

## Overview

### Motivation

Generalize to novel physical manipulation tasks with *compositional* structure

### Idea 1: Entity Abstraction

*Symmetric* local processing of entities, rather than global processing of scenes, enables knowledge about an entity in one context to directly transfer to modeling the same entity in different contexts.

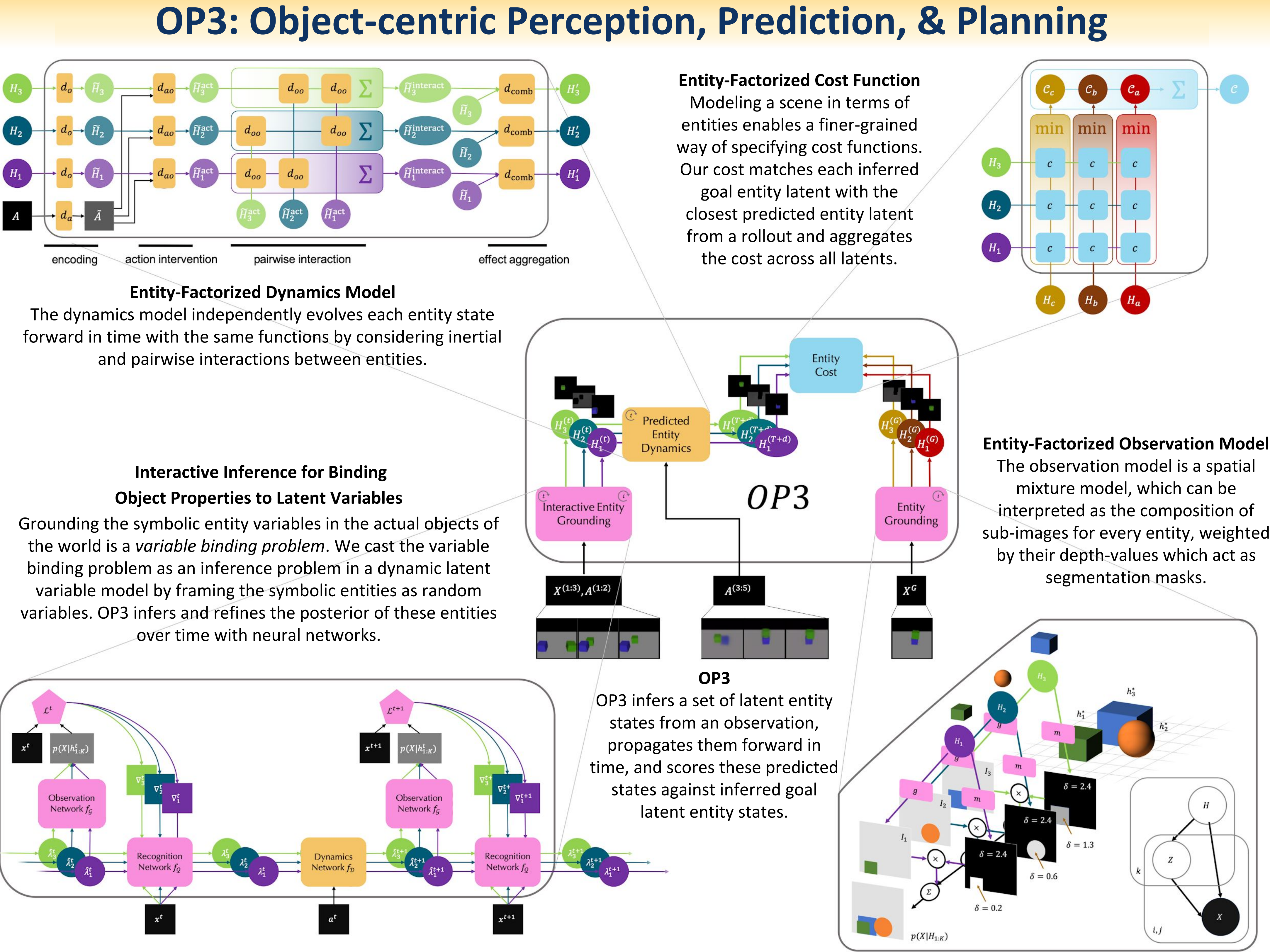
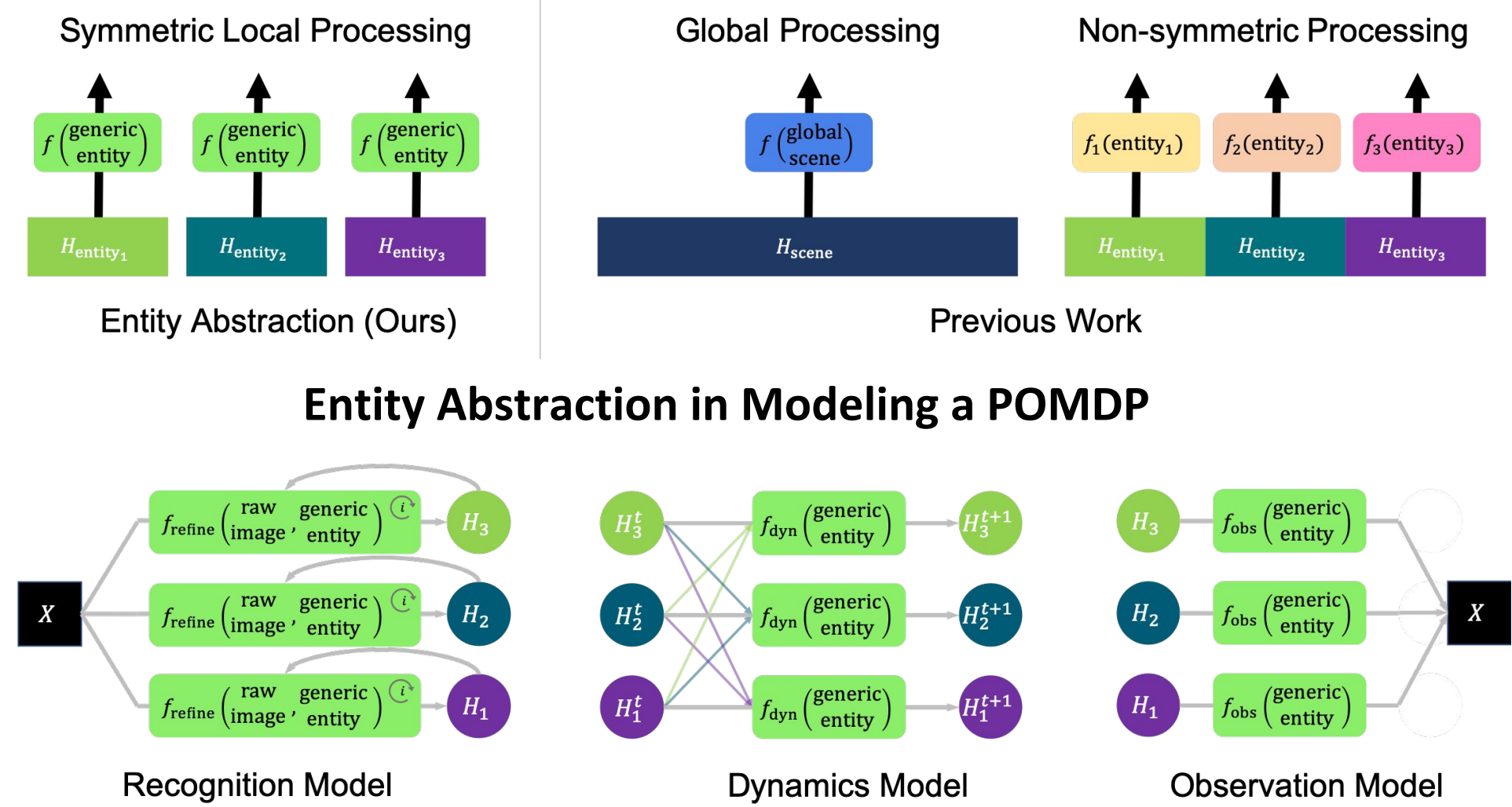
### Idea 2: Interactive Inference

