# Deep Learning and Optimization
## Unpacking Transformers, LLMs and Diffusion

# Session 4

olivier.koch@ensae.fr

slack #ensae-dl-2025

Key ingredients to practical deep learning (activation, regularization, normalization, residual networks, etc.)

Recurrent networks struggle to learn long-term dependencies (vanishing gradients) and are slow/inefficient to train (sequential, not parallel, processing)

We learned tensor-based DL and applied it to MNIST.

| | Session | Date | Content |
|---|---|---|---|
| Foundations | 1 | Jan, 28 | Intro to DL<br>TP: micrograd |
| | 2 | Feb, 4 | Fundamentals I: inductive bias, loss functions<br>TP: bigram, MLP for next character prediction |
| | 3 | Feb, 11 | Fundamentals II: DL architectures<br>TP: tensor-based models |
| Applications | 4 | Feb, 18 | Attention & Transformers<br>TP: GPT from scratch |
| | 5 | Feb, 25 | DL for Computer vision<br>TP: convnets on CIFAR-10 |
| | 6 | Mar, 11 | VAE and Diffusion<br>TP: diffusion from scratch<br>Quiz / Exam |

The restaurant refused to serve me a ham sandwich because it only cooks vegetarian food. In the end, they just gave me two slices of bread. Their ambiance was just as good as the food and the service.
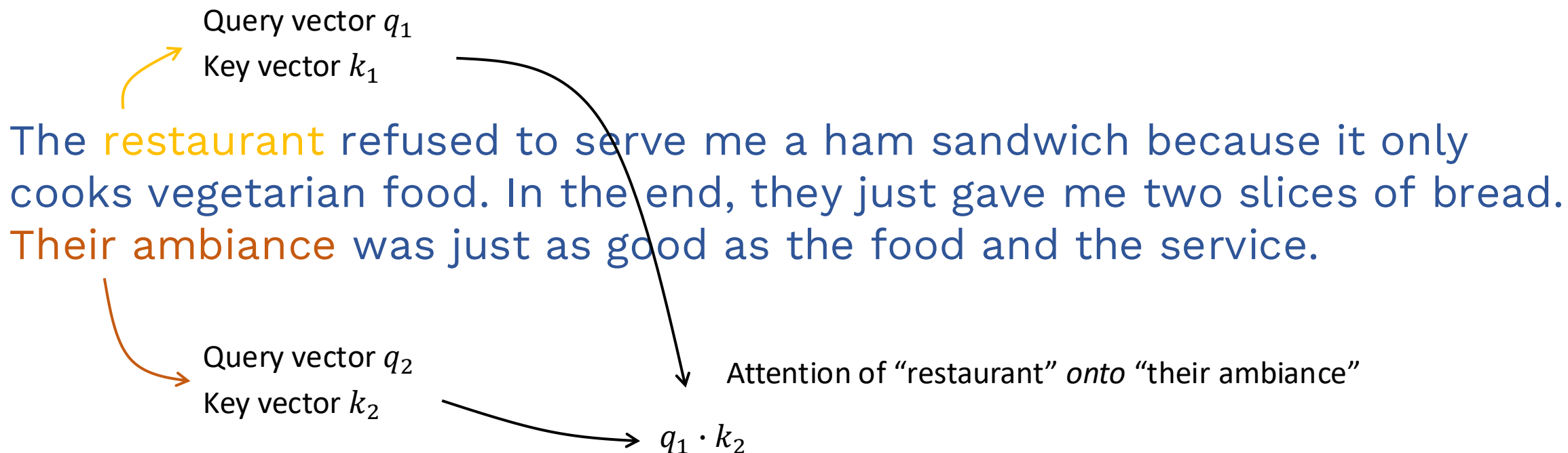
Query vector $q_1$

Key vector $k_1$

The restaurant refused to serve me a ham sandwich because it only cooks vegetarian food. In the end, they just gave me two slices of bread. Their ambiance was just as good as the food and the service.

Query vector $q_2$

Key vector $k_2$

Query vector $q_1$
Key vector $k_1$

The restaurant refused to serve me a ham sandwich because it only cooks vegetarian food. In the end, they just gave me two slices of bread. Their ambiance was just as good as the food and the service.
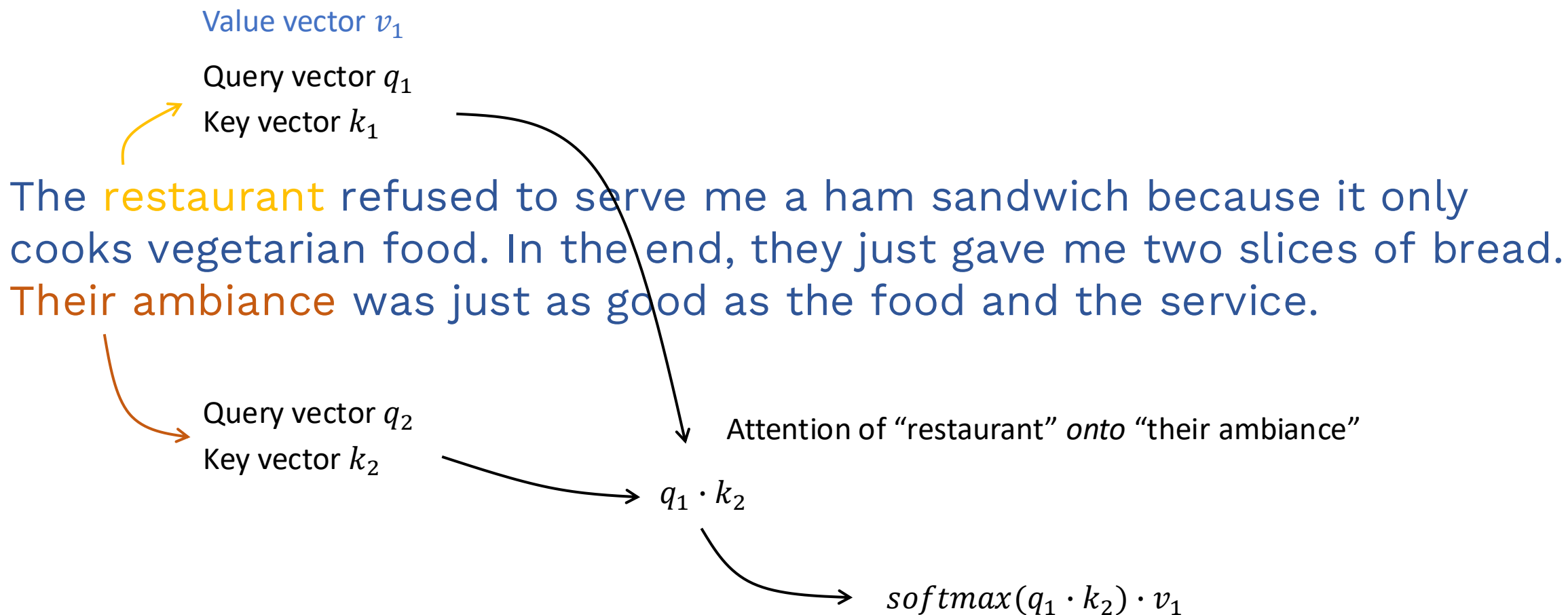
Query vector $q_2$
Key vector $k_2$

Attention of "restaurant" *onto* "their ambiance"

$q_1 \cdot k_2$

Value vector $v_1$

Query vector $q_1$

Key vector $k_1$

The restaurant refused to serve me a ham sandwich because it only cooks vegetarian food. In the end, they just gave me two slices of bread. Their ambiance was just as good as the food and the service.

