

Long Text Sequences

Tasks In NLP:

Long Text Sequences

Tasks In NLP:



Writing Books

Long Text Sequences

Tasks In NLP:



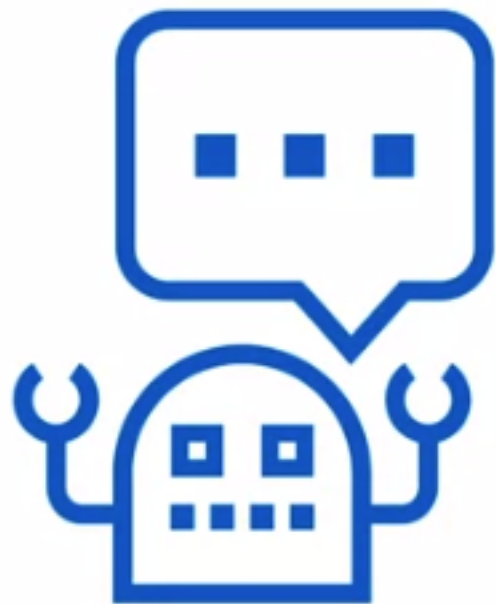
Writing Books



Chatbots

Chatbots

Context Windows:

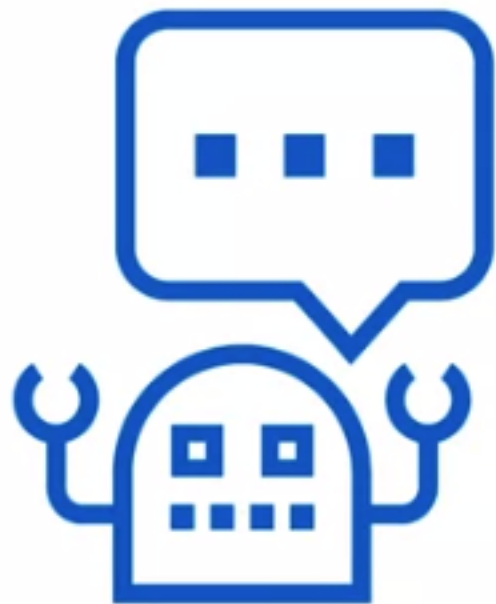


Chatbots

Context Windows:

User 1: What's for dinner?

Chatbot: Who's cooking, you or me?



Chatbots

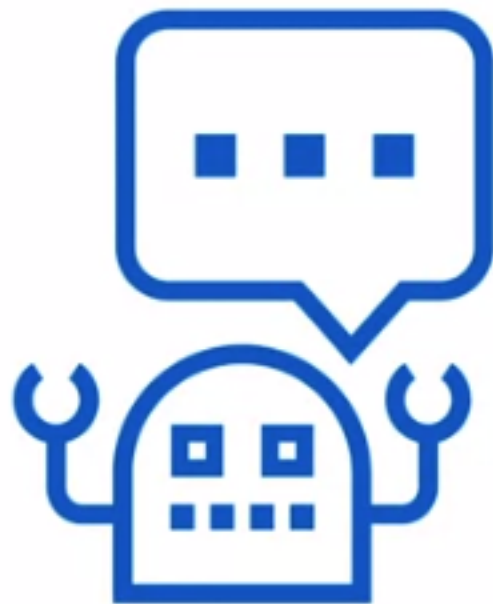
Context Windows:

User 1: What's for dinner?

Chatbot: Who's cooking, you or me?

User 1: Hey now chatbot.

Chatbot: I hope it's not hay, that's
what horses eat.



Transformer Issues

Transformer Issues

- Attention on sequence of length L takes L^2 time and memory

Transformer Issues

- Attention on sequence of length L takes L^2 time and memory

Transformer Issues

- Attention on sequence of length L takes L^2 time and memory
- N layers take N times as much memory
GPT-3 has 96 layers and new models will have more

Transformer Issues

- Attention on sequence of length L takes L^2 time and memory

$L=100$ $L^2 = 10K$ (0.001s at 10M ops/s)

$L=1000$ $L^2 = 1M$ (0.1s at 10M ops/s)

$L=10000$ $L^2 = 100M$ (10s at 10M ops/s)

$L=100000$ $L^2 = 10B$ (1000s at 10M ops/s)

- N layers take N times as much memory

GPT-3 has 96 layers and new models will have more

Transformer Issues

- Attention on sequence of length L takes L^2 time and memory

$L=100$ $L^2 = 10K$ (0.001s at 10M ops/s)

$L=1000$ $L^2 = 1M$ (0.1s at 10M ops/s)

$L=10000$ $L^2 = 100M$ (10s at 10M ops/s)

$L=100000$ $L^2 = 10B$ (1000s at 10M ops/s)

- N layers take N times as much memory

GPT-3 has 96 layers and new models will have more

Attention Complexity

- Attention: $\text{softmax}(QK^T)V$

Attention Complexity

- Attention: $\text{softmax}(QK^T)V$
- Q, K, V are all $[L, d_{\text{model}}]$

Attention Complexity

- Attention: $\text{softmax}(QK^T)V$
- Q, K, V are all $[L, d_{\text{model}}]$
- QK^T is $[L, L]$

Attention Complexity

- Attention: $\text{softmax}(QK^T)V$
- Q, K, V are all $[L, d_{\text{model}}]$
- QK^T is $[L, L]$
- Save compute by using area of interest for large L

Memory with N Layers

- Activations need to be stored for backprop

Memory with N Layers

- Activations need to be stored for backprop
- Big models are getting bigger

Memory with N Layers

- Activations need to be stored for backprop
- Big models are getting bigger
- Compute vs memory tradeoff

What does Attention do?

Select Nearest Neighbors (K,Q) and return corresponding V

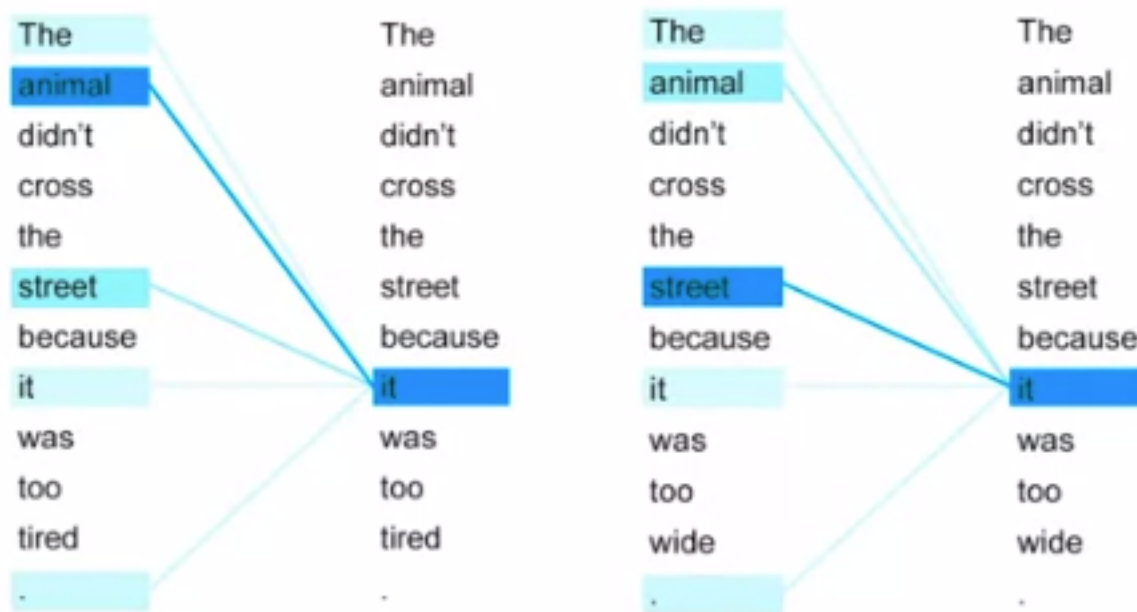


image ©
[\(Transformer: A Novel Neural Network Architecture for Language Understanding.\)](#)

What does Attention do?

Select Nearest Neighbors (K,Q) and return corresponding V

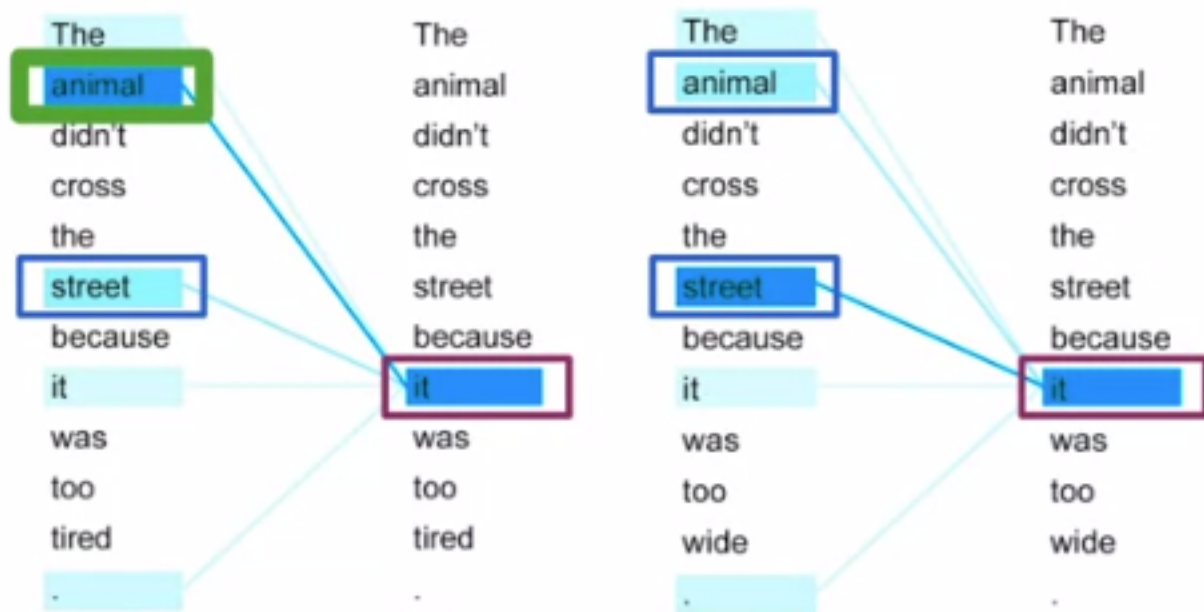


image ©
[\(Transformer: A Novel Neural Network Architecture for Language Understanding.\)](#)

What does Attention do?

Select Nearest Neighbors (K,Q) and return corresponding V

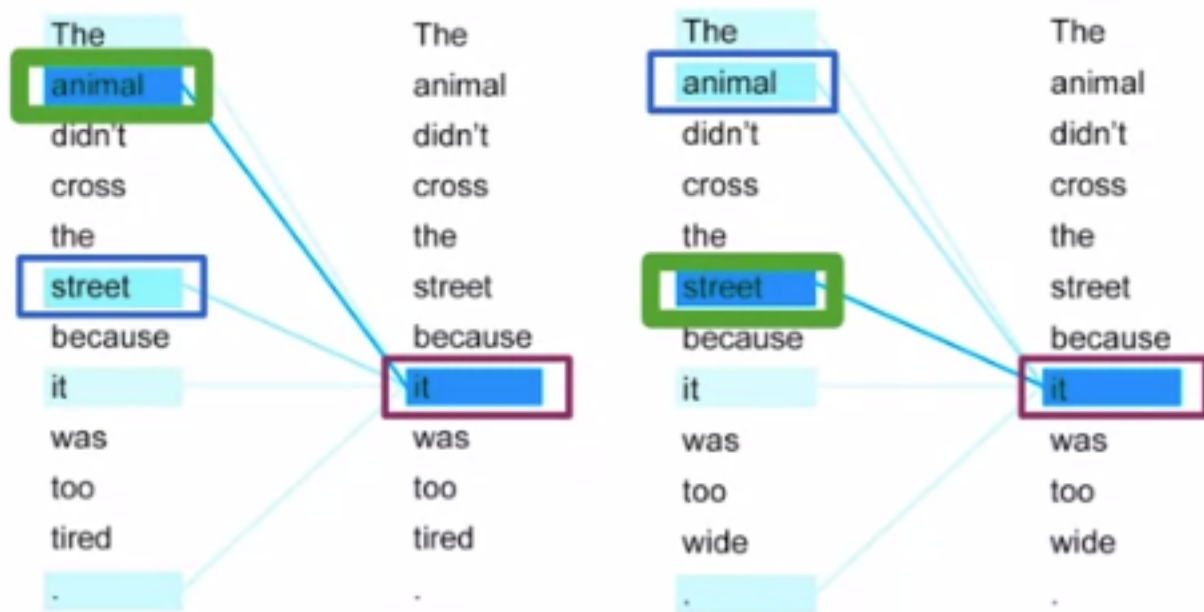


image ©
[\(Transformer: A Novel Neural Network Architecture for Language Understanding.\)](#)

Nearest Neighbors

Course:

Natural Language Processing with Classification and Vector Spaces

Lessons:

- KNN
- Hash Tables and Hash Functions
- Locality Sensitive Hashing
- Multiple Planes

Nearest Neighbors

Course:

Natural Language Processing with Classification and Vector Spaces

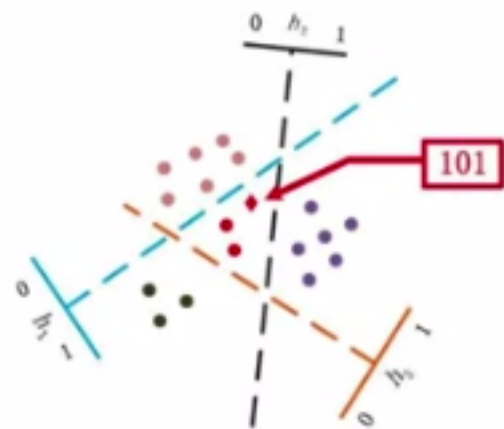
Lessons:

- KNN
- Hash Tables and Hash Functions
- Locality Sensitive Hashing
- Multiple Planes

Nearest Neighbors

Compute the nearest neighbor to q among vectors $\{k_1, \dots, k_n\}$

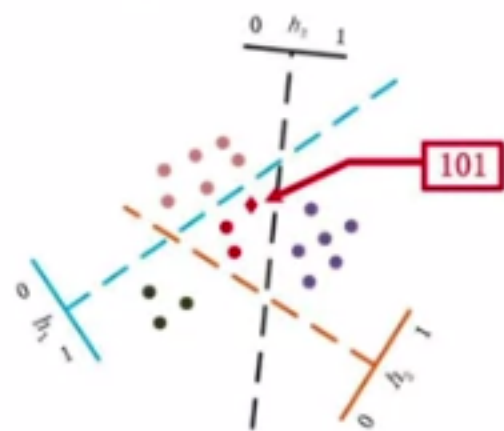
- Attention computes $d(q, k_i)$ for i from 1 to n which can be slow
- Faster *approximate* uses locality sensitive hashing (LSH)



Nearest Neighbors

Compute the nearest neighbor to q among vectors $\{k_1, \dots, k_n\}$

- Attention computes $d(q, k_i)$ for i from 1 to n which can be slow
- Faster *approximate* uses locality sensitive hashing (LSH)
- Locality sensitive: if q is close to k_i :
 $\text{hash}(q) == \text{hash}(k_i)$



Nearest Neighbors

Compute the nearest neighbor to q among vectors $\{k_1, \dots, k_n\}$

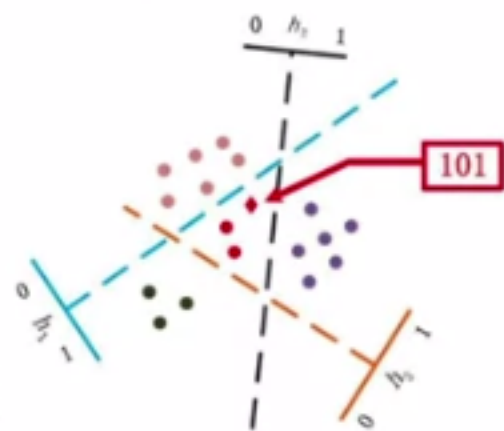
- Attention computes $d(q, k_i)$ for i from 1 to n which can be slow
- Faster *approximate* uses locality sensitive hashing (LSH)

- Locality sensitive: if q is close to k_i :

$$\text{hash}(q) == \text{hash}(k_i)$$

- Achieve by randomly cutting space

$$\text{hash}(x) = \text{sign}(xR) \quad R: [d, n_hash_bins]$$

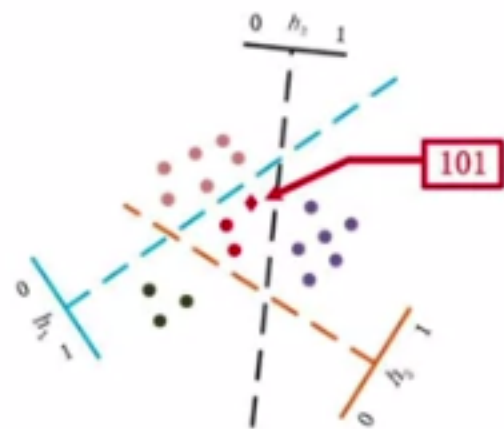


LSH Attention

Standard Attention:

$$A(Q, K, V) = \text{softmax}(QK^T)V$$

LSH Attention:



