



GenParticles: Probabilistic Particle-Based Modeling for Object-Centric Motion

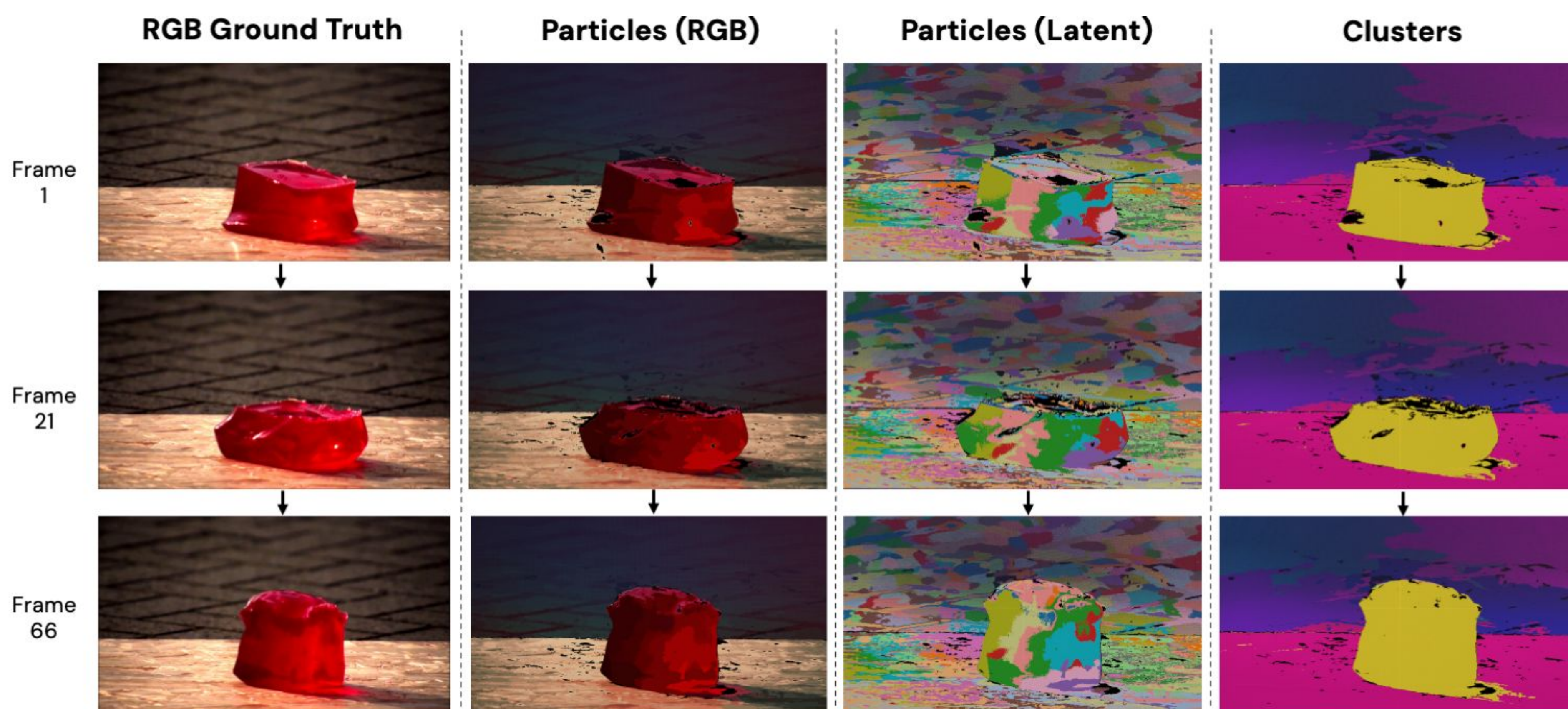
Arijit Dasgupta*, Eric Li*, Mathieu Huot, William T. Freeman, Vikash Mansinghka, Joshua B. Tenenbaum

**Structured World Models
for Robotic Manipulation**

RSS 2025 | Los Angeles, USA | June 21, 2025

Introduction

- Object **motion** and **structure** inference from visual input is key for robotic manipulation in dynamic, disturbed or cluttered scenes.
- Humans use motion cues to infer object identity and deformation, inspiring models for flexible, persistent representations of objects.



