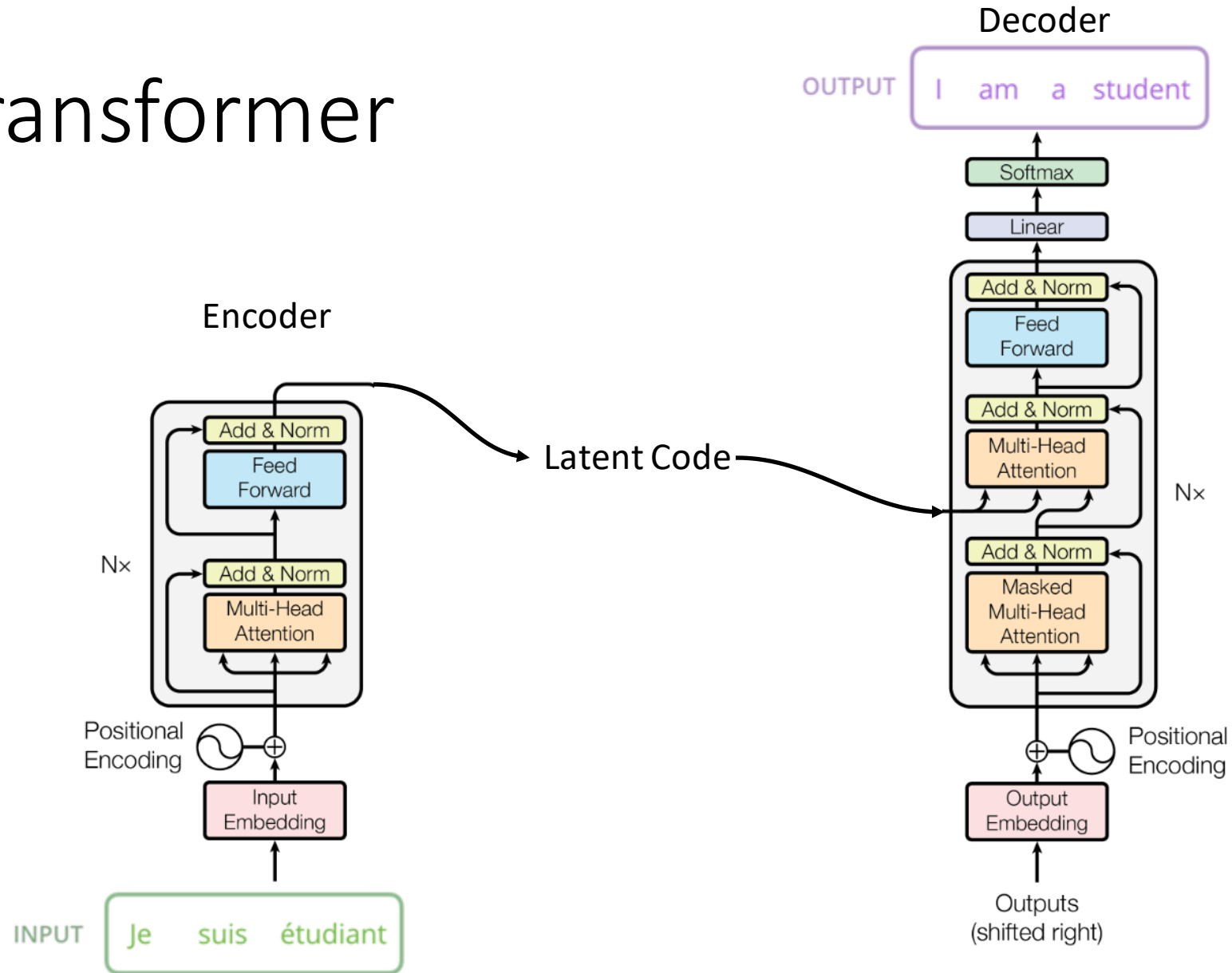


Neural Computation

