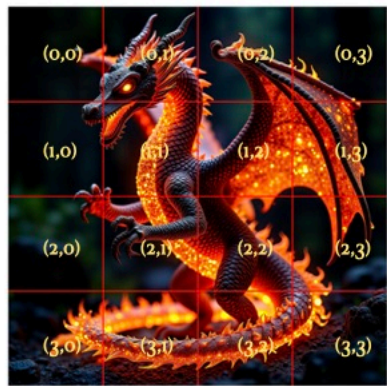


High-Res  $X$   
Partitioned to Local Windows

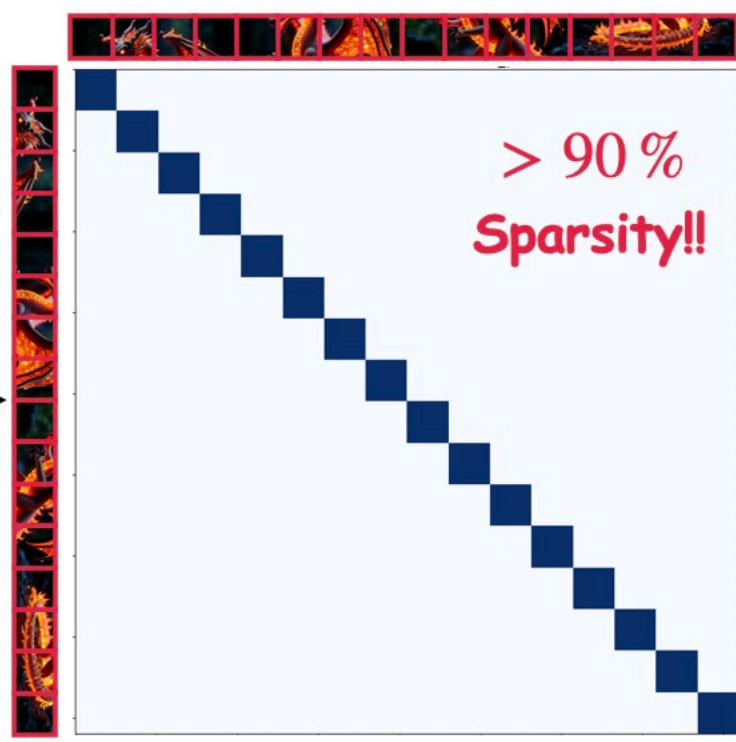


Low-Res  $X_{lr}$

Hilbert-Order  
Permutation

Guidance Via  
Attention

Position Scaling



Local Self-Attention with  
Contiguous Mask Layout