To ground symbolic entity variables in actual objects, use iterative inference in a entity-factorized dynamic latent variable model.

### Contribution

A factorized model-based reinforcement learning framework with entity variables inferred directly from visual observations and actions without any object-level supervision.

## Entity Abstraction



## Interactive Inference

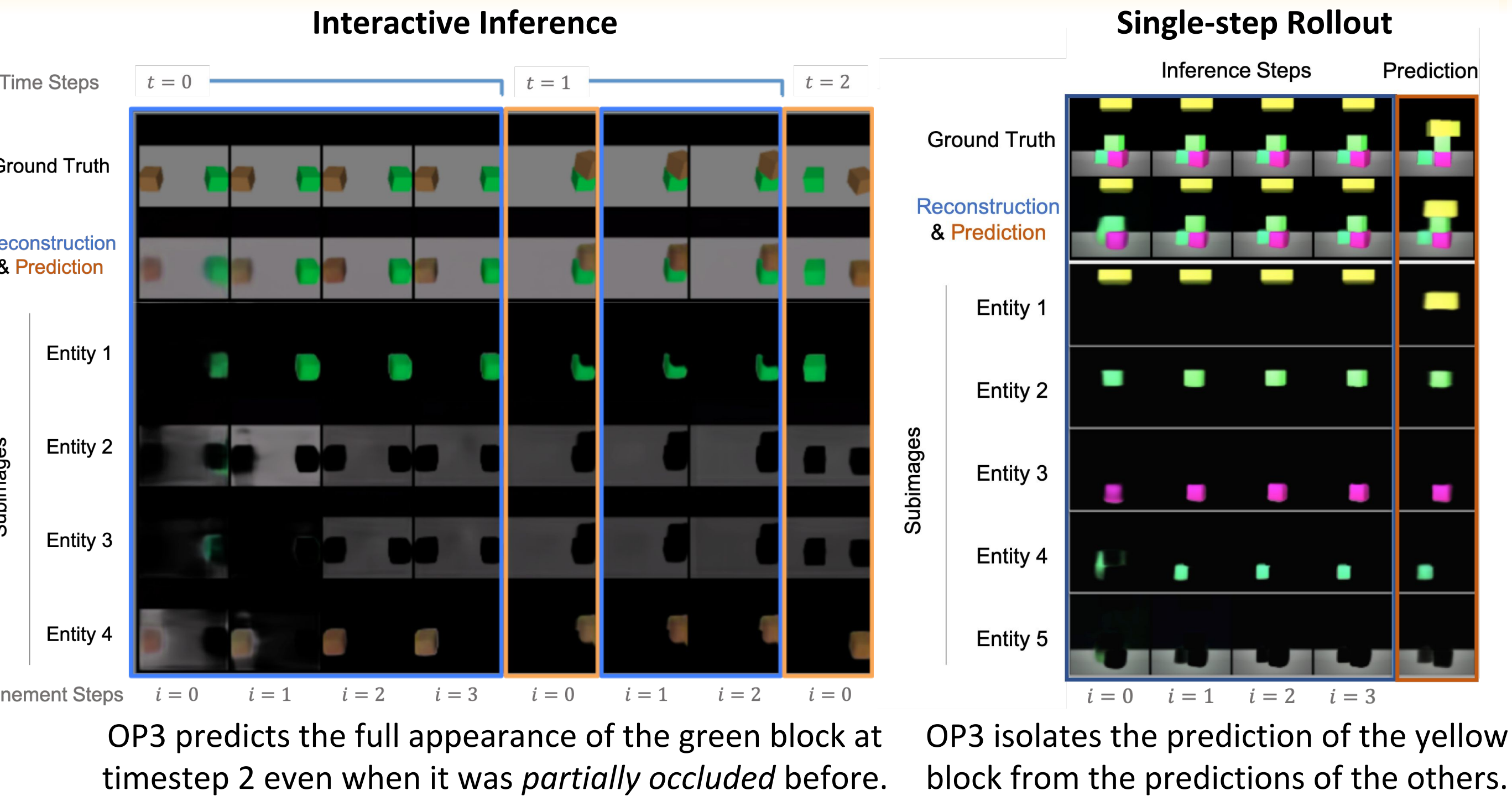
### Generative Distribution

$$p\left(X^{(0:T)}, H_{1:K}^{(0:T)} \mid a^{(0:T-1)}\right) = p\left(H_{1:K}^{(0)}\right) \prod_{t=1}^T p\left(H_{1:K}^{(t)} \mid H_{1:K}^{(t-1)}, a^{(t-1)}\right) \prod_{t=0}^T p\left(X^{(t)} \mid H_{1:K}^{(t)}\right)$$

