



SELF ATTENTION IN TRANSFORMERS

Consider two text statements:

- 1) money bank grows
- 2) river bank flows

Here, the meaning of word bank is different in both statements which the static word embedding won't be able to capture. So, we'll need contextual word embeddings.

So, we'll represent the two statements in this way -

money bank grows

river bank flows

money: $0.7 \text{ money} + 0.2 \text{ bank} + 0.1 \text{ grows}$
bank: $0.25 \text{ money} + 0.7 \text{ bank} + 0.05 \text{ grows}$
grows: $0.1 \text{ money} + 0.2 \text{ bank} + 0.7 \text{ grows}$

river: $0.8 \text{ river} + 0.15 \text{ bank} + 0.05 \text{ flows}$
bank: $0.2 \text{ river} + 0.78 \text{ bank} + 0.02 \text{ flows}$
flows: $0.4 \text{ river} + 0.01 \text{ bank} + 0.59 \text{ flows}$

Now here as we can see, each word of a statement is now depending on other words of their respective statement. This helps retain contextual information of the statements. But machine doesn't understand words so we need to represent each word as it's embedding form.

⇒

$$e_{\text{money}}^{(\text{new})} = 0.7 e_{\text{money}} + 0.2 e_{\text{bank}} + 0.1 e_{\text{grows}}$$

e_{money} : Embedding of word money

$$e_{\text{bank}}^{(\text{new})} = 0.25 e_{\text{money}} + 0.7 e_{\text{bank}} + 0.05 e_{\text{grows}}$$

e_{bank} : Embedding of word bank

$$e_{\text{grows}}^{(\text{new})} = 0.1 e_{\text{money}} + 0.2 e_{\text{bank}} + 0.7 e_{\text{grows}}$$

e_{grows} : Embedding of word grows

Embeddings are vector representation of words & can be n dimensional



The nos. before each word embedding represents the weight that describes the similarity between each word embedding.
i.e. $e_{\text{money}}^{(\text{new})} = (e_{\text{money}} \cdot e_{\text{money}}^T) e_{\text{money}} + (e_{\text{money}} \cdot e_{\text{bank}}^T) e_{\text{bank}} + (e_{\text{money}} \cdot e_{\text{grows}}^T) e_{\text{grows}}$

$$e_{\text{bank}}^{(\text{new})} = (e_{\text{bank}} \cdot e_{\text{money}}^T) e_{\text{money}} + (e_{\text{bank}} \cdot e_{\text{bank}}^T) e_{\text{bank}} + (e_{\text{bank}} \cdot e_{\text{grows}}^T) e_{\text{grows}}$$

$$e_{\text{grows}}^{(\text{new})} = (e_{\text{grows}} \cdot e_{\text{money}}^T) e_{\text{money}} + (e_{\text{grows}} \cdot e_{\text{bank}}^T) e_{\text{bank}} + (e_{\text{grows}} \cdot e_{\text{grows}}^T) e_{\text{grows}}$$

* The similarity b/w each word embedding is calculated using simple dot product b/w them.

$$\Rightarrow e_{\text{money}}^{(\text{new})} = s_{11} e_{\text{money}} + s_{12} e_{\text{bank}} + s_{13} e_{\text{grows}}$$

$$e_{\text{bank}}^{(\text{new})} = s_{21} e_{\text{money}} + s_{22} e_{\text{bank}} + s_{23} e_{\text{grows}}$$

