# Perception and Generation of Physical Interactions

Srinath Sridhar

GAMES Seminar

August 12, 2021

BROWN

# AI for Logical Reasoning



*Deep Blue versus Garry Kasparov (1997)*



*AlphaGo versus Lee Sedol (2016)*

# AI for Physical Motion*

*Boston Dynamics (2020)*

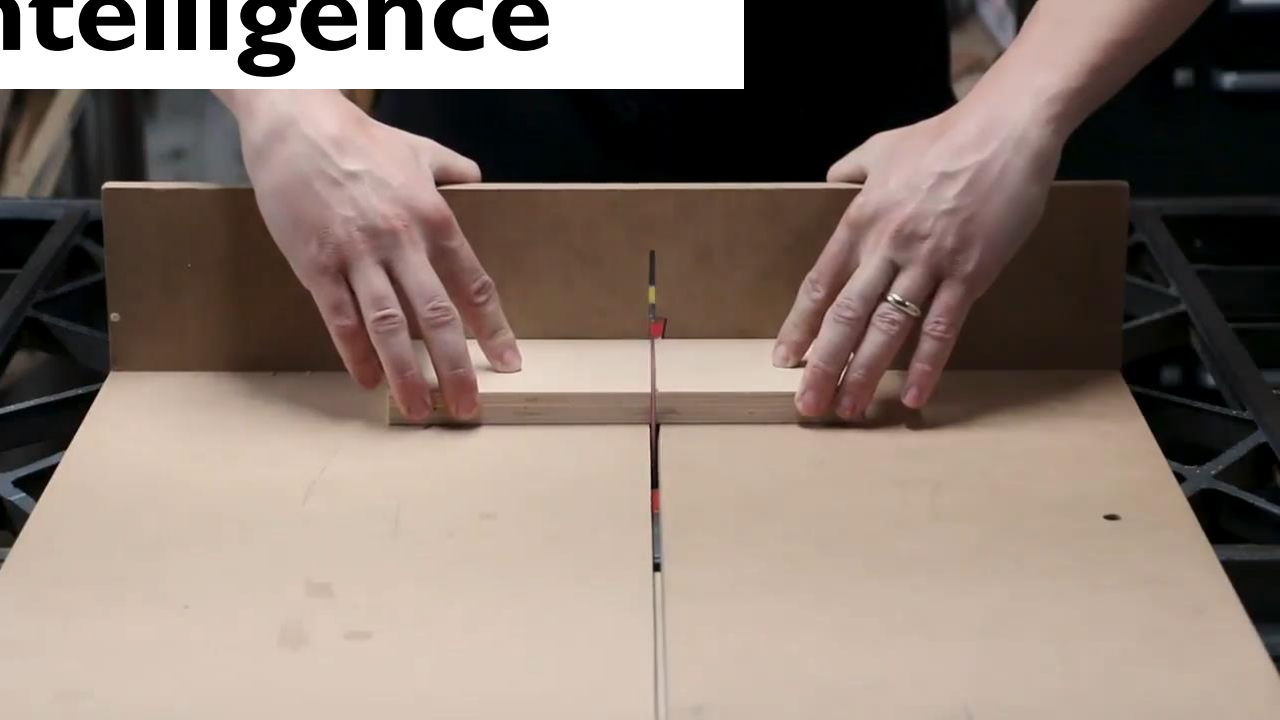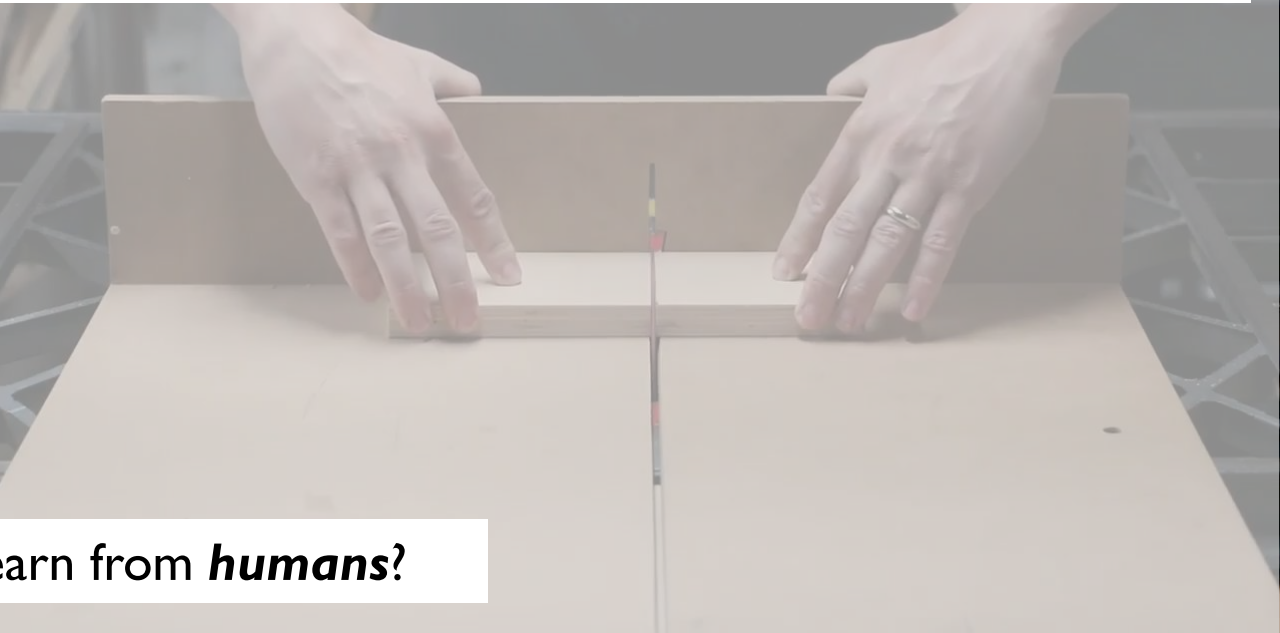*Boston Dynamics (2017)*

Physical Intelligence

Learn physical intelligence by **observing humans**

Hmmm, but why learn from *humans*?

