# Visualize a batch of data

C
=64

B
=16

T =32

1 3 8 39 25 40 64 54 ... 12 18 4 7 16 9 58 14
20 38 49 35 50 60 34 ... 10 8 5 4 18 9 5 13 2

45 37 62 2 26 33 56 5 ... 1 52 11 9 4 6 28 0

B: batch_size
T: block_size
C: n_embd

C
=64

1.2 15 458 0.1 52 111
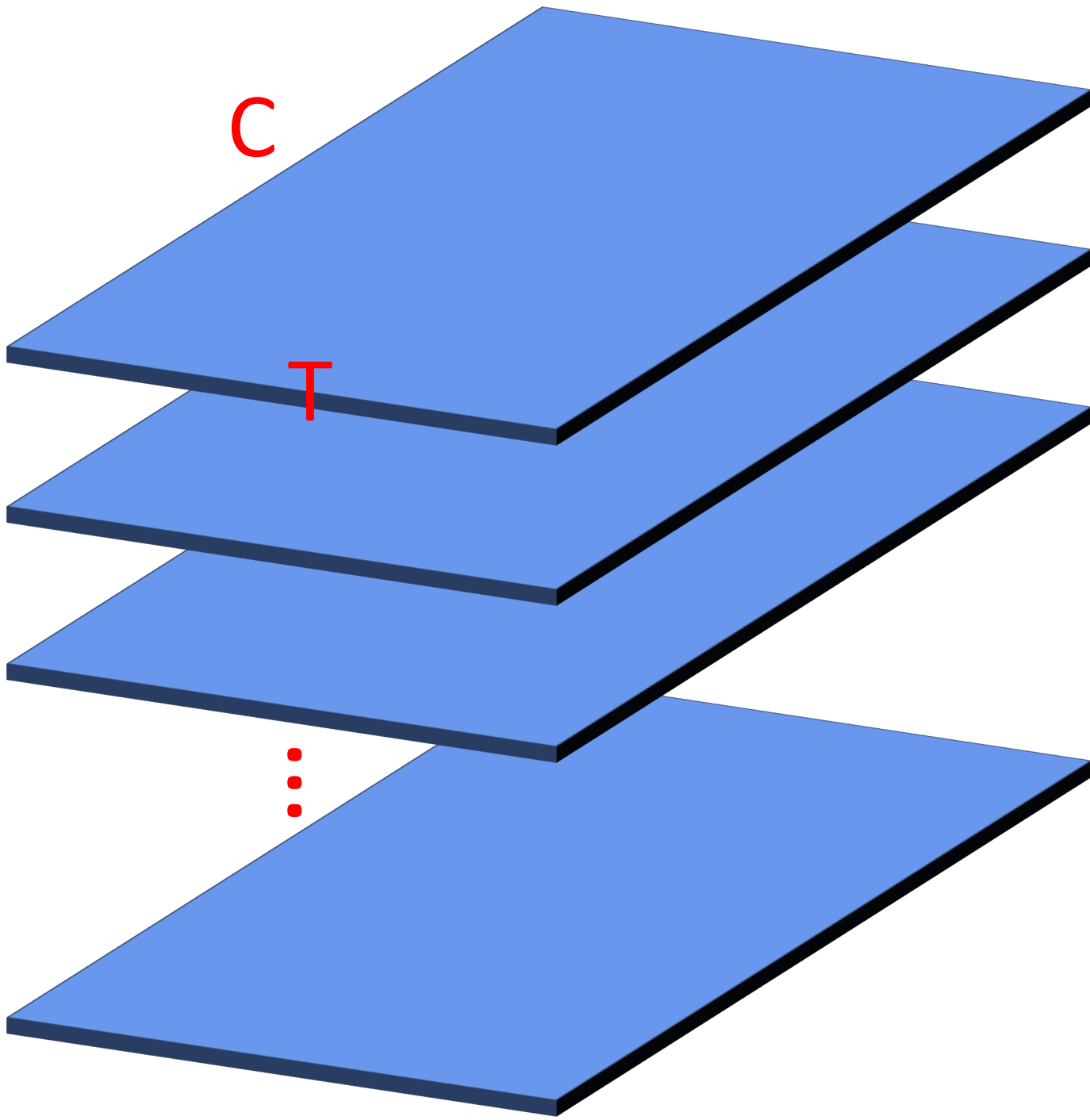14 14.7 56 36.5 5 14.5
9 2.5 3.4 10 2.2 45 67 54 ...
11.5 35 1.2 22 35 88 5.2 ...
2 35 6.7 8.8 90 121 5 12
5.4 3.3 13 27 45.7 56.8 ...