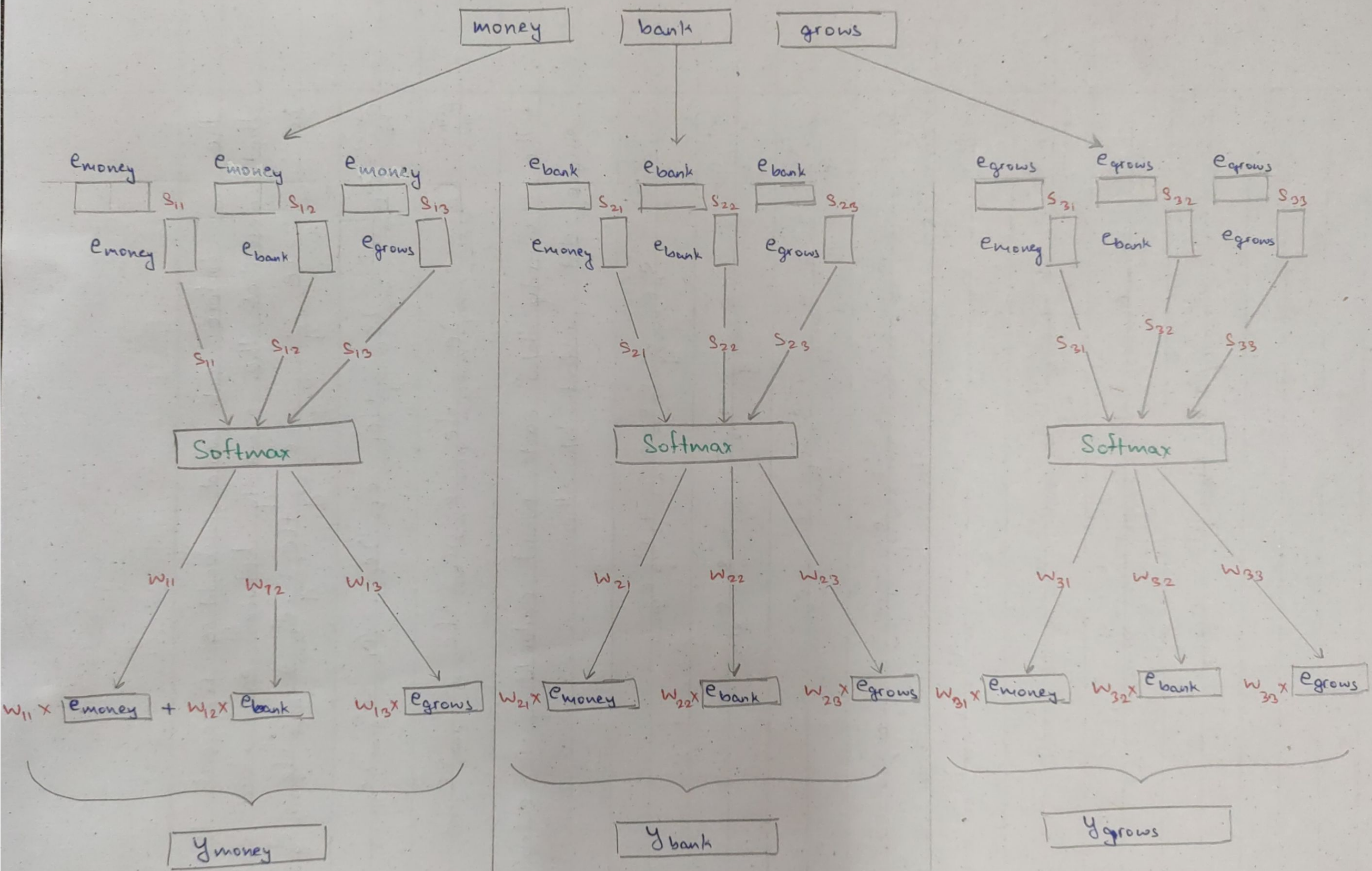
$$e_{\text{grows}}^{(\text{new})} = s_{31} e_{\text{money}} + s_{32} e_{\text{bank}} + s_{33} e_{\text{grows}}$$

$s_{11}, s_{12}, s_{13}, s_{21}, s_{22}, s_{23}, s_{31}, s_{32}, s_{33}$ are not normalized & are normalized using a softmax function resulting in normalized weights of $w_{11}, w_{12}, w_{13}, w_{21}, w_{22}, w_{23}, w_{31}, w_{32}, w_{33}$

$$\Rightarrow e_{\text{money}}^{(\text{new})} = w_{11} e_{\text{money}} + w_{12} e_{\text{bank}} + w_{13} e_{\text{grows}}$$

$$e_{\text{bank}}^{(\text{new})} = w_{21} e_{\text{money}} + w_{22} e_{\text{bank}} + w_{23} e_{\text{grows}}$$

$$e_{\text{grows}}^{(\text{new})} = w_{31} e_{\text{money}} + w_{32} e_{\text{bank}} + w_{33} e_{\text{grows}}$$





This whole process consists of three steps, 1 softmax operation & 2 dot product operations. Learnable parameters can be introduced in dot product operations.

Observing the diagram we can see that each word embedding vector is performing different roles at each step.

For word embedding vector e_{money} at first operation of dot product it is querying for the similarity of b/w it & all word embedding vectors & also returning the similarity during the dot product. In the last dot product it is being used as value to output final contextual word embedding vector for the word money. Same goes with other word embedding vectors.

∴, we can see that e_{money} & others, are being playing 3 roles in the whole process - query, key, & value. All 3 roles are being performed by the same word embedding, e_{money} for instance, which won't be able to do it properly. Instead, it should 3 separate vectors for performing those 3 roles, q_{money} for query, k_{money} for key, & v_{money} for value.

