

Physical Scene Understanding with Compositional Structure

Jiajun Wu (MIT CSAIL)

GAMES Webinar, Jan 10, 2019



Advisors and Collaborators



Josh Tenenbaum



Bill Freeman



Antonio Torralba



Pushmeet Kohli

Tianfan Xue
Ilker Yildirim
Joseph Lim
Jun-Yan Zhu
Katie Bouman
Chengkai Zhang
Yunzhu Li
Harry Hsu
Jiancheng Liu

Andrew Owens
Xiuming Zhang
Zhoutong Zhang
Wenzhen Yuan
Erika Lu
Zhijian Liu
Chen Sun
Yuandong Tian
Andrew Spielberg

Michael Janner
Anurag Ajay
Nima Fazeli
Maria Bauza
James Traer
Hongyi Zhang
Yuanming Hu
Tejas Kulkarni
Kevin Ellis

David Zheng
Shaoxiong Wang
Xingyuan Sun
Qiujia Li
Zhengjia Huang
Yifan Wang
Shunyu Yao
Andrew Luo
Yonglong Tian

Josh McDermott
Leslie Kaelbling
Alberto Rodriguez
Kevin Murphy
Ted Adelson
Russ Tedrake
Daniela Rus
Wojciech Matusik
Daniel Ritchie

What can we learn from this video?



Human Physical Scene Understanding

What can we see in this video?

I. Scene structure (**perception**)

- Object appearance (geometry, texture)
- Physical properties (e.g., mass)

II. Interactions and events (**physics**)

- Collision, rolling, etc.

III. Concepts and regularity (**reasoning**)

- Balls can roll, but not blocks
- Blocks are of the same size and shape
- Blocks are lined up in a row



sphere
orange
plastic



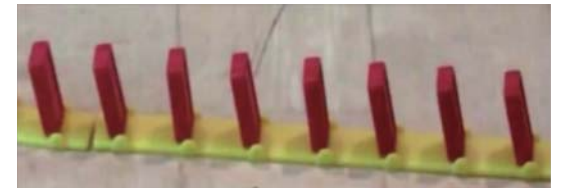
hedge
white/yellow
wooden



collisions



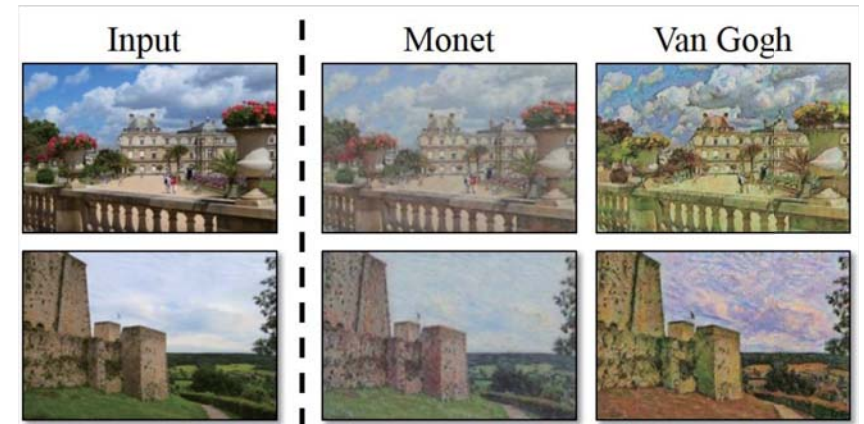
rolling



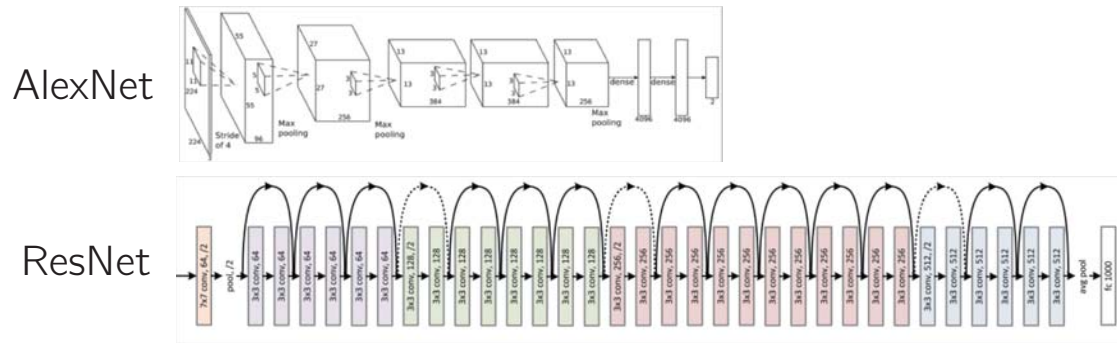
Current Machine Scene Understanding



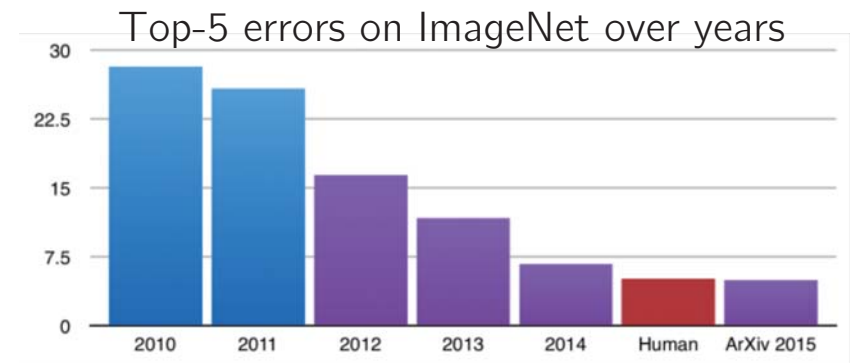
Recognition



Synthesis



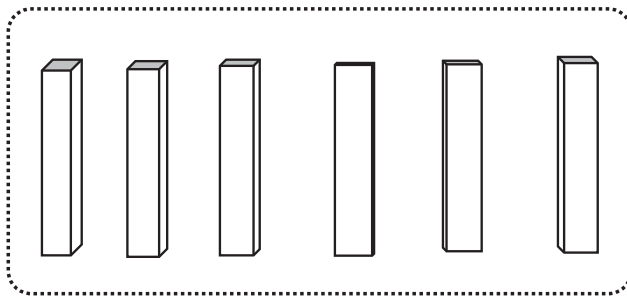
Deep Learning Models



Performance

Image Credit: DeepLab, Chen et al., 2018; CycleGAN, Zhu et al., 2017

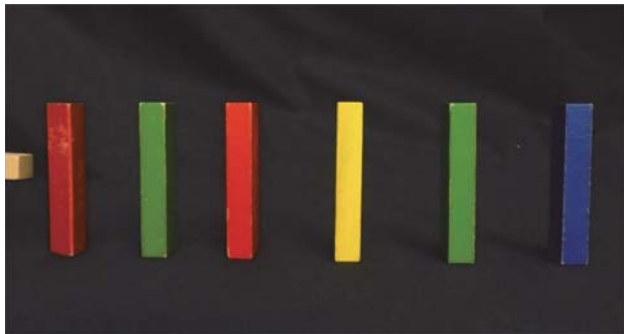
Modeling the Physical World



World State (t-1)

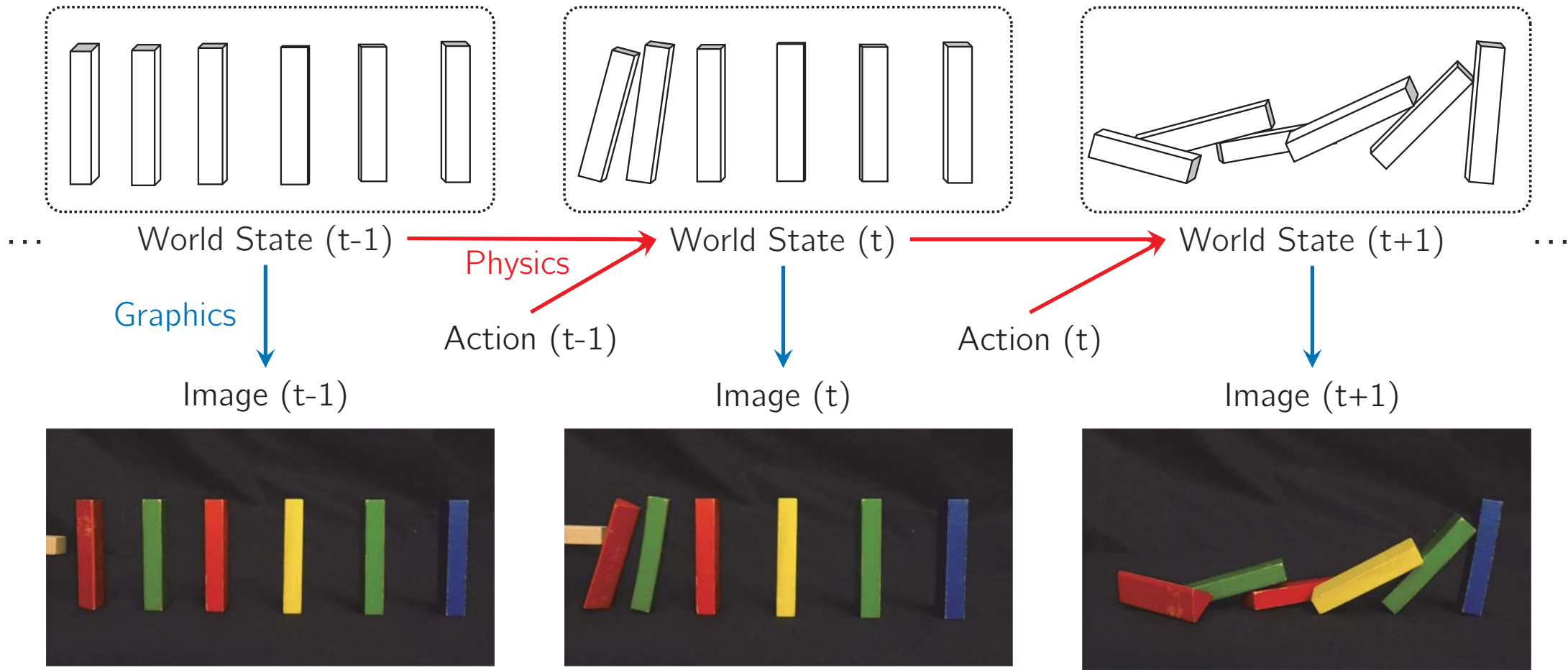
Graphics

Image (t-1)

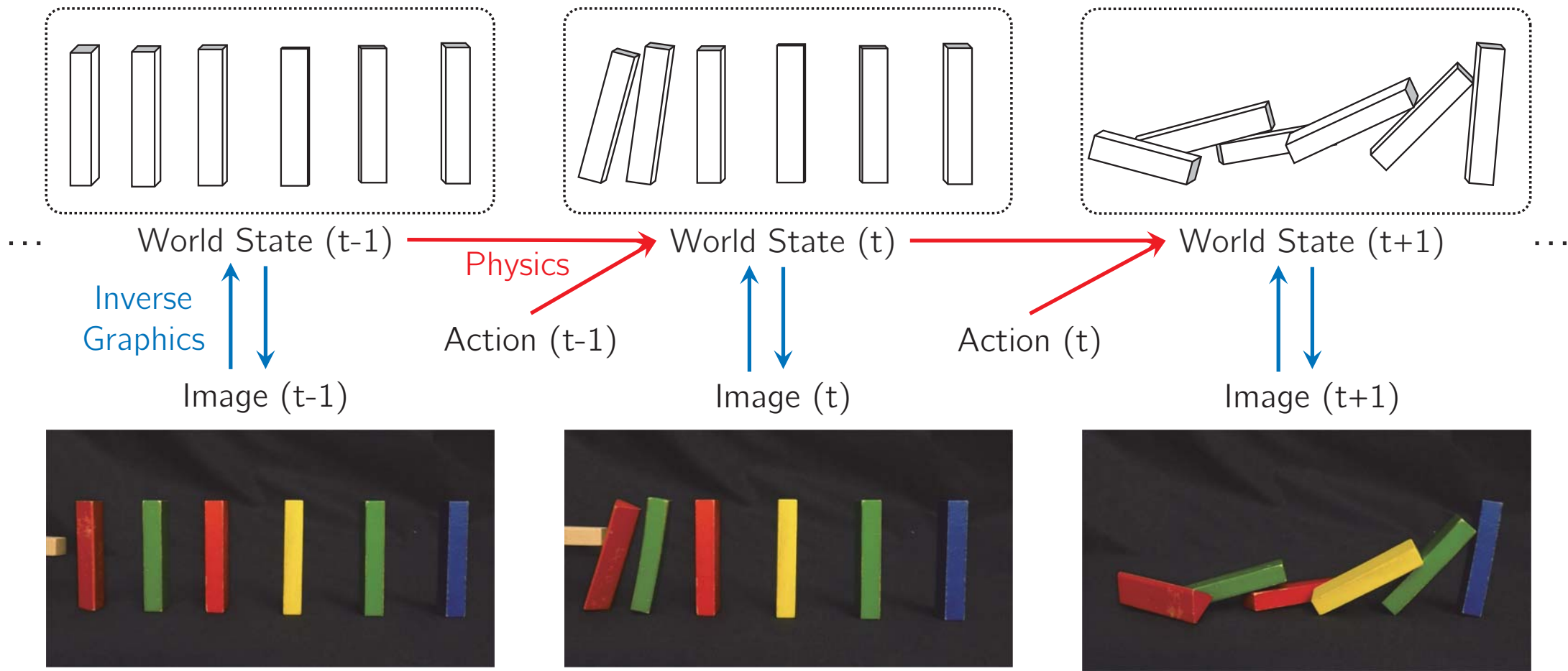


- Object Intrinsic
 - Geometry
 - Physical properties
- Object Extrinsic
 - Position
 - Velocity
- Scene Descriptions
 - Lighting
 - Camera parameters

Modeling the Physical World

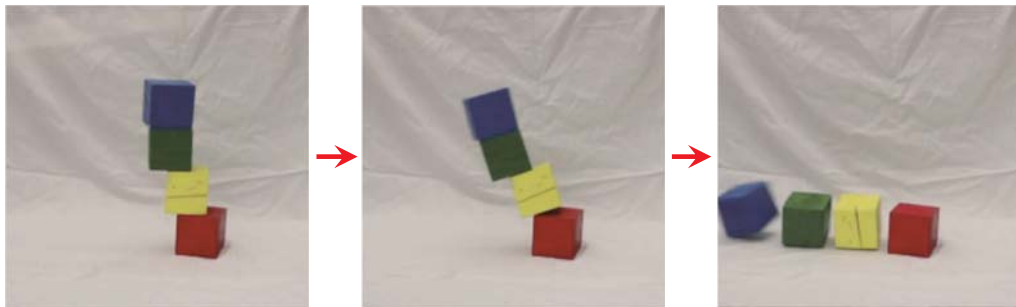
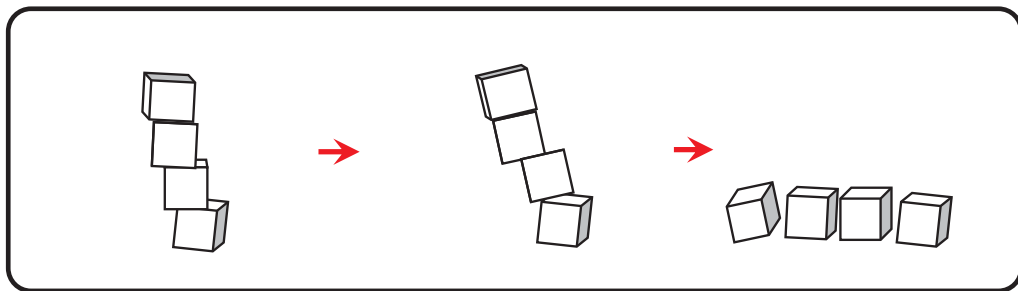


Modeling the Physical World



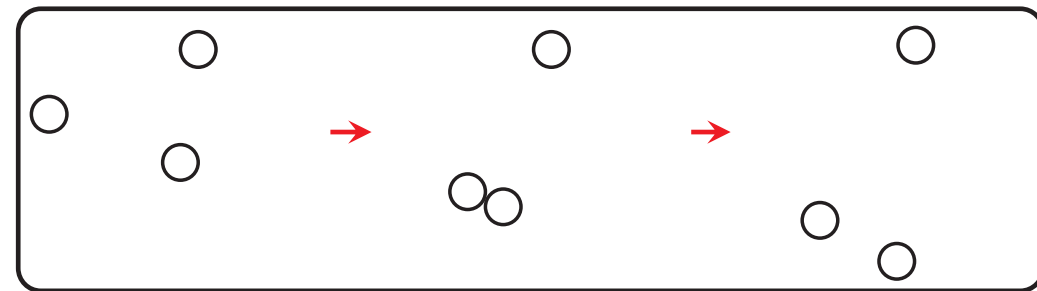
Physical World Representations are Universal