1 3 8 39 25 40 64 54 ... 12 18 4 7 16 9 58 14

B
=16

20 38 49 35 50 60 34 ... 10 8 5 4 18 9 5 13 2

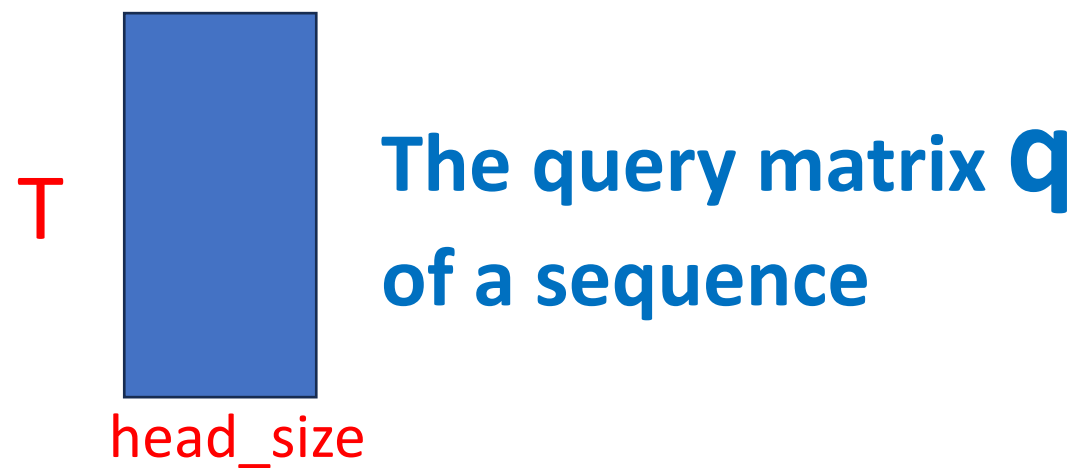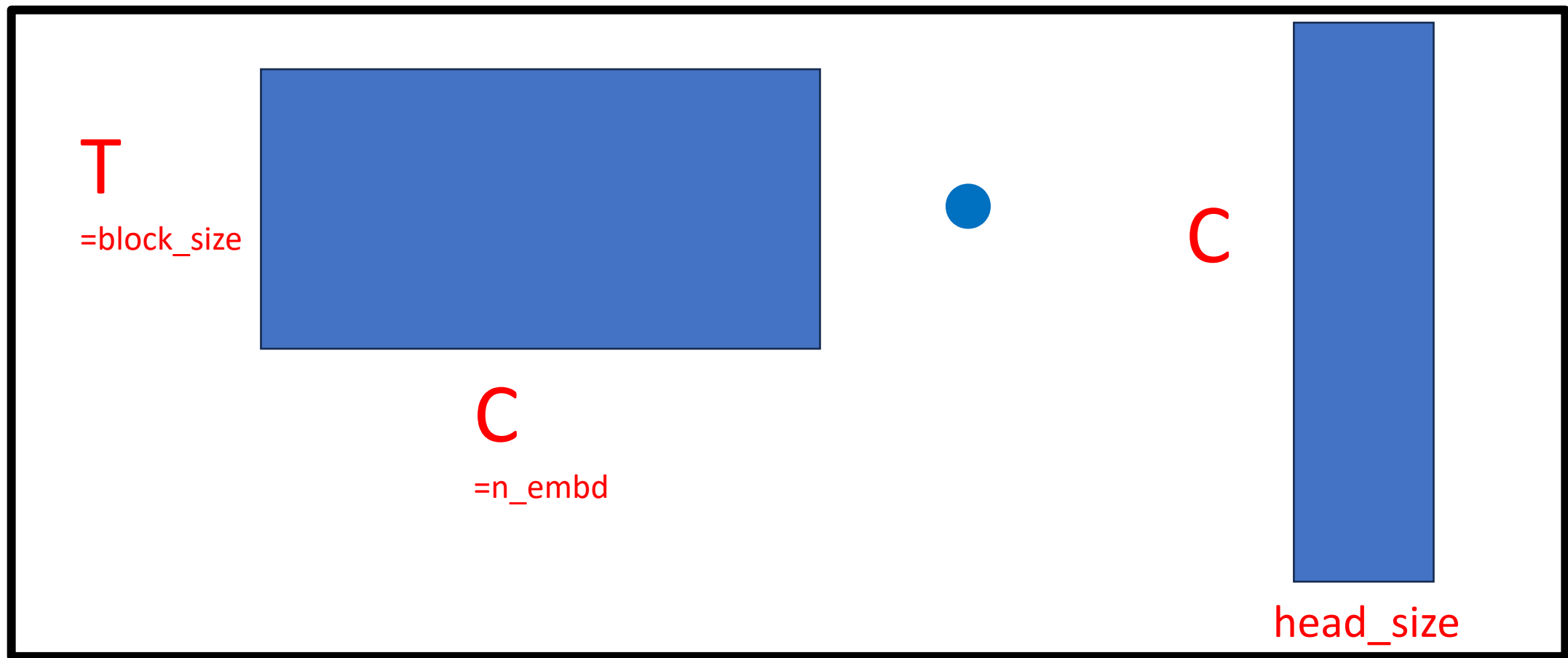45 37 62 2 26 33 56 5 ... 1 52 11 9 4 6 28 0

