

Embodied Vision

Physical Modeling: Rigid-body Motion

Tsung-Wei Ke

Spring 2026



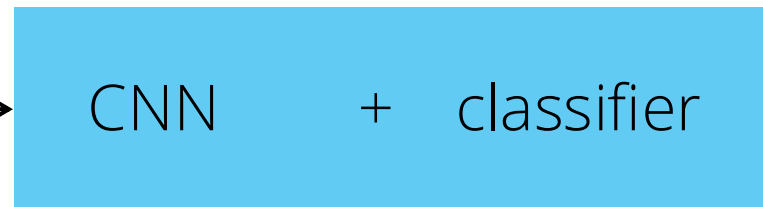
Disclaimer

- This lecture borrows contents heavily from
 - [Physics-based Animation](#) by David I.W Levin at University of Toronto
 - [Physics-based Animation of Solids and Fluids](#) by Minchen Li at CMU

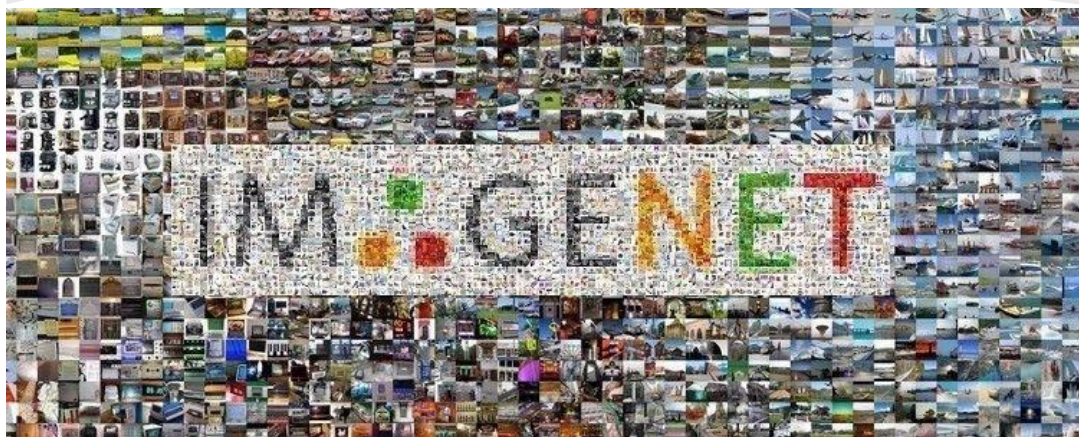
What We've Discussed Still Follows the Data-Driven Paradigm

- We try to follow the same successful story: all you need is to collect more data and labels

179	211	214	209	201	187	192	212	221	231
197	225	223	210	159	147	182	209	219	235
215	230	228	183	103	105	170	206	221	233
232	238	230	155	89	169	173	209	226	231
240	245	224	133	109	198	200	202	229	232
234	244	218	122	66	80	109	187	229	230
218	243	209	124	62	133	82	183	227	224
200	241	207	124	58	55	88	172	211	215
182	237	212	148	100	141	149	190	208	210
170	225	213	163	148	188	221	196	206	210

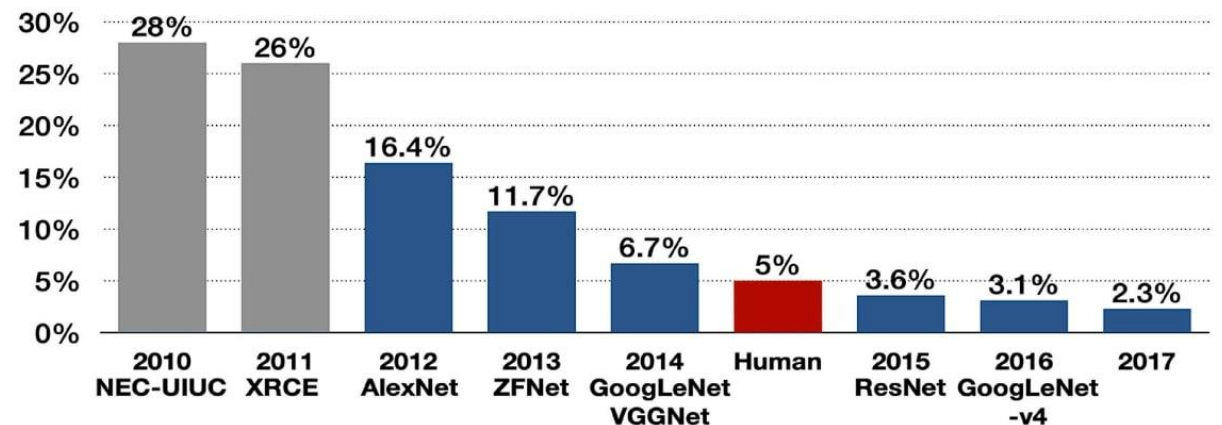


Labels



1,000,000 images with 1,000,000 labels

Top-5 error



Can We Solve Physics with the Data-Driven Paradigm?

- How to annotate the water wave motion?



Can We Solve Physics with the Data-Driven Paradigm?

- How to annotate the elastic-object or fluid motion?



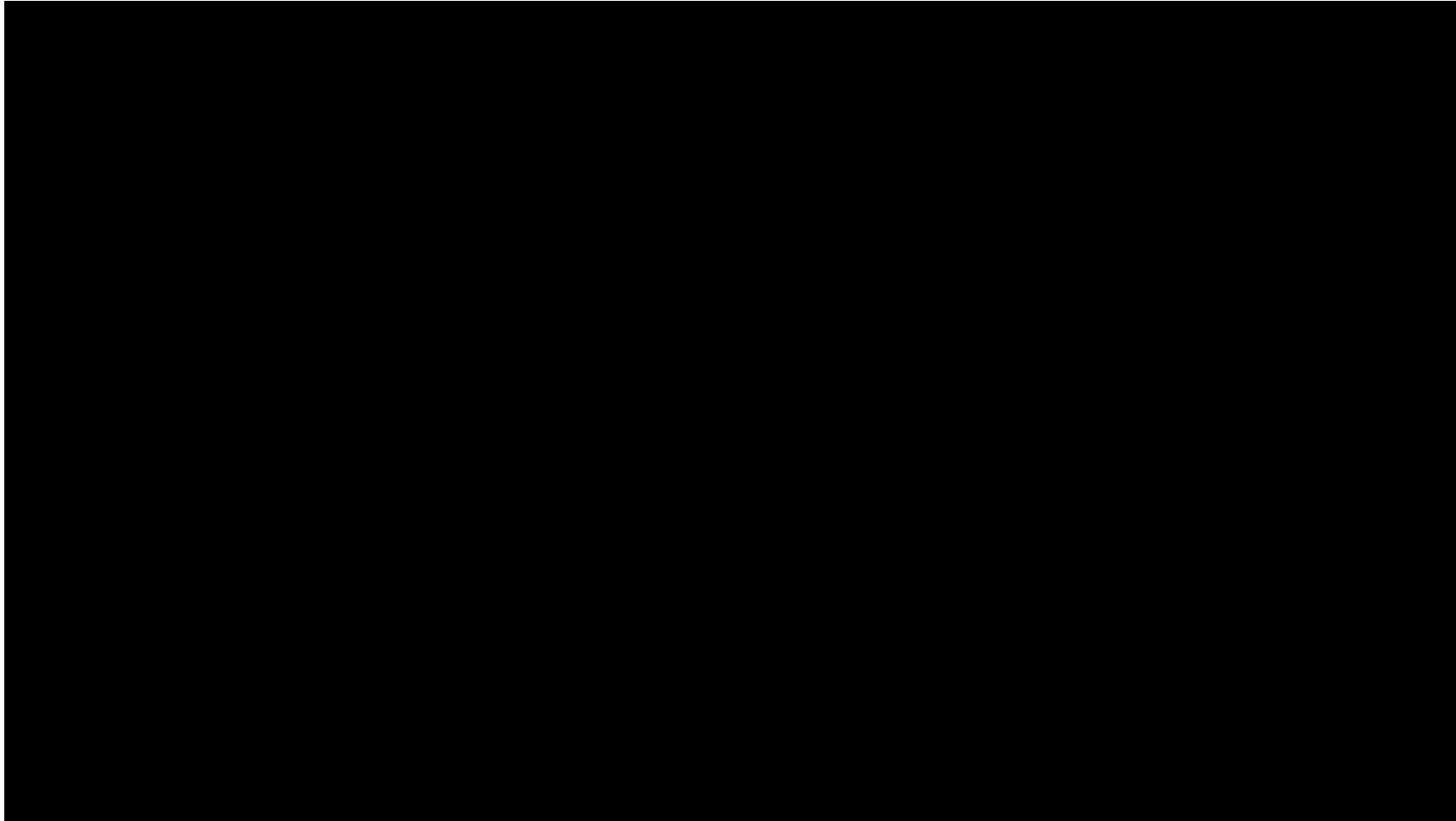
Can We Solve Physics with the Data-Driven Paradigm?

- How to annotate the motion of snow?



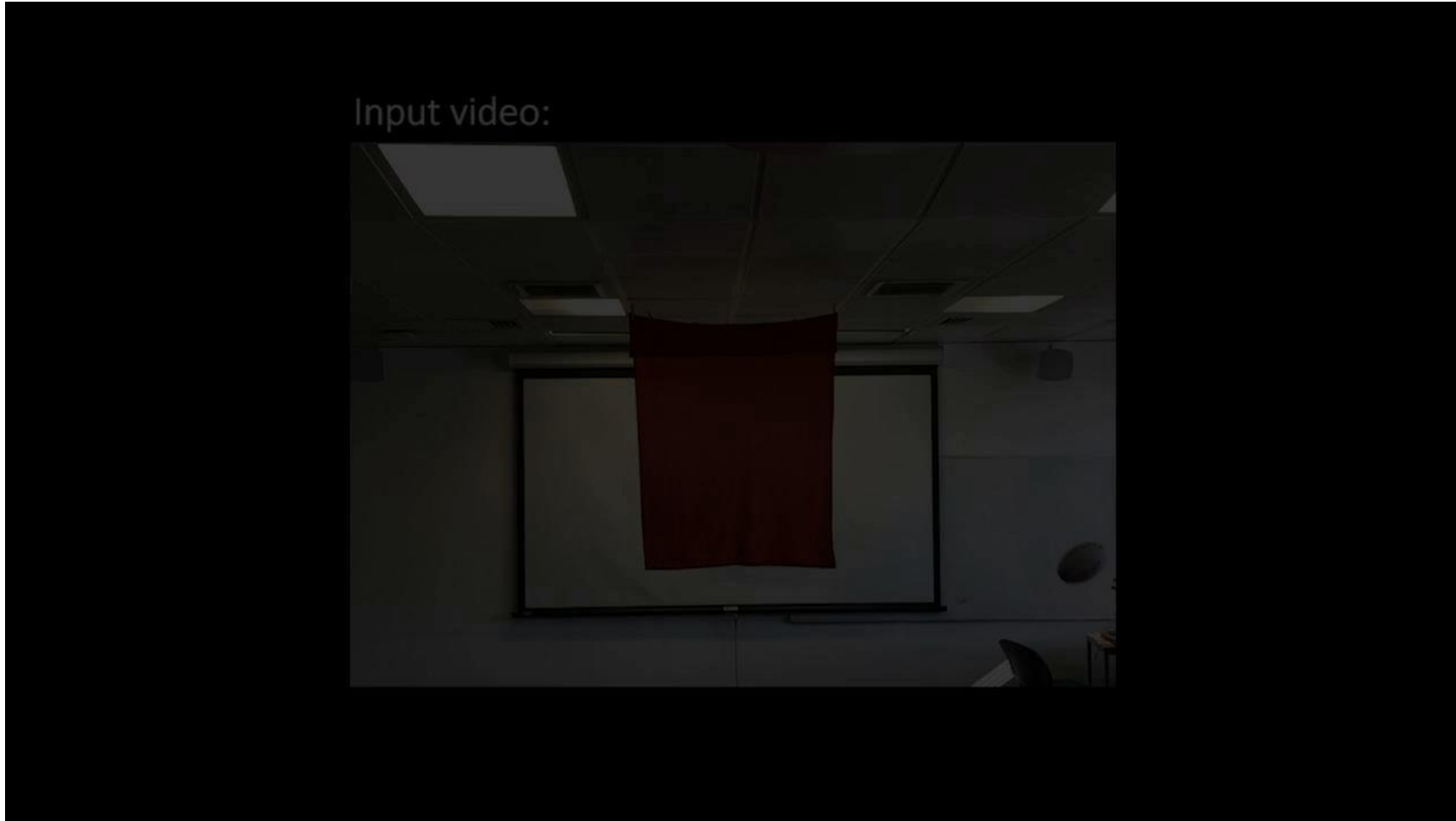
Can We Solve Physics with the Data-Driven Paradigm?

- How to annotate and model the motion with friction



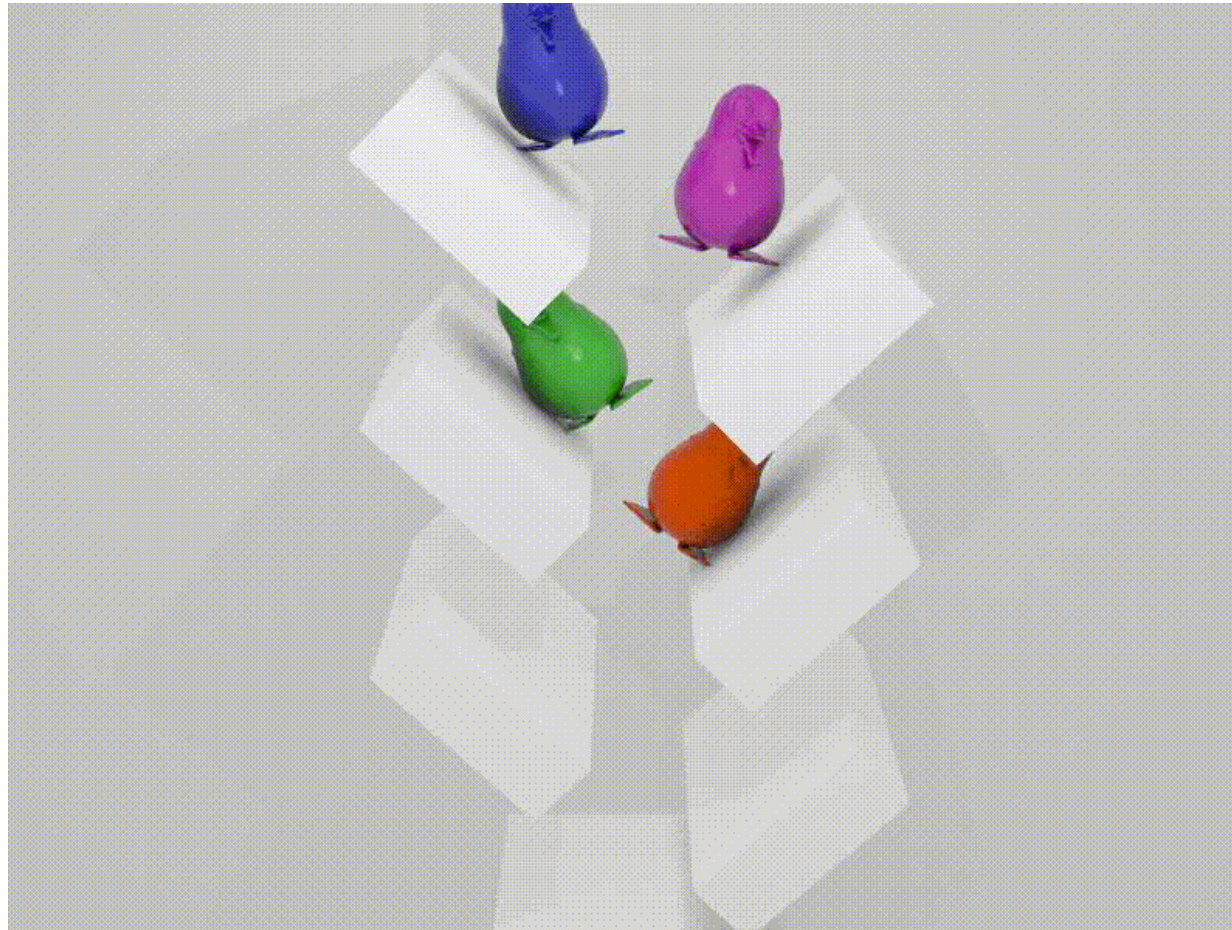
Can We Solve Physics with the Data-Driven Paradigm?

