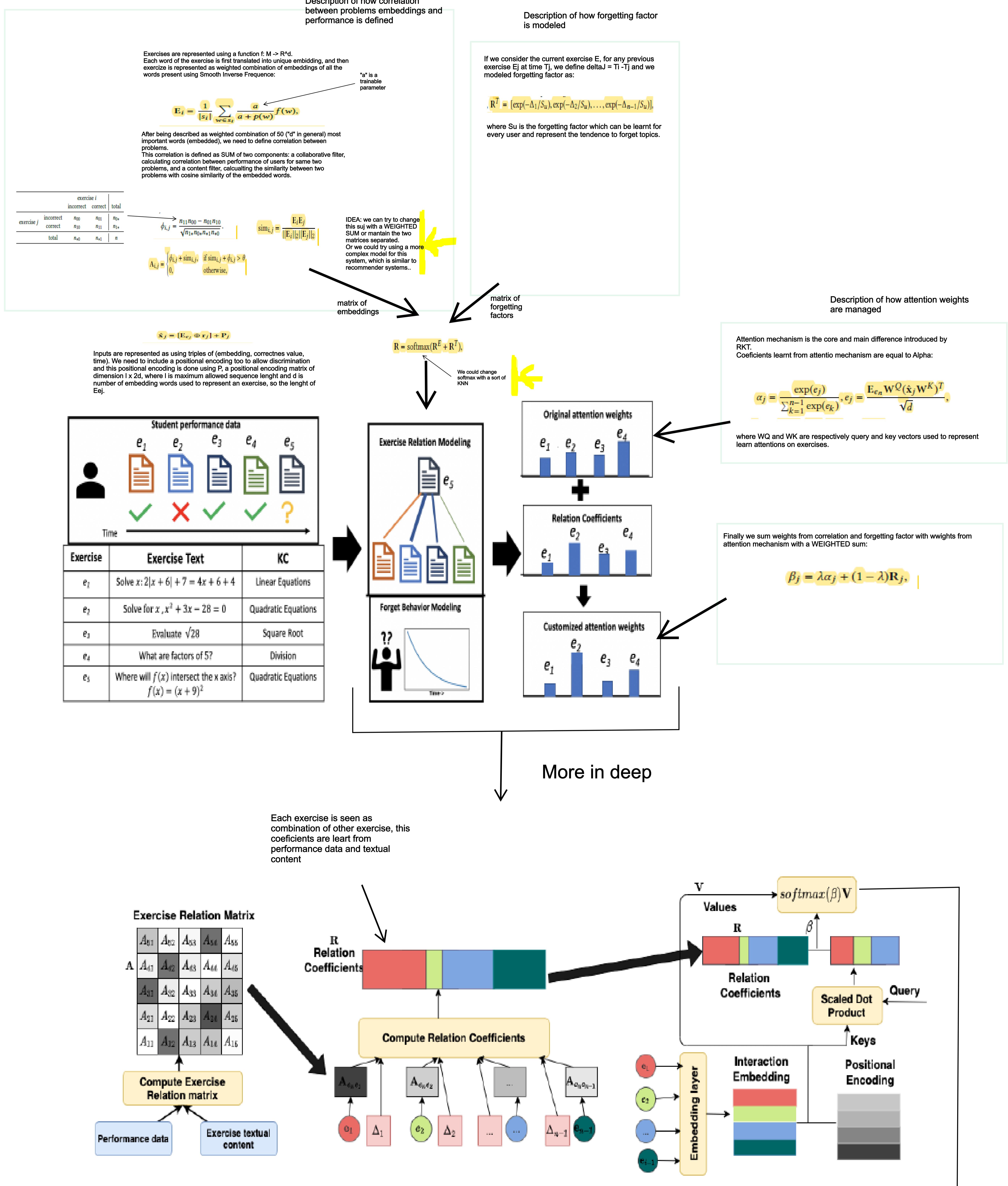


Relation Aware Self-attention Knowledge Tracing



In the end we obtain outputs from weights as:

$$\mathbf{o} = \sum_{j=1}^{n-1} \beta_j \hat{x}_j \mathbf{W}^V,$$

We add Always residual network and apply Normalization and Dropout

And this outputs are then given as input to a Feed Forward Neural Network, to insert possibility of relations between latent dimensions, introducing non-linearities:

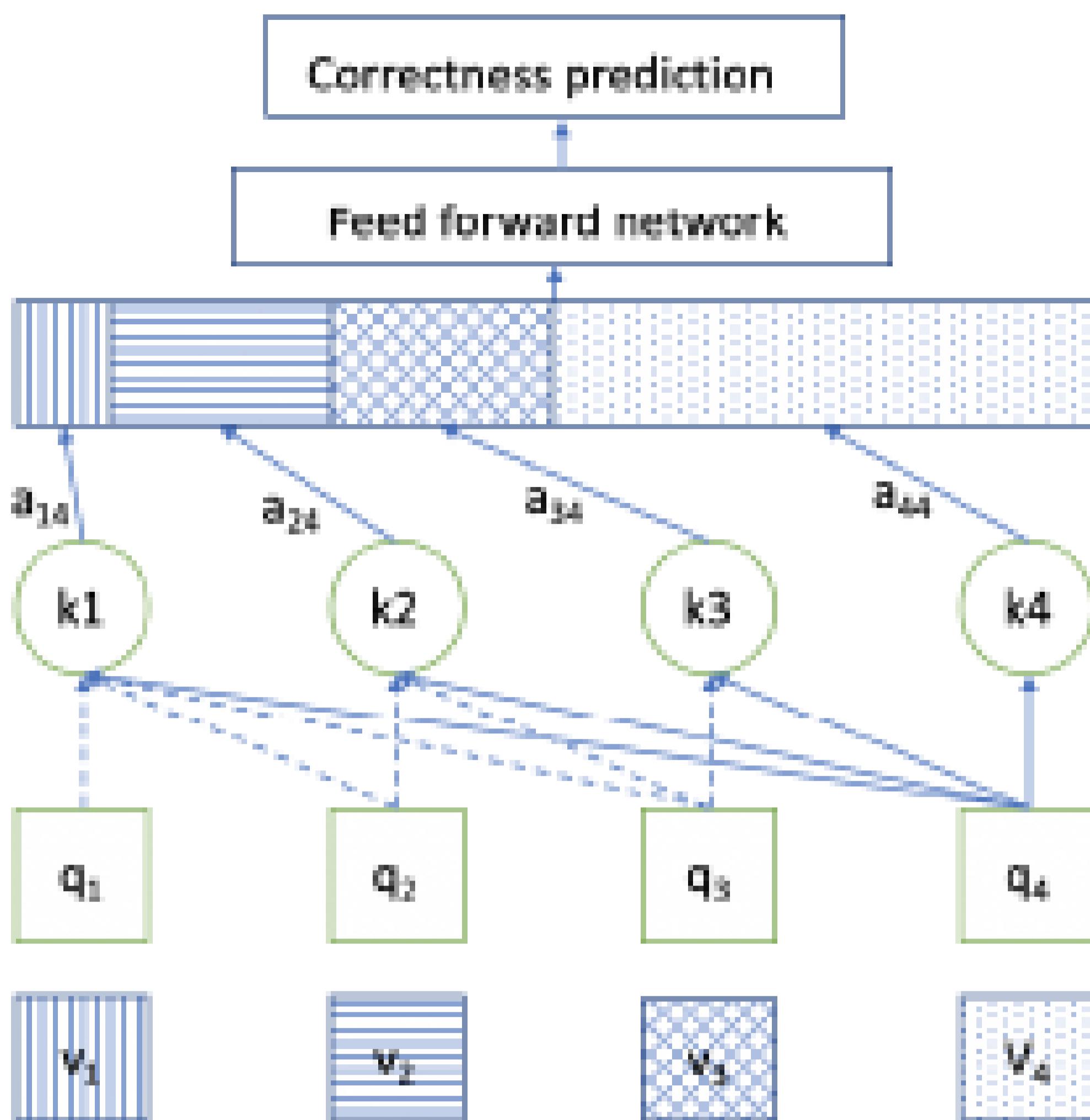
$$\mathbf{F} = \text{ReLU}(\mathbf{o} \mathbf{W}^{(1)} + \mathbf{b}^{(1)}) \mathbf{W}^{(2)} + \mathbf{b}^{(2)},$$