T =32

B: batch_size
T: block_size
C: n_embd

C

T
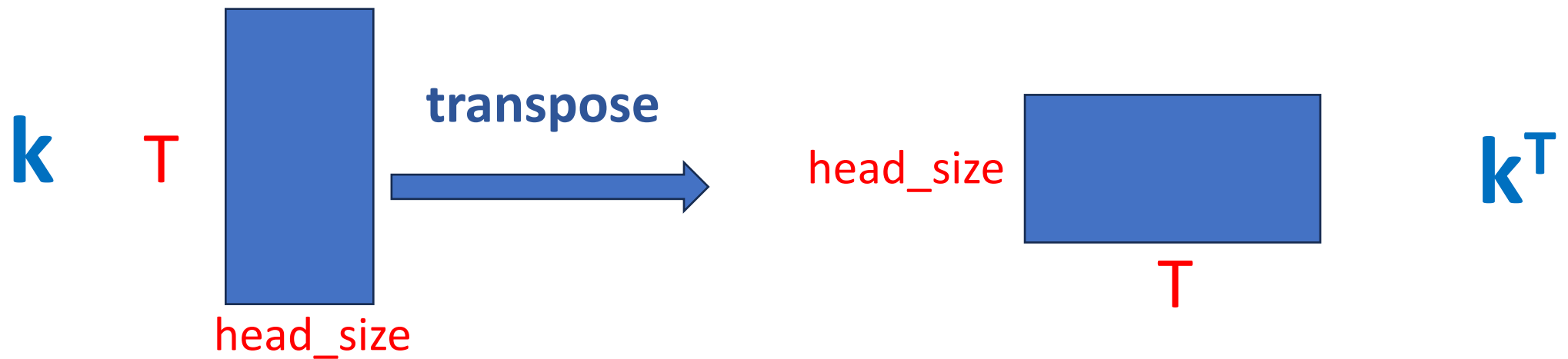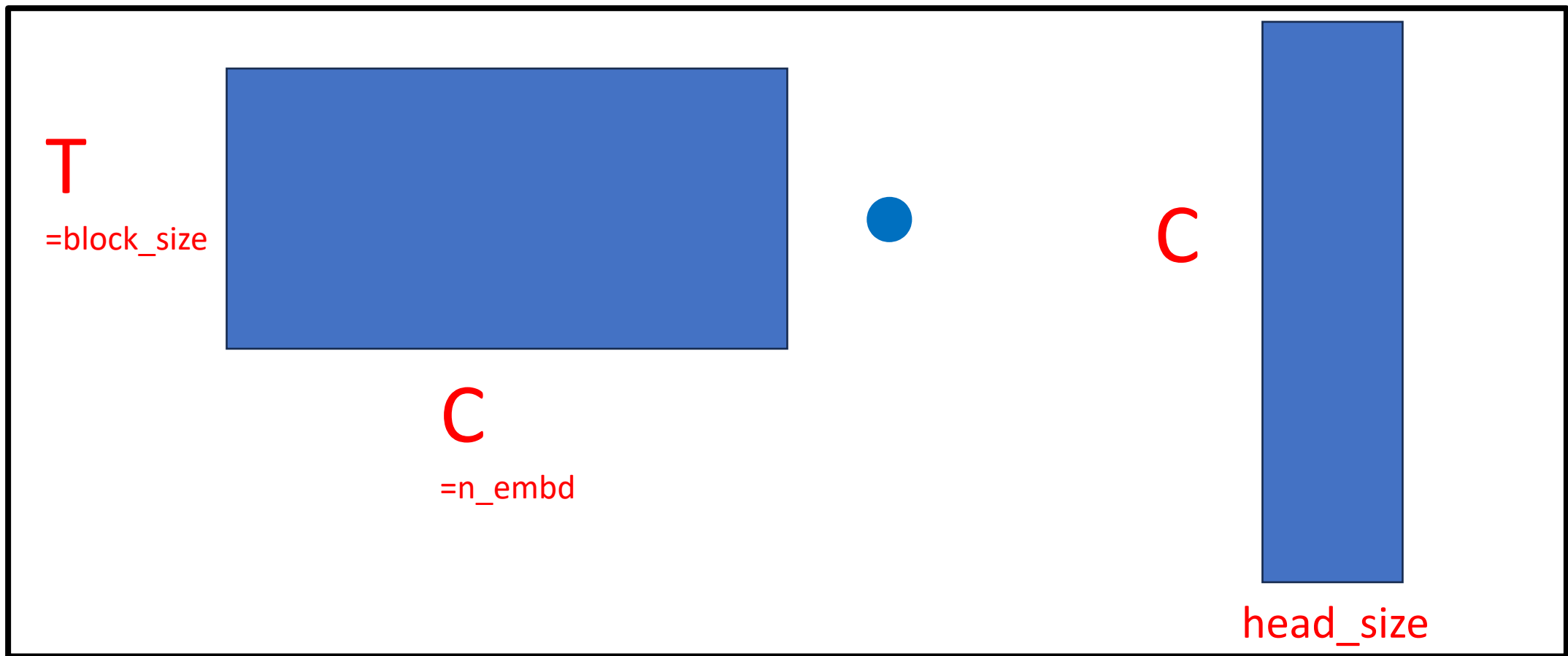
16 slates stacked together

each slate is a sequence of chars in the batch and is a matrix of shape (T, C)