# Applications



*Interactive Household Robots*
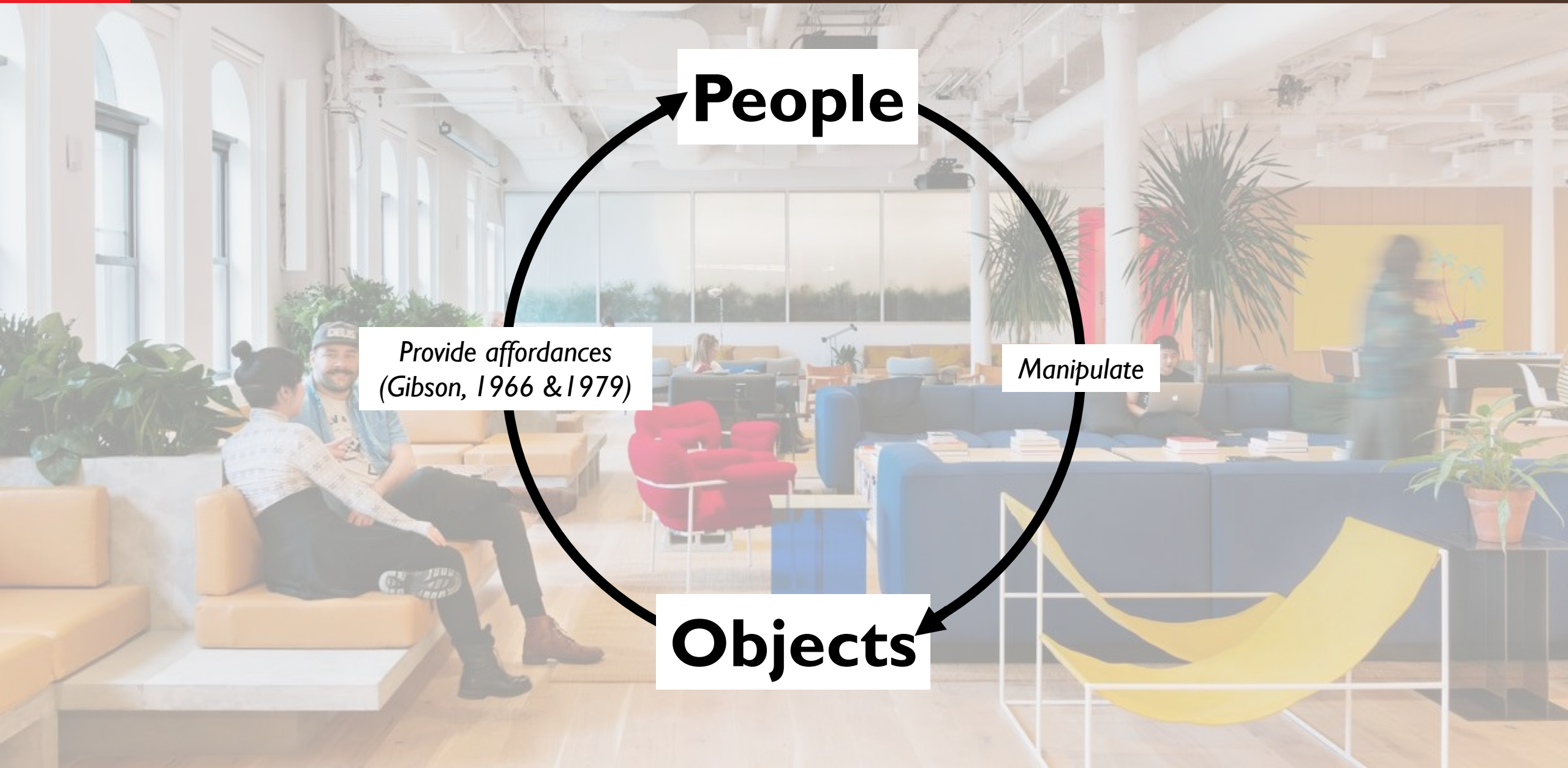
*Interactive Mixed Reality*

**Understanding = Perception + Generation**

*What I cannot create, I do not understand.*
- Richard Feynman

Shotton et al. 2013

Baak et al. 2011

Wei et al. 2016



Girshick et al. 2011
Ganapathi et al. 2012
Ma and Wu 2014
Newell et al. 2016
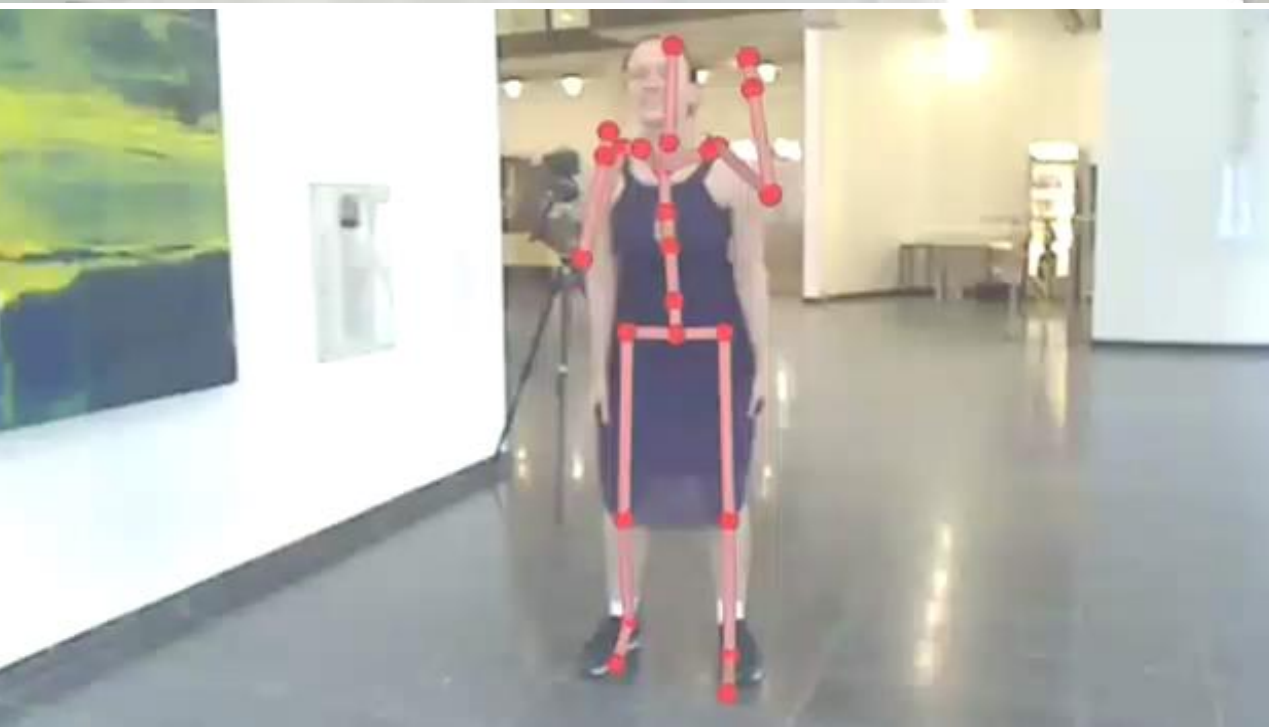Tompson et al. 2014
Insafutdinov et al. 2016
Cao et al. 2017
…

## Limitations

Cao et al. 2019
Zheng et al., 2019
Zhang et al. 2020
Kanazawa et al. 2018
…

- Limited to 2D
- Depth-sensor based
- Lacking generalizability

11

- Uses a single RGB camera / community videos
- Works for diverse scenes and subjects
- Fast (>40 FPS)