World States



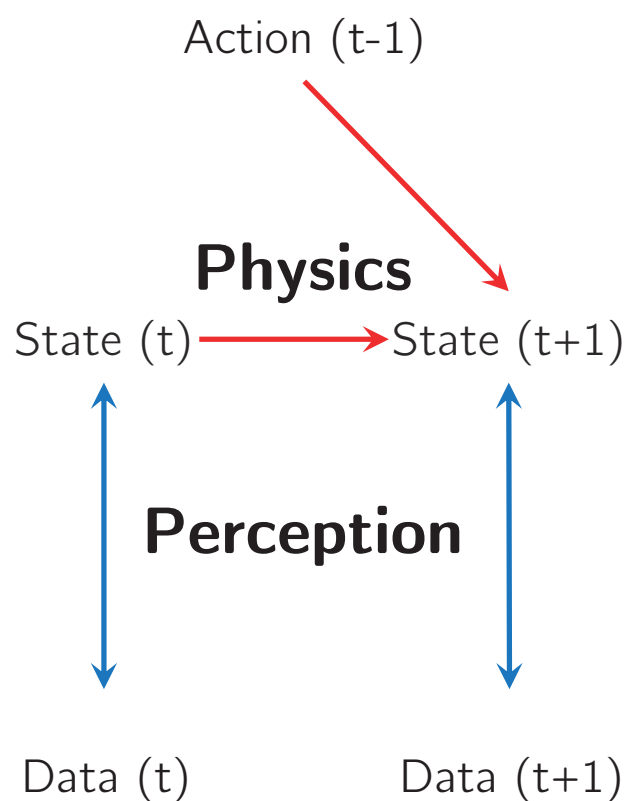
Visual Observation

World States



Visual Observation

Cognitive Science Meets Machine Scene Understanding



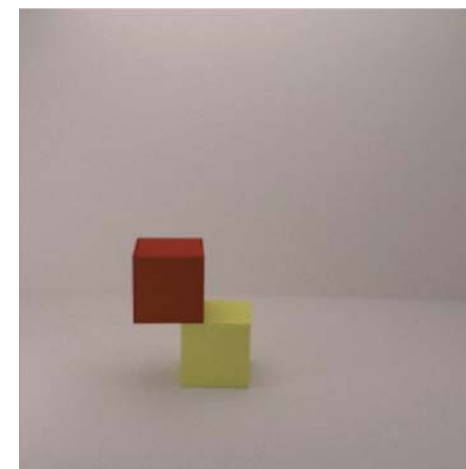
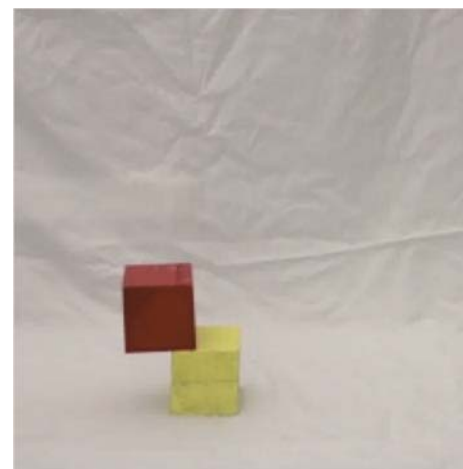
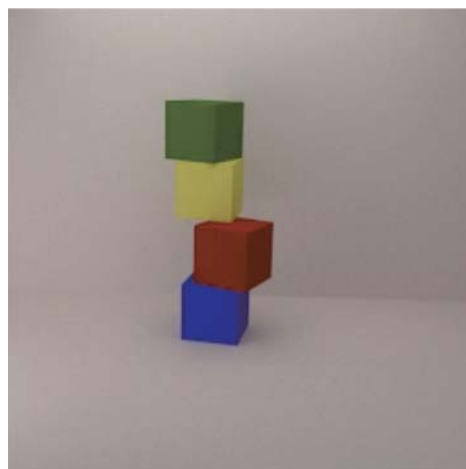
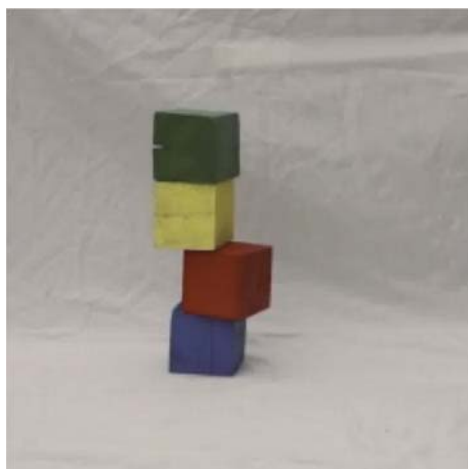
Causal structure and cognitive science insights provide guidance on building machine scene understanding models:

- When and where to use top-down simulation engines vs. bottom-up neural networks?
- What training targets to use for neural networks?
- What intermediate representations to use?
- What training data to use?

Research in machine intelligence helps to stimulate research in human cognition and neuroscience:

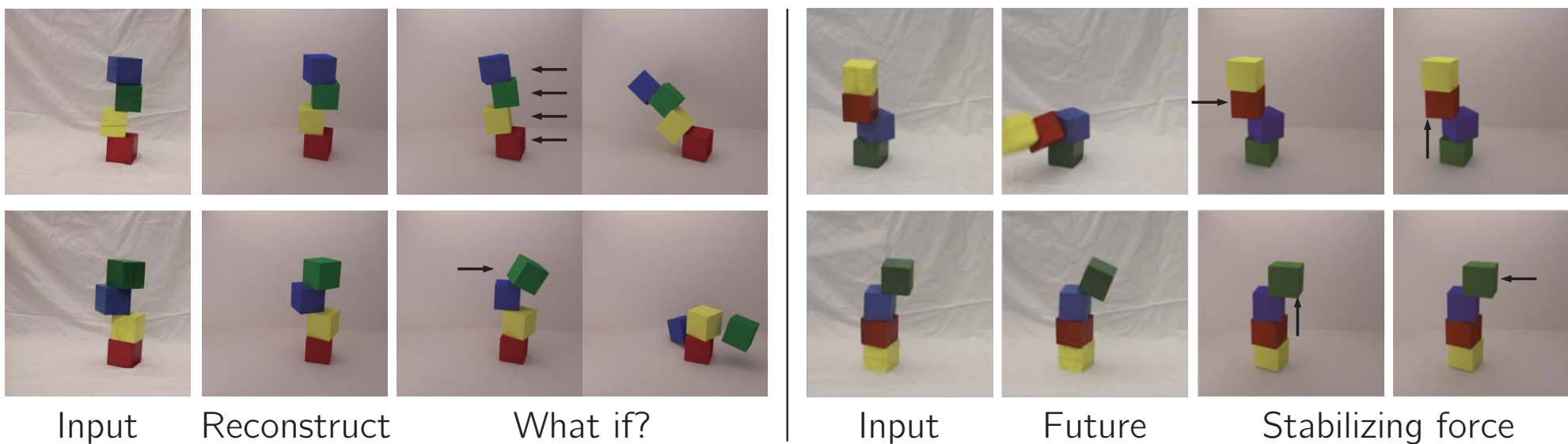
- Computational models for human behaviors;
- Algorithms and representations in the brain.

Learning to See Physics via Visual De-animation



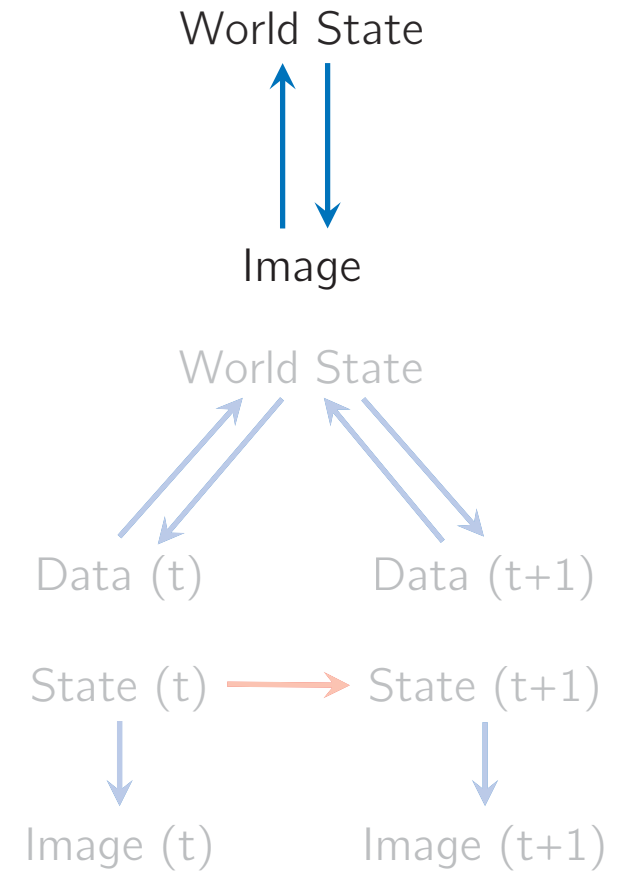
Wu, Lu, Kohli, Freeman, Tenenbaum. NeurIPS'17

Learning to See Physics via Visual De-animation



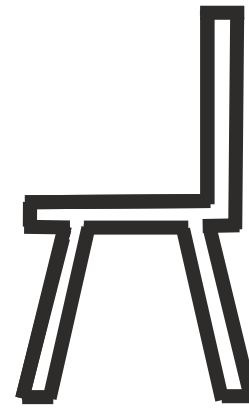
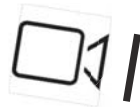
Physical Scene Understanding

- Learning to invert a graphics engine
- Learning to invert a physics engine
- Learning simulation engines themselves



3D Reconstruction

Forward: image formation



Inverse: shape estimation



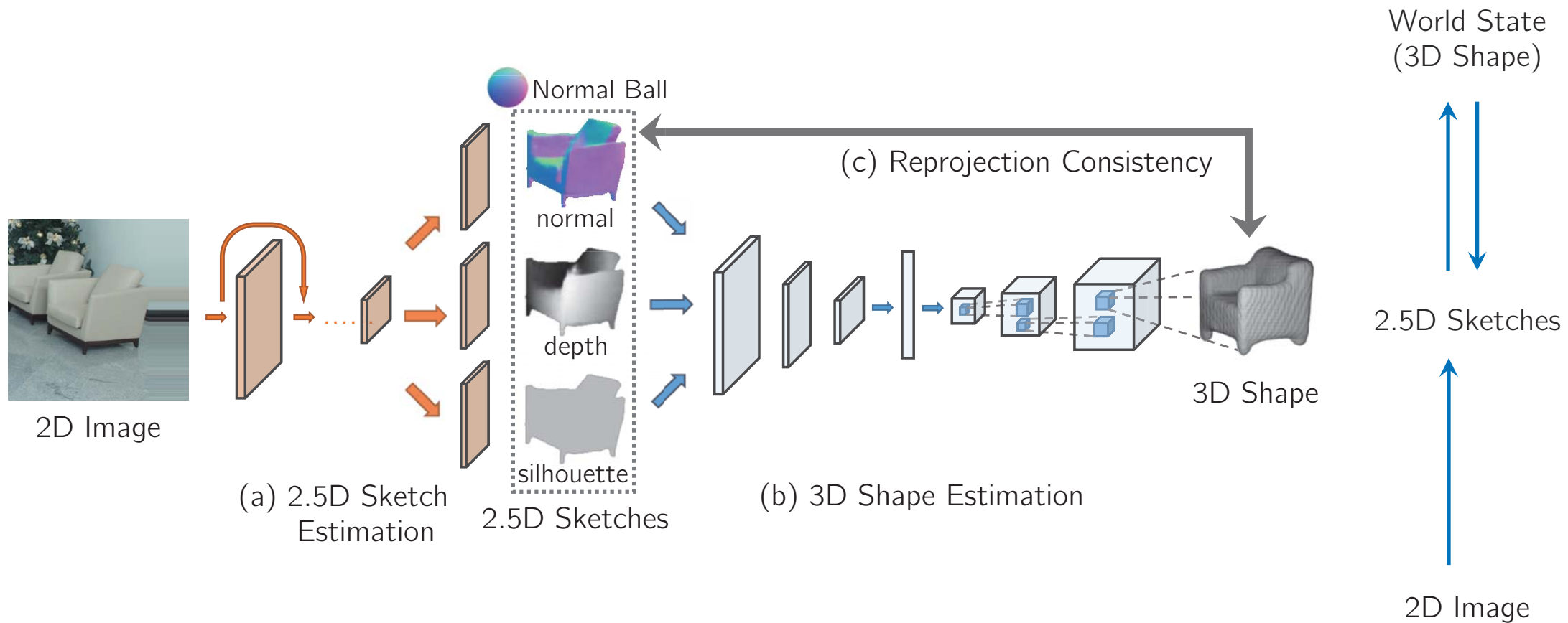
Depth Estimation



Shape Completion

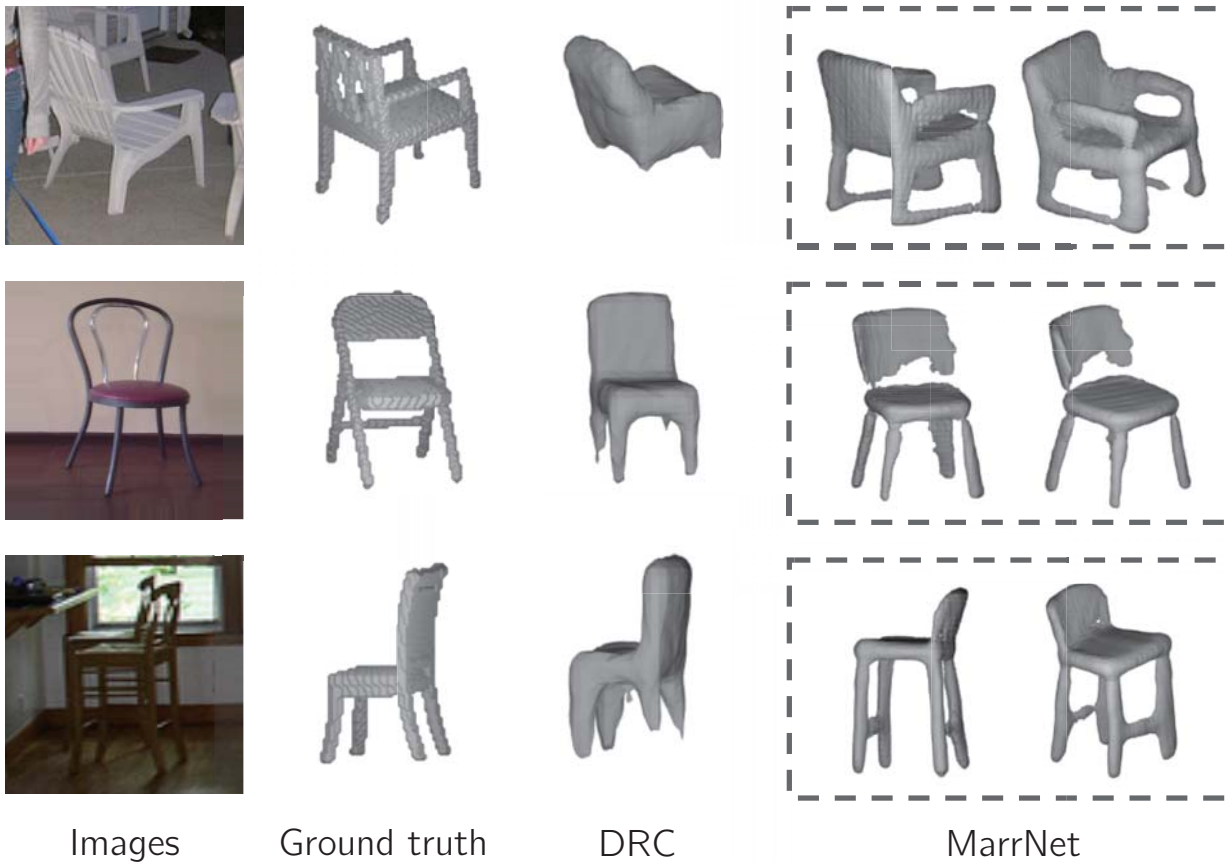


MarrNet: 3D Reconstruction via 2.5D Sketches



Wu*, Wang*, Xue, Sun, Freeman, Tenenbaum. NeurIPS'17

Comparisons



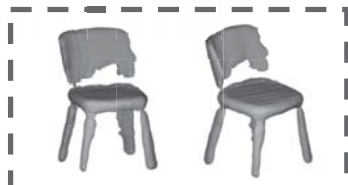
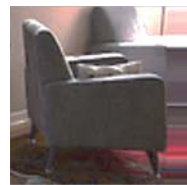
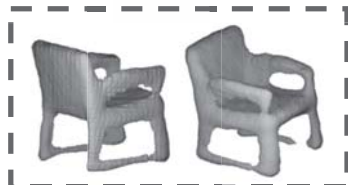
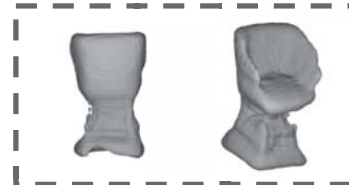
Methods	IoU
DRC 3D [CVPR '17]	0.34
MarrNet	0.38