# Visualize a single head

T
=block_size

C

C
=n_embd

head_size

T

head_size

The query matrix **q** of a sequence

$$\mathbf{q}^{chr1} \begin{bmatrix} q_1^{chr1} & q_2^{chr1} & \cdots & q_{16}^{chr1} \\ q_1^{chr2} & q_2^{chr2} & \cdots & q_{16}^{chr2} \\ \cdots & \cdots & \cdots & \cdots \\ q_1^{chr32} & q_2^{chr32} & \cdots & q_{16}^{chr32} \end{bmatrix}$$

$$\mathbf{k}^{\text{chr1}} \qquad \mathbf{k}^{\text{chr2}} \qquad \mathbf{\ldots} \qquad \mathbf{k}^{\text{chr32}}$$

$$\begin{bmatrix} k_1^{chr1} & k_1^{chr2} & \ldots & k_1^{chr32} \\ k_2^{chr1} & k_2^{chr2} & \ldots & k_2^{chr32} \\ \ldots & \ldots & \ldots & \ldots \\ k_{16}^{chr1} & k_{16}^{chr2} & \ldots & k_{16}^{chr32} \end{bmatrix}$$

$$\mathbf{q}^{\text{chr1}} \quad \mathbf{q}^{\text{chr2}} \quad \ldots \quad \mathbf{q}^{\text{chr32}}$$

$$
\begin{array}{c}
\mathbf{q}^{\text{chr1}} \\
\mathbf{q}^{\text{chr2}} \\
\ldots \\
\mathbf{q}^{\text{chr32}}
\end{array}
\begin{bmatrix}
q_1^{chr1} & q_2^{chr1} & \cdots & q_{16}^{chr1} \\
q_1^{chr2} & q_2^{chr2} & \cdots & q_{16}^{chr2} \\
\cdots & \cdots & \cdots & \cdots \\
q_1^{chr32} & q_2^{chr32} & \cdots & q_{16}^{chr32}
\end{bmatrix}
*
\begin{matrix}
\mathbf{k}^{\text{chr1}} & \mathbf{k}^{\text{chr2}} & \mathbf{...} & \mathbf{k}^{\text{chr32}} \\
\begin{bmatrix}
k_1^{chr1} & k_1^{chr2} & \cdots & k_1^{chr32} \\
k_2^{chr1} & k_2^{chr2} & \cdots & k_2^{chr32} \\
\cdots & \cdots & \cdots & \cdots \\
k_{16}^{chr1} & k_{16}^{chr2} & \cdots & k_{16}^{chr32}
\end{bmatrix}
\end{matrix}
$$

$$(\mathsf{T}, \mathsf{T})$$

$$\begin{bmatrix} \mathbf{q}^{chr1} * \mathbf{k}^{chr1} & \mathbf{q}^{chr1} * \mathbf{k}^{chr2} & \dots & \mathbf{q}^{chr1} * \mathbf{k}^{chr32} \\ \mathbf{q}^{chr2} * \mathbf{k}^{chr1} & \mathbf{q}^{chr2} * \mathbf{k}^{chr2} & \dots & \mathbf{q}^{chr2} * \mathbf{k}^{chr32} \\ \dots & \dots & \dots & \dots \\ \mathbf{q}^{chr32} * \mathbf{k}^{chr1} & \mathbf{q}^{chr32} * \mathbf{k}^{chr2} & \dots & \mathbf{q}^{chr32} * \mathbf{k}^{chr32} \end{bmatrix}$$

$$(\mathsf{T}, \mathsf{T})$$

$$
\begin{bmatrix}
\mathbf{q}^{chr1} * \mathbf{k}^{chr1} & \mathbf{q}^{chr1} * \mathbf{k}^{chr2} & \dots & \mathbf{q}^{chr1} * \mathbf{k}^{chr32} \\
\mathbf{q}^{chr2} * \mathbf{k}^{chr1} & \mathbf{q}^{chr2} * \mathbf{k}^{chr2} & \dots & \mathbf{q}^{chr2} * \mathbf{k}^{chr32} \\
\dots & \dots & \dots & \dots \\
\mathbf{q}^{chr32} * \mathbf{k}^{chr1} & \mathbf{q}^{chr32} * \mathbf{k}^{chr2} & \dots & \mathbf{q}^{chr32} * \mathbf{k}^{chr32}
\end{bmatrix}
$$

$$(T, T)$$

$$\begin{bmatrix} \mathbf{q}^{chr1} * \mathbf{k}^{chr1} & \mathbf{-inf} & \ldots & \ldots & \mathbf{-inf} \\ \mathbf{q}^{chr2} * \mathbf{k}^{chr1} & \mathbf{q}^{chr2} * \mathbf{k}^{chr2} & \mathbf{-inf} & \ldots & \mathbf{-inf} \\ \ldots & \ldots & \ldots & \ldots & \ldots \\ \mathbf{q}^{chr32} * \mathbf{k}^{chr1} & \mathbf{q}^{chr32} * \mathbf{k}^{chr2} & \ldots & \ldots & \mathbf{q}^{chr32} * \mathbf{k}^{chr32} \end{bmatrix}$$