**VNect.** Mehta, Sridhar, Sotnychenko, Rhodin, Shafiei, Seidel, Xu, Casas, Theobalt. **SIGGRAPH 2017**

# Sufficient for understanding?

**Input Video**

**Monocular Total Capture**



*Xiang et al., CVPR 2019*

VIBE

Animation

Alternate View

Physical Interactions

Kocabas et al., CVPR 2020

# Generation: Human Motion Prior

## HuMoR:
### 3D Human Motion Model for Robust Pose Estimation

Davis **Rempe**, Tolga **Birdal,** Aaron **Hertzmann**, Jimei **Yang,** Srinath **Sridhar,** Leonidas **Guibas**
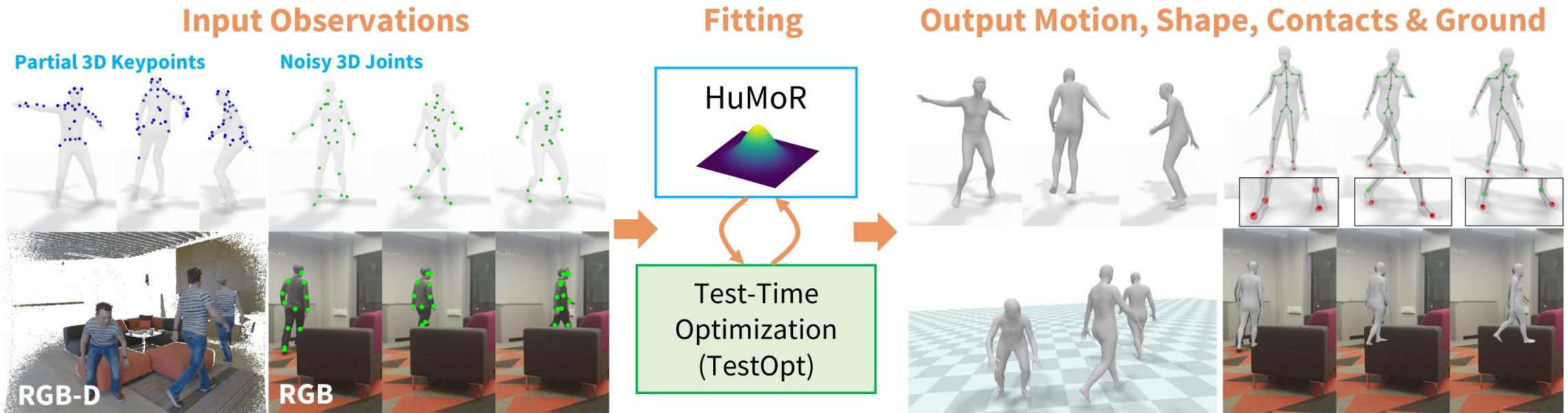
ICCV 2021 (Oral)

**geometry.stanford.edu/projects/humor/**

**HuMoR**. Rempe, Birdal, Hertzmann, Yang, Sridhar, Guibas. **ICCV 2021 [oral presentation]**

1. Learned **generative model** of plausible 3D motion (HuMoR)

2. Time-time optimization (TestOpt) using **HuMoR as a prior**

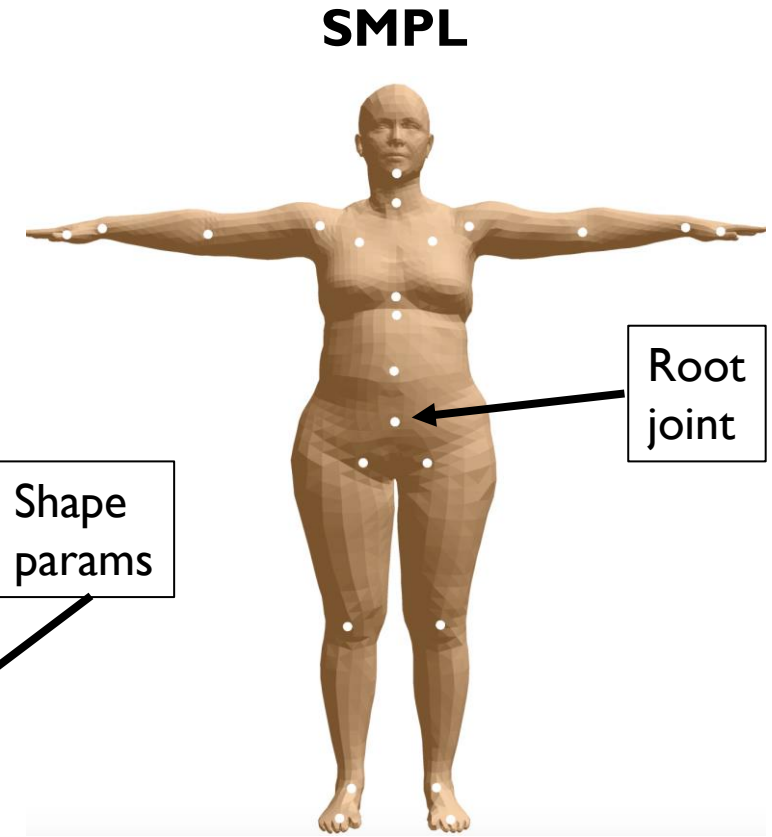$$\mathbf{x} = \begin{bmatrix} \mathbf{r} & \dot{\mathbf{r}} & \Phi & \dot{\Phi} & \Theta & \mathbf{J} & \dot{\mathbf{J}} \end{bmatrix}$$

**SMPL Root**

**SMPL Body**

**Joints**

**Position/Vel** $\in \mathbb{R}^3$

**Rot/Vel** $\in \mathbb{R}^3$

**Joint Angles** $\in \mathbb{R}^{3 \times 21}$

**Joint Pos/Vel** $\in \mathbb{R}^{3 \times 22}$

**SMPL**

Root joint

Shape params

*Over-parameterization* of joint positions:
(i)   Implicit through SMPL  $\mathbf{J}^{\mathrm{SMPL}} = M(\mathbf{r}, \Phi, \Theta, \beta)$
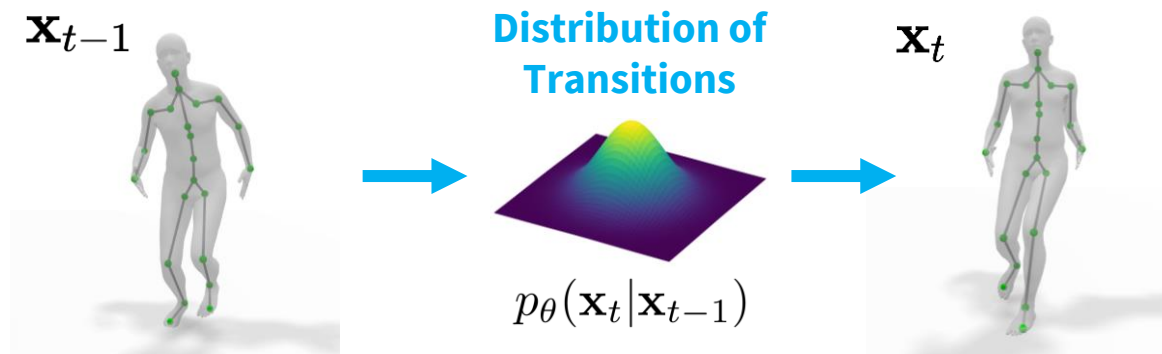(ii)  Explicit from state  $\mathbf{J}$

*Loper et al., SIGGRAPH Asia 2015*

- For **state** $\mathbf{x}_t$ at time $t$

$$p_\theta(\mathbf{x}_0, \mathbf{x}_1, \ldots, \mathbf{x}_T)$$

$$= p_\theta(\mathbf{x}_0) \prod_{t=1}^{T} \underbrace{p_\theta(\mathbf{x}_t \mid \mathbf{x}_{t-1})}_{\text{HuMoR}}$$