Intersection over Union (IoU)

	DRC 3D	MarrNet	GT
DRC 3D	50	26	17
MarrNet	74	50	42
GT	83	58	50

Percentages of users that preferred the left approach to the top one

Results on PASCAL 3D+



Images

MarrNet

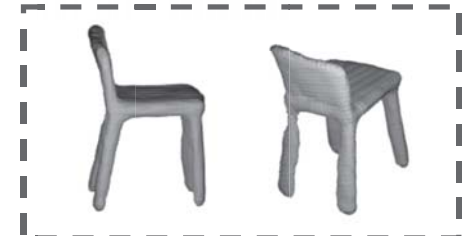
Images

MarrNet

Images

MarrNet

Results on IKEA



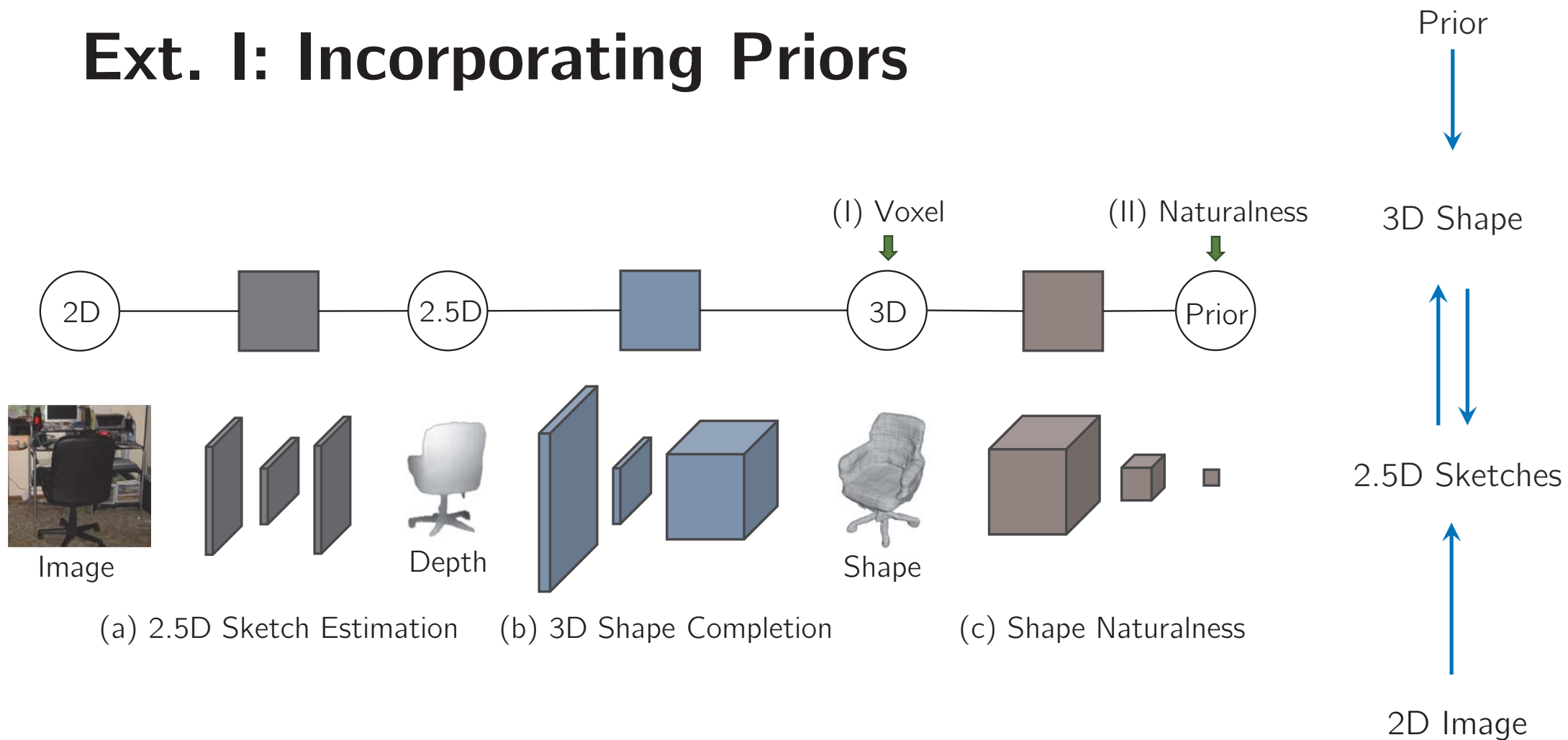
Images Ground truth

MarrNet

Images Ground truth

MarrNet

Ext. I: Incorporating Priors

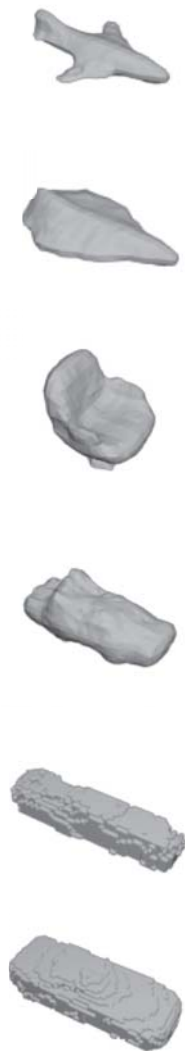




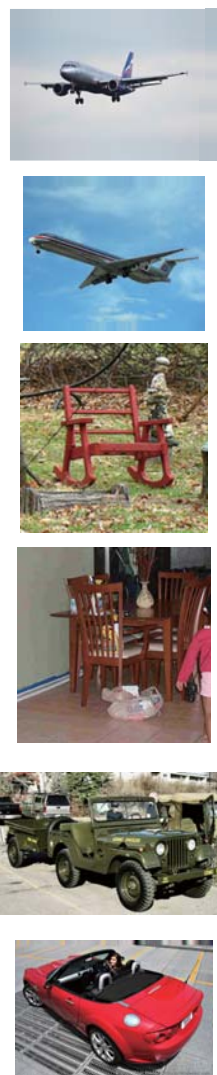
Input



ShapeHD



Best alternative



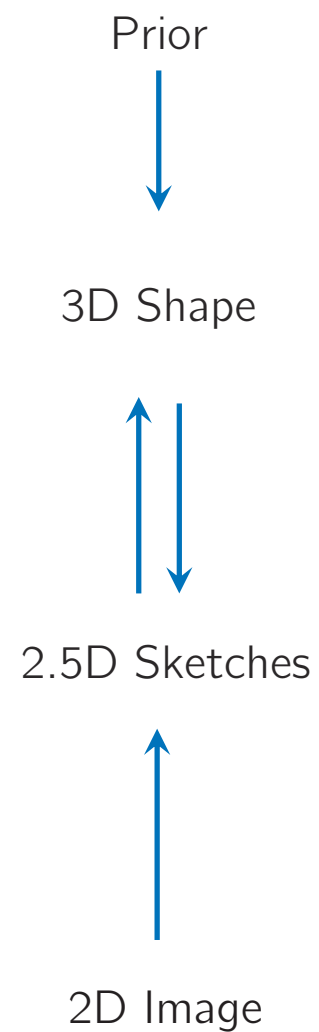
Input



ShapeHD

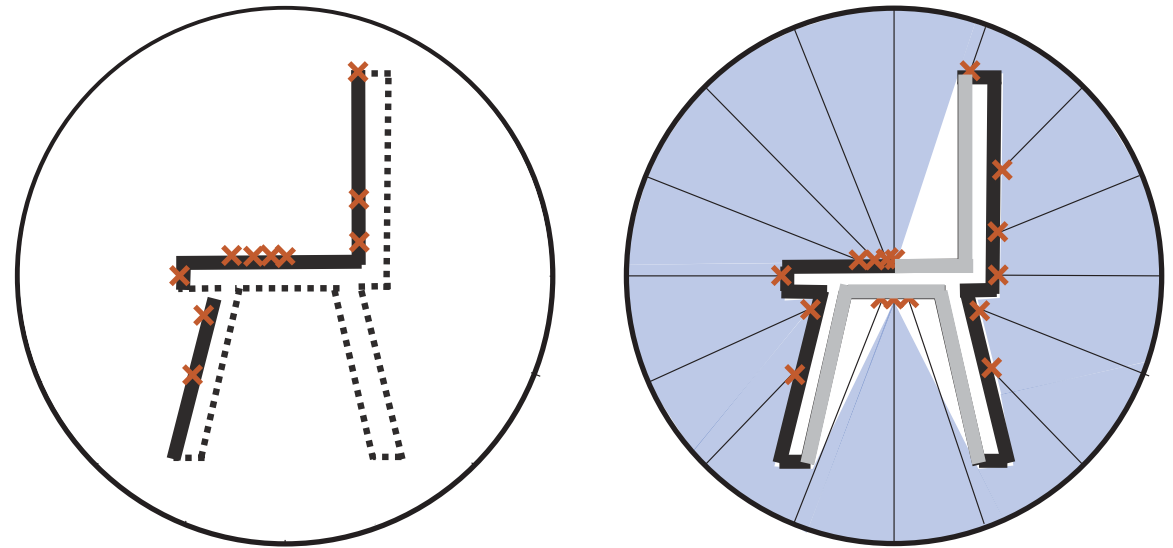
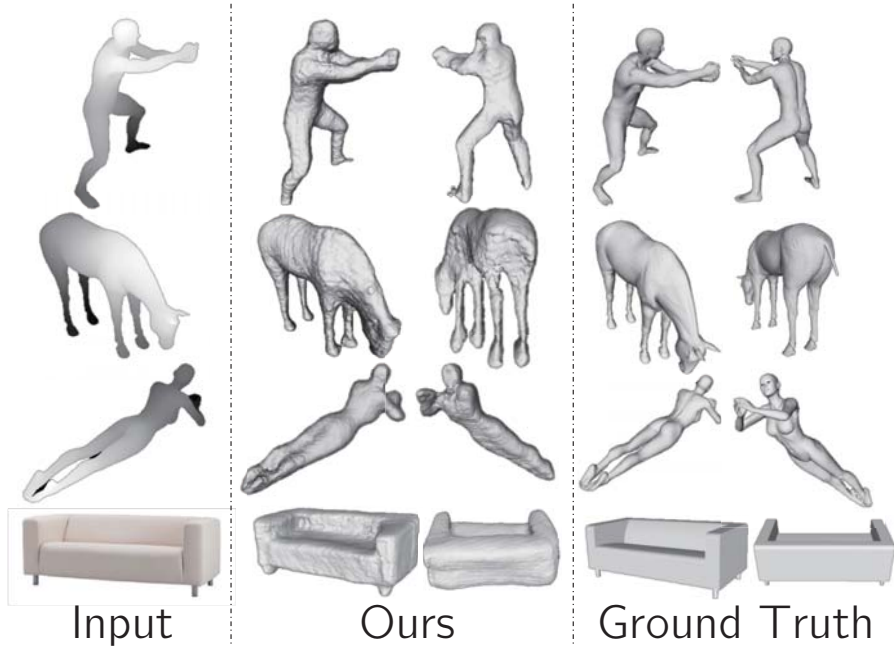
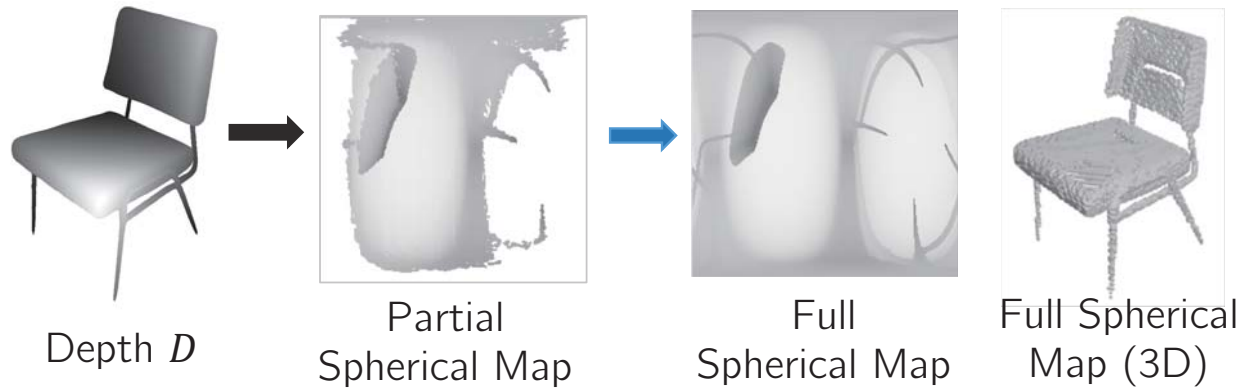
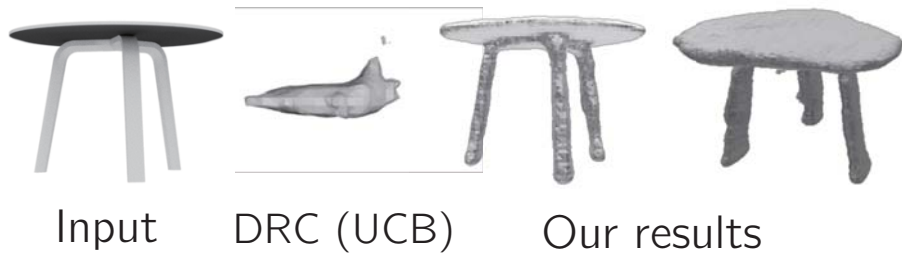


Best alternative



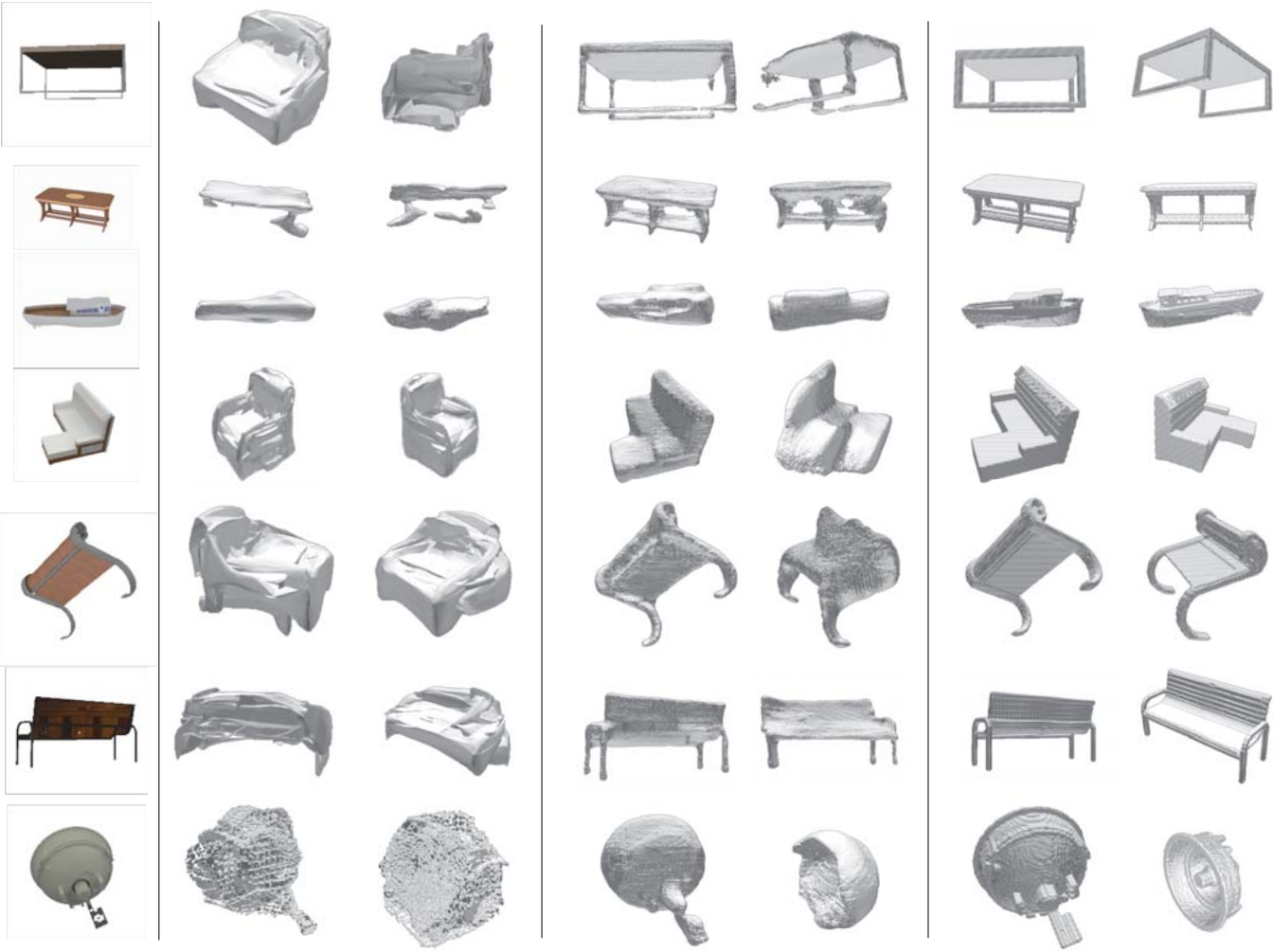
Ext. II: Generalization to Unseen Classes