$$(T, T)$$

$$\begin{bmatrix} \mathbf{q}^{chr1} * \mathbf{k}^{chr1} & \mathbf{-inf} & \dots & \dots & \mathbf{-inf} \\ \mathbf{q}^{chr2} * \mathbf{k}^{chr1} & \mathbf{q}^{chr2} * \mathbf{k}^{chr2} & \mathbf{-inf} & \dots & \mathbf{-inf} \\ \dots & \dots & \dots & \dots & \dots \\ \mathbf{q}^{chr32} * \mathbf{k}^{chr1} & \mathbf{q}^{chr32} * \mathbf{k}^{chr2} & \dots & \dots & \mathbf{q}^{chr32} * \mathbf{k}^{chr32} \end{bmatrix}$$
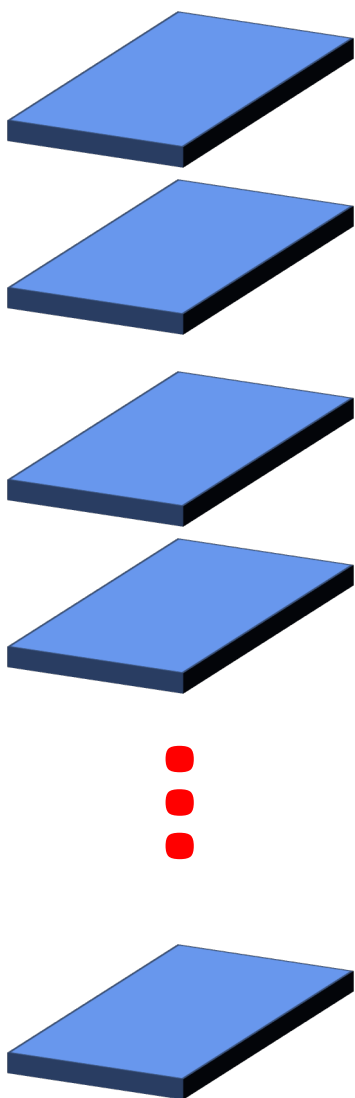
**Softmax to each row**

$$\begin{bmatrix} w_{\mathbf{q}^{chr1}*\mathbf{k}^{chr1}} & 0 & \dots & \dots & 0 \\ w_{\mathbf{q}^{chr2}*\mathbf{k}^{chr1}} & w_{\mathbf{q}^{chr2}*\mathbf{k}^{chr2}} & 0 & \dots & 0 \\ \dots & \dots & \dots & \dots & 0 \\ w_{\mathbf{q}^{chr32}*\mathbf{k}^{chr1}} & w_{\mathbf{q}^{chr32}*\mathbf{k}^{chr2}} & \dots & \dots & w_{\mathbf{q}^{chr32}*\mathbf{k}^{chr32}} \end{bmatrix}$$

(T, T)

sum up to 1

sum up to 1

sum up to 1

sum up to 1

$$(\text{T, T})$$

$$
\begin{bmatrix}
1 & 0 & 0 & \dots & 0 \\
0.22 & 0.78 & 0 & \dots & 0 \\
\dots & \dots & \dots & \dots & \dots \\
0.02 & 0.12 & 0.24 & \dots & 0.45
\end{bmatrix}
\quad
\begin{aligned}
&\text{sum up to 1} \\
&\text{sum up to 1} \\
&\text{sum up to 1} \\
&\text{sum up to 1}
\end{aligned}
$$

$$
\underbrace{\begin{bmatrix} w_{\mathbf{q}^{chr1}*\mathbf{k}^{chr1}} & 0 & \dots & \dots & 0 \\ w_{\mathbf{q}^{chr2}*\mathbf{k}^{chr1}} & w_{\mathbf{q}^{chr2}*\mathbf{k}^{chr2}} & 0 & \dots & 0 \\ \dots & \dots & \dots & \dots & 0 \\ w_{\mathbf{q}^{chr32}*\mathbf{k}^{chr1}} & w_{\mathbf{q}^{chr32}*\mathbf{k}^{chr2}} & \dots & \dots & w_{\mathbf{q}^{chr32}*\mathbf{k}^{chr32}} \end{bmatrix}}_{(T,\ T)} \bullet \underbrace{\begin{bmatrix} v_1^{chr1} & v_2^{chr1} & \dots & v_{16}^{chr1} \\ v_1^{chr2} & v_2^{chr2} & \dots & v_{16}^{chr2} \\ \dots & \dots & \dots & \dots \\ v_1^{chr32} & v_2^{chr32} & \dots & v_{16}^{chr32} \end{bmatrix}}_{(T,\ head\_size)}
$$

16 slates stacked together

each slate is of shape (T, head_size) and is a matrix of weighted values of one sequence in the batch