Learn the *plausibility* of a transition, *i.e.*, **distribution of dynamics**



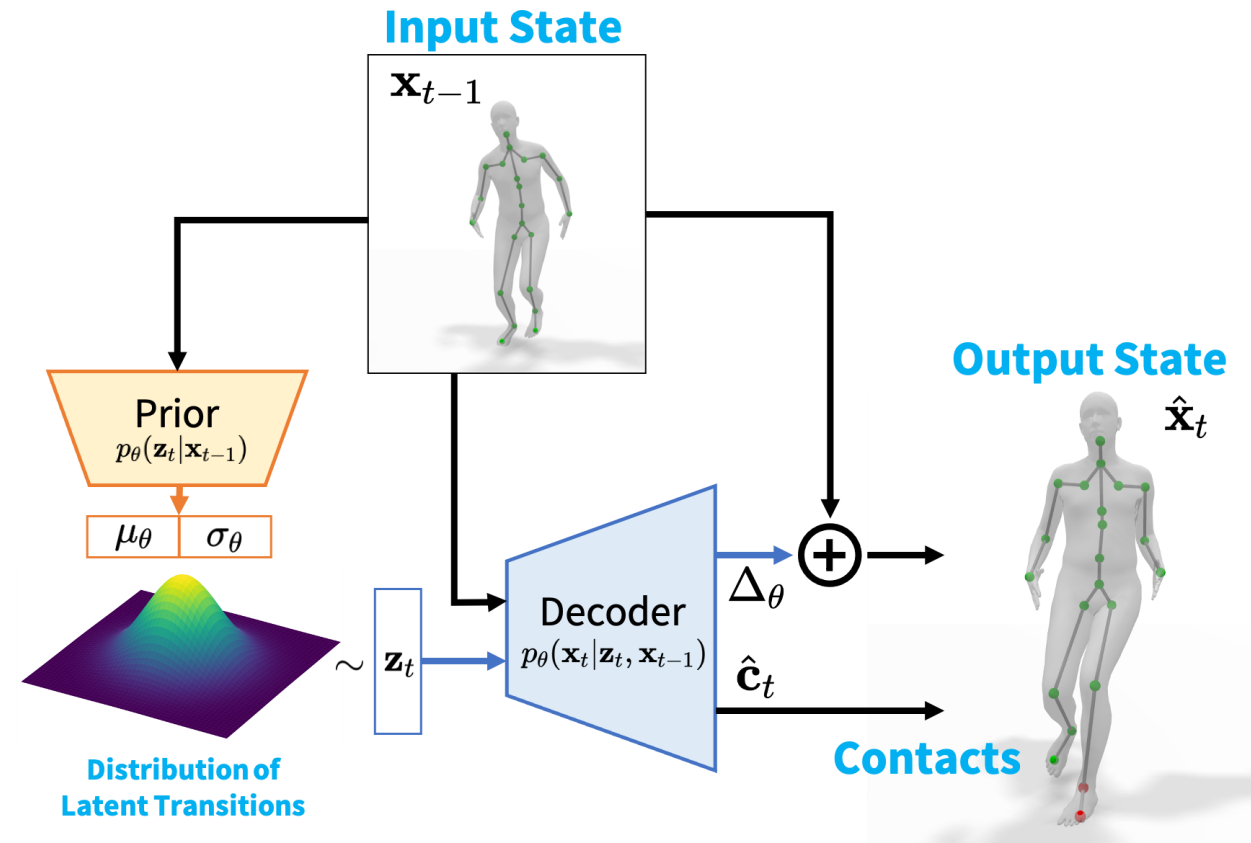$\mathbf{x}_{t-1}$     Distribution of Transitions     $\mathbf{x}_t$

$$p_\theta(\mathbf{x}_t \mid \mathbf{x}_{t-1})$$

- Use **latent variable model**

- Conditional VAE

**Generation:**

- Conditional Prior

- Decoder

**Outputs:**

    *Change* in state $\Delta_\theta$
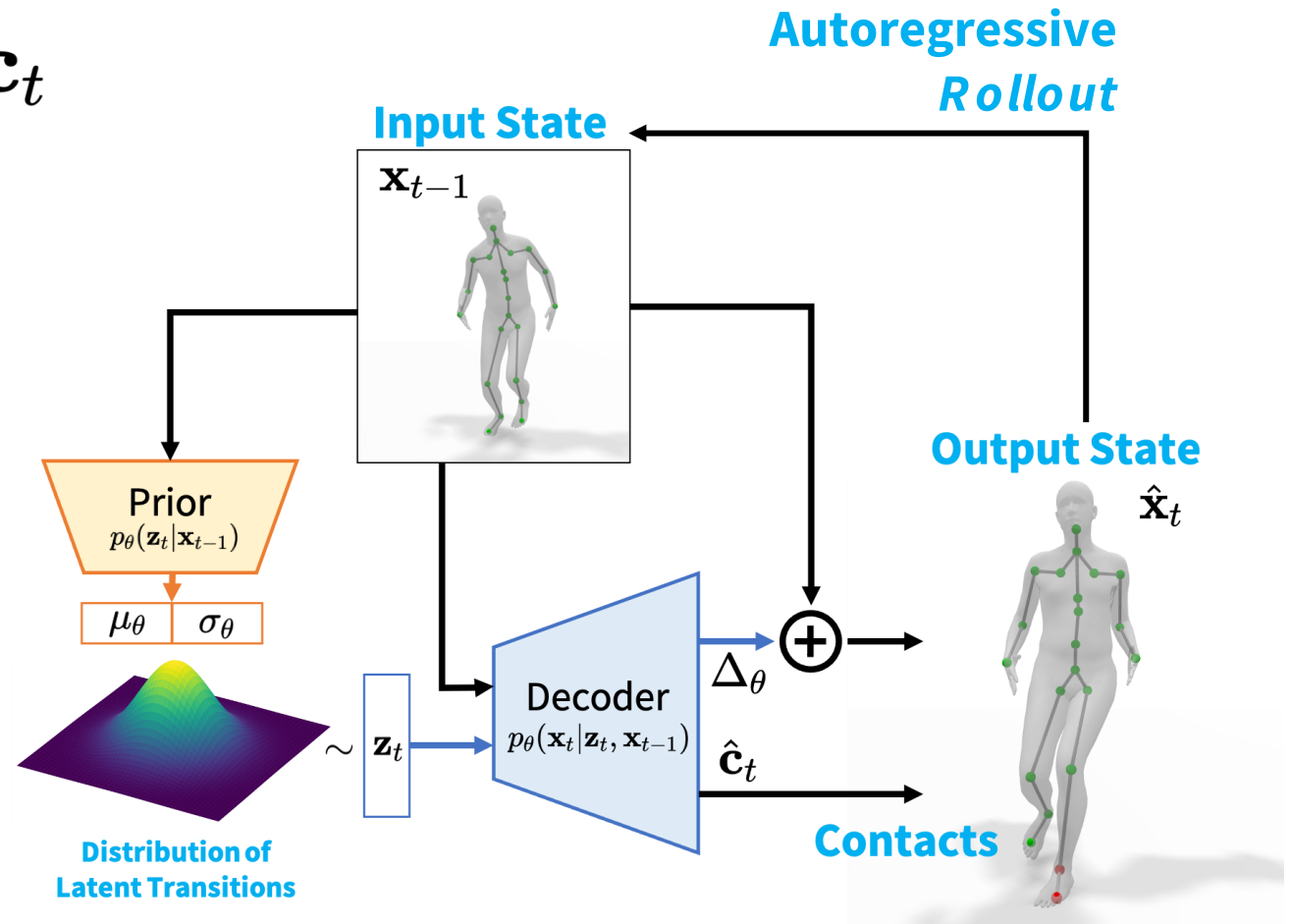    Ground *contact* classification $\mathbf{c}_t$

