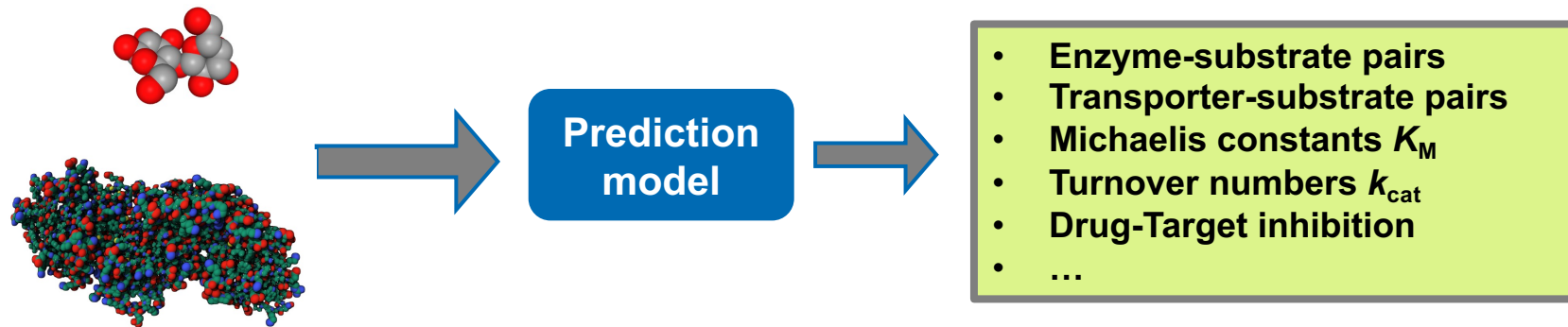
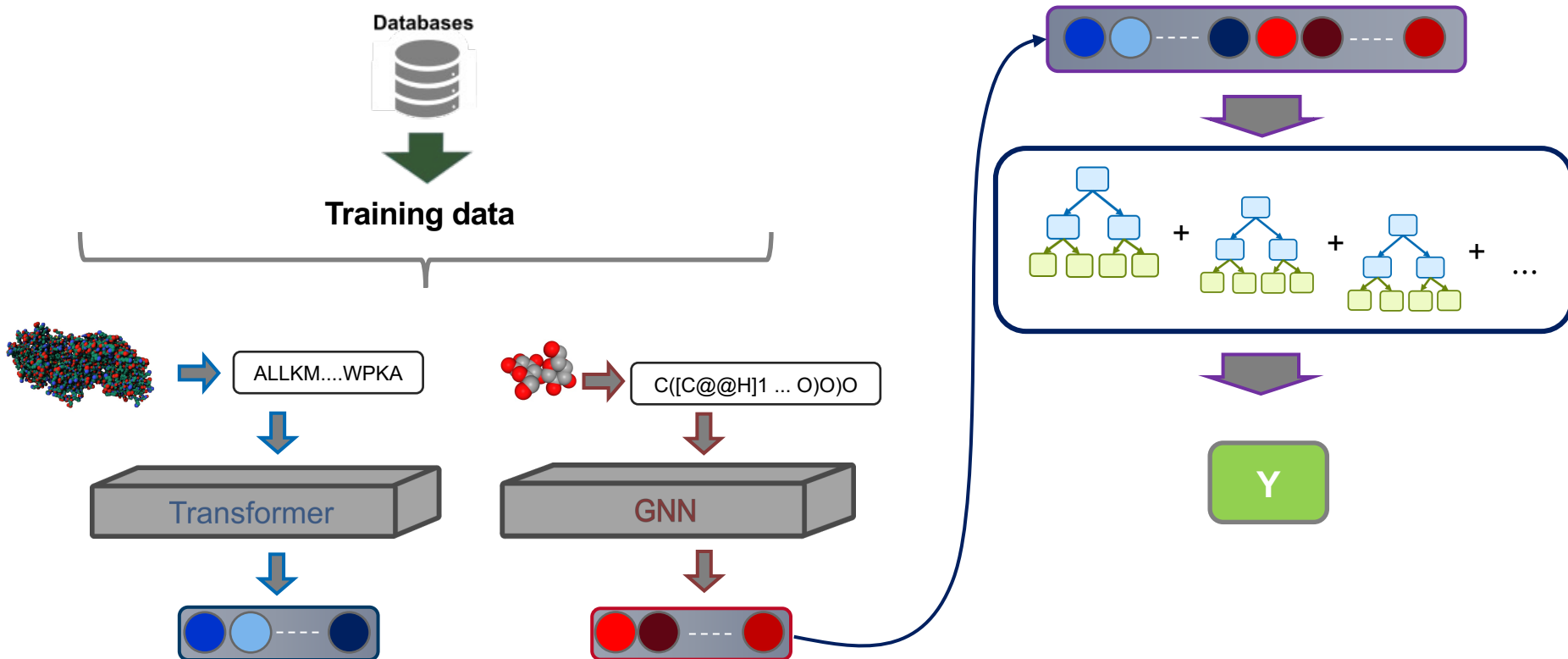