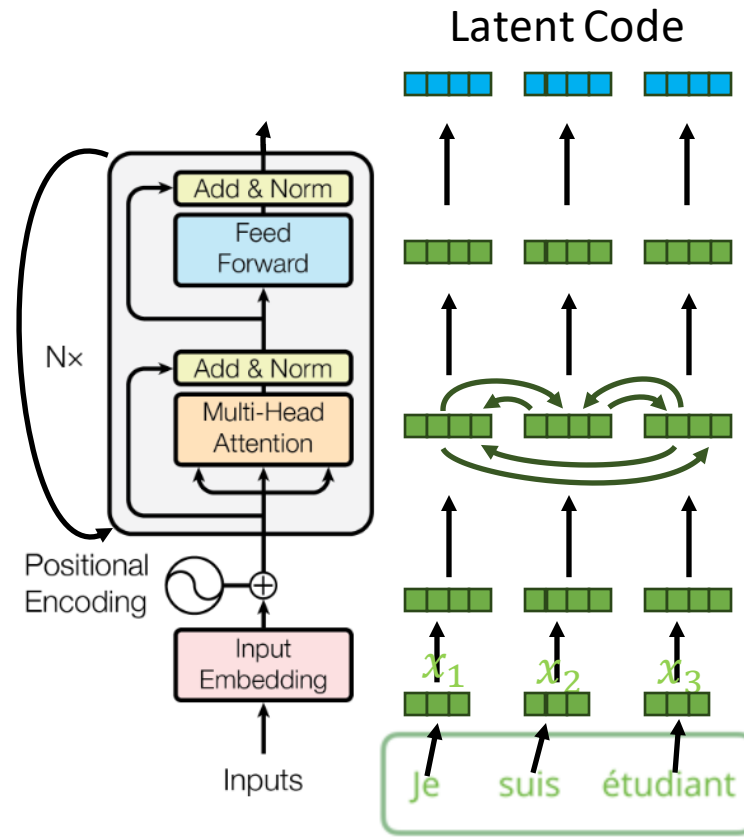
The Decoder - Part 2

Putting Everything Together

