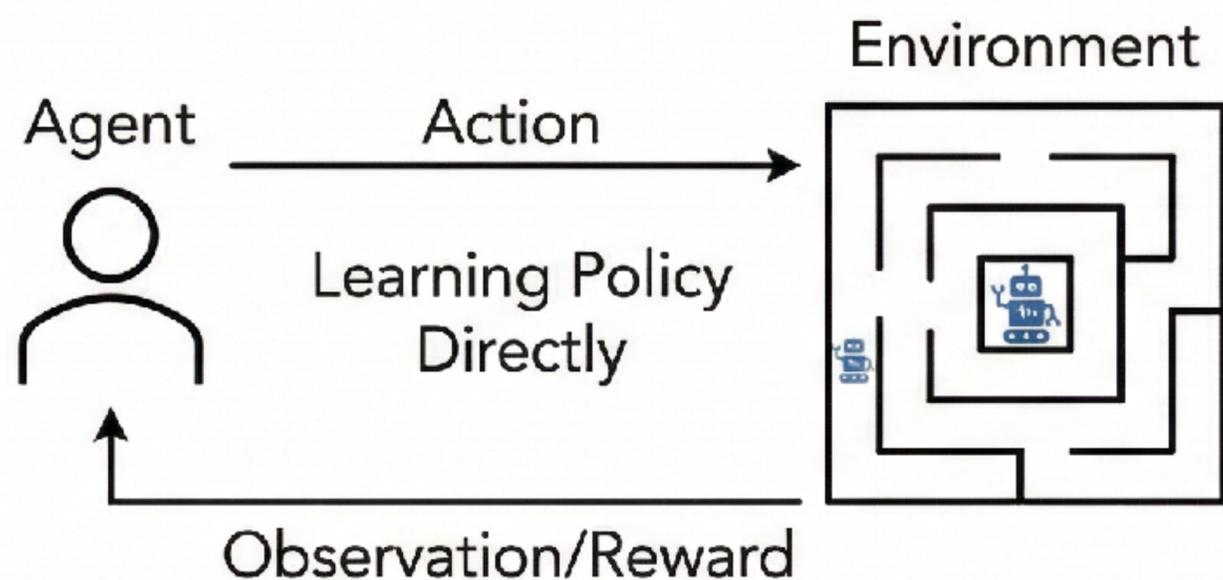# CSCI1470

**Deep Learning**

**Randall Balestriero**
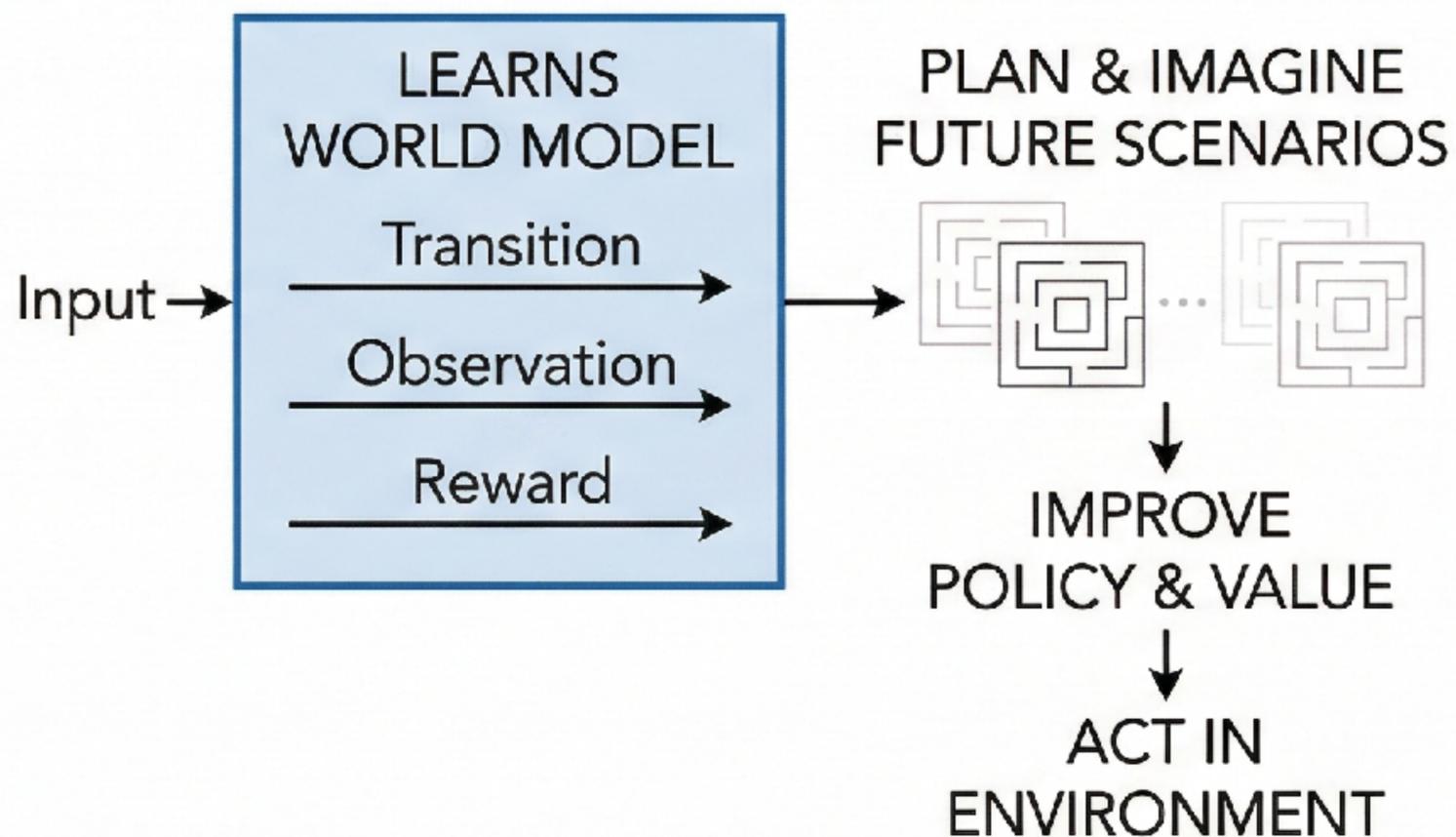
# Recap

# MODEL-BASED REINFORCEMENT LEARNING: THE CONCEPT
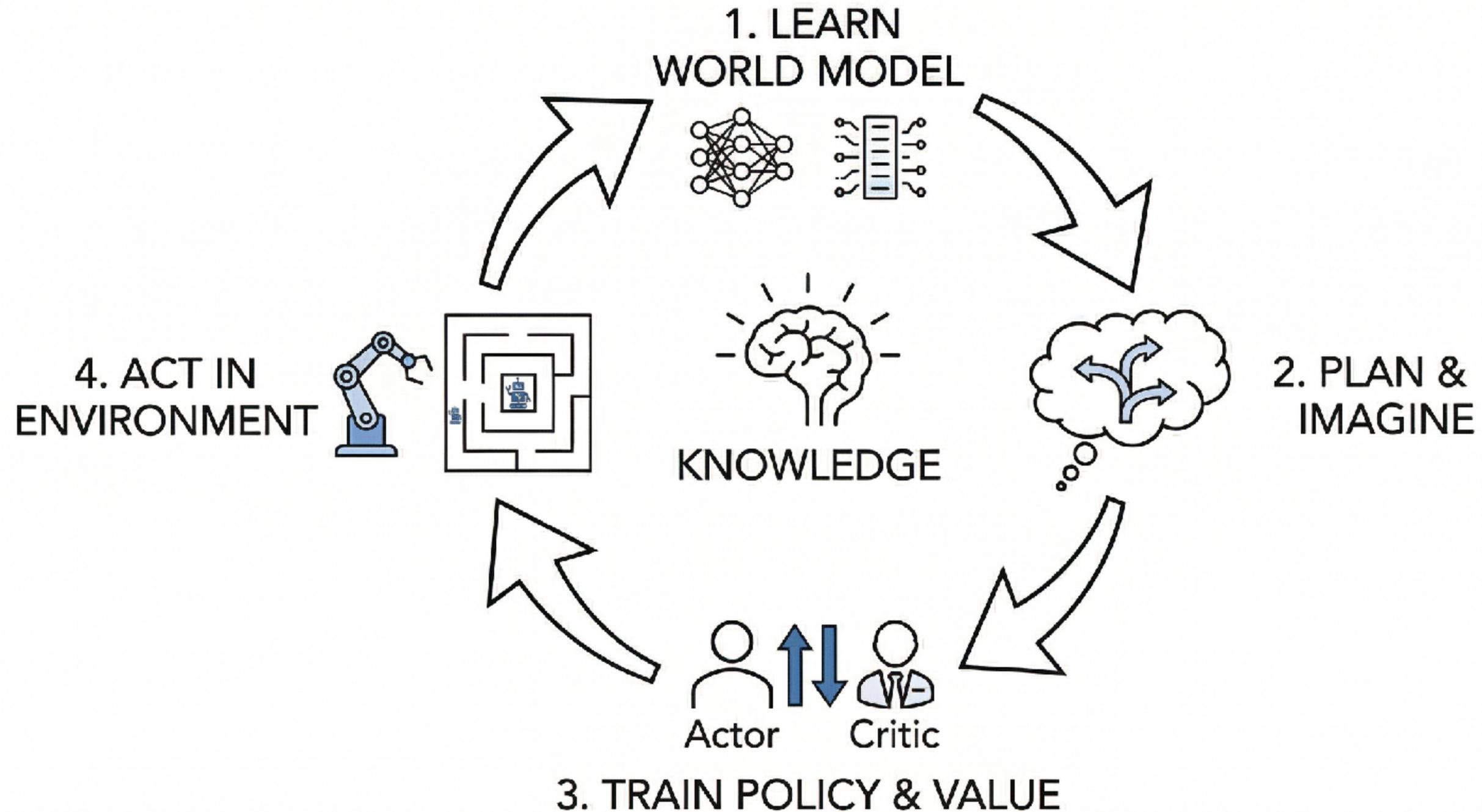
# THE MODEL-BASED RL CYCLE

# EXAMPLE: PIXEL-SPACE WORLD MODEL (DREAMER-v3)



Observation $o_t$ (pixels)

ENCODER

Latent $z_t$

$$\begin{Bmatrix} \vec{v}_t \\ \vec{v}_t, \\ |\vec{v}|_{in}^l \end{Bmatrix}$$

Latent $z_t$

DECODER

Reconstruction $\hat{o}_t$

Reconstruction $\hat{o}_t$

CRITICAL PIXEL RECONSTRUCTION LOSS

REWARD MODEL → Scalar value

Action $a_t$ → TRANSITION MODEL → Predicted Latent $z_{t+1}$

# DreamerV4!

## 1. TRANSFORMER-BASED TRANSITION MODEL

Latent $z_t$ → TRANSFORMER

Action $a_t$ → TRANSFORMER

→ Predicted Latent $z'_{t+1}$

*Models longer horizons.

*Scalable and computationally efficient.

## 2. HYBRID LATENT SPACE LEARNING

Pixel reconstruction Loss

DECODER → $\hat{o}_t$ (Reconstruction)

Initial training provides visual details via reconstruction.

$o_t$

Observations

$o_{t+1}$

ENCODER → Latent $z_t$

→ Encoded Latent $z_{t+1}$

PREDICTOR → Predicted Latent $z'_{t+1}$

**Stage 1: Reconstruction-Based (Initial)**

LATENT PREDICTION LOSS With Shortcut Forcing

Later prediction handles function without a decoder.

**Stage 2: Prediction-Based (Planning)**

*Initial training captures visual detail via pixel reconstruction.

