

BERT:

Pre-training of Deep Bidirectional Transformers for Language Understanding

Vasiliki Kougia

Instructor: John Pavlopoulos

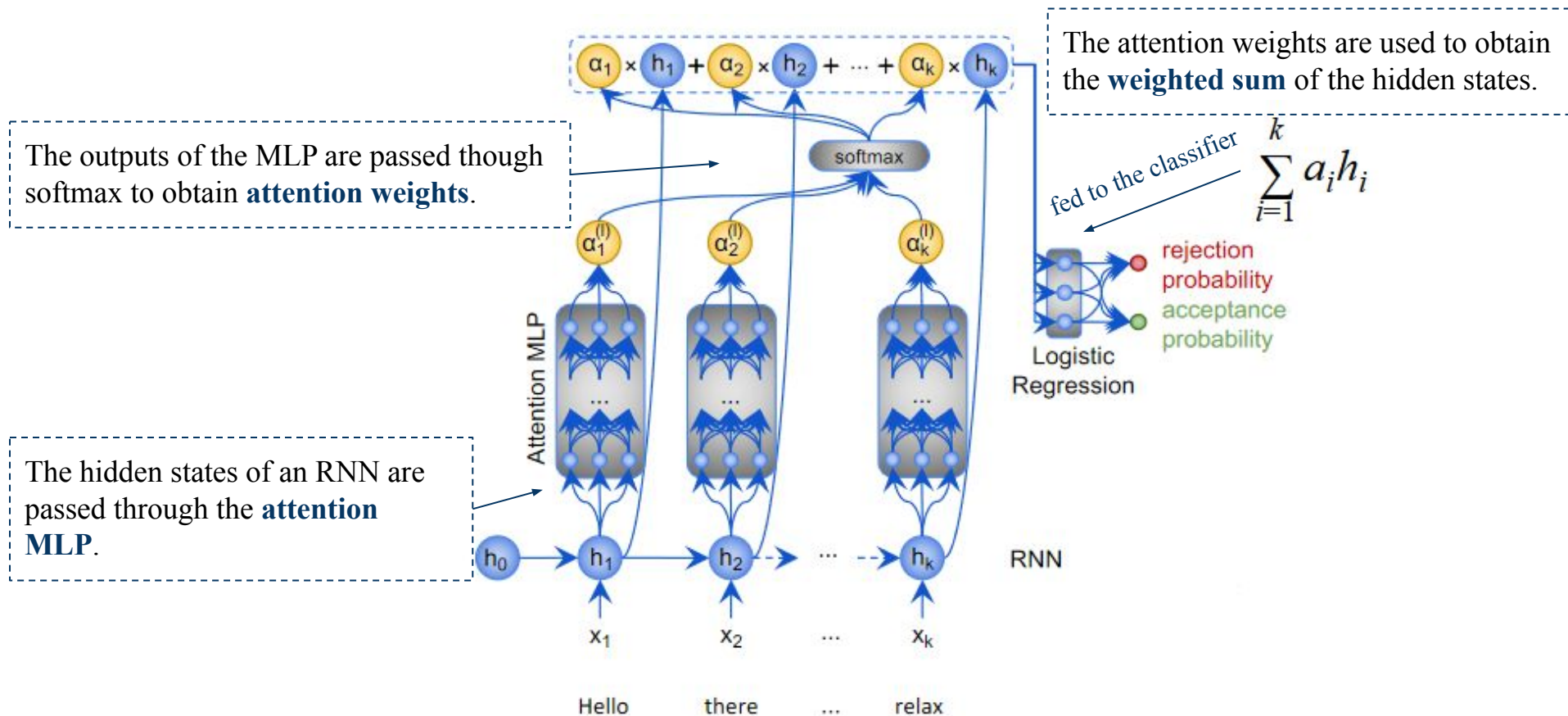
28/4/2021



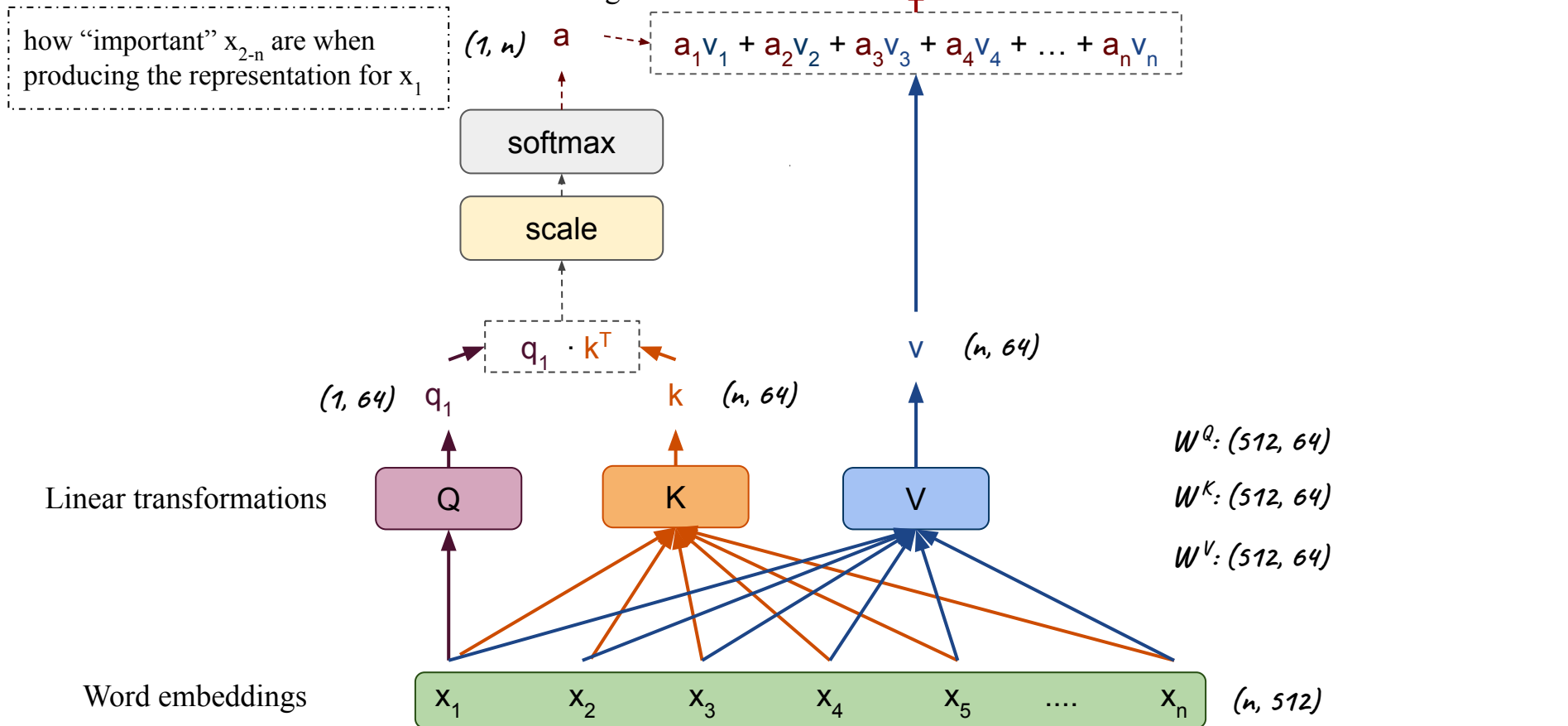