Generative Particle Model

Algorithm 1 Generative Particle Model

```

1: Input:
2:    $K, L, N$  ▷ Number of clusters, particles, and observed data points
3:   Priors:  $\alpha, \beta, (\mu^{\mathcal{H}}, \sigma_{\mu^{\mathcal{H}}}^2, \Psi^{\mathcal{H}}, \nu^{\mathcal{H}}), (\Psi^{\mathcal{B}}, \nu^{\mathcal{B}}), \sigma_V^2, (\Psi^{\mathcal{V}}, \nu^{\mathcal{V}})$ 
4:   Sample cluster weights:  $\pi^{\mathcal{H}} \sim \text{Dir}(\alpha)$ 
5:   Sample particle weights:  $\pi^{\mathcal{B}} \sim \text{Dir}(\beta)$ 
6:   for  $k = 1$  to  $K$  do
7:     Sample cluster covariance:  $\Sigma_k^{\mathcal{H}} \sim \mathcal{W}^{-1}(\Psi^{\mathcal{H}}, \nu^{\mathcal{H}})$ 
8:     Sample cluster mean:  $\mu_k^{\mathcal{H}} \sim \mathcal{N}(\mu^{\mathcal{H}}, \sigma_{\mu^{\mathcal{H}}}^2 \mathbf{I})$ 
9:     Sample cluster translation:  $\mathbf{t}_k \sim \text{DiscreteNormal}(\mathbf{0}, s^2 \mathbf{I})$ 
10:    Sample cluster rotation:  $\mathbf{R}_k \sim \text{DiscreteVMF}(\kappa^{\text{vmf}}, \theta_{\max})$ 
11:  end for
12:  for  $\ell = 1$  to  $L$  do
13:    Sample cluster assignment:  $z_{\ell}^{\mathcal{H}} \sim \text{Cat}(\pi^{\mathcal{H}})$ 
14:    Let  $k = z_{\ell}^{\mathcal{H}}$ 
15:    Sample particle covariance:  $\Sigma_{\ell}^{\mathcal{B}} \sim \mathcal{W}^{-1}(\Psi^{\mathcal{B}}, \nu^{\mathcal{B}})$ 
16:    Sample particle mean:  $\mu_{\ell}^{\mathcal{B}} \sim \mathcal{N}(\mu_k^{\mathcal{H}}, \Sigma_k^{\mathcal{H}})$ 
17:    Compute cluster-induced velocity:  $\bar{\mathbf{v}}_{\ell} = \mathbf{t}_k + (\mathbf{R}_k - \mathbf{I})(\mu_{\ell}^{\mathcal{B}} - \mu_k^{\mathcal{H}})$ 
18:    Sample particle velocity mean:  $\mathbf{v}_{\ell} \sim \mathcal{N}(\bar{\mathbf{v}}_{\ell}, \sigma_V^2 \mathbf{I})$ 
19:    Sample particle velocity covariance:  $\Sigma_{\ell}^{\mathcal{V}} \sim \mathcal{W}^{-1}(\Psi^{\mathcal{V}}, \nu^{\mathcal{V}})$ 
20:  end for
21:  for  $n = 1$  to  $N$  do
22:    Sample particle assignment:  $z_n^{\mathcal{B}} \sim \text{Cat}(\pi^{\mathcal{B}})$ 
23:    Let  $\ell = z_n^{\mathcal{B}}$ 
24:    Sample data point position:  $\mathbf{x}_n \sim \mathcal{N}(\mu_{\ell}^{\mathcal{B}}, \Sigma_{\ell}^{\mathcal{B}})$ 
25:    Sample data point velocity:  $\mathbf{v}_n \sim \mathcal{N}(\mathbf{v}_{\ell}, \Sigma_{\ell}^{\mathcal{V}})$ 
26:  end for