### Inference Algorithm

**Algorithm** Interactive Inference: Timestep  $t$

**Input:** observation  $x^{(t)}$ , action  $a^{(t)}$ , previous entity states  $h_{1:K}^{(t-1)}$   
**Predict**  $\lambda_k^{(t,0)} \leftarrow f_\lambda\left(h_k^{(t-1)}, h_{[\neq k]}^{(t-1)}, a^{(t)}\right)$  for each entity  $k$   
**for**  $i = 0$  **to**  $M - 1$  **do**  
    Sample  $h_k^{(t,i)} \sim \mathcal{N}\left(\lambda^{(t,i)}\right)$  for each entity  $k$   
    Evaluate  $\mathcal{L}^{(t,i)} \approx \log \mathcal{Q}\left(x^{(t)} \mid h_{1:K}^{(t,i)}\right) - D_{KL}\left(\mathcal{N}\left(\lambda_{1:K}^{(t,i)}\right) \parallel \mathcal{N}\left(\lambda_{1:K}^{(t,0)}\right)\right)$   
    Calculate  $\nabla_{\lambda_k} \mathcal{L}^{(t,i)}$  for each entity  $k$   
    Assemble auxiliary inputs  $\beta_k$  for each entity  $k$   
    Update  $\lambda_k^{(t,i+1)} \leftarrow f_\lambda\left(x^{(t)}, \nabla_{\lambda_k} \mathcal{L}^{(t,i)}, \lambda_k^{(t,i)}, \beta_k^{(t,i)}\right)$  for each entity  $k$   
**end for**  
**return**  $\lambda^{(t,M)}$



## Planning Algorithm

### Algorithm OBJECT-CENTRIC-PLANNING

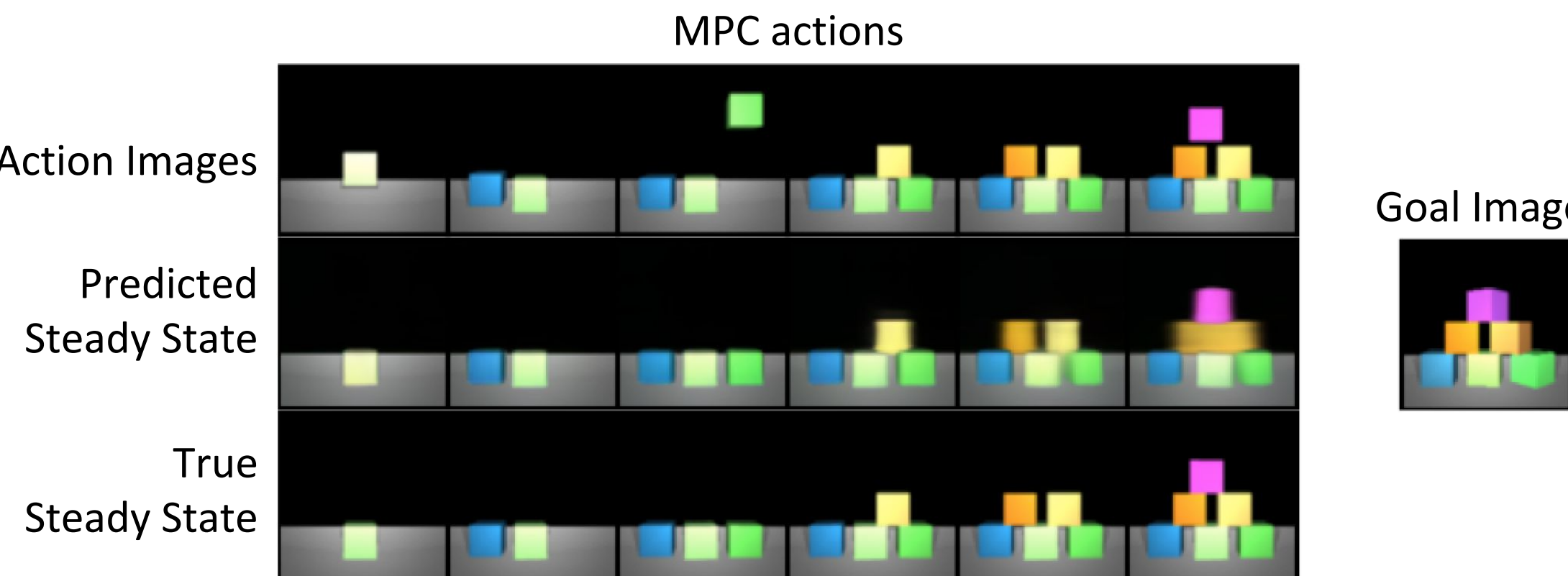
- Input:**  $x^{(0:\tau)}, a^{(0:\tau-1)}, x^{(G)}$
- Initialize:**  $\lambda_{1:K}^{(0)}$
- $h_{1:K}^{(G)} \leftarrow \text{GOAL INFERENCE}(\lambda_{1:K}^{(0)}, x^{(G)})$
- $\lambda_{1:K}^{(\tau)}, h_{1:K}^{(\tau)} \leftarrow \text{STATE ACQUISITION}(\lambda_{1:K}^{(0)}, x^{(0:\tau)}, a^{(0:\tau-1)})$
- for**  $t \leftarrow \tau$  **to**  $T$  **do**
- $a^t \leftarrow \text{ACTION SELECTION}(h_{1:K}^{(t)}, h_{1:K}^{(G)})$
- $x^{t+1} \leftarrow \text{ENVIRONMENT-STEP}(a^t)$
- $\lambda_{1:K}^{(t+1)}, h_{1:K}^{(t+1)} \leftarrow \text{STATE ACQUISITION}(\lambda_{1:K}^{(t)}, a^{(t)}, x^{(t+1)})$
- end for**
- return**  $a^{(\tau:T)}$

