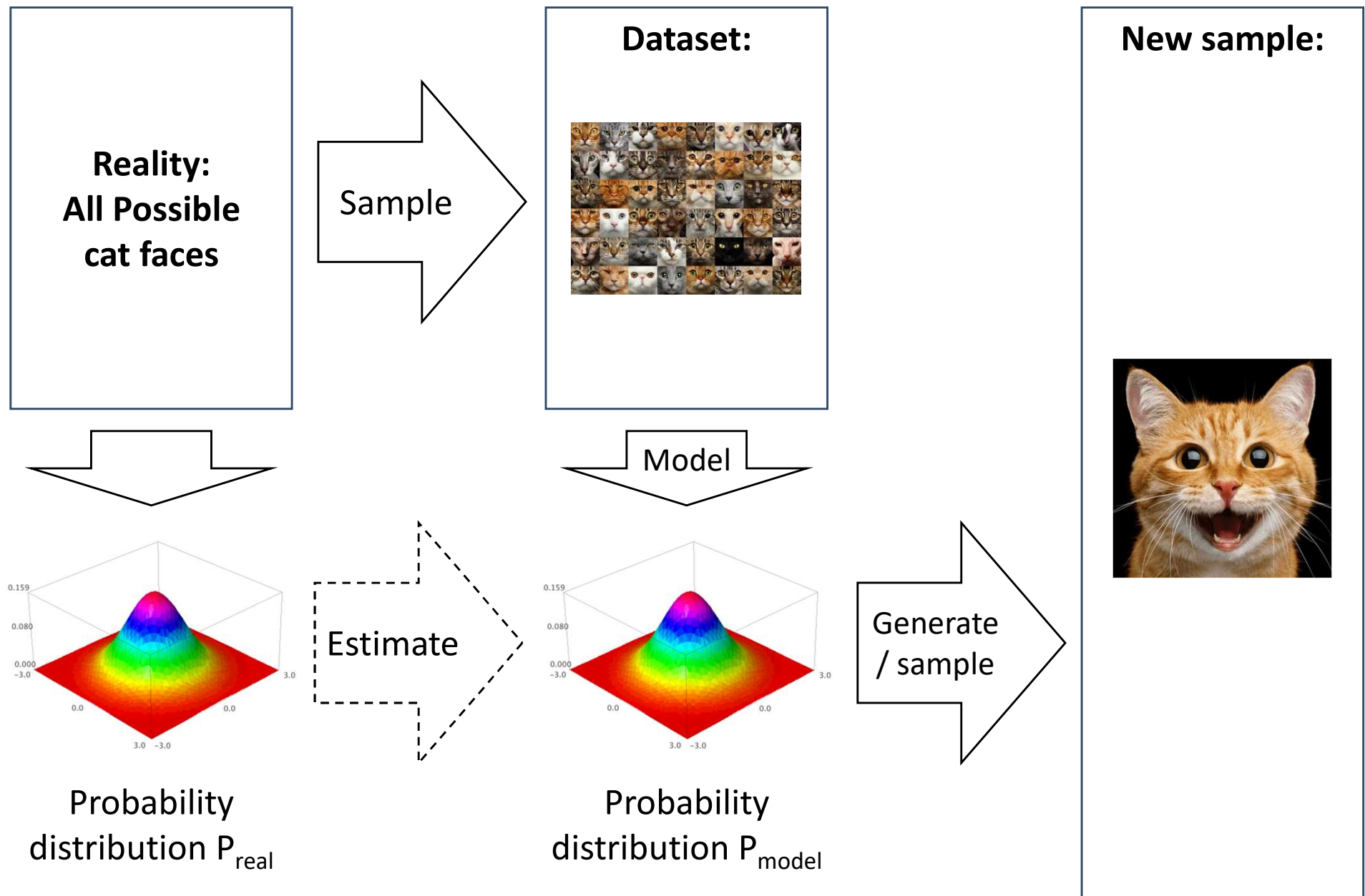- **Hidden Markov Models (HMM)**

Discriminative classifiers:
- **Logistic regression**
- **Support Vector Machines**
- **Traditional neural networks**
- **k-Nearest Neighbors**
- **Conditional Random Fields (CRF)s**