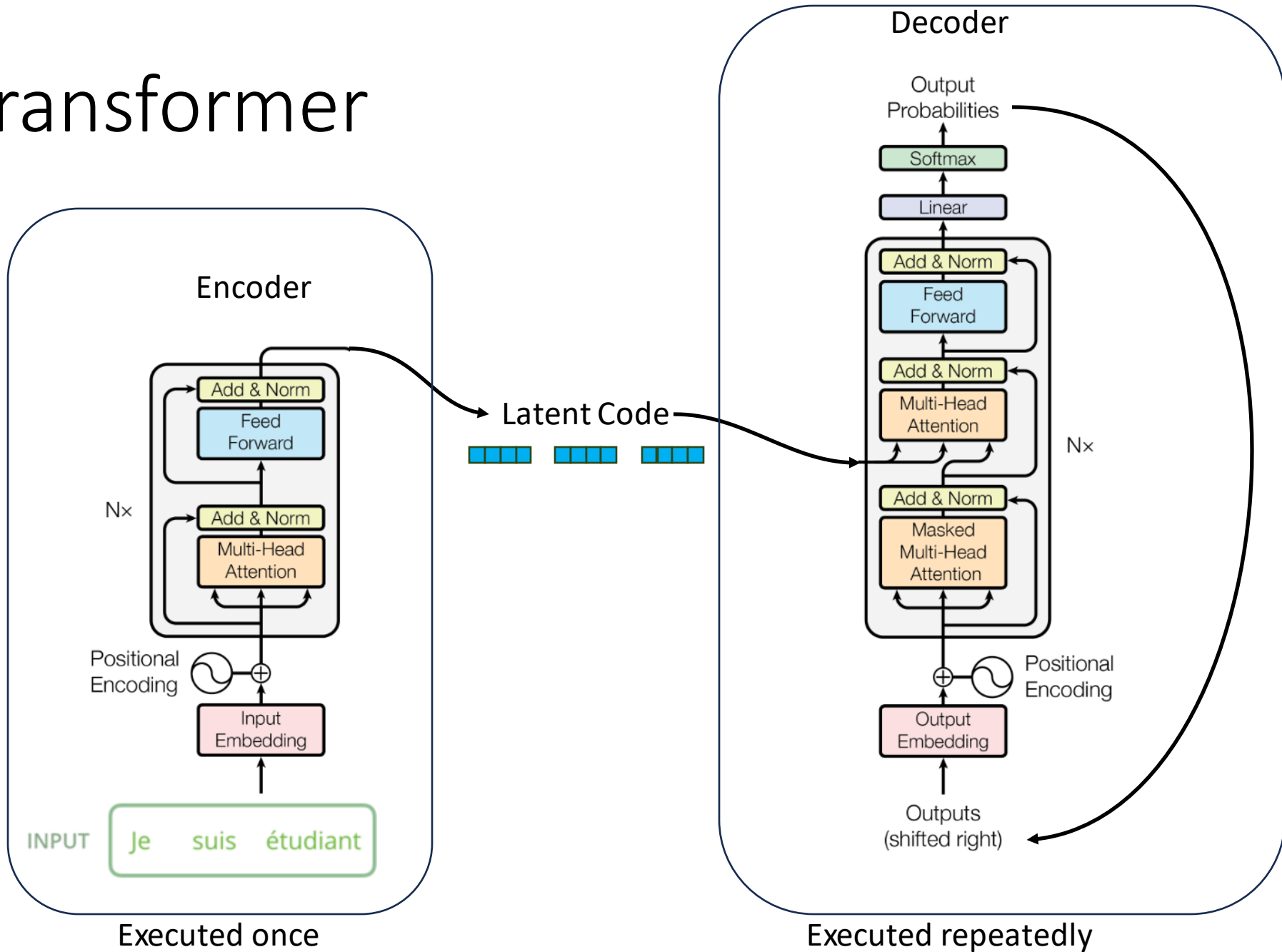
The Transformer



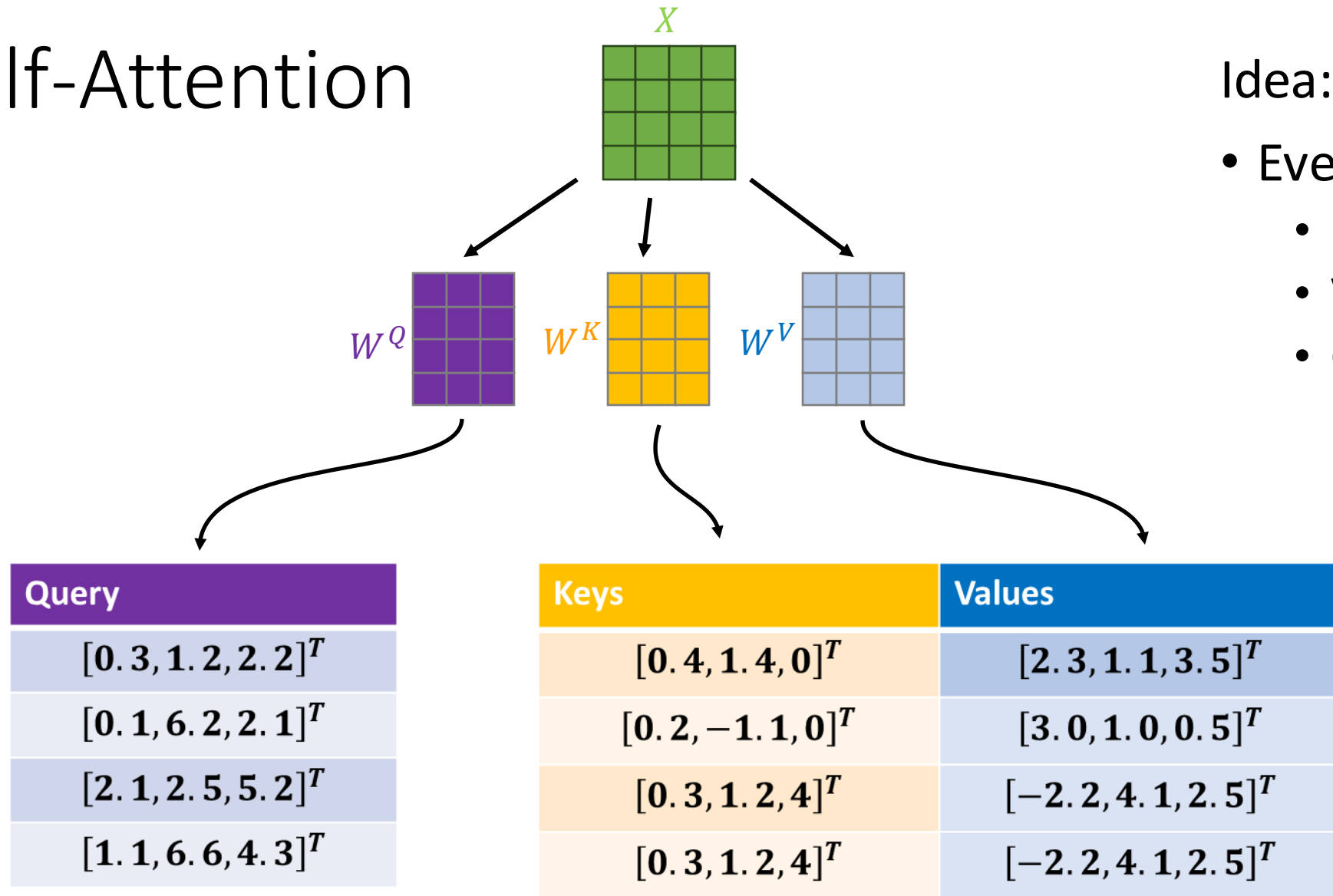