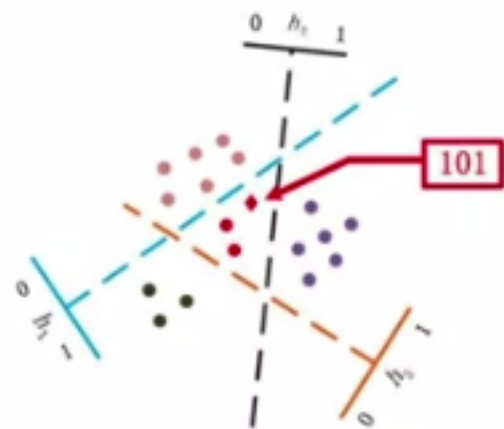
LSH Attention

Standard Attention:

$$A(Q, K, V) = \text{softmax}(QK^T)V$$

LSH Attention:

- Hash Q and K



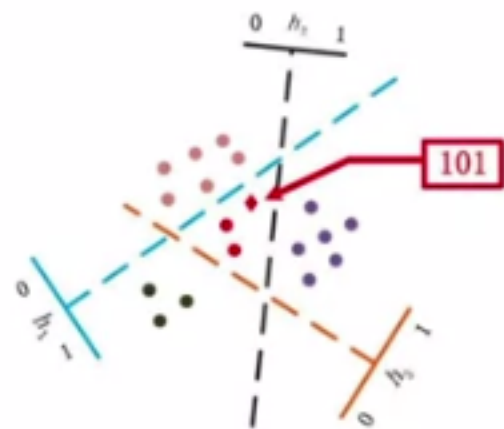
LSH Attention

Standard Attention:

$$A(Q, K, V) = \text{softmax}(QK^T)V$$

LSH Attention:

- Hash Q and K
- Standard attention within same-hash bins



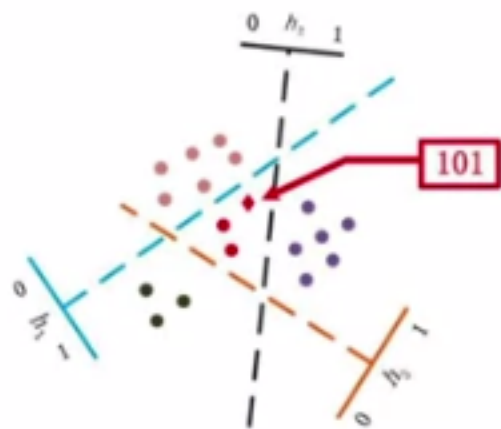
LSH Attention

Standard Attention:

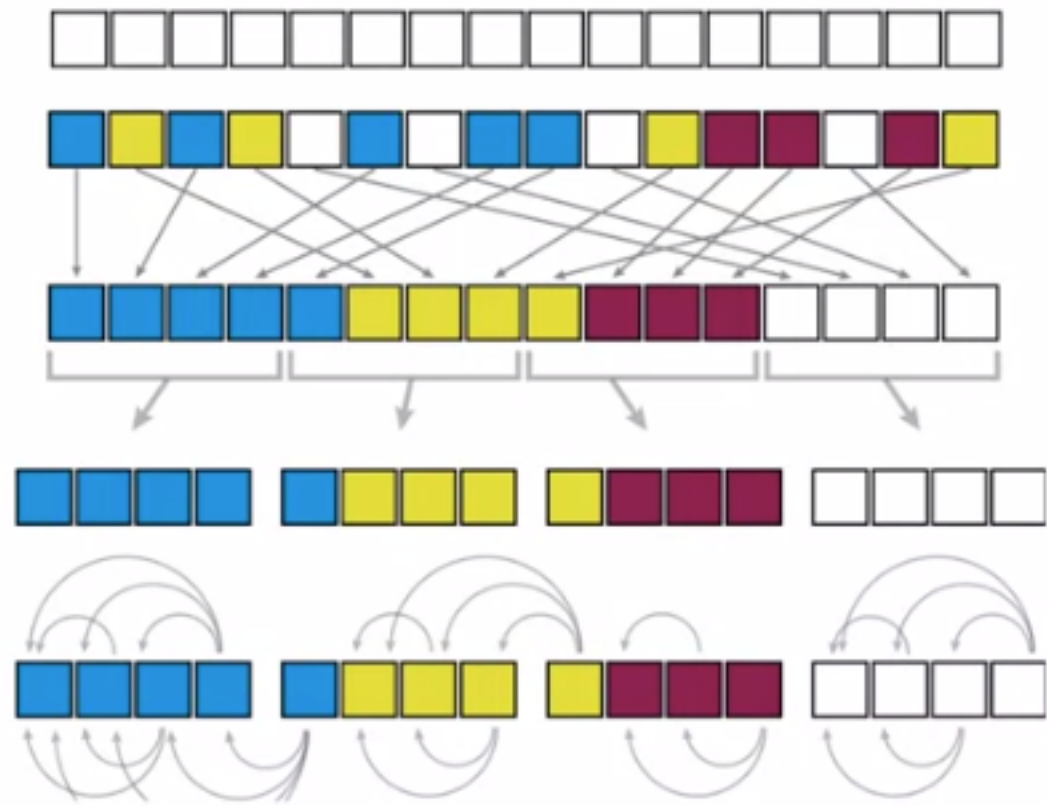
$$A(Q, K, V) = \text{softmax}(QK^T)V$$

LSH Attention:

- Hash Q and K
- Standard attention within same-hash bins
- Repeat a few times to increase probability of key in the same bin

