

Visual scale S_v

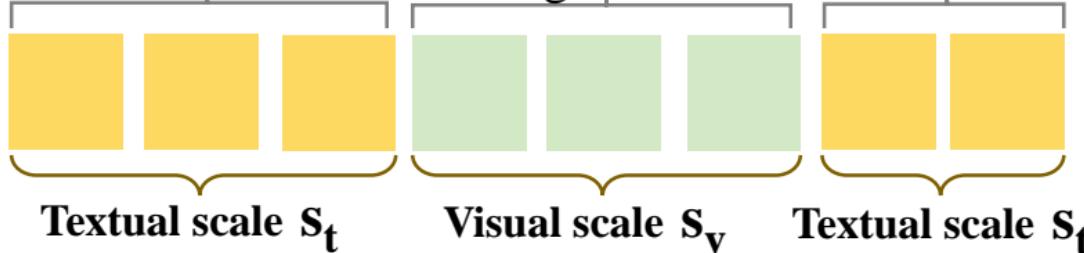
Textual scale S_t



Unified Modality-decoupled Token Representation

Attention-Invariant Flexible Switch (AIFS)

Mixed Multi-modal Token Representation
text1 image1 text2



Modality-specific Static Quantization (MSQ)

Unified token index

unified casual mask

naive casual mask