The Encoder



The Transformer



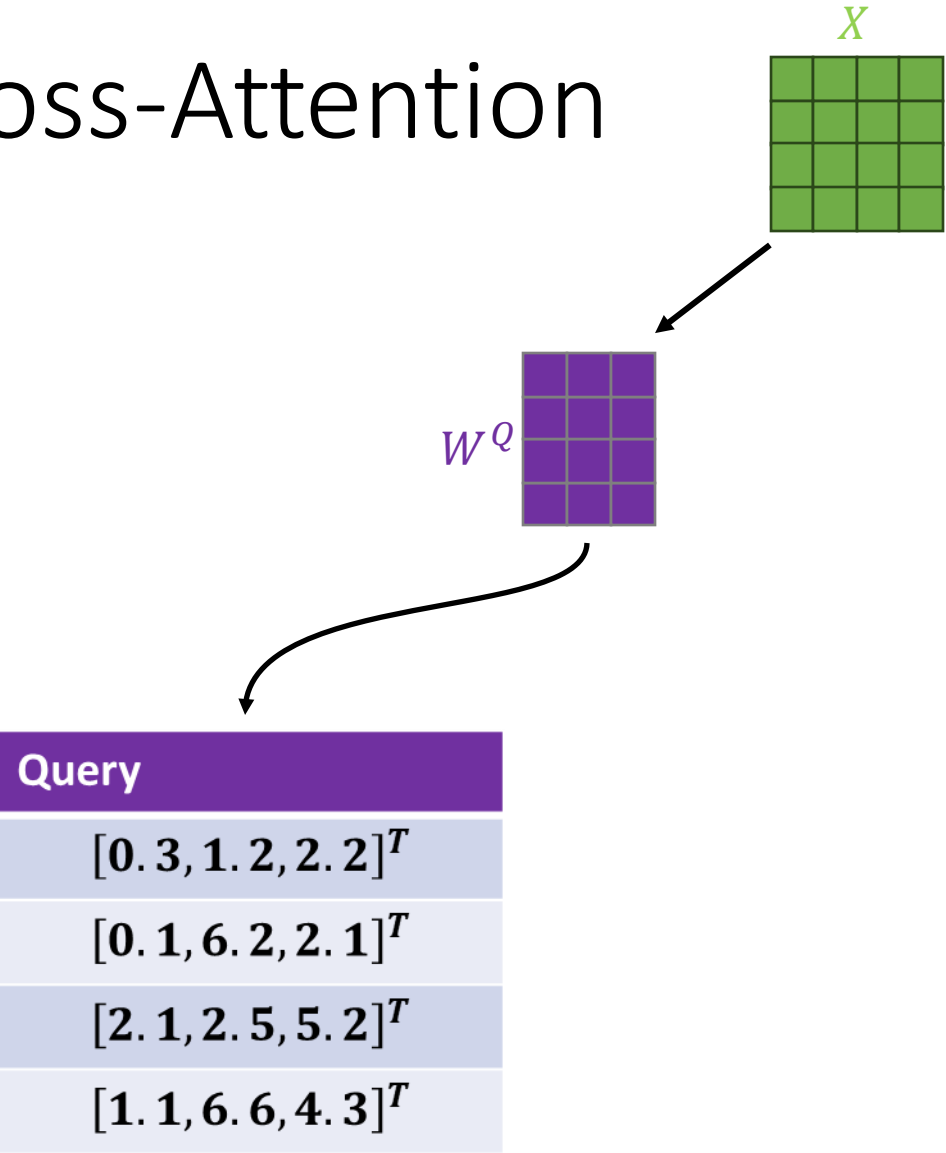