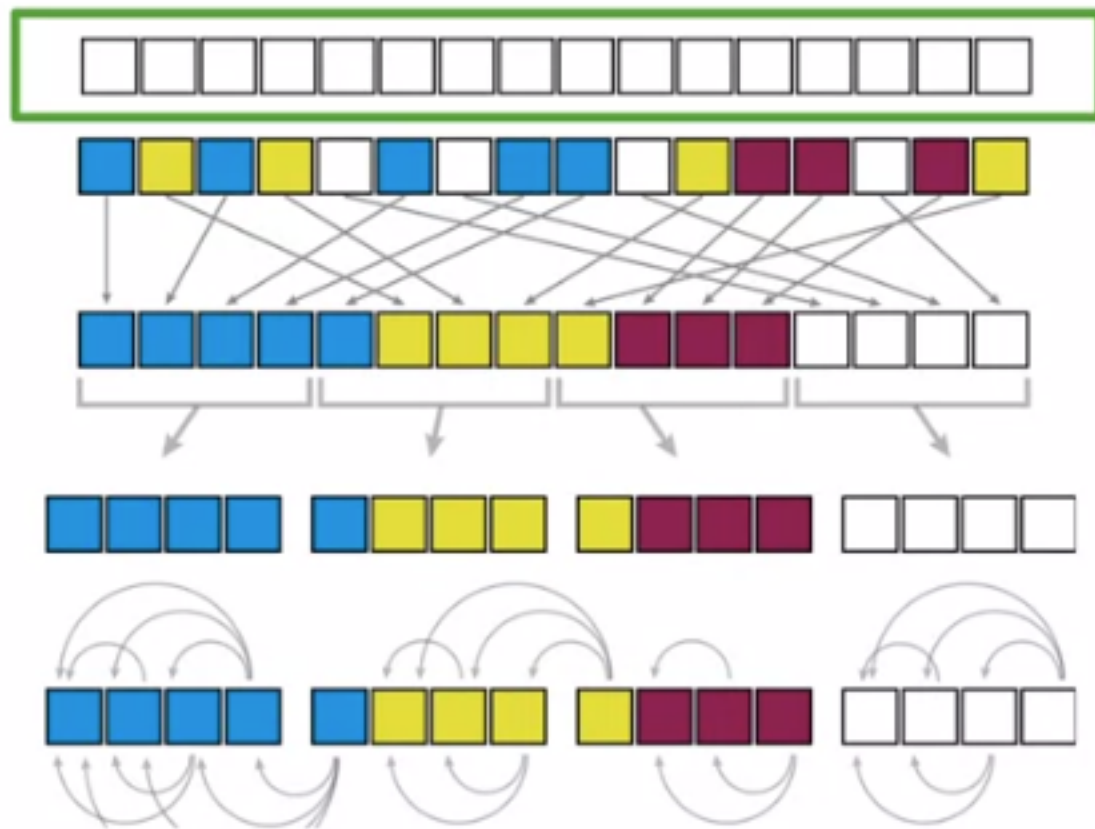


LSH Attention



LSH Attention

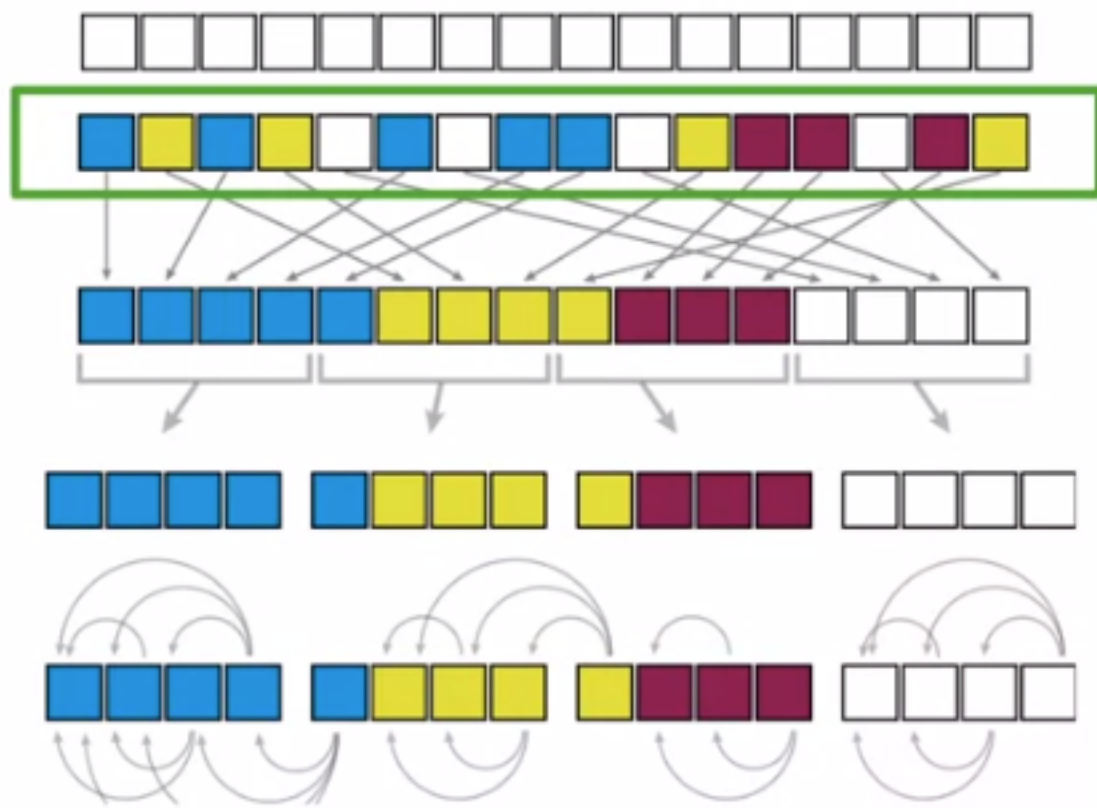
Sequence of Queries = Keys



LSH Attention

Sequence of Queries = Keys

LSH bucketing

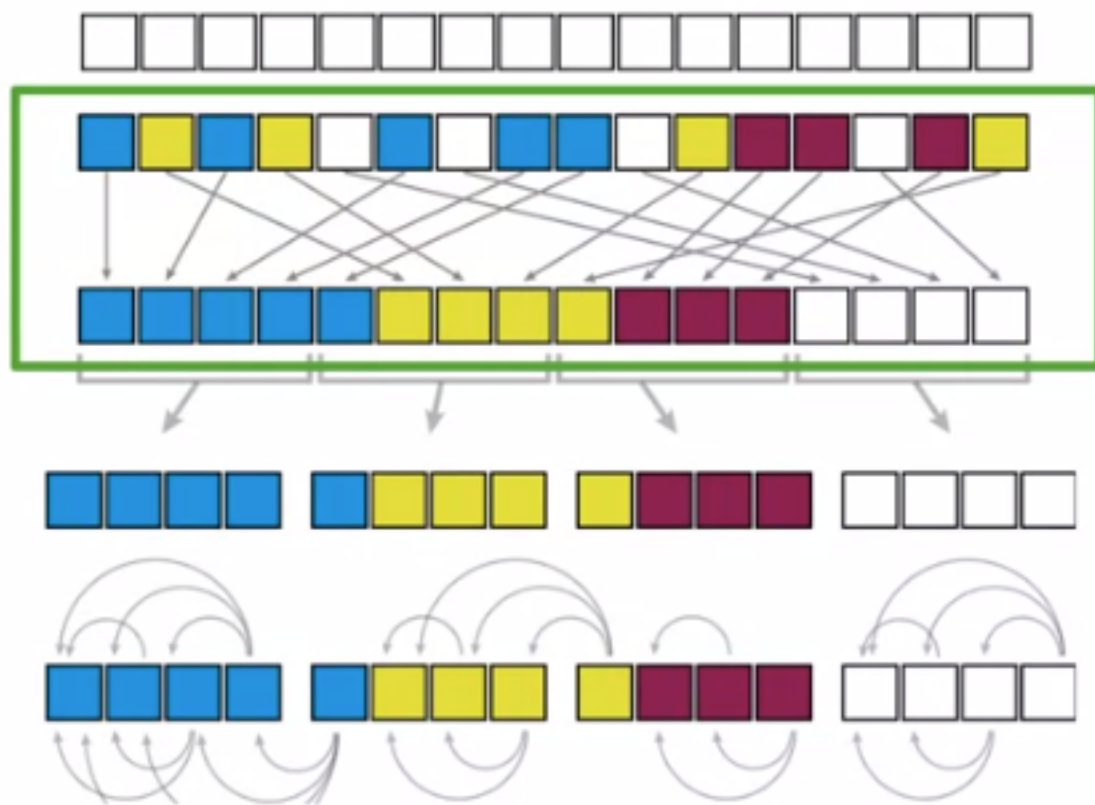


LSH Attention

Sequence of Queries = Keys

LSH bucketing

Sort by LSH bucket



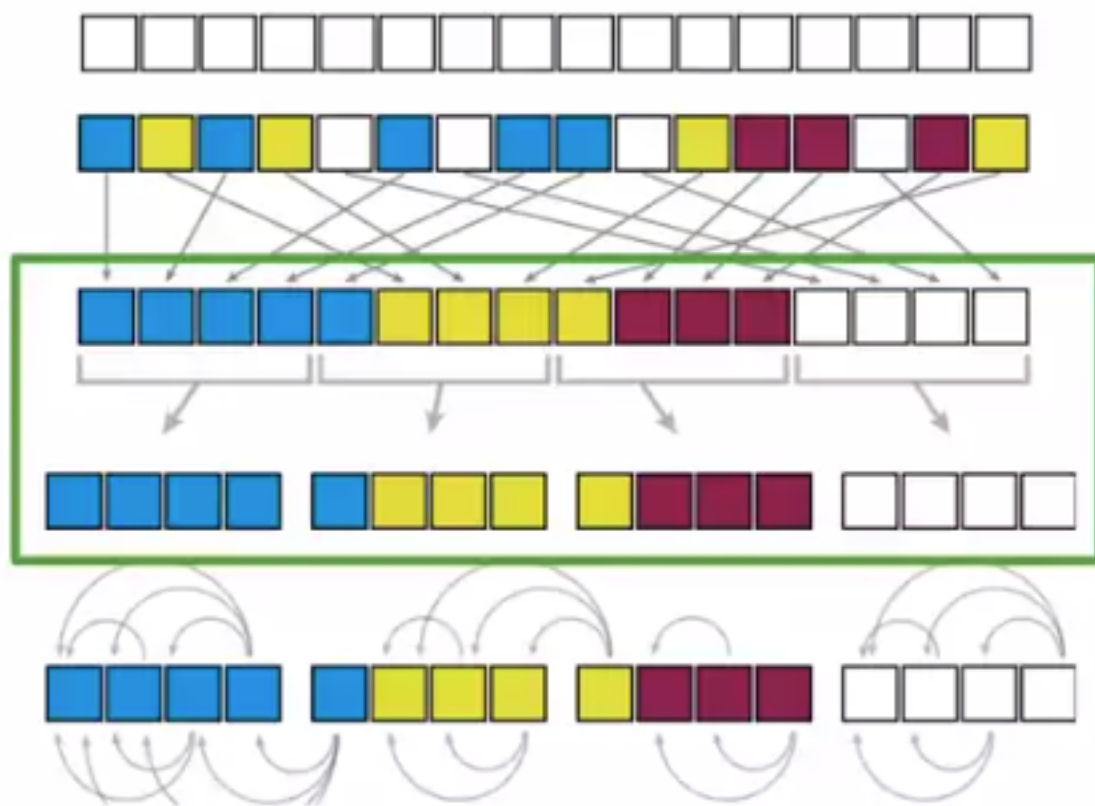
LSH Attention

Sequence of Queries = Keys

LSH bucketing

Sort by LSH bucket

Chunk sorted sequence
to parallelize



LSH Attention

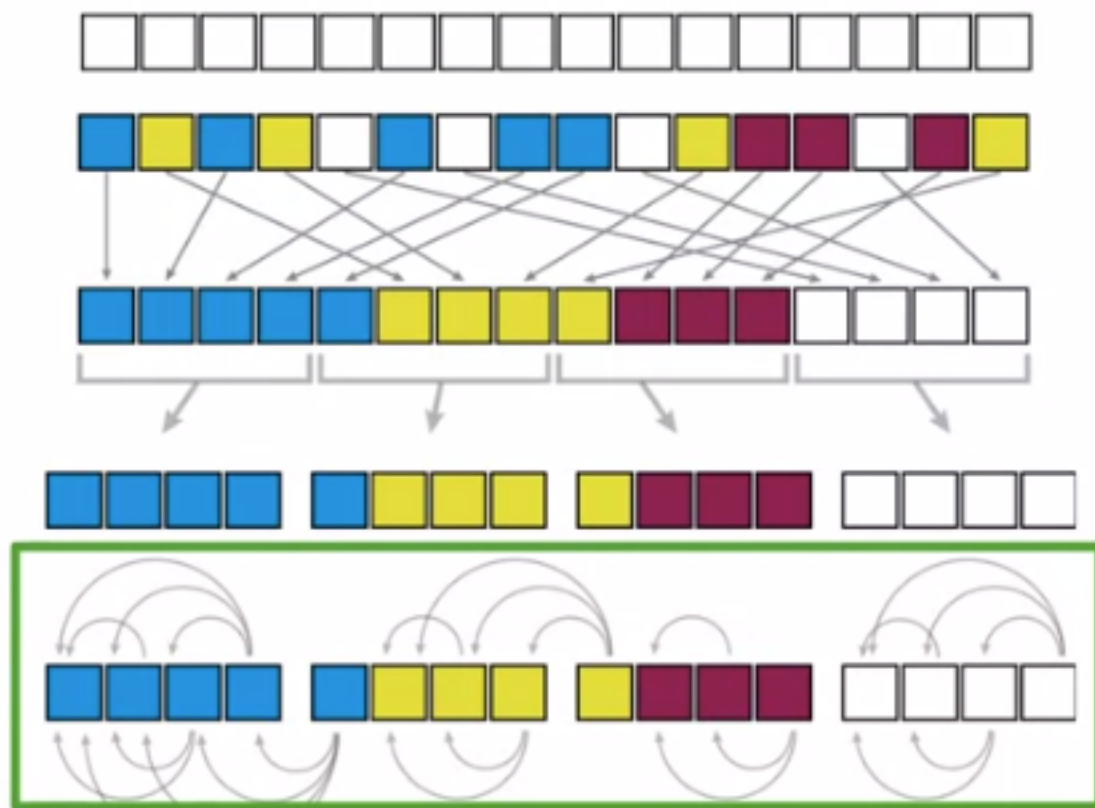
Sequence of Queries = Keys

LSH bucketing

Sort by LSH bucket

Chunk sorted sequence
to parallelize

Attend within same bucket of
own chunk and previous chunk



Memory Efficiency



Memory Efficiency



input



$L = 1$ million tokens

Memory Efficiency



$d_{\text{model}} = 512$



input

$L = 1 \text{ million tokens}$

Memory Efficiency



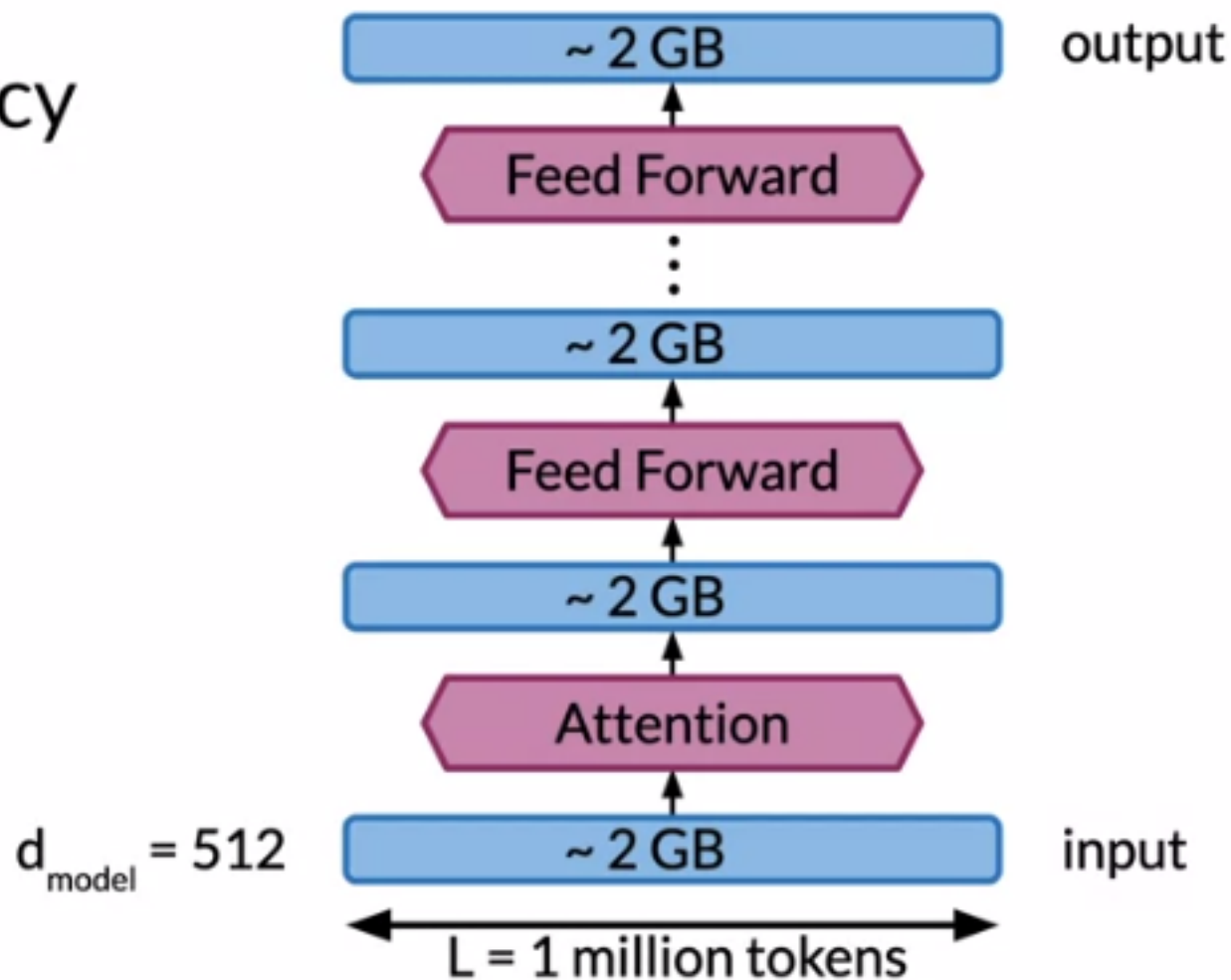
$d_{\text{model}} = 512$

~ 2 GB

input

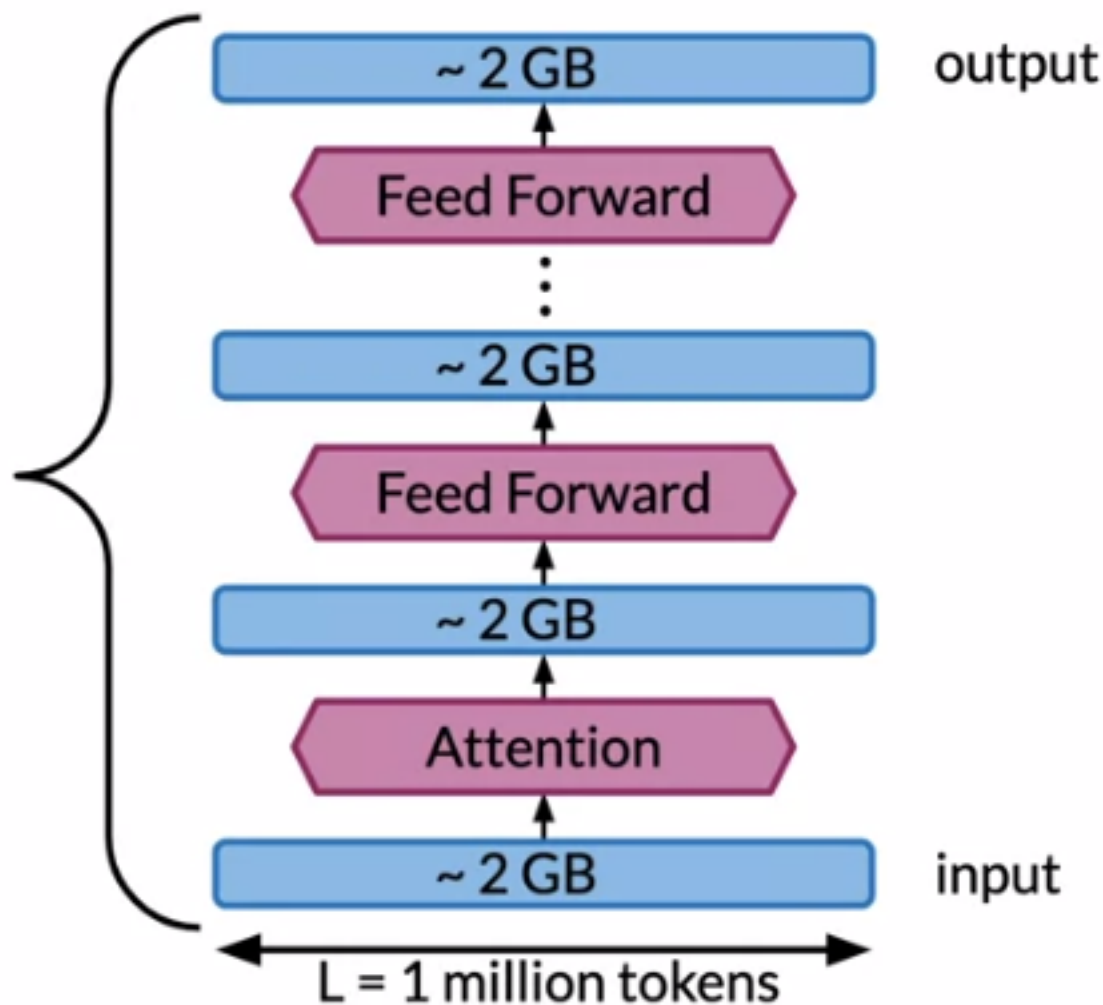
$L = 1 \text{ million tokens}$

Memory Efficiency



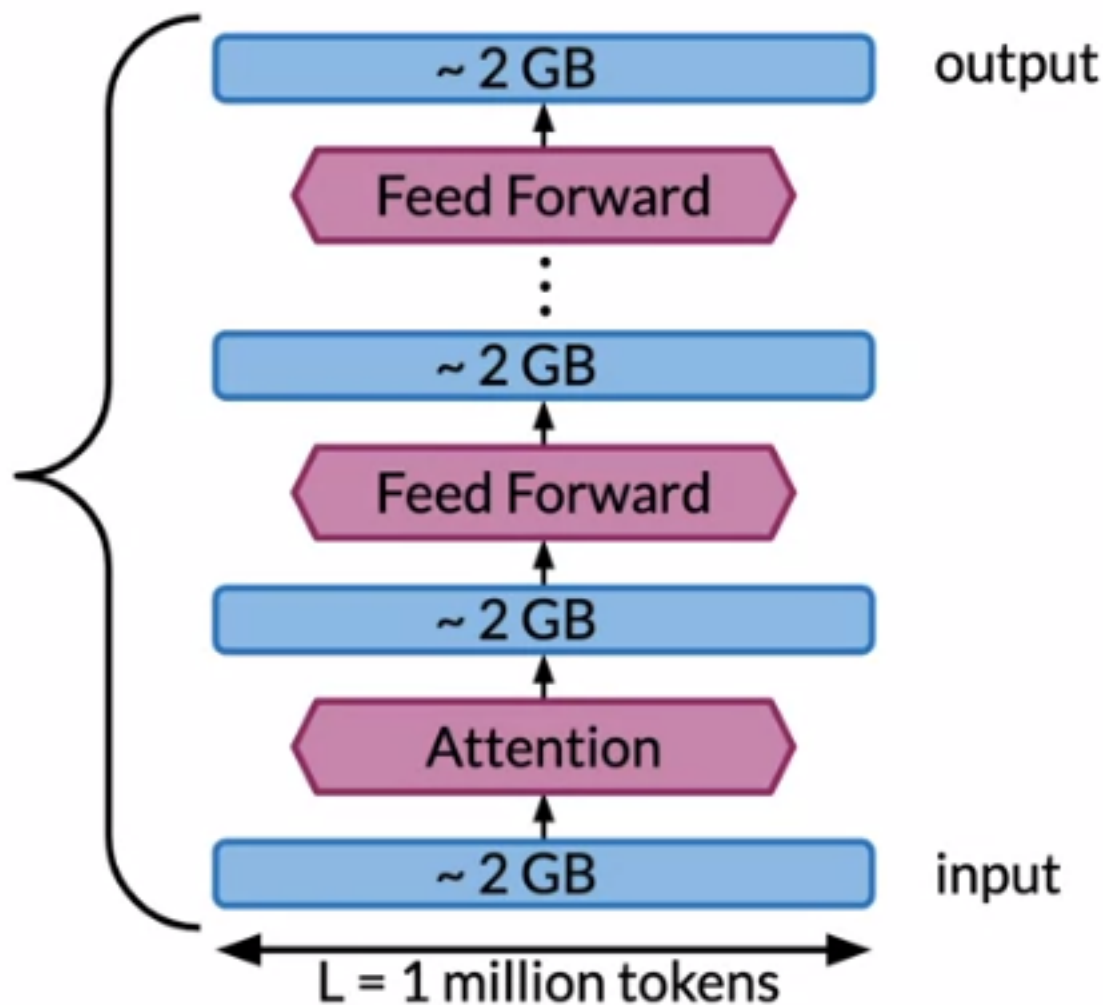
Memory Efficiency

12 x Attention
12 x Feed-Forward



Memory Efficiency

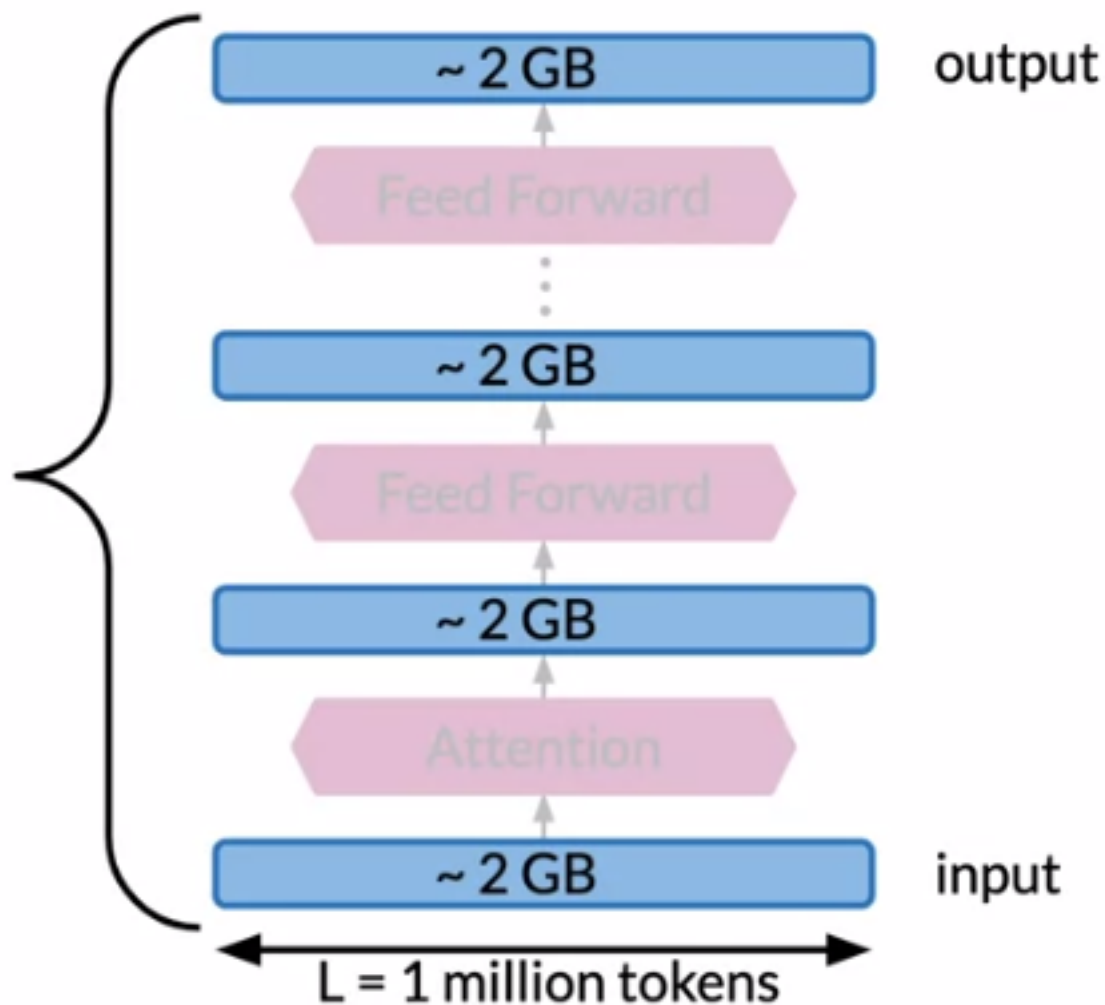
12 x Attention
12 x Feed-Forward



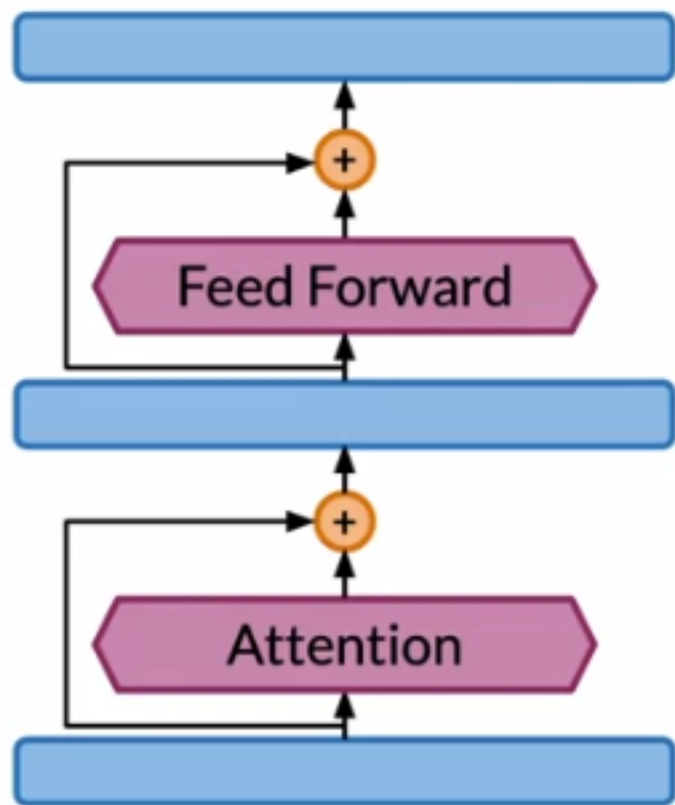
Memory Efficiency

12 x Attention
12 x Feed-Forward

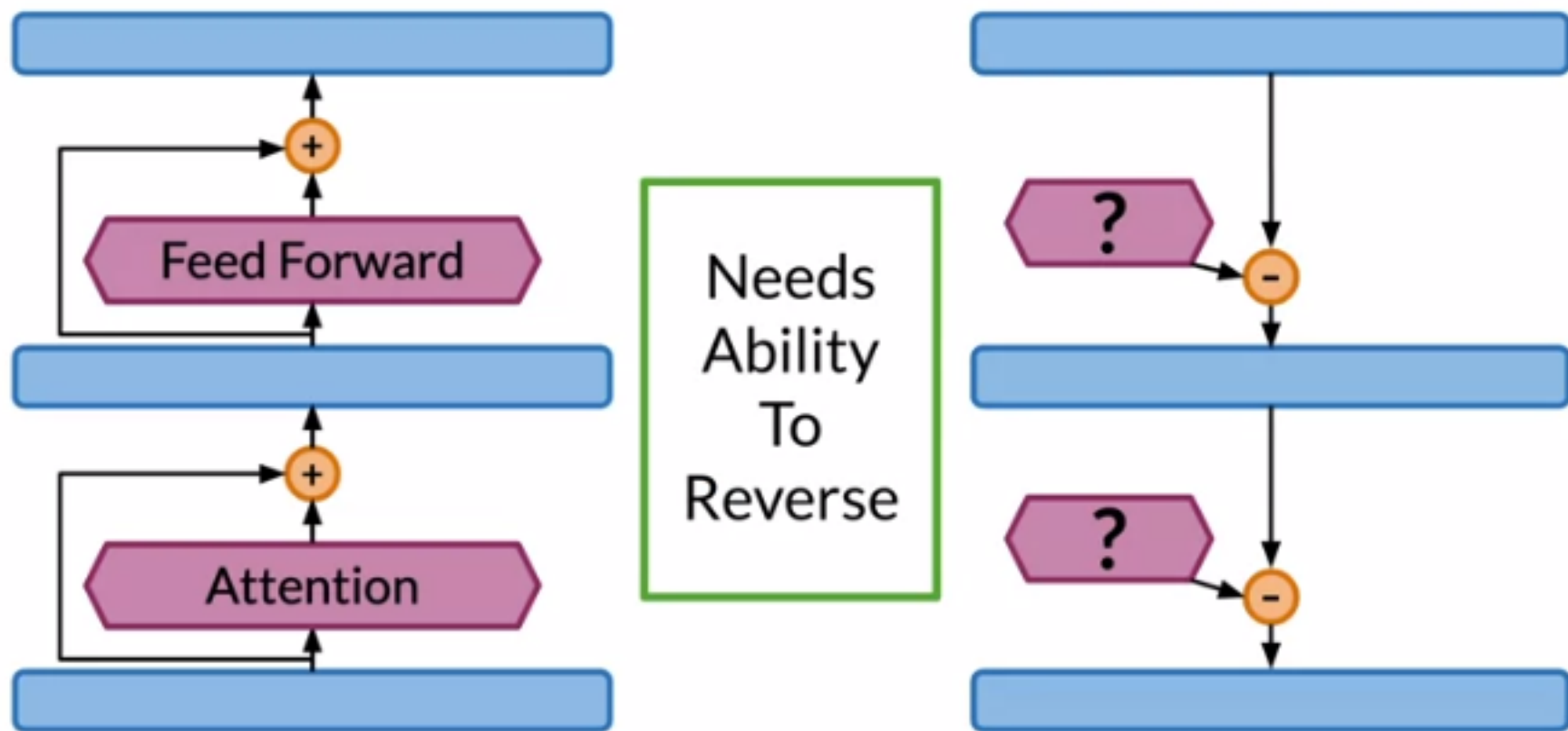
50 GB total



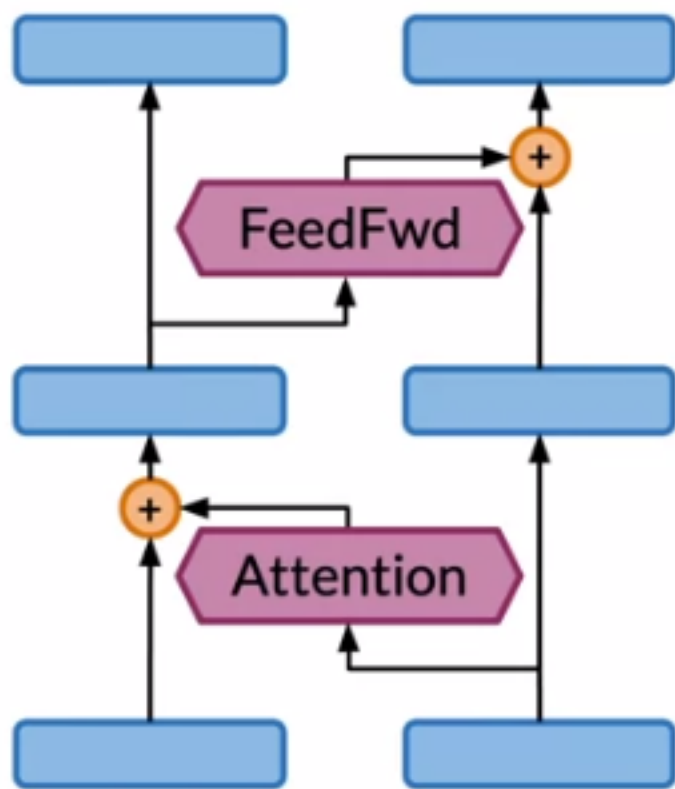
Residual Blocks in Transformer



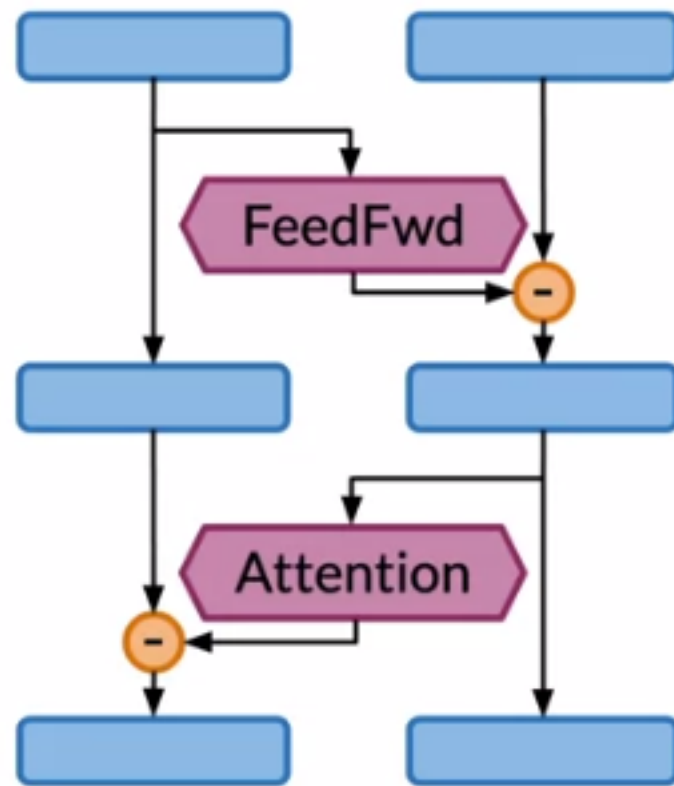
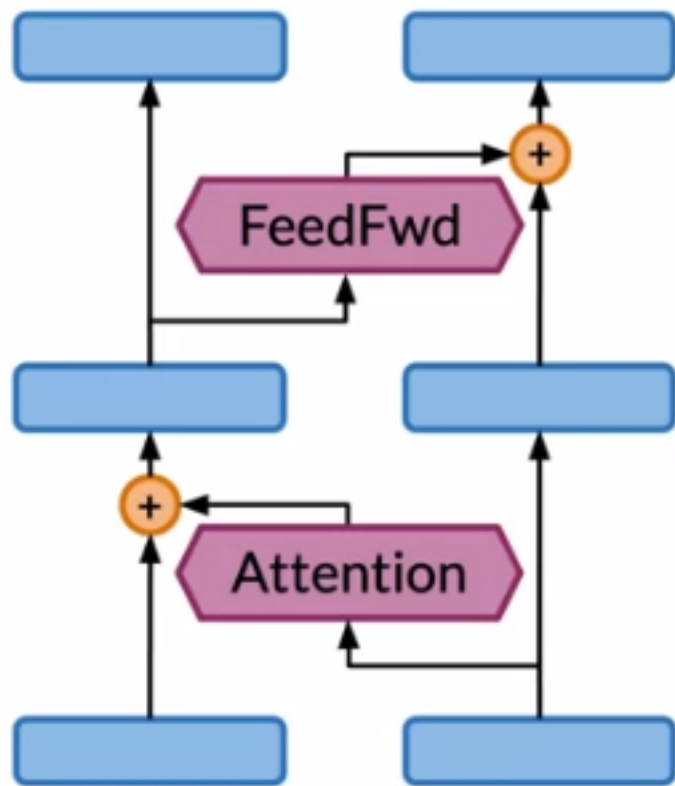
Residual Blocks in Transformer