**Goal Inference:** Iterative inference on the static goal image to infer goal latent entities  
**State Acquisition:** Interactive inference on a sequence of observed images and actions to estimate the latent entities  
**Action Selection:** Run CEM by rolling out many proposed actions sequences and scoring them with the cost function

## Experiments

### Single-Step Block-Stacking (Predict Fall)

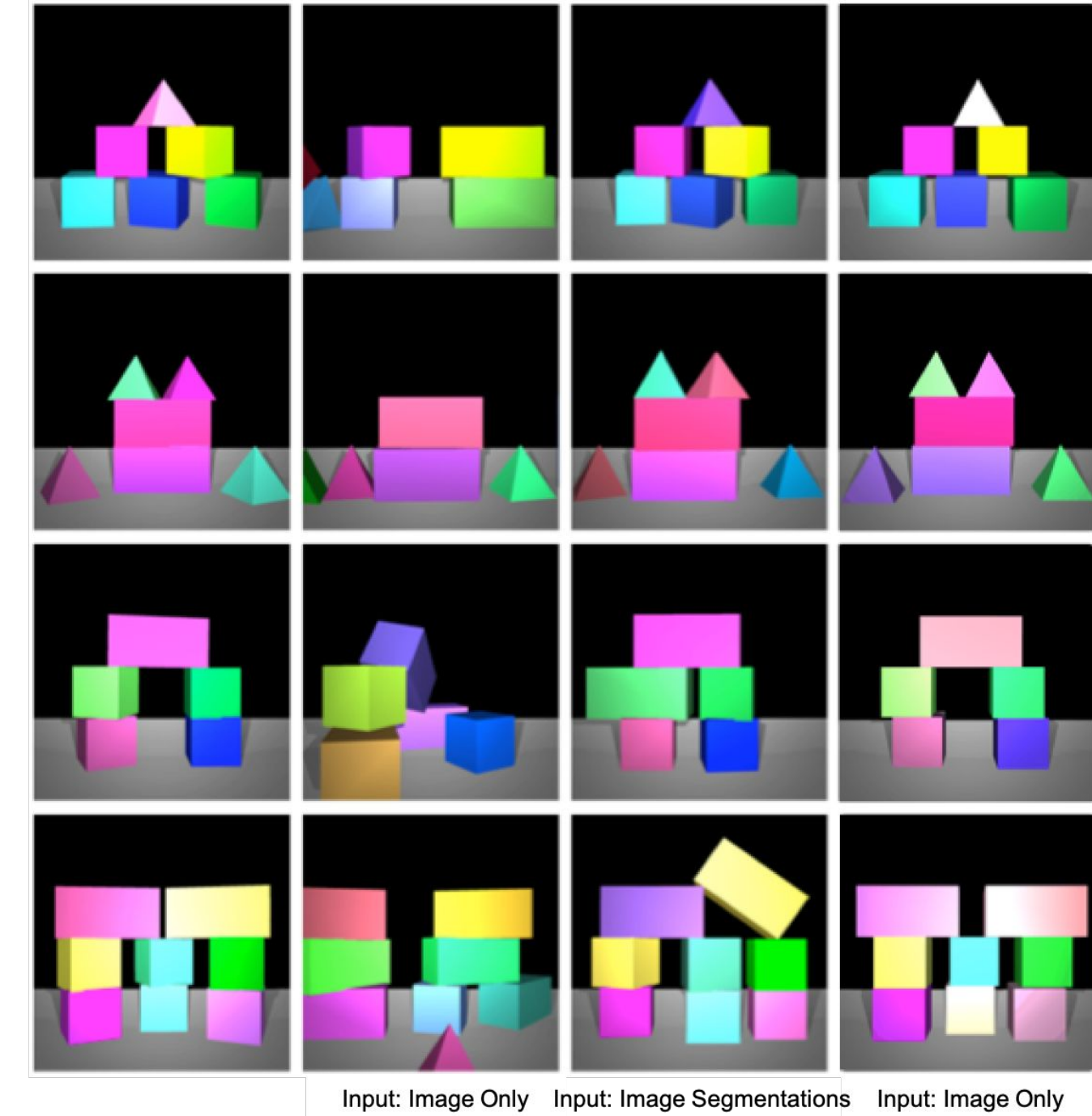
OP3 can generalize to building complex block structures by planning in object-space from raw pixel input. We plan by first inferring the hidden states of each proposed action image, predicting the resultant steady state using our dynamics model, and then comparing the hidden states with the goal states using our cost function.



