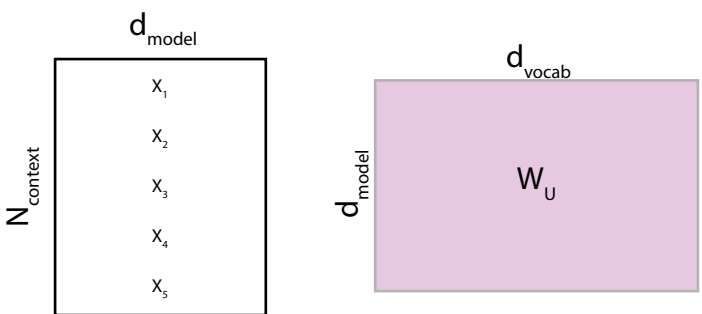


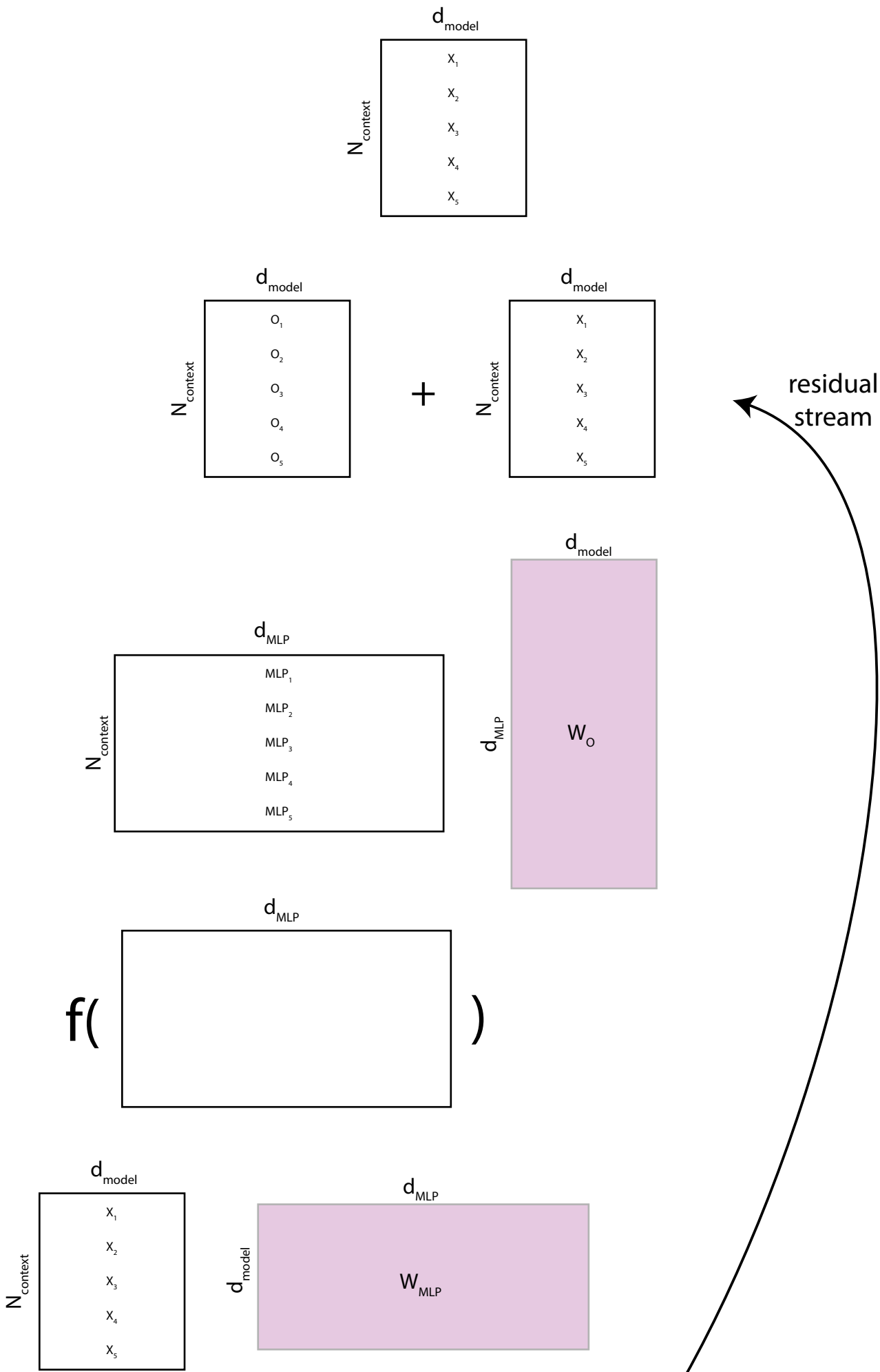
We take the final residual stream values and linearly transform them back to $d_{vocab'}$ and then take a softmax row-wise. Now, for every position in the context window, we have a probability distribution over our vocabulary.