Autoregressive **sampling**

Deterministic **rollout**

    Motion "parameters" $\mathbf{x}_0, \mathbf{z}_{1:T}$
    give $\mathbf{x}_t = \mathbf{x}_{t-1} + \Delta_\theta(\mathbf{z}_t, \mathbf{x}_{t-1})$
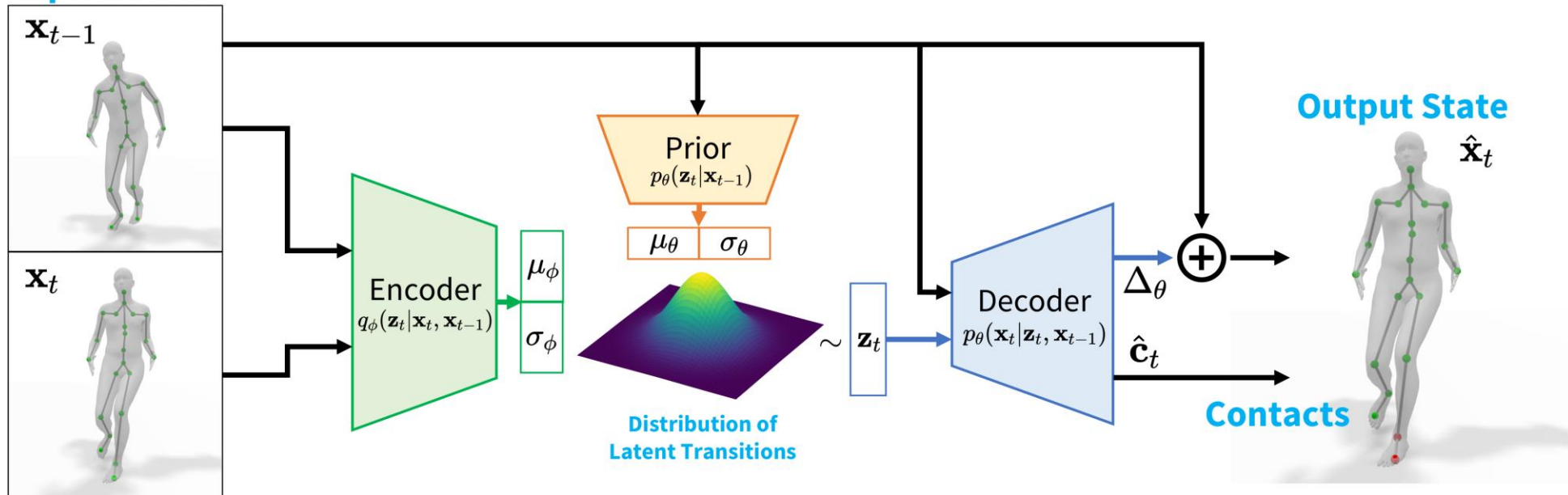    Have *prior* on $\mathbf{z}_{1:T}$

- Encoder for training on AMASS [Mahmood et al., ICCV 2019]

- Loss based on lower bound

  - Reconstruction, KL, consistency

$$\log p_\theta(\mathbf{x}_t|\mathbf{x}_{t-1}) \geq \mathbb{E}_{q_\phi}[\log p_\theta(\mathbf{x}_t|\mathbf{z}_t, \mathbf{x}_{t-1})]$$
$$- D_{\mathrm{KL}}(q_\phi(\mathbf{z}_t|\mathbf{x}_t, \mathbf{x}_{t-1}) \parallel p_\theta(\mathbf{z}_t|\mathbf{x}_{t-1}))$$

