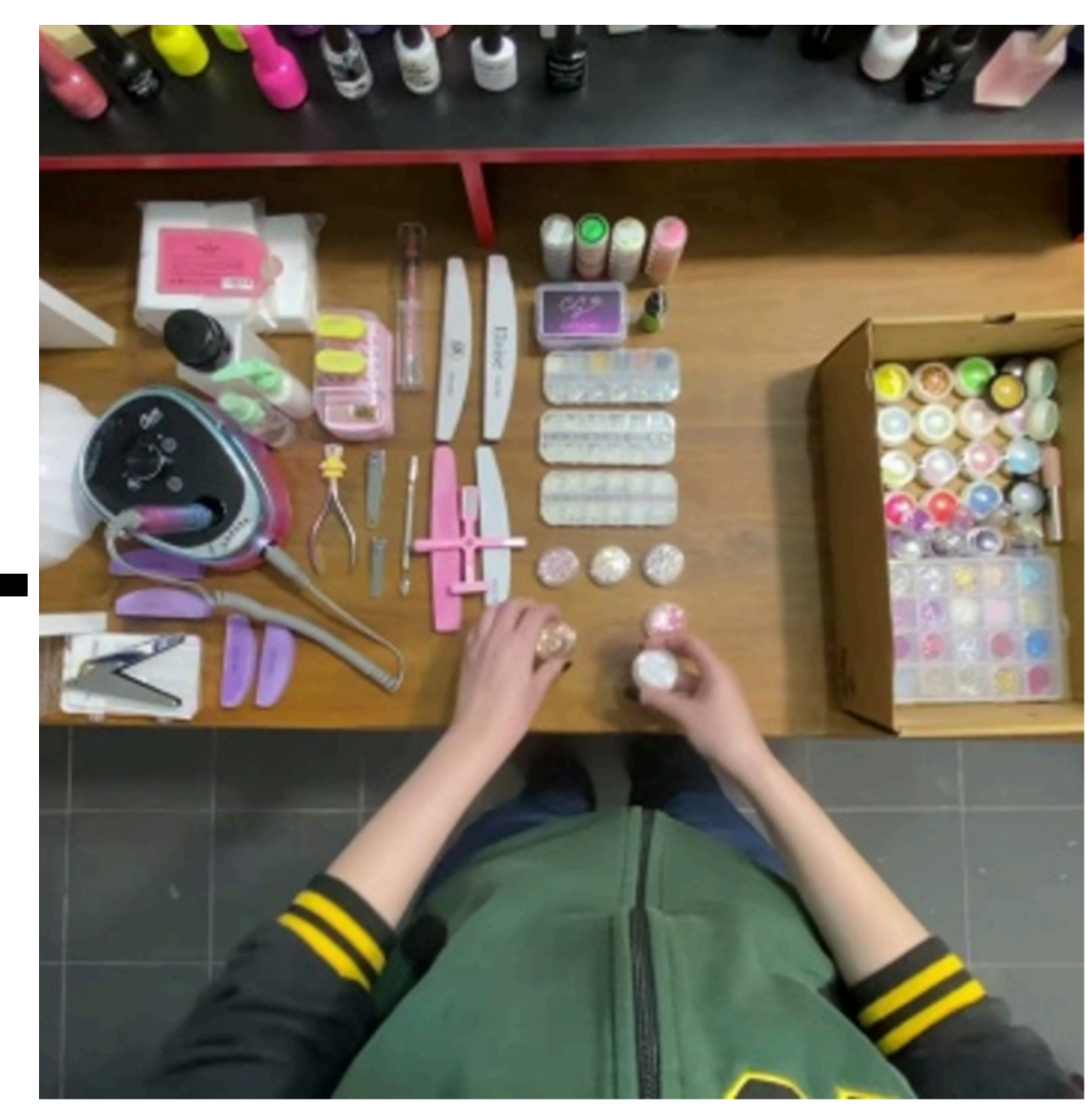


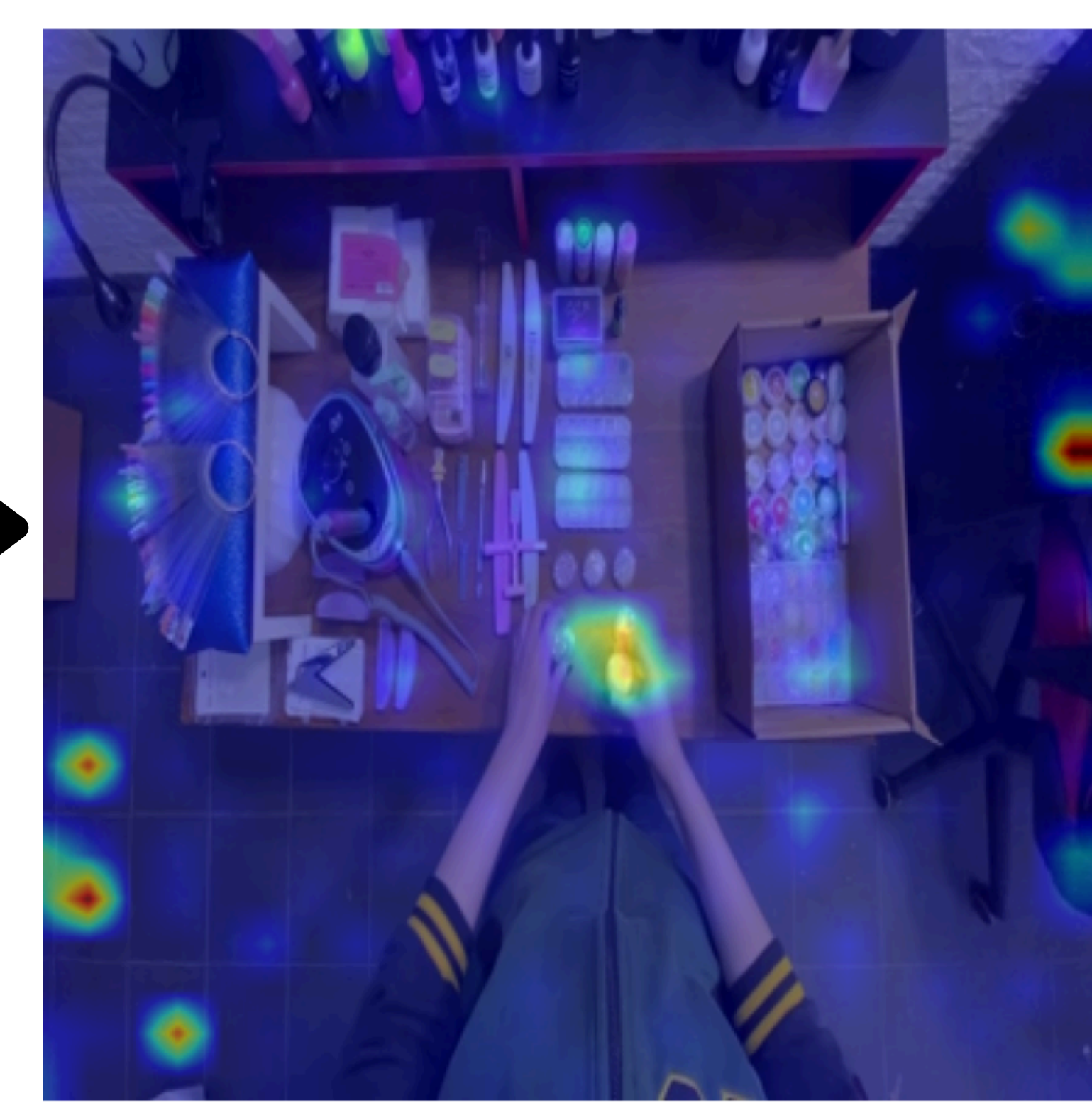
Action-Centric Features

“Organize supplies...”



