**Inference T=1**

$$Q \quad \times \quad K^T \quad = \quad QK^T \quad \times \quad V \quad = \quad \text{Attention}$$

| Token 1 |

Token 1

$q_1 k_1$

| Token 1 |

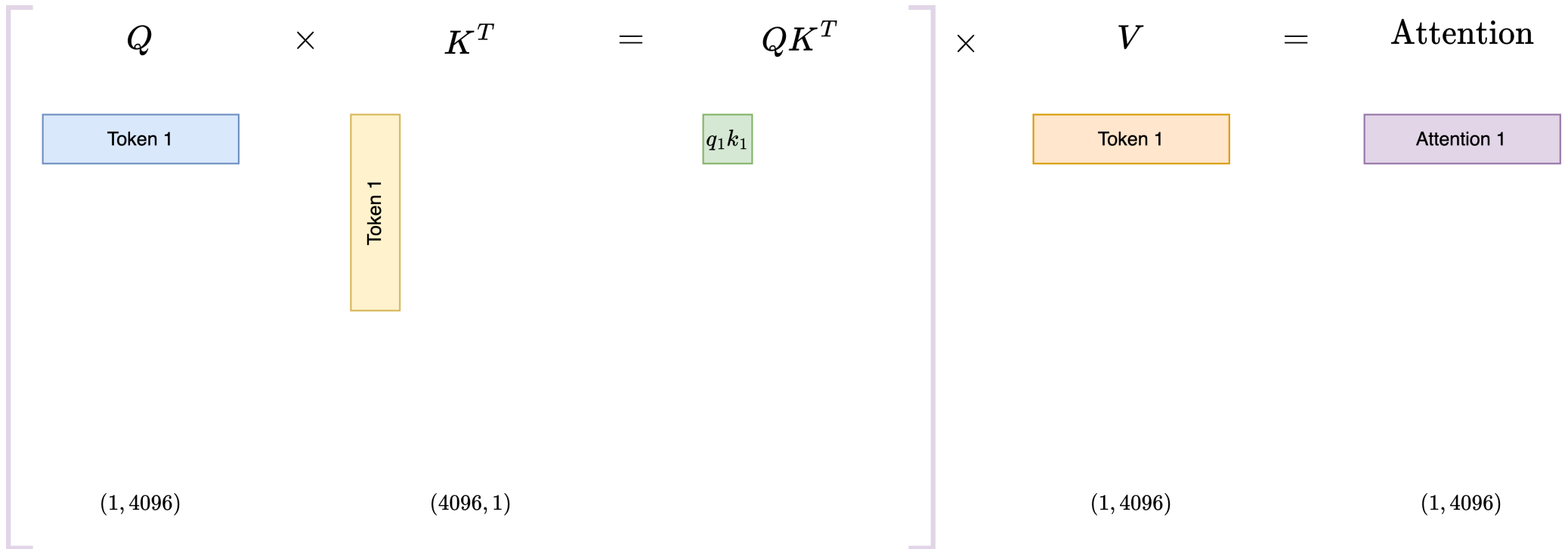| Attention 1 |

$$(1, 4096) \qquad (4096, 1) \qquad (1, 4096) \qquad (1, 4096)$$

$$\text{Attention}(Q, K, V) = \text{Softmax}\left(\frac{QK^T}{\sqrt{d_K}}\right)V$$

**Inference T=2**

$$Q \quad \times \quad K^T \quad = \quad QK^T \quad \times \quad V \quad = \quad \text{Attention}$$

| Token 1 |
|---------|
| Token 2 |

Token 1 | Token 2

| $q_1 k_1$ | $q_1 k_2$ |
|-----------|-----------|
| $q_2 k_1$ | $q_2 k_2$ |

| Token 1 |
|---------|
| Token 2 |

| Attention 1 |
|-------------|
| Attention 2 |

$(2, 4096)$        $(4096, 2)$        $(2, 4096)$        $(2, 4096)$

$$\text{Attention}(Q, K, V) = \text{Softmax}\left(\frac{QK^T}{\sqrt{d_K}}\right)V$$

2

**Inference T=3**

$$Q \quad \times \quad K^T \quad = \quad QK^T \quad \times \quad V \quad = \quad \text{Attention}$$

| Token 1 |
|---------|
| Token 2 |
| Token 3 |

Token 1   Token 2   Token 3

| $q_1k_1$ | $q_1k_2$ | $q_1k_3$ |
|----------|----------|----------|
| $q_2k_1$ | $q_2k_2$ | $q_2k_3$ |
| $q_3k_1$ | $q_3k_2$ | $q_3k_3$ |

| Token 1 |
|---------|
| Token 2 |
| Token 3 |

| Attention 1 |
|-------------|
| Attention 2 |
| Attention 3 |

$(3, 4096)$      $(4096, 3)$      $(3, 4096)$      $(3, 4096)$

$$\text{Attention}(Q, K, V) = \text{Softmax}\left(\frac{QK^T}{\sqrt{d_K}}\right)V$$