Predicting protein-small molecule interactions

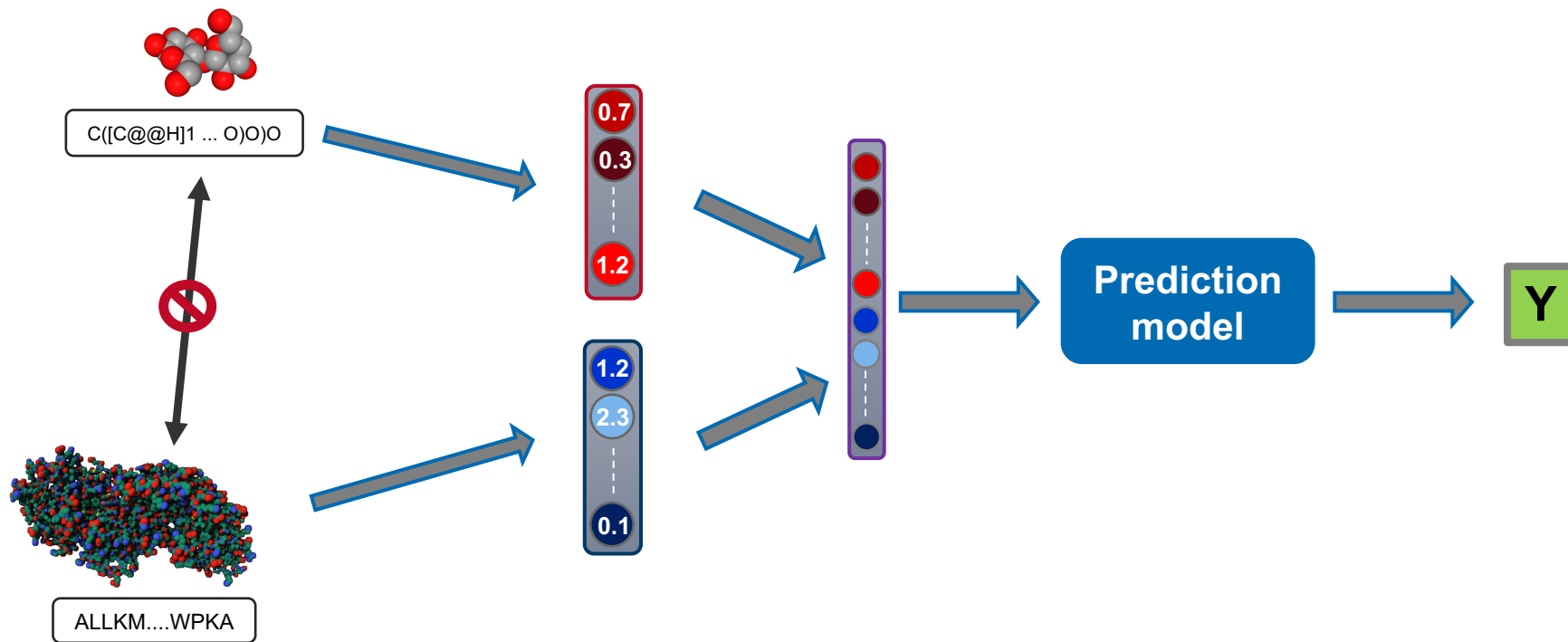
Predicting protein-small molecule interactions