Trained on chairs, planes, cars
Tested on tables



Zhang*, Zhang*, Zhang, Tenenbaum, Freeman, Wu. NeurIPS'18

Generalization to Novel Classes (Table, Boat, Sofa, Bench, Lamp)



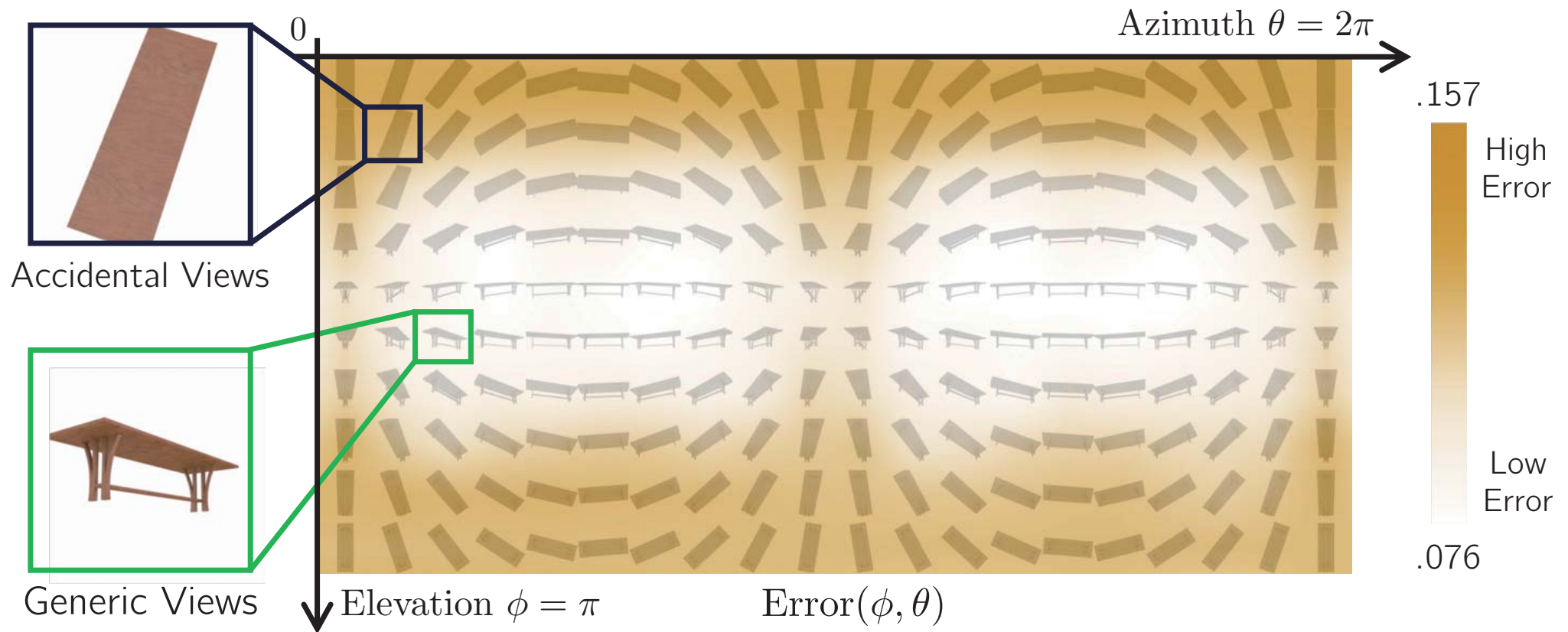
Input

Best Baseline

Ours

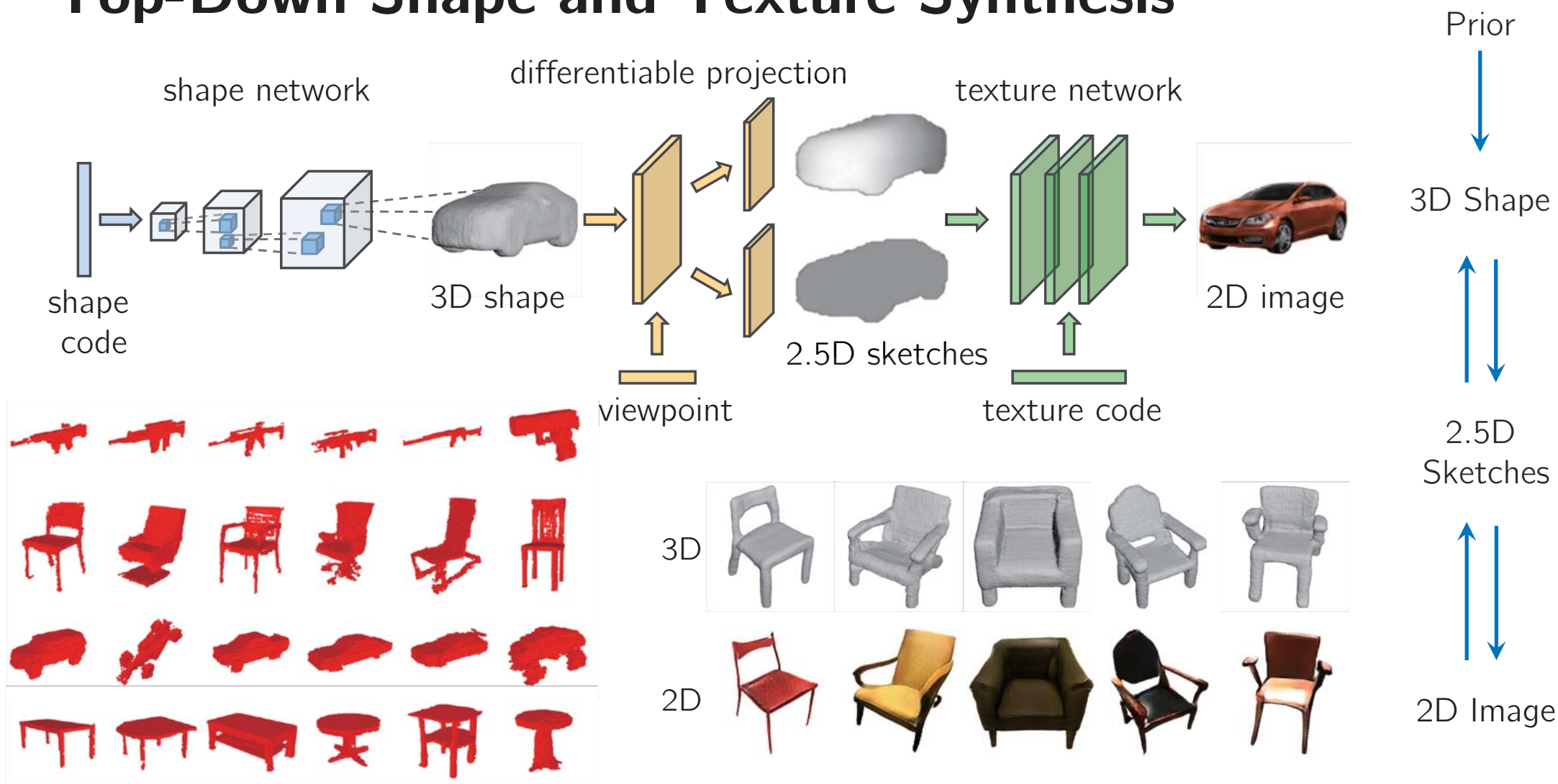
Ground Truth

Canonical Viewpoints in Generalization



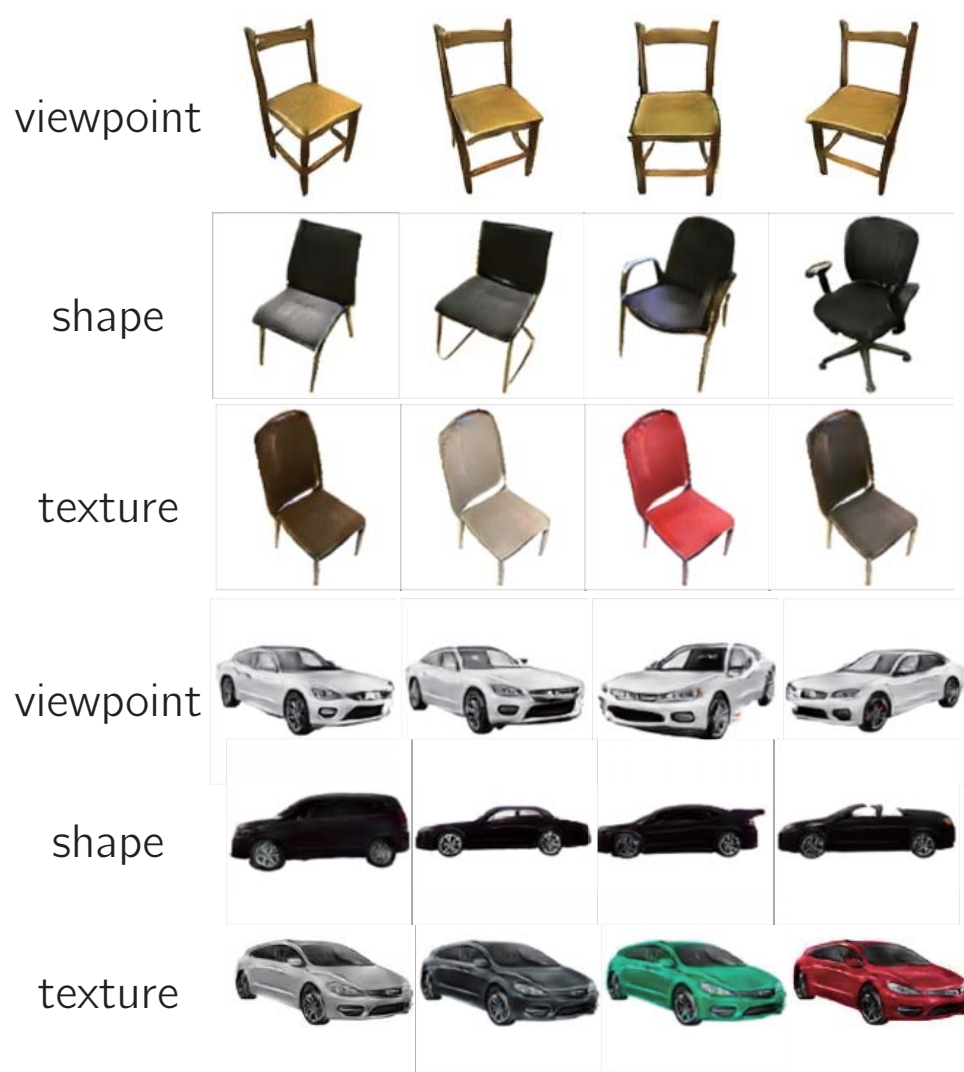
Accidental views \rightarrow large errors

Top-Down Shape and Texture Synthesis



Wu*, Zhang*, Xue, Freeman, Tenenbaum. NeurIPS'16

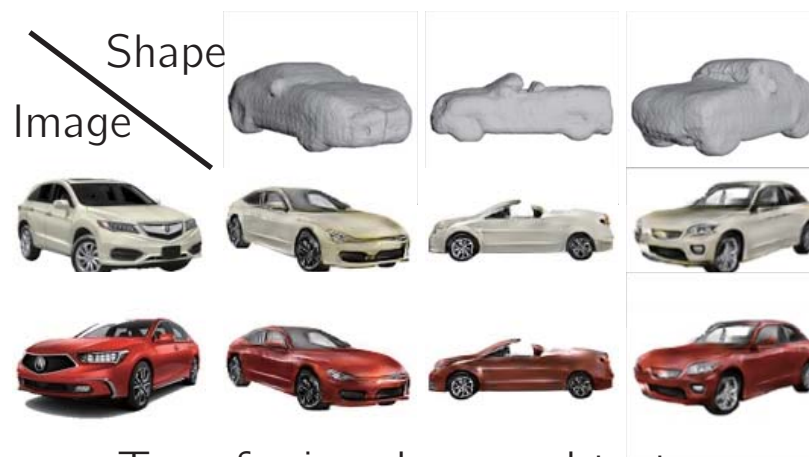
Zhu, Zhang, Zhang, Wu, Torralba, Tenenbaum, Freeman. NeurIPS'18