Query vector $q_2$

Key vector $k_2$
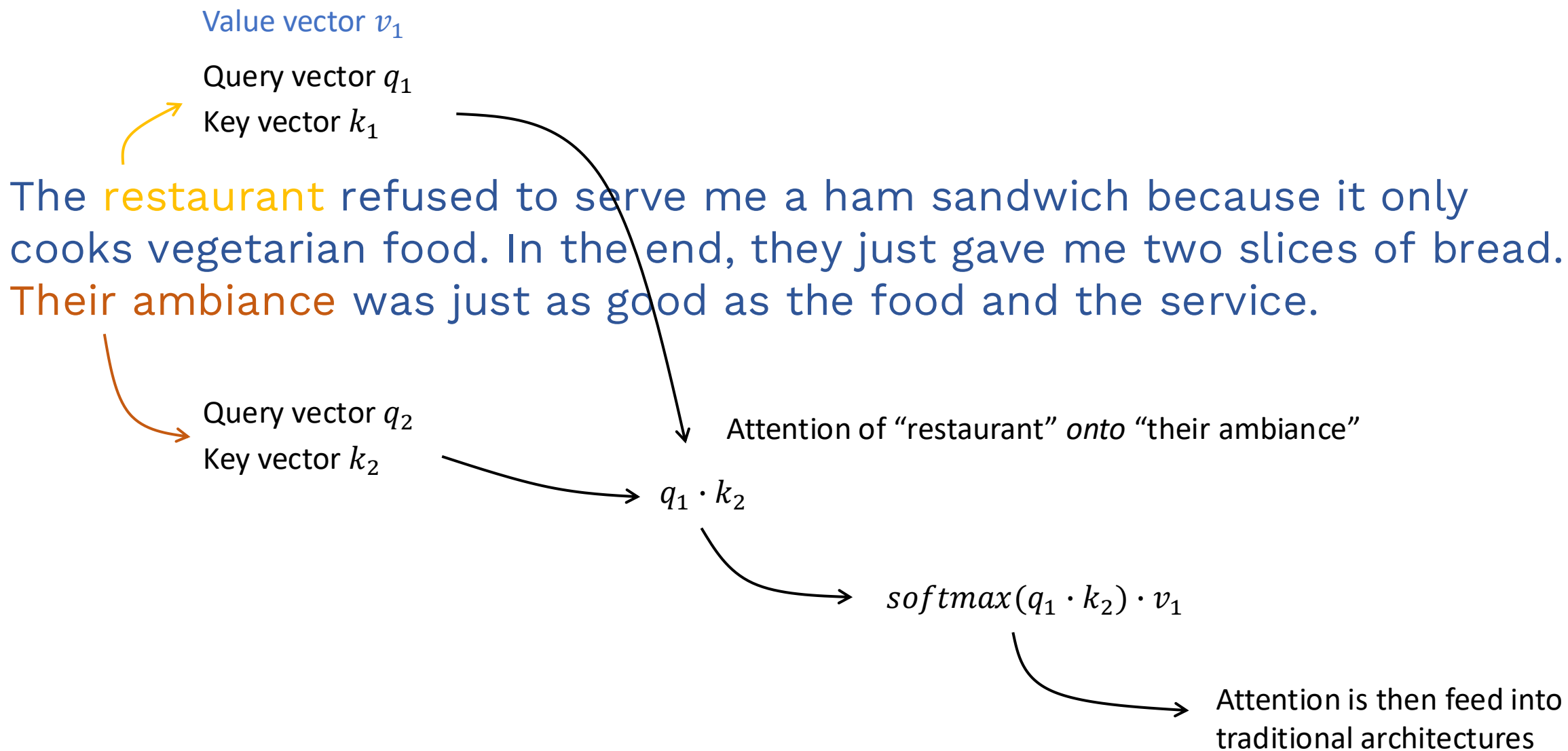
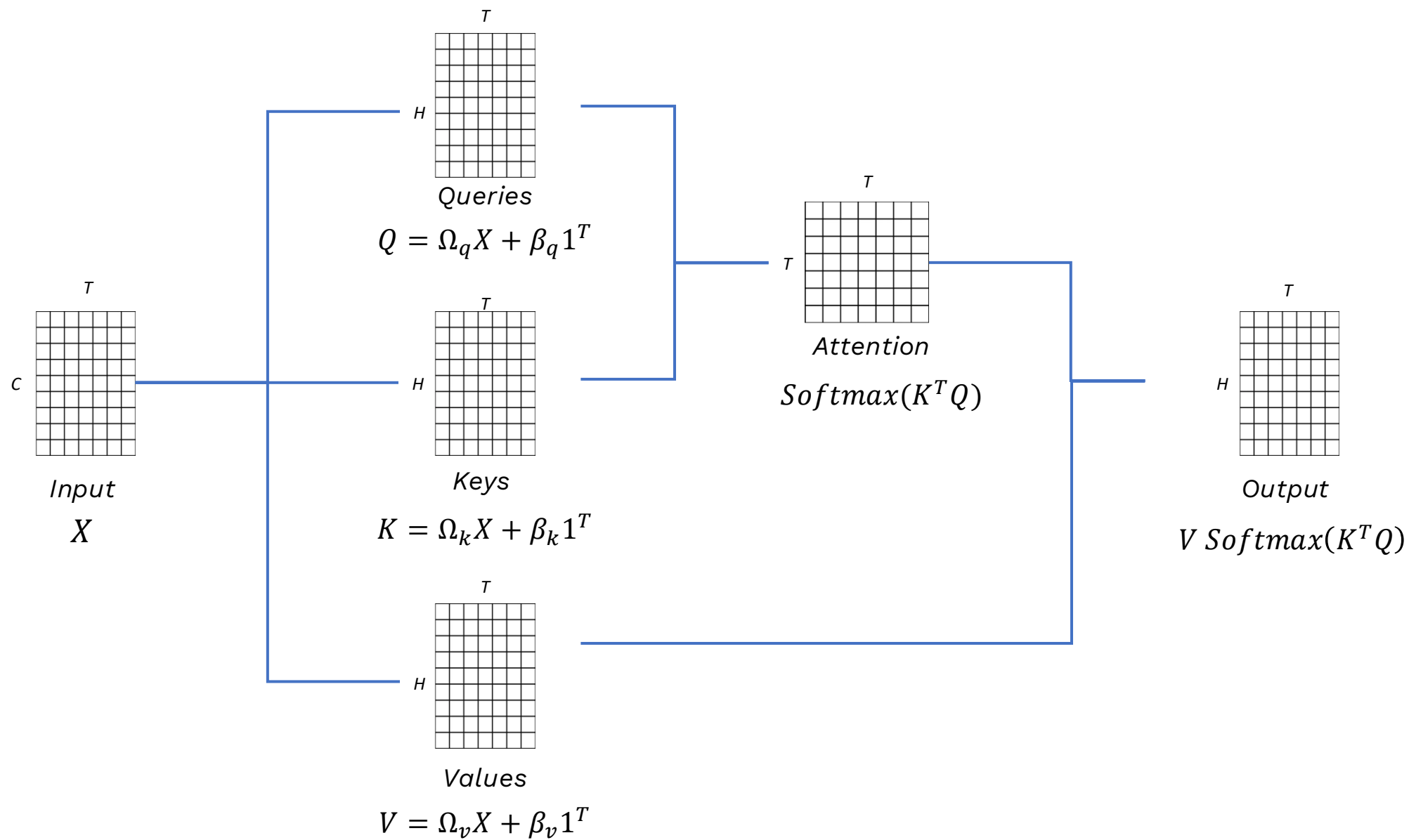Attention of "restaurant" *onto* "their ambiance"

$q_1 \cdot k_2$

$softmax(q_1 \cdot k_2) \cdot v_1$

Value vector $v_1$

Query vector $q_1$

Key vector $k_1$

The restaurant refused to serve me a ham sandwich because it only cooks vegetarian food. In the end, they just gave me two slices of bread. Their ambiance was just as good as the food and the service.

Query vector $q_2$

Key vector $k_2$

Attention of "restaurant" *onto* "their ambiance"

$q_1 \cdot k_2$

$softmax(q_1 \cdot k_2) \cdot v_1$

Attention is then feed into traditional architectures

# Attention and Transformers



$T$

$H$

Queries

$$Q = \Omega_q X + \beta_q 1^T$$

$T$

$C$

Input

$X$

$H$

Keys

$$K = \Omega_k X + \beta_k 1^T$$

$T$

$H$

Values

$$V = \Omega_v X + \beta_v 1^T$$

$T$

$T$

Attention

$$Softmax(K^T Q)$$

$T$

$H$

Output

$$V \, Softmax(K^T Q)$$

# Attention and Transformers

Number of input vectors
(**block size**)

Vector dimension
(e.g. **embedding size**)

Input
$X$

Queries
$$Q = \Omega_q X + \beta_q 1^T$$

Keys
$$K = \Omega_k X + \beta_k 1^T$$

Values
$$V = \Omega_v X + \beta_v 1^T$$

Attention
$$Softmax(K^T Q)$$

Output
$$V\, Softmax(K^T Q)$$

# Attention and Transformers



**Head size**

$T$

$H$

Queries

$$Q = \Omega_q X + \beta_q 1^T$$

Q is a linear operation on X

$\Omega_q$ is of size $H \times C$

$T$

$C$

Input

$X$

$T$

$H$

Keys

$$K = \Omega_k X + \beta_k 1^T$$

$T$

$T$

Attention

$$Softmax(K^T Q)$$

$T$

$H$

Output

$$V\ Softmax(K^T Q)$$

$T$

$H$

Values

$$V = \Omega_v X + \beta_v 1^T$$

# Attention and Transformers



**Head size**

$T$

$H$ Queries

$Q = \Omega_q X + \beta_q 1^T$

Q is a linear operation on X

$\Omega_q$ is of size $H \times C$

$T$

$C$ $H$ Keys

$K = \Omega_k X + \beta_k 1^T$

Input

$X$

$T$

$T$ Attention

$Softmax(K^T Q)$

$T$

$H$ Output

$V\, Softmax(K^T Q)$

$T$

$H$ Values

$V = \Omega_v X + \beta_v 1^T$

Attention weights (wei)