$\text{softmax}(\frac{1}{N_{\text{context}}} \begin{bmatrix} \vdots \\ \vdots \\ \vdots \end{bmatrix})$



More layers of Multiheaded Attention and MLPs

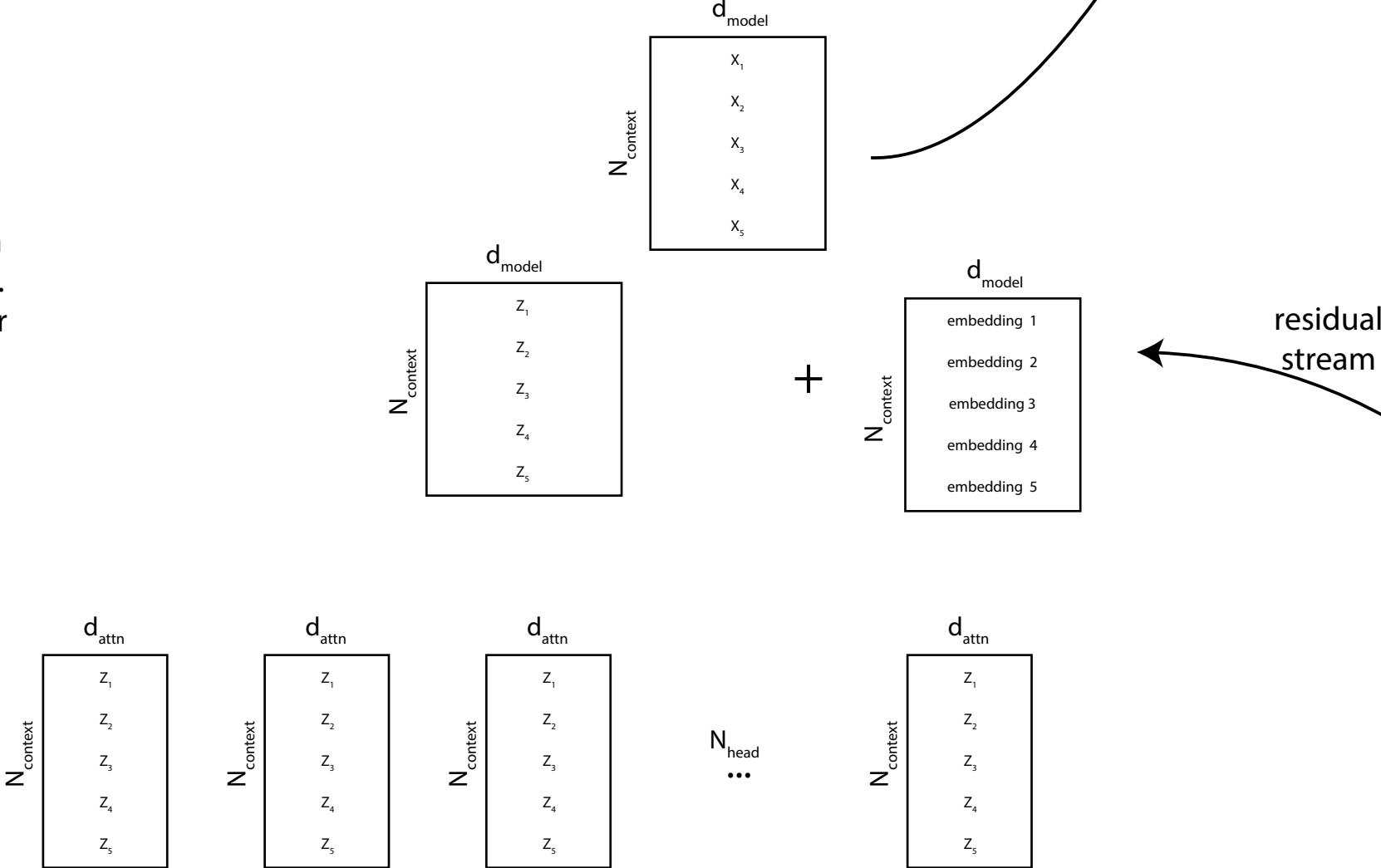
We take the residual stream values and run them each individually through a MLP layer, which is just a linear transform (to a higher dimension) followed by a nonlinearity. We then use another matrix to go back to the residual stream dimension, d_{model} . We add our result back to the residual stream.