- How to model the rigid-body motion with collision



Can We Solve Physics with the Data-Driven Paradigm?

- How to model the elastic-object motion with collision and friction



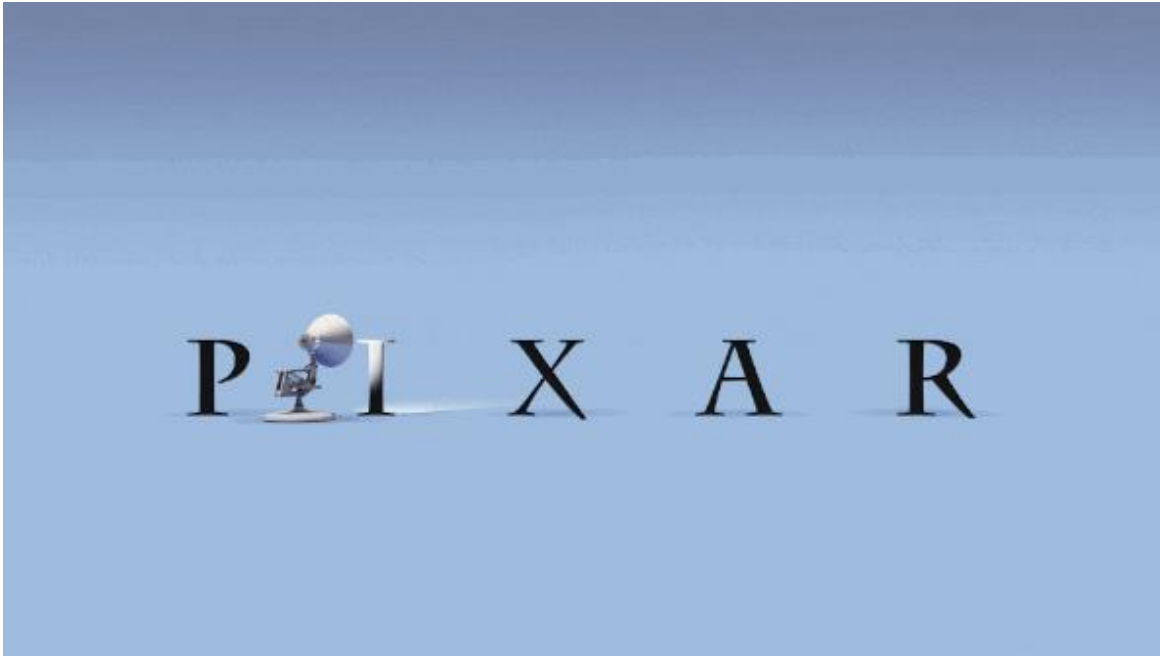
Can We Solve Physics with the Data-Driven Paradigm?

- Supervised learning methods need:
 - Data
 - Labels ← It's almost impossible to label particle motion of monochromatic objects
- Unsupervised learning methods need only data, such as video generative models. However, the dynamics of generated videos are often physically incorrect.

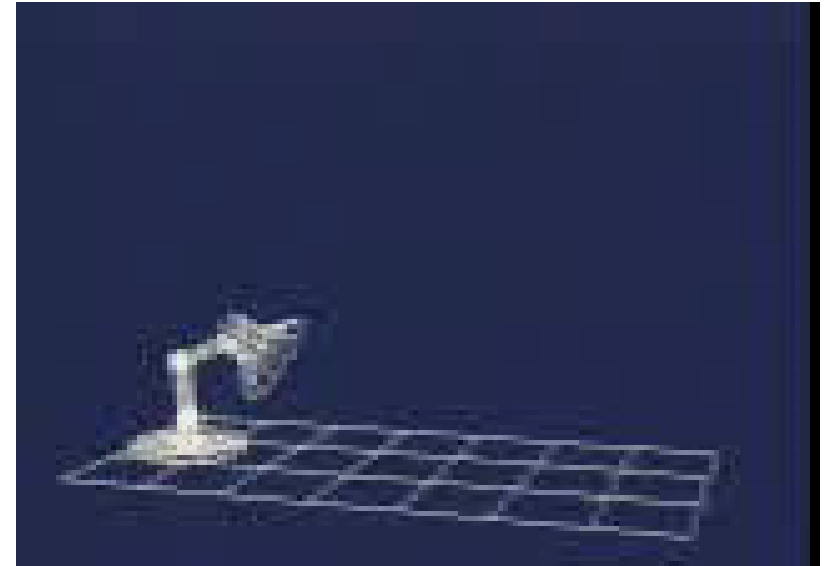


Computer Graphics Has Solutions: Physical-Based Simulation

- Idea: Simulate motions with physic laws



Video source: Pixar



Spacetime Constraints. Witkins and Kass.

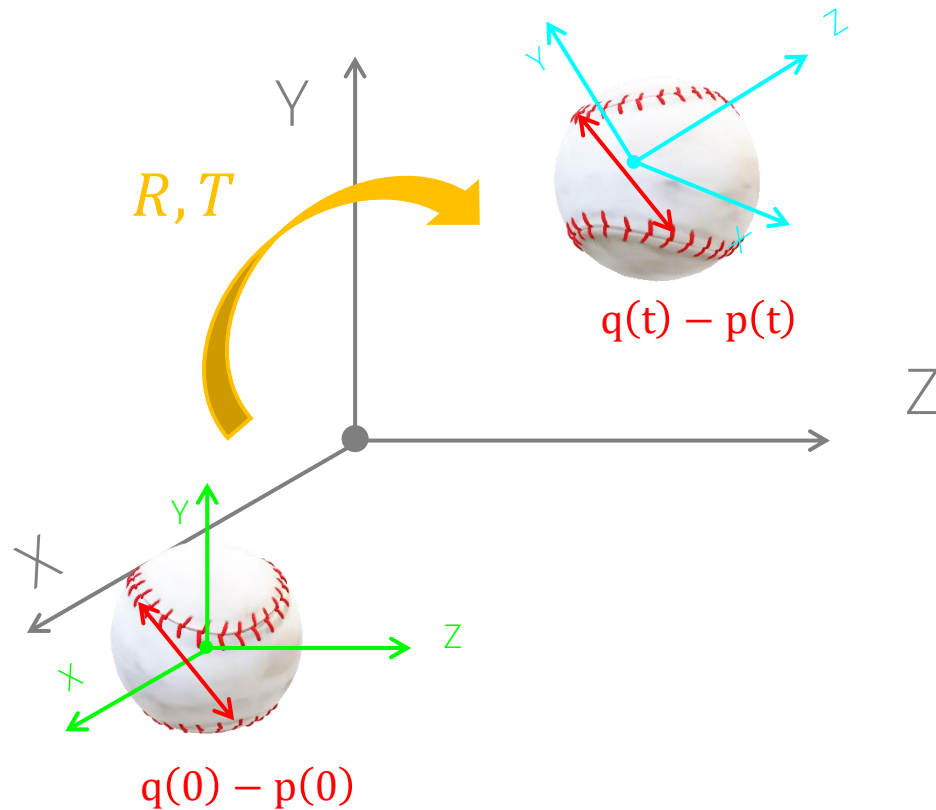
Content

- Rigid-body Modeling
 - Physical law of rigid bodies
 - Formulation of articulated objects
- Dynamic Modeling of Rigid Bodies
 - Video generation of rigid-body motions
 - Generation of 3D articulated objects

Content

- Rigid-body Modeling
 - Physical law of rigid bodies
 - Formulation of articulated objects
- Dynamic Modeling of Rigid Bodies
 - Video generation of rigid-body motions
 - Generation of 3D articulated objects

Rigid-Body Motions



- Rigid body: A rigid body is a collection of particles s.t. the distance between any 2 particles are fixed:

$$\|q(t) - p(t)\| = \|q(0) - p(0)\|$$

- Every particle motion can be described by a rotation matrix \mathbf{R} and a translation vector \mathbf{T}

Rigid-Body Physic Law

Rigid-body physic law:

$$\frac{d}{dt} \mathbf{q}(t) = \frac{d}{dt} \begin{bmatrix} \mathbf{t}(t) \\ \mathbf{R}(t) \\ \boldsymbol{\nu}(t) \\ \boldsymbol{\omega}(t) \end{bmatrix} = \begin{bmatrix} \mathbf{v}(t) \\ \boldsymbol{\omega}(t) \times \mathbf{R}(t) \\ \frac{\mathbf{F}(t)}{M} \\ \mathbf{I}(t)^{-1}(\boldsymbol{\tau} - \boldsymbol{\omega}(t) \times \mathbf{I}(t)\boldsymbol{\omega}(t)) \end{bmatrix}$$

Simulation:

$$\mathbf{q}(t) = \mathbf{q}(0) + \int_0^t \frac{d}{dt} \mathbf{q}(t) dt = \mathbf{q}(0) + \sum_{i=1}^T \frac{d}{dt} \mathbf{q}(t)|_{t=t_i} \Delta t$$

We will cover in the next lecture

Kinematics of Rotation $\frac{d\mathbf{R}(t)}{dt} = \boldsymbol{\omega}(t) \times \mathbf{R}(t)$

- Let \mathbf{r}_{body} be a constant, fixed point on the rigid body relative to its center of mass. Its position in the fixed world frame is $\mathbf{r}_{\text{world}}(t) = \mathbf{R}(t)\mathbf{r}_{\text{body}}$
- If the body is rotating with an angular velocity $\boldsymbol{\omega}(t)$, the linear velocity of that point in the world frame is $\dot{\mathbf{r}}_{\text{world}}(t) = \boldsymbol{\omega}(t) \times \mathbf{r}_{\text{world}}(t)$
- We then have

$$\begin{aligned}\frac{d\mathbf{r}_{\text{world}}}{dt} &= \frac{d\mathbf{R}(t)\mathbf{r}_{\text{body}}}{dt} = \boldsymbol{\omega}(t) \times \mathbf{r}_{\text{world}}(t) = \boldsymbol{\omega}(t) \times \mathbf{R}(t)\mathbf{r}_{\text{body}} \\ &\Rightarrow \dot{\mathbf{R}}(t)\mathbf{r}_{\text{body}} = \boldsymbol{\omega}(t) \times \mathbf{R}(t)\mathbf{r}_{\text{body}} \\ &\Rightarrow \frac{d\mathbf{R}(t)}{dt} = \boldsymbol{\omega}(t) \times \mathbf{R}(t)\end{aligned}$$

Dynamics of Rotation $\frac{d\omega(\mathbf{t})}{dt} = \mathbf{I}(\mathbf{t})^{-1}(\boldsymbol{\tau} - \omega(\mathbf{t}) \times \mathbf{I}(\mathbf{t})\omega(\mathbf{t}))$

- Angular momentum in the world frame is the product of the world-frame inertia tensor $\mathbf{I}(\mathbf{t})$ and the angular velocity $\omega(\mathbf{t})$: $\mathbf{L}(\mathbf{t}) = \mathbf{I}(\mathbf{t})\omega(\mathbf{t})$
- The net torque $\boldsymbol{\tau}$ equals the rate of change of angular momentum

$$\boldsymbol{\tau} = \frac{d\mathbf{L}(\mathbf{t})}{dt} = \dot{\mathbf{I}}(\mathbf{t})\omega(\mathbf{t}) + \mathbf{I}(\mathbf{t})\dot{\omega}(\mathbf{t})$$

Dynamics of Rotation $\frac{d\omega(t)}{dt} = \mathbf{I}(t)^{-1}(\boldsymbol{\tau} - \omega(t) \times \mathbf{I}(t)\omega(t))$

- Angular momentum in the world frame is the product of the world-frame inertia tensor $\mathbf{I}(t)$ and the angular velocity $\omega(t)$: $\mathbf{L}(t) = \mathbf{I}(t)\omega(t)$
- The net torque $\boldsymbol{\tau}$ equals the rate of change of angular momentum

$$\boldsymbol{\tau} = \frac{d\mathbf{L}(t)}{dt} = \dot{\mathbf{I}}(t)\omega(t) + \mathbf{I}(t)\dot{\omega}(t)$$

- The inertia tensor changes as the object rotates in the world frame

$$\mathbf{I}(t) = \mathbf{R}(t)\mathbf{I}_{\text{body}}\mathbf{R}(t)^\top \Rightarrow \frac{d\mathbf{I}(t)}{dt} = \frac{d\mathbf{R}(t)}{dt}\mathbf{I}_{\text{body}}\mathbf{R}(t)^\top + \mathbf{R}(t)\mathbf{I}_{\text{body}}\frac{d\mathbf{R}(t)^\top}{dt}$$

$$\Rightarrow \dot{\mathbf{I}}(t) = \omega(t) \times \mathbf{R}(t)\mathbf{I}_{\text{body}}\mathbf{R}(t)^\top + \mathbf{R}(t)\mathbf{I}_{\text{body}}(\omega(t) \times \mathbf{R}(t))^\top$$

$$\Rightarrow \dot{\mathbf{I}}(t) = \omega(t) \times \mathbf{I}(t) - \mathbf{R}(t)\mathbf{I}_{\text{body}}\mathbf{R}(t)^\top \times \omega(t)$$

$$\Rightarrow \dot{\mathbf{I}}(t) = \omega(t) \times \mathbf{I}(t) - \mathbf{I}(t) \times \omega(t)$$

Dynamics of Rotation $\frac{d\omega(t)}{dt} = I(t)^{-1}(\tau - \omega(t) \times I(t)\omega(t))$

- The inertia tensor changes as the object rotates in the world frame

$$\dot{I}(t) = \omega(t) \times I(t) - I(t) \times \omega(t)$$

$$\Rightarrow \dot{I}(t)\omega(t) = \omega(t) \times I(t)\omega(t) - I(t) \times \omega(t)\omega(t)$$

- Then we have

$$\tau = \dot{I}(t)\omega(t) + I(t)\dot{\omega}(t) = \omega(t) \times I(t)\omega(t) + I(t)\dot{\omega}(t)$$