An *action image* depicts how an action intervenes on the state by raising a block in the air. OP3 is trained to predict the steady-state outcome of dropping the block.

SAVP (baseline, no factorization)	O2P2 (oracle, pre-segmented input)	OP3 (ours) (raw image input)
24% accuracy	76% accuracy	<b>82% accuracy</b>

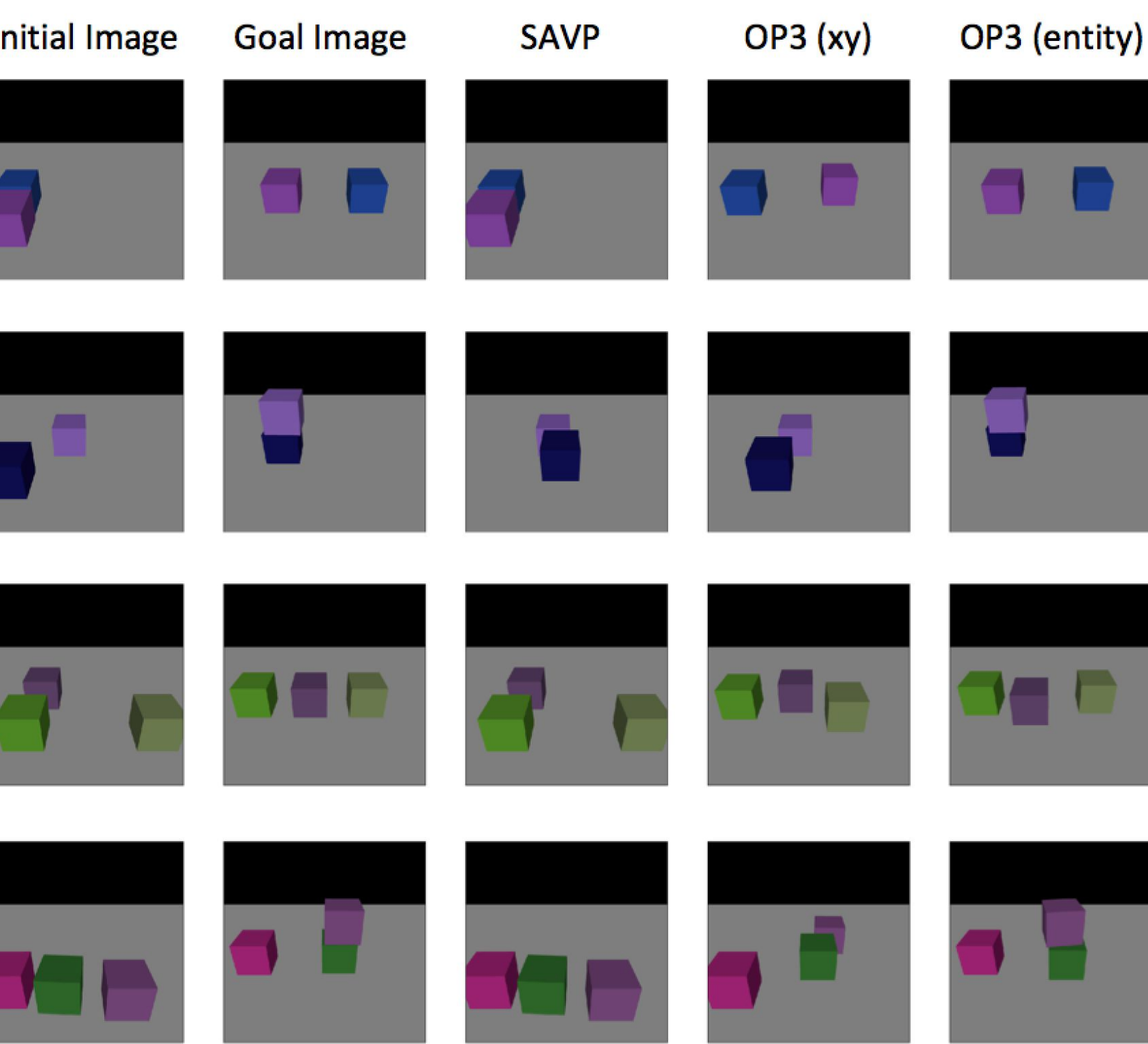
Goal Image   SAVP   O2P2   OP3 (ours)