Self-Attention



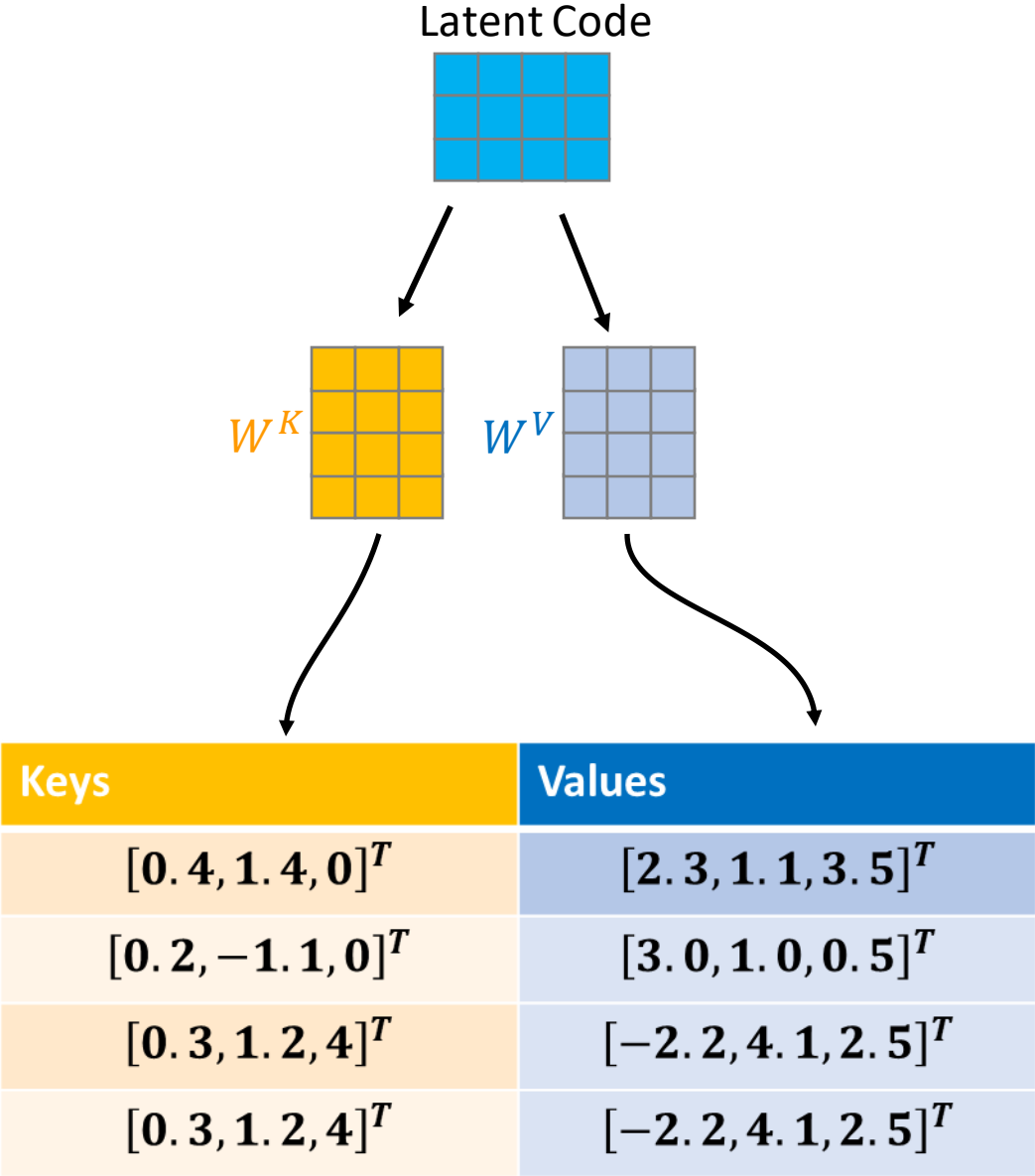
Idea:

- Every token makes:
 - Key
 - Value
 - Query

Cross-Attention

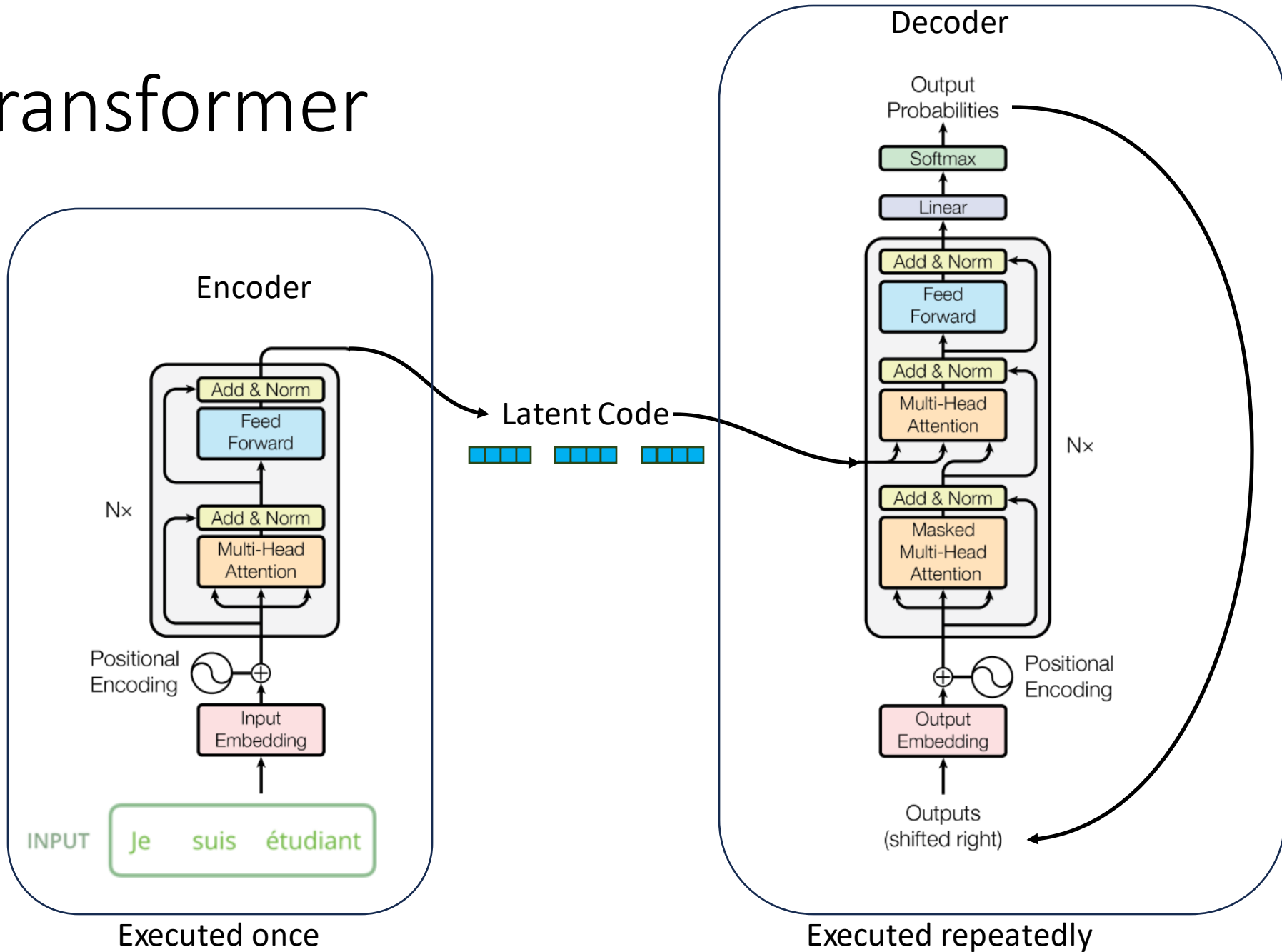


From Decoder



From Encoder

The Transformer



Which word in our vocabulary
is associated with this index?

Get the index of the cell
with the highest value
(**argmax**)

log_probs



am

5

Softmax

logits



Linear

Decoder stack output



Which word in our vocabulary
is associated with this index?

am

Get the index of the cell
with the highest value
(argmax)

