

Inference on Nano GPT - for Beginners

Pre-Requisites

Auto Regressive task ?

Time Vs Space Complexity

Inference and setting it up

**Kv cacheing (although there is some issues in the code
here)**

Auto Regressive Task

The current output is depends upon the series of input which the model have seen through which can be neatly said in the equation given below

$$P(x_t|x_0, \dots, x_{t-1}) = \prod_{i=1}^t p(x_i|x_{i-1})$$

But in Language Modeling to pick a particular token Transformer does this by adding decoder casual mask

$$token = K(P(x_t|x_0, \dots, x_{t-1}))$$

$$\begin{cases} \text{argmax if } K = \textit{greedysampling} \\ \text{Multinomial(topk) if } K = \textit{multinomalsampling} \end{cases}$$

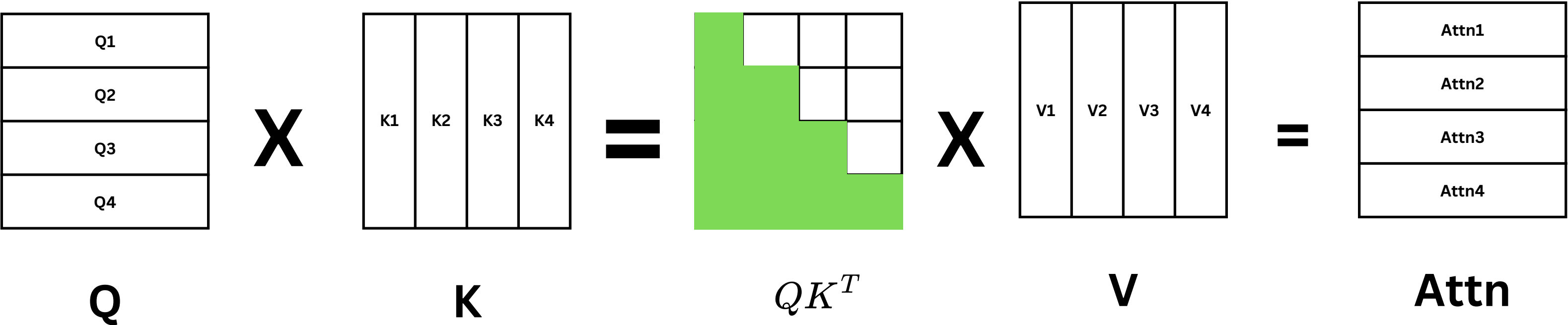
Time Vs Space complexity

- **Bubble sort: Time= $O(n^2)$, Space= $O(1)$**
 - **Merge sort: Time = $O(n \log n)$ space $O(\log n)$**
 - **Heap Sort : Time= $O(n \log n)$ space $O(n)$**
-
- We could observe that with Increasing space complexity we could kind of mitigate the time if correctly used.
 - if space could be increased reasonable amount Time complexity could be kind of reduced