**Results**  
OP3 was only trained on up to five objects but can generalize to nine objects and new configurations during testing. OP3 achieves three times the accuracy of a state-of-the-art video prediction model without entity abstraction (SAVP) and outperforms an oracle model (O2P2) that receives the object segmentations.

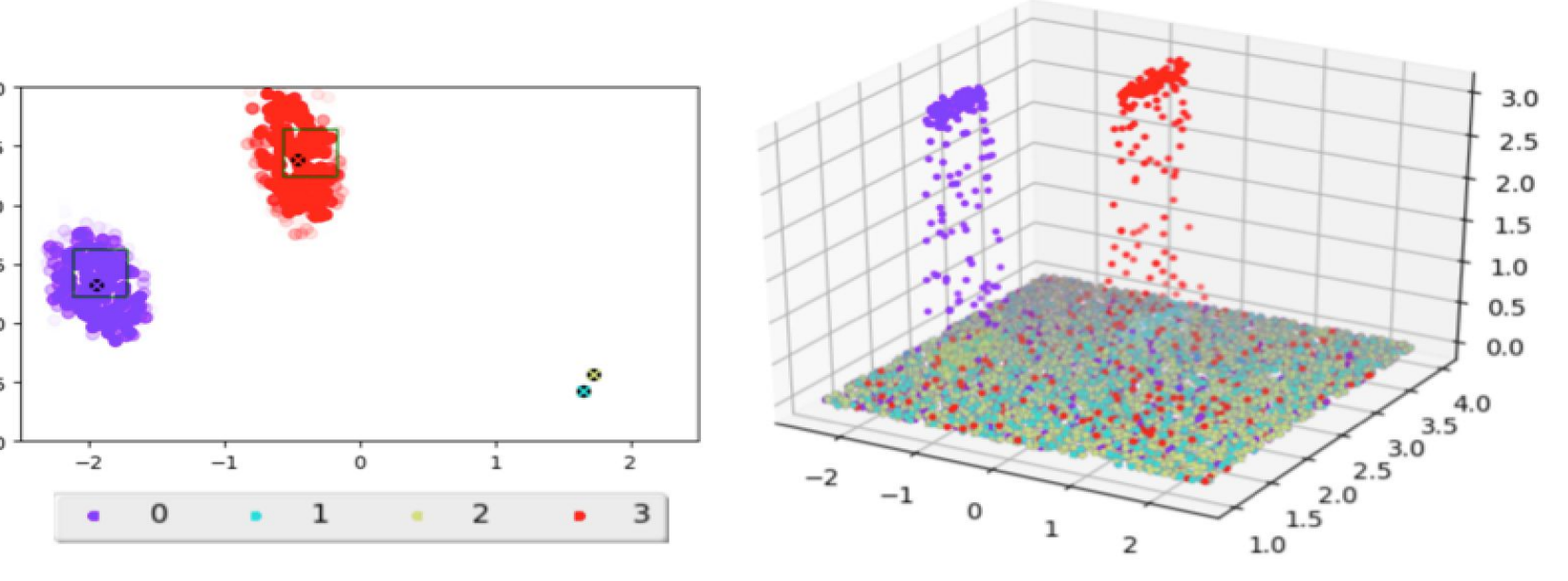
### Multi-Step Block-Stacking (Pick/Place)

On a sparse block-stacking task, OP3 can plan over the space of objects multiple steps into the future.