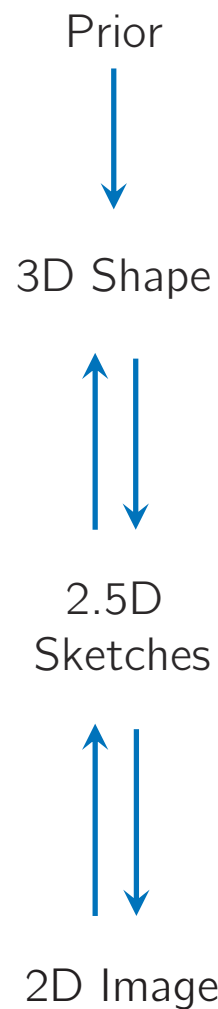
Editing viewpoint, shape, and texture



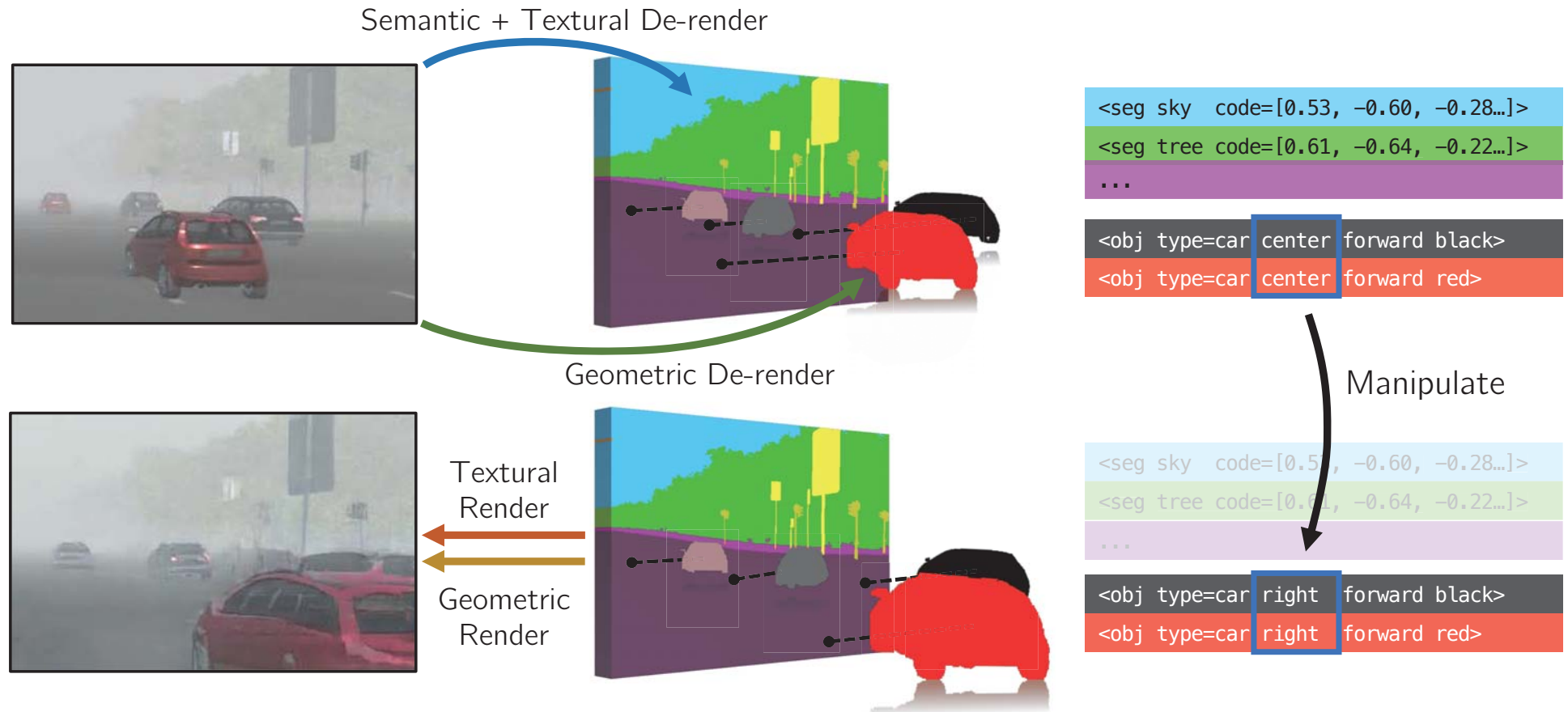
Interpolation in the latent space



Transferring shape and texture



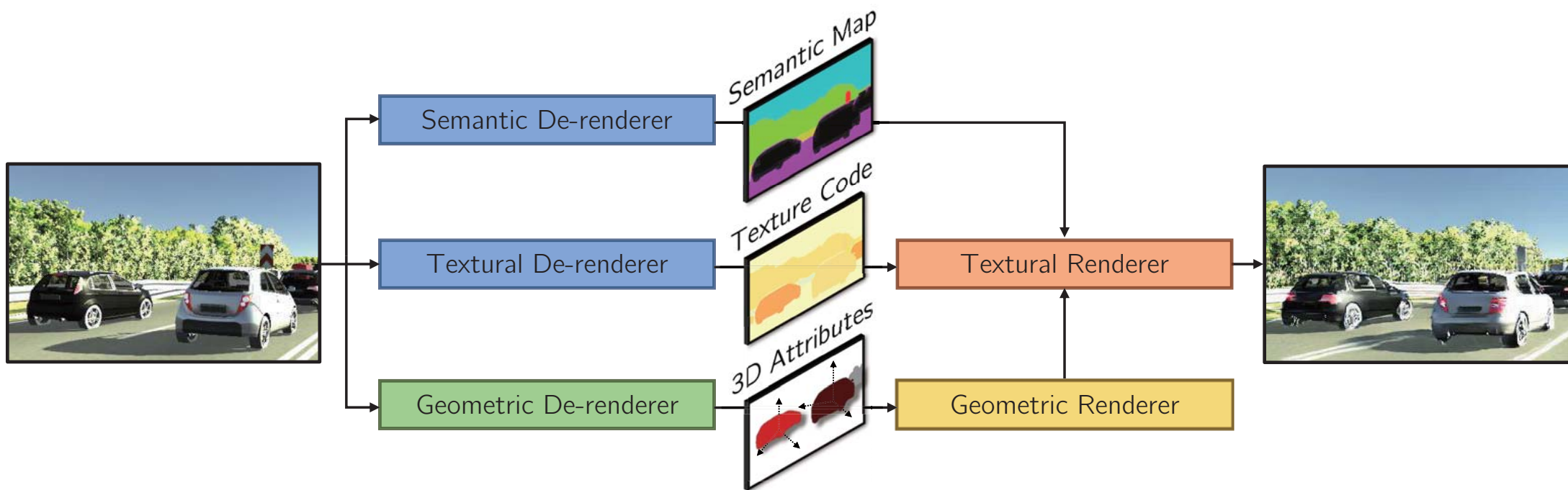
Ext. IV: Extension to Scenes



Goal: Recovering a structured, 3D-aware scene representation.

The structured representation allows re-rendering and editing the image.

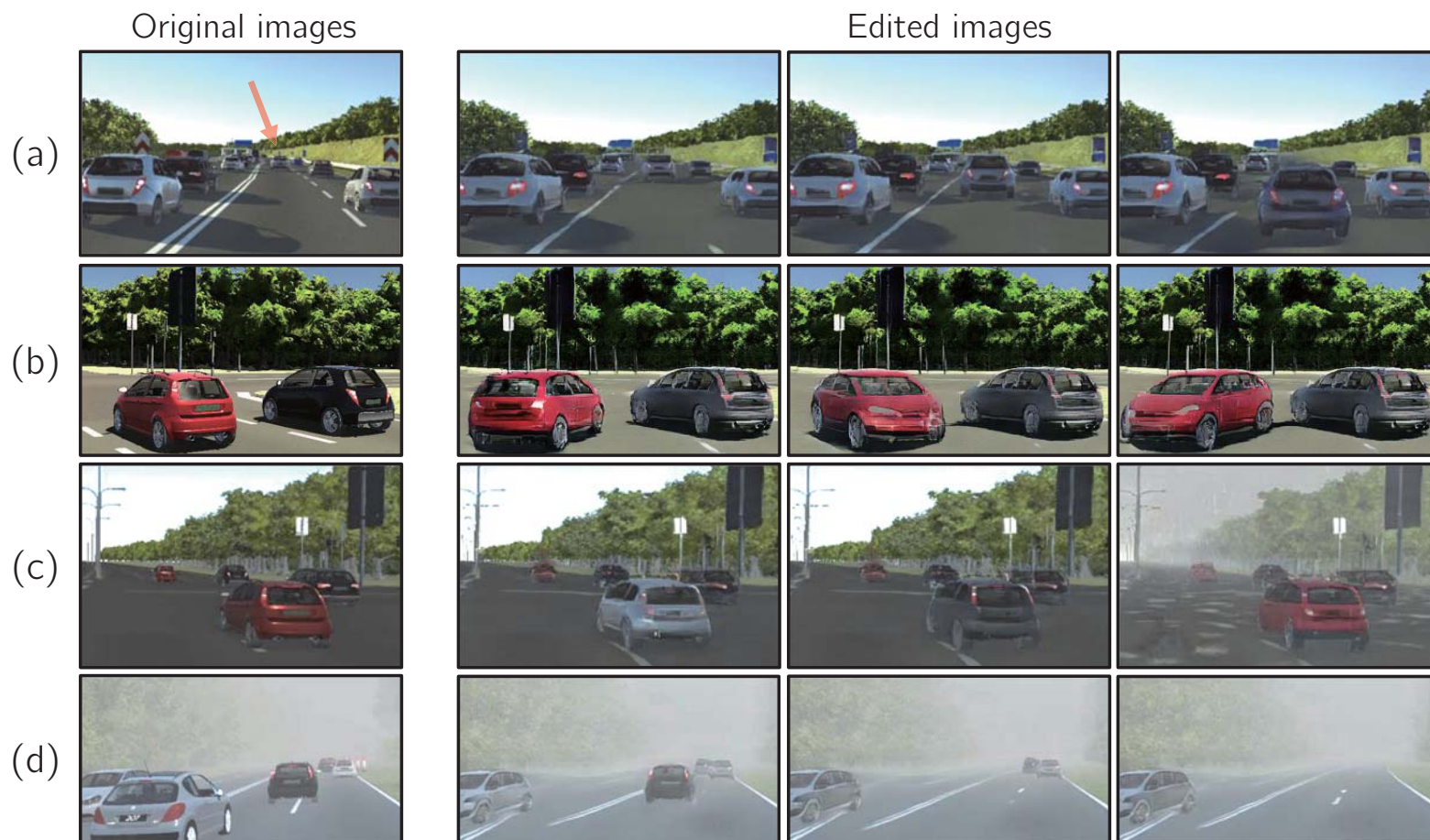
3D Disentangled Scene Representation



Disentangled model for the scene's semantics, texture, and object geometry and 6DOF pose.

Yao*, Hsu*, Zhu, Wu, Torralba, Freeman, Tenenbaum. NeurIPS'18

Image Editing on Virtual KITT



Yao*, Hsu*, Zhu, Wu, Torralba, Freeman, Tenenbaum. NeurIPS'18

Image Editing on CityScapes (Real Images)

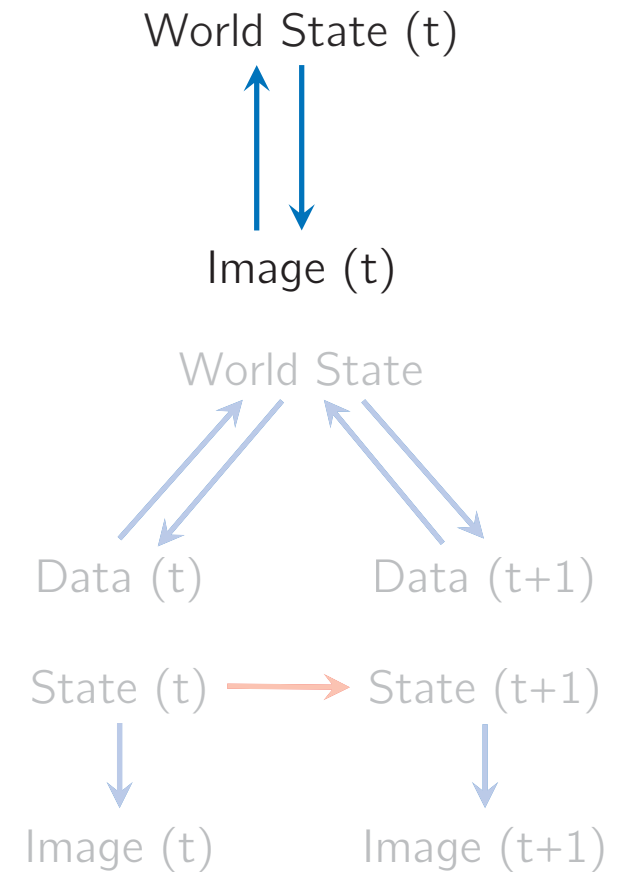
Original images

Edited images



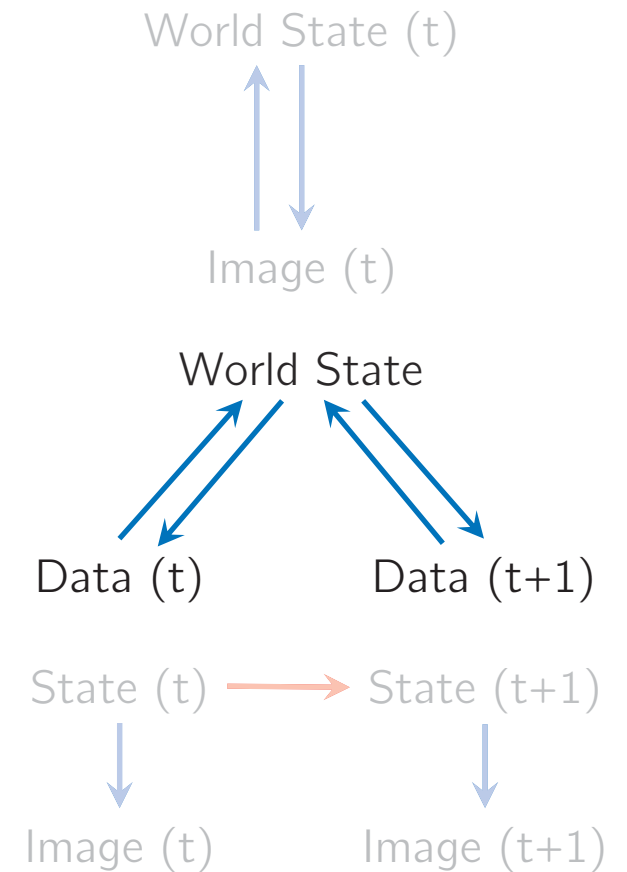
Physical Scene Understanding

- Learning to invert a graphics engine
 - Inferring fine object geometry
 - Learning structured shape representations (shape + texture)
 - Beyond single object, learning scene representations
- Learning to invert a physics engine
- Learning simulation engines themselves

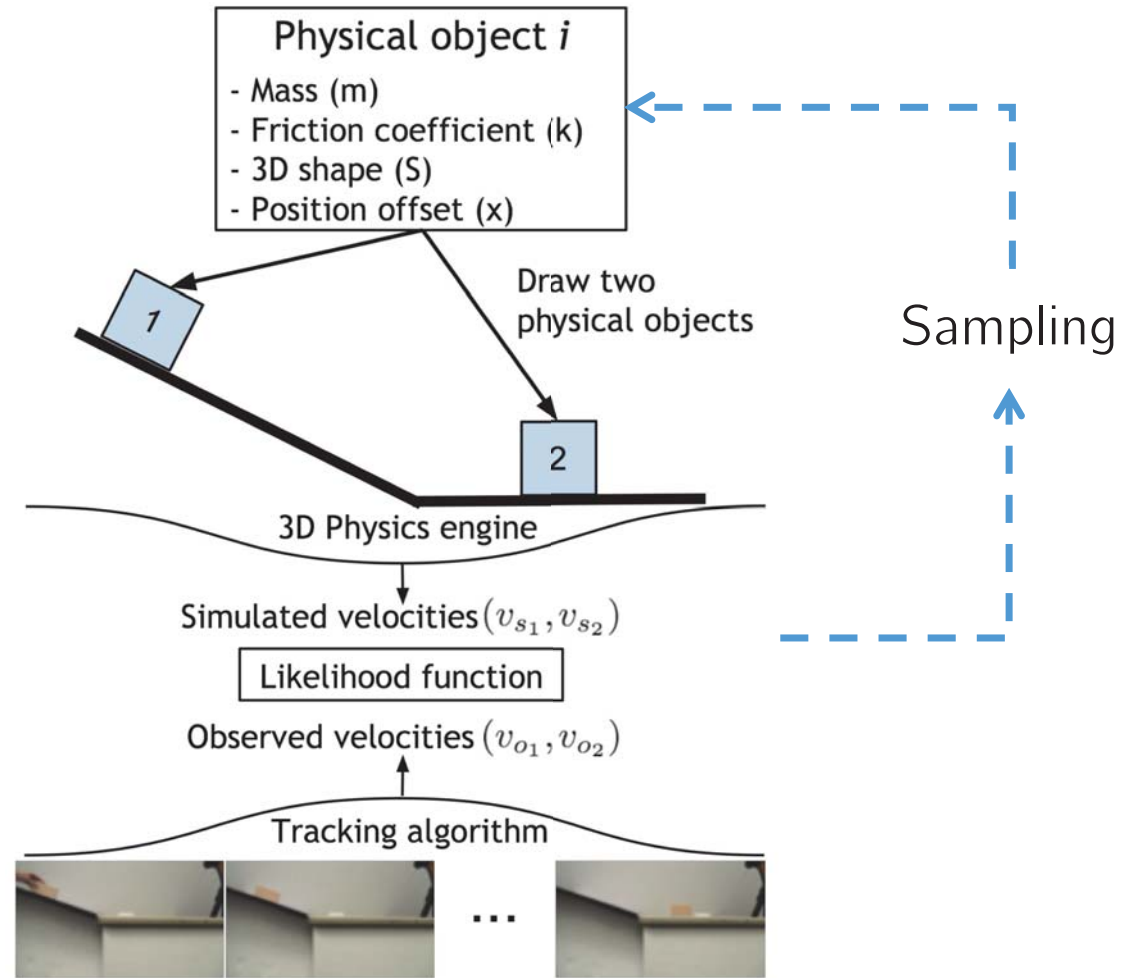
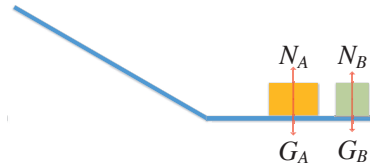
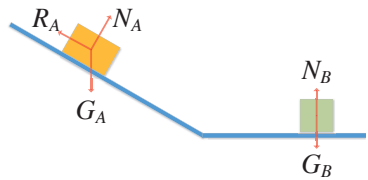


Physical Scene Understanding

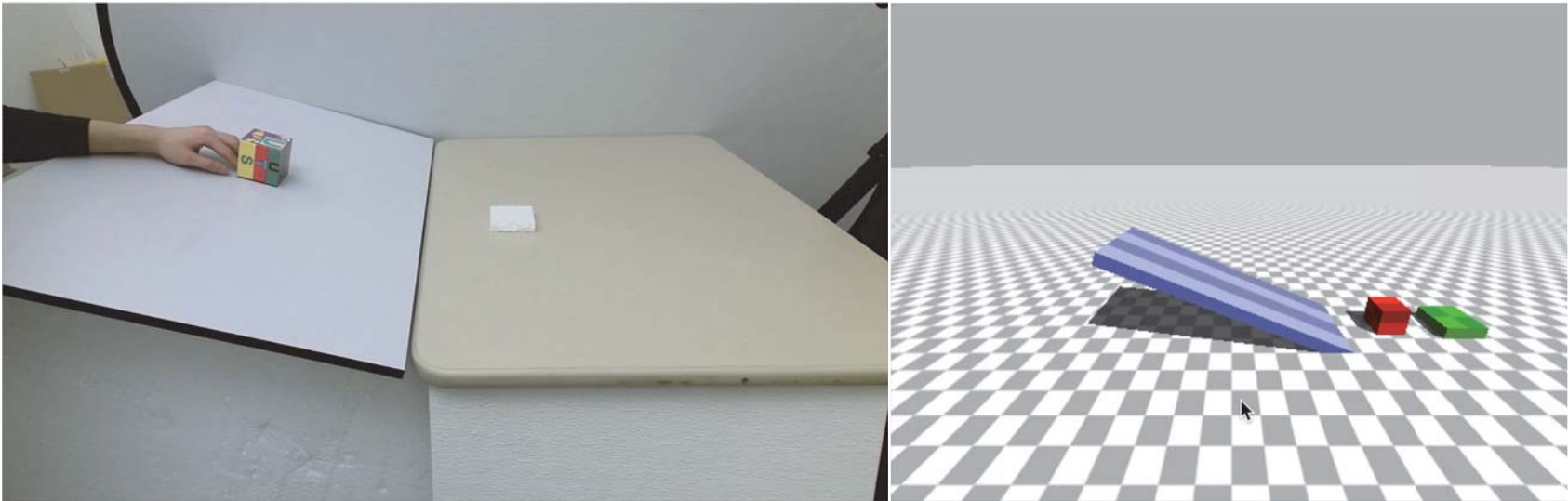
- Learning to invert a graphics engine
 - Inferring fine object geometry
 - Learning structured shape representations (shape + texture)
 - Beyond single object, learning scene representations
- Learning to invert a physics engine
- Learning simulation engines themselves