Self Attention

Source: Tom Yeh, 2024

# Attention and Transformers

Scaled Dot-Product Self-Attention

$$Sa[X] = V \cdot Softmax\left[\frac{K^T Q}{\sqrt{D_q}}\right]$$

Attention Is All You Need, A. Vaswani et al, NeurIPS 2017

Scaled Dot-Product Self-Attention

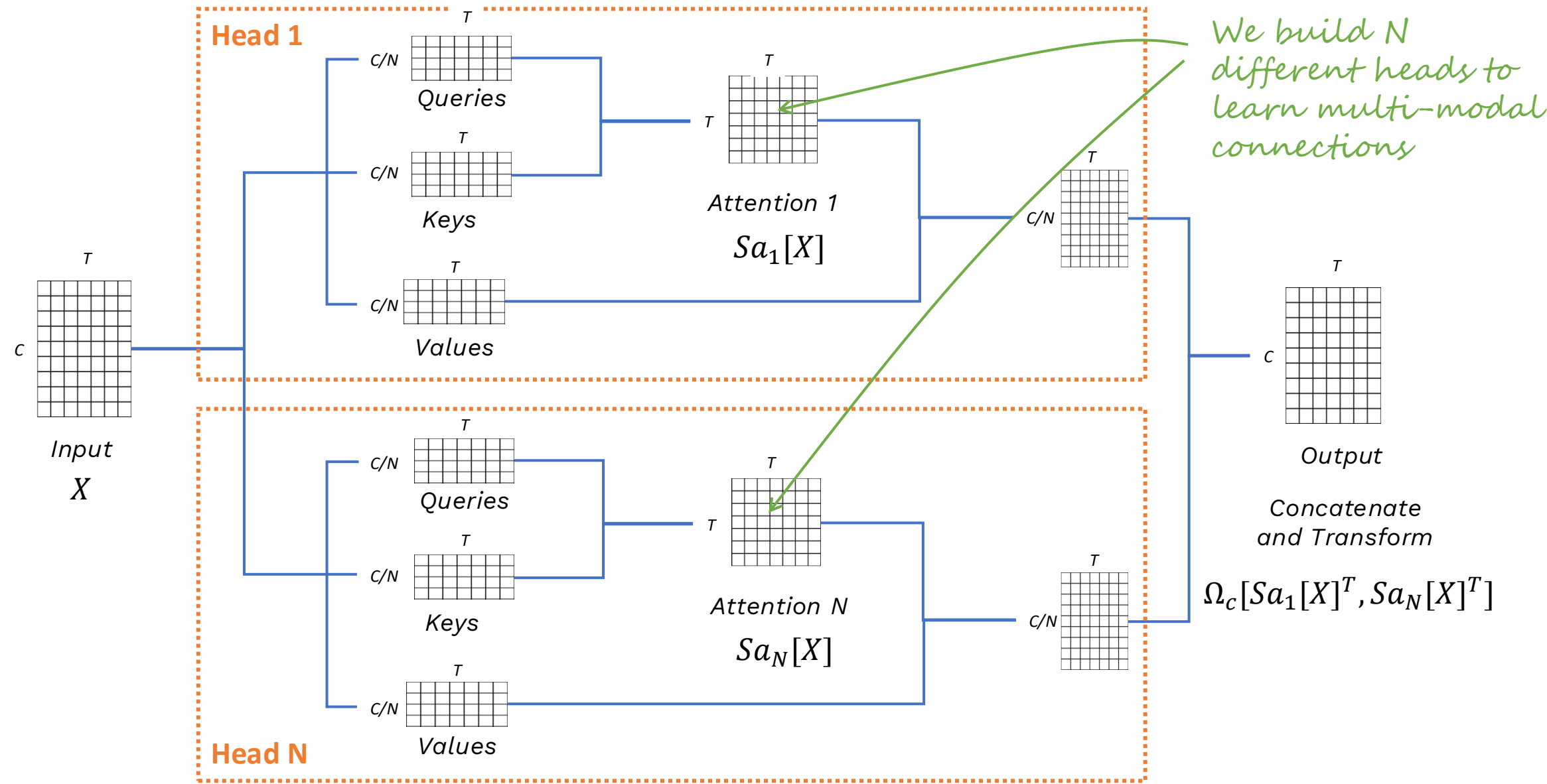$$Sa[X] = V \cdot Softmax\left[\frac{K^T Q}{\sqrt{D_q}}\right]$$

*Quadratic in X!*

$$Sa[X] = V \cdot Softmax\left[\frac{X^T \Omega_K^T \Omega_Q X}{\sqrt{D_q}}\right]$$

# Multi-head attention (for N heads)

# Multi-head attention (for N heads)



Head 1

$T$

$C/N$ Queries
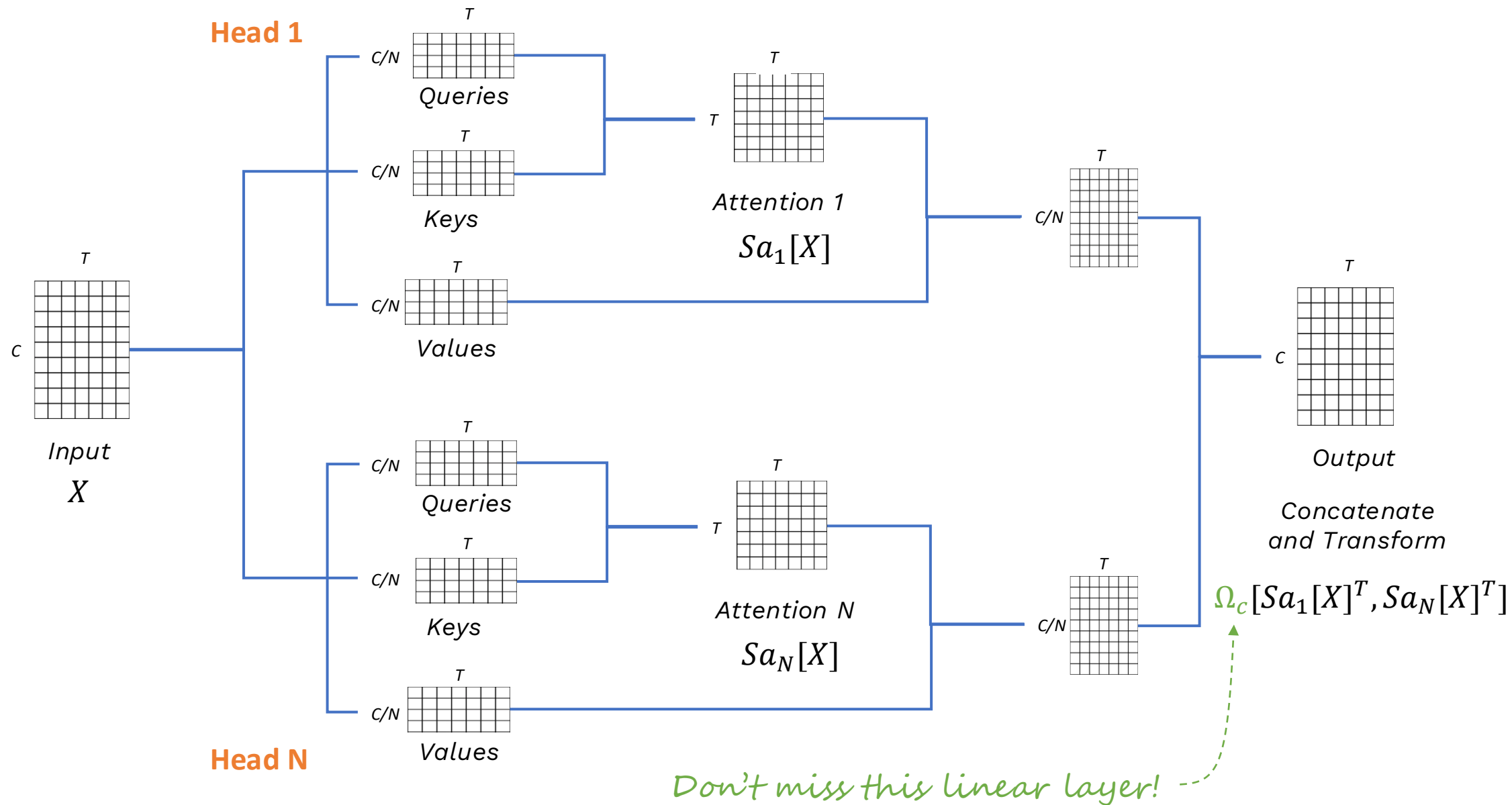
$C/N$ Keys

$C/N$ Values

$T$ Attention 1 $Sa_1[X]$

$C/N$

We build N different heads to learn multi-modal connections

Input $X$

$C$

$T$

Head N

$C/N$ Queries

$C/N$ Keys

$C/N$ Values

$T$ Attention N $Sa_N[X]$

$C/N$

Output

Concatenate and Transform

$\Omega_c[Sa_1[X]^T, Sa_N[X]^T]$

# Multi-head attention (for N heads)

"The cat sat on the mat"

Head 0: syntactic — "sat" attends to "cat" (subject-verb)