Reversible Residual Blocks



Reversible layers



Reversible layers equations

Standard Transformer:

$$y_a = x + \text{Attention}(x)$$

$$y_b = y_a + \text{FeedFwd}(y_a)$$

Reversible layers equations

Standard Transformer:

$$y_a = x + \text{Attention}(x)$$

$$y_b = y_a + \text{FeedFwd}(y_a)$$

Reversible:

$$y_1 = x_1 + \text{Attention}(x_2)$$

$$y_2 = x_2 + \text{FeedFwd}(y_1)$$

Reversible layers equations

Standard Transformer:

$$y_a = x + \text{Attention}(x)$$

$$y_b = y_a + \text{FeedFwd}(y_a)$$

Reversible:

$$y_1 = x_1 + \text{Attention}(x_2)$$

$$y_2 = x_2 + \text{FeedFwd}(y_1)$$

Recompute x_1, x_2 from y_1, y_2 :

$$x_1 = y_1 - \text{Attention}(x_2)$$

$$x_2 = y_2 - \text{FeedFwd}(y_1)$$

Reversible layers equations

Standard Transformer:

$$y_a = x + \text{Attention}(x)$$

$$y_b = y_a + \text{FeedFwd}(y_a)$$

Reversible:

$$y_1 = x_1 + \text{Attention}(x_2)$$

$$y_2 = x_2 + \text{FeedFwd}(y_1)$$

Recompute x_1, x_2 from y_1, y_2 :

$$x_1 = y_1 - \text{Attention}(x_2)$$

$$x_2 = y_2 - \text{FeedFwd}(y_1)$$

Reversible layers equations

Standard Transformer:

$$y_a = x + \text{Attention}(x)$$

$$y_b = y_a + \text{FeedFwd}(y_a)$$

Reversible:

$$y_{\boxed{1}} = x_{\boxed{1}} + \text{Attention}(x_{\boxed{2}})$$

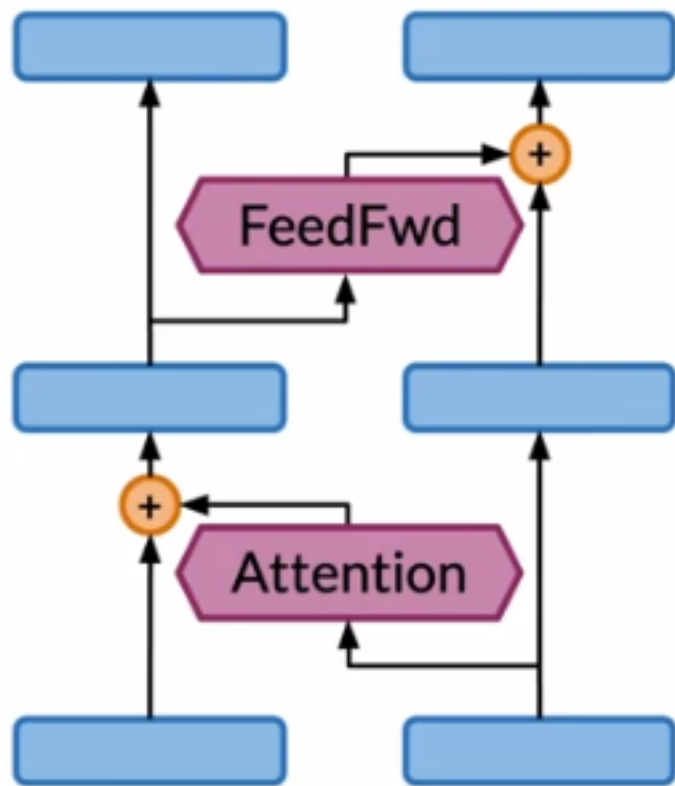
$$y_{\boxed{2}} = x_{\boxed{2}} + \text{FeedFwd}(y_{\boxed{1}})$$

Recompute x_1, x_2 from y_1, y_2 :

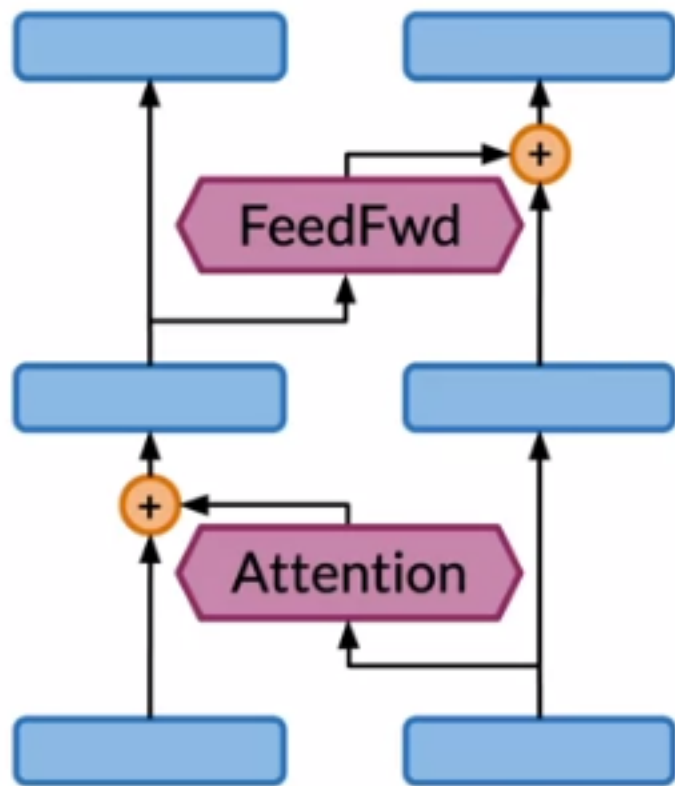
$$x_1 = y_1 - \text{Attention}(x_2)$$

$$x_2 = y_2 - \text{FeedFwd}(y_1)$$

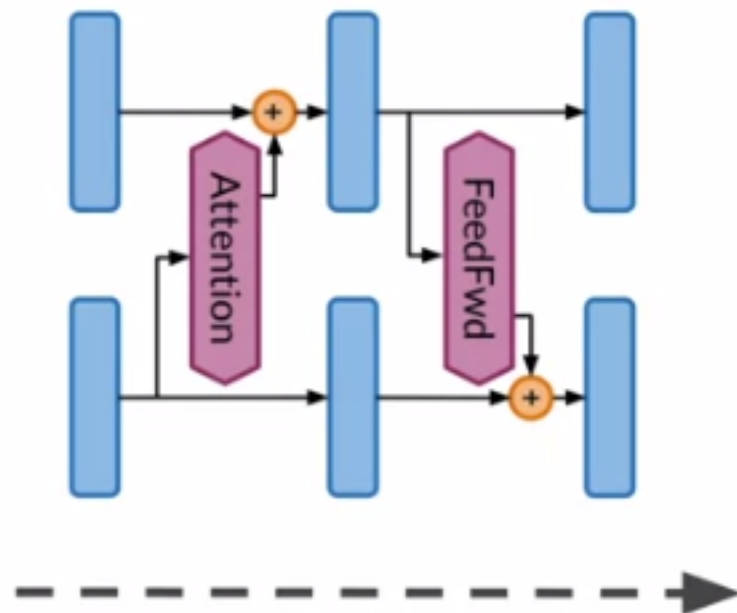
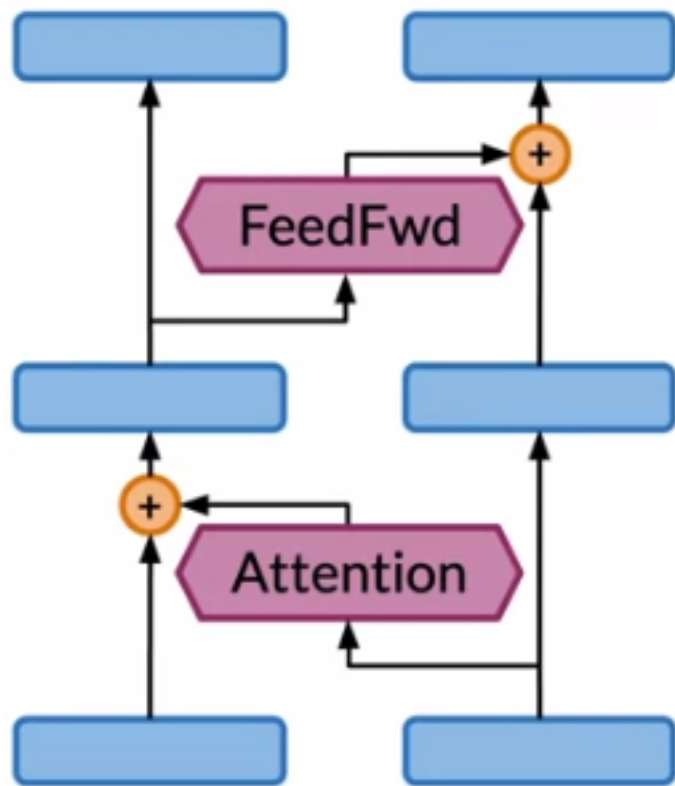
Reversible layers equations