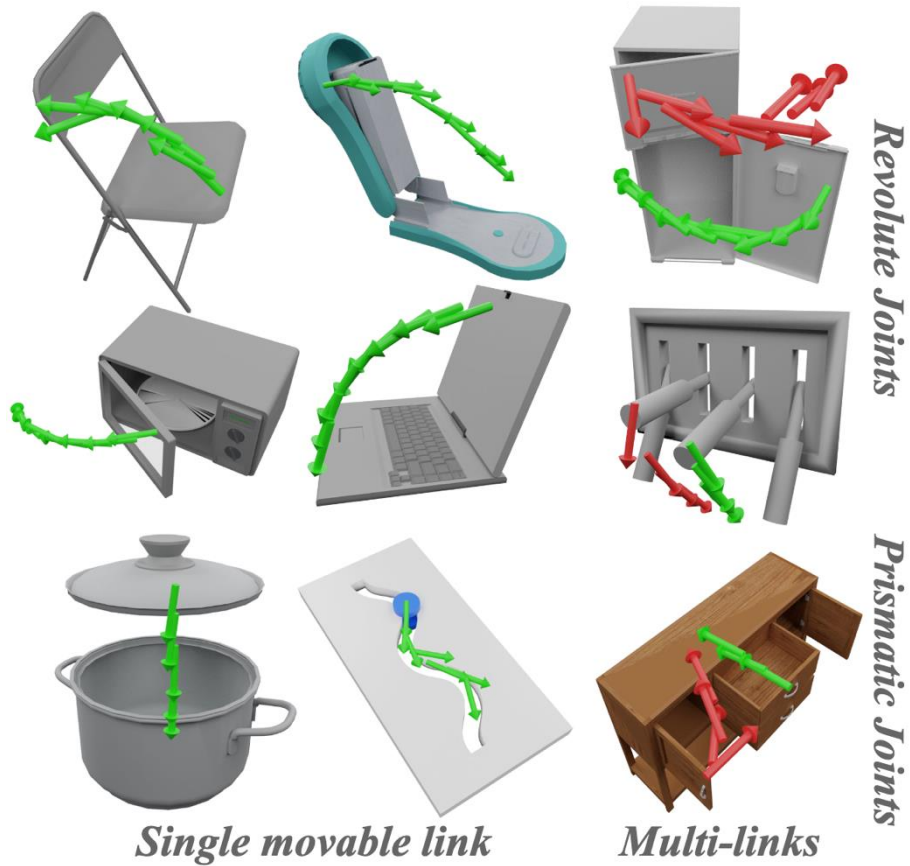
$$\Rightarrow \dot{\omega}(t) = I(t)^{-1}(\tau - \omega(t) \times I(t)\omega(t))$$

Check “A Mathematical Introduction to Robotic Manipulation” by Richard M. Murray, Zexiang Li and S. Shankar Sastry for more details

Content

- Rigid-body Modeling
 - Physical law of rigid bodies
 - Formulation of articulated objects
- Dynamic Modeling of Rigid Bodies
 - Video generation of rigid-body motions
 - Generation of 3D articulated objects

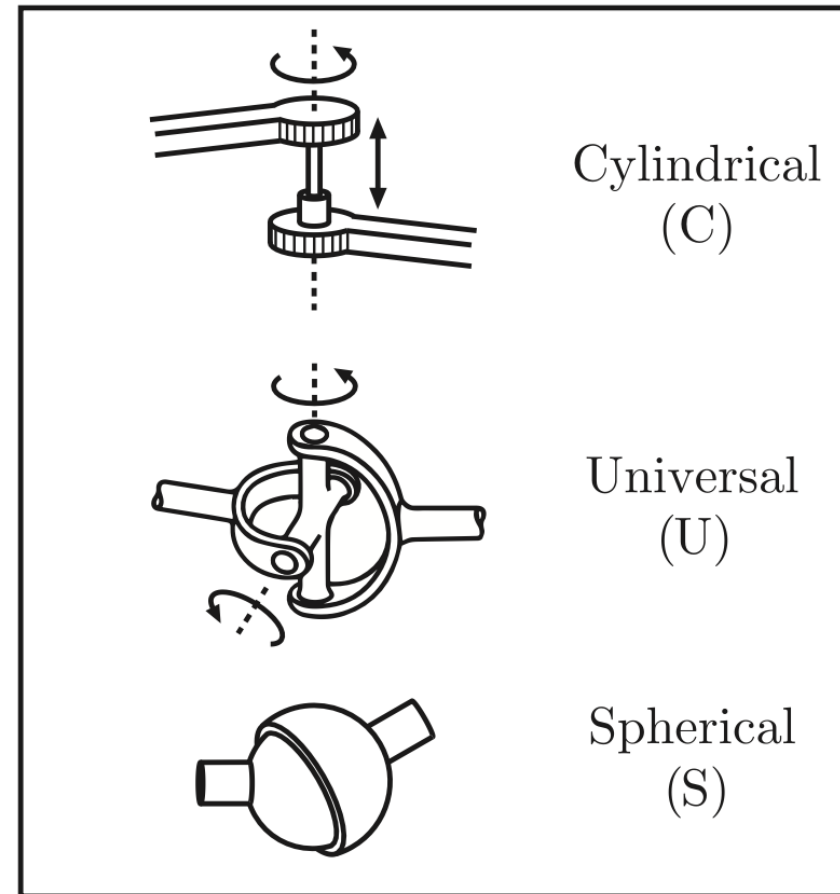
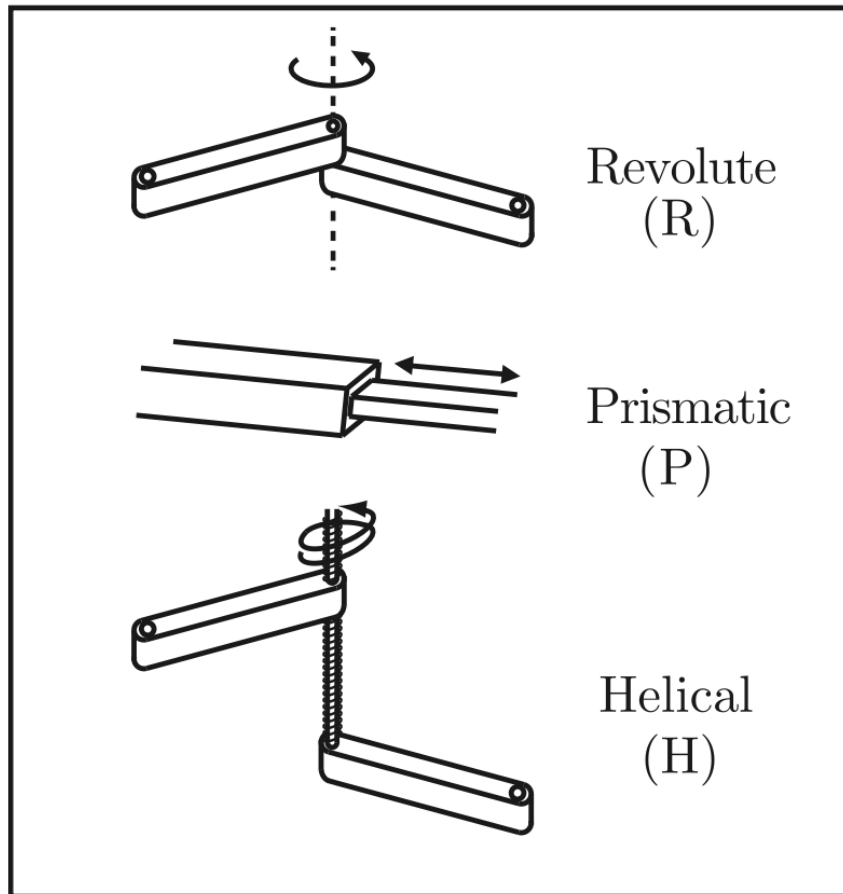
What are Articulated Objects



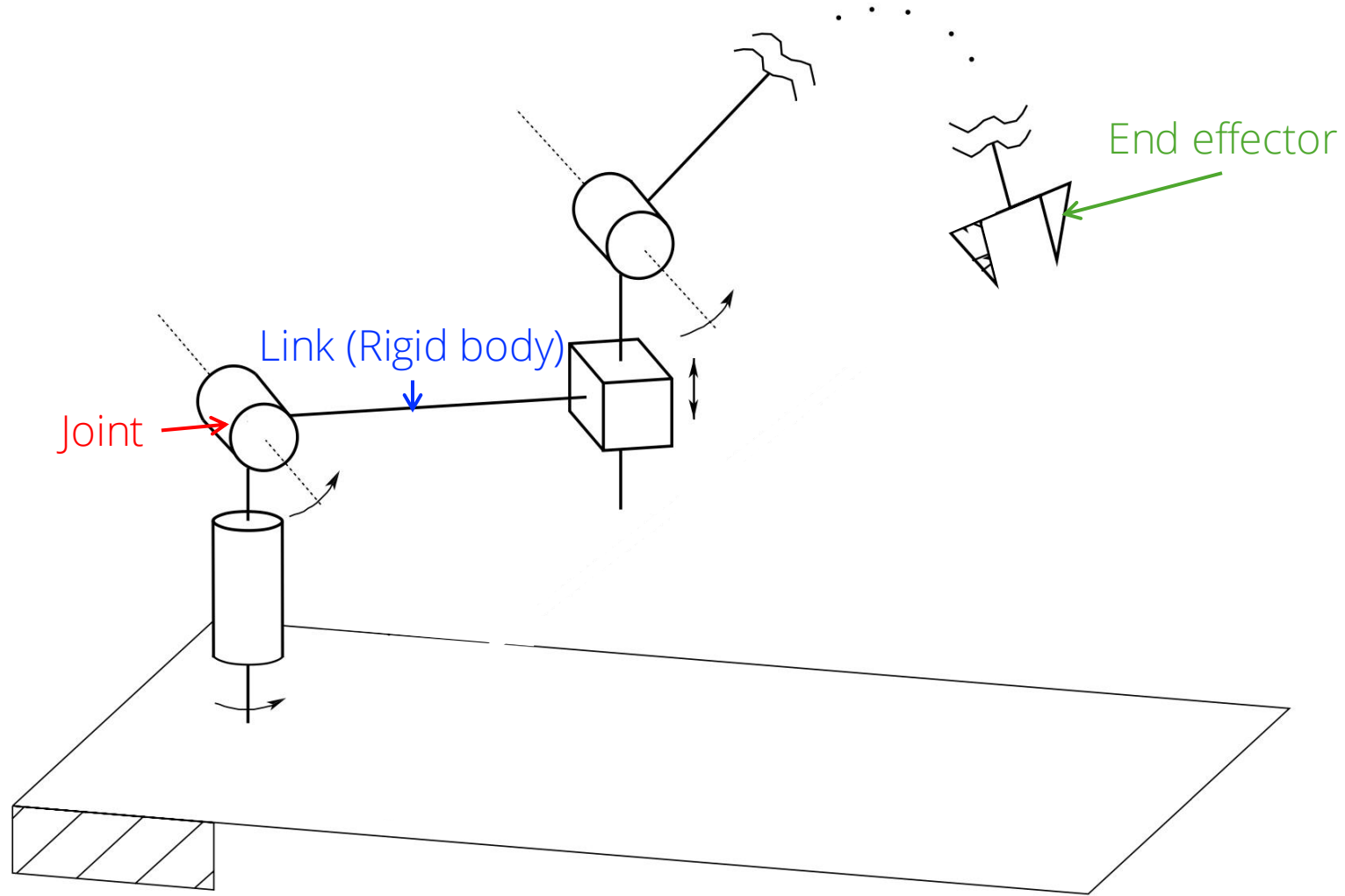
- Articulated objects: objects composed of multiple rigid bodies which move relative to one another
- Their relative motions are characterized by the “joint”
- We can describe the motion of one part by chaining all relative motions of its root parts

Typical Joints

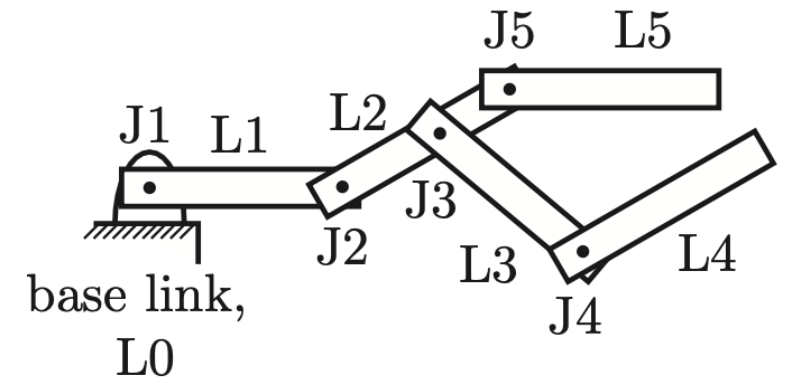
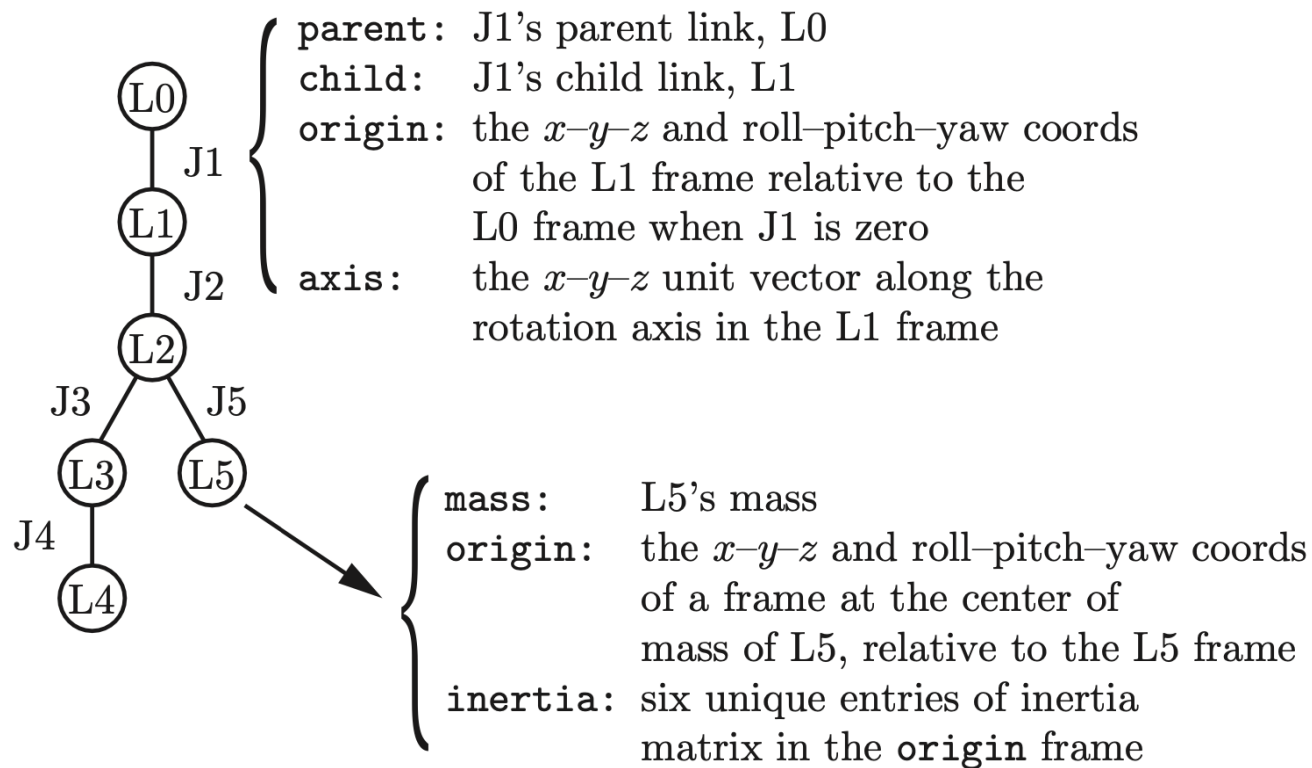
Relative motions of parts can be described by the configuration of each joint (e.g. the angle of the joint)