Modern NLP reading group



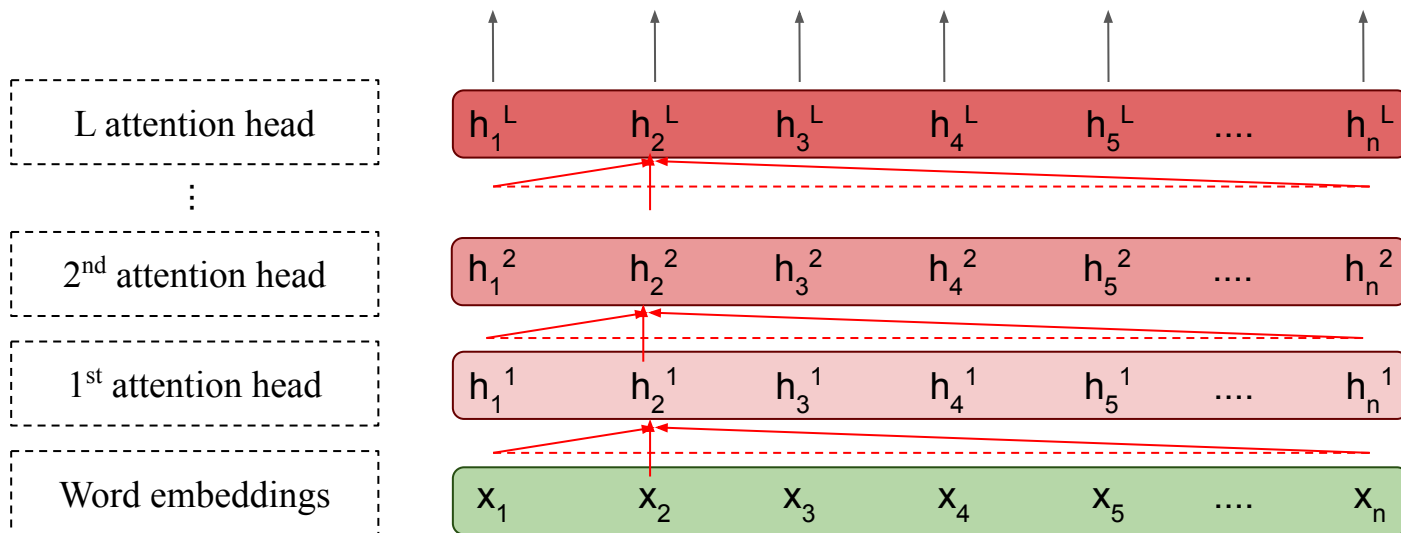
Self-attention in RNNs



Query-key-value self-attention



Transformers - the bigger picture



For position 2