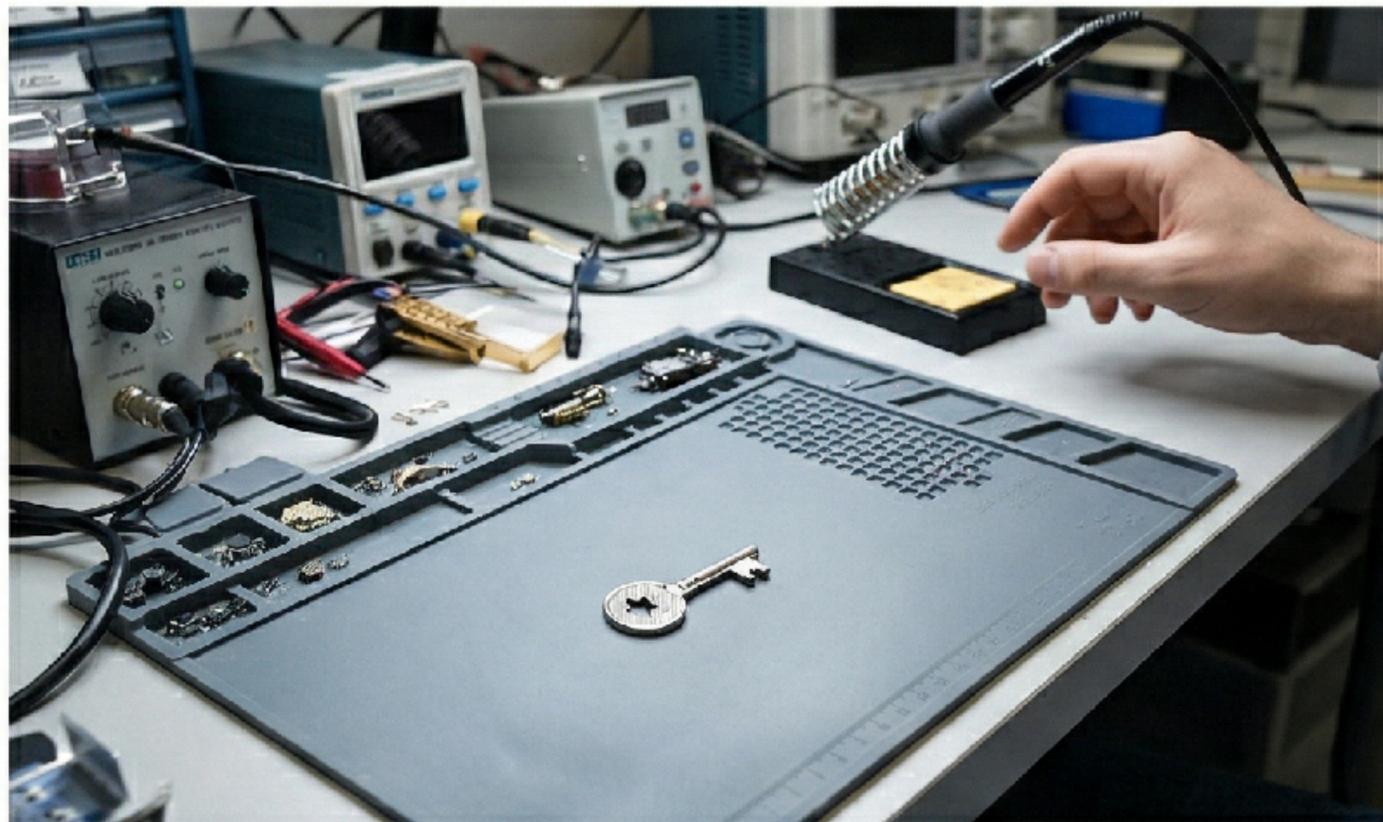*Later training predicts functional structure without a decoder.

## TAKEAWAY: Data-Efficient & Scalable Real-Time Performance and Control from Pixels.
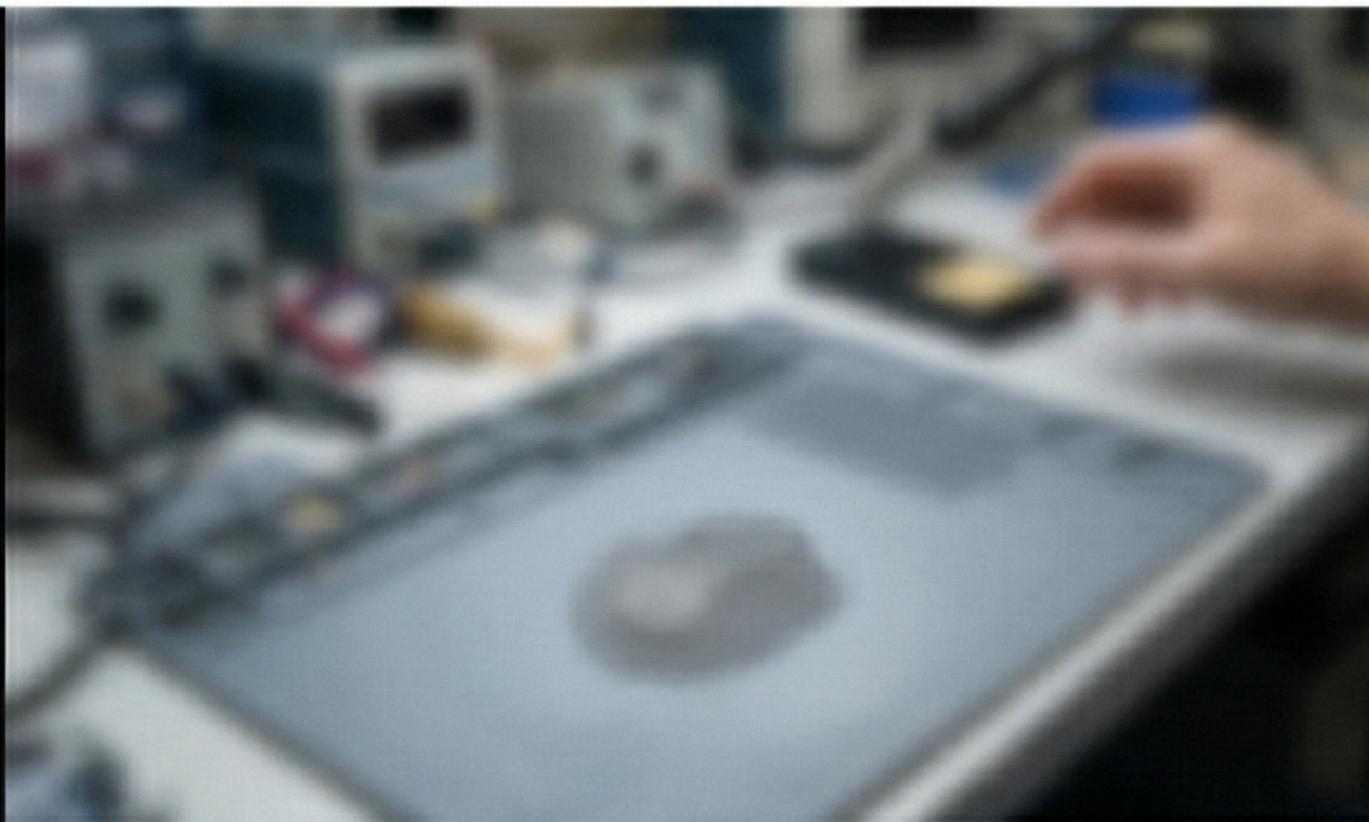
• Scales to large unlabeled datasets.