Common Robots are Articulated Objects



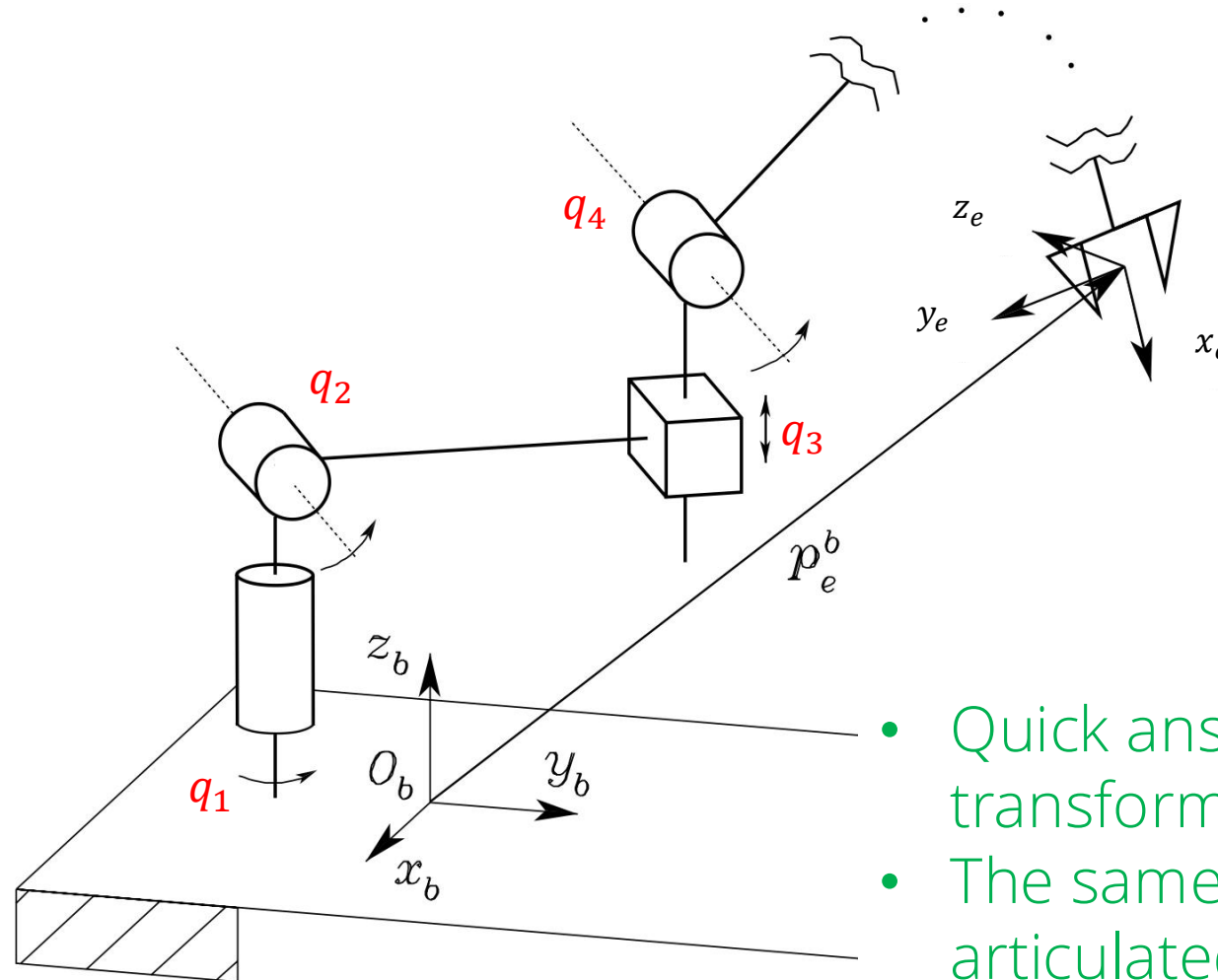
Articulated Objects are a Tree of Rigid Bodies



We need to generate:

1. The connectivity graph of parts
2. The geometry and appearance of parts
3. The kinematics of connection (positions, directions and types of joints)

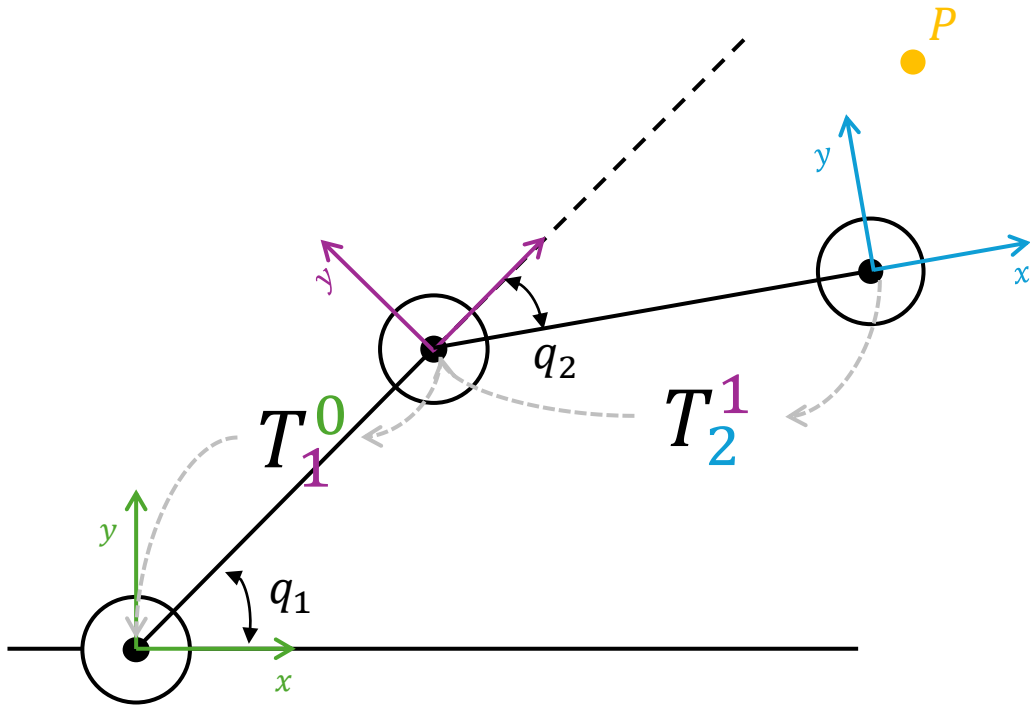
What is the Pose of a Robot's Part given Its Configuration?



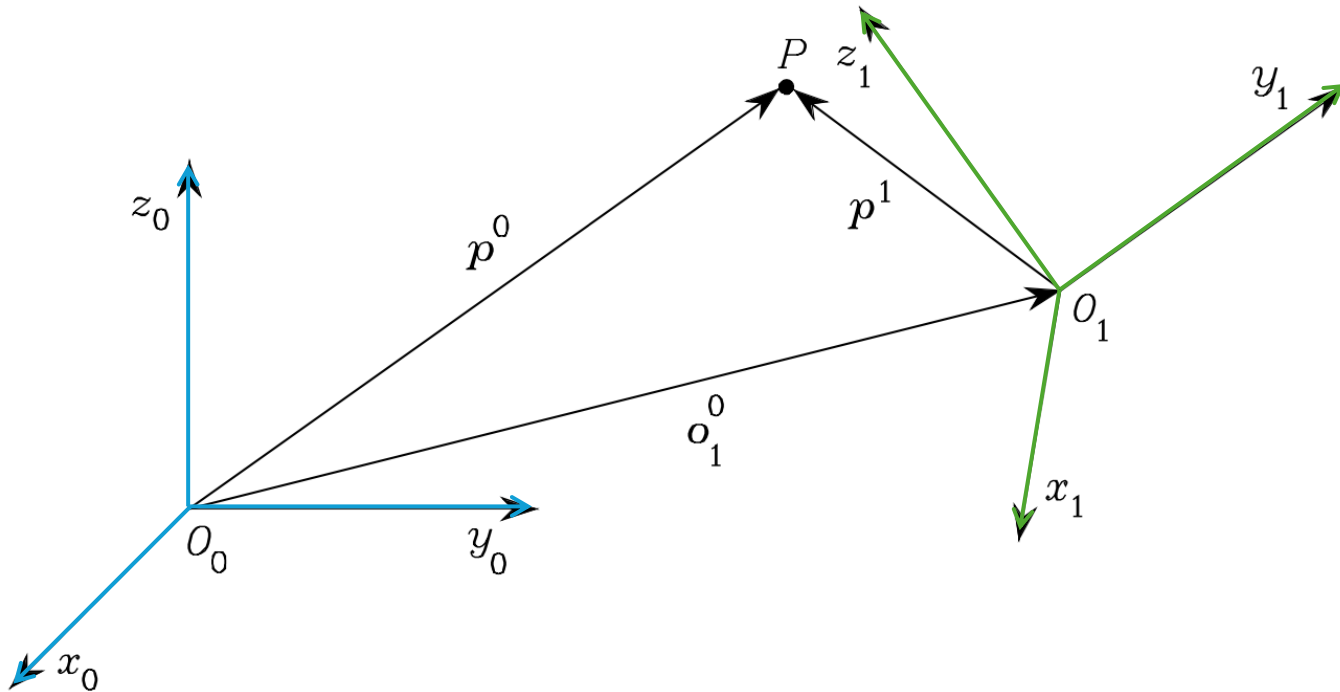
- Quick answer: chaining relative transformation across root parts
- The same idea applies to all articulated objects

How to Calculate the Pose of a Part given the Configuration?

- We can consider each joint applies coordinate transformation



Rigid-Body Transformation



The coordinate of point
P in frame O_0

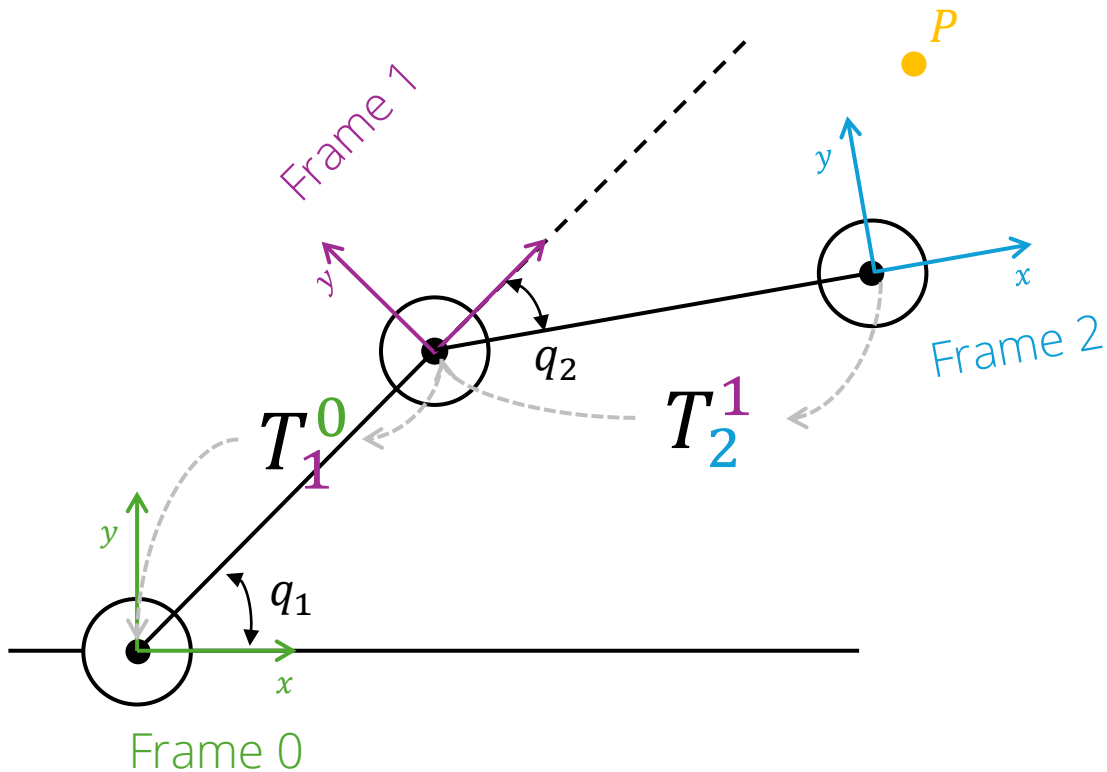
The coordinate of point
P in frame O_1

$$\begin{bmatrix} p^0 \\ 1 \end{bmatrix} = \begin{bmatrix} R_1^0 & o_1^0 \\ \mathbf{0}^T & 1 \end{bmatrix} \begin{bmatrix} p^1 \\ 1 \end{bmatrix}$$

Spatial Transformation
between frame O_1 and O_0

How to Calculate the Pose of a Part given the Configuration?