5

log_probs



↑



↑

logits



↑



↑

Decoder stack output



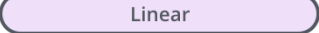
↑



↑



↑



↑



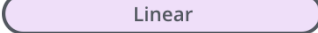
↑



↑



↑



↑



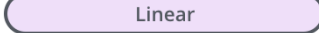
↑



↑



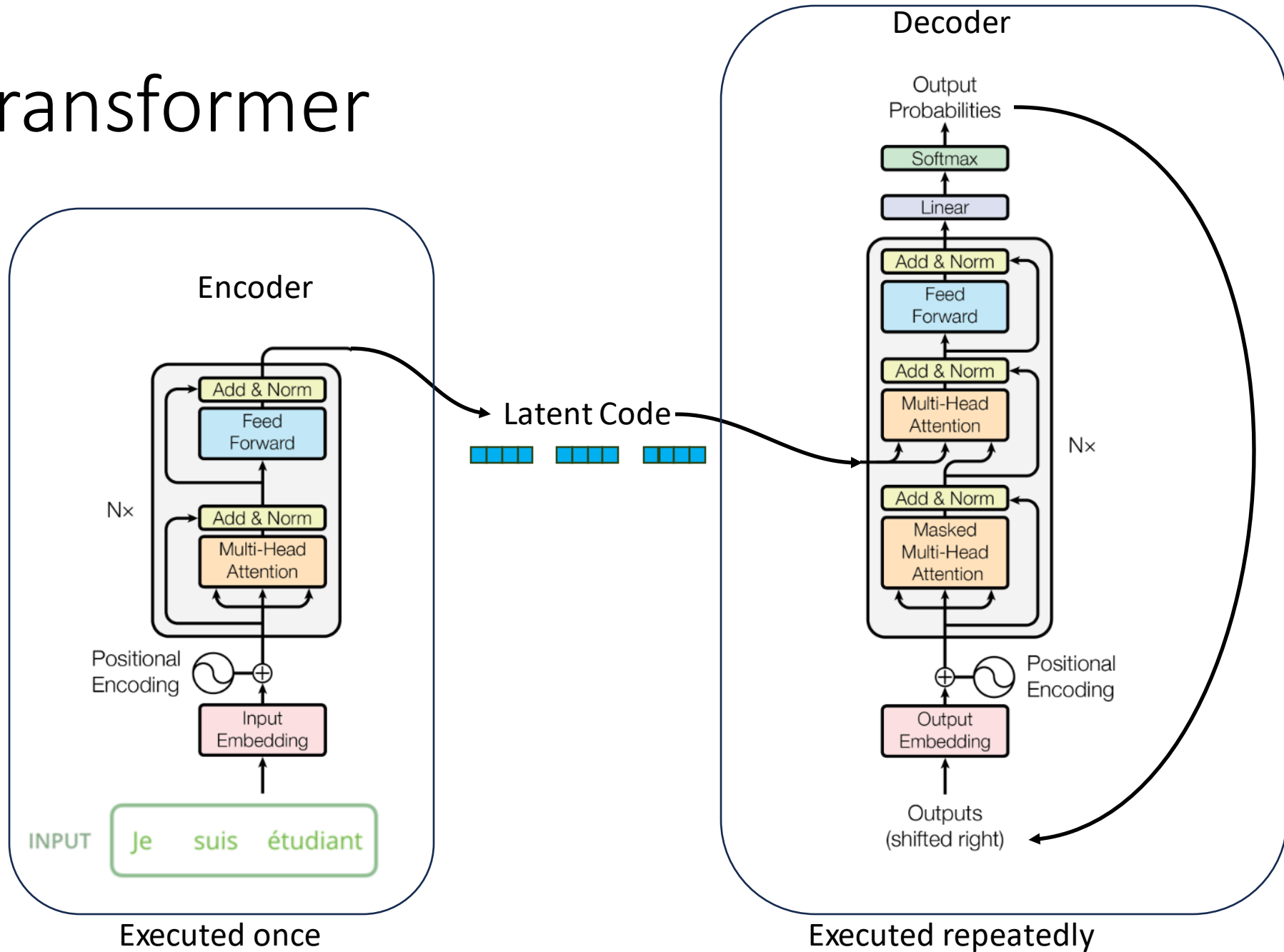
↑



↑

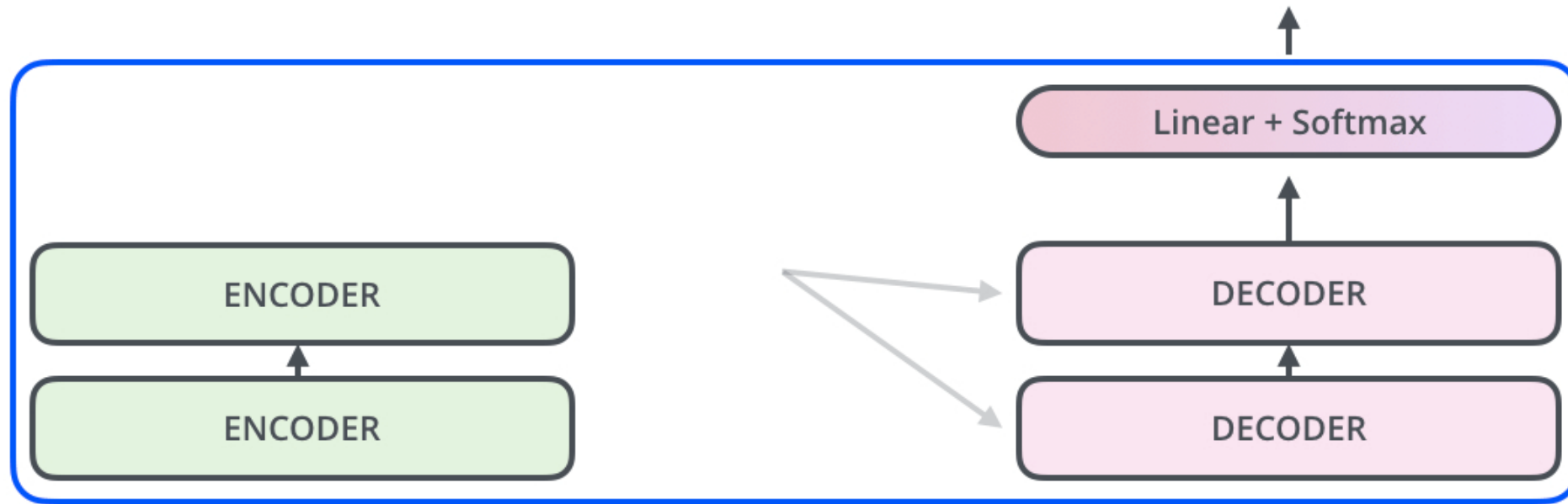


The Transformer

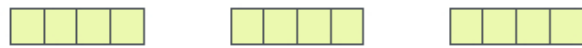


Decoding time step: 1 2 3 4 5 6