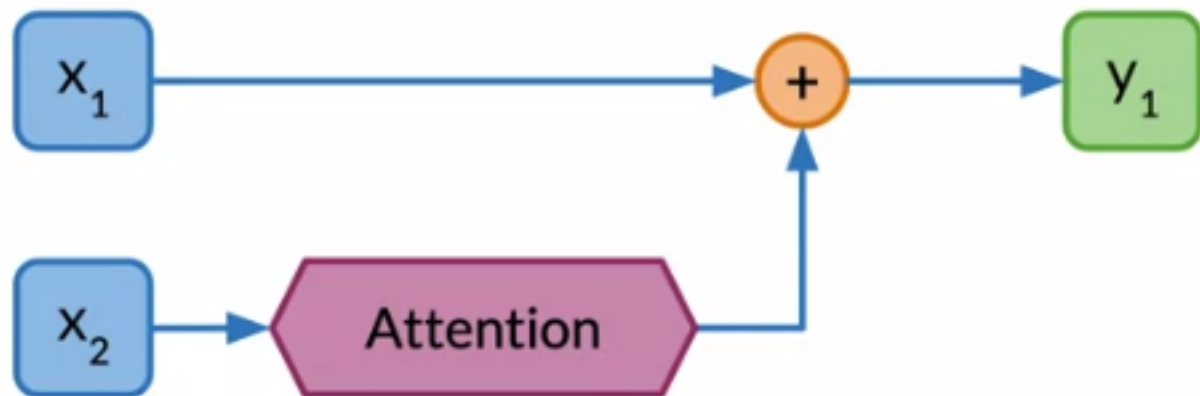
Reversible layers equations



Reversible layers equations



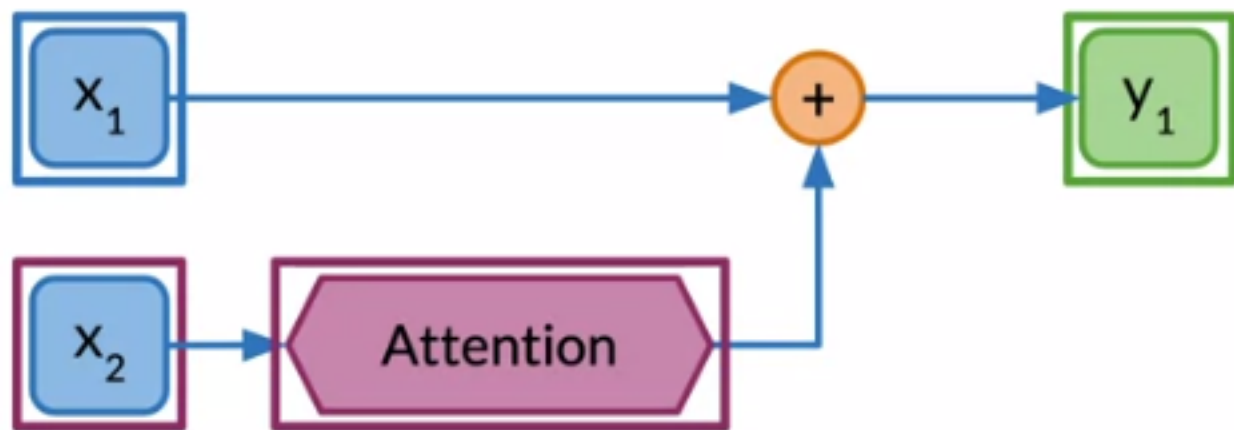
Reversible layers equations



$$y_1 = x_1 + \text{Attention}(x_2)$$

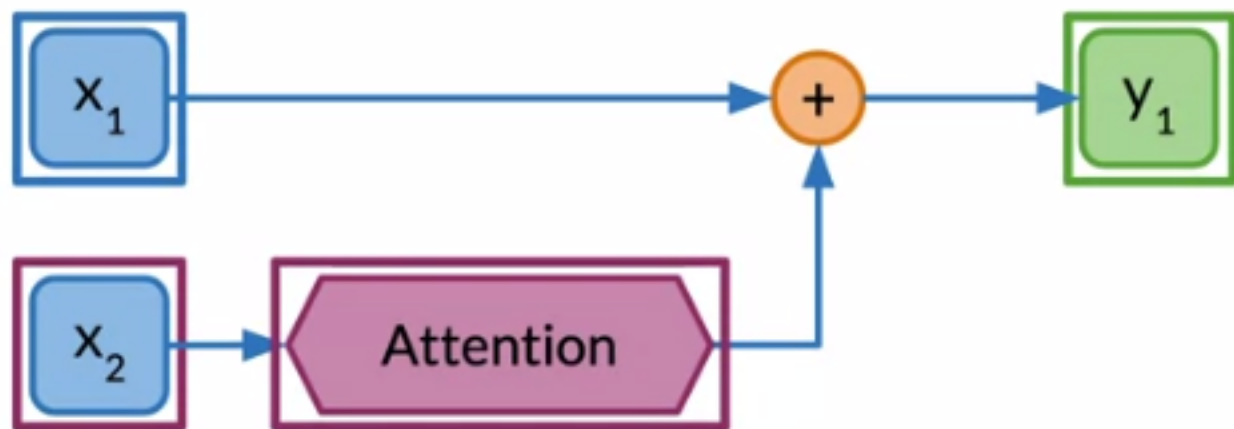
$$y_2 = x_2 + \text{FeedFwd}(y_1)$$

Reversible layers equations



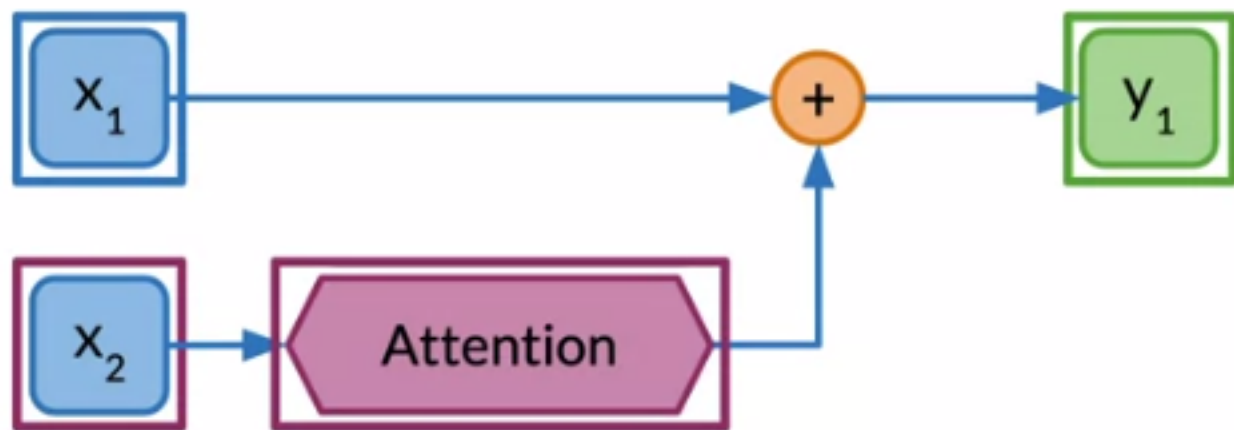
$$y_1 = x_1 + \text{Attention}(x_2)$$
$$y_2 = x_2 + \text{FeedFwd}(y_1)$$

Reversible layers equations



$$y_1 = x_1 + \text{Attention}(x_2)$$
$$y_2 = x_2 + \text{FeedFwd}(y_1)$$

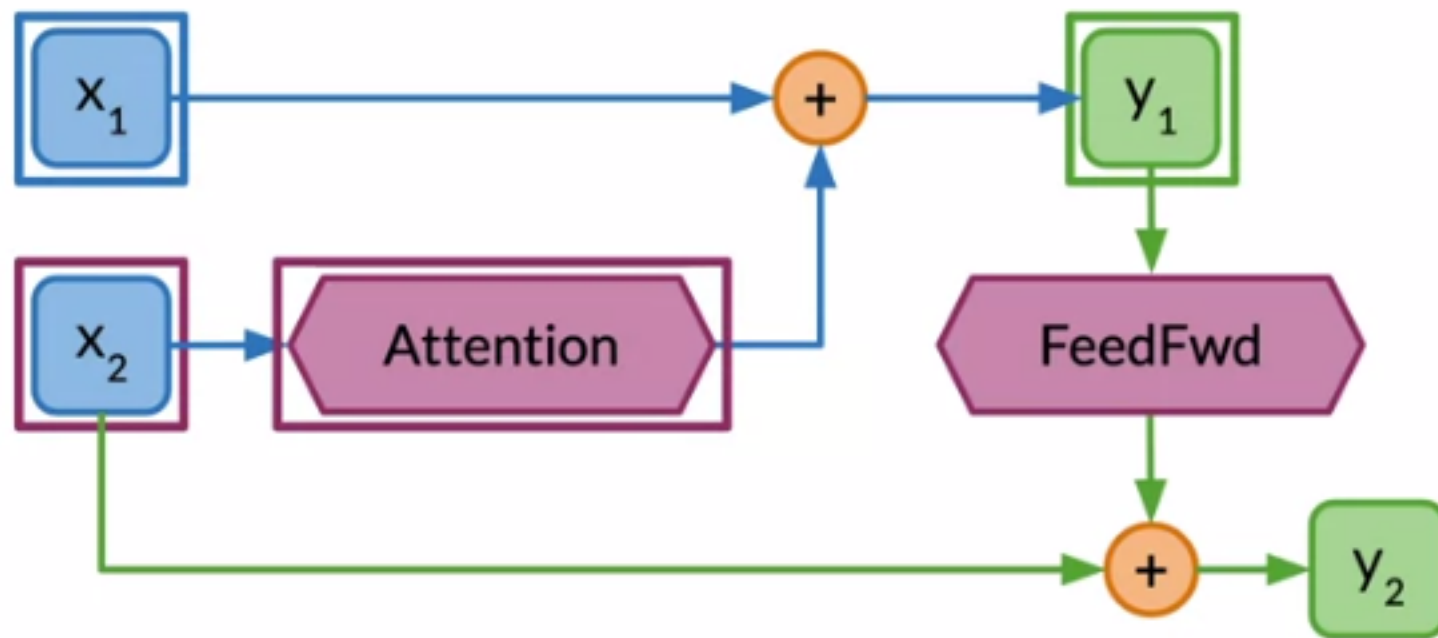
Reversible layers equations



$$y_1 = x_1 + \text{Attention}(x_2)$$
$$y_2 = x_2 + \text{FeedFwd}(y_1)$$

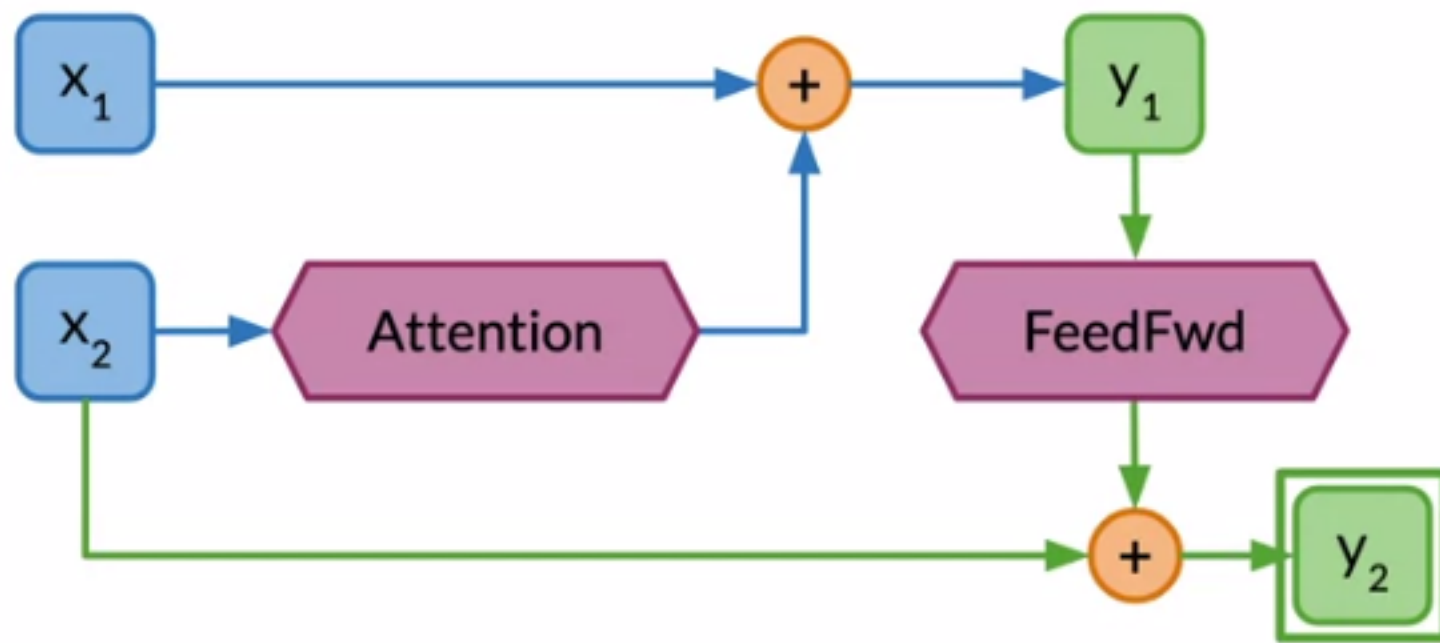
A red arrow points from the circled y_1 in the first equation to the $\text{FeedFwd}(y_1)$ term in the second equation.