- We can consider each joint applies coordinate transformation

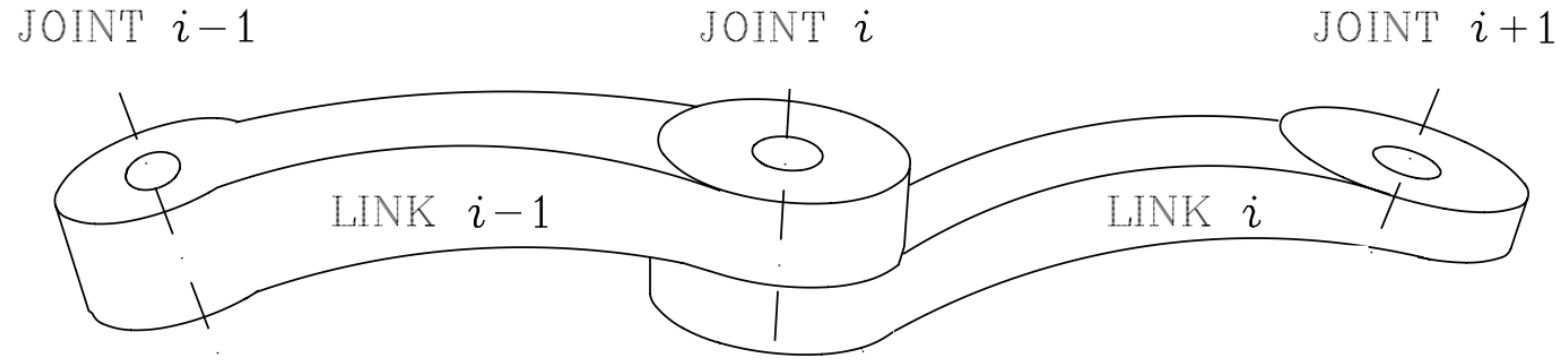


Configuration of Link n

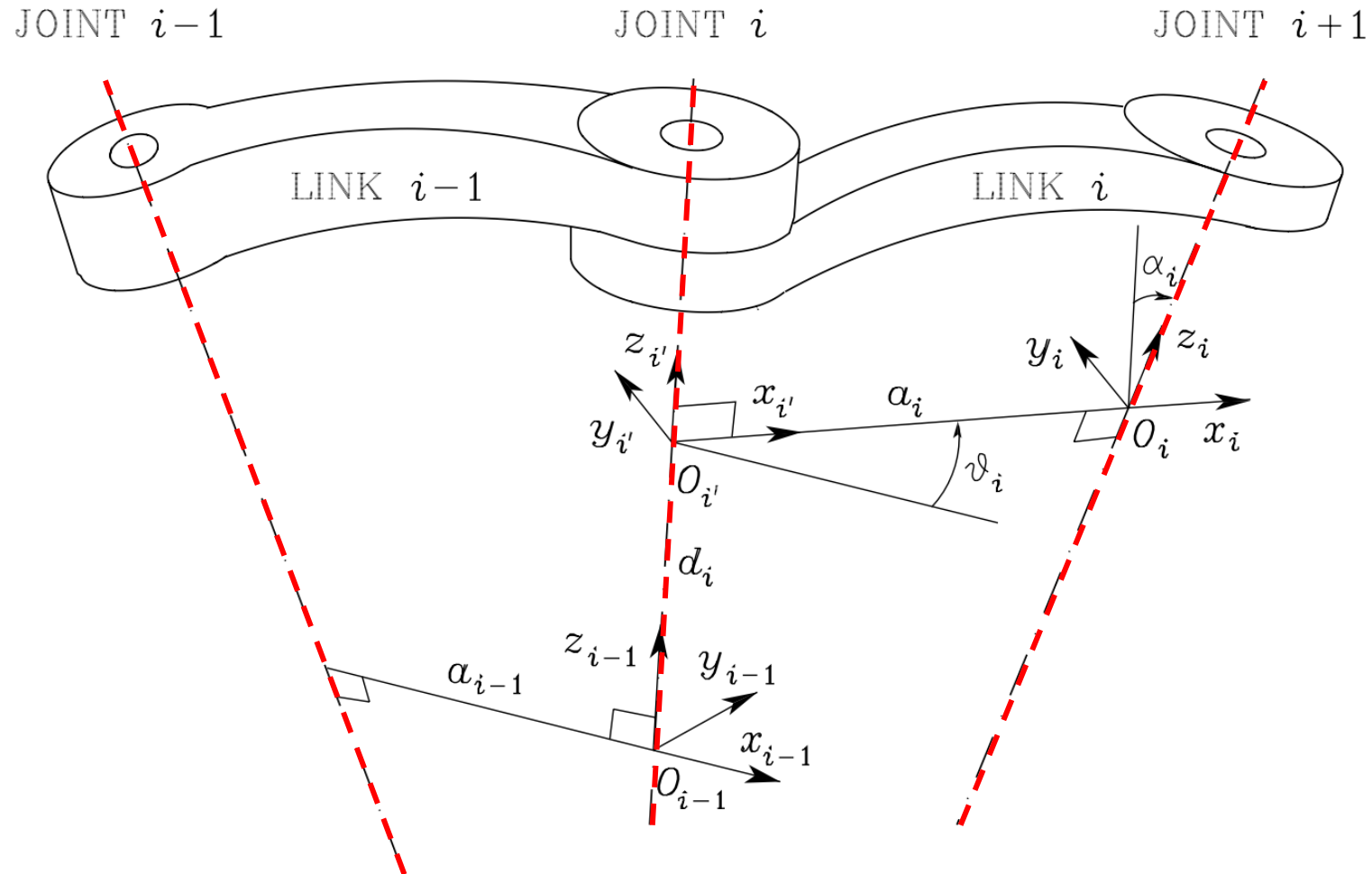
$$T_n^W(\mathbf{q}) = T_0^W T_1^0(q_1) \dots T_n^{n-1}(q_n)$$

Coordinate transformation from
Frame 0 to the world frame

How to Determine Frames Attached to the Two Links?

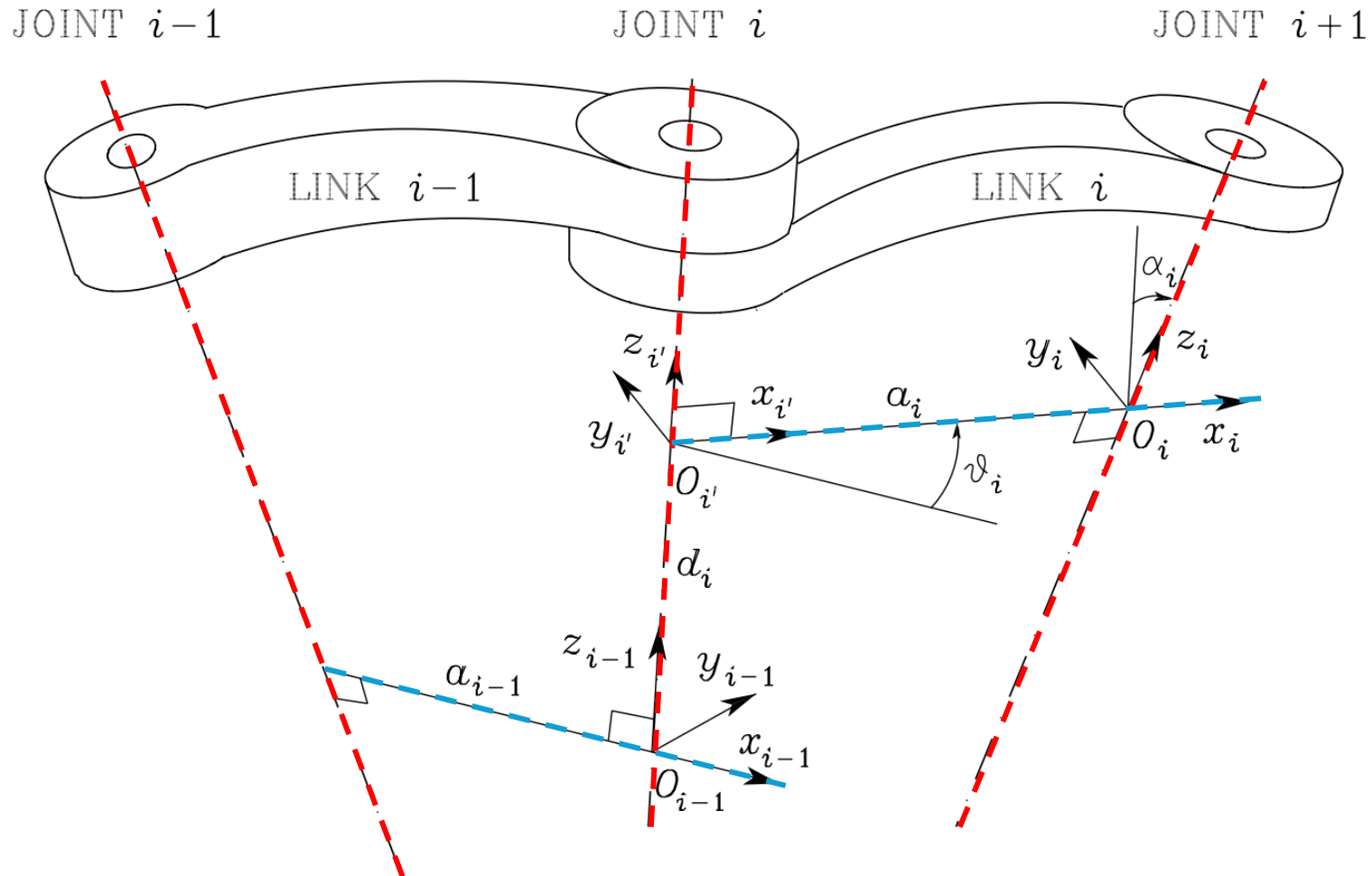


Decide Frames with Denavit–Hartenberg Convention



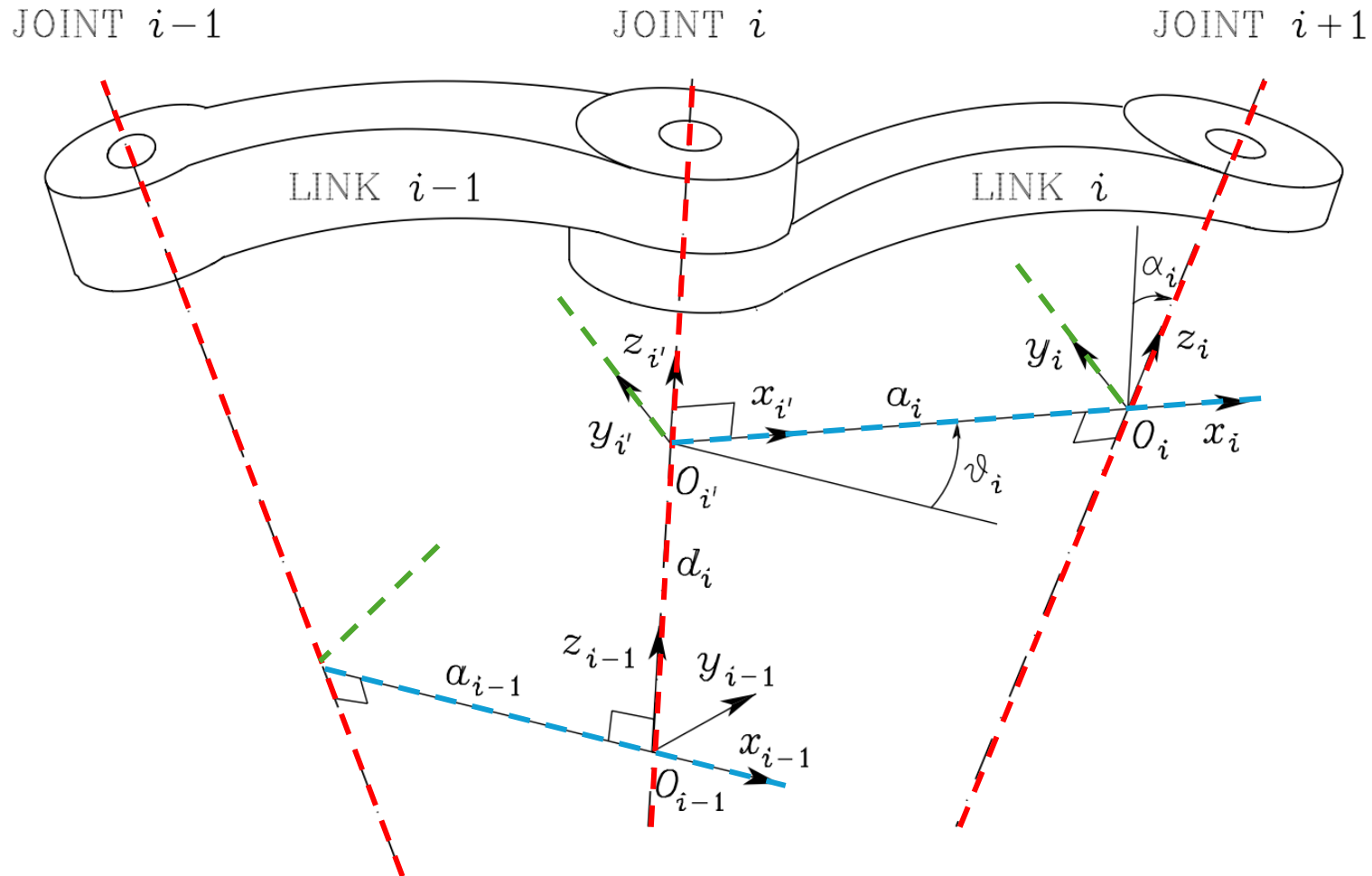
1. Choose axis z_i along the axis of Joint $i + 1$

Decide Frames with Denavit–Hartenberg Convention



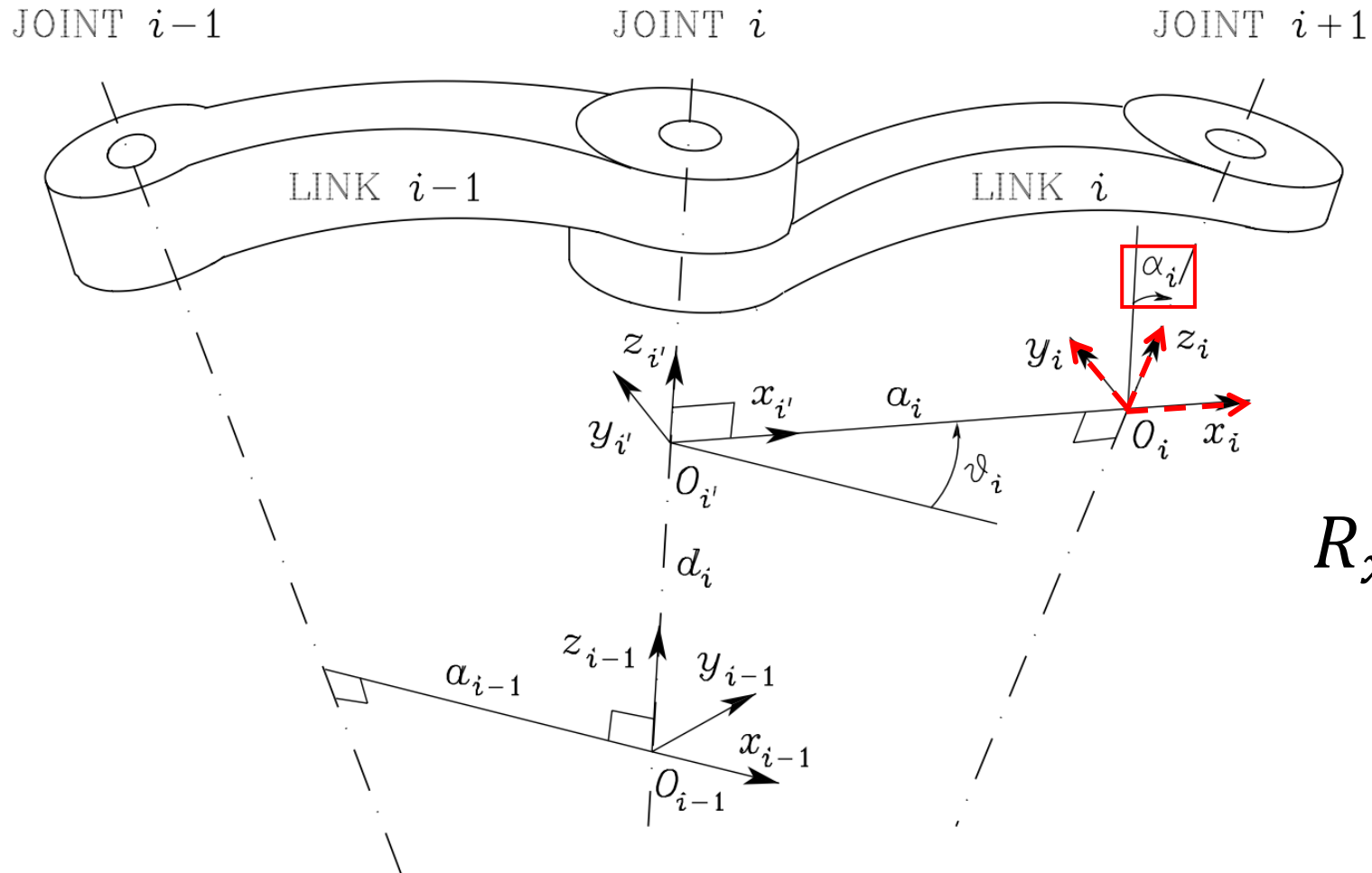
1. Choose axis z_i along the axis of Joint $i + 1$
2. Choose axis x_i along the direction of the common normal (vector of minimum distance between axis z_i and z_{i-1})