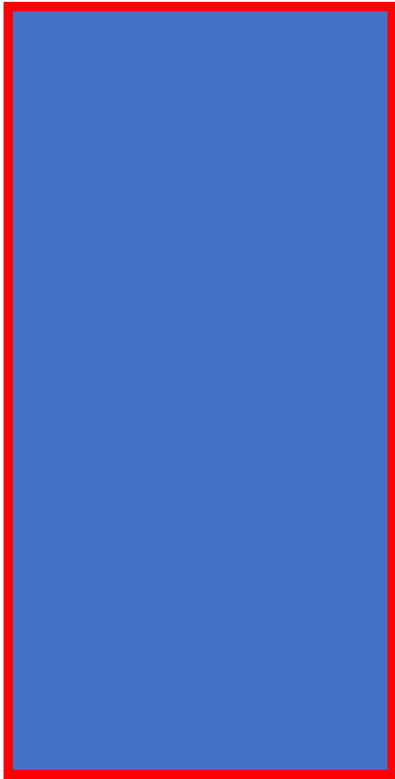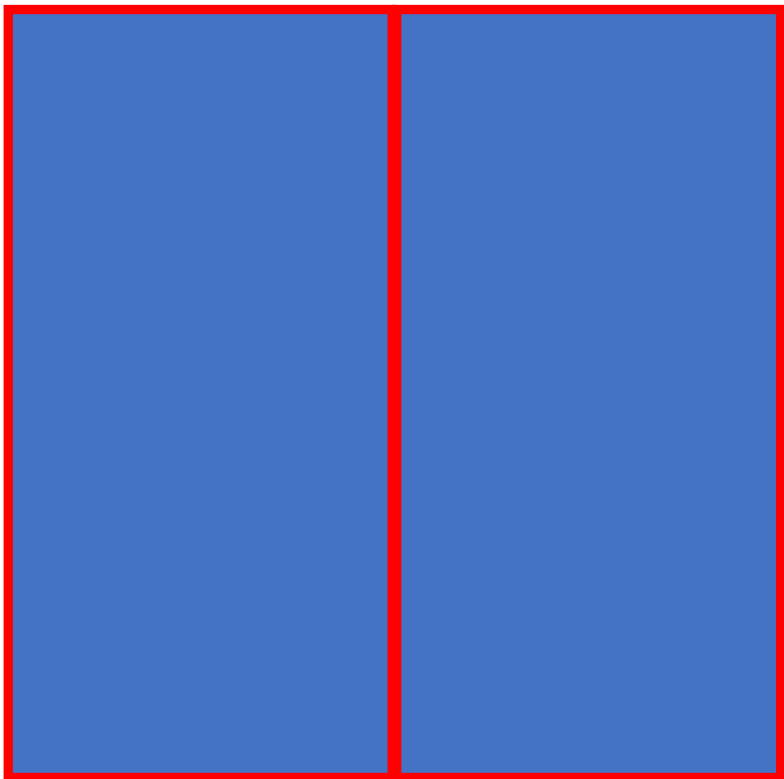The cuboid formed by these slates is of shape (B, T, head_size)