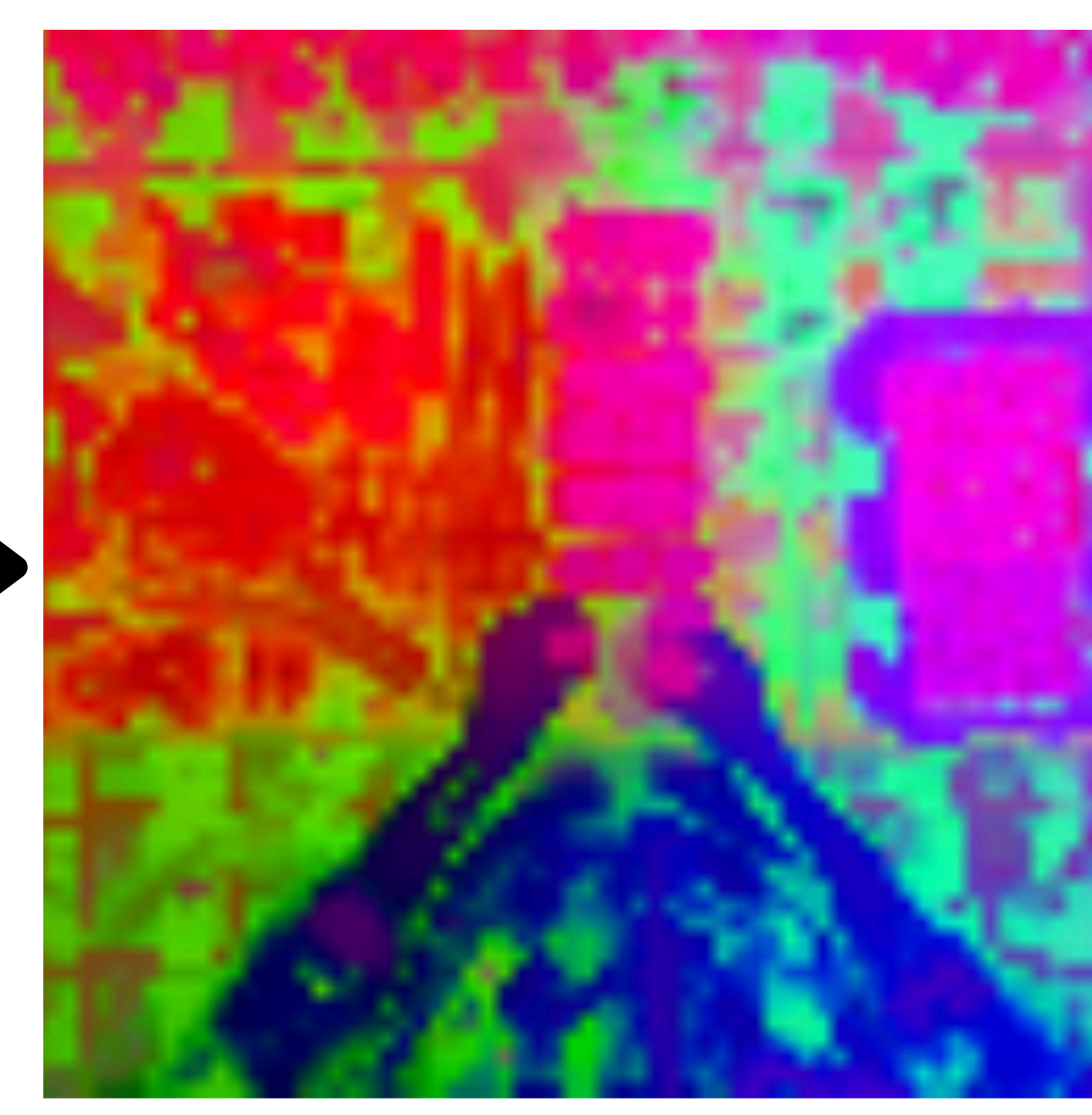
SigLIP

Semantic Meaning ✓
Action-Centric ✗



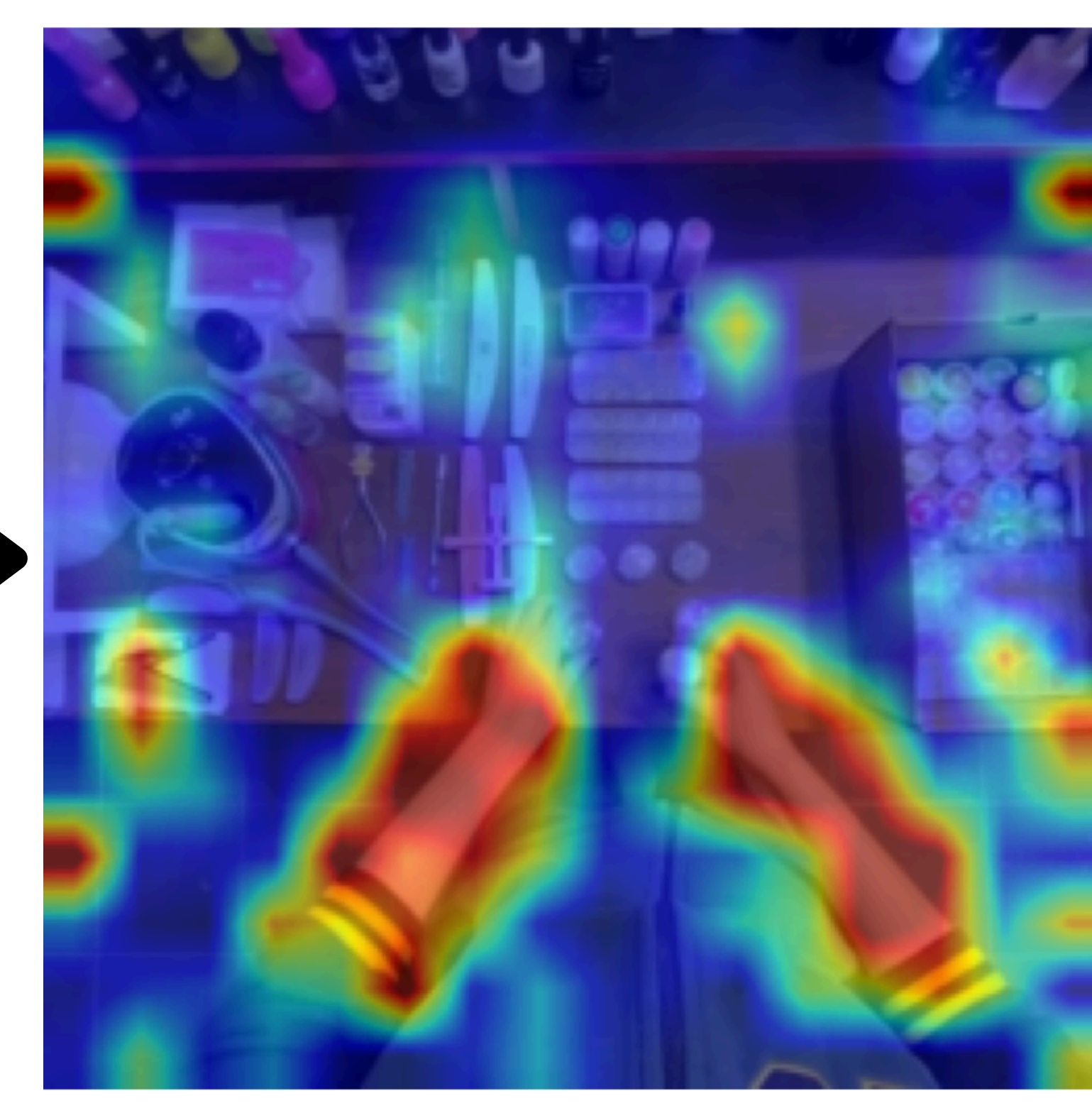
DINOv2

Spatial Detail ✓