**Unseen Body Shapes**

Subject 1

Subject 2

Subject 3

Subject 4

**Diverse Samples**

Sample 1

Sample 2
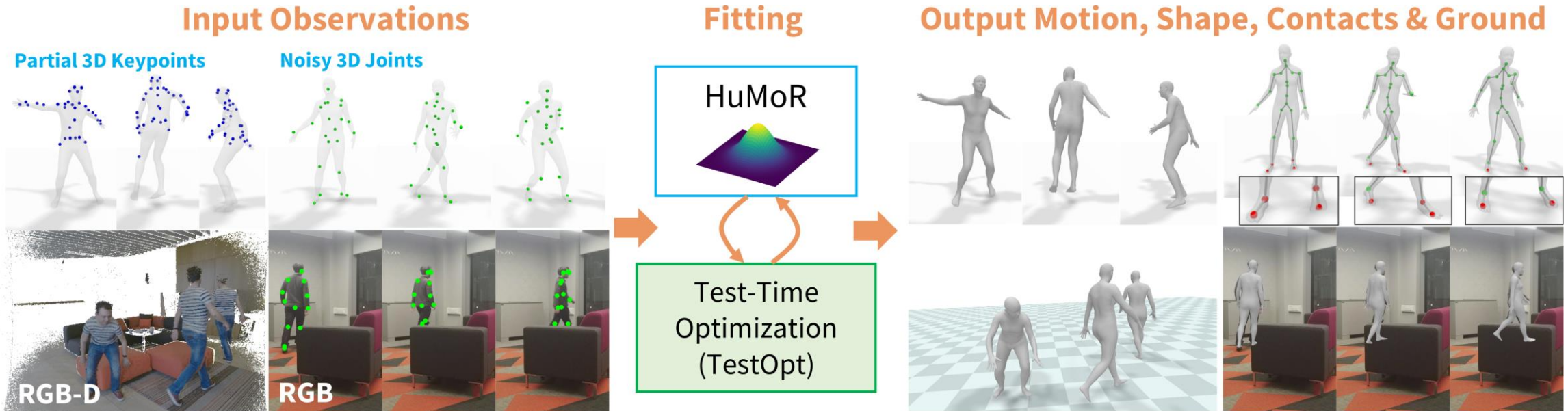
Sample 3

Sample 4

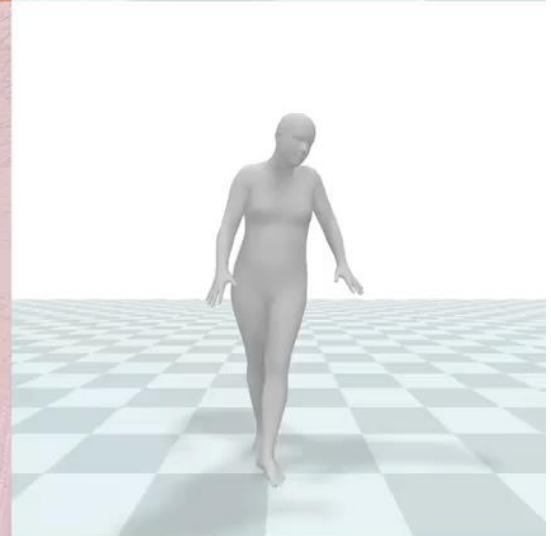Partial 3D Keypoints: **Sequence 1**

Observations & Ground Truth

Output

Input | Motion & Shape | Alternate View

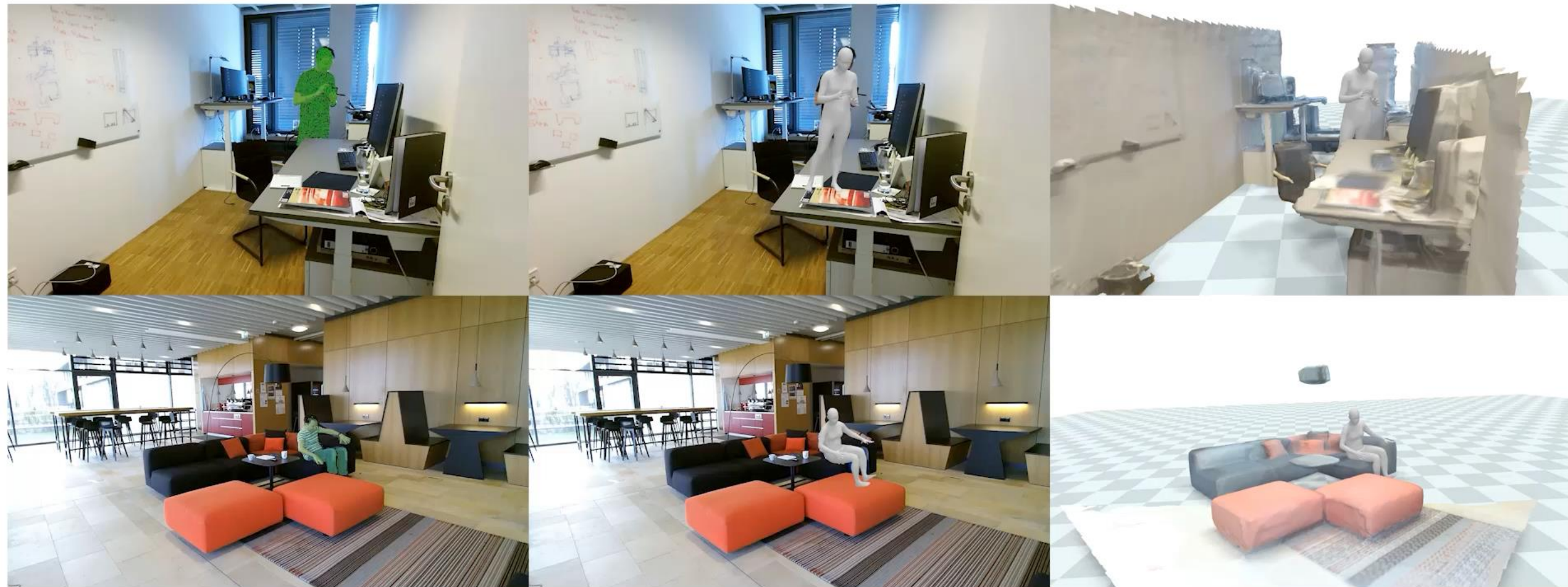Input

Ground Contacts

Alternate View

# Fitting to 2D Joints + 3D (RGB-D)

Input

Motion & Shape

Ground Plane

Input | Final

Input | Final

*Tsuchida et al., ISMIR 2019*

**People**

**Objects**

*Provide affordances
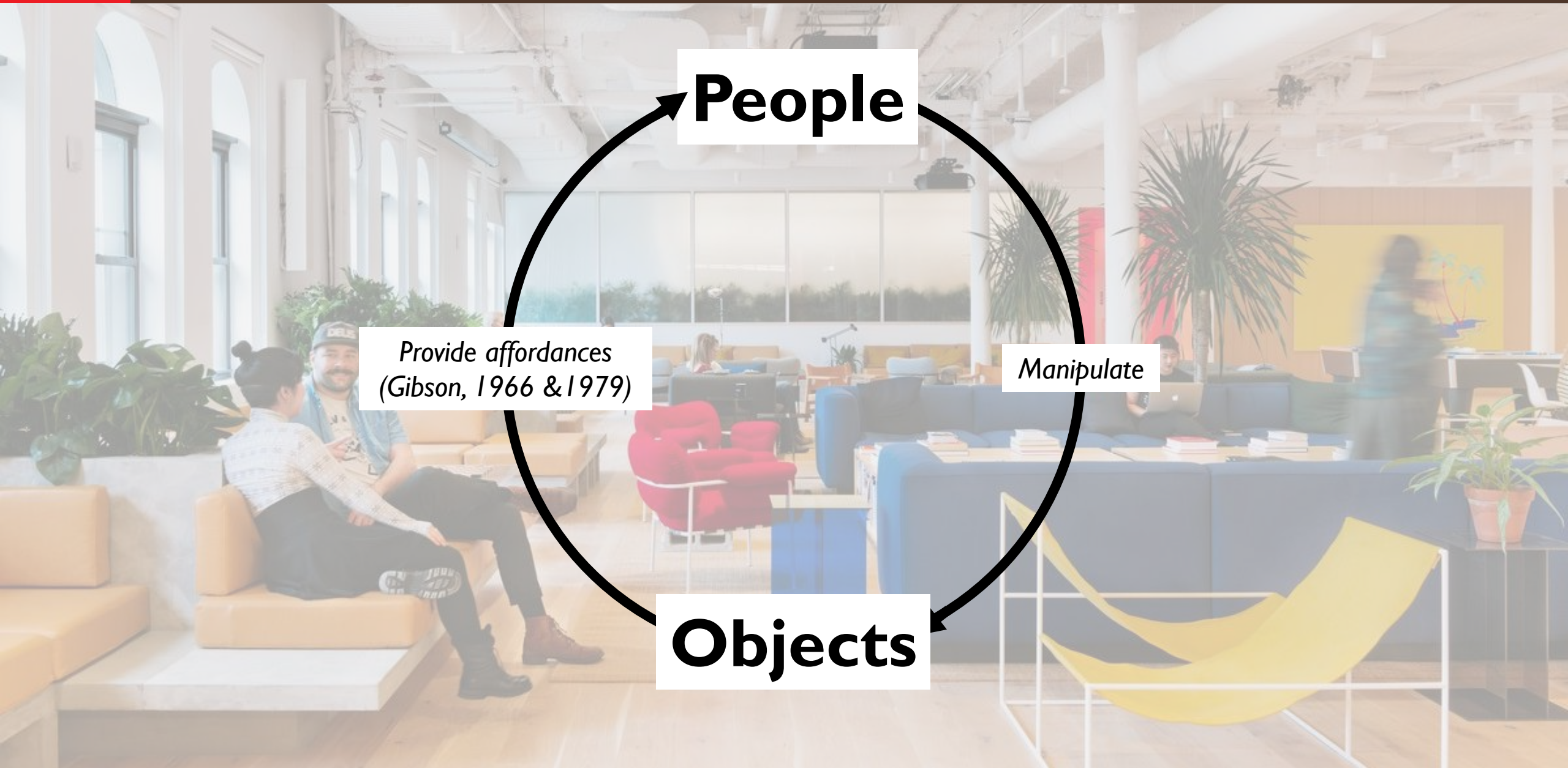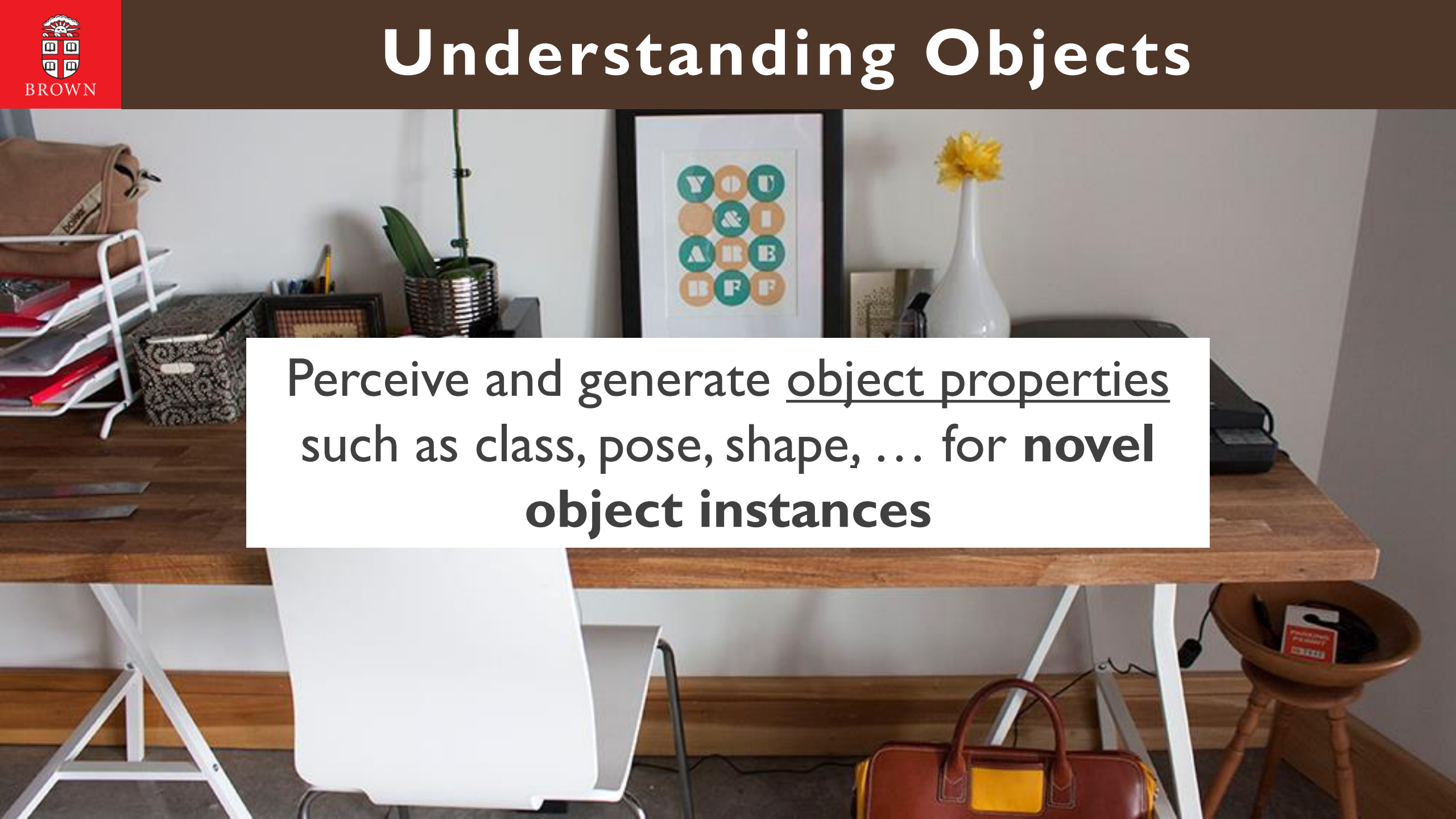(Gibson, 1966 &1979)*

