# Transformers as pure LMs: improving the past context



**Tr LMs**

**RNN LMs**

**FF LMs** (4-gram)

Token embeddings

| [/s] | il | ne | faut | pas | dire | que | le | roi | est | un | fou |

Position embeddings

+ + + + + + + + + + + +

| 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 |

Also: relax causality, recompute past representations after each new word

# Faster, Better Encoder-Decoder + Attention : Transformer



In the decoder: cross-attention with respect to the last encoder layer