Decide Frames with Denavit–Hartenberg Convention



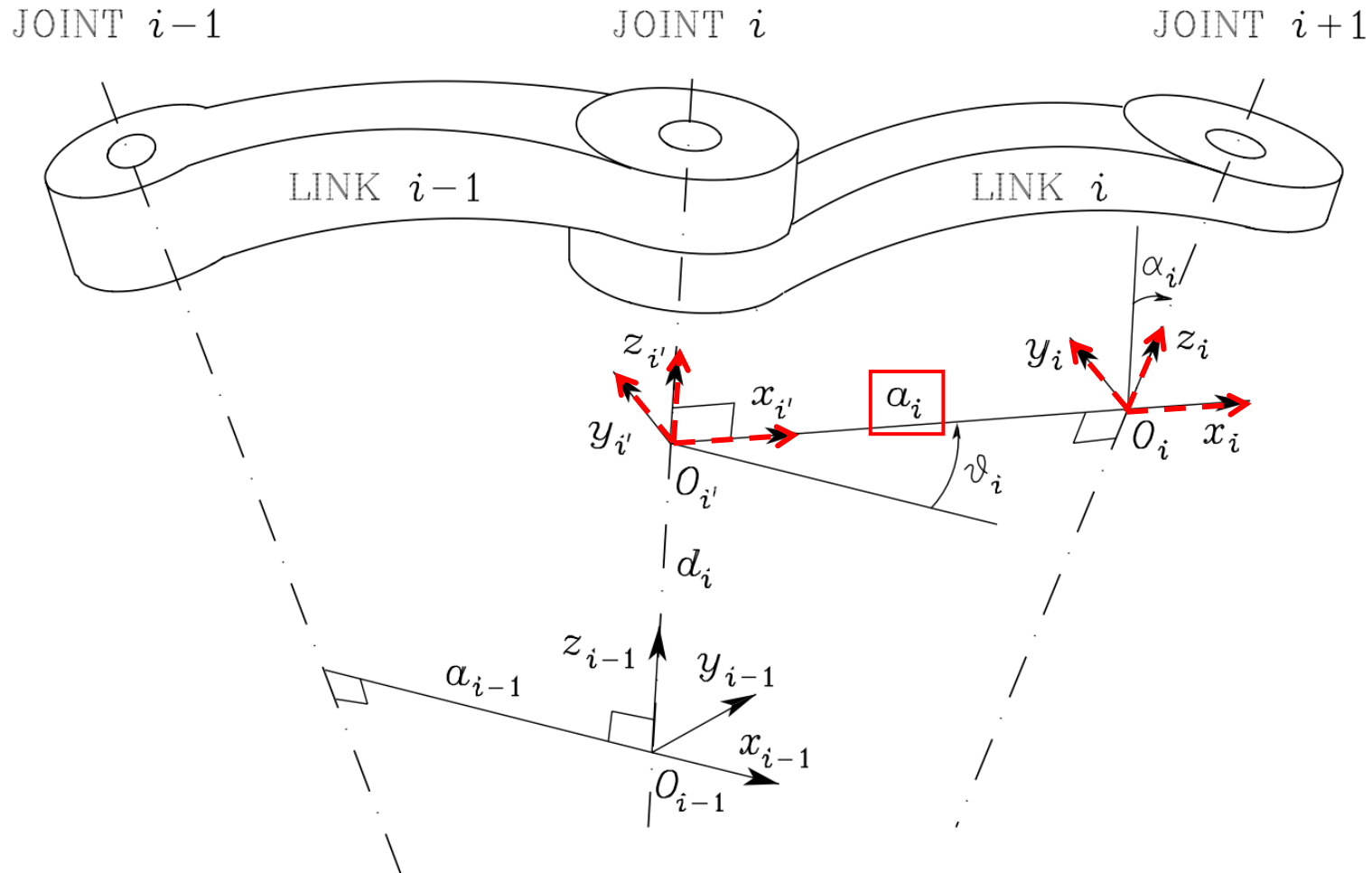
1. Choose axis z_i along the axis of Joint $i + 1$
2. Choose axis x_i along the direction of the common normal (vector of minimum distance between axis z_i and z_{i-1})
3. Axis y_i is decided by right-hand rule

What are the Spatial Transformations under Denavit–Hartenberg Convention?



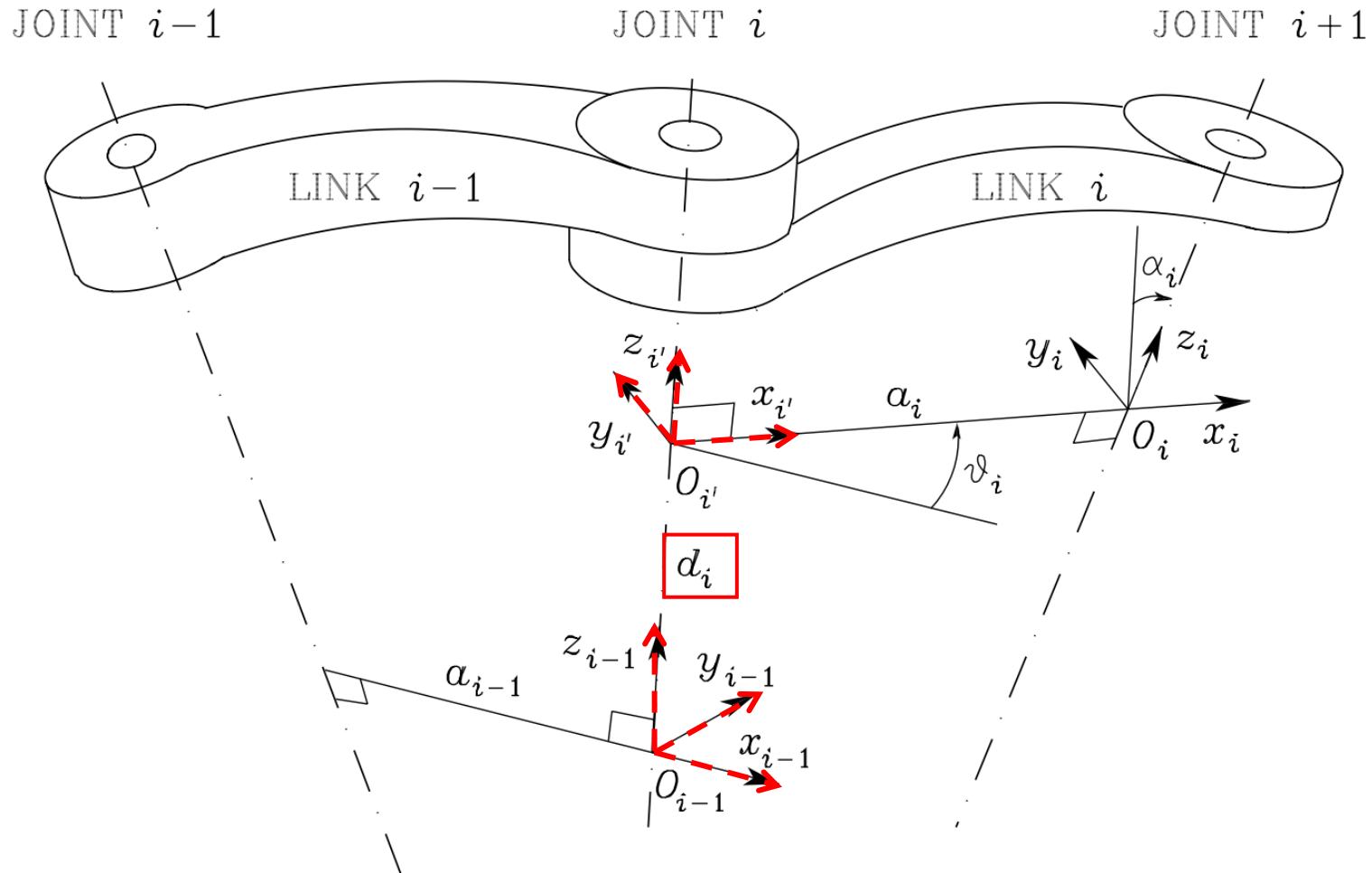
$$R_{x, \alpha_i} = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & c_{\alpha_i} & -s_{\alpha_i} & 0 \\ 0 & s_{\alpha_i} & c_{\alpha_i} & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix}$$

What are the Spatial Transformations under Denavit–Hartenberg Convention?



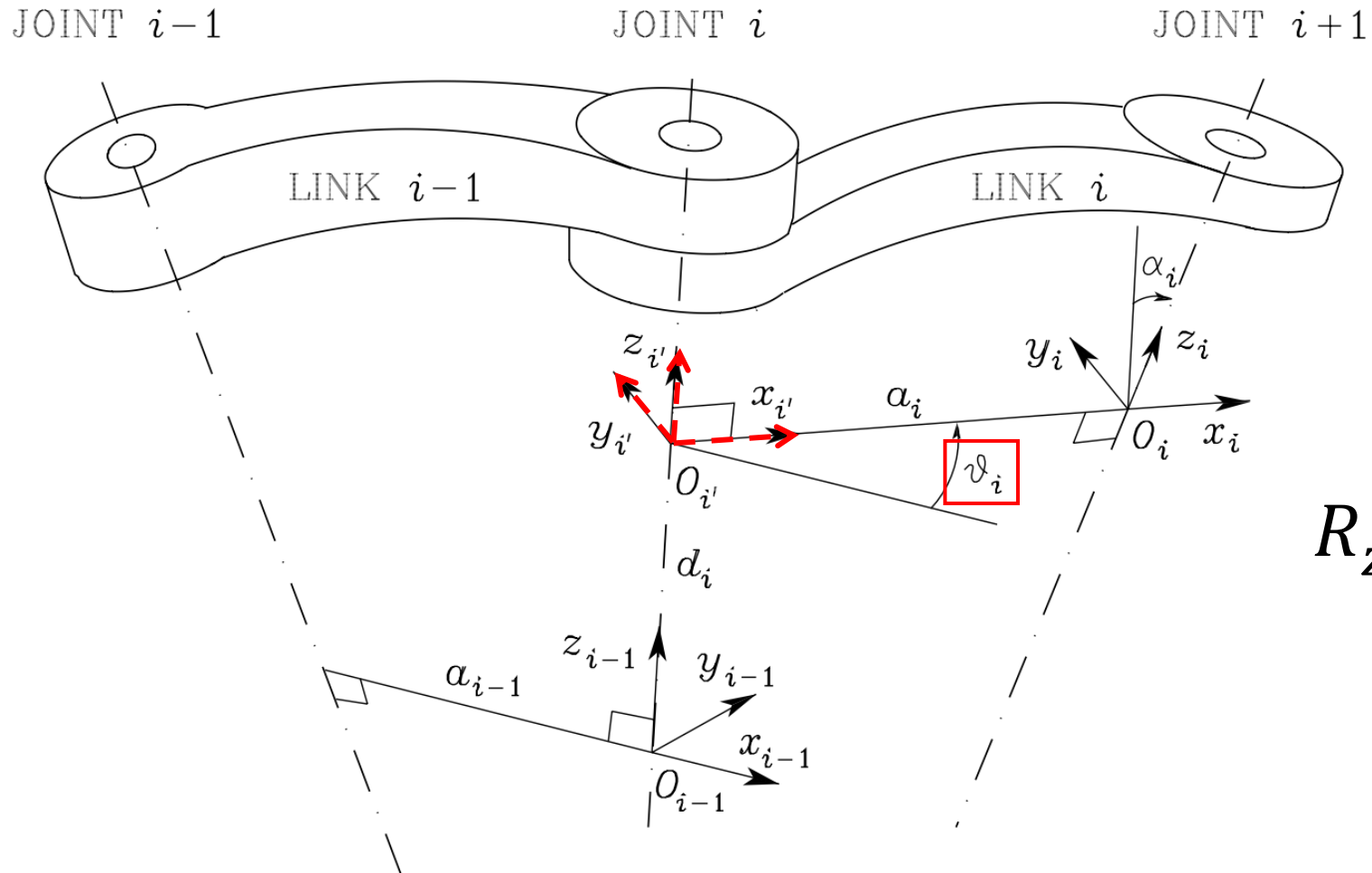
$$T_i^{i'} = \begin{bmatrix} 1 & 0 & 0 & a_i \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix}$$

What are the Spatial Transformations under Denavit–Hartenberg Convention?



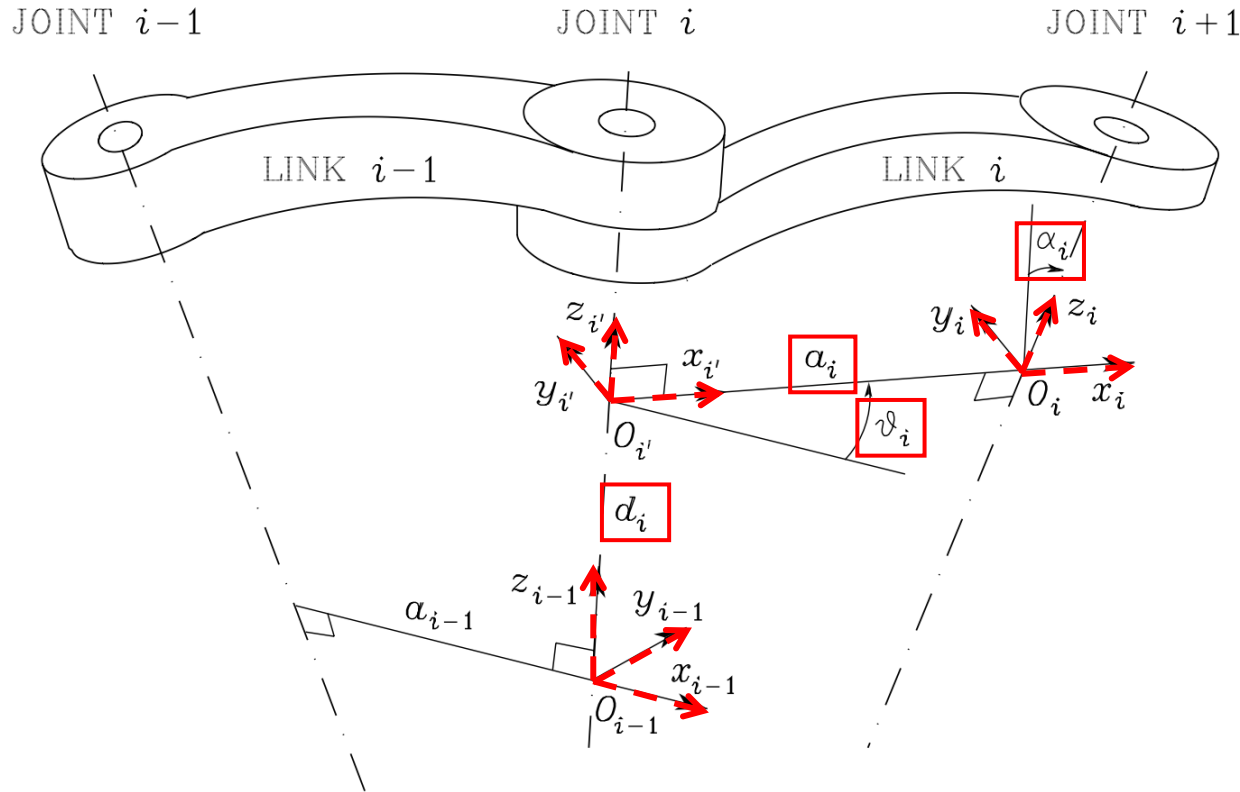
$$T_{i'}^{i-1} = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & d_i \\ 0 & 0 & 0 & 1 \end{bmatrix}$$

What are the Spatial Transformations under Denavit–Hartenberg Convention?



$$R_{z, \vartheta_i} = \begin{bmatrix} C_{\vartheta_i} & -S_{\vartheta_i} & 0 & 0 \\ S_{\vartheta_i} & C_{\vartheta_i} & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix}$$

What are the Spatial Transformations under Denavit–Hartenberg Convention?



$$A_i^{i-1} = R_{z, \vartheta_i} T_{i'}^{i-1} T_i^{i'} R_{x, \alpha_i}$$

$$= \begin{bmatrix} c_{\vartheta_i} & -s_{\vartheta_i} c_{\alpha_i} & s_{\vartheta_i} s_{\alpha_i} & a_i c_{\vartheta_i} \\ s_{\vartheta_i} & c_{\vartheta_i} c_{\alpha_i} & -c_{\vartheta_i} s_{\alpha_i} & a_i s_{\vartheta_i} \\ 0 & s_{\alpha_i} & c_{\alpha_i} & d_i \\ 0 & 0 & 0 & 1 \end{bmatrix}$$

Forward Kinematics: Calculate the Pose of a Part given the Configuration

- We know:

Transformation from frame i to frame $i - 1$ ← $A_i^{i-1} = \begin{bmatrix} c_{\vartheta_i} & -s_{\vartheta_i}c_{\alpha_i} & s_{\vartheta_i}s_{\alpha_i} & a_i c_{\vartheta_i} \\ s_{\vartheta_i} & c_{\vartheta_i}c_{\alpha_i} & -c_{\vartheta_i}s_{\alpha_i} & a_i s_{\vartheta_i} \\ 0 & s_{\alpha_i} & c_{\alpha_i} & d_i \\ 0 & 0 & 0 & 1 \end{bmatrix}$

Transformation from frame 0 to the world frame

Transformation from the frame n to the world frame

← $A_n^w(\mathbf{q}) = A_0^w A_1^0(q_1) \dots A_n^{n-1}(q_n)$

Configuration of joint n

$$\mathbf{p}_n^w = A_0^w A_1^0(q_1) \dots A_n^{n-1}(q_n) \mathbf{p}$$

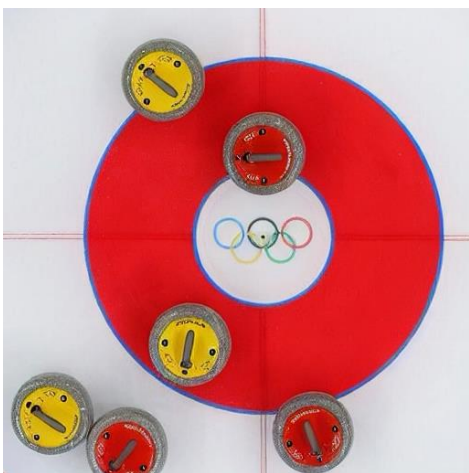
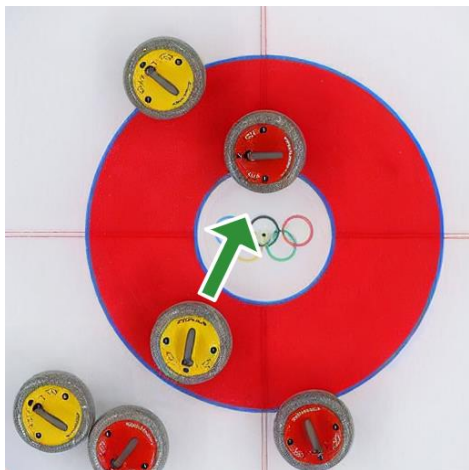
Content

- Rigid-body Modeling
 - Physical law of rigid bodies
 - Formulation of articulated objects
- Dynamic Modeling of Rigid Bodies
 - Video generation of rigid-body motions
 - Generation of 3D articulated objects

Physically Accurate Video Generator



Idea: Integrate Video Generative Model with a Rigid-Body Motion Simulator

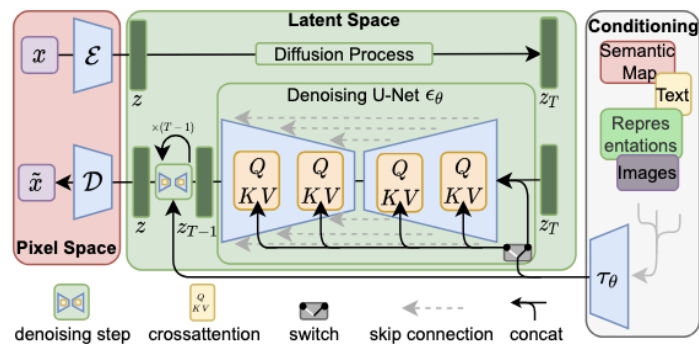


Rigid-body physic law:

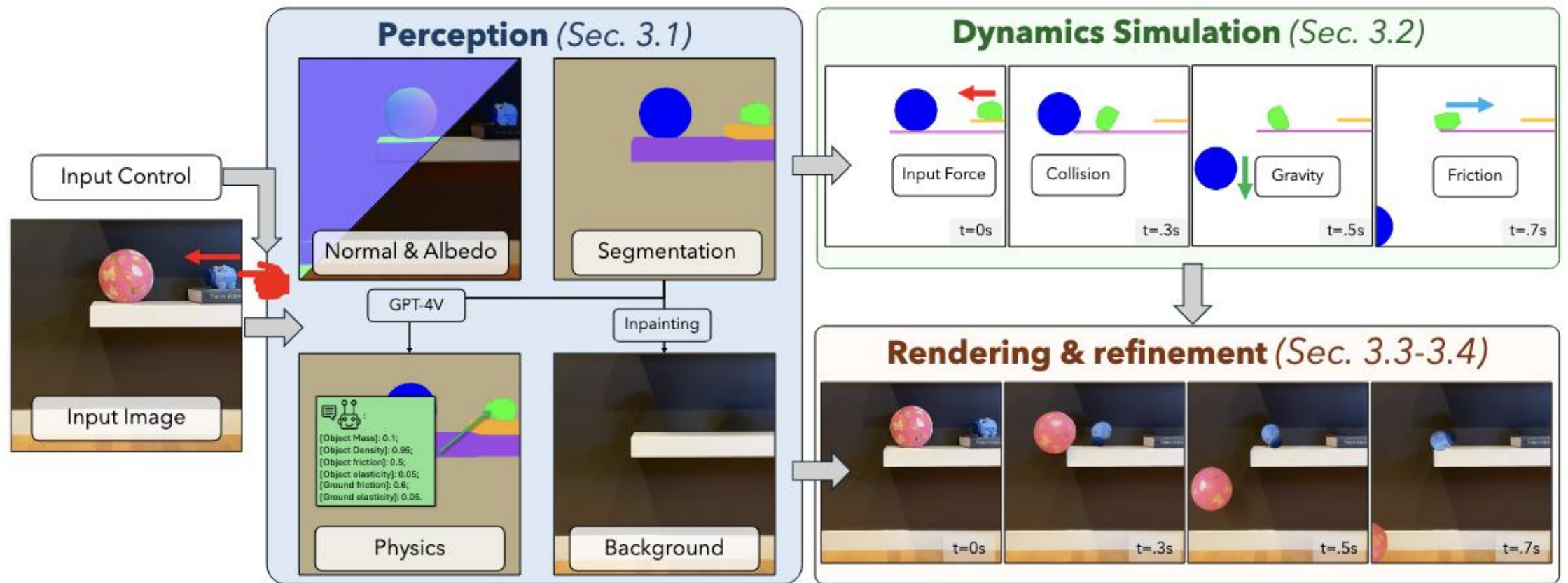
$$\frac{d}{dt} \mathbf{q}(t) = \frac{d}{dt} \begin{bmatrix} \mathbf{t}(t) \\ \mathbf{R}(t) \\ \boldsymbol{\nu}(t) \\ \boldsymbol{\omega}(t) \end{bmatrix} = \begin{bmatrix} \mathbf{v}(t) \\ \boldsymbol{\omega}(t) \times \mathbf{R}(t) \\ \frac{\mathbf{F}(t)}{M} \\ \mathbf{I}(t)^{-1}(\boldsymbol{\tau} - \boldsymbol{\omega}(t) \times \mathbf{I}(t)\boldsymbol{\omega}(t)) \end{bmatrix}$$

Simulation:

$$\mathbf{q}(t) = \mathbf{q}(0) + \int_0^t \frac{d}{dt} \mathbf{q}(t) dt = \mathbf{q}(0) + \sum_{i=1}^T \frac{d}{dt} \mathbf{q}(t)|_{t=t_i} \Delta t$$



PhysGen: Rigid-Body Physics-Grounded Image-to-Video Generation



Results: Interactive Video Generation

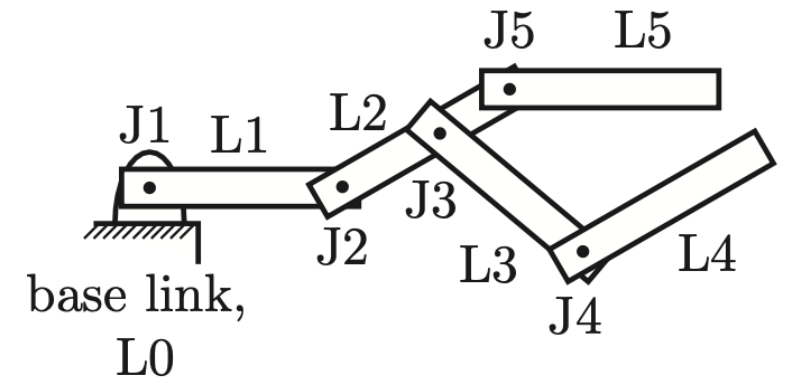
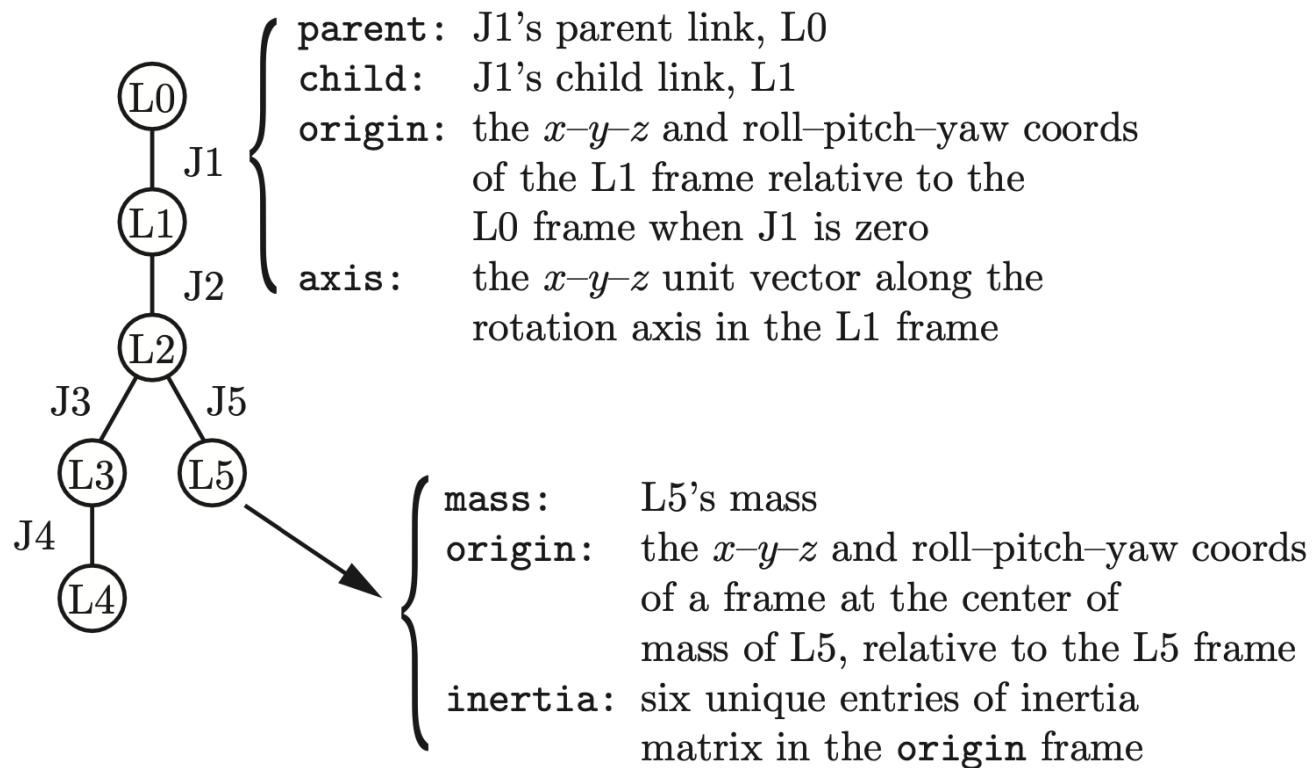
- Check [interactive demo](#)



Content

- Rigid-body Modeling
 - Physical law of rigid bodies
 - Formulation of articulated objects
- Dynamic Modeling of Rigid Bodies
 - Video generation of rigid-body motions
 - Generation of 3D articulated objects

Articulated Objects are a Tree of Rigid Bodies



We need to generate:

1. The connectivity graph of parts
2. The geometry and appearance of parts
3. The kinematics of connection (positions, directions and types of joints)

An Easier Problem: Reconstruct 3D Articulated Objects from an Input Image



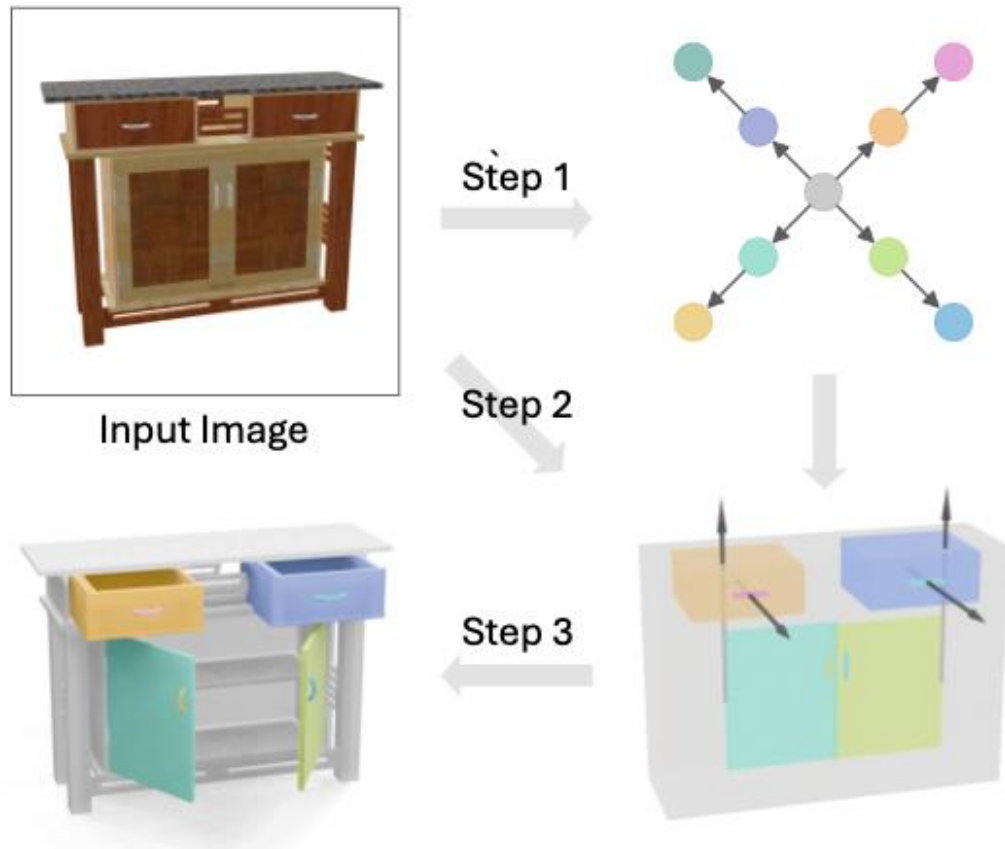
reconstruct



We need to reconstruct:

1. The connectivity graph of parts
2. The geometry and appearance of parts
3. The kinematics of connection (positions, directions and types of joints)