```

Clusters: Globally coherent, rigidly moving particle groups

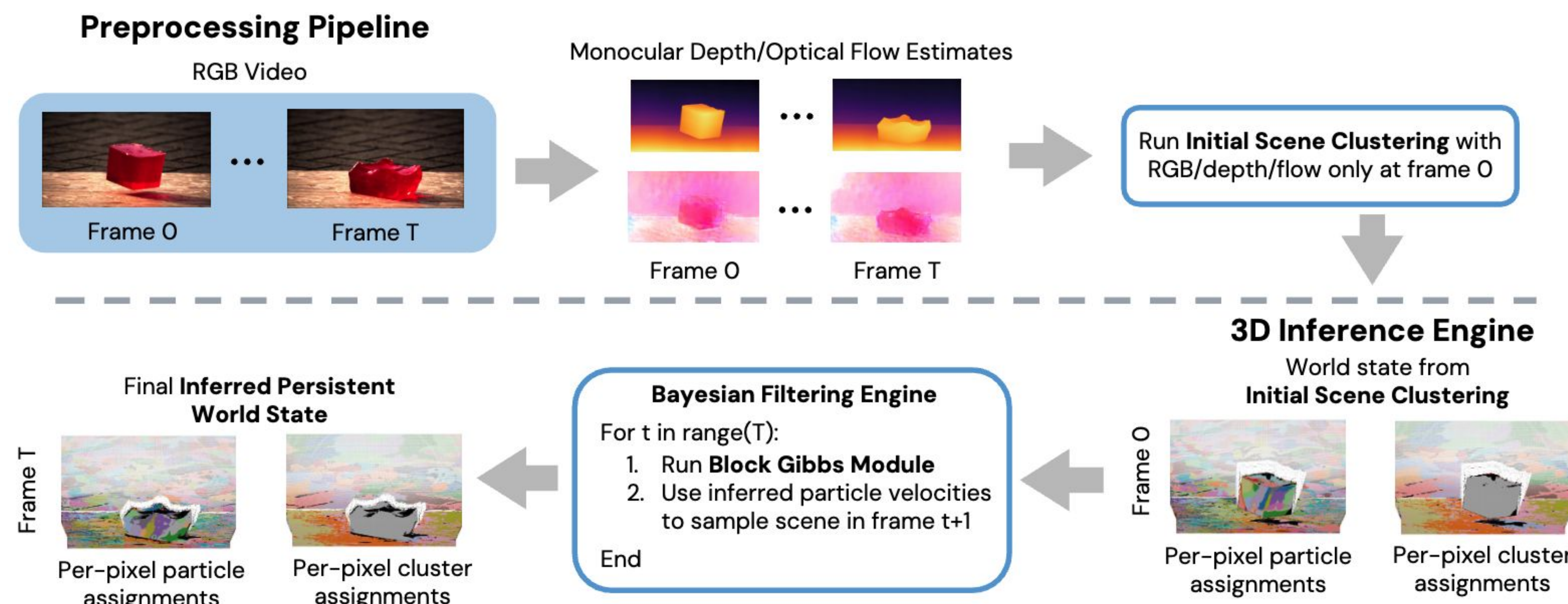
Particles: Localized components assigned to clusters, capturing spatial variability and approximate local translational motion

Data Points: point-level positions and instantaneous velocity

*Equal contribution

Approximate Inference

GenParticles performs sequential inference using a blocked Gibbs sampler that propagates *inferred particles from the last time-step* forward, assigns data points by spatial proximity, updates particle parameters, then infers cluster parameters in a strict order to ground higher-level structure while preserving object consistency with fixed assignments and covariances.



Object Persistence in Video

- We evaluate **GenParticles** on 33 single-object DAVIS videos by measuring object particle persistence, which is the average percentage of particles remaining inside the segmentation mask over time relative to their initialization, running on a single NVIDIA L4 24GB GPU.
- We compare against two *state-of-the-art* particle video baselines, CoTracker3 and SpaTracker, which run offline using full video access and structured 25×25 particle grids to represent objects, constrained by the same GPU memory limits for fair comparison.