Head 1: positional — "sat" attends to nearby words

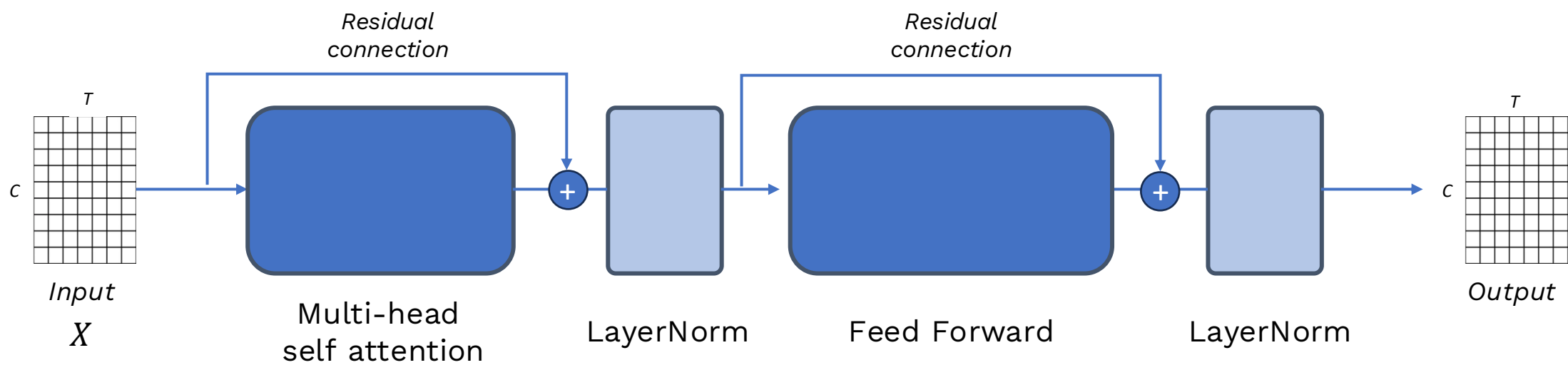Head 2: semantic — "sat" attends to "mat" (related concepts)
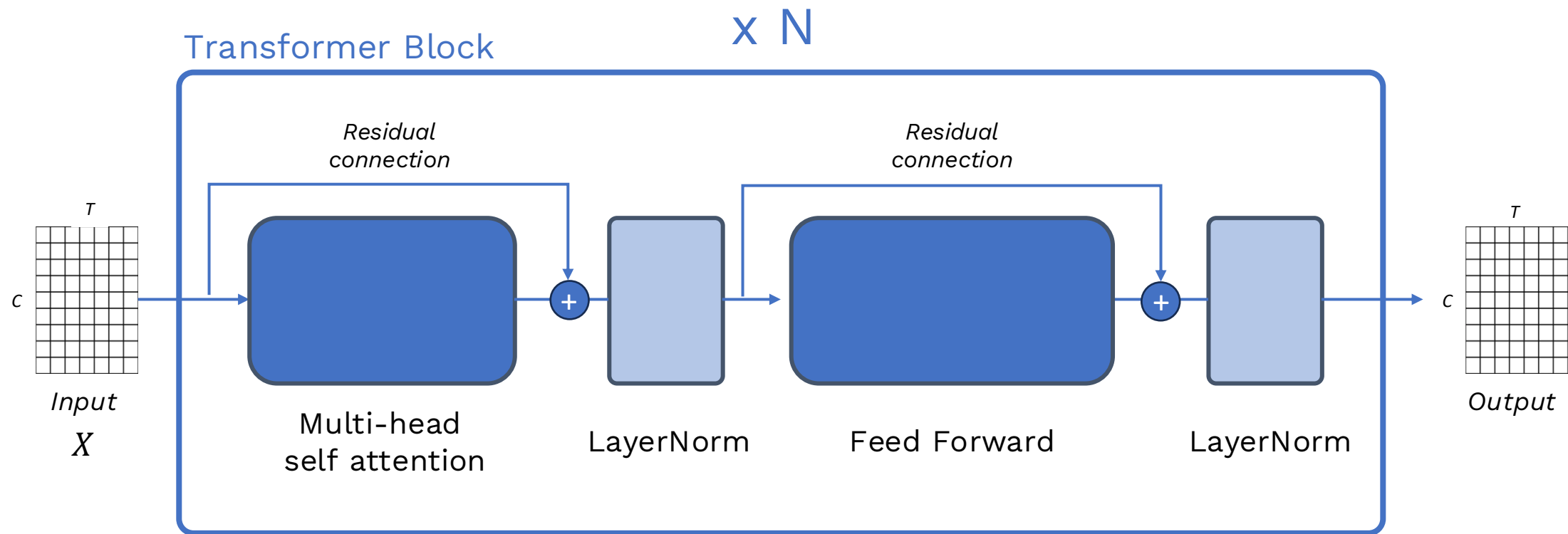
# Multi-head attention (for N heads)



**Head 1**

T

C/N — Queries (T)

C/N — Keys (T)

T — Attention 1 $Sa_1[X]$ (T)

C/N (T)

C — Input $X$ (T)

C/N — Values (T)

**Head N**

C/N — Queries (T)

C/N — Keys (T)

T — Attention N $Sa_N[X]$ (T)

C/N (T)

C/N — Values (T)

Output (T)

Concatenate and Transform

$\Omega_c[Sa_1[X]^T, Sa_N[X]^T]$

*Don't miss this linear layer!*

# Attention and Transformers



Attention Is All You Need, A. Vaswani et al, NeurIPS 2017

# Attention and Transformers



Residual connection

Residual connection

T

c

Input
$X$

Multi-head
self attention

+

LayerNorm

Feed Forward

+

LayerNorm

T

c

Output

Two linear layers with a 4X
expansion factor

Attention Is All You Need, A. Vaswani et al, NeurIPS 2017

# Attention and Transformers



x N

Transformer Block

Residual connection

Residual connection

$T$

$c$

Input

$X$

Multi-head self attention

LayerNorm

Feed Forward

LayerNorm

$T$

$c$

Output

Attention Is All You Need, A. Vaswani et al, NeurIPS 2017

**Figure 12.9** The input embedding matrix $\mathbf{X} \in \mathbb{R}^{D \times N}$ contains $N$ embeddings of length $D$ and is created by multiplying a matrix $\mathbf{\Omega}_e$ containing the embeddings for the entire vocabulary with a matrix containing one-hot vectors in its columns that correspond to the word or sub-word indices. The vocabulary matrix $\mathbf{\Omega}_e$ is considered a parameter of the model and is learned along with the other parameters. Note that the two embeddings for the word an in $\mathbf{X}$ are the same.
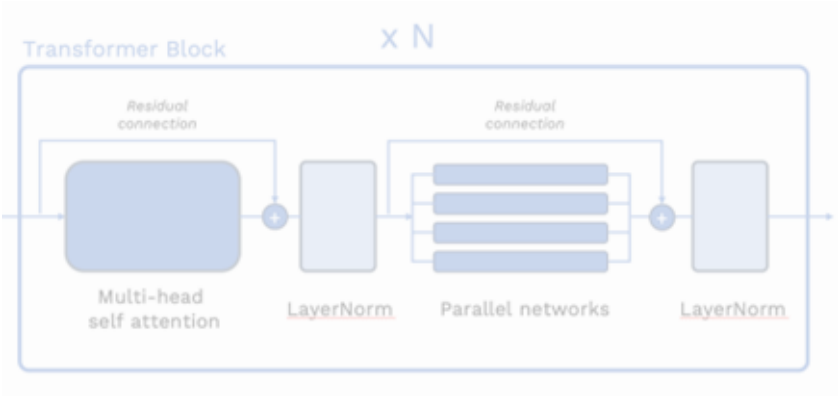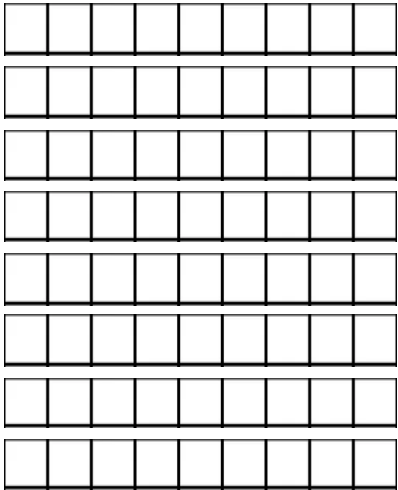
# Pretraining for BERT-like encoder

**Token embeddings**

She
<mask>
reading
a
<mask>
in
the
library

**Transformer Block**

x N

Residual connection

Residual connection

Multi-head self attention

+

LayerNorm

Parallel networks

+

LayerNorm

Linear + softmax

Probability distribution over vocabulary

...

Encoder: all tokens interact with each other

# Fine-tuning to specific tasks: review prediction

**Token embeddings**

<mask>
The
soup
had
a
terrible
taste
and

Transformer Block

x N

Residual connection

Residual connection

Multi-head self attention

LayerNorm

Parallel networks

LayerNorm

Linear + softmax

Probability of a positive review

...