$$h_i^j = \sum_{r=1}^n \text{softmax}(q_i^j k_r^{jT}) v_r^j$$

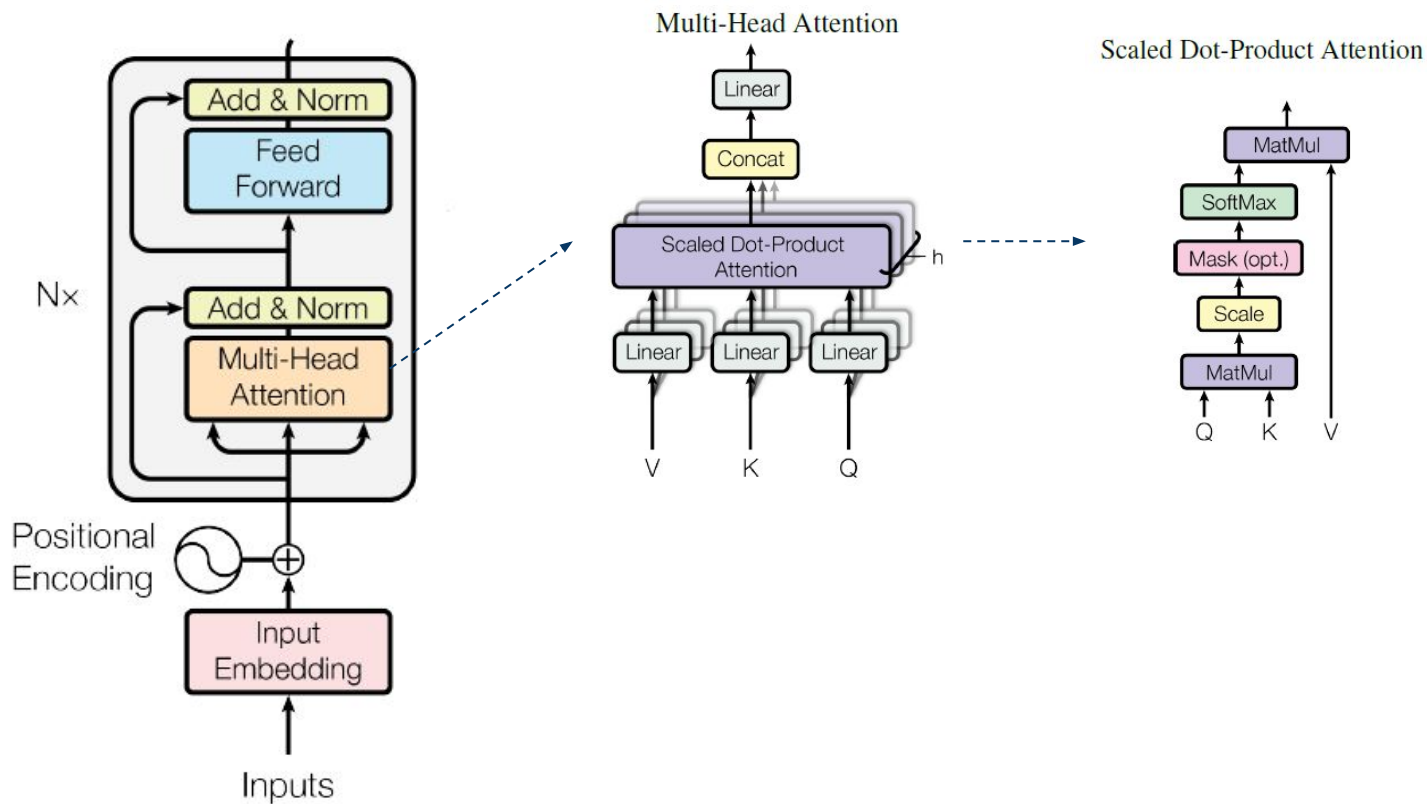
$$q_i^j = W^{Q,j} h_i^{j-1}$$

$$k_r^j = W^{K,j} h_r^{j-1}$$

$$v_r^j = W^{V,j} h_r^{j-1}$$