Galileo



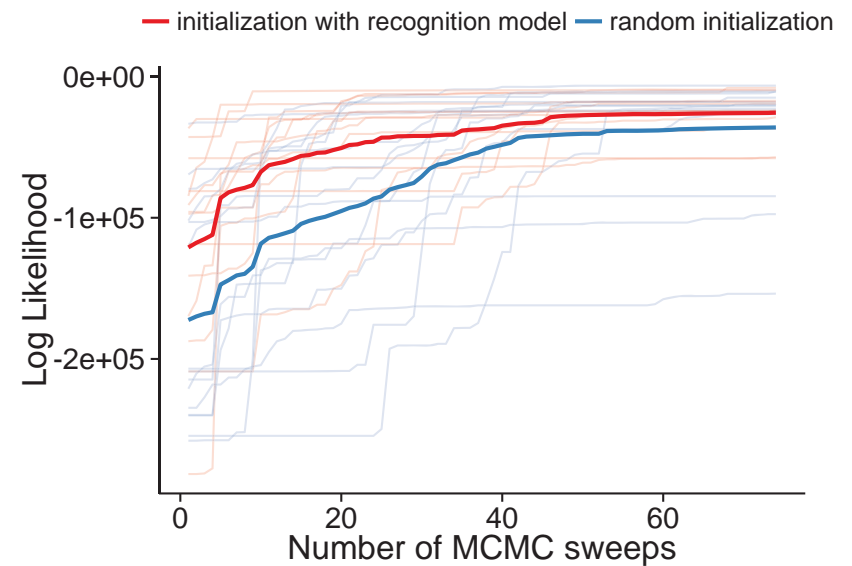
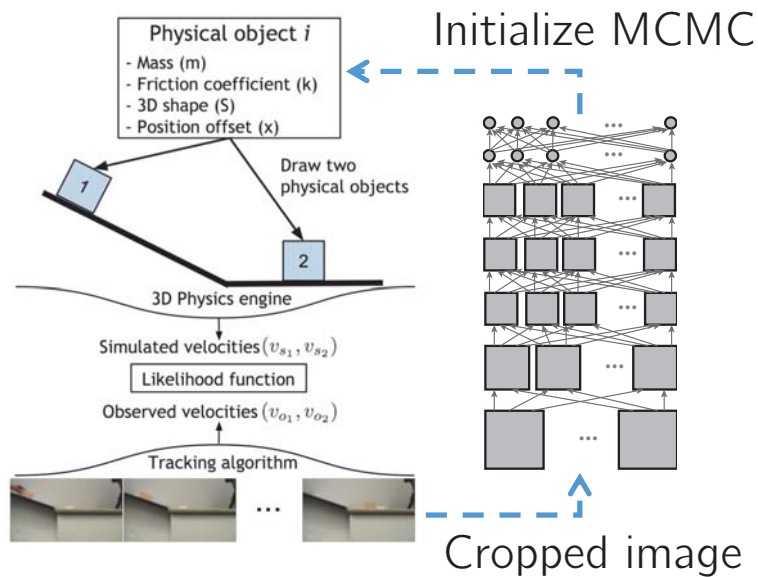
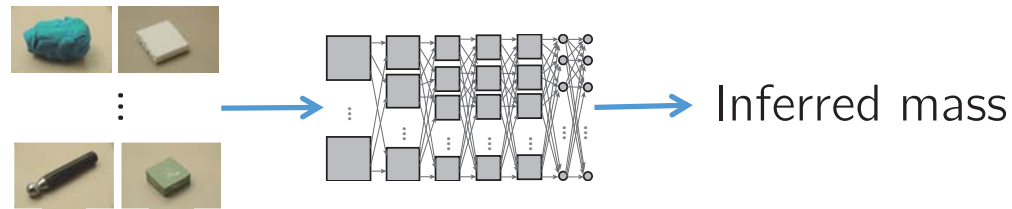
Results



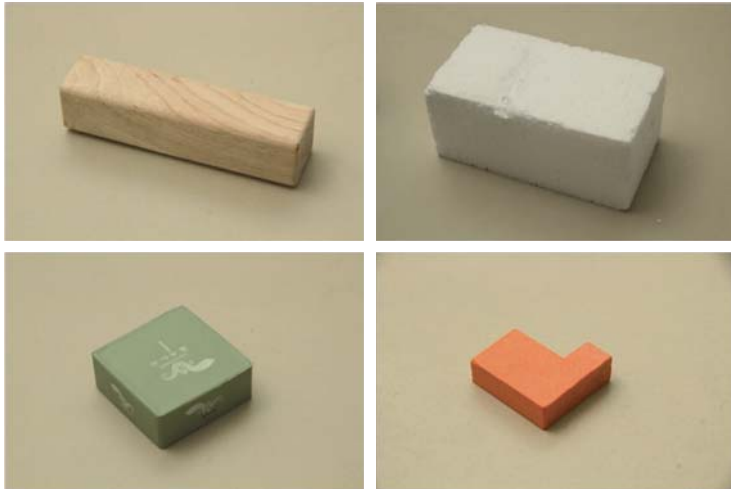
Wu*, Yildirim*, Lim, Freeman, Tenenbaum. NeurIPS'15

Generative + Recognition Model

If the model has prior knowledge like humans do...



We've seen...



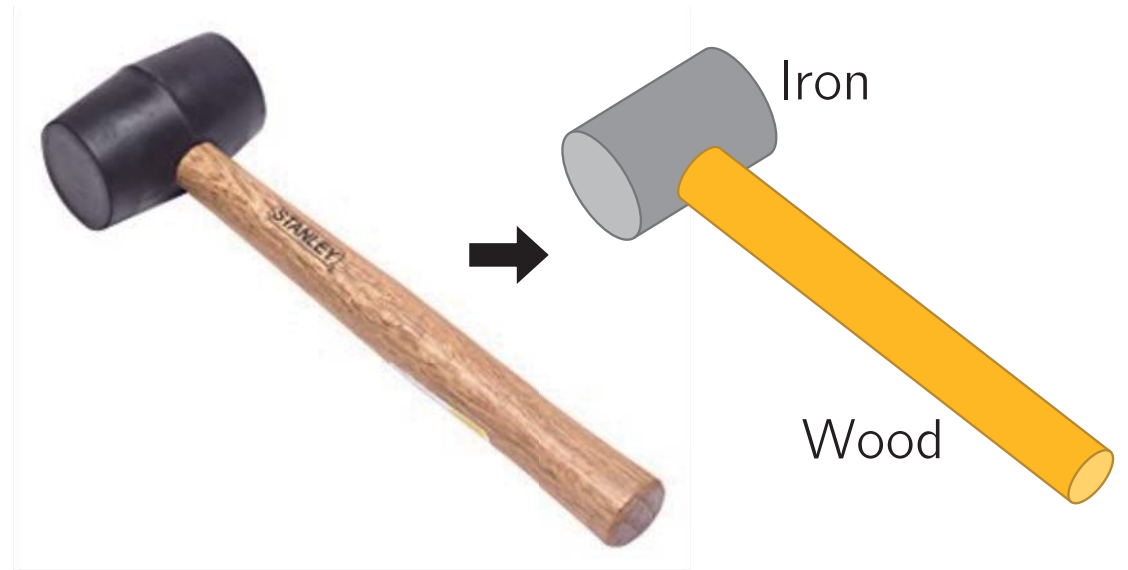
What about?



Learning Shape Abstractions



Physical Primitive Decomposition



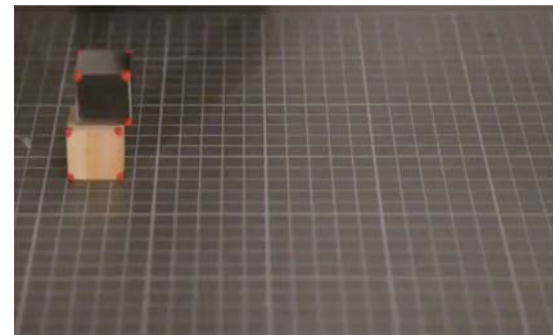
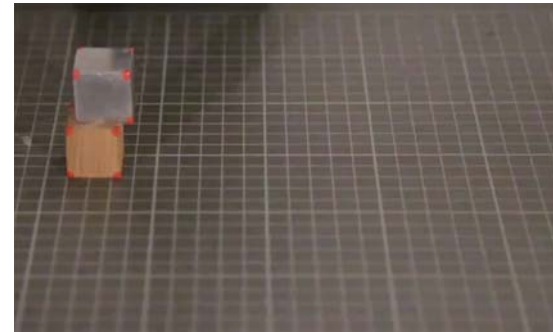
Appearance + Physics

Aluminum
(2.87g/ml)

Oak
(0.67g/ml)

Steel
(7.74g/ml)

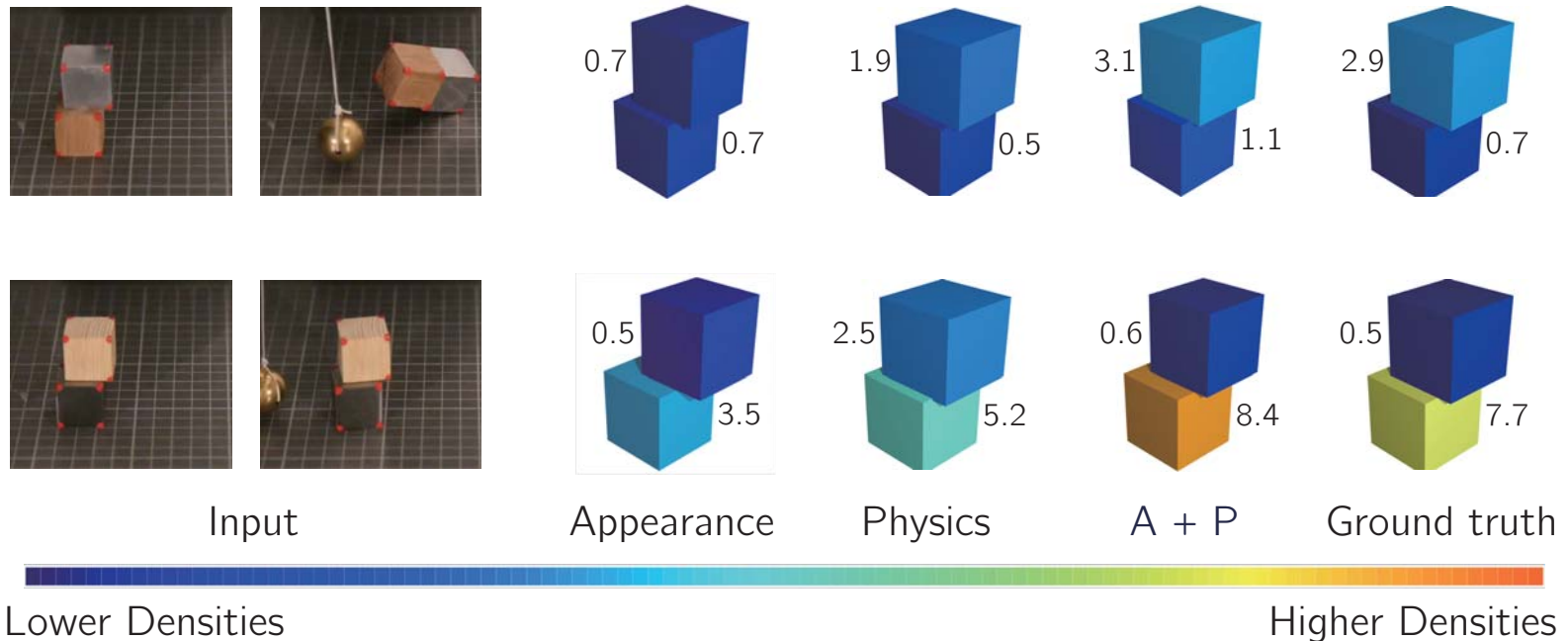
Pine
(0.48g/ml)



Visual Appearance
(Very Similar)

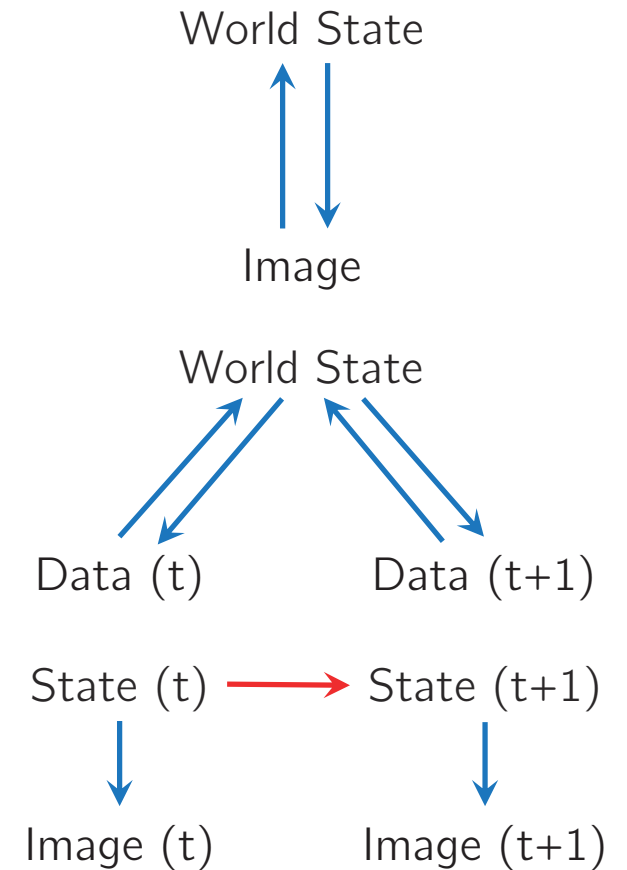
Physics Trajectory
(Very Different)

Physical Primitive Decomposition



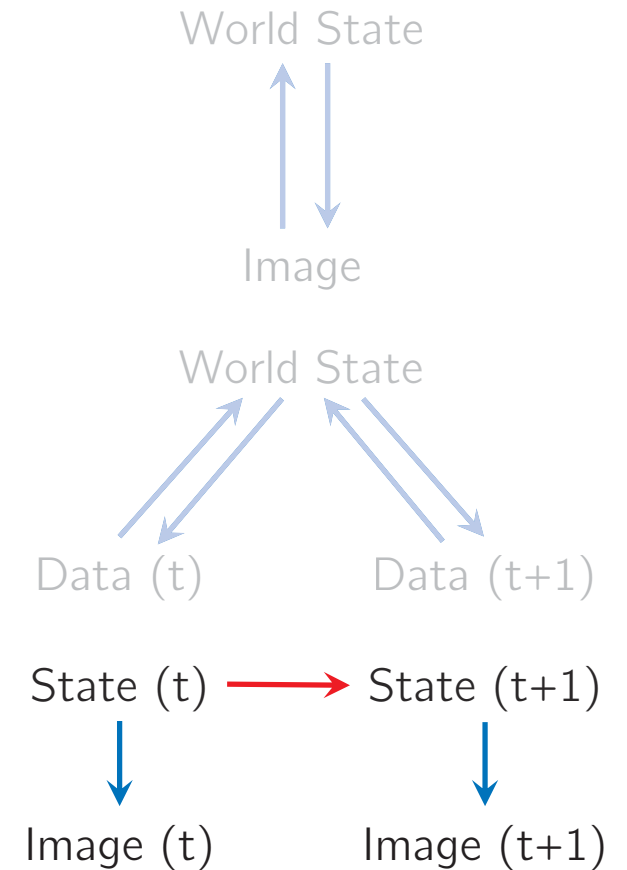
Physical Scene Understanding

- Learning to invert a graphics engine
 - Inferring fine object geometry
 - Learning structured shape representations (shape + texture)
 - Beyond single object, learning scene representations
- Learning to invert a physics engine
 - Inferring object physical properties
 - Joint modeling of object shape and physics
- Learning simulation engines themselves