# Fine-tuning to specific tasks: text classification

Token embeddings

<mask>

We
arrived
in
Venice
with
Jennie
on

Transformer Block                    x N

Residual
connection

Residual
connection

Multi-head
self attention          LayerNorm    Parallel networks    LayerNorm

Linear
+
softmax

Distribution over
classes

country     person     color     animal

# Pretraining for GPT-like decoder

Token embeddings

Linear + softmax

Probability distribution over vocabulary

She
was
reading
a
book
in
the
library

Transformer Block

x N

Residual connection

Residual connection

Multi-head self attention

LayerNorm

Parallel networks

LayerNorm

Task: predict future tokens

Decoder: tokens can only interact with past tokens (masked attention)
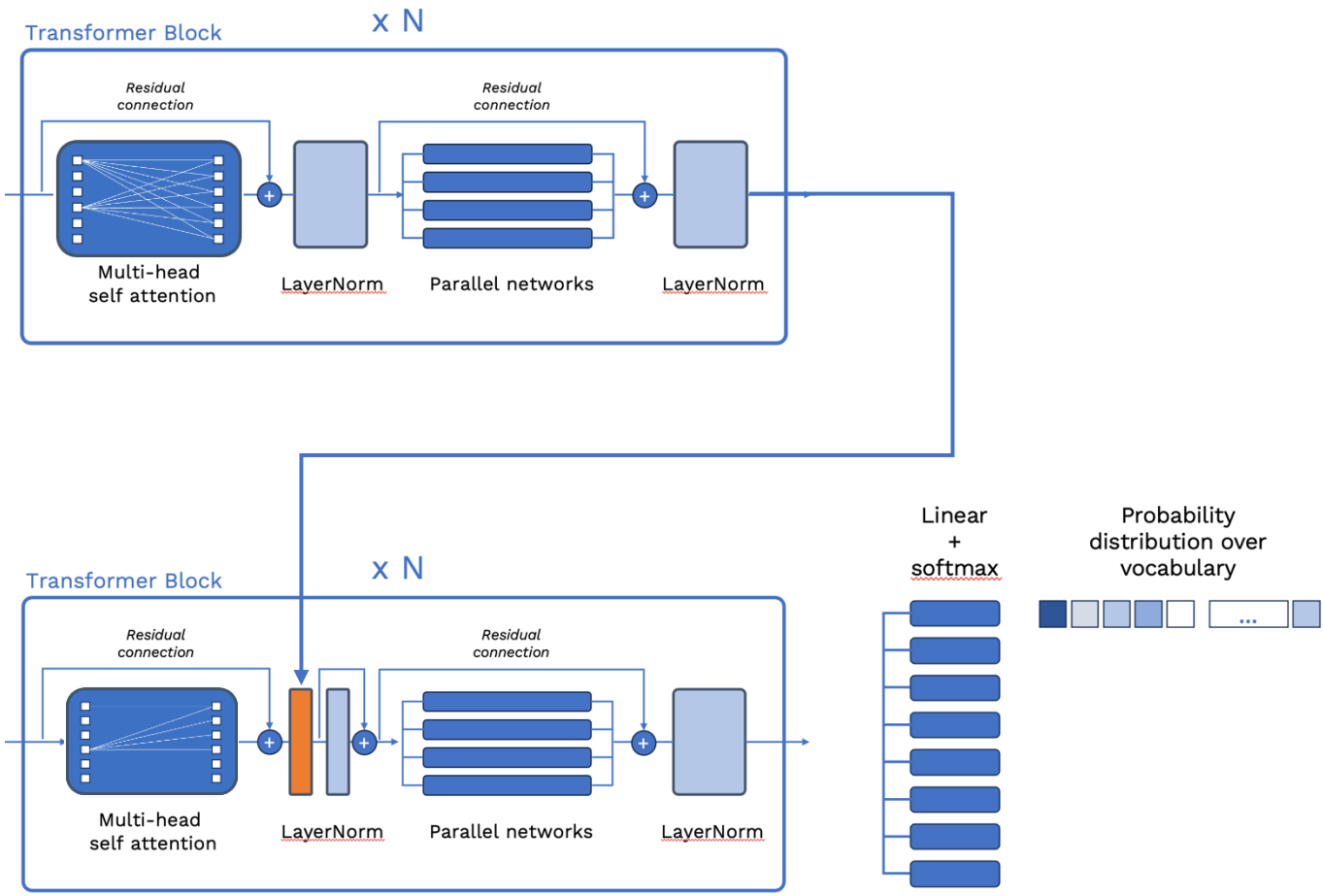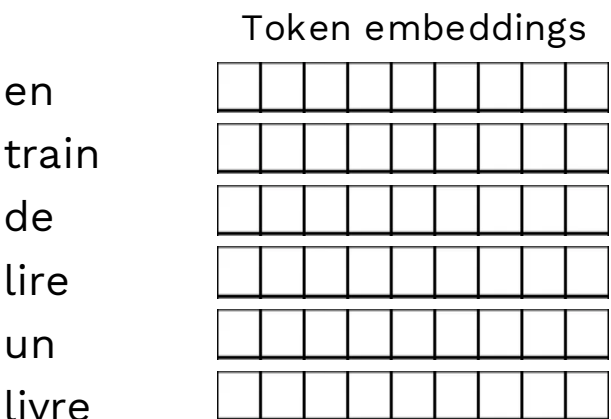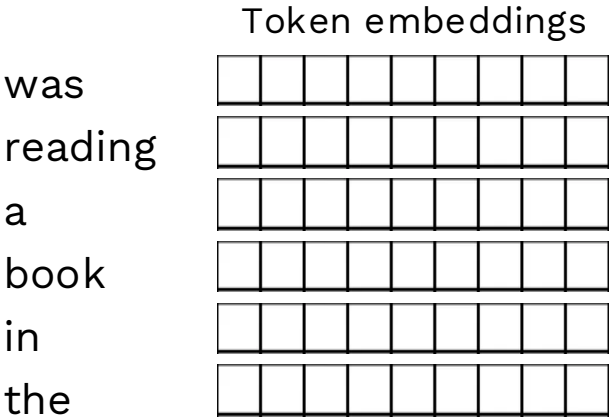
# Encoder-decoder architecture for translation with cross-attention

# Cross-attention



Decoder $X_{dec}$

Encoder $X_{inc}$

Queries

Keys

Values

Attention

$Softmax(K^T Q)$

Output

$V \, Softmax(K^T Q)$

# The original Transformer architecture



Encoder

Decoder

Figure 1: The Transformer - model architecture.

# Positional encoding



Figure 1: The Transformer - model architecture.

Encodes the relative position of tokens in the block

Vocabulary embedding table

*embedding size x vocab size*

Position embedding table

*embedding size x block size*

# Practical 4: Let's build a GPT-like encoder!

Token embeddings

She
was
reading
a
book
in
the
library

Task: predict future tokens

Transformer Block                                    x N

Residual connection                    Residual connection

Multi-head self attention    LayerNorm    Parallel networks    LayerNorm

Linear + softmax

Probability distribution over vocabulary

...

Decoder: tokens can only interact with past tokens (masked attention)

# Dataset generation

First, you know Caius Marcius is chief enemy to the people.

18        47    56      57        22       13

# Dataset generation

|       |    |    |    |    |    |    |
|-------|----|----|----|----|----|----|
|       | 18 | 47 | 56 | 57 | 22 | 13 |
|       | 22 | 31 | 82 | 16 | 46 | 81 |
| Batch | 46 | 36 | 32 | 82 | 10 | 11 |
|       | 74 | 59 | 82 | 91 | 60 | 38 |
|       | 27 | 21 | 37 | 26 | 42 | 18 |

# Dataset generation

| | | | X | | | | | | Y | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | 18 | 47 | 56 | 57 | 22 | 13 | 47 | 56 | 57 | 22 | 13 | 42 |
| | 22 | 31 | 82 | 16 | 46 | 81 | 31 | 82 | 16 | 46 | 81 | 32 |
| Batch | 46 | 36 | 32 | 82 | 10 | 11 | 36 | 32 | 82 | 10 | 11 | 69 |
| | 74 | 59 | 82 | 91 | 60 | 38 | 59 | 82 | 91 | 60 | 38 | 41 |
| | 27 | 21 | 37 | 26 | 42 | 18 | 21 | 37 | 26 | 42 | 18 | 77 |

# Dataset generation

Example 1

X                                                    Y

| 18 | 47 | 56 | 57 | 22 | 13 | | 47 | 56 | 57 | 22 | 13 | 42 |
|----|----|----|----|----|----|---|----|----|----|----|----|----|
| 22 | 31 | 82 | 16 | 46 | 81 | | 31 | 82 | 16 | 46 | 81 | 32 |
| Batch 46 | 36 | 32 | 82 | 10 | 11 | | 36 | 32 | 82 | 10 | 11 | 69 |
| 74 | 59 | 82 | 91 | 60 | 38 | | 59 | 82 | 91 | 60 | 38 | 41 |
| 27 | 21 | 37 | 26 | 42 | 18 | | 21 | 37 | 26 | 42 | 18 | 77 |

# Dataset generation

Example 2

X

Y

| | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| **18** | **47** | 56 | 57 | 22 | 13 | 47 | **56** | 57 | 22 | 13 | 42 |
| 22 | 31 | 82 | 16 | 46 | 81 | 31 | 82 | 16 | 46 | 81 | 32 |
| 46 | 36 | 32 | 82 | 10 | 11 | 36 | 32 | 82 | 10 | 11 | 69 |
| 74 | 59 | 82 | 91 | 60 | 38 | 59 | 82 | 91 | 60 | 38 | 41 |
| 27 | 21 | 37 | 26 | 42 | 18 | 21 | 37 | 26 | 42 | 18 | 77 |

Batch

# Dataset generation

Example 3

X
Y

| Batch | | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 18 | 47 | 56 | 57 | 22 | 13 | 47 | 56 | 57 | 22 | 13 | 42 |
| 22 | 31 | 82 | 16 | 46 | 81 | 31 | 82 | 16 | 46 | 81 | 32 |
| 46 | 36 | 32 | 82 | 10 | 11 | 36 | 32 | 82 | 10 | 11 | 69 |
| 74 | 59 | 82 | 91 | 60 | 38 | 59 | 82 | 91 | 60 | 38 | 41 |
| 27 | 21 | 37 | 26 | 42 | 18 | 21 | 37 | 26 | 42 | 18 | 77 |

# Dataset generation

X      Example 7      Y

| | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| 18 | 47 | 56 | 57 | 22 | 13 | 47 | 56 | 57 | 22 | 13 | 42 |
| **22** | 31 | 82 | 16 | 46 | 81 | **31** | 82 | 16 | 46 | 81 | 32 |
| 46 | 36 | 32 | 82 | 10 | 11 | 36 | 32 | 82 | 10 | 11 | 69 |
| 74 | 59 | 82 | 91 | 60 | 38 | 59 | 82 | 91 | 60 | 38 | 41 |
| 27 | 21 | 37 | 26 | 42 | 18 | 21 | 37 | 26 | 42 | 18 | 77 |

Batch

# Dataset generation

|  | X |  |  |  |  |  | Y |  |  |  |  |  |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
|  | 18 | 47 | 56 | 57 | 22 | 13 | 47 | 56 | 57 | 22 | 13 | 42 |
|  | 22 | 31 | 82 | 16 | 46 | 81 | 31 | 82 | 16 | 46 | 81 | 32 |
| Batch | 46 | 36 | 32 | 82 | 10 | 11 | 36 | 32 | 82 | 10 | 11 | 69 |
|  | 74 | 59 | 82 | 91 | 60 | 38 | 59 | 82 | 91 | 60 | 38 | 41 |
|  | 27 | 21 | 37 | 26 | 42 | 18 | 21 | 37 | 26 | 42 | 18 | 77 |