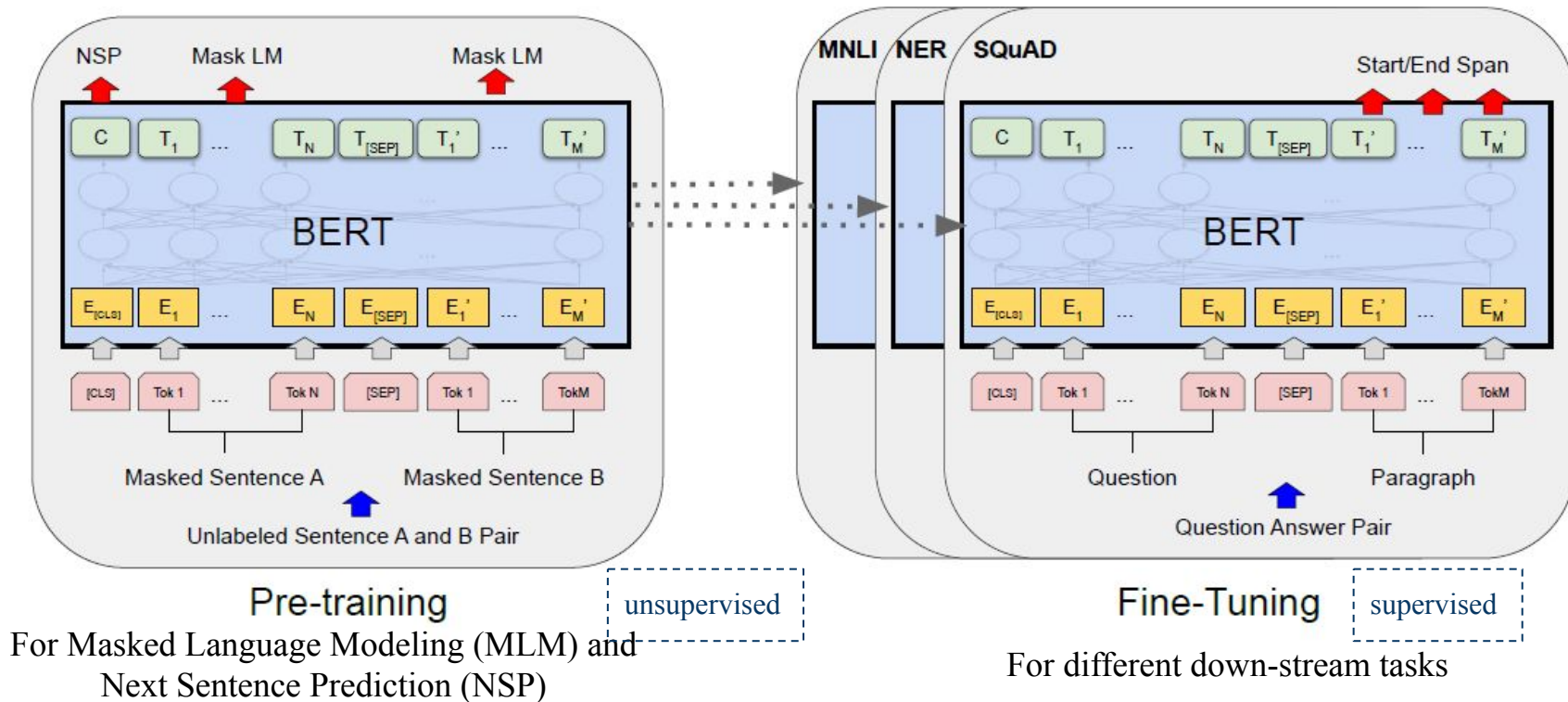
Different weight matrices at each head

Transformer encoder block



BERT

- Consists of stacked transformer encoders
- Pre-trained on a huge corpus



BERT - text classification