3

$$Q \quad \times \quad K^T \quad = \quad QK^T \quad \times \quad V \quad = \quad \text{Attention}$$

| Token 1 |
|---|
| Token 2 |
| Token 3 |
| Token 4 |

Token 1 | Token 2 | Token 3 | Token 4

| $q_1k_1$ | $q_1k_2$ | $q_1k_3$ | $q_1k_4$ |
|---|---|---|---|
| $q_2k_1$ | $q_2k_2$ | $q_2k_3$ | $q_2k_4$ |
| $q_3k_1$ | $q_3k_2$ | $q_3k_3$ | $q_3k_4$ |
| $q_4k_1$ | $q_4k_2$ | $q_4k_3$ | $q_4k_4$ |

| Token 1 |
|---|
| Token 2 |
| Token 3 |
| Token 4 |

| Attention 1 |
|---|
| Attention 2 |
| Attention 3 |
| Attention 4 |

$$(4, 4096) \qquad (4096, 4) \qquad\qquad (4, 4096) \qquad (4, 4096)$$
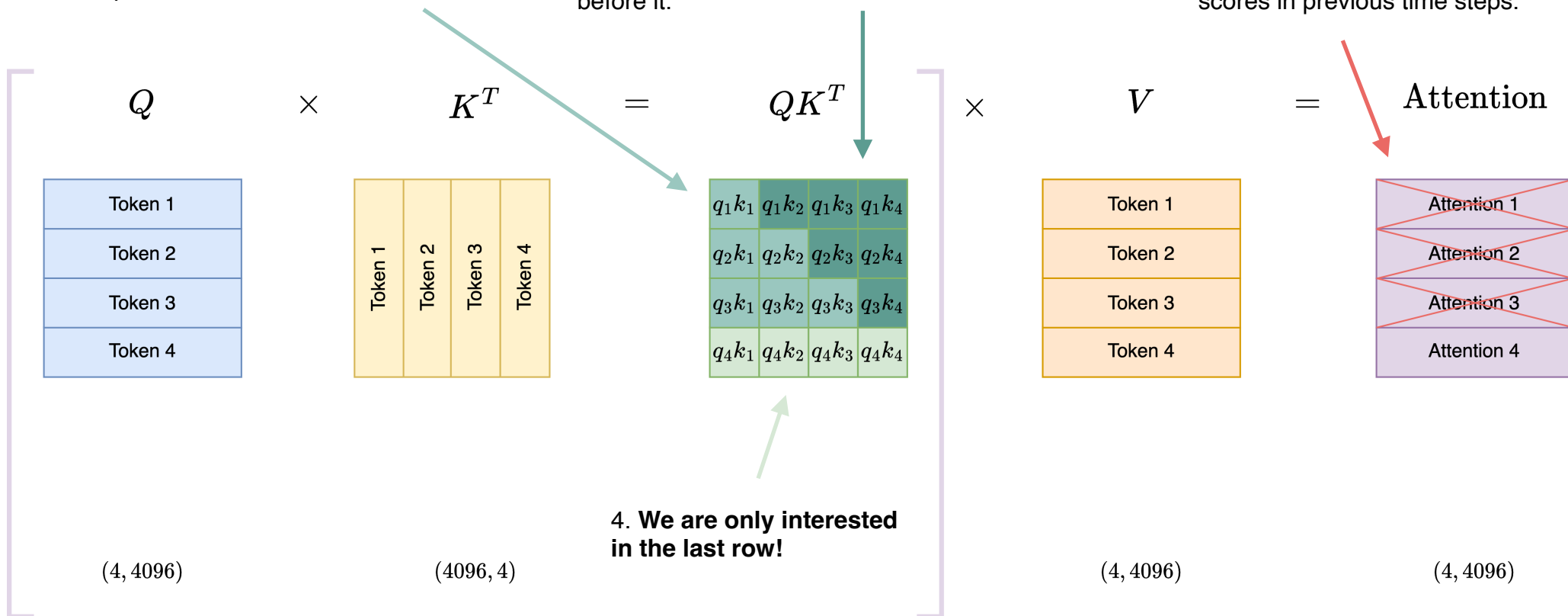
$$\text{Attention}(Q, K, V) = \text{Softmax}\left(\frac{QK^T}{\sqrt{d_K}}\right)V$$

1. We already computed these dot products in the previous time steps. **Can we cache them?**

2. Since the model is causal, **we don't care about the attention of a token with its successors**, but only with the tokens before it.

3. **We don't care about these,** as we want to predict the next token and we have already predicted these attention scores in previous time steps.

$Q$ $\times$ $K^T$ $=$ $QK^T$ $\times$ $V$ $=$ Attention

| Token 1 |
| Token 2 |
| Token 3 |
| Token 4 |

Token 1 | Token 2 | Token 3 | Token 4

| $q_1k_1$ | $q_1k_2$ | $q_1k_3$ | $q_1k_4$ |
| $q_2k_1$ | $q_2k_2$ | $q_2k_3$ | $q_2k_4$ |
| $q_3k_1$ | $q_3k_2$ | $q_3k_3$ | $q_3k_4$ |
| $q_4k_1$ | $q_4k_2$ | $q_4k_3$ | $q_4k_4$ |

| Token 1 |
| Token 2 |
| Token 3 |
| Token 4 |

| Attention 1 |
| Attention 2 |
| Attention 3 |
| Attention 4 |

$(4, 4096)$ $\qquad$ $(4096, 4)$

4. **We are only interested in the last row!**

$(4, 4096)$ $\qquad$ $(4, 4096)$

$$\text{Attention}(Q, K, V) = \text{Softmax}\left(\frac{QK^T}{\sqrt{d_K}}\right)V$$