Example 30

# Dataset generation

|       | X   |     |     |     |     |     |
|-------|-----|-----|-----|-----|-----|-----|
|       | 18  | 47  | 56  | 57  | 22  | 13  |
|       | 22  | 31  | 82  | 16  | 46  | 81  |
| Batch | 46  | 36  | 32  | 82  | 10  | 11  |
|       | 74  | 59  | 82  | 91  | 60  | 38  |
|       | 27  | 21  | 37  | 26  | 42  | 18  |

# Dataset generation

X

| 18 | 47 | 56 | 57 | 22 | 13 |
|----|----|----|----|----|----|
| 22 | 31 | 82 | 16 | 46 | 81 |
| 46 | 36 | 32 | 82 | 10 | 11 |
| 74 | 59 | 82 | 91 | 60 | 38 |
| 27 | 21 | 37 | 26 | 42 | 18 |

Batch

(B)

Embedding ("channel")

(C)

Time (T)

# Dataset generation

A batch is a 3D tensor

X

| 18 | 47 | 56 | 57 | 22 | 13 |
|----|----|----|----|----|----|
| 22 | 31 | 82 | 16 | 46 | 81 |
| 46 | 36 | 32 | 82 | 10 | 11 |
| 74 | 59 | 82 | 91 | 60 | 38 |
| 27 | 21 | 37 | 26 | 42 | 18 |

Batch

(B)

Time (T)

B

C

T

Embedding ("channel")

(C)

# Tensor computation

```python
ones = torch.zeros(2, 2) + 1

twos = torch.ones(2, 2) * 2

threes = (torch.ones(2, 2) * 7 - 1) / 2

fours = twos ** 2

sqrt2s = twos ** 0.5
```

powers2 = twos ** torch.tensor([[1, 2], [3, 4]])

fives = ones + fours

dozens = threes * fours

tensor([[ 2.,  4.],
        [ 8., 16.]])

# Tensor computation

```
a = torch.rand((2,4,3))

a.transpose()
```

# Tensor computation

a = torch.rand((2,4,3))

a.transpose()

TypeError: transpose() received an invalid combination of arguments

# Tensor computation

```
a = torch.rand((2,4,3))

a.transpose(-2, -1)

a.shape
```

```
torch.Size([2, 3, 4])
```

# Tensor computation

```
a = torch.rand(2, 3)
b = torch.rand(3, 2)

print(a * b)
```

# Tensor computation

```
a = torch.rand(2, 3)
b = torch.rand(3, 2)

print(a * b)
```

RuntimeError: The size of tensor a (3) must match the size of tensor b (2) at non-singleton dimension 1

a = torch.rand(2, 3)
b = torch.rand(3, 2)

print(a @ b)

This works!

@ is for matrix multiplication

* is for element-wise multiplication

# Tensor broadcasting

```
a = torch.rand(2, 3)
b = torch.rand(1, 3)

print(a * b)

This works!
```

# Tensor broadcasting

```
a = torch.tensor([[1, 2, 1], [2, 5, 1]])
b = torch.ones(1, 3) + 1

print(a * b)

tensor([[ 2.,  4.,  2.],
        [ 4., 10.,  2.]])
```

# Tensor broadcasting

Brodcasting rules:

Comparing the dimension sizes of the two tensors, going from last to first:

Each dimension must be equal, or

One of the dimensions must be of size 1, or

The dimension does not exist in one of the tensors

```
a = torch.rand(5, 4, 3)
b = torch.rand(1, 3, 6)

print(a @ b)
```

# Tensor computation

```
a = torch.rand(5, 4, 3)
b = torch.rand(1, 3, 6)

print(a @ b)


This works!
```

```
a = torch.rand(1, 5, 4, 3)
b = torch.rand(3, 1, 3, 6)

print(a @ b)
```

# Tensor computation

```
a = torch.rand(1, 5, 4, 3)
b = torch.rand(3, 1, 3, 6)

print(a @ b)
```
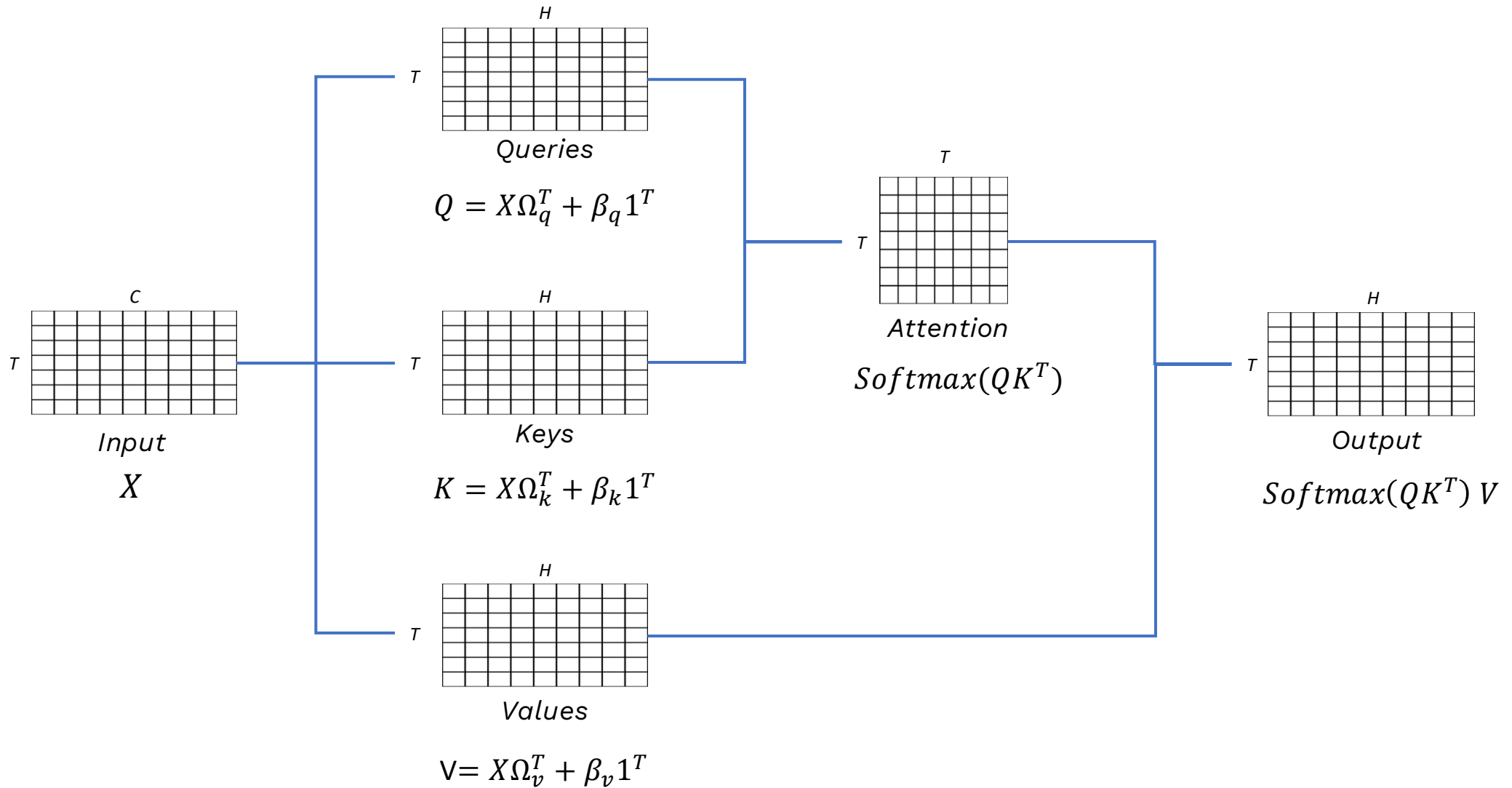
This works!

# Tensor computation
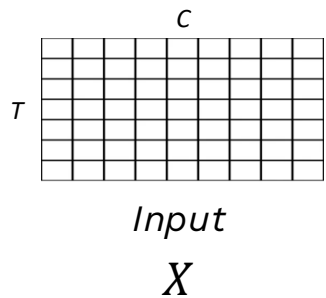
xbow = wei @ x  # (B, T, T) x (B, T, C) --> (B, T, C)

$H$

Queries

$$Q = X\Omega_q^T + \beta_q 1^T$$

$C$

Input

$X$

$H$

Keys

$$K = X\Omega_k^T + \beta_k 1^T$$

$T$

Attention

$$Softmax(QK^T)$$

$H$

Output

$$Softmax(QK^T)V$$

$H$

Values

$$V = X\Omega_v^T + \beta_v 1^T$$

$c$

$T$

*Input*

$X$

This is because pytorch is channel-last for memory optimization.

