

(a). From video prediction to visual editing

Video-generation WAM

Observation

Predict full future videos



✗ redundant visual details

✗ higher latency

✗ wrong generation

ImageWAM (Ours)

Observation

Task-relevant state



+ "Scan barcode on the object."



✓ task-relevant modification

✓ single edited state

✓ lower latency

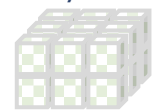
(b). ImageWAM backbone

Observation



VAE

Noisy Latent



Instruction

"Scan barcode on the object."

Tokenizer

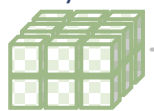
Image Editing Backbone

KV

VAE



Noisy Latent



Action Expert

Robot Actor

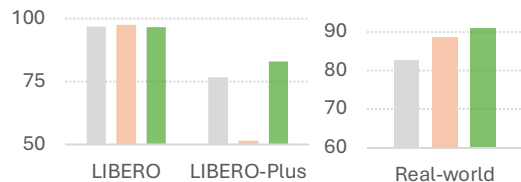
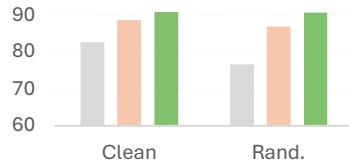


(c). Better performance

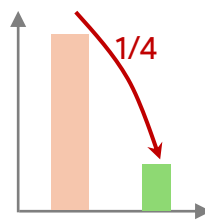
■ VLA ■ WAM ■ ImageWAM



Success Rate (%)



Latency



FLOPs