where \mathbf{W} are weight matrices and $\mathbf{b}^{(1)}$ and $\mathbf{b}^{(2)}$ are Bias vectors.
We add Always residual network and apply Normalization and Dropout

In the end we apply a sigmoid function to a linear layer of \mathbf{F} matrix obtained to generate probability of answering correctly:

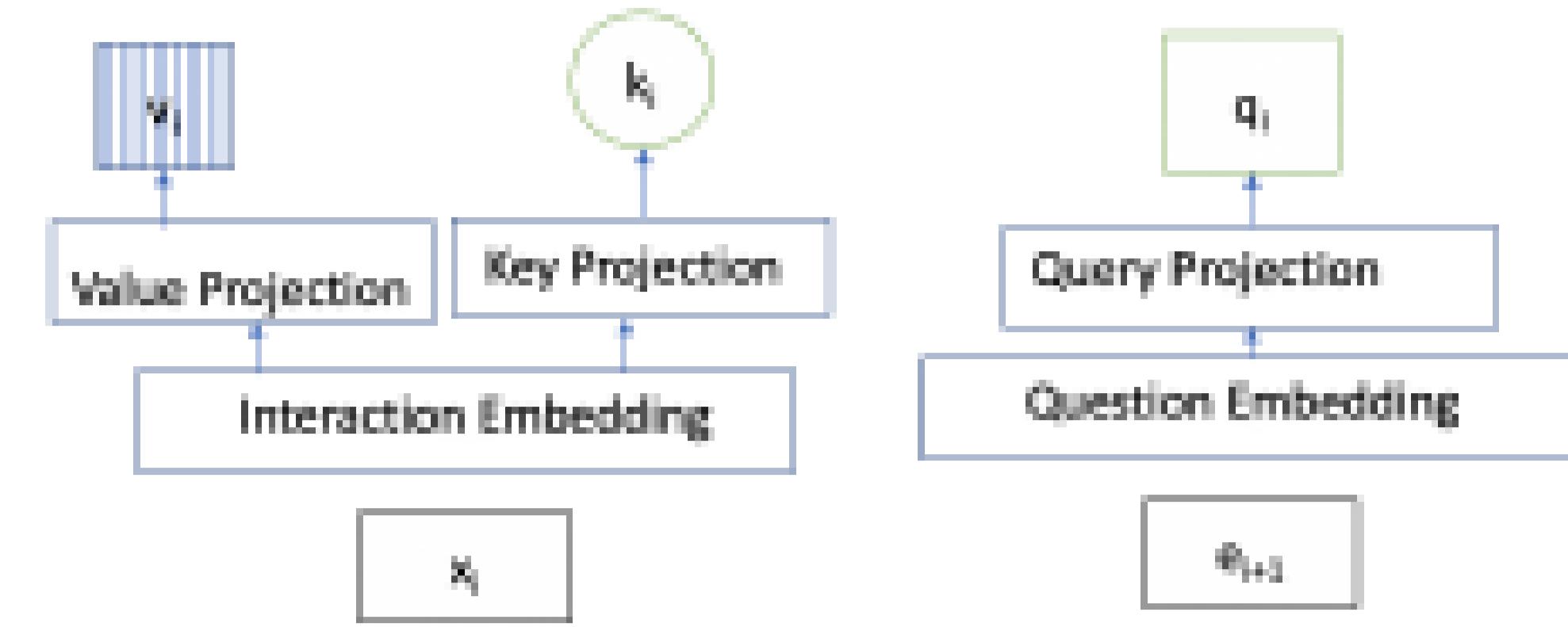
$$p = \sigma(\mathbf{F} \mathbf{W} + \mathbf{b}),$$

