(a) PerceiverIO

(b) Transolver

($B$ blocks)

(c) LNO

($B$ blocks)

(d) FLARE

($B$ blocks)

Token sequences

Input/output ($N$ tokens)  Latent ($M$ tokens)

Attention Operators

Self attention ($M \rightarrow M$)  Cross attention ($N \rightarrow M$)

Encoder ($N \rightarrow M$)/ decoder ($M \rightarrow N$)

Feature projector  Cross attn projector