Action-Centric ✗

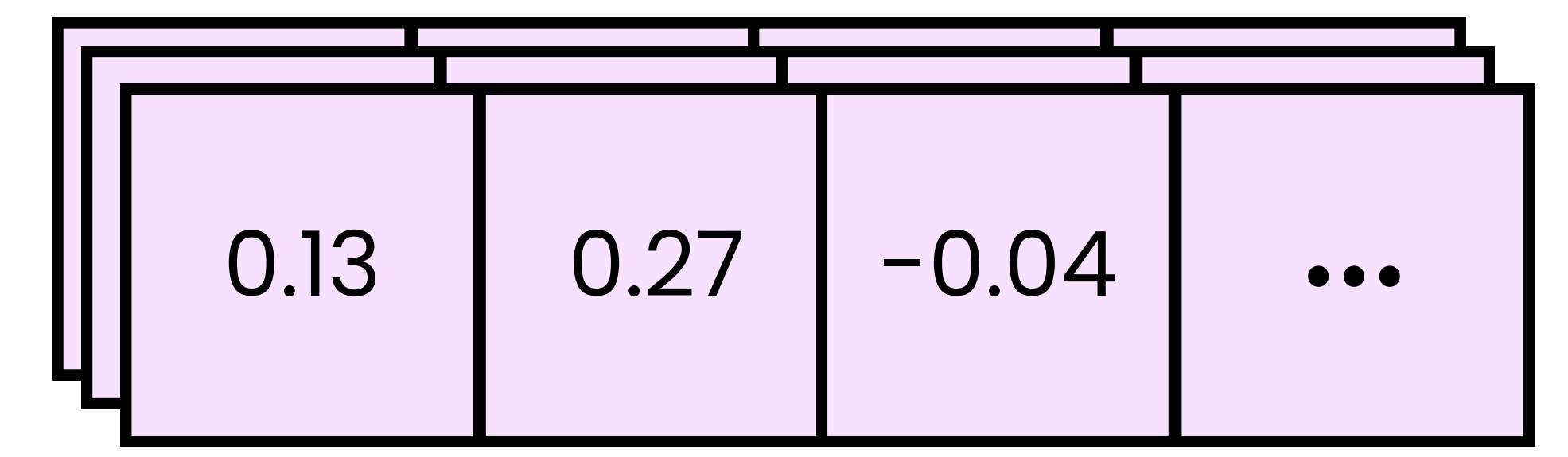


CAIP

Manipulation-relevant features ✓
Action-Centric ✓