SAKT- Self-Attentive model for Knowledge Tracing



In the end we take outputs of Feed Forward network to calculate prediciton as:

$$p_i = \text{Sigmoid}(\mathbf{F}_i \mathbf{w} + \mathbf{b}),$$



Input X_i is composed of a couple (exercise, result) and we have a sequence of inputs X_i which constitutes the rel input of the system and the new exercise we want to predict the result E_{i+1} .

The past interactions are embedded using a trained interaction embedding matrix M of dimensions $2E \times d$, to express each exercise as a weighted combination of some terms (embedded concepts). Exercise to predict instead will be embedded with a trained exercise matrix of dimensions $E \times d$. This will lead to a matrix M and a vector E , we can use them as inputs to a SELF ATTENTIVE layer.

$$\mathbf{Q} = \hat{\mathbf{E}}\mathbf{W}^Q, \mathbf{K} = \hat{\mathbf{M}}\mathbf{W}^K, \mathbf{V} = \hat{\mathbf{M}}\mathbf{W}^V,$$

Where W_q , W_k and W_v are respectively attention matrices for query, key and value and are learnt during training.

In the end we use Q , K and V as input to calculate attention:

$$\text{Attention}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = \text{softmax}\left(\frac{\mathbf{Q}\mathbf{K}^T}{\sqrt{d}}\right)\mathbf{V}.$$

We use Multi Head Attention too, meaning that we calculate different W_q , W_k , W_v to allow detection of different templates.

In the end we take the outputs from different attention heads and use as inputs for a Feed Forward Neural Network to allow non linearities and relations between different heads.

$$\mathbf{F} = \text{FFN}(\mathbf{S}) = \text{ReLU}(\mathbf{S}\mathbf{W}^{(1)} + \mathbf{b}^{(1)})\mathbf{W}^{(2)} + \mathbf{b}^{(2)},$$

We use Normalization and add Residual Connections too.

KT-XL: A Knowledge Tracing Model for Predicting Learning Performance Based on Transformer-XL