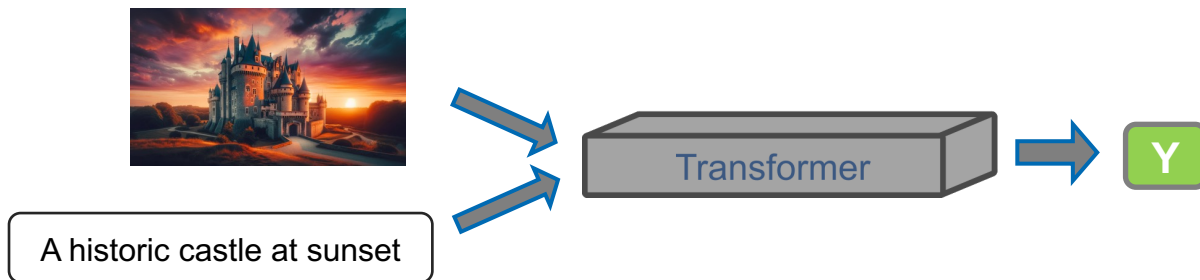
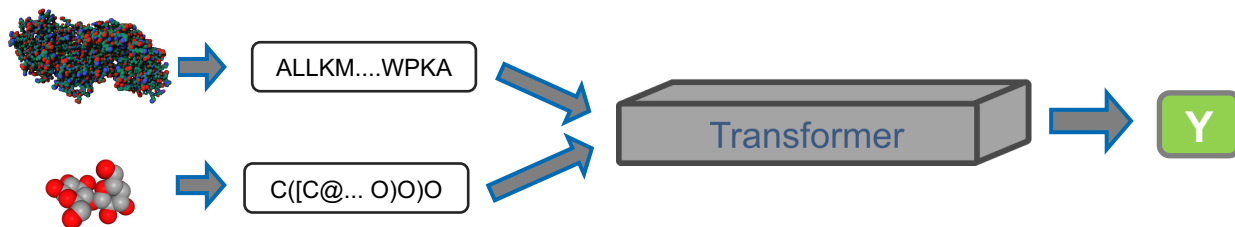
Typical Prediction Pipeline



Missing Information Exchange



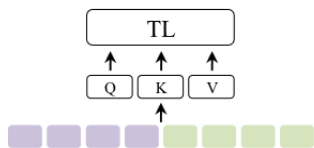
Multimodal Transformer Networks (1)



Multimodal Transformer Networks (2)

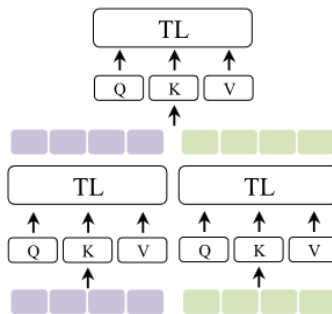
■ ■ ■ ■ = 1st input type/modality (e.g., text or protein)

■ ■ ■ ■ = 2nd input type/modality (e.g., image or small molecule)



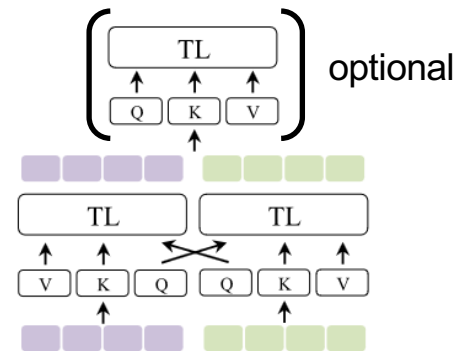
Concatenating both initial input sequences

- Generating a joint vocabulary for tokens of both input types
- Concatenated sequences is input for single Transformer



Concatenating sequence embeddings after separate processing

- Each input sequence is input of separate Transformer layers
- Resulting embeddings are used as the input of shared Transformer layers



Cross-attention between both sequences

- Using embeddings from other input type to calculate attention scores
- No self-attention

ProSmith – Protein-Small Molecule Transformer