DAVIS Video	GenParticles (Ours)	CoTracker3	SpaTracker	DAVIS Video	GenParticles (Ours)	CoTracker3	SpaTracker
boat	100.00 ± 0.00	90.69	92.14	bus	93.33 ± 0.52	82.16	82.79
car-turn	100.00 ± 0.00	97.15	99.89	dance-jump	93.09 ± 1.80	88.46	85.21
drift-chicane	100.00 ± 0.00	74.84	51.28	dog	92.78 ± 2.45	96.34	97.09
car-roundabout	99.79 ± 0.25	96.50	97.78	dance-twirl	92.64 ± 1.12	83.87	90.83
flamingo	99.53 ± 0.42	80.83	90.98	mallard-water	92.56 ± 1.00	97.51	94.28
breakdance-flare	99.51 ± 0.25	82.19	98.74	goat	92.26 ± 2.82	88.73	88.61
camel	99.40 ± 0.53	93.45	96.34	koala	91.46 ± 1.03	63.95	57.65
cows	99.28 ± 0.89	94.25	93.02	lucia	87.73 ± 2.18	93.77	97.72
rallye	98.96 ± 0.66	100.00	80.00	dog-agility	84.78 ± 1.50	78.65	74.20
rollerblade	97.96 ± 1.92	95.78	94.64	libby	83.35 ± 3.19	77.78	79.23
rhino	97.50 ± 0.39	88.19	87.87	parkour	78.57 ± 1.72	84.96	76.00
blackswan	96.89 ± 0.69	99.91	100.00	mallard-fly	69.23 ± 3.03	71.01	85.27
bear	96.72 ± 1.29	95.87	94.28	drift-turn	64.75 ± 2.76	96.46	91.47
elephant	96.71 ± 1.11	89.82	90.14	drift-straight	63.75 ± 12.47	94.17	74.00
breakdance	96.31 ± 1.21	69.06	94.83	varanus-cage	51.93 ± 9.78	67.19	64.93
hike	94.36 ± 2.56	92.67	94.46	soccerball	11.85 ± 0.22	88.94	87.22
car-shadow	94.11 ± 0.63	96.08	97.28	Median Accuracy	94.11	89.82	90.98

- GenParticles outperforms SpaTracker and CoTracker3 on **20 of 33** sequences, showing strong tracking but limited handling of occlusion and out-of-frame motion due to the absence of a dynamics prior.

Blocked Gibbs Sampling

Latent Assignments

$$p(z_n^{\mathcal{B}} = \ell | \dots) \propto \pi_{\ell}^{\mathcal{B}} \cdot \mathcal{N}(\mathbf{x}_n | \mu_{\ell}^{\mathcal{B}}, \Sigma_{\ell}^{\mathcal{B}}) \cdot \mathcal{N}(\mathbf{v}_n | \mathbf{v}_{\ell}, \Sigma_{\ell}^{\mathcal{V}})$$

$$p(z_{\ell}^{\mathcal{H}} = k | \dots) \propto \pi_k^{\mathcal{H}} \cdot \mathcal{N}(\mu_{\ell}^{\mathcal{B}} | \mu_k^{\mathcal{H}}, \Sigma_k^{\mathcal{H}}) \cdot \mathcal{N}(\mathbf{v}_{\ell} | \bar{\mathbf{v}}_{\ell, k}, \sigma_V^2 \mathbf{I})$$