Physical Scene Understanding

- Learning to invert a graphics engine
 - Inferring fine object geometry
 - Learning structured shape representations (shape + texture)
 - Beyond single object, learning scene representations
- Learning to invert a physics engine
 - Inferring object physical properties
 - Joint modeling of object shape and physics
- Learning simulation engines themselves



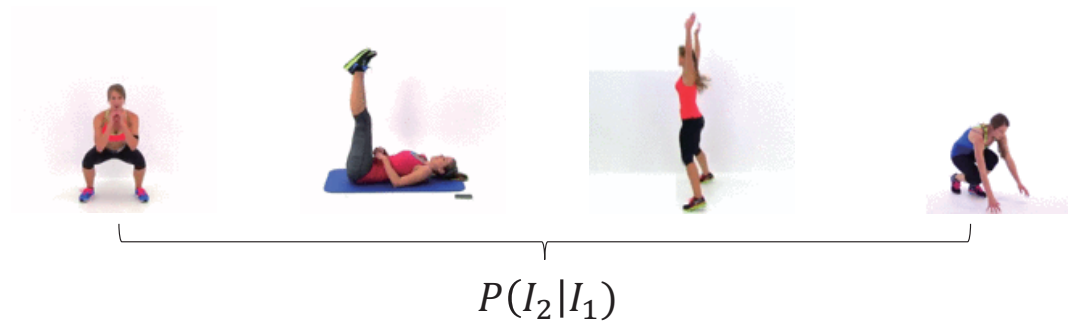
Key Features on Dynamics Modeling

- Depending on visual content
- Modeling uncertainty



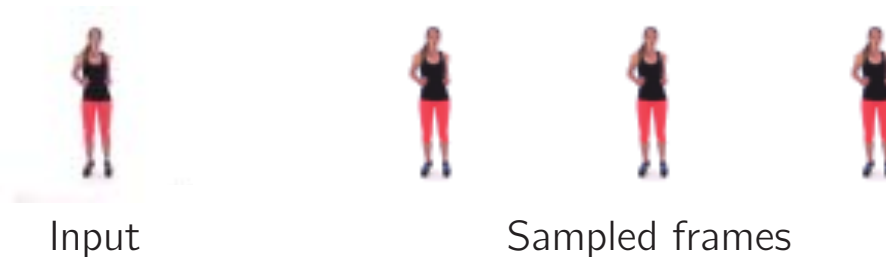
Visual Dynamics

- Two temporally-consecutive frames



$P(I_2|I_1)$: Probabilistic distribution of the second frame conditioned on the first frame

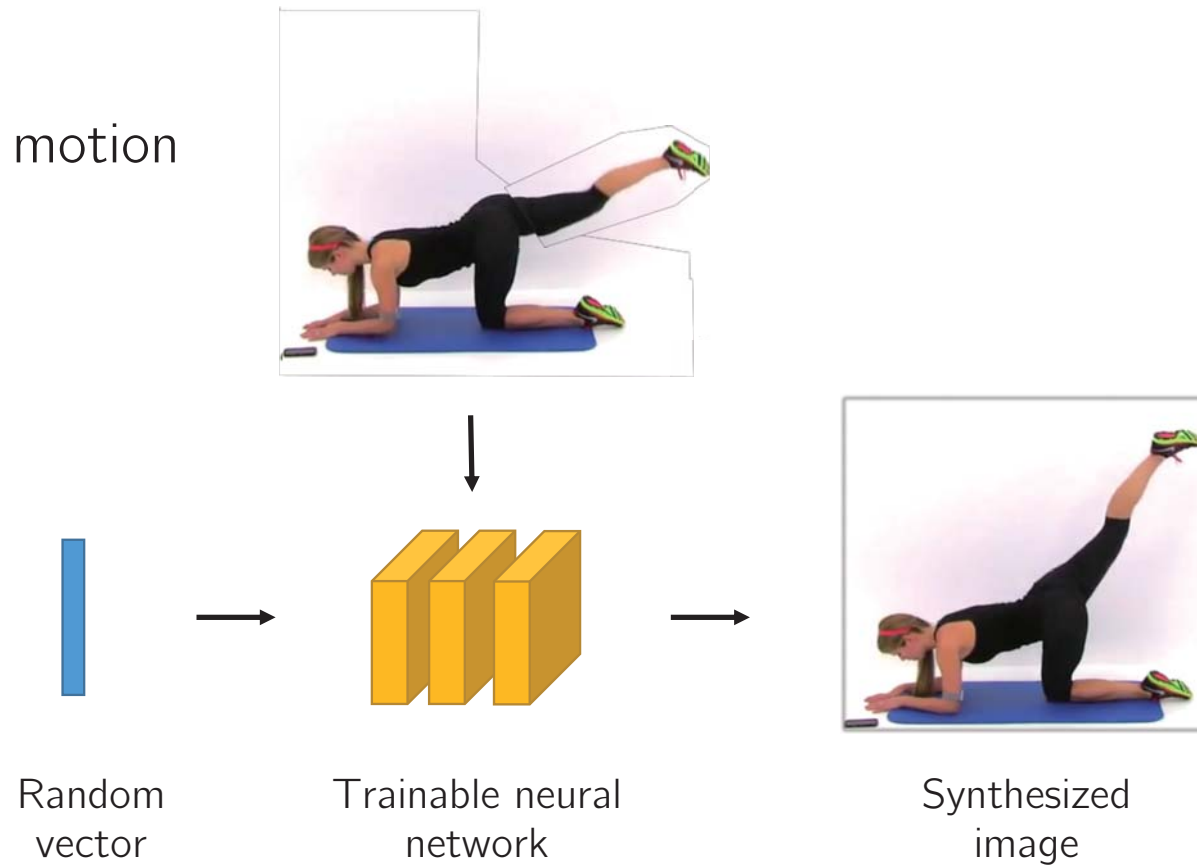
- Prediction



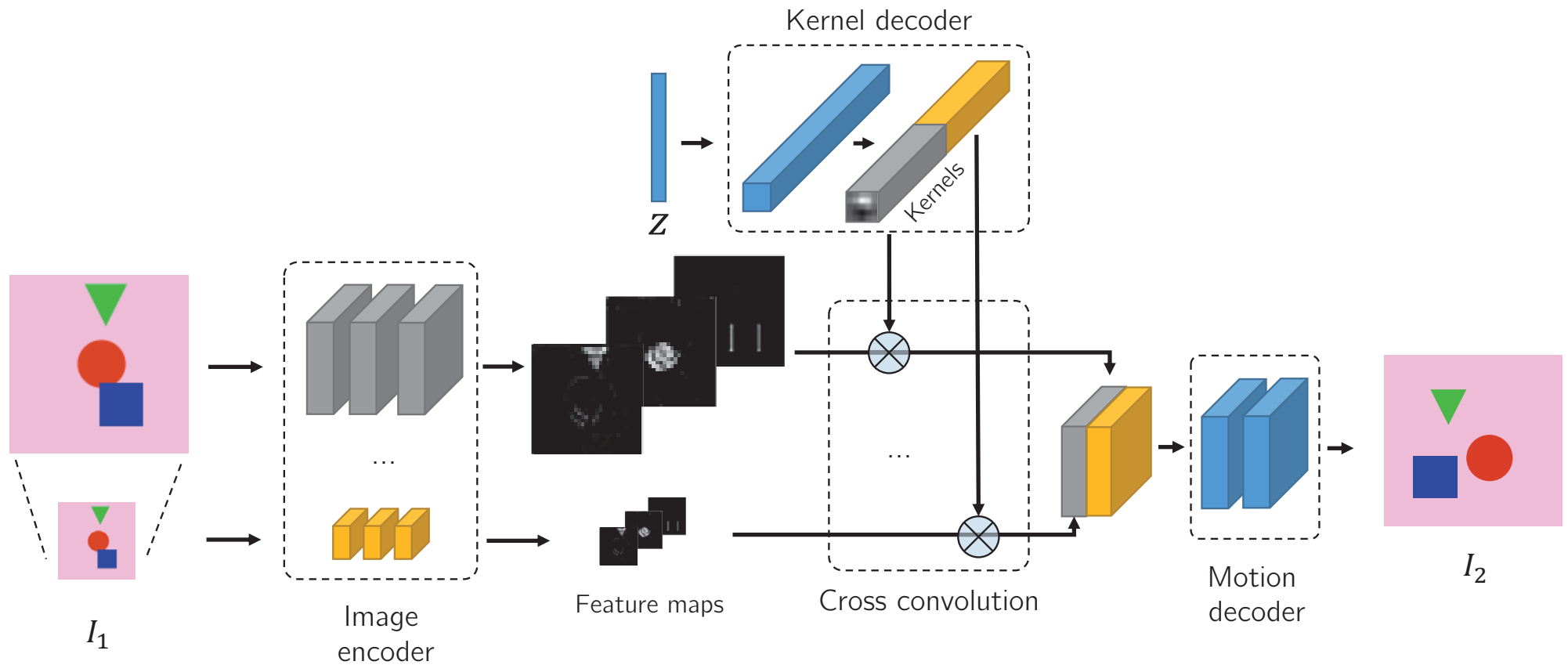
Xue*, Wu*, Bouman, Freeman. NeurIPS'16, TPAMI'18

Decomposing Objects into Independently Movable Parts

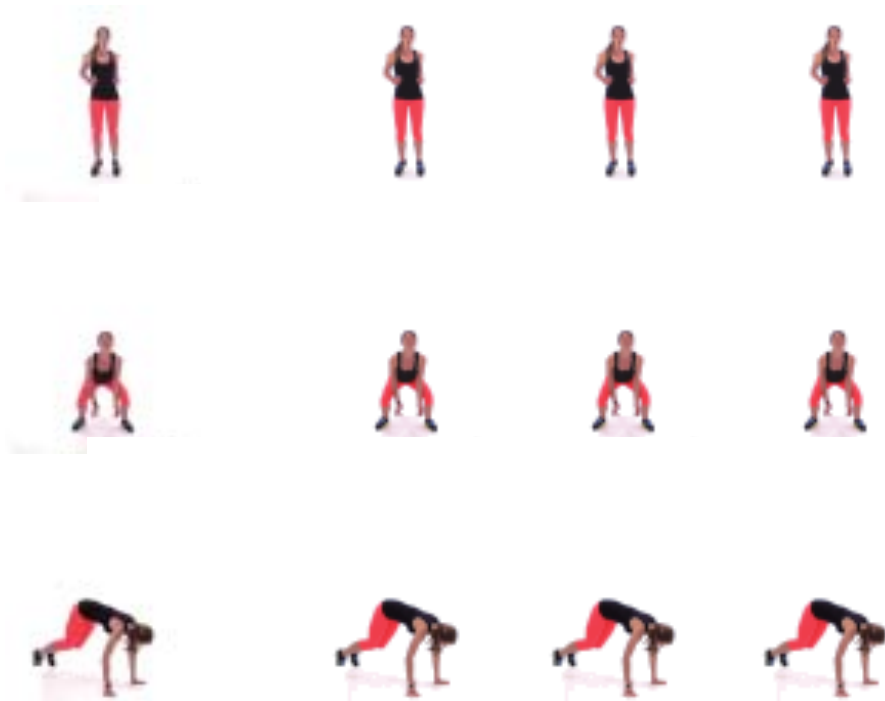
- Identify movable segments
- Model their dynamics
- Combine the sampled motion



Layered Cross-Convolutional Networks

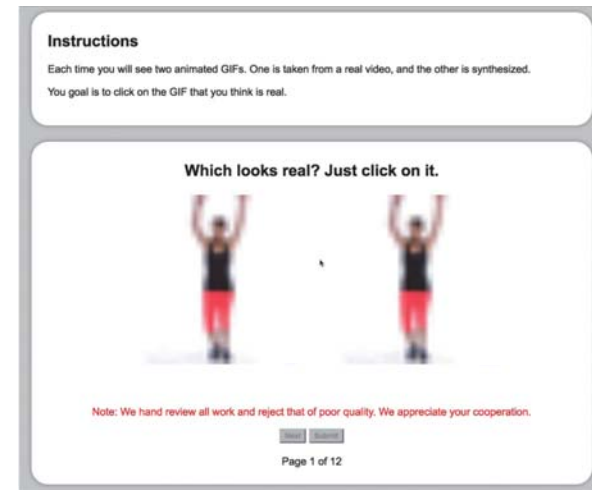


Results on Real Videos



Input

Synthesized next frames

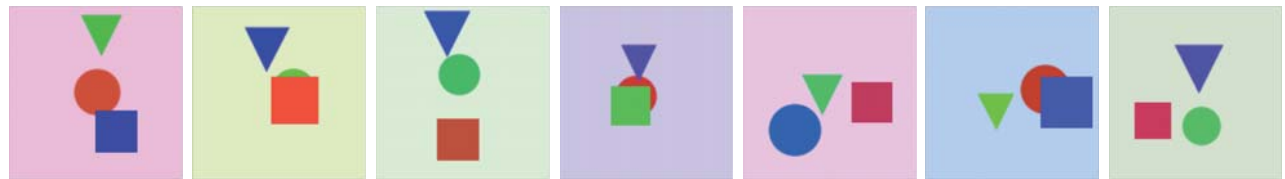


	% of synthetic labeled as real
Transfer flow	25.5
Ours	31.3

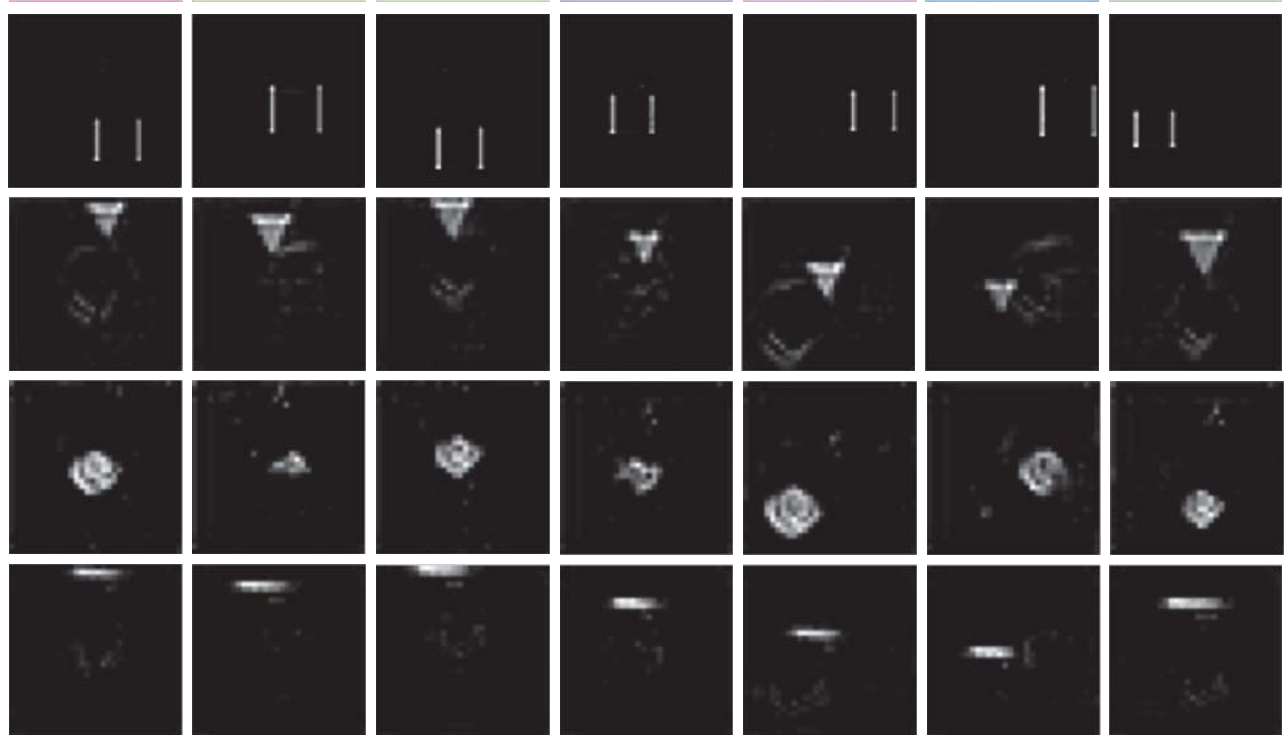
Visualize learned features



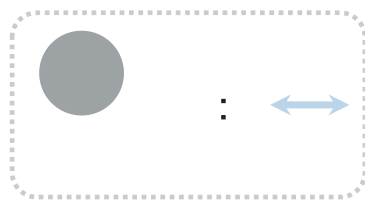
Input



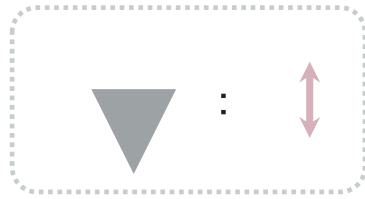
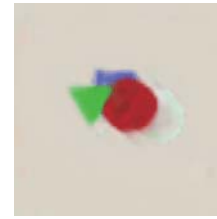
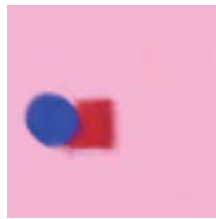
Feature maps