Prediction Layer

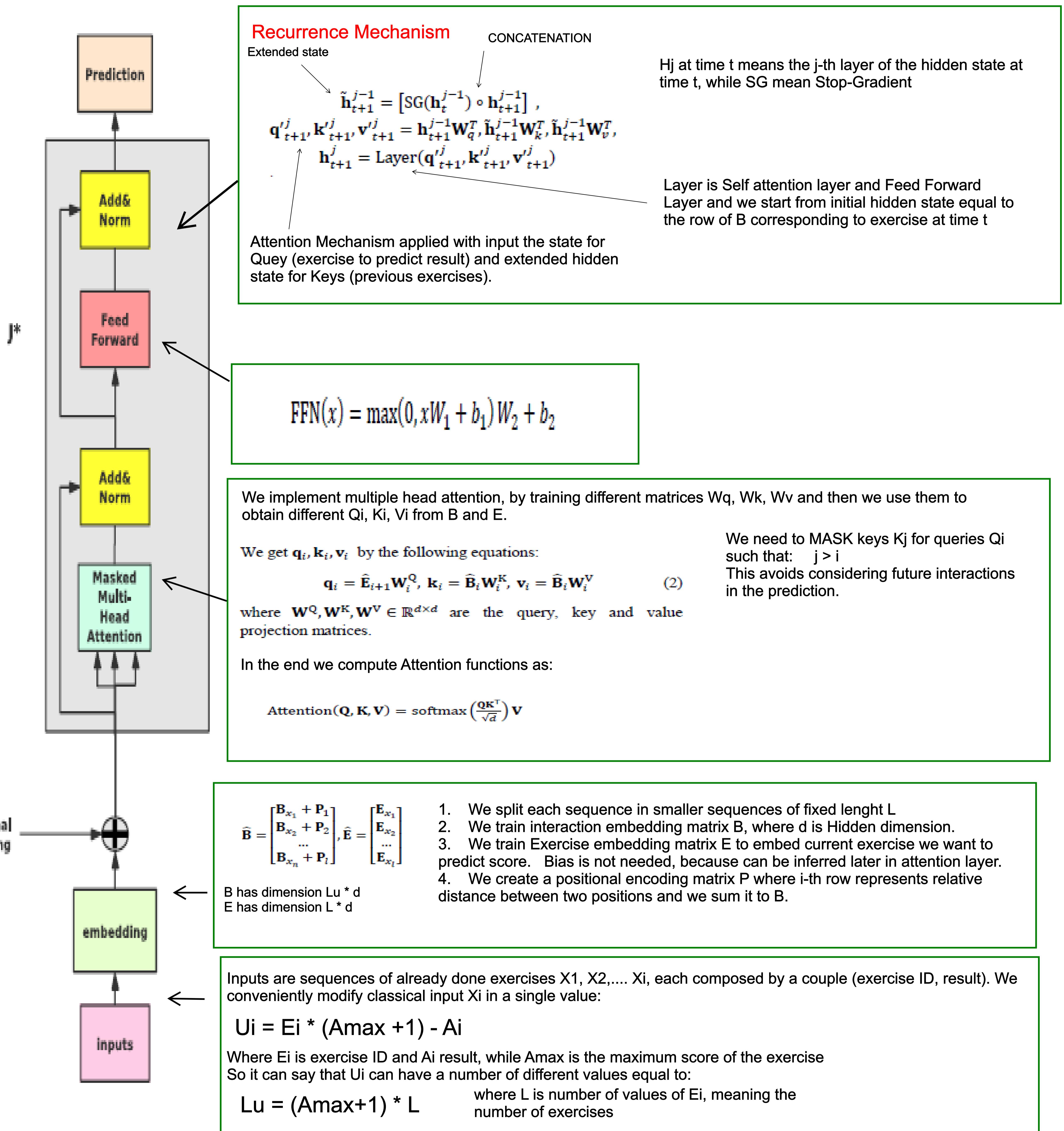
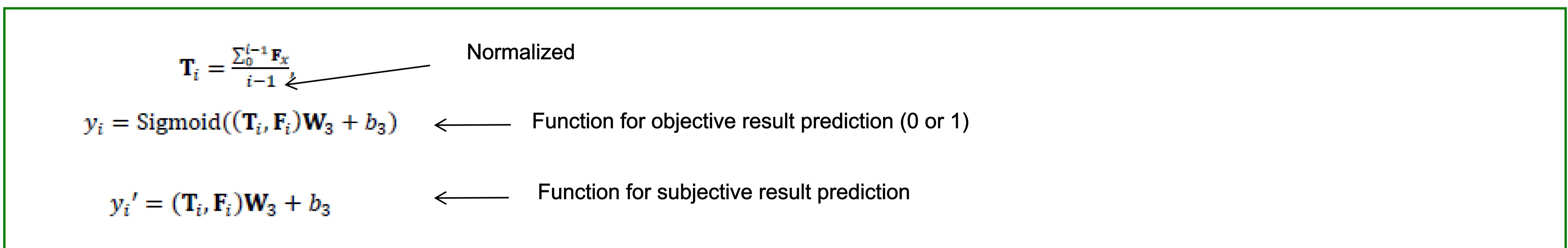


Figure 1: The architecture of KT-XL