Mixture weights are updated via *Dirichlet-Categorical Conjugacy*

Particle Updates

Velocity and Spatial Covariances via *Normal Inverse-Wishart Conjugacy*

$$\Sigma_{\ell}^{\mathcal{B}}, \Sigma_{\ell}^{\mathcal{V}}$$

Velocities and Positions via *Normal-Normal Conjugacy*

$$\mu_{\ell}^{\mathcal{B}}, \mathbf{v}_{\ell}$$

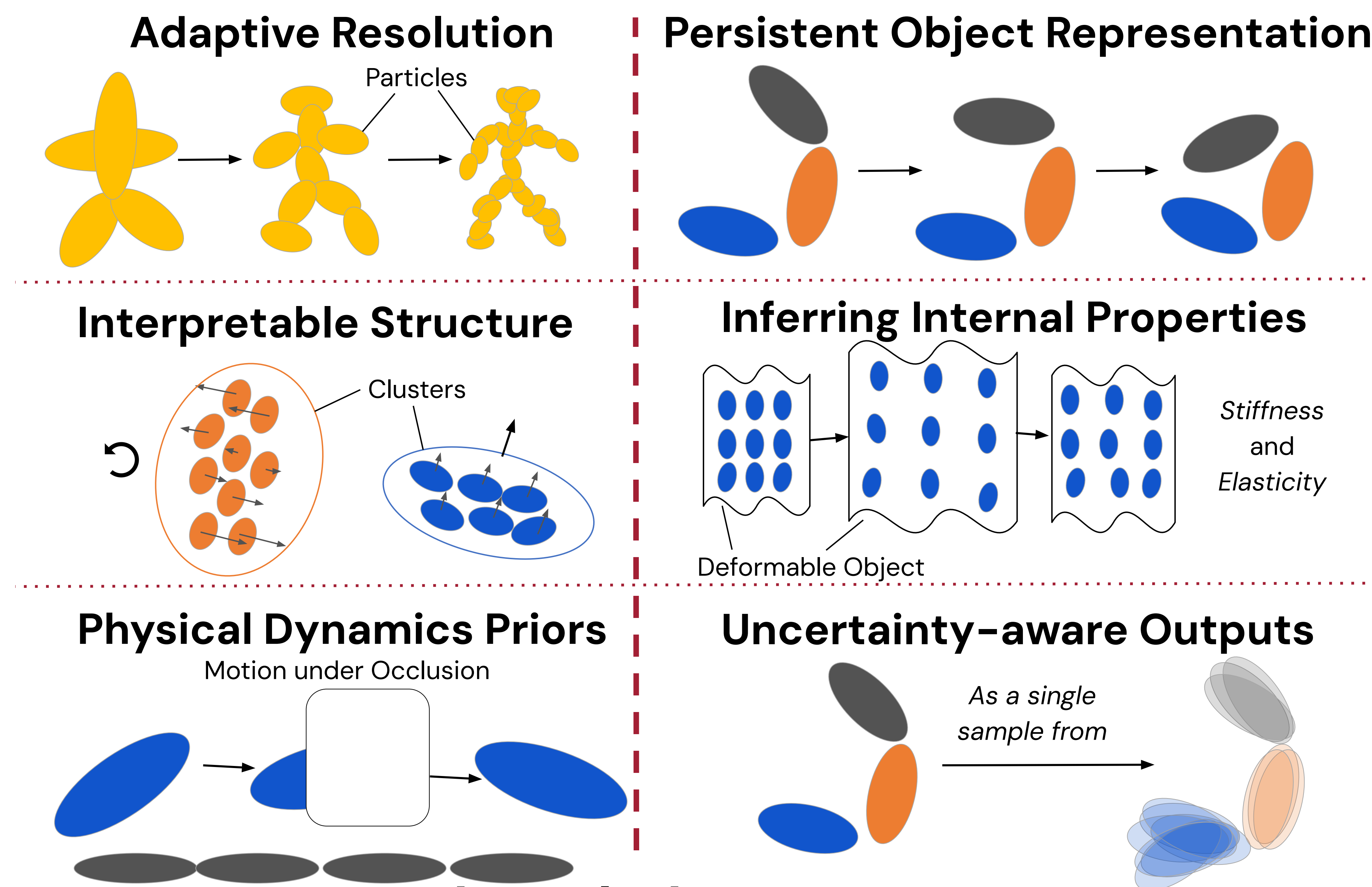
Cluster Updates

Rotations and Translations are sampled via full enumeration

$$\mathbf{R}_k, \mathbf{t}_k$$

All inference moves are functionally parallelized and matrix operations are vectorized via **GenJAX** and **JAX**

Future Applications for Robotic Manipulation



Acknowledgements

This work was supported in part by CoCoSys, one of seven centers in JUMP 2.0, Semiconductor Research Corporation (SRC) program sponsored by DARPA