

Attention

Discussion highlights of Tuesday 9/11/2021

Typical attention operation:

- Here we start by $qkv = \text{concat}(Q, K, V)$ all together.
- Keys, Queries, Values are $K = XW^K$, $Q = XW^Q$, $V = XW^V$
- For the attention heatmap $\mathcal{A}(K, Q) = \text{softmax}\left(QK^T / \sqrt{d_K}\right)$
- Then attention scores are $\text{Attention}(K, Q, V) = \mathcal{A}(K, Q) \cdot V$

$$X \in \mathbb{R}^{1,197,768}$$

$$qkv \in \mathbb{R}^{1,197,3 \times 768} = \mathbb{R}^{1,197,2304}$$

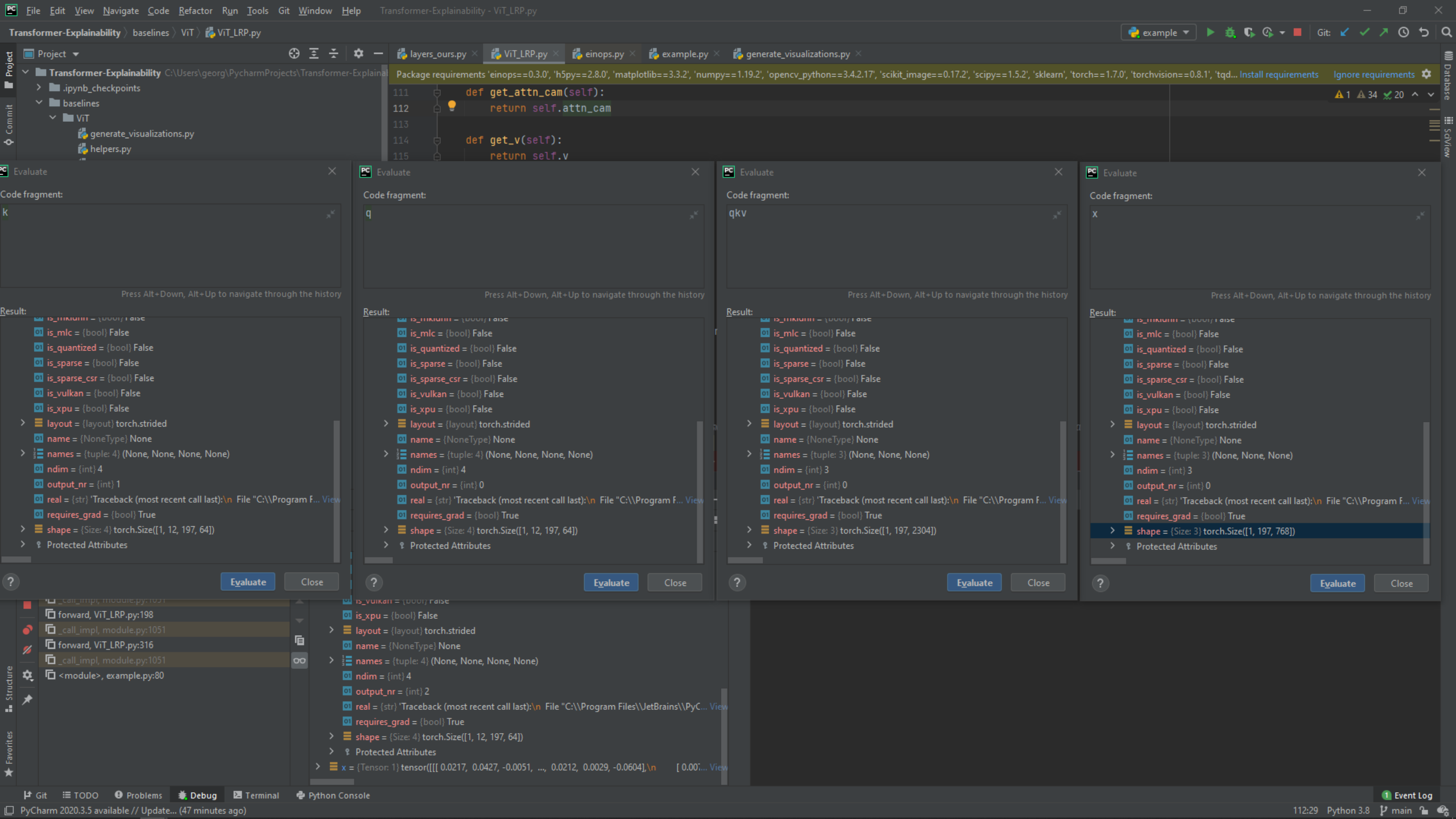
$$q \in \mathbb{R}^{1,12,197,64}$$

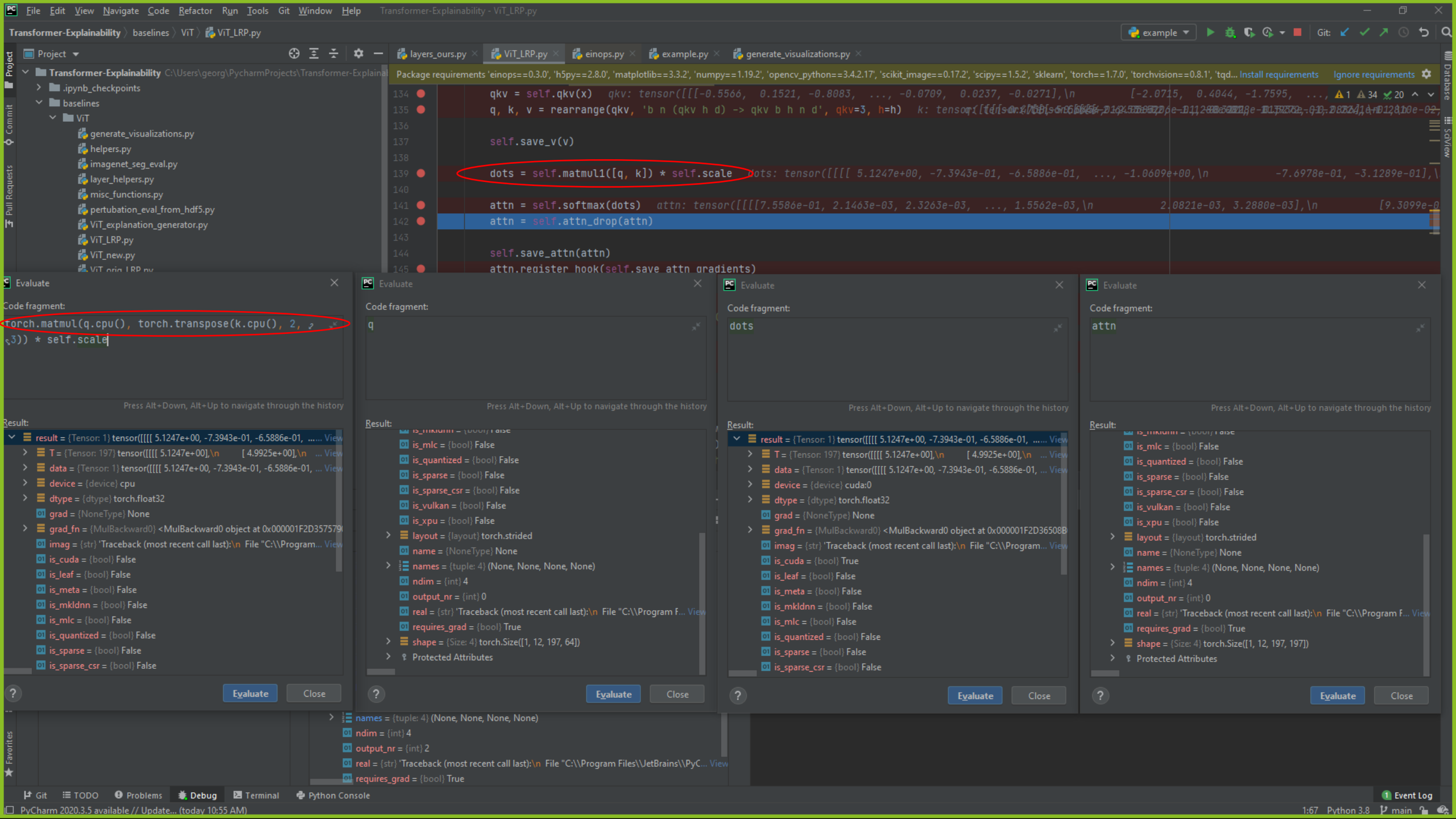
$$k \in \mathbb{R}^{1,12,197,64}$$

$$v \in \mathbb{R}^{1,12,197,64}$$

$$768 = 12 \times 64$$

$$\sqrt{d_K} = \sqrt{\frac{768}{12}} = 0.125$$





Four screenshots of the PyCharm 'Evaluate' window showing the results of different code fragments. Each window has a 'Code fragment' field at the top and a 'Result' field below it. The first window shows the result of a matrix multiplication and transpose operation. The second window shows the result of a softmax operation. The third window shows the result of a dot product operation. The fourth window shows the result of an attention operation. Each result is a dictionary containing various attributes like 'T', 'data', 'device', 'dtype', 'grad', 'imag', 'is_cuda', 'is_leaf', 'is_meta', 'is_mkldnn', 'is_mlc', 'is_quantized', 'is_sparse', and 'is_sparse_csr'.