Reversible layers equations



$$y_1 = x_1 + \text{Attention}(x_2)$$
$$y_2 = x_2 + \text{FeedFwd}(y_1)$$

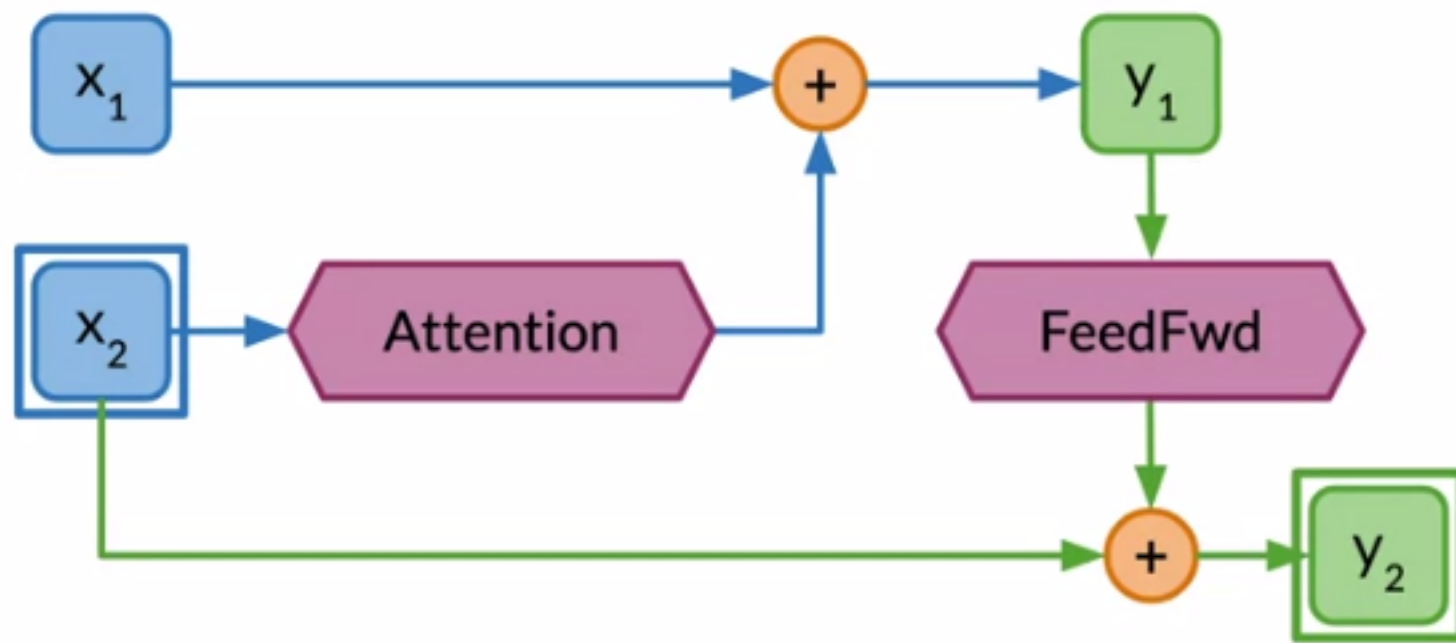
Reversible layers equations



$$y_1 = x_1 + \text{Attention}(x_2)$$

$$\boxed{y_2} = x_2 + \text{FeedFwd}(y_1)$$

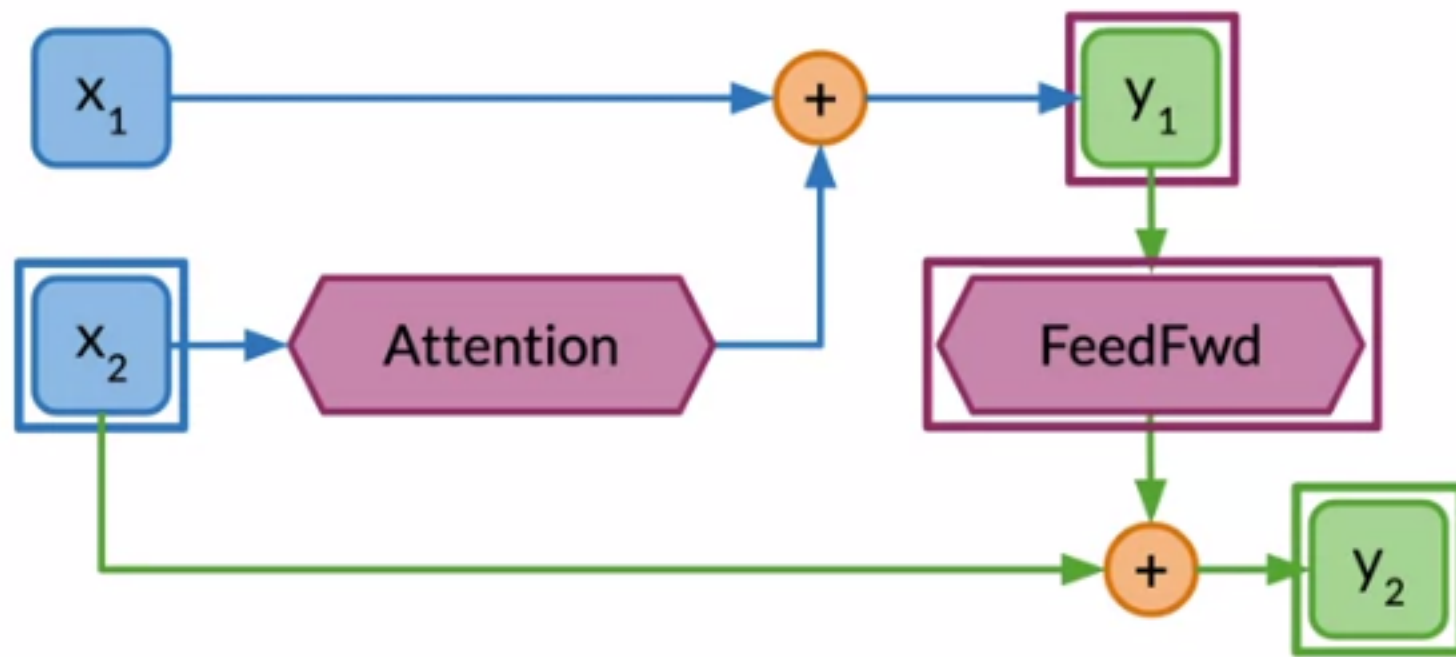
Reversible layers equations



$$y_1 = x_1 + \text{Attention}(x_2)$$

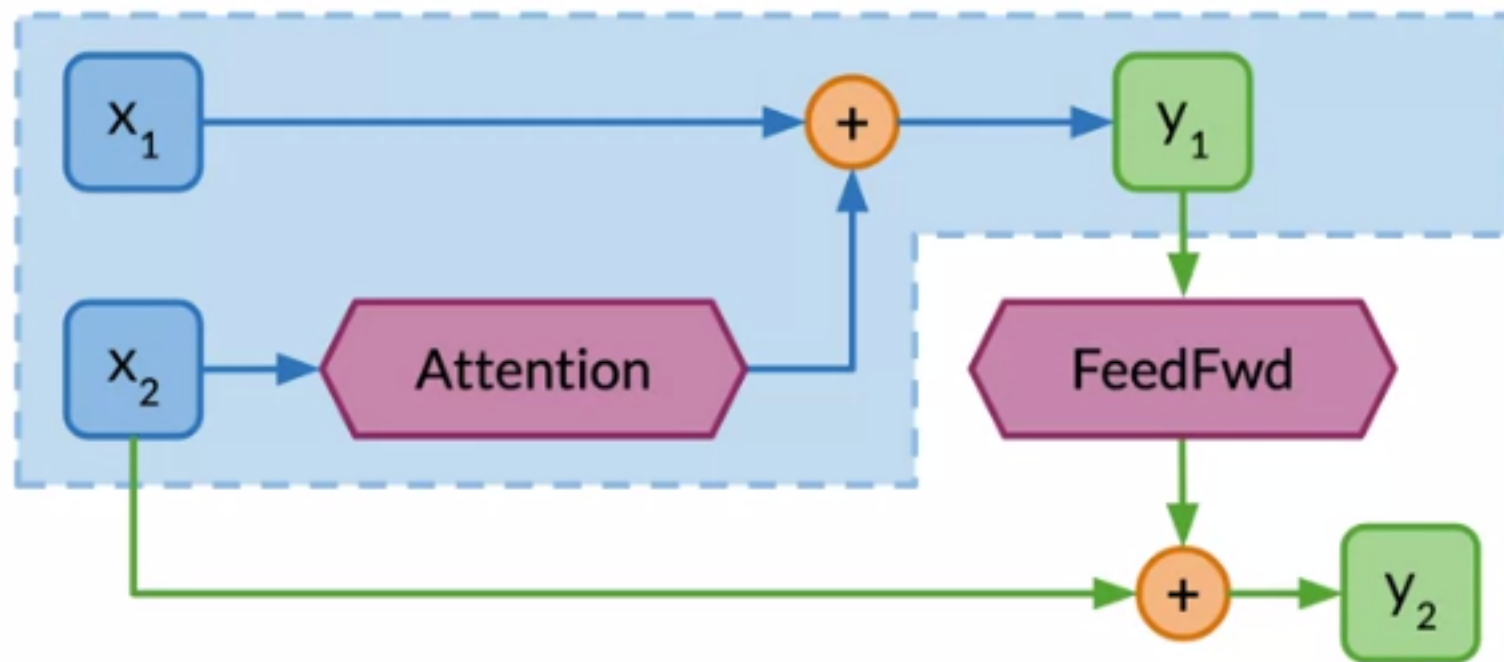
$$y_2 = x_2 + \text{FeedFwd}(y_1)$$

Reversible layers equations



$$y_1 = x_1 + \text{Attention}(x_2)$$
$$y_2 = x_2 + \text{FeedFwd}(y_1)$$

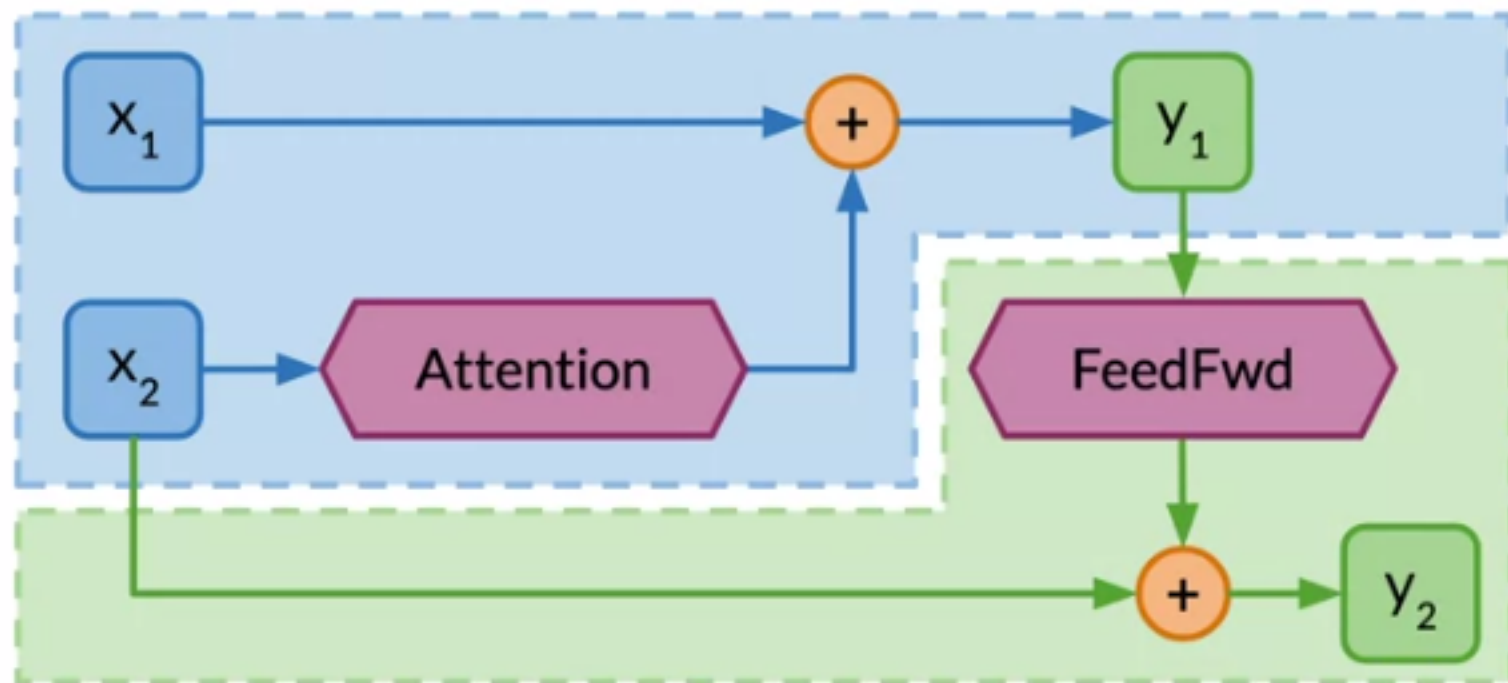
Reversible layers equations



$$y_1 = x_1 + \text{Attention}(x_2)$$

$$y_2 = x_2 + \text{FeedFwd}(y_1)$$

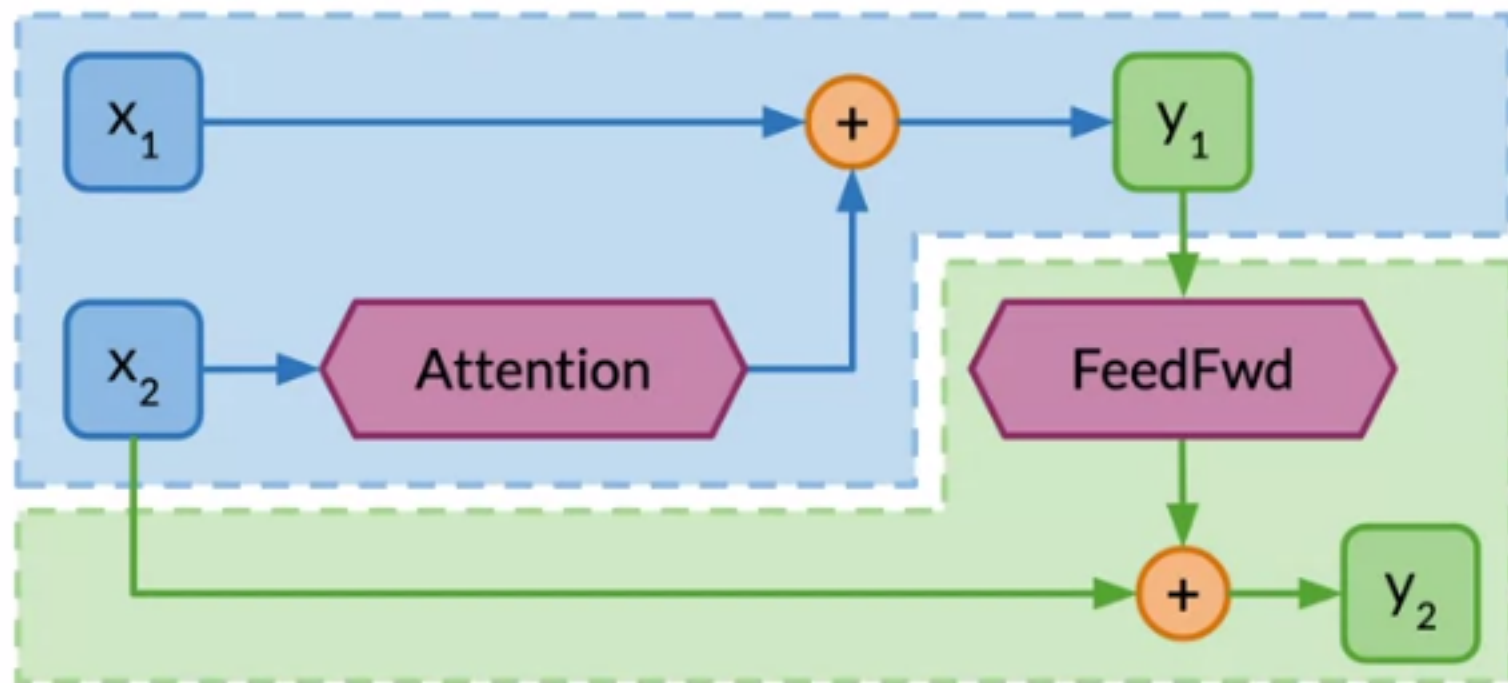
Reversible layers equations



$$y_1 = x_1 + \text{Attention}(x_2)$$

$$y_2 = x_2 + \text{FeedFwd}(y_1)$$