# In a pytorch implementation, the last two dimensions are inverted!

$c$

$T$

Input

$X$

A linear layer nn.Linear(in, out) implements:
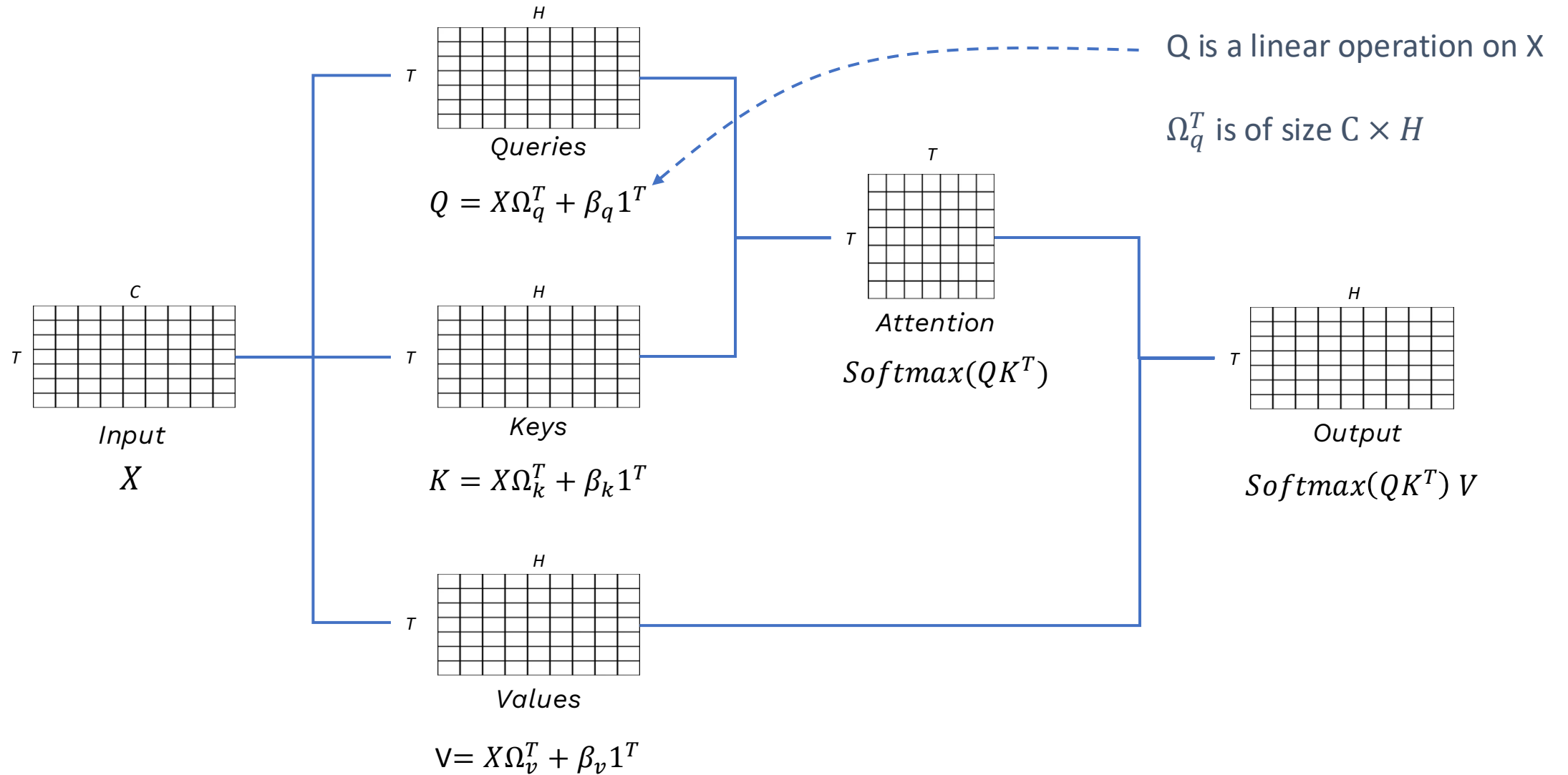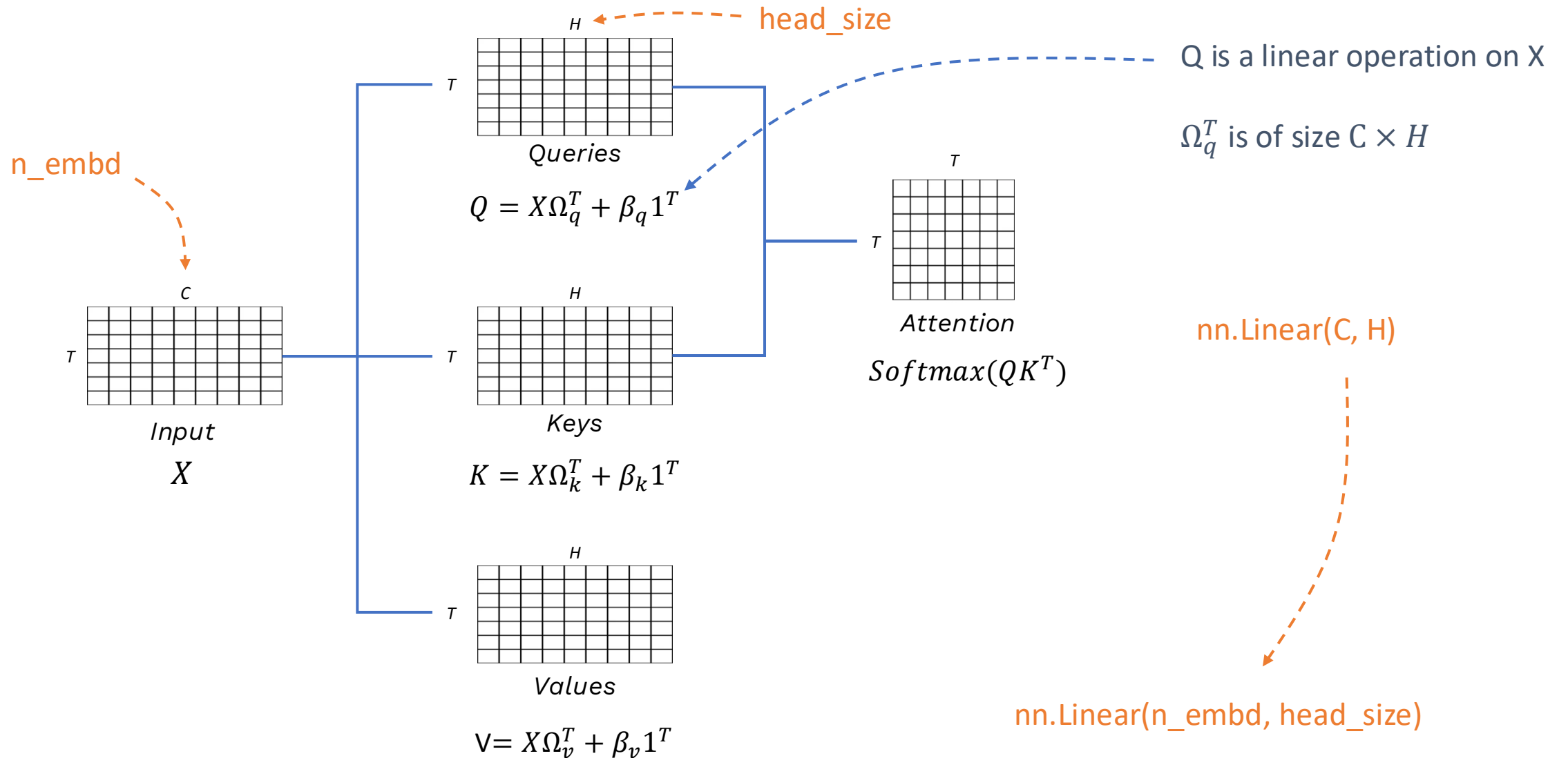
$$y = x.A^T + b$$

and not

$$y = A.x + b$$

Therefore the shape of A is (out, in)