# Generative AI Model: the Idea



**Reality:**
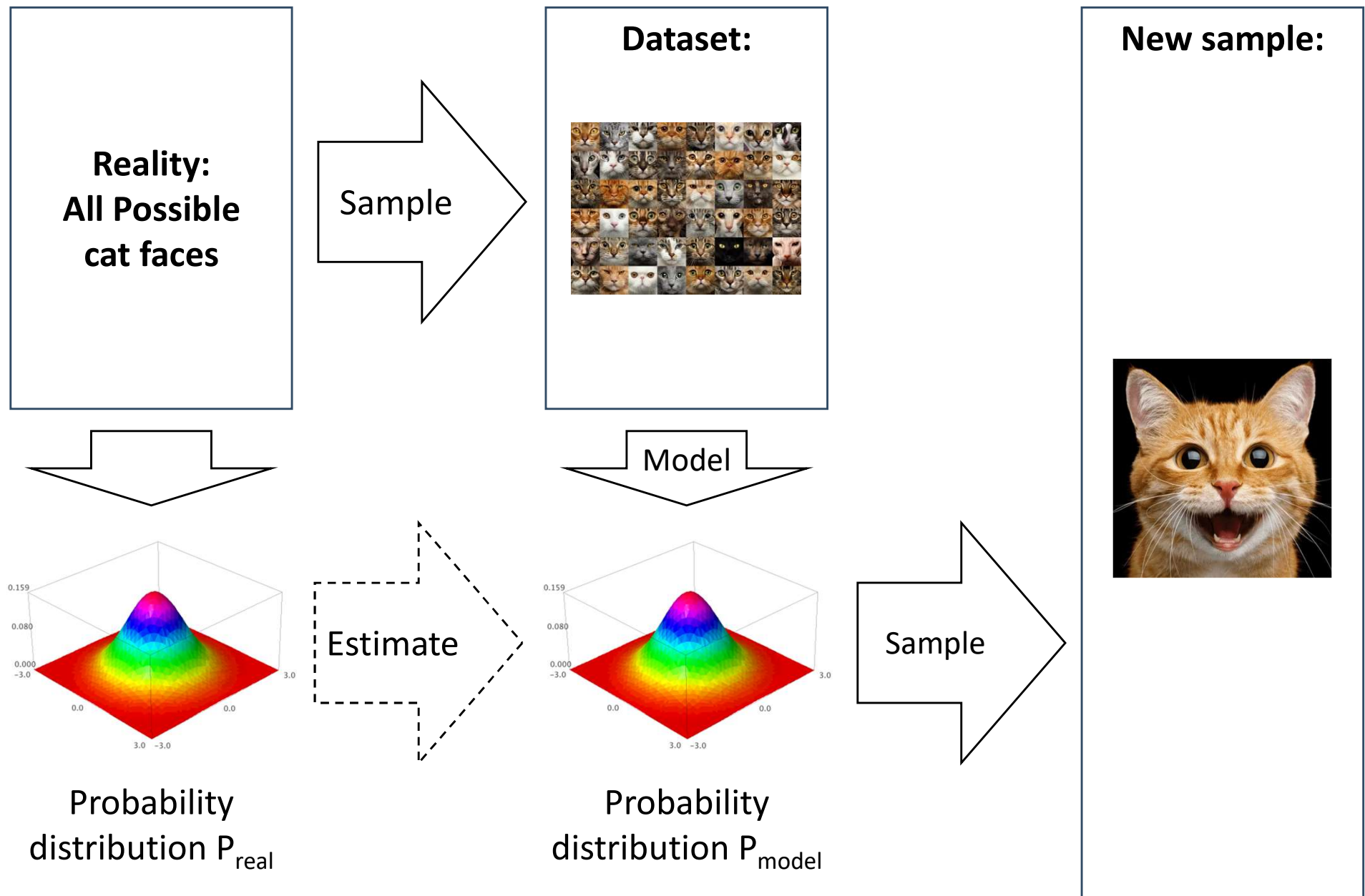**some sample**
**space**
**(possibly**
**infinite)**

Sample

**Dataset:**
**Subset of**
**"Reality"**
**(finite)**

Model

**New sample**

Probability
distribution $P_{real}$

Estimate

Probability
distribution $P_{model}$

Generate
/ sample

# Generative AI Model: the Idea



Reality:
All Possible
cat faces

Sample →

Dataset:

Model

New sample:

Estimate →

Generate / sample

Probability distribution $P_{real}$

Probability distribution $P_{model}$

# Taxonomy of Generative AI Models

GAN – Generative Adversarial Network

# Explicit Density



Reality: All Possible cat faces → Sample → Dataset: → Model → New sample:

Probability distribution $P_{real}$ → Estimate → Probability distribution $P_{model}$ → Sample

# Implicit Density



**Reality:
All Possible
cat faces**

Sample

**Dataset:**

Model

Probability
distribution $P_{real}$

**Stochastic
generation
process**
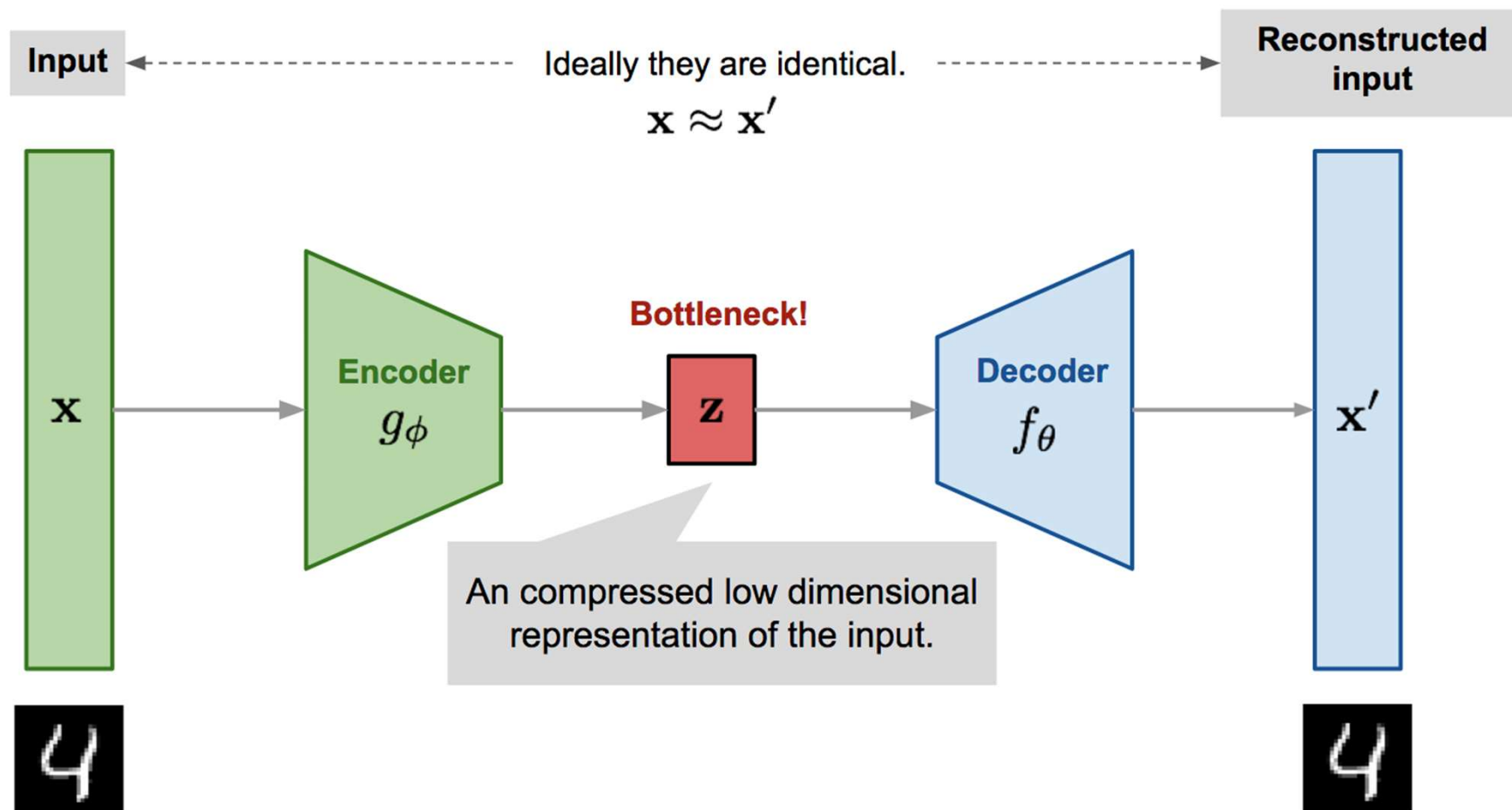
Generate

**New sample:**

# Tractable vs. Approximate Density

**Tractable density models place constraints on the model architecture so that the density function has a form that makes it easy to calculate.**