Fast real-time (20 FPS+)

Offline ur-picking context (Minecraft, from <IMAGE 6)

# FOOD FOR THOUGHT



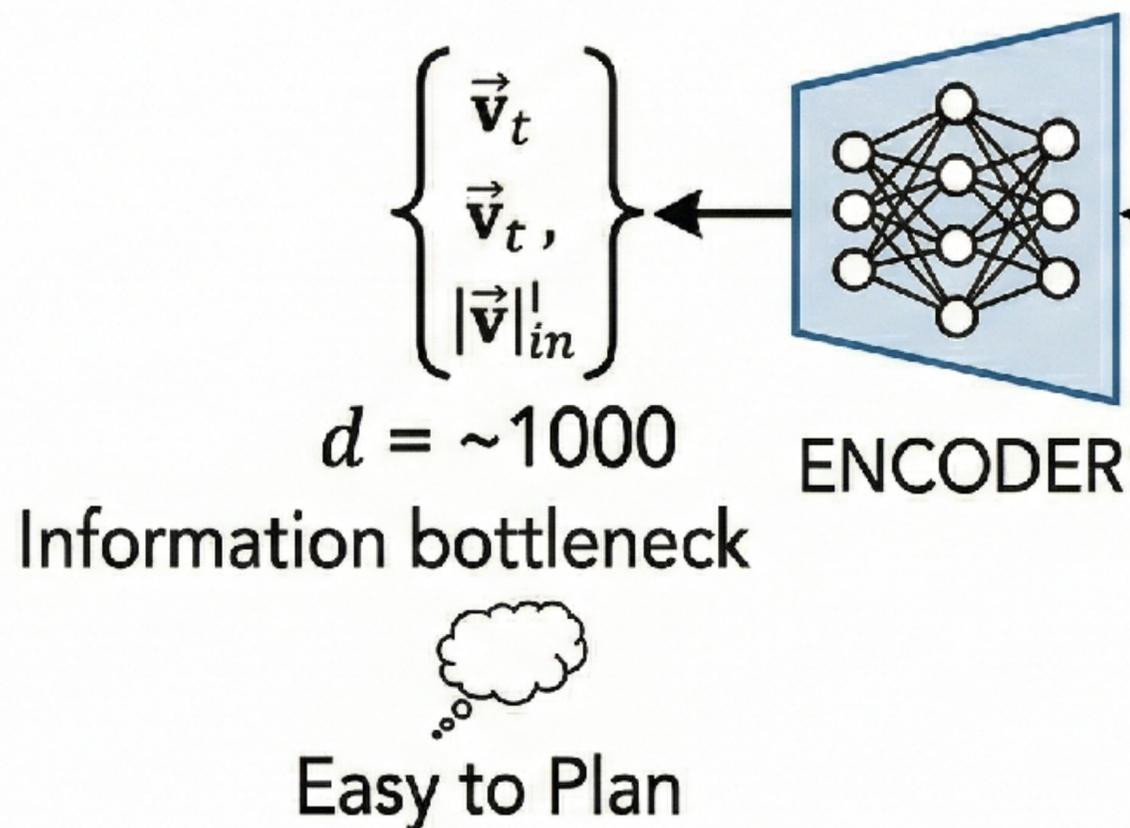THE REALITY $(o_t)$       DECODED IMAGINATION $(\hat{o}_t)$
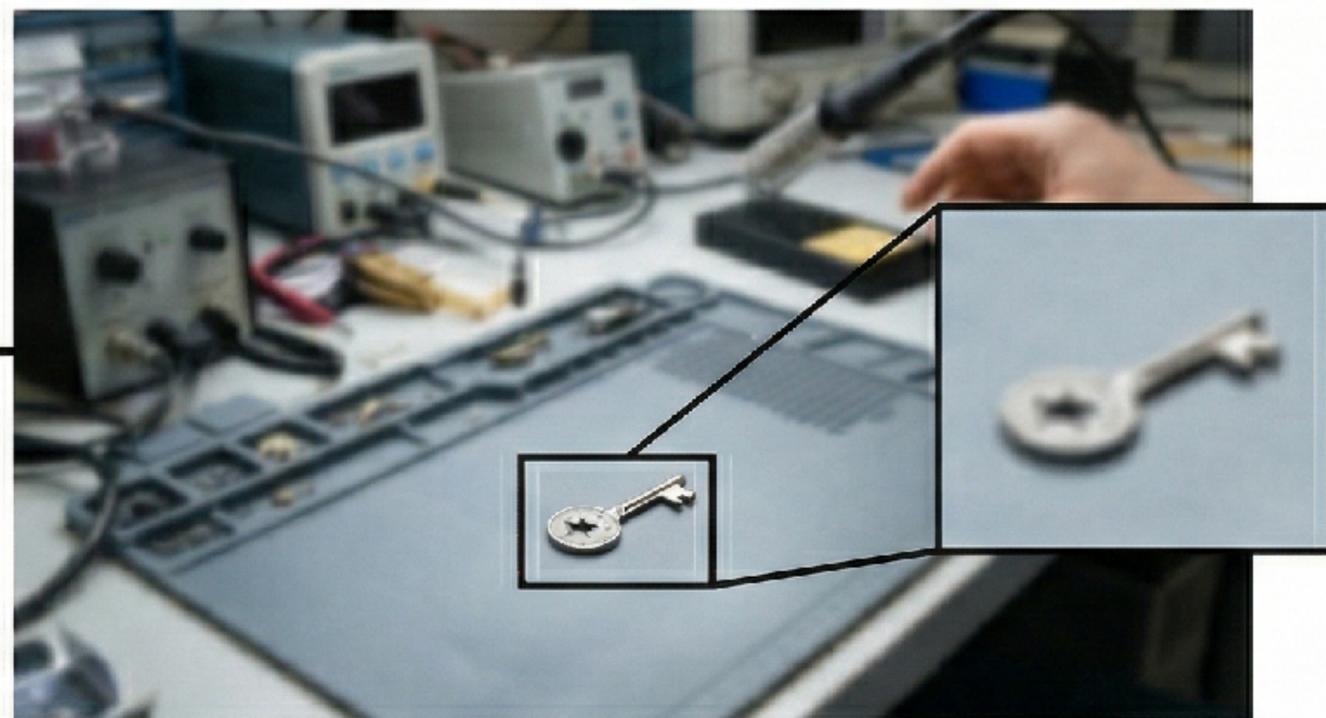
## IS RECONSTRUCTING EVERY PIXEL A GOOD USE OF COMPUTATION?

Wasteful? Irrelevant details? Wrong metric?