*Manipulate*
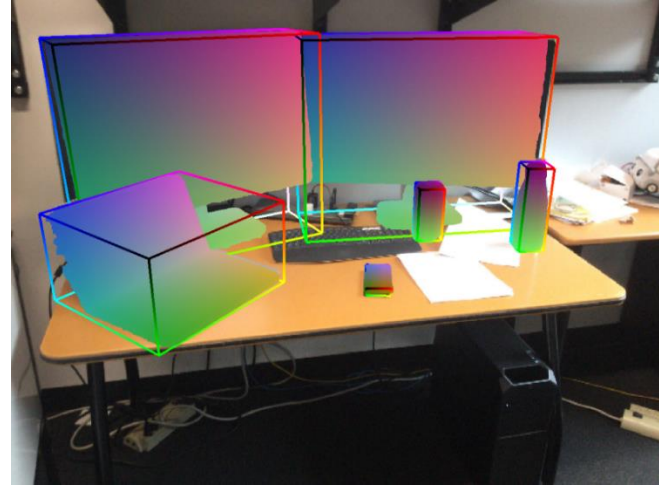
Perceive and generate <u>object properties</u> such as class, pose, shape, … for **novel object instances**

RGB(-D)

- 6 DoF Pose
  - 3D position
  - 3D orientation
- Object dimensions

## Limitations

- Limited category-level reasoning
- Datasets
- Generalizability

**6 DoF Pose**

- Brachmann et al. 2014
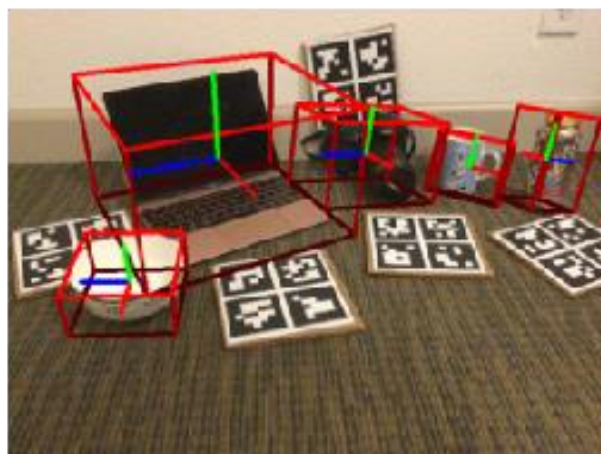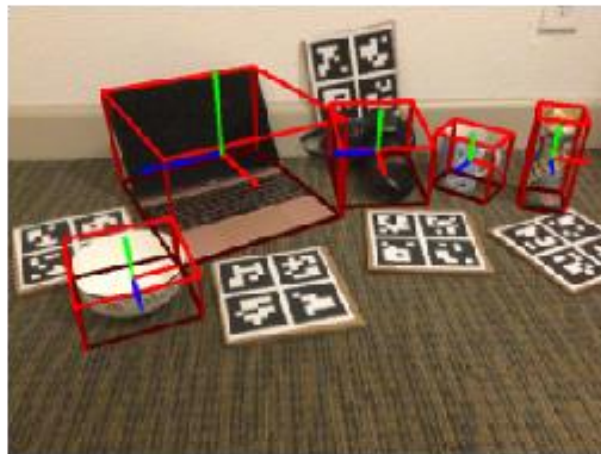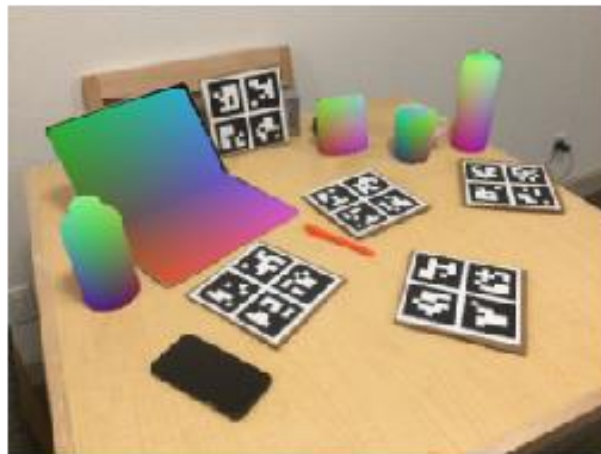- Kehl et al. 2017
- Xiang et al. 2017
- Krull et al. 2016
- Doumanoglou et al. 2016 …

**3D Object Detection**

- Gupta et al. 2013, 2014
- Engelcke et al. 2017
- Song et al. 2016
- Qi et al. 2018
- Zhou et al. 2017 …

Ground Truth

Prediction

***Pix2Surf.*** Lei, Sridhar, Guerrero, Sung, Mitra, Guibas. **ECCV 2020**

**Canonicalize**

- Position
- Orientation
- Size
- Articulation

RGB colors represent *XYZ* coordinates of shape.

(0, 1, 1)  (0, 1, 0)

(1, 1, 1)  (1, 1, 0)

(0, 0, 0)

(1, 0, 0)

Motion over time. How do we
reconstruct?

Flow in space (CNF)

Flow in time (Latent ODE)

**CaSPR**: Learning <u>Ca</u>nonical <u>S</u>patiotemporal <u>P</u>oint Cloud <u>R</u>epresentations

Davis Rempe, Tolga Birdal, Yongheng Zhao, Zan Gojcic, Srinath Sridhar, Leonidas Guibas

NeurIPS 2020 (Spotlight)

**geometry.stanford.edu/projects/caspr/**

Goal:

Learn a **representation** from point cloud inputs that can capture and generate **spatio-temporal** changes in **object properties**.