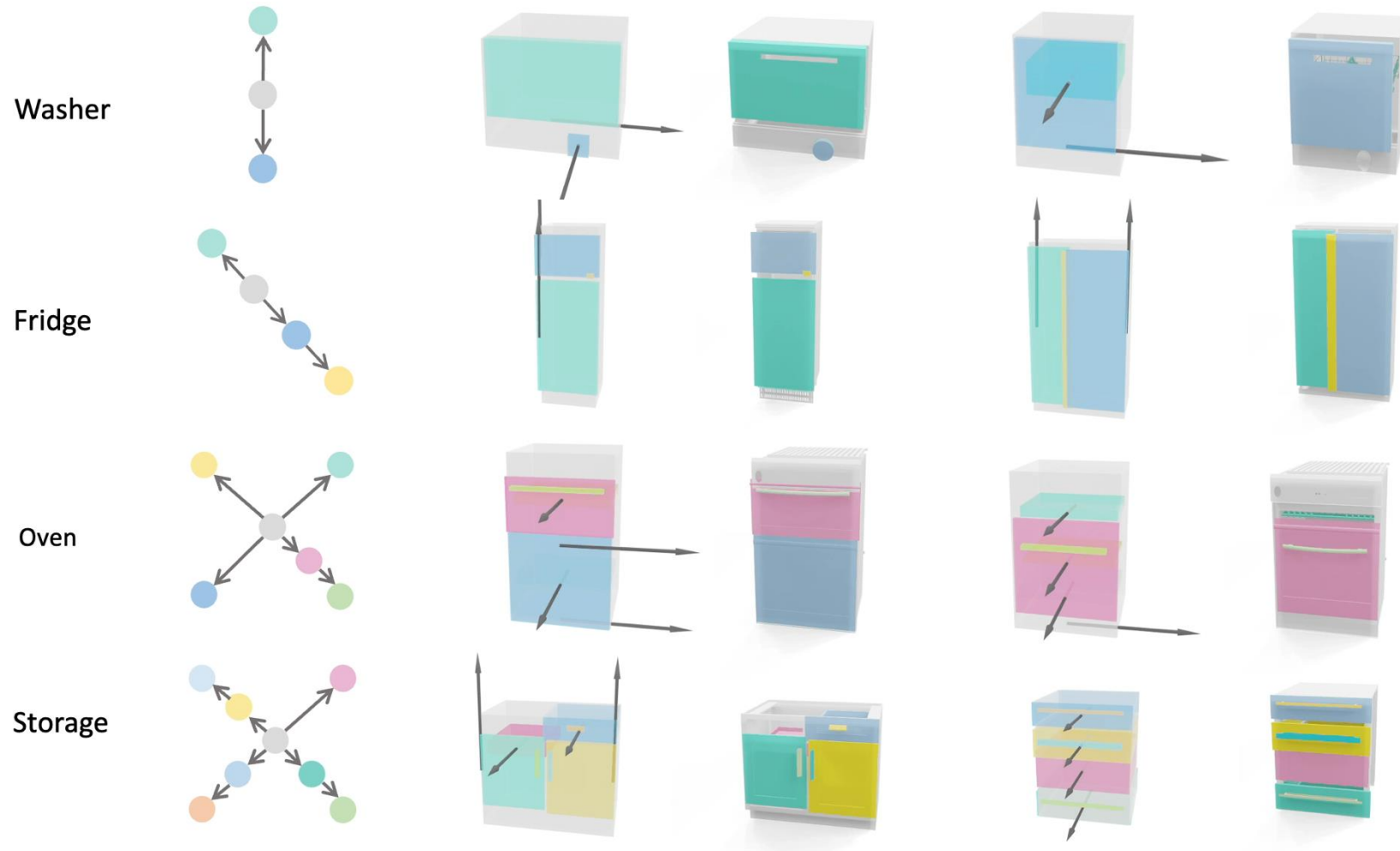
Idea: Breakdown Reconstruction Pipeline



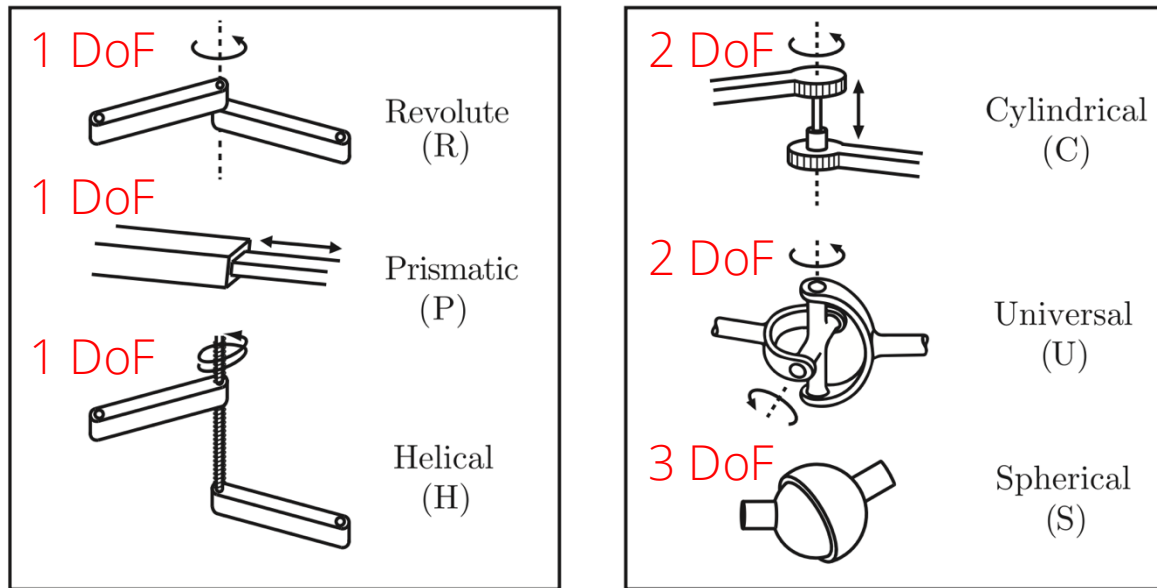
Reconstruction pipeline:

1. Infer the connectivity graph of parts from input image
2. Infer the kinematics of connection (positions, directions and types of joints)
3. Reconstruct the geometry and appearance of parts

How to Represent the Connective Graph?



How to Represent the Kinematics of Connections?



- For simplification, current work only consider a limited set of connections:
- Articulation type: fixed, revolute, prismatic, continuous, and screw joints
 - Kinematic parameters:
 - The location and rotation axis
 - The motion range

How to Represent the Geometry and Appearance of Parts?



For each part, we need to denote:

- Size of each part with bounding box
- Geometry and appearance of each part with structured latents or object code

SINGAPO: Single Image Controlled Generation of Articulated Parts in Objects

Jiayi Liu¹, Denys Iliash¹, Angel X. Chang^{1,2}, Manolis Savva¹, Ali Mahdavi-Amiri¹

¹Simon Fraser University, ²Canada-CIFAR AI Chair, Amii

 ICLR 2025

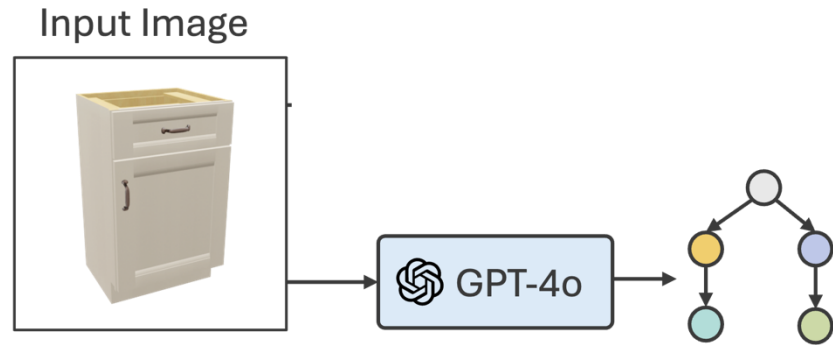
 arXiv

 Code

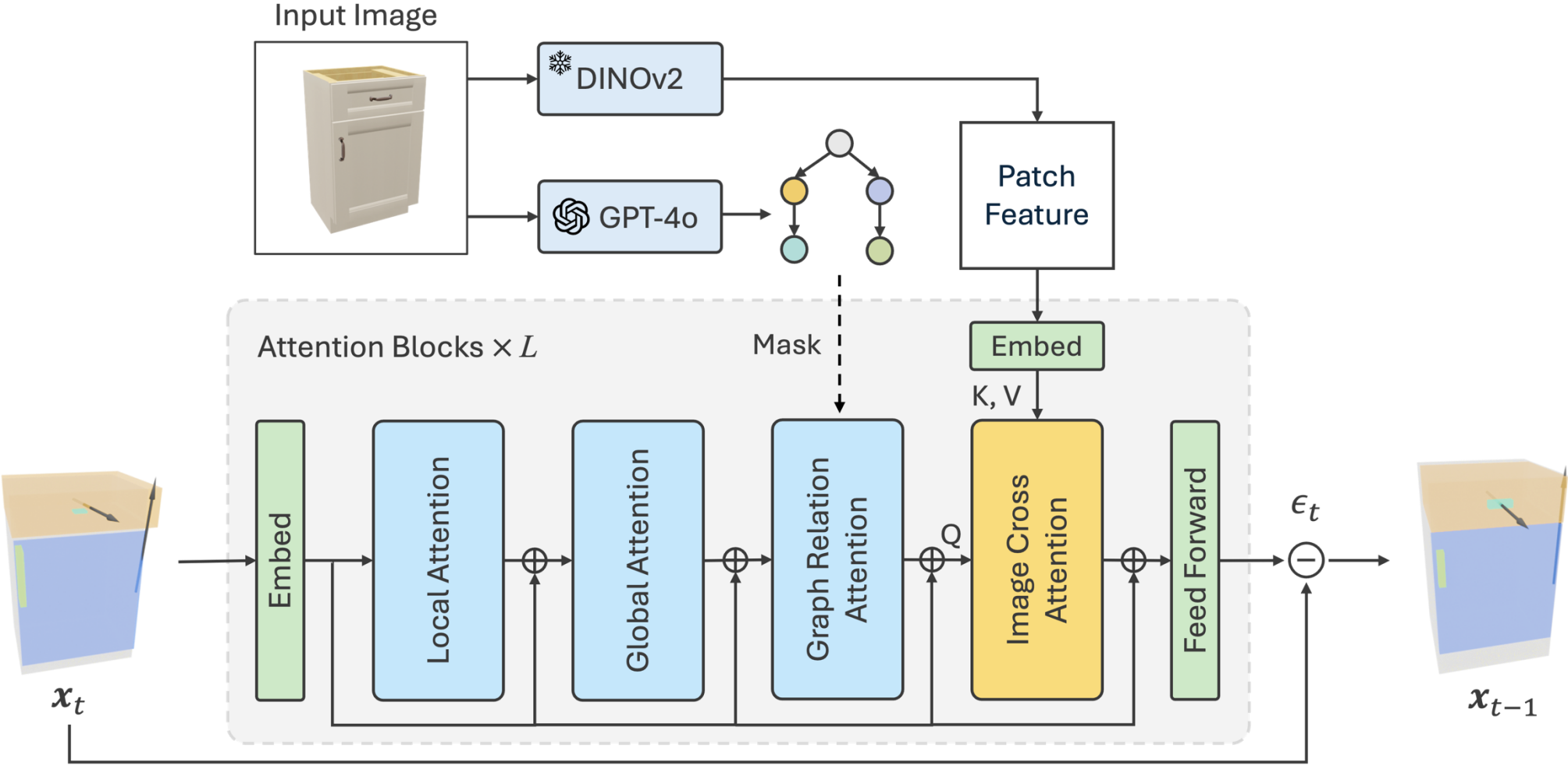
Reconstruction pipeline:

1. Infer the connectivity graph of parts from input image **with VLMs**
2. Infer the kinematics of connection (positions, directions and types of joints) **by learning a transformer**
3. Reconstruct the geometry and appearance of parts **by retrieving 3D assets from a dataset**

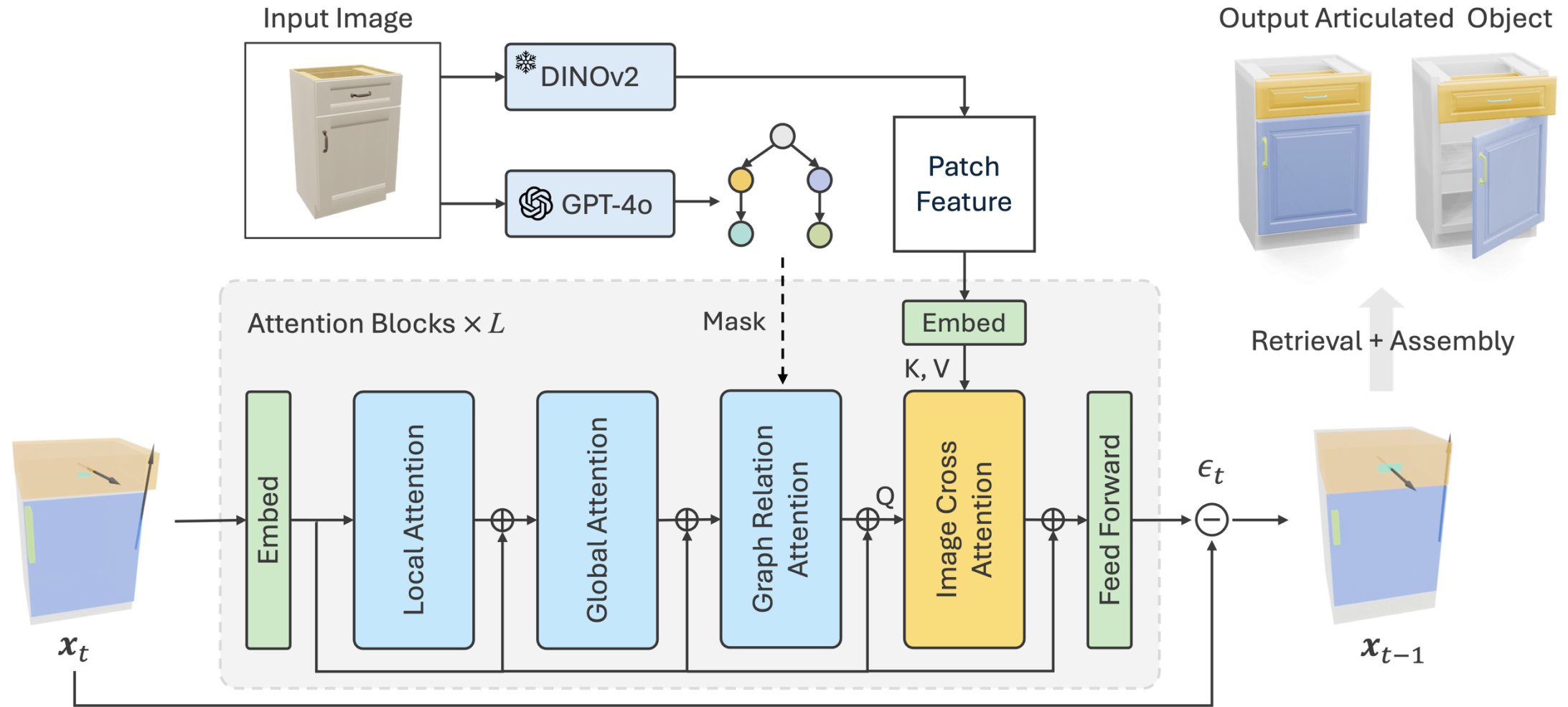
Step 1: Infer the Connectivity Graph of Parts from Input Image with VLMs



Step 2: Infer the Kinematics of Connection by Learning a Transformer



Step 3: Retrieve and Assemble 3D Assets from a Dataset





ArtFormer: Controllable Generation of Diverse 3D Articulated Objects

Jiayi Su^{1,*,\dagger} Youhe Feng^{2,*} Zheng Li⁴ Jinhua Song¹ Yangfan He^{5,6} Botao Ren³ Botian Xu^{3,\ddagger}

¹Xiamen University Malaysia ²Renmin University of China ³Tsinghua University

⁴Southern University of Science and Technology ⁵University of Minnesota-Twin Cities

⁶Henan RunTai Digital Technology Group Co., Ltd.

* Equal Contribution. ^{\dagger}CST2209162@xmu.edu.my ^{\ddagger}btx0424@outlook.com