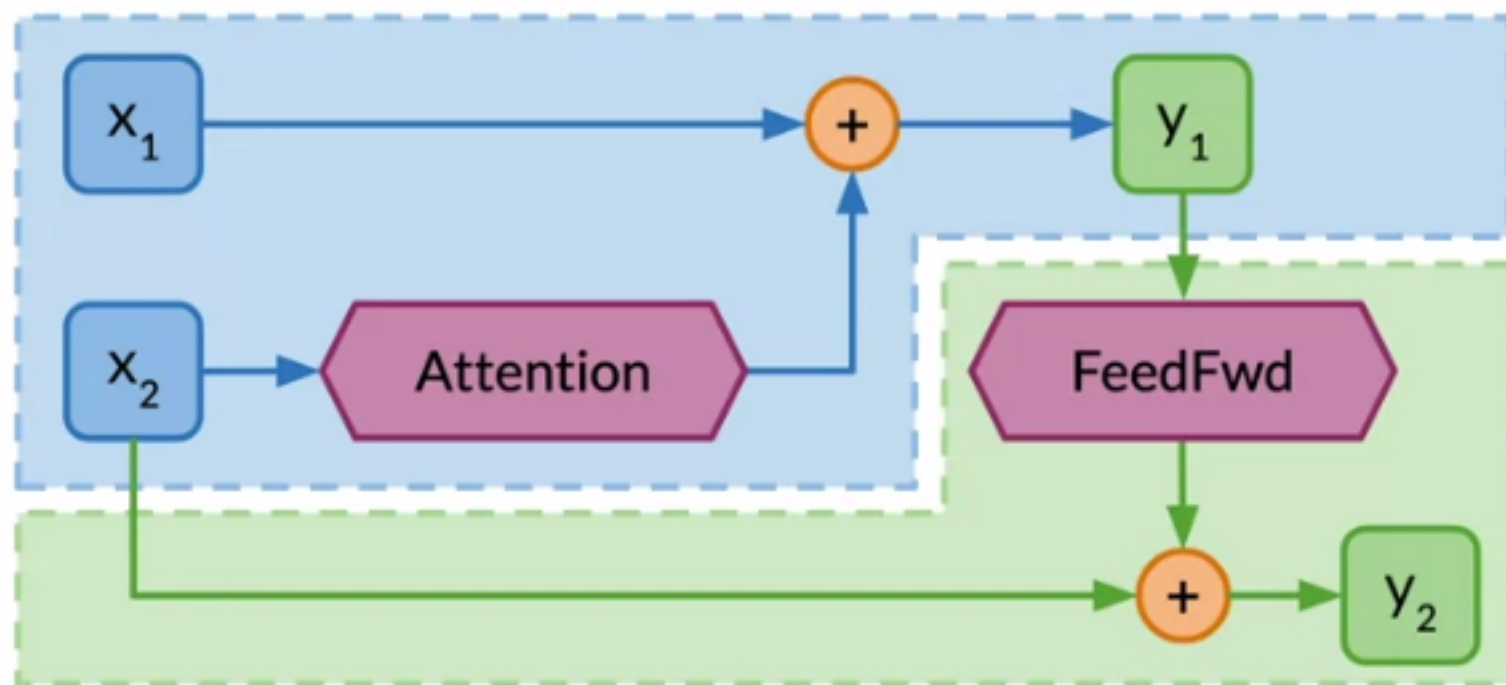
Reversible layers equations



➔

$$y_1 = x_1 + \text{Attention}(x_2)$$
$$y_2 = x_2 + \text{FeedFwd}(y_1)$$

Reversible layers equations



➔

$$y_1 = x_1 + \text{Attention}(x_2)$$

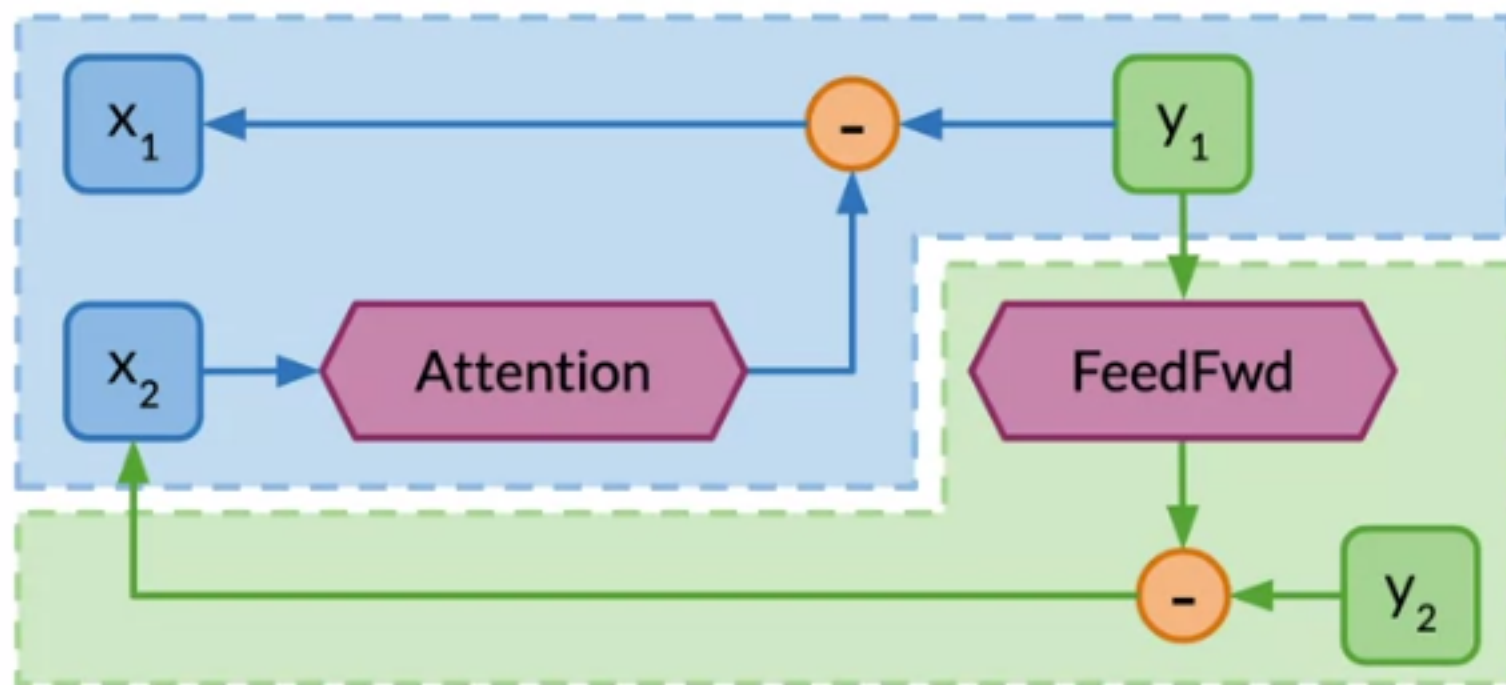
$$y_2 = x_2 + \text{FeedFwd}(y_1)$$

➔

$$x_2 = y_2 - \text{FeedFwd}(y_1)$$

$$x_1 = y_1 - \text{Attention}(x_2)$$

Reversible layers equations



➔

$$y_1 = x_1 + \text{Attention}(x_2)$$

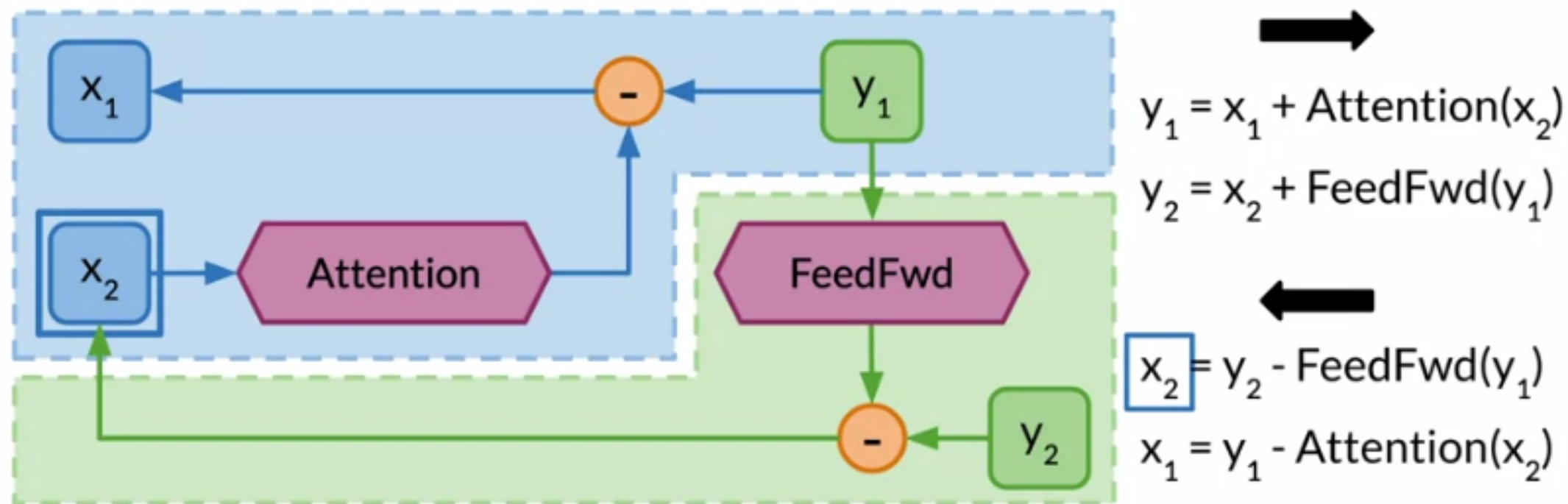
$$y_2 = x_2 + \text{FeedFwd}(y_1)$$

➔

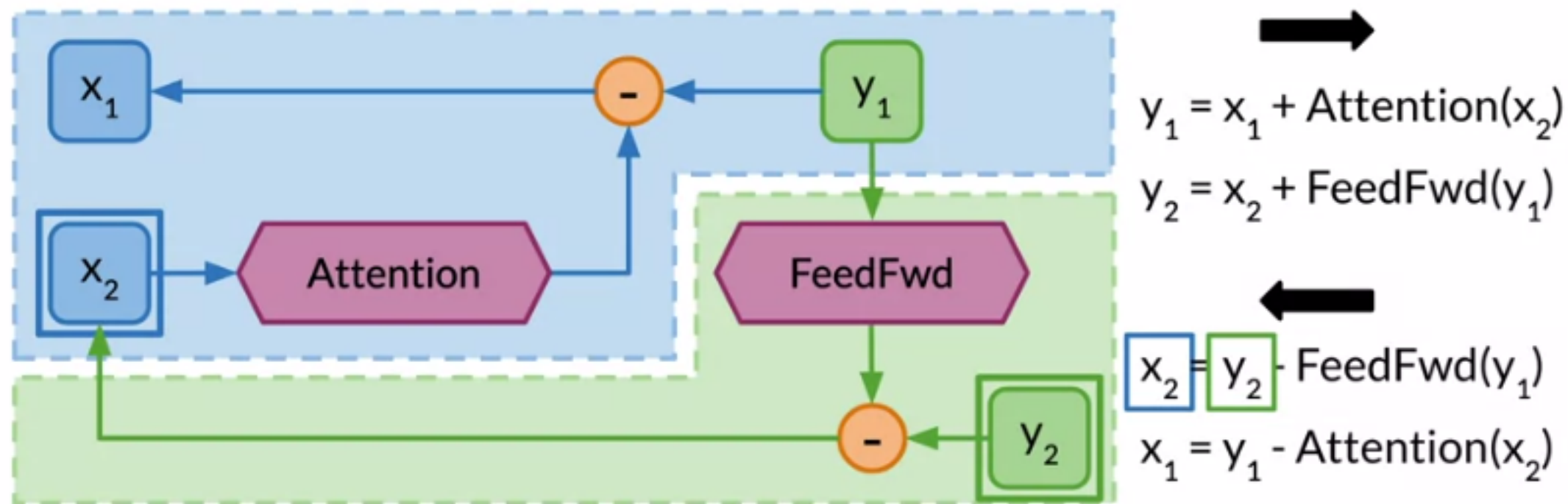
$$x_2 = y_2 - \text{FeedFwd}(y_1)$$

$$x_1 = y_1 - \text{Attention}(x_2)$$

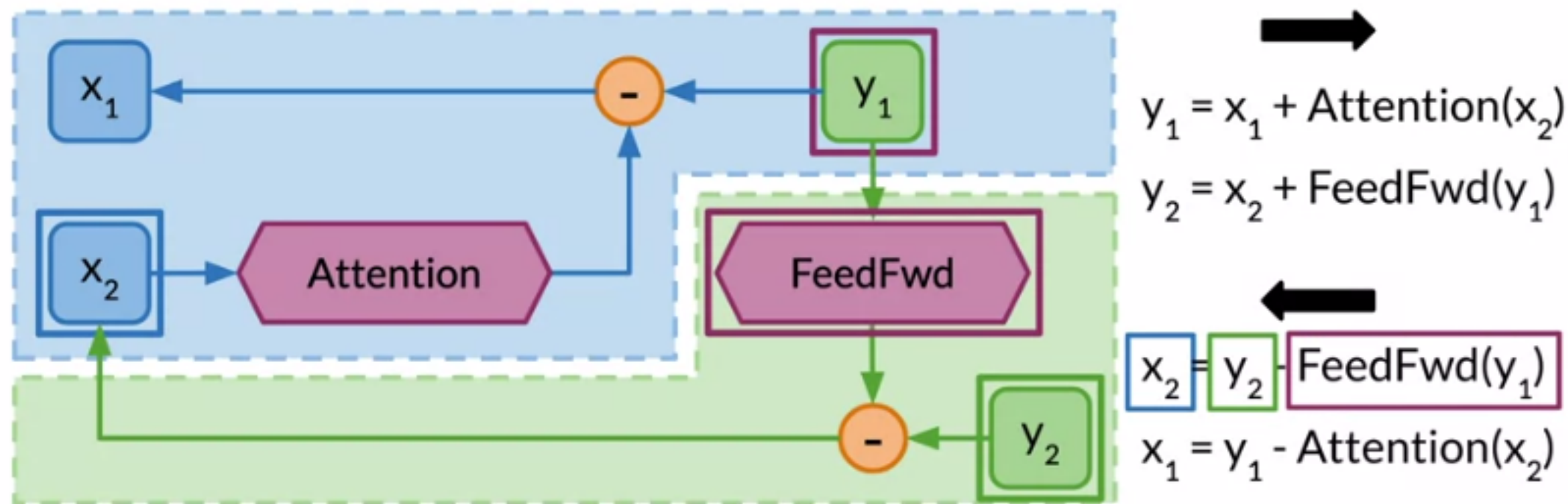
Reversible layers equations



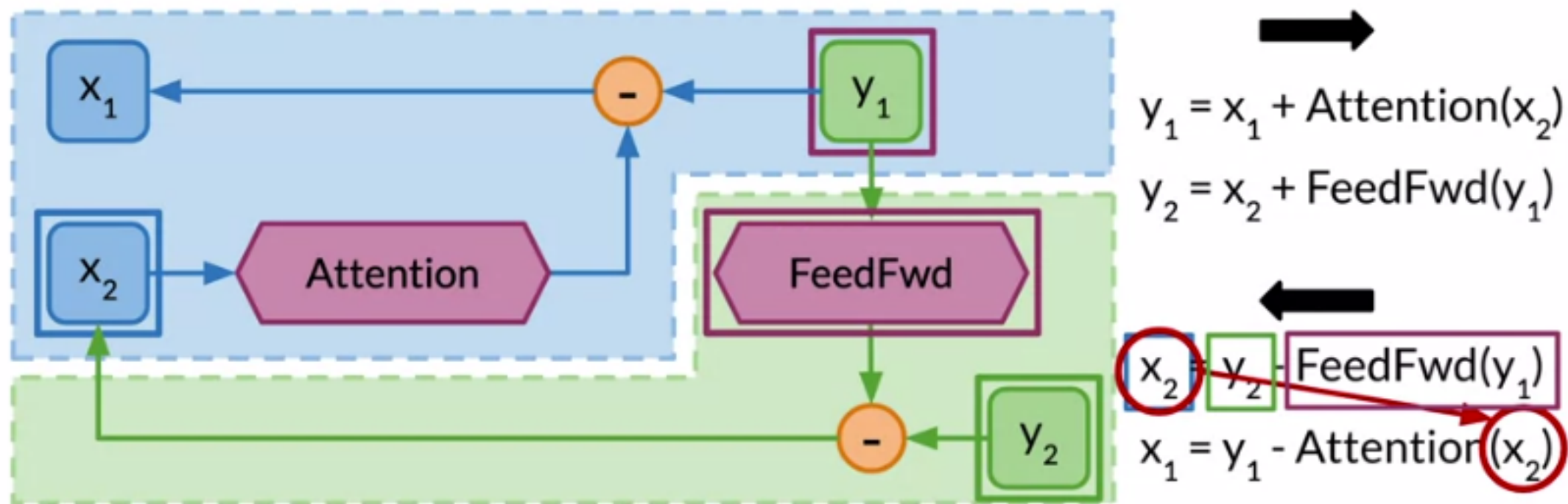
Reversible layers equations



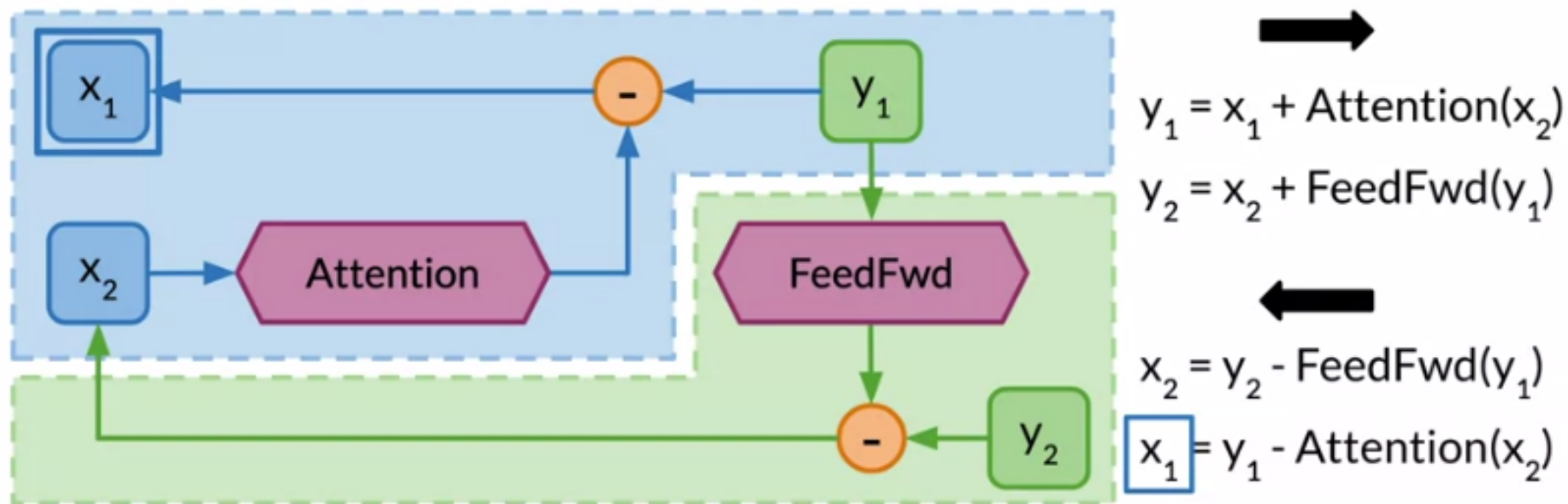
Reversible layers equations



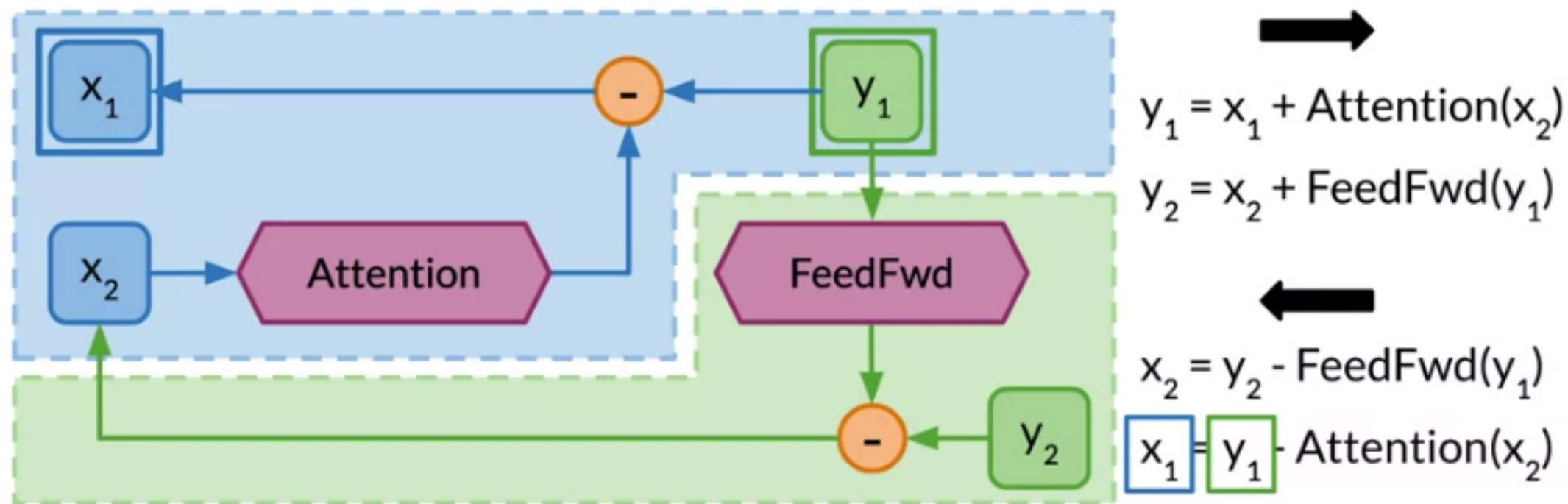