# In a pytorch implementation, the last two dimensions are inverted!



Q is a linear operation on X

$\Omega_q^T$ is of size $C \times H$

Queries

$Q = X\Omega_q^T + \beta_q 1^T$

Attention

$Softmax(QK^T)$

Input

$X$

Keys

$K = X\Omega_k^T + \beta_k 1^T$

Output

$Softmax(QK^T)V$

Values

$V = X\Omega_v^T + \beta_v 1^T$

# In a pytorch implementation, the last two dimensions are inverted!



head_size

Q is a linear operation on X

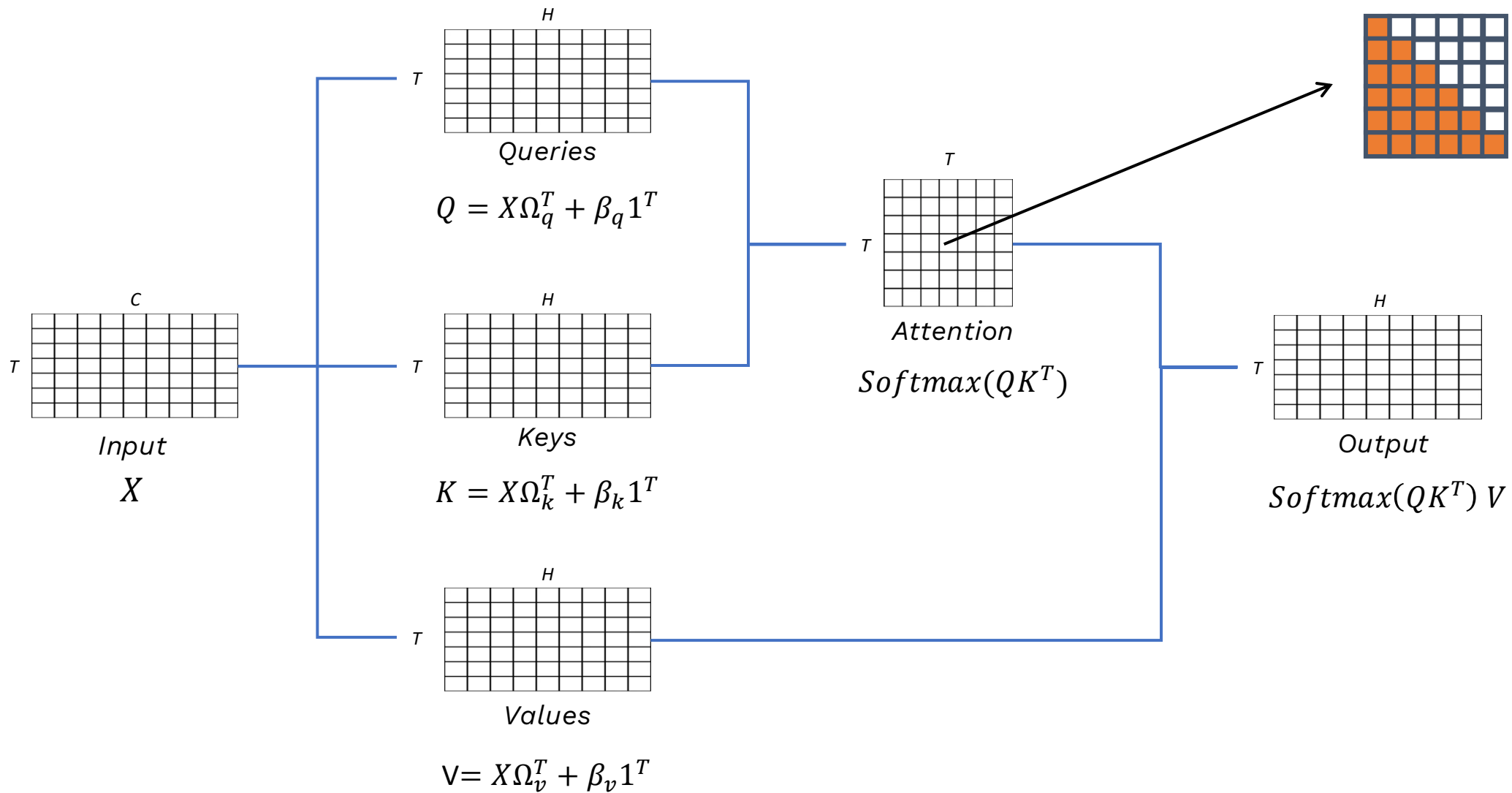$\Omega_q^T$ is of size $C \times H$

$H$

$T$

Queries

$$Q = X\Omega_q^T + \beta_q 1^T$$

n_embd

$C$

$T$

Input

$X$

$H$

$T$

Keys

$$K = X\Omega_k^T + \beta_k 1^T$$

$T$

$T$

Attention

$$Softmax(QK^T)$$

nn.Linear(C, H)

$H$

$T$

Values

$$V = X\Omega_v^T + \beta_v 1^T$$

nn.Linear(n_embd, head_size)

# Masked self-attention



Attention weights (wei)

**Queries**

$$Q = X\Omega_q^T + \beta_q 1^T$$

**Keys**

$$K = X\Omega_k^T + \beta_k 1^T$$

**Values**

$$V = X\Omega_v^T + \beta_v 1^T$$

**Input**

$$X$$

**Attention**

$$Softmax(QK^T)$$

**Output**

$$Softmax(QK^T)\,V$$

# Masked self-attention



$H$

$T$

Queries

$$Q = X\Omega_q^T + \beta_q 1^T$$

$C$

$T$

Input

$X$

$H$

$T$

Keys

$$K = X\Omega_k^T + \beta_k 1^T$$

$H$

$T$

Values

$$V = X\Omega_v^T + \beta_v 1^T$$

$T$

$T$

Attention

$$Softmax(QK^T)$$

$H$

$T$

Output

$$Softmax(QK^T)\,V$$

# Masked self-attention

```
tril = torch.tril(torch.ones(T,T))
```

```
tensor([[1., 0., 0.],
        [1., 1., 0.],
        [1., 1., 1.]])
```

# Masked self-attention

```
tril = torch.tril(torch.ones(T,T))

wei = torch.zeros((T,T))

wei = wei.masked_fill(tril == 0, float('-inf'))

tensor([[0., -inf, -inf],
        [0., 0., -inf],
        [0., 0., 0.]])
```

# Masked self-attention

```
tril = torch.tril(torch.ones(T,T))

wei = torch.zeros((T,T))

wei = wei.masked_fill(tril == 0, float('-inf’))


tensor([[0., -inf, -inf],
        [0., 0., -inf],
        [0., 0., 0.]])

wei = F.softmax(wei, dim=-1)
```

# Masked self-attention

```
tril = torch.tril(torch.ones(T,T))

wei = torch.zeros((T,T))

wei = wei.masked_fill(tril == 0, float('-inf'))
```
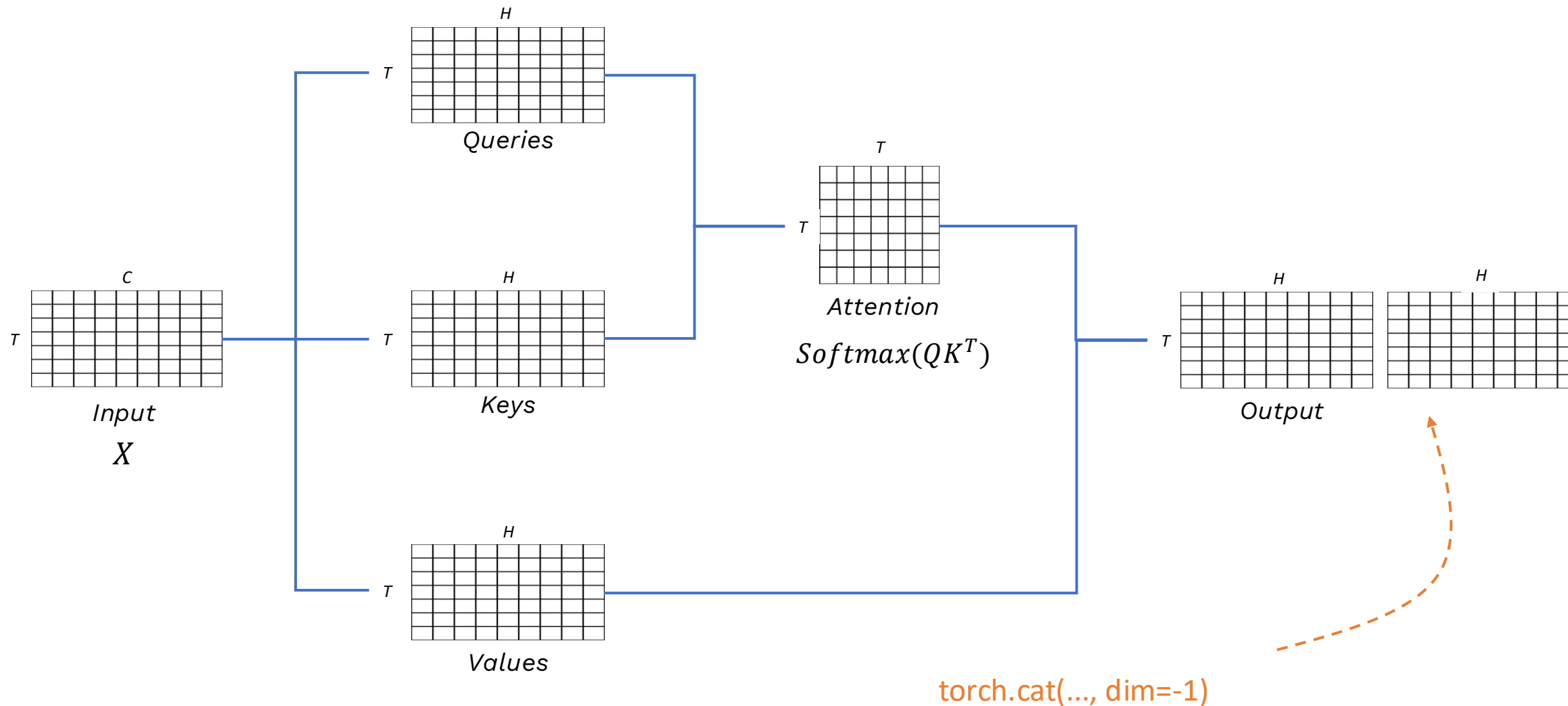
```
tensor([[0., -inf, -inf],
        [0., 0., -inf],
        [0., 0., 0.]])
```

```
wei = F.softmax(wei, dim=-1)
```

```
tensor([[1.0000, 0.0000, 0.0000],
        [0.5000, 0.5000, 0.0000],
        [0.3333, 0.3333, 0.3333]])
```

# Beware of concatenation in multi-head

# Getting started

```python
class Head(nn.Module):
    """ one head of self-attention """

    def __init__(self, head_size):
        super().__init__()
        self.key = nn.Linear(..., ..., bias=False)

        ...

    def forward (self, x):
        B, T, C = x.shape
        k = self.key(x) # (B,T,C)

        ...
```

# Getting started

```python
class Head(nn.Module):
    """ one head of self-attention """

    def __init__(self, head_size):
        super().__init__()
        self.key = nn.Linear(..., ..., bias=False)
        ...

    def forward (self, x):
        B, T, C = x.shape
        k = self.key(x) # (B,T,C)
        q = ...
        # compute self attention scores (affinities)
        wei = ...
        wei = wei.masked_fill(self.tril[:T, :T] == 0, float('-inf'))
        wei = F.softmax(wei, dim=-1)
        ...
```

# Practical 4 summary

1. Self-attention by hand

2. Self-attention in pytorch

3. GPT piece-by-piece

4. GPU goes rrr!

Dataset: Shakespeare's corpus (input.txt)