# THE CORE CHALLENGE: COMPARING IMAGE OBSERVATIONS



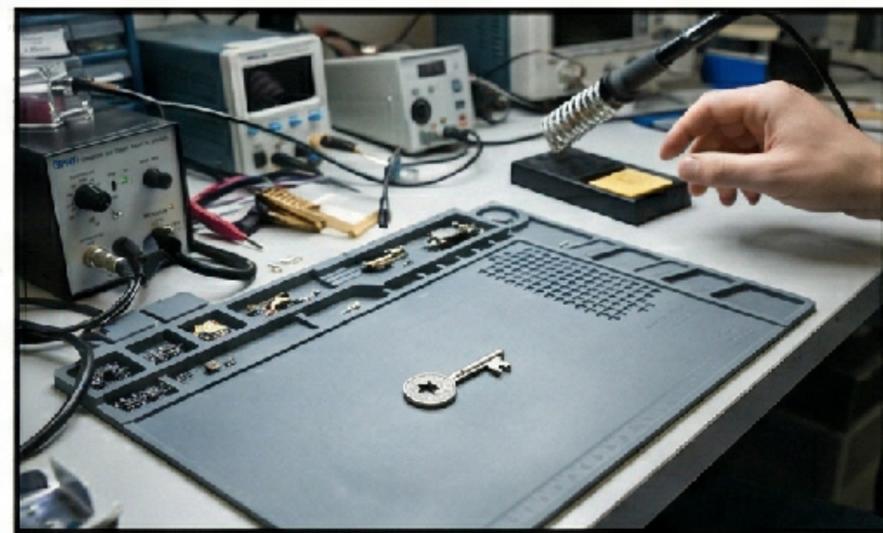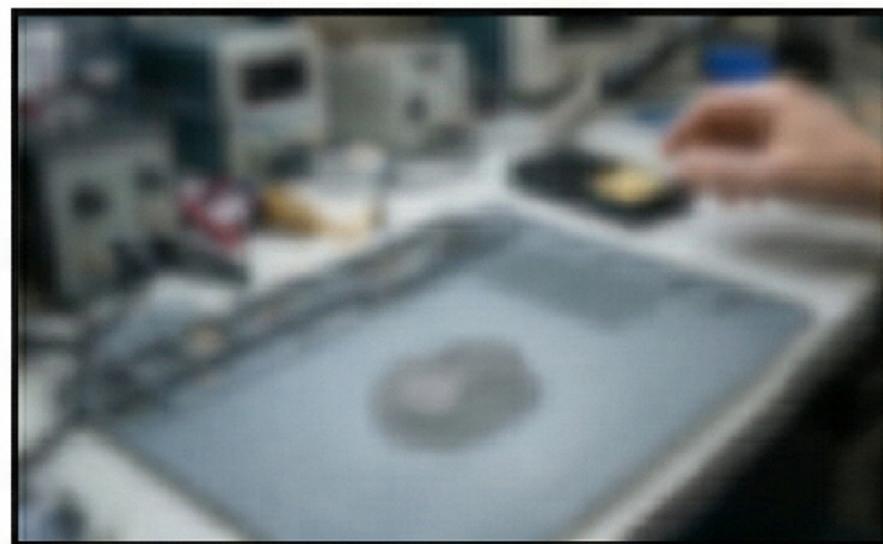Original Observation $o_t$

Reconstruction $\hat{o}_t$ (MSE ≈ 50)

High Visual Similarity (Shift)

❌

**PIXEL MSE IS POOR PROXY PERCEPTION**

❌

Low Visual Similarity (Blur)
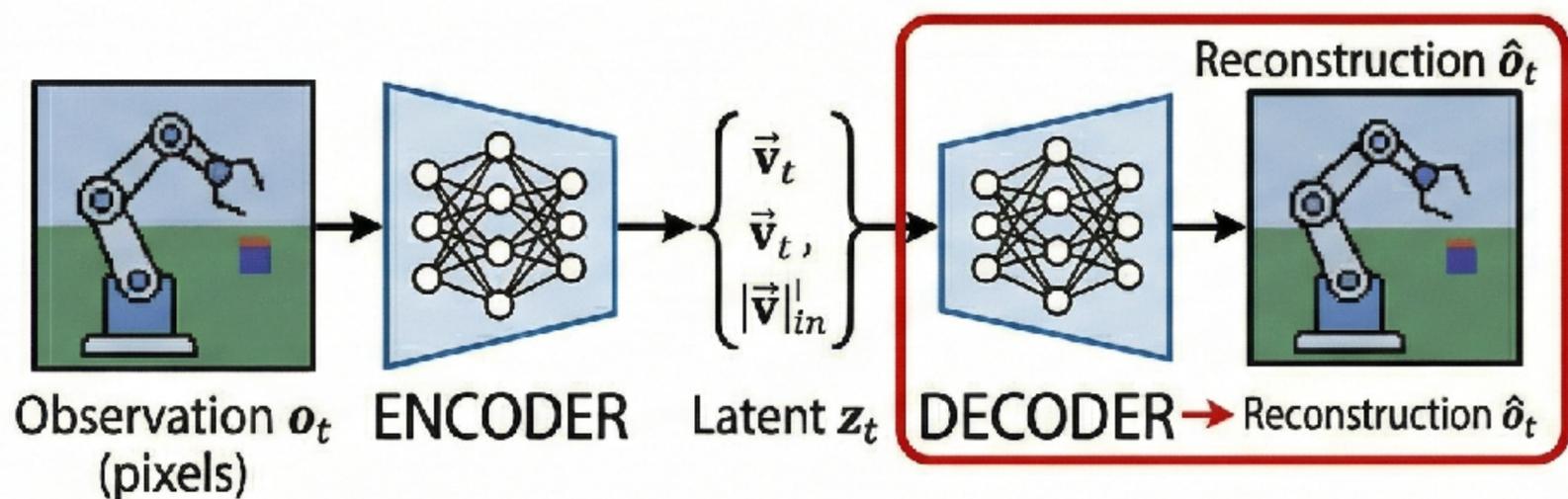
Reconstruction $\hat{o}_t$ (MSE ≈ 50)