Whats that with transformers

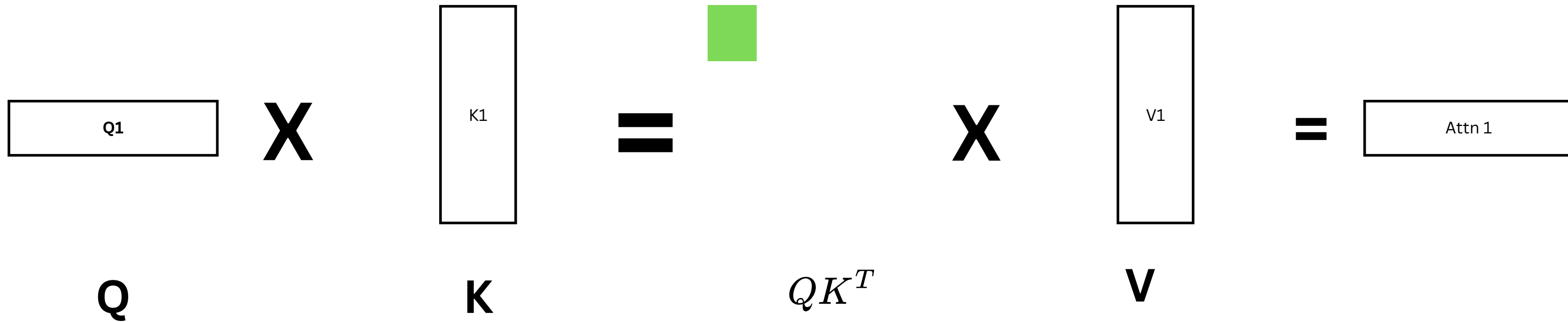
Inference on Transformers

$$SM \left(\frac{QK^T}{\sqrt{d}} \right) V$$



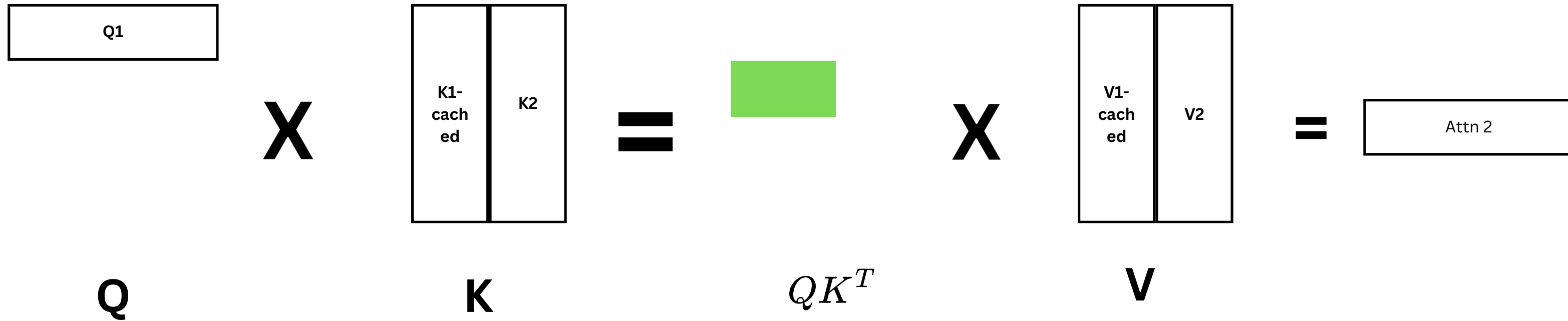
KV cacheing

step 1



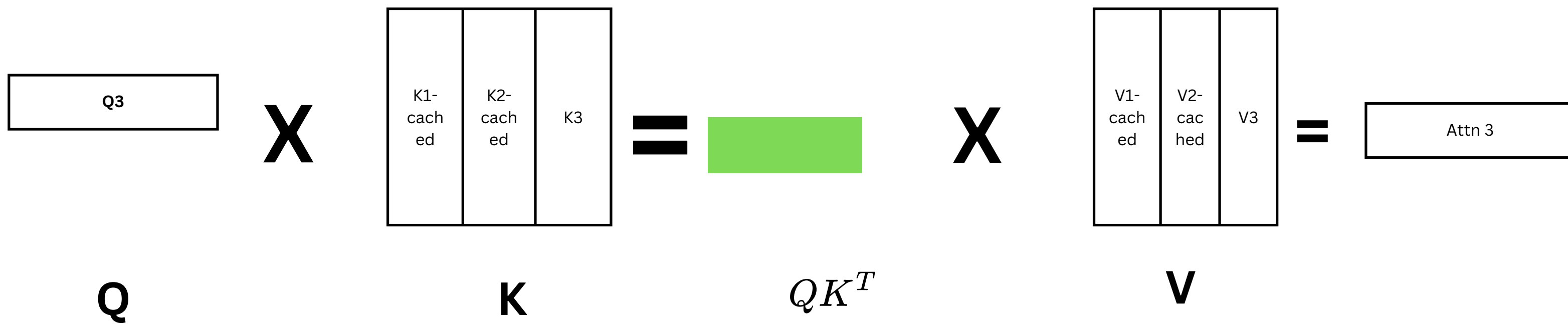
KV cacheing

step 1



KV cacheing

step 1



KV cacheing