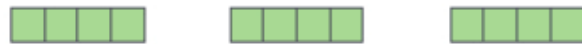
OUTPUT



EMBEDDING
WITH TIME
SIGNAL



EMBEDDINGS

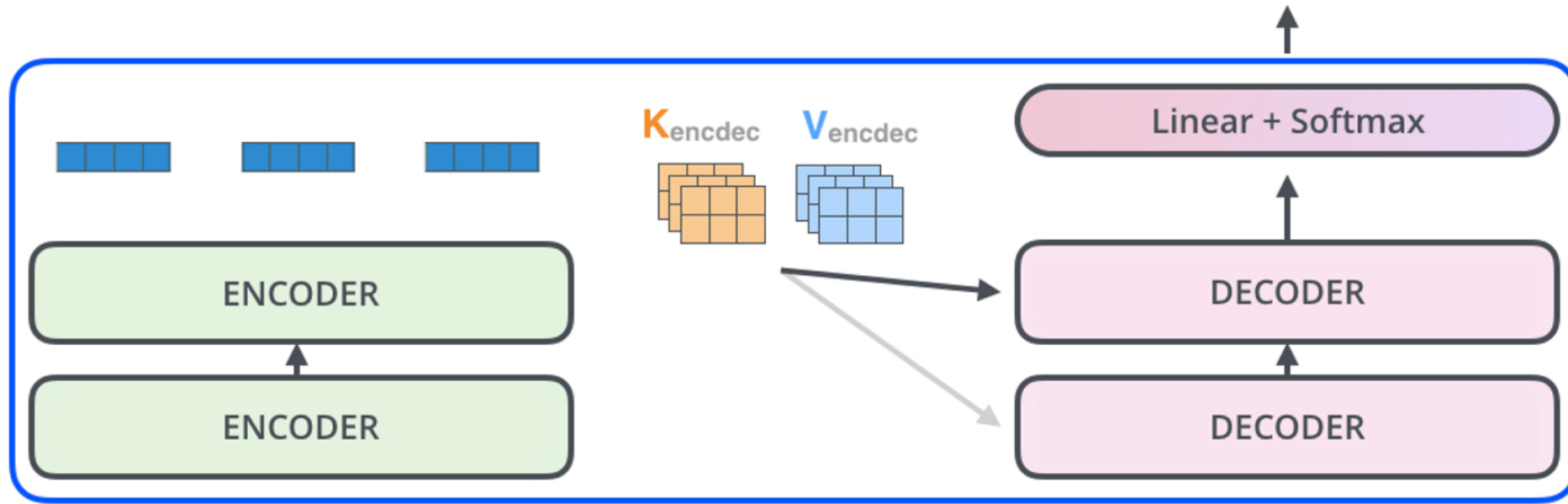


INPUT

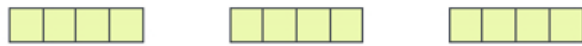
Je suis étudiant

Decoding time step: 1 2 3 4 5 6

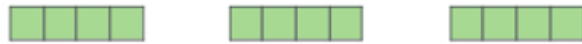
OUTPUT |



EMBEDDING
WITH TIME
SIGNAL



EMBEDDINGS



INPUT Je suis étudiant

Decoding time step: 1 2 3 4 5 6

OUTPUT |