# SUMMARY: THE PARADIGM SHIFT

**PIXEL-SPACE WORLD MODEL**

Observation $o_t$ (pixels) → ENCODER → Latent $z_t$ $\left\{\begin{array}{l}\vec{\mathbf{v}}_t \\ \vec{\mathbf{v}}_t, \\ |\vec{\mathbf{v}}|_{in}\end{array}\right\}$ → DECODER → Reconstruction $\hat{o}_t$

Reconstruction $\hat{o}_t$
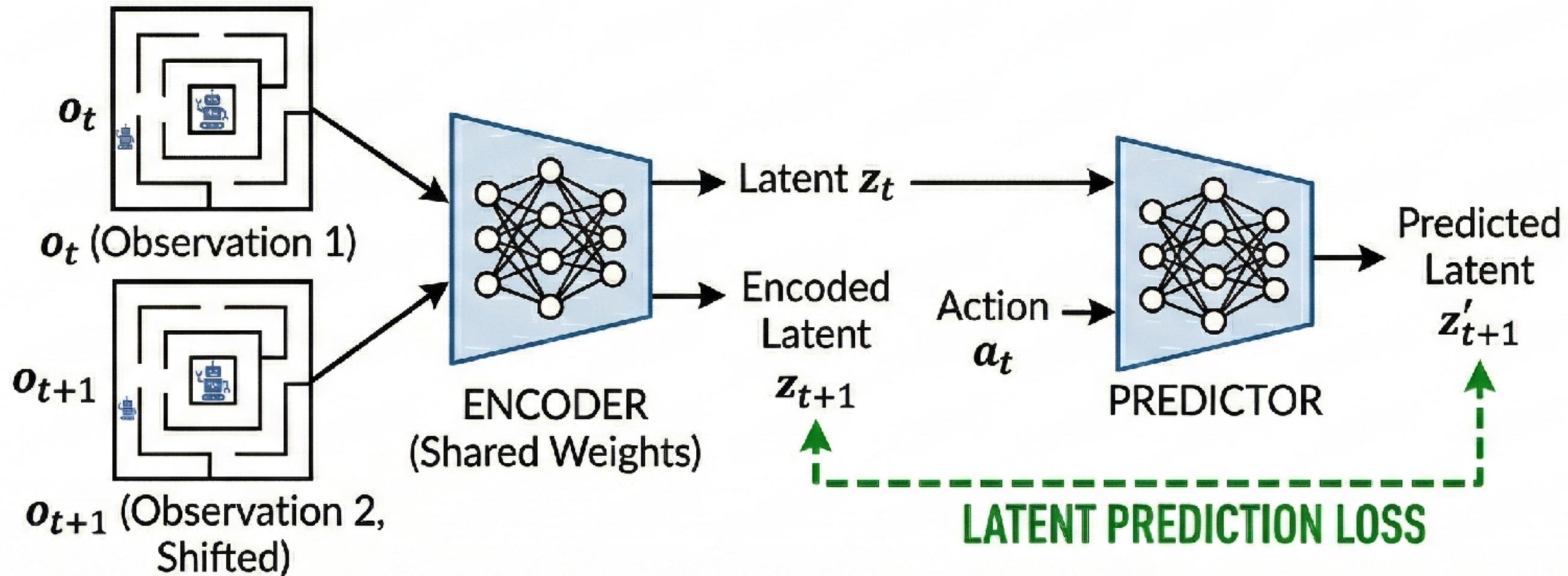
- High Reconstruct Cost
- Metric = MSE/Pixels
- Focus on *Appearance*
- Data Hungry

**RECONSTRUCTION-FREE WORLD MODEL**

Object 3 icons    Abstraction icon

- Zero Reconstruct Cost
- Metric = Dynamics/Function
- Focus on *Structure*
- Data Efficient

# Thank you!
# See you Friday!