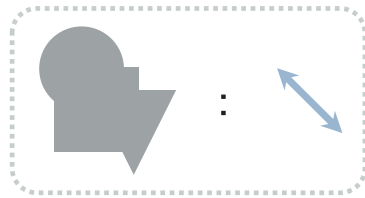
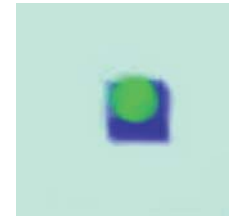
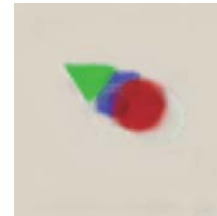
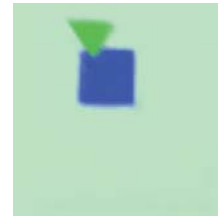
Interpretable Latent Representations



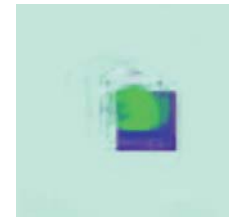
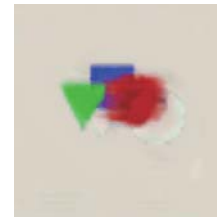
D69



D80



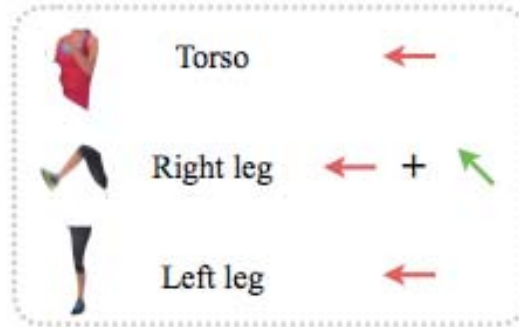
D121



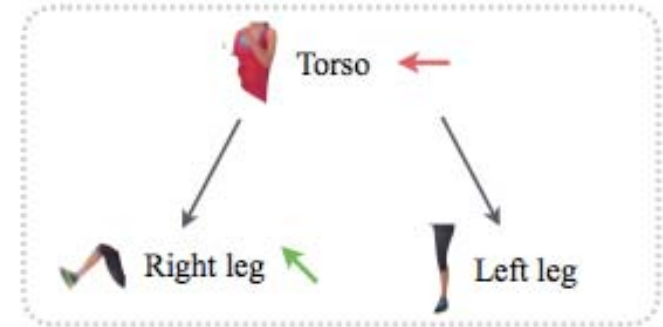
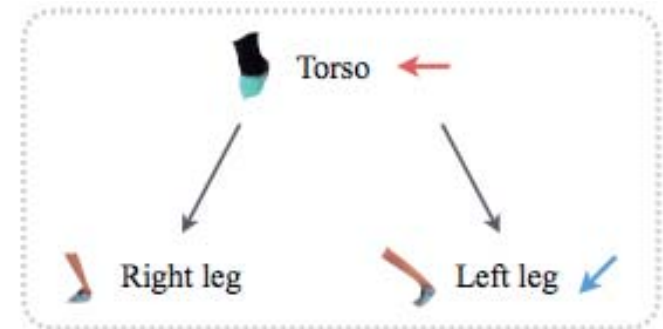
Ext. I: Unsupervised Structure Discovery



(a) Pairs of images



(b) Global motions of parts



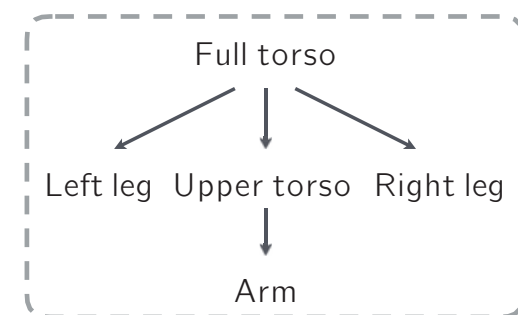
(c) Local motions of parts

Ext. I: Unsupervised Structure Discovery



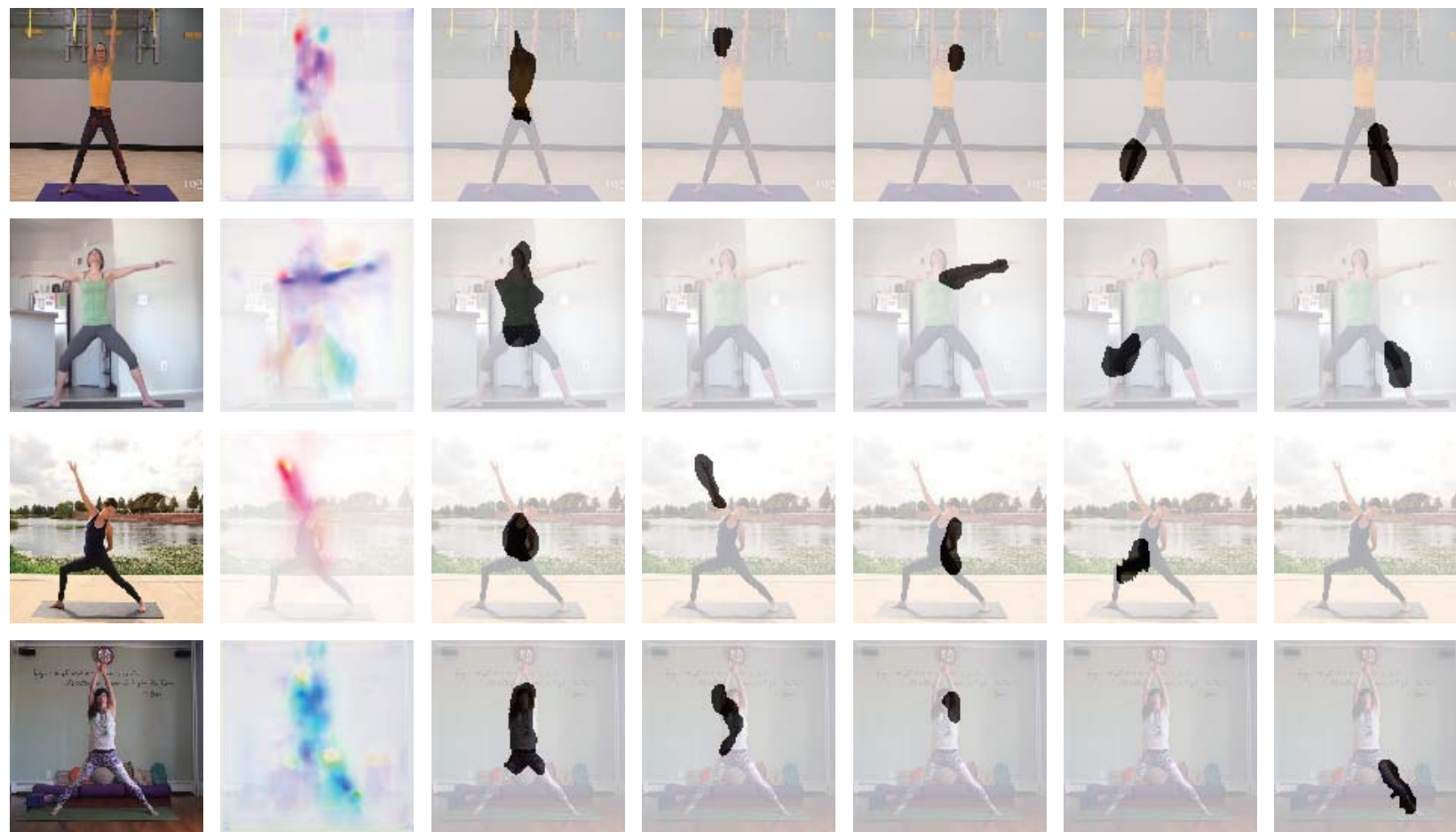
Input frame Flow field (a) Full torso (b) Upper torso (c) Arm (d) Right leg (e) Left leg

	(a)	(b)	(c)	(d)	(e)
(a)	1				
(b)	1	1			
(c)		1	1		
(d)	1			1	
(e)	1				1

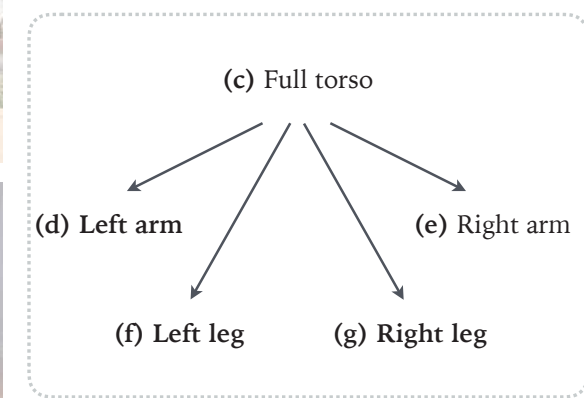


Hierarchical tree structure

Ext. I: Unsupervised Structure Discovery



	(c)	(d)	(e)	(f)	(g)
(c)	1				
(d)	1	1			
(e)	1		1		
(f)	1			1	
(g)	1				1

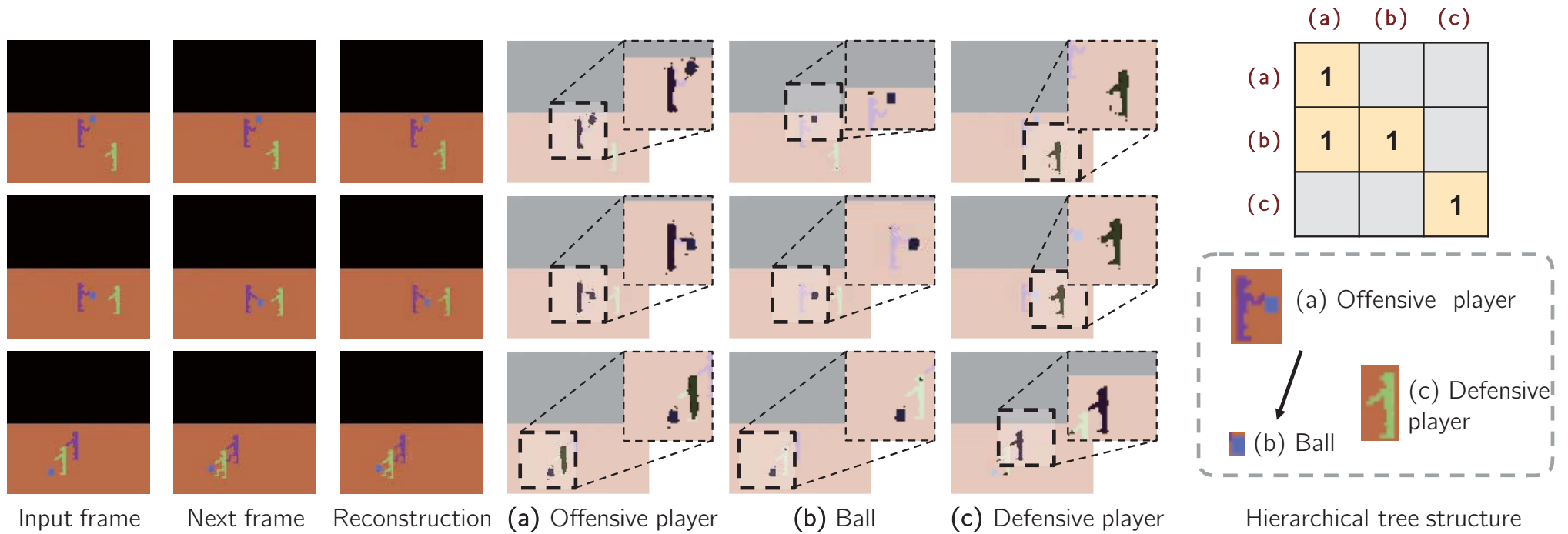


(a) Input frame (b) Flow field (c) Full torso (d) Left arm (e) Right arm (f) Left leg (g) Right leg

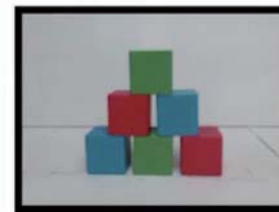
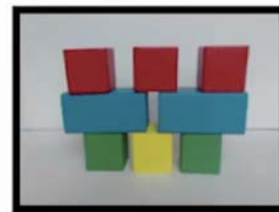
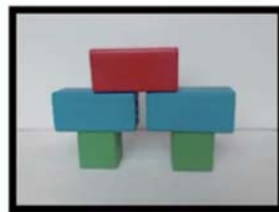
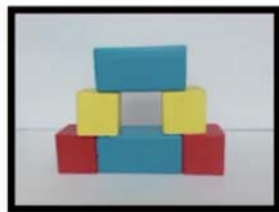
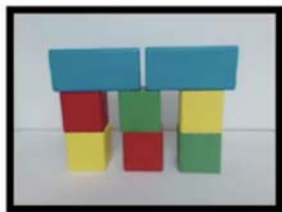
(h) Hierarchical tree structure

Xu*, Liu*, Sun, Murphy, Freeman, Tenenbaum, Wu. ICLR'19

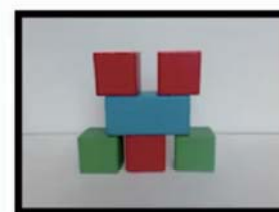
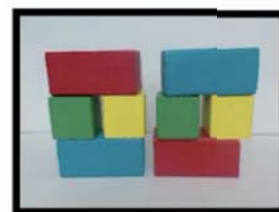
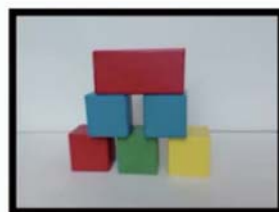
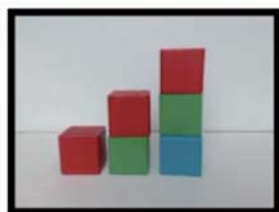
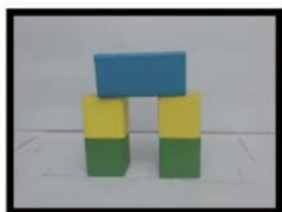
Ext. I: Unsupervised Structure Discovery



Ext II: Planning and Control



goal images

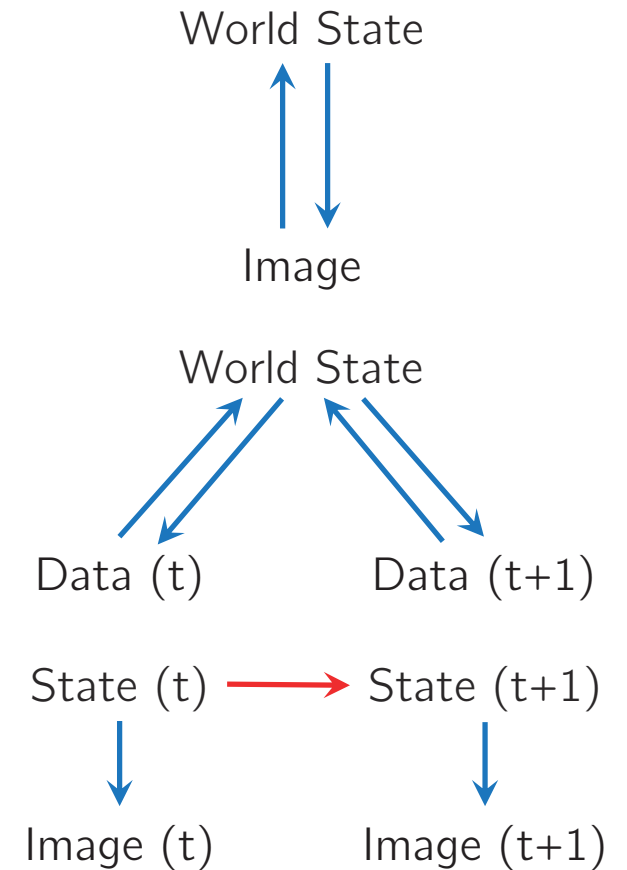


Ext II: Planning and Control



Physical Scene Understanding

- Learning to invert a graphics engine
 - Inferring fine object geometry
 - Learning structured shape representations (shape + texture)
 - Beyond single object, learning scene representations
- Learning to invert a physics engine
 - Inferring object physical properties
 - Joint modeling of object shape and physics
- Learning simulation engines themselves
 - Learning object dynamics in the pixel space
 - Modeling object dynamics for control



Physical Scene Understanding with Compositional Structure

Goal

- Explaining and reasoning about data

Approach

- Leveraging causal structure to integrate generative, forward models with efficient inference algorithms.

Advantages

- 1. Guiding and facilitating model design.
- 2. Allowing learning with little or no supervision.
- 3. Offering rich generalization power.