Multilayer Perceptron

Do the attention mechanism in parallel N_{head} times, with different W_Q, W_K and W_V learned matrices. Concatenate the results together so that we are back in d_{model} dimensions. We then add this back to the matrix we sent into the attention mechanism (the residual stream).

Multihead Head Attention



Single Head Attention

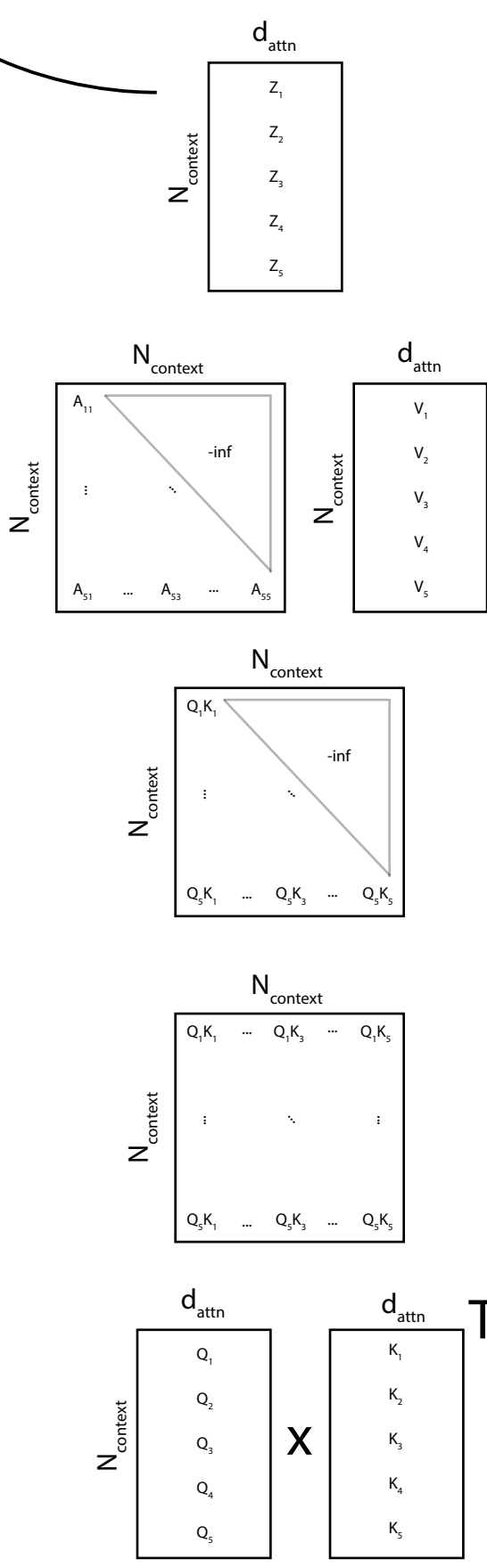
We use the attention patterns to take weighted averages of our V-vectors. We call the result Z.