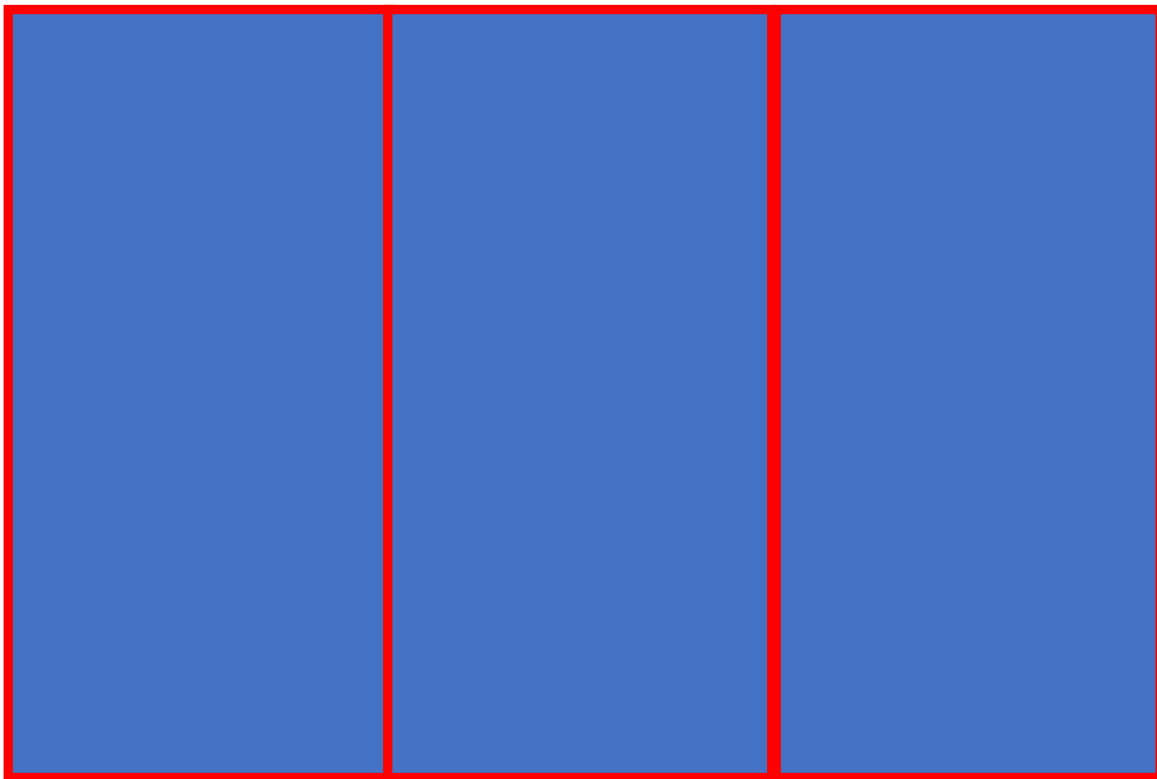
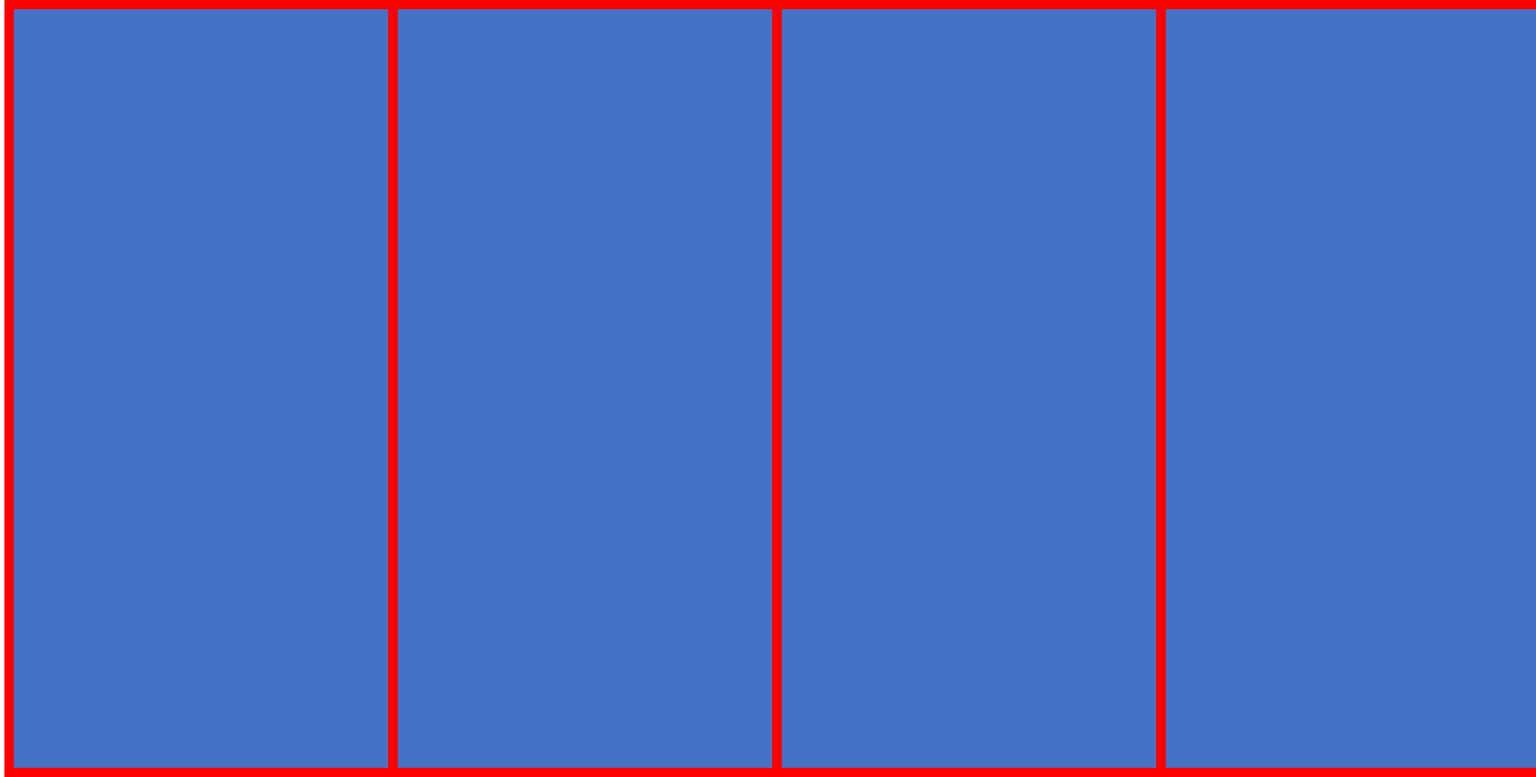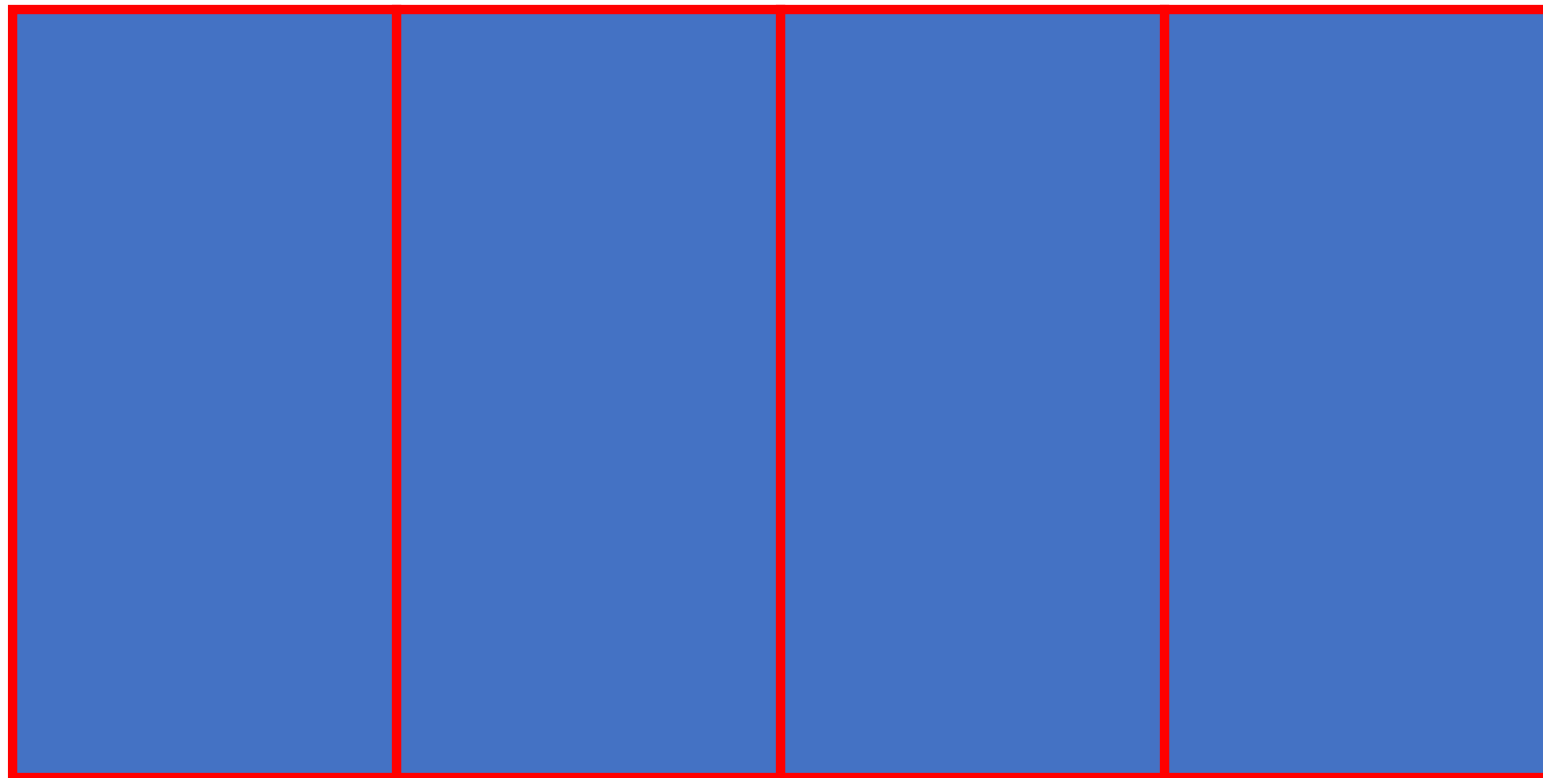head_size