- Gives **spatiotemporal continuity**

- **Latent ODE** allows "querying" any intermediate timestamp

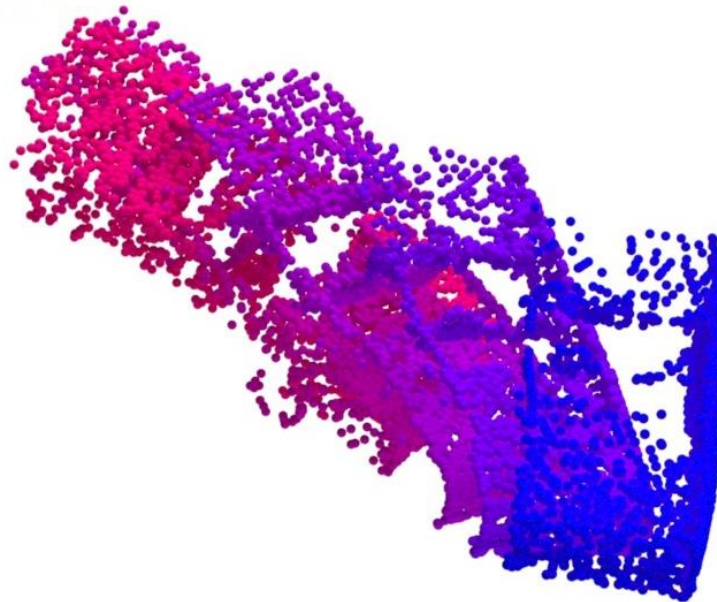- Continuous Normalizing Flow (CNF) allows dense spatial sampling



Flow in space (CNF)

Flow in time (Latent ODE)

Input

Predicted Canonicalization

Predicted Reconstruction

Ground Truth Canonicalized Input

Predicted Spatiotemporal Interpolation

Normalizing Flow Mapped Gaussian

# Perception and Generation of Physical Interactions

Srinath Sridhar

GAMES Seminar

August 12, 2021

- Object detection
  - Faster R-CNN (Ren et al. 2015)
- Instance segmentation
  - FCN (Long et al. 2015)
- Hand segmentation
  - FCN (Long et al. 2015)
- Action
  - Detection and segmentation
  - Wang et al. 2019 (Ours)

Reconstructed

Generated

## Perception: "what is"



learn to digitize

### Requires

- 3D Understanding

- Deeper knowledge of human skills

- Expressive models of objects

- Properties other than shape & appearance

Leo Guibas    Christian Theobalt    Antti Oulasvirta    Hans-Peter Seidel    Niloy Mitra

Karlin Bark / Lee Beckwith / Florian Bernard / Rishabh Bhandari / Sebastian Boring / Sofien Bouaziz / Dan Casas / Anna Maria Feit / Paul Guerrero / Aaron Hertzmann / Judy Hoffman / Jingwei Huang / Krishna Murthy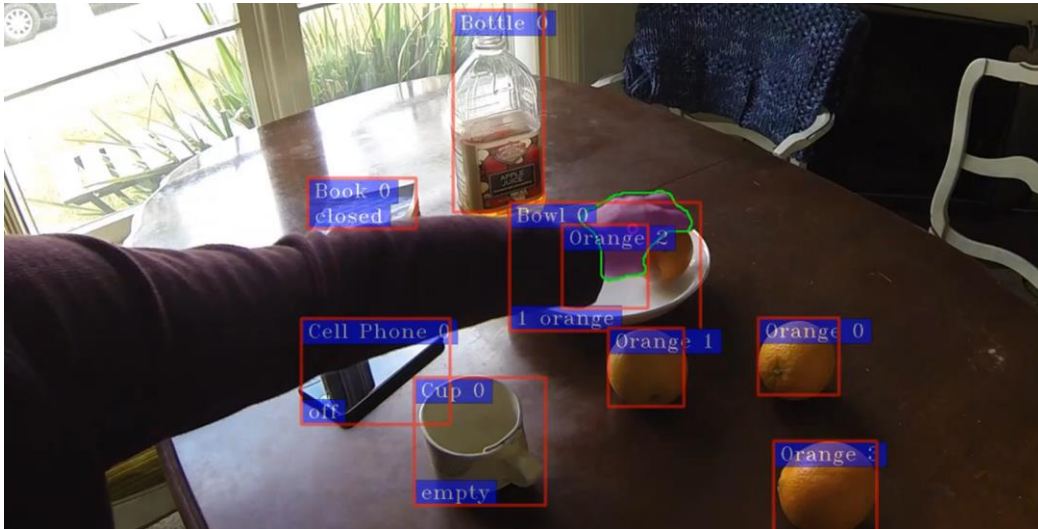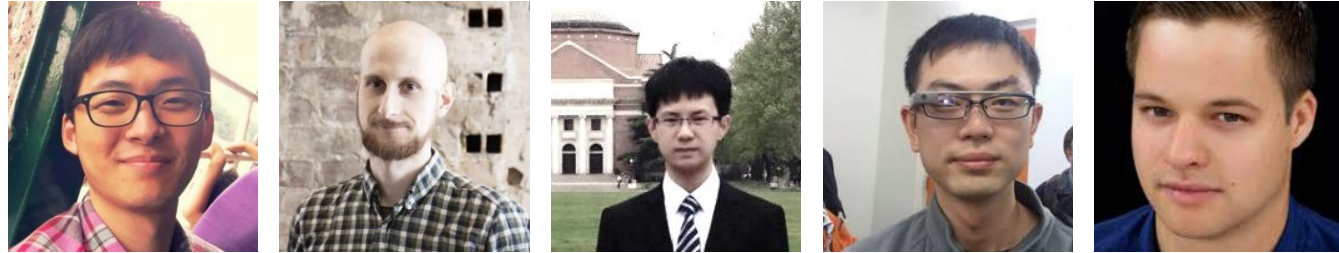 Jatavallabhula / Vladimir Kim / K. Madhava Krishna / Jiahui Lei / Or Litany / Anders Markussen / Dushyant Mehta / Ari Morcos / Franziska Mueller / Victor Ng-Thow-Hing / Soeren Pirk / Gerard Pons-Moll / Davis Rempe / Helge Rhodin / Ozan Sener / Mohammad Shafiei / Rahul Sajnani / AadilMehdi Sanchawala / Shuran Song / Oleksandr Sotnychenko / Minhyuk Sung / Cuong Tran / Julien Valentin / He Wang / Weipeng Xu / Jimei Yang / Zhangsihao Yang / Ersin Yumer / Michael Zollhöfer

BROWN    Stanford    mpii max planck institut informatik

NSF    AIR FORCE RESEARCH LABORATORY

Google    TOYOTA RESEARCH INSTITUTE    aws

**srinathsridhar.com**

Brown
Visual Computing

PhD Students

Kefan Chen
Rao Fu

MS Students

Sijie Ding
Trevor Houchens
Aparna Natarajan

Undergrad

Qiuhong (Anna) Wei
Yiheng Xie

Visitors

Radhika Dua
Shivam Duggal

Jivitesh Jain
Rahul Sajnani
Apoorve Singhal

**Always looking for motivated students!**

Topics:

3D computer vision, deep learning, robotics, graphics, …

srinath@brown.edu

@drsrinathsridha

*ivl.cs.brown.edu*