We then softmax each row in order to make every row a probability distrution. We call this the attention pattern.

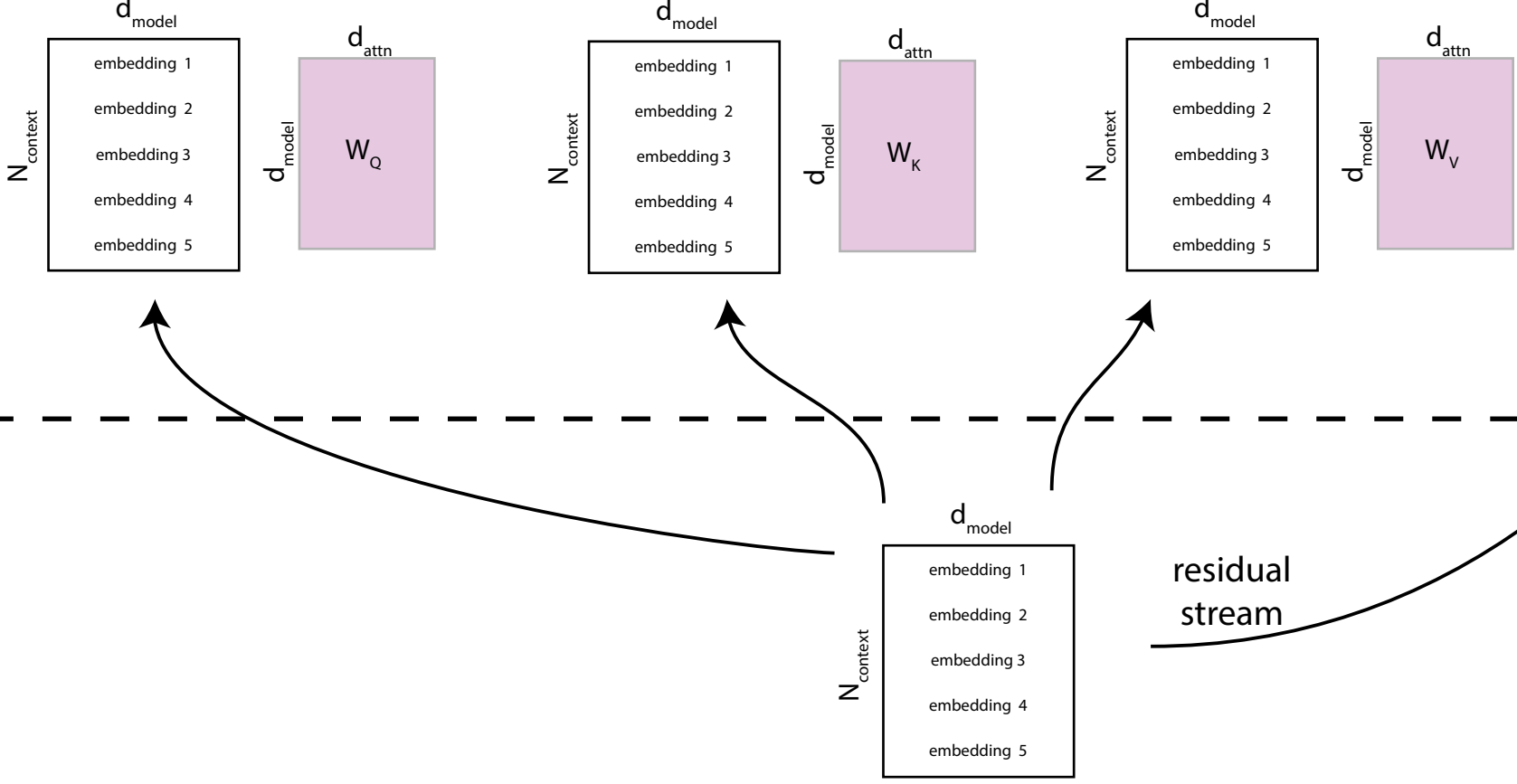
We don't want a position i to be able to get information for predicting the next token from K's from the future, so we make all positions $j > i$ to -inf.