B

T

(B, T, head_size)

# Visualize a multi-head

T

head_size * 4

C = head_size * num_heads

T

C

(32, 64)

$$
\begin{bmatrix}
a^{h1}_{chr1_{v1}} & \cdots & a^{h1}_{chr1_{v16}} \\
a^{h1}_{chr2_{v1}} & \cdots & a^{h1}_{chr2_{v16}} \\
\cdots & \cdots & \cdots \\
a^{h1}_{chr32_{v1}} & \cdots & a^{h1}_{chr32_{v16}}
\end{bmatrix}
\quad
\begin{bmatrix}
a^{h2}_{chr1_{v1}} & \cdots & a^{h2}_{chr1_{v16}} \\
a^{h2}_{chr2_{v1}} & \cdots & a^{h2}_{chr2_{v16}} \\
\cdots & \cdots & \cdots \\
a^{h2}_{chr32_{v1}} & \cdots & a^{h2}_{chr32_{v16}}
\end{bmatrix}
\quad
\begin{bmatrix}
a^{h3}_{chr1_{v1}} & \cdots & a^{h3}_{chr1_{v16}} \\
a^{h3}_{chr2_{v1}} & \cdots & a^{h3}_{chr2_{v16}} \\
\cdots & \cdots & \cdots \\
a^{h3}_{chr32_{v1}} & \cdots & a^{h3}_{chr32_{v16}}
\end{bmatrix}
\quad
\begin{bmatrix}
a^{h4}_{chr1_{v1}} & \cdots & a^{h4}_{chr1_{v16}} \\
a^{h4}_{chr2_{v1}} & \cdots & a^{h4}_{chr2_{v16}} \\
\cdots & \cdots & \cdots \\
a^{h4}_{chr32_{v1}} & \cdots & a^{h4}_{chr32_{v16}}
\end{bmatrix}
$$

HEAD 1  HEAD 2  HEAD 3  HEAD 4

C

T

16 slates stacked together

each slate is an output from a multi-head layer and is a matrix of shape (T, C)