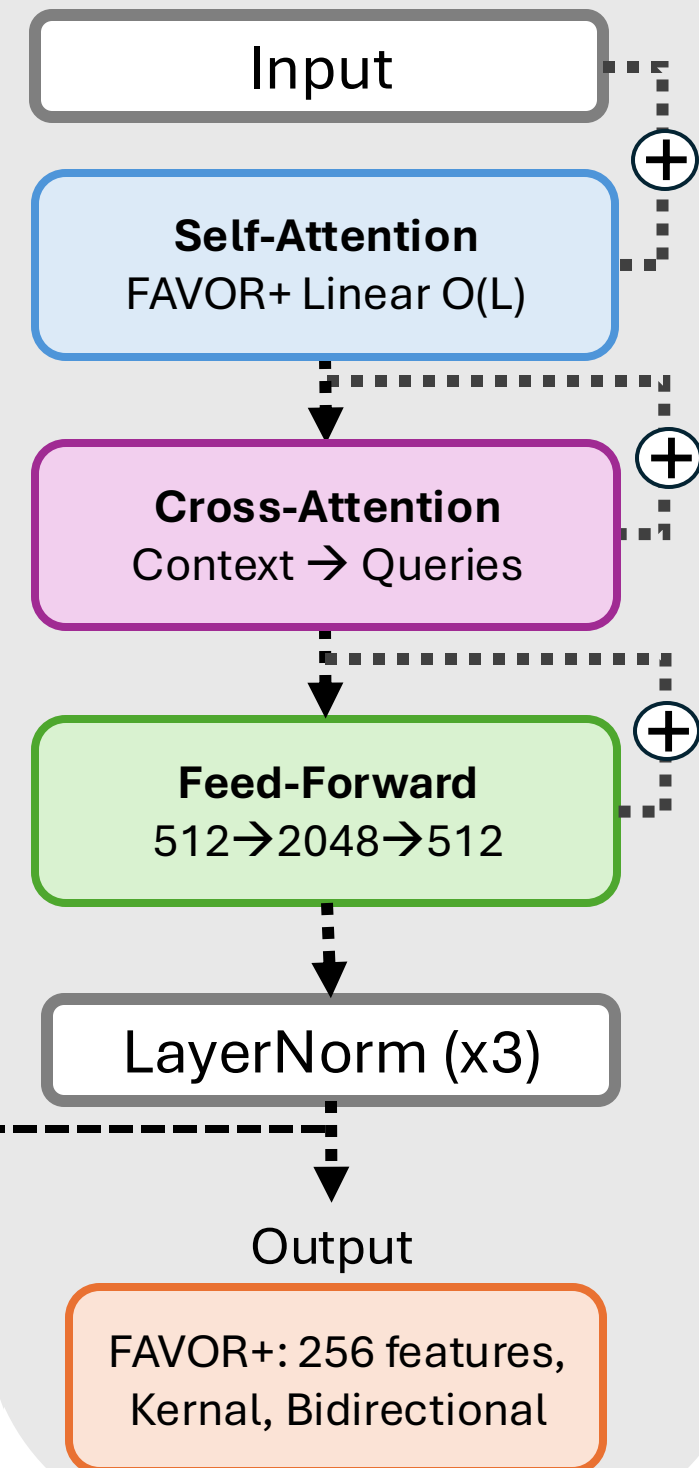


### Performer Block (N=6)



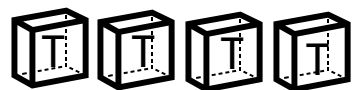
## Temporal Multi-Sensory Encoder

**Context Tokens**  
( $B, T_{obs}, 1024$ )

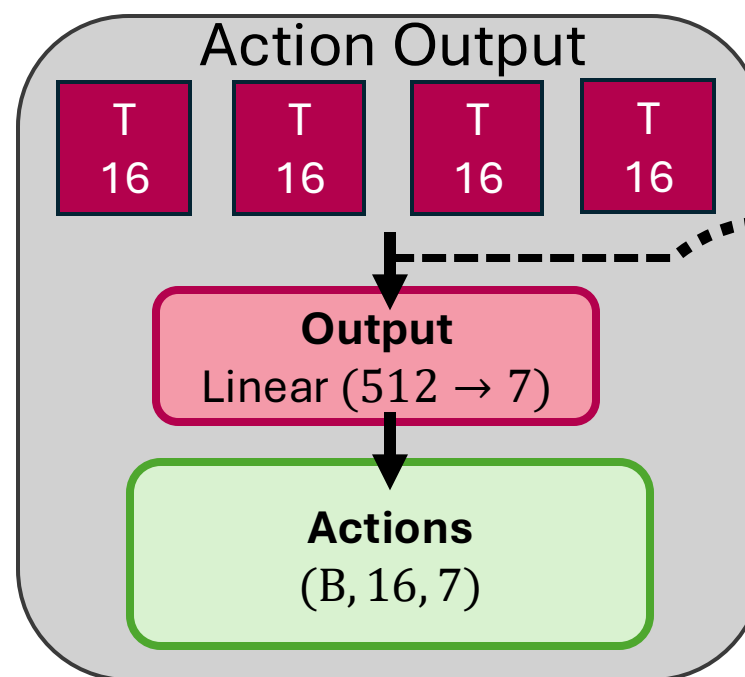
**Query Tokens**  
( $16, 512$ ) Learned  
 $T_{pred} = 16$

**Noise**  
( $B, 64$ )  $\rightarrow$  ( $B, 16, 512$ )  
Projection + Expand

**Pos Encoding**  
( $16, 512$ ) Sinusoidal  
Temporal ordering



## Performer Transformer



**RS-IMLE Loss**  
Batch-global  $\epsilon$ -rejection  
Robust distance EMA  $\epsilon$   
calibration

**Top K Loss**  
Diversity  
Weight: 0.02  
Huber loss

**Top Loss**  
 $L_{hard} + \lambda L_{soft}$   
 $\lambda = 0.02$   
End-to-end

### Training Objectives