Think of $Q_i K_j$ as describing how important K_j is for predicting the i-th positions next token, given Q_i .

Take an outer product of the Q and K vectors, so that our new matrix gives the dot product between every pair of Q and K vectors.

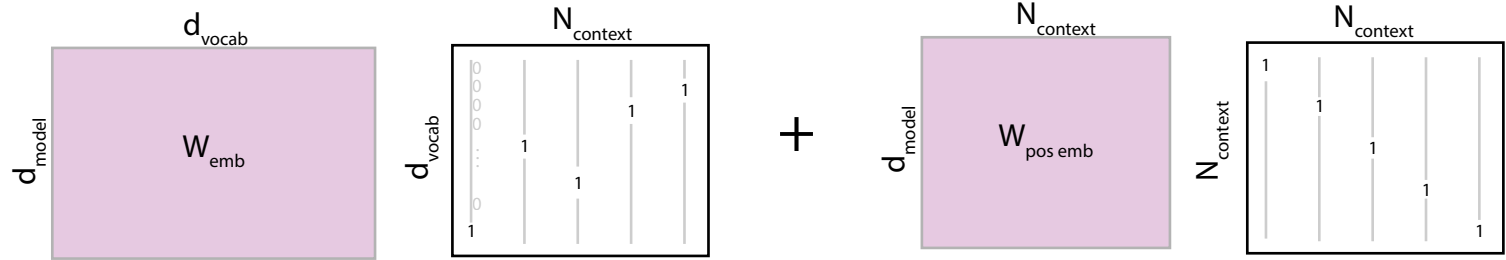


Convert the embeddings into 3 different attn dimensional vectors using $W_{Q'}, W_K$ and $W_{V'}$.



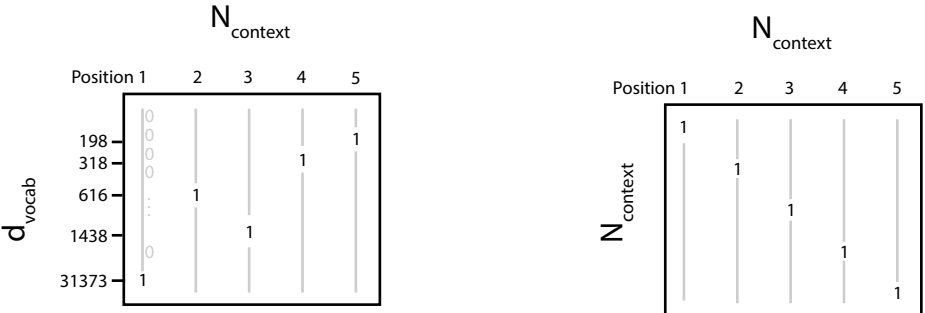
Embedding

Use W_{emb} and $W_{\text{pos emb}}$ to convert the token and positional one hot encodings into the model dimensions. We add the token and positional embeddings together to get a final embedding. The final embedding enters the residual stream.



One Hot Encoding

One hot encode the string of integers such that the j-th column is a vector of 0s with a 1 in the x_j -th position. We do the same for the positions, j , which just gives us an $N_{\text{context}} \times N_{\text{context}}$ identity matrix.



Tokenization

Tokenize the input string into a string of integers according to a predefined vocabulary. Call this string x_j with j between 1 and N_{context}

hello my name is

"hello my name is"