This is called separation of concerns as each concerned role has now it's specific vector for it which will improve performance.

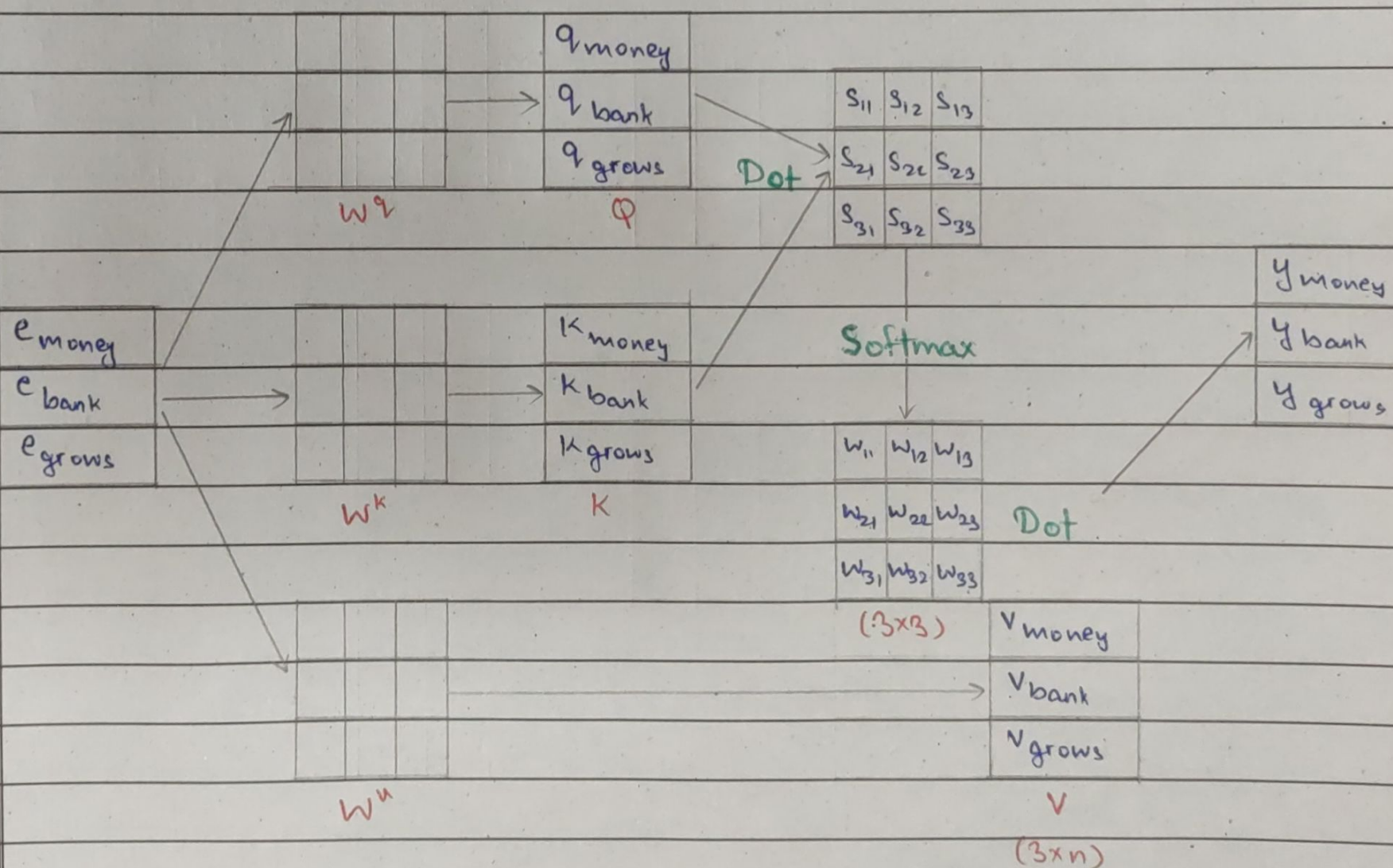
Now, to generate 3 new vectors from another vector we can use linear transformation where we can multiply the word embedding vector with 3 different weight matrices, namely W^q for query, W^k for key, & W^v for values respectively.



The values of the matrices can be randomly initialized which will get updated to best values while training & backpropagation. And we will do this for all word embedding vectors.

So now, our model architecture can also work parallelly & has learnable parameters to generate Task specific Contextual Word Embeddings.

So now, our final model architecture looks like this -



* The value of W^Q , W^K , W^V will be same for all word embeddings