Reversible layers equations



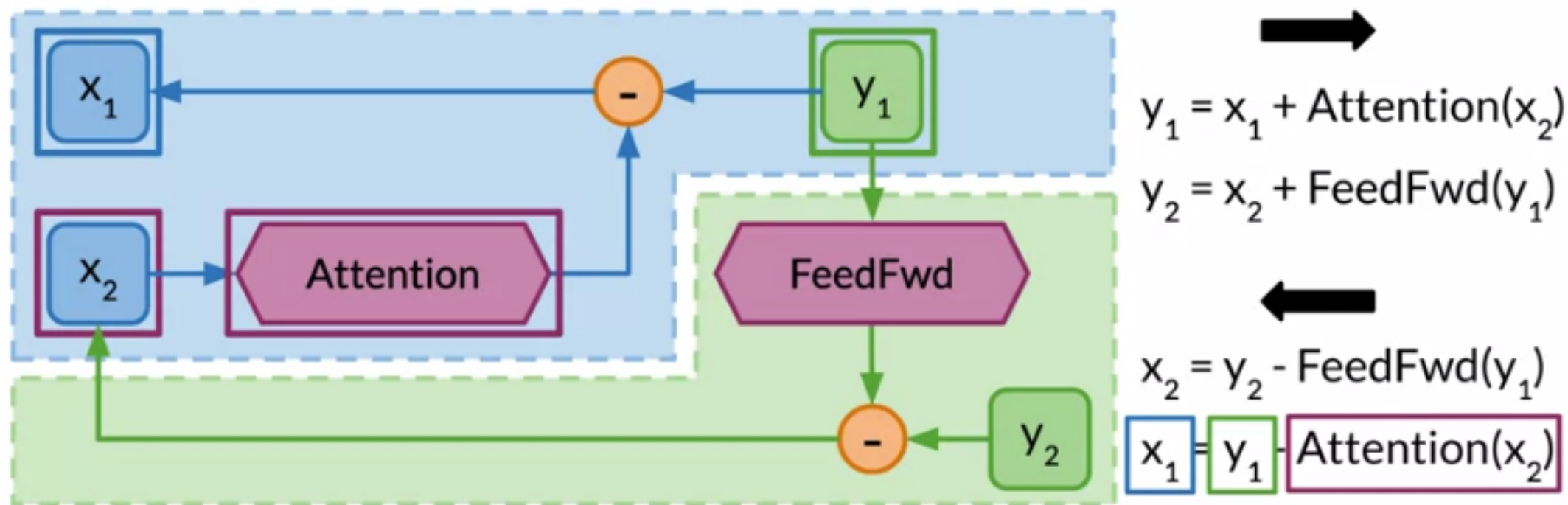
Reversible layers equations



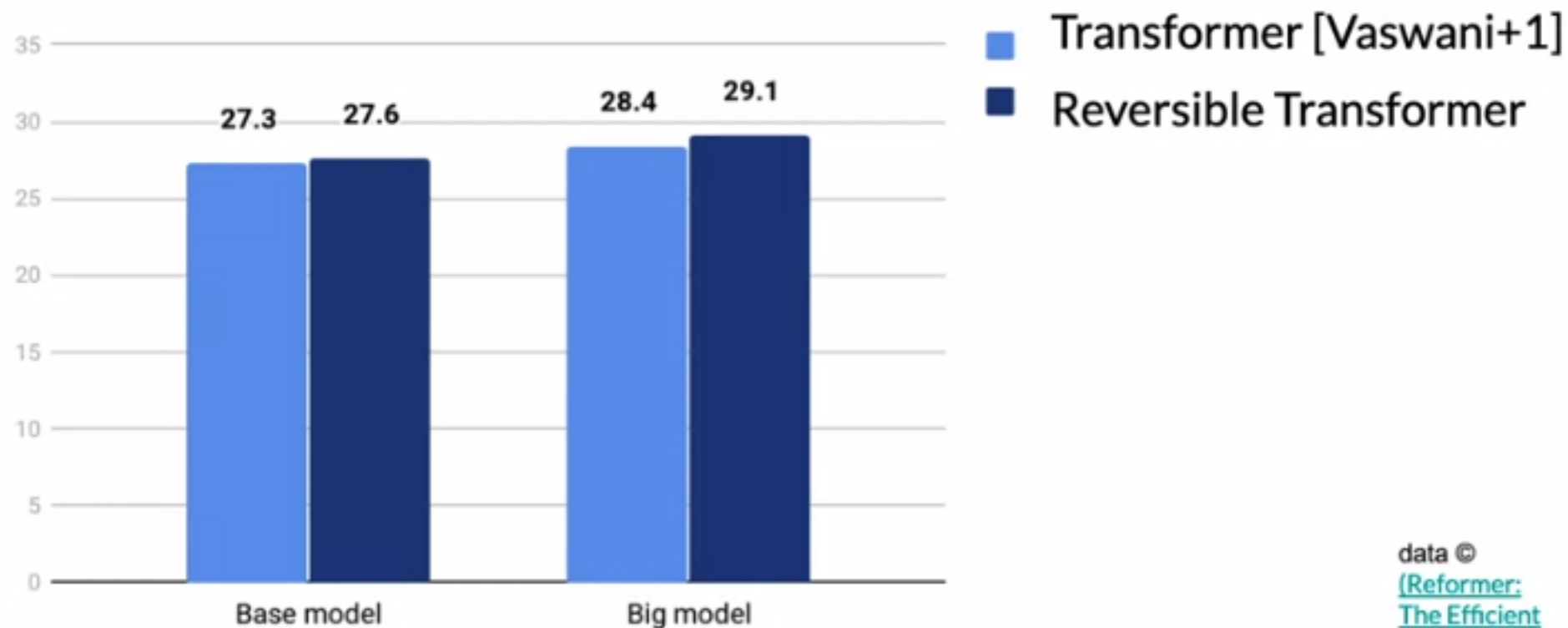
Reversible layers equations



Reversible layers equations



Reversible Transformer: BLEU Scores



data ©
[\(Reformer:
The Efficient
Transformer\)](#)

Reformer

The Reversible Transformer



$L = 1$ million tokens



1 GPU
(16 GB)

Reformer

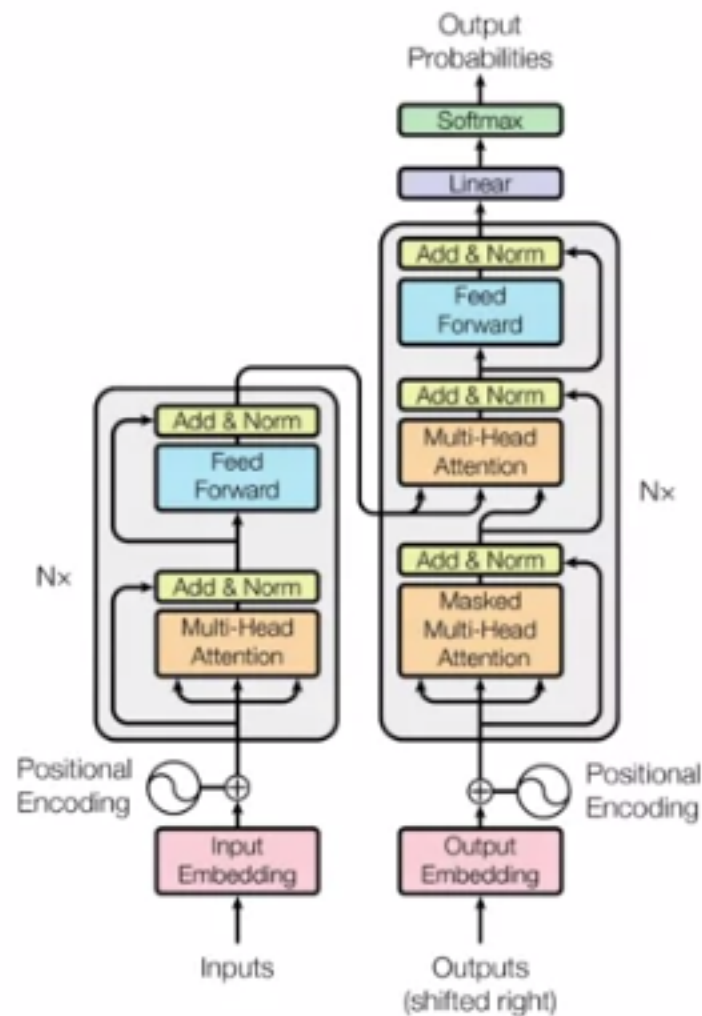


image ©
[\(Attention Is All You Need\)](#)

Reformer

- LSH Attention
- Reversible Layers

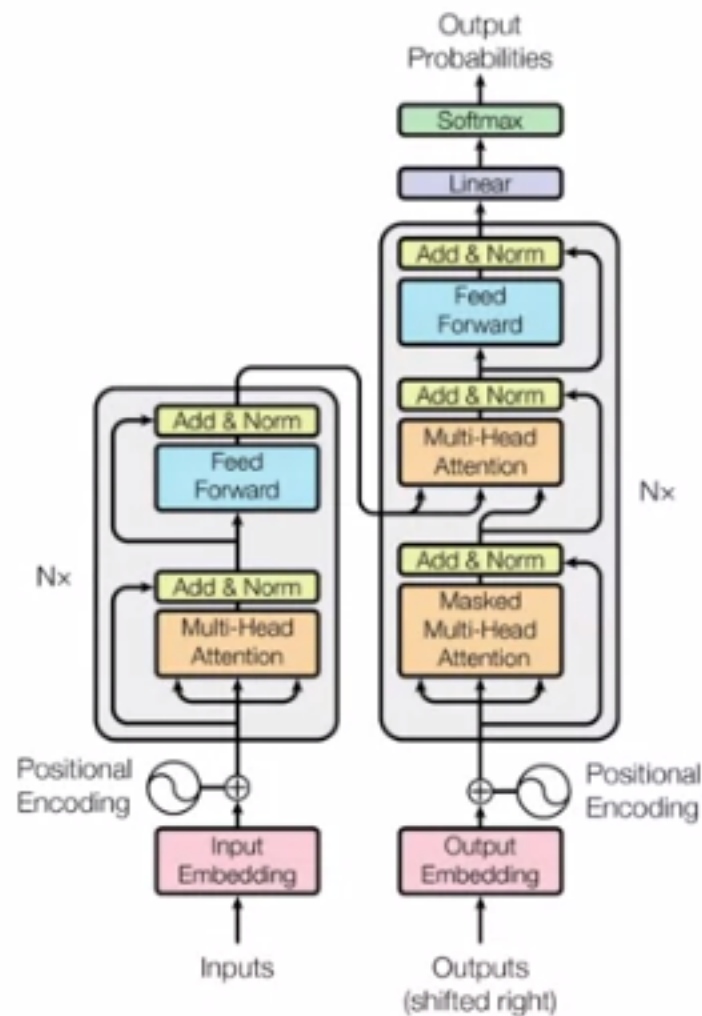


image ©
[\(Attention Is All You Need\)](#)

Reformer

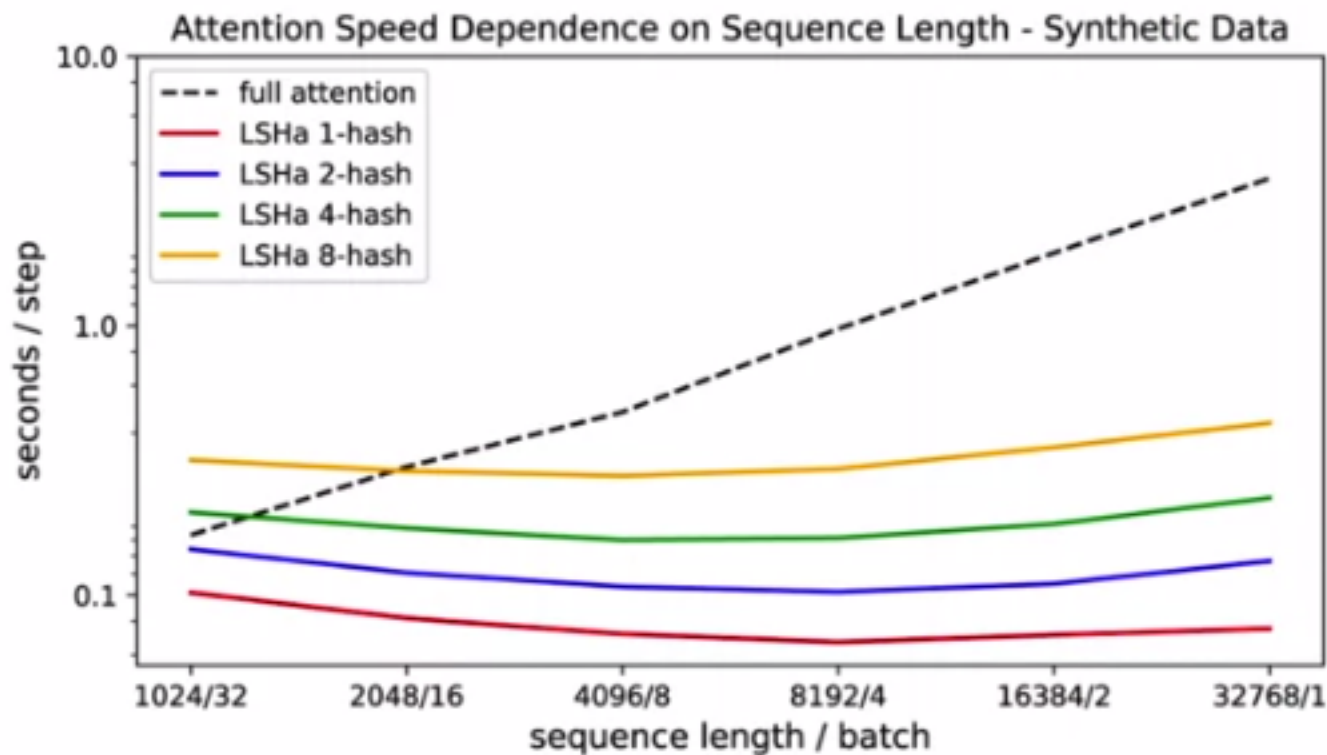
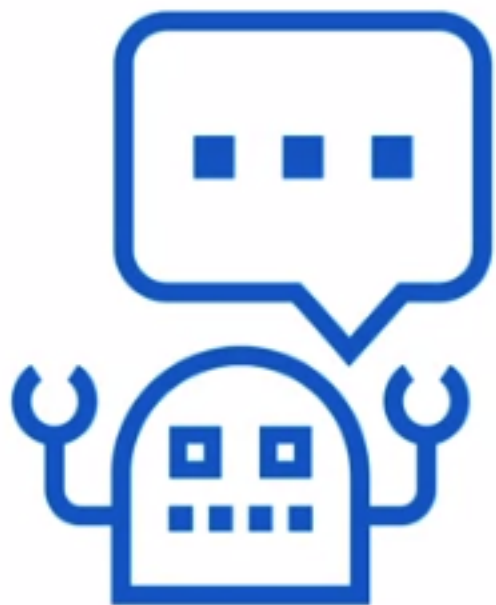
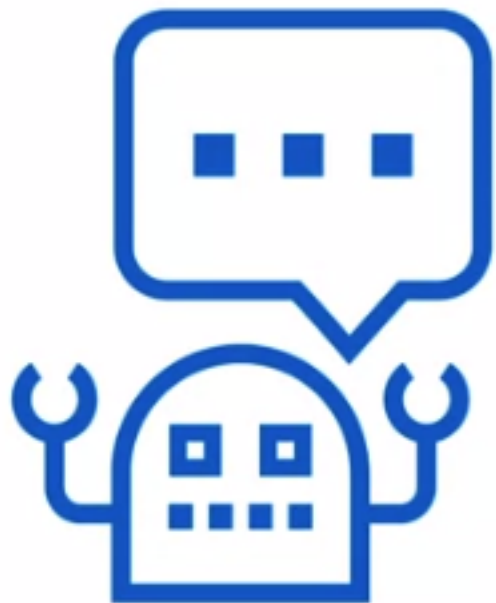


image ©
[\(Reformer:
The Efficient
Transformer\)](#)

Chatbot



Chatbot



- Reformer
- MulitiWOZ dataset
- Trax