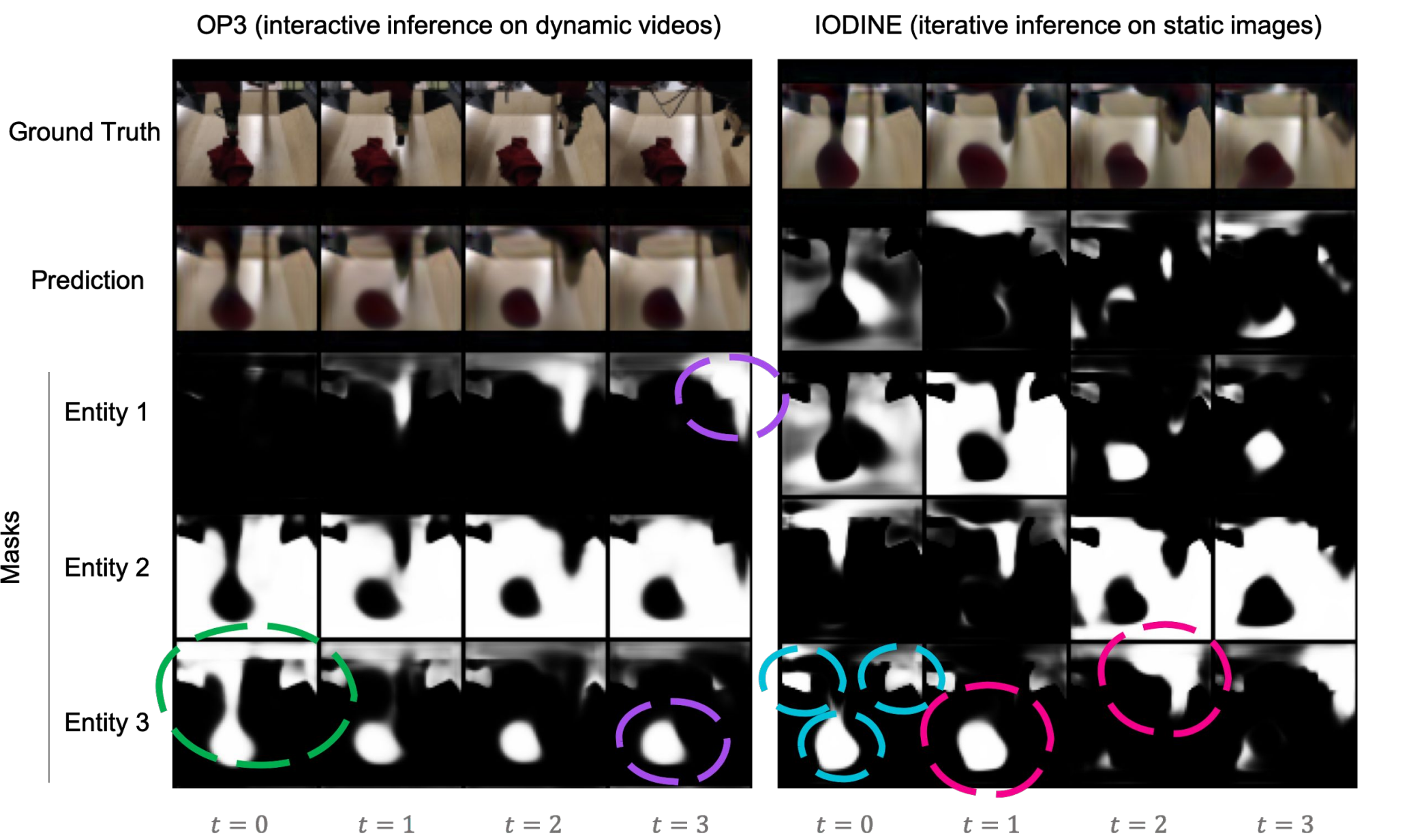
#Blocks	SAVP	OP3 (xy)	OP3 (entity)
1	54%	73%	<b>91%</b>
2	28%	55%	<b>80%</b>
3	28%	41%	<b>55%</b>

OP3 (xy): action space is (pick\_xy, place\_xy)  
OP3 (entity): OP3 has access to *pointers* to each entity, enabling an *entity-centric action space* (entity\_id, place\_xy)



### Real World Evaluation

We evaluate how well OP3 can disambiguate objects on a robotic pushing task with clutter and occlusions.



Initially, both OP3 (green circle) and IODINE (cyan circles) both disambiguate objects via color segmentation. As time progresses, OP3 uses temporal continuity and interactive feedback to disambiguate latents (purple), whereas applying IODINE on a per-frame basis cannot do so.

## Conclusion

- OP3 integrates graphical models, symbolic computation, and neural networks in a model-based reinforcement learner
- Models as compositions of *locally-scoped* functions
- Symbolic grounding as variable binding as posterior inference*
- Modeling entities and their interactions provides significant generalization improvement in *combinatorially complex* tasks

Challenge	Solution
Combinatorial Generalization	Entity Abstraction
Grounding Entity Variables	Interactive Inference