# Multimodal VQA Architecture

*CLIP (Frozen) + Trainable Projection + Question-Type-Specific Heads*