**Approximate density models use variety of techniques to approximate the density function:**
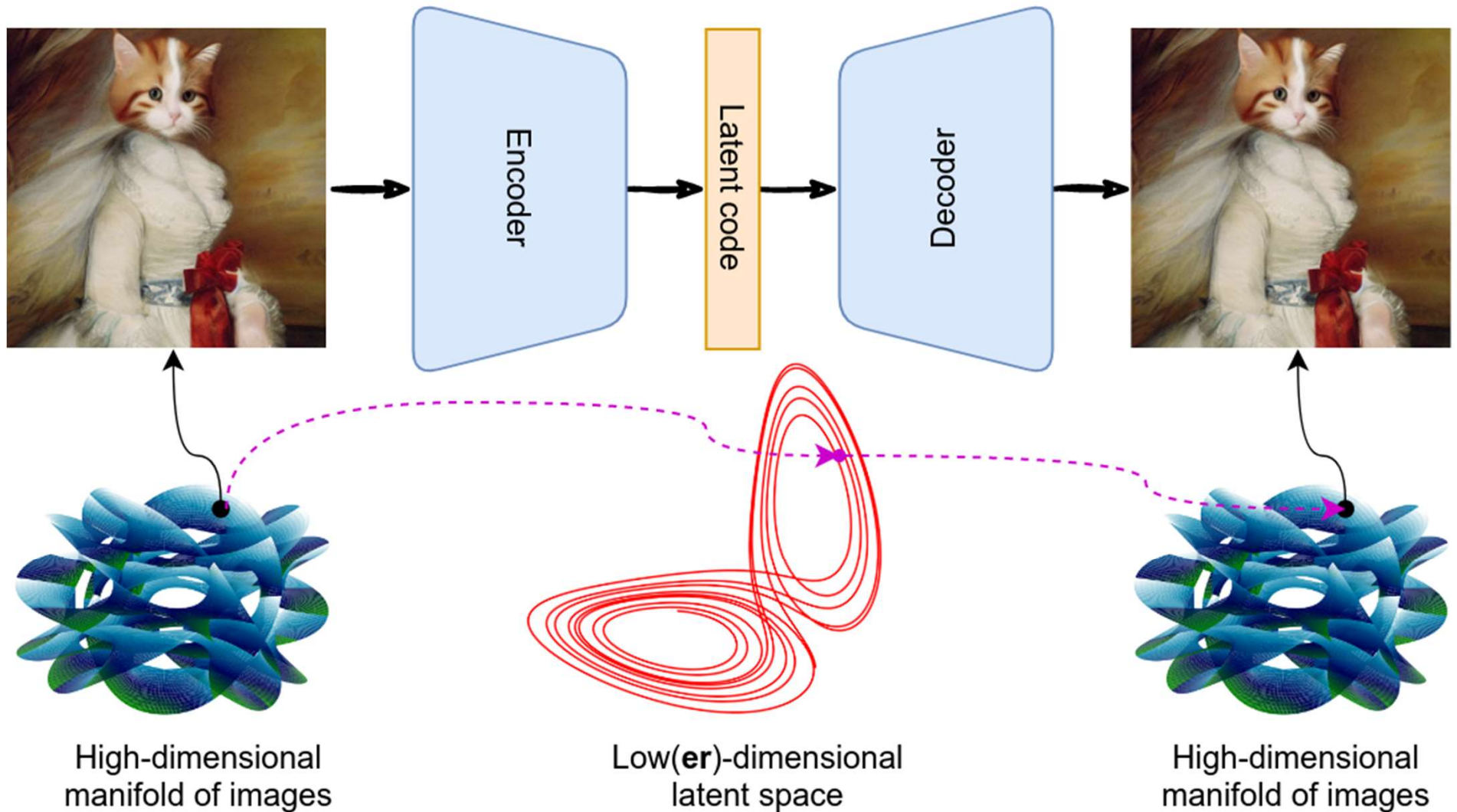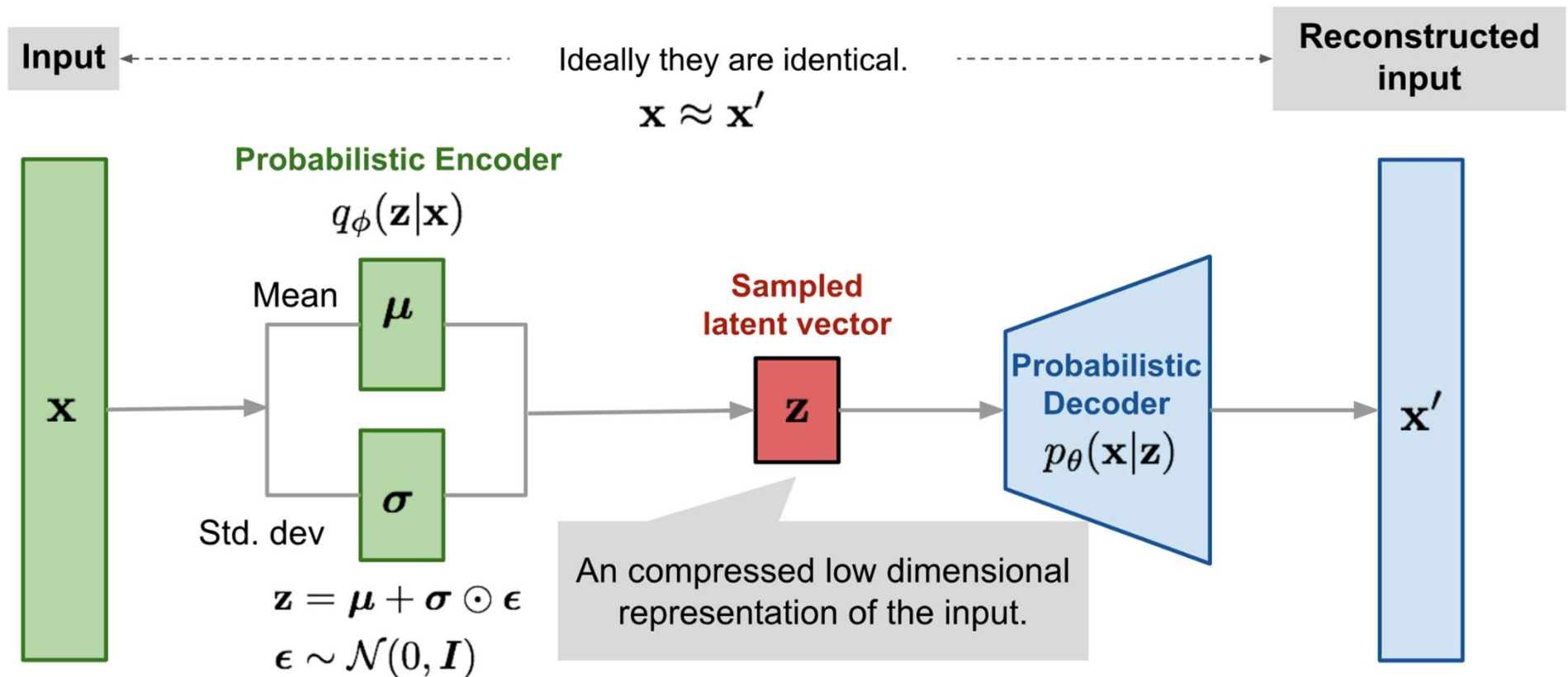
- **latent vectors**
- **denoising**