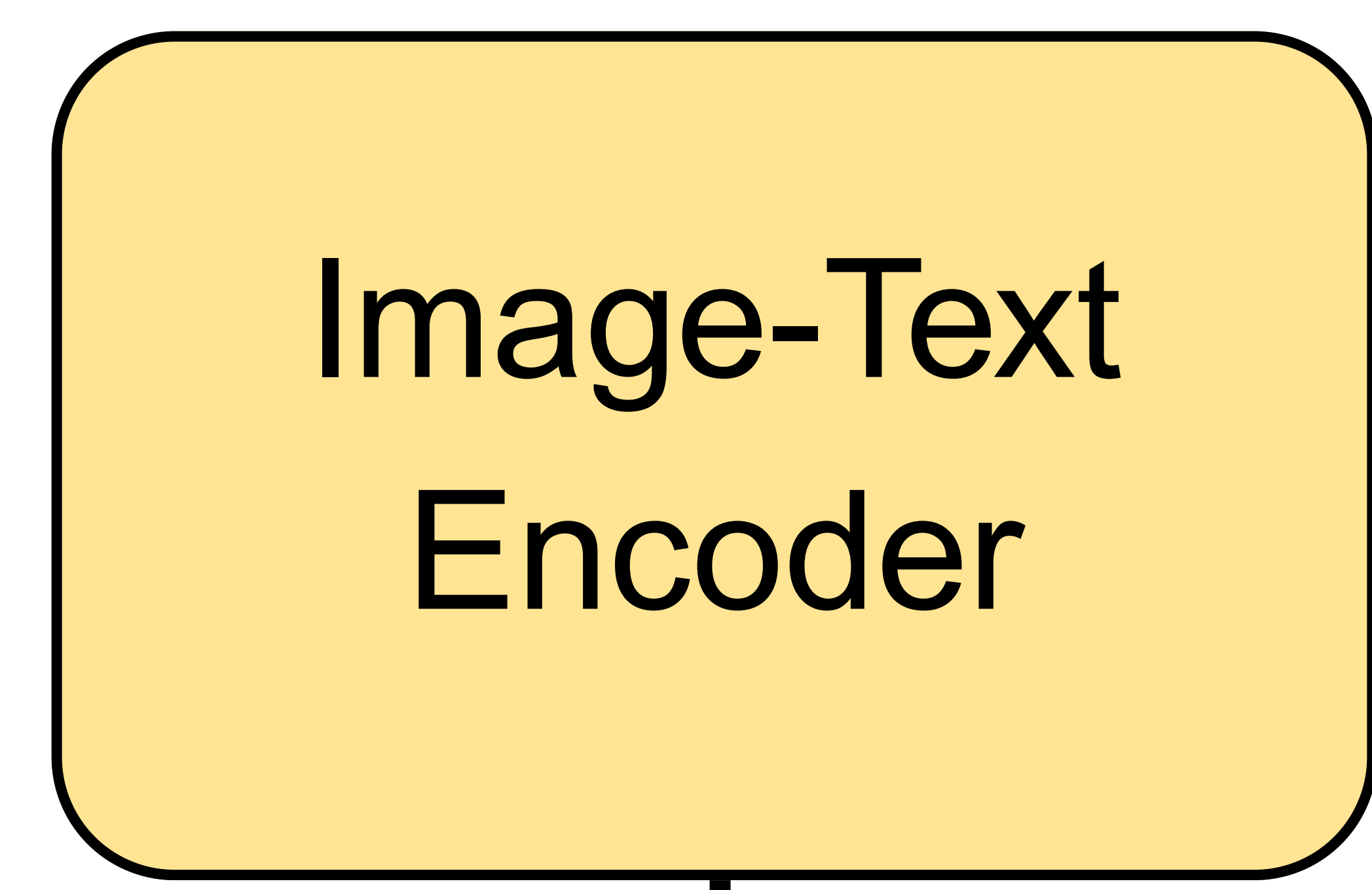
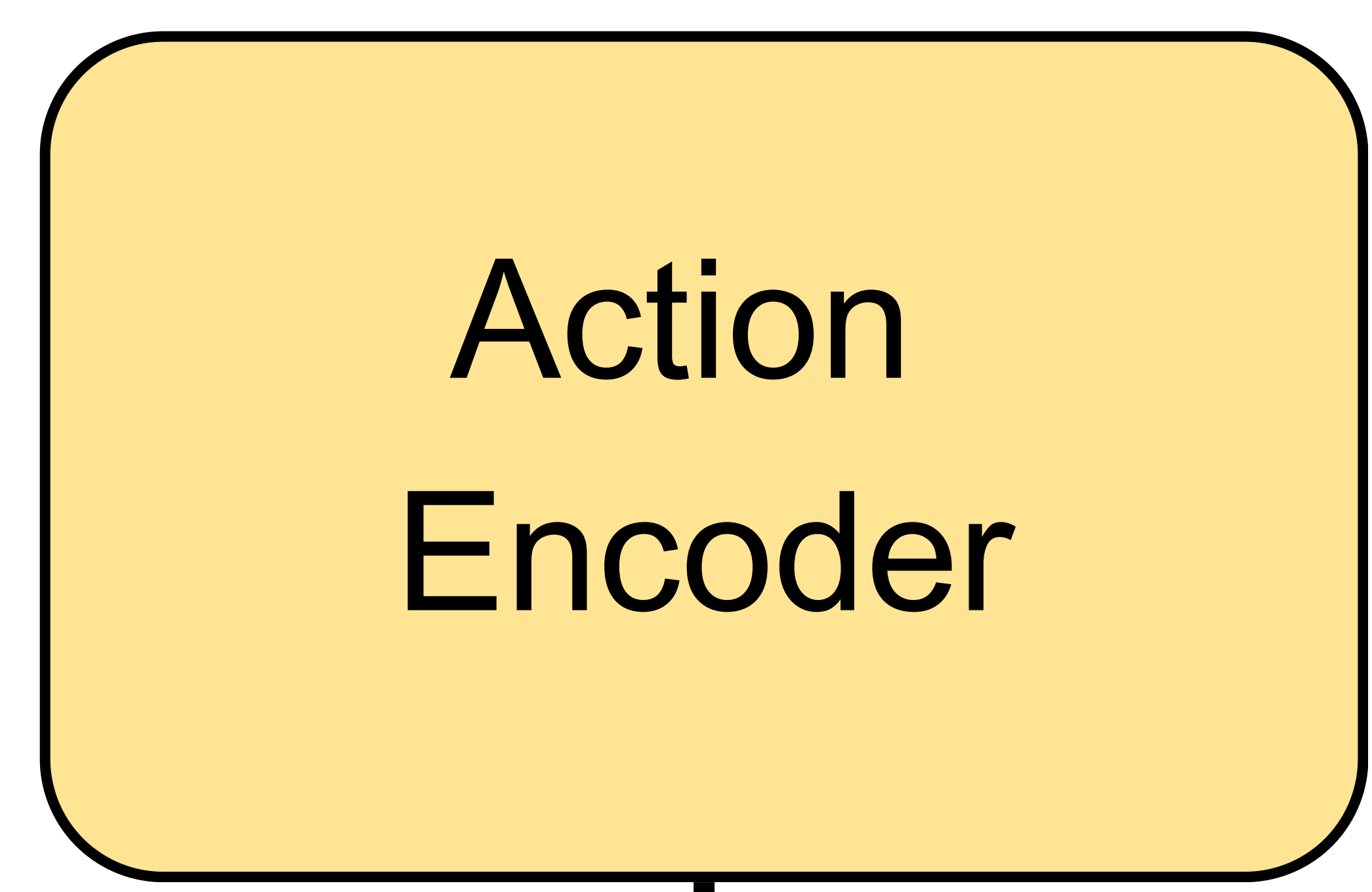
Contrastive Action-Image Pre-training



Hand Pose Actions



“Vacuum floor...”



Action Embedding

Image-Text Embedding

	z_1^i	z_2^i	z_3^i
z_1^a	+	-	-
z_2^a	-	+	-
z_3^a	-	-	+

Contrastive Training

Downstream Policy Performance