# Autoencoder Model



Source: https://lilianweng.github.io/posts/2018-08-12-vae/

# Latent Space



High-dimensional manifold of images

Low(**er**)-dimensional latent space

High-dimensional manifold of images

# Variational Autoencoder Model



*Source: https://lilianweng.github.io/posts/2018-08-12-vae/*

# Autoregressive Model (GPT-3)

**What is it?**

Generative Pre-trained Transformer 3 (GPT-3) is an **autoregressive language model that uses deep learning to produce human-like text**. It is the third-generation language prediction model in the GPT-n series (and the successor to GPT-2) created by OpenAI, a San Francisco-based artificial intelligence research laboratory.
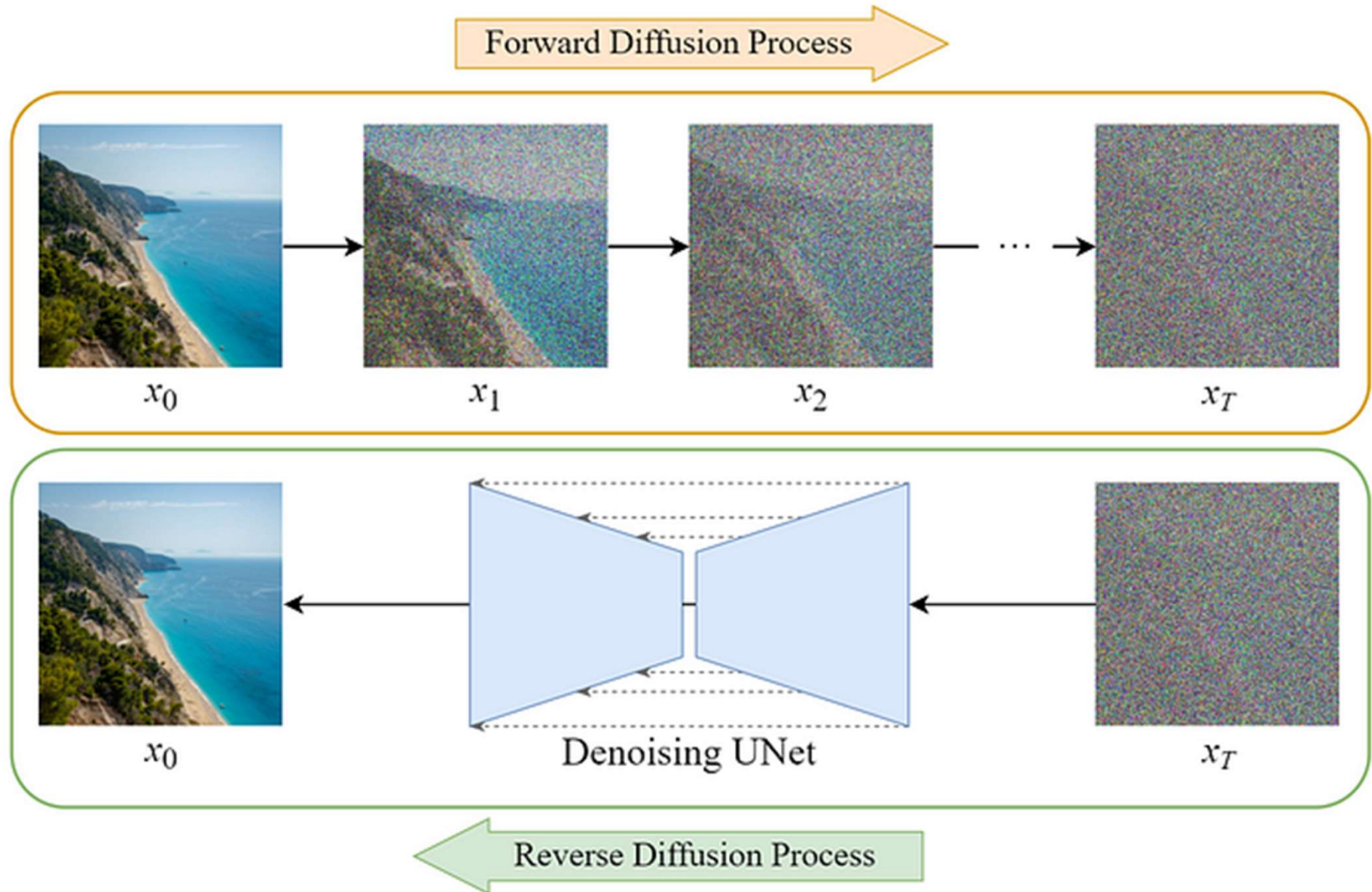
Size:
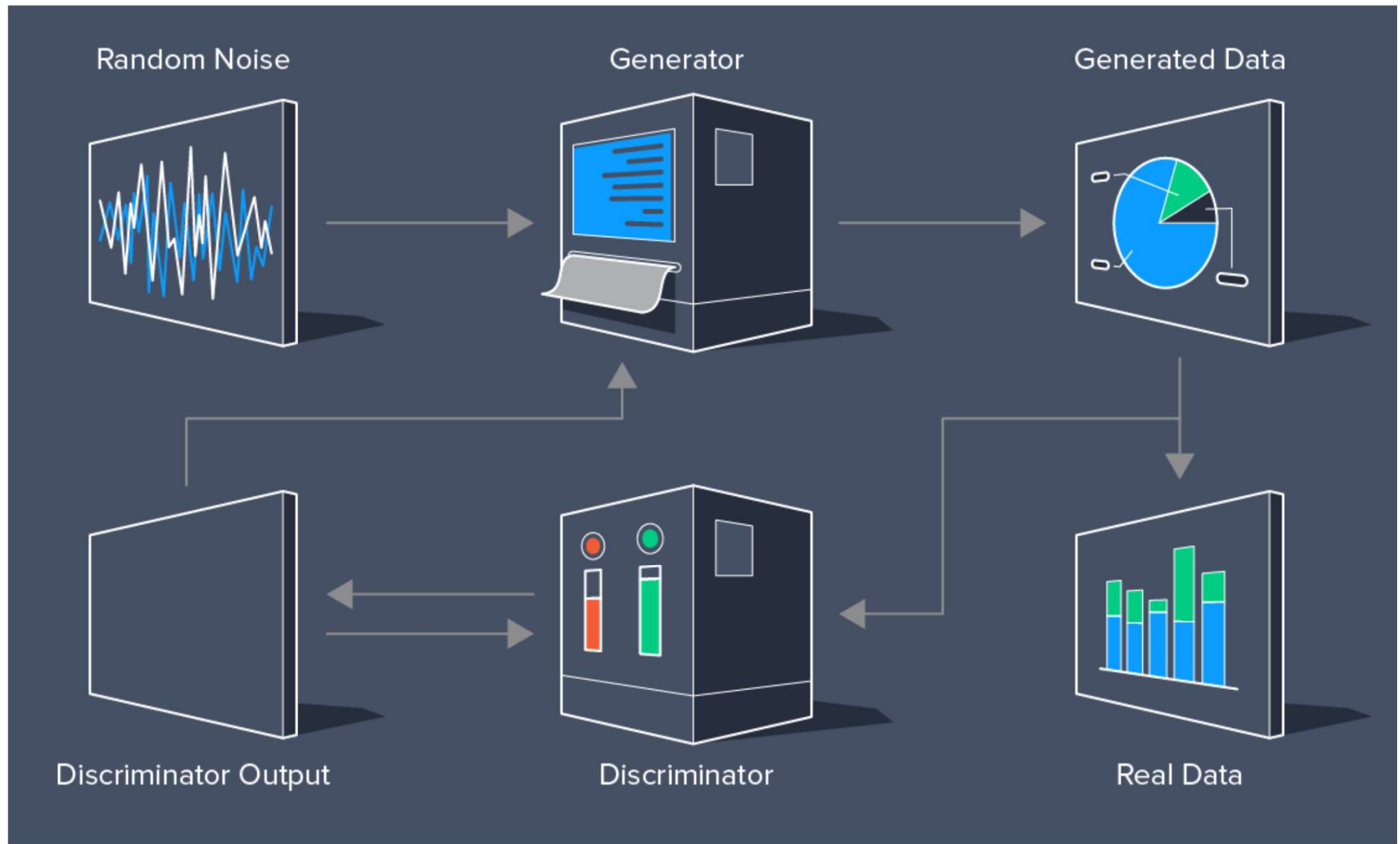
**175 billion machine learning parameters**

~45 GB

*Source: Wikipedia*

# Diffusion Model



*Source: https://medium.com/@steinsfu/stable-diffusion-clearly-explained-ed008044e07e*

# Generative Adversarial Network



*Source: https://www.toptal.com/machine-learning/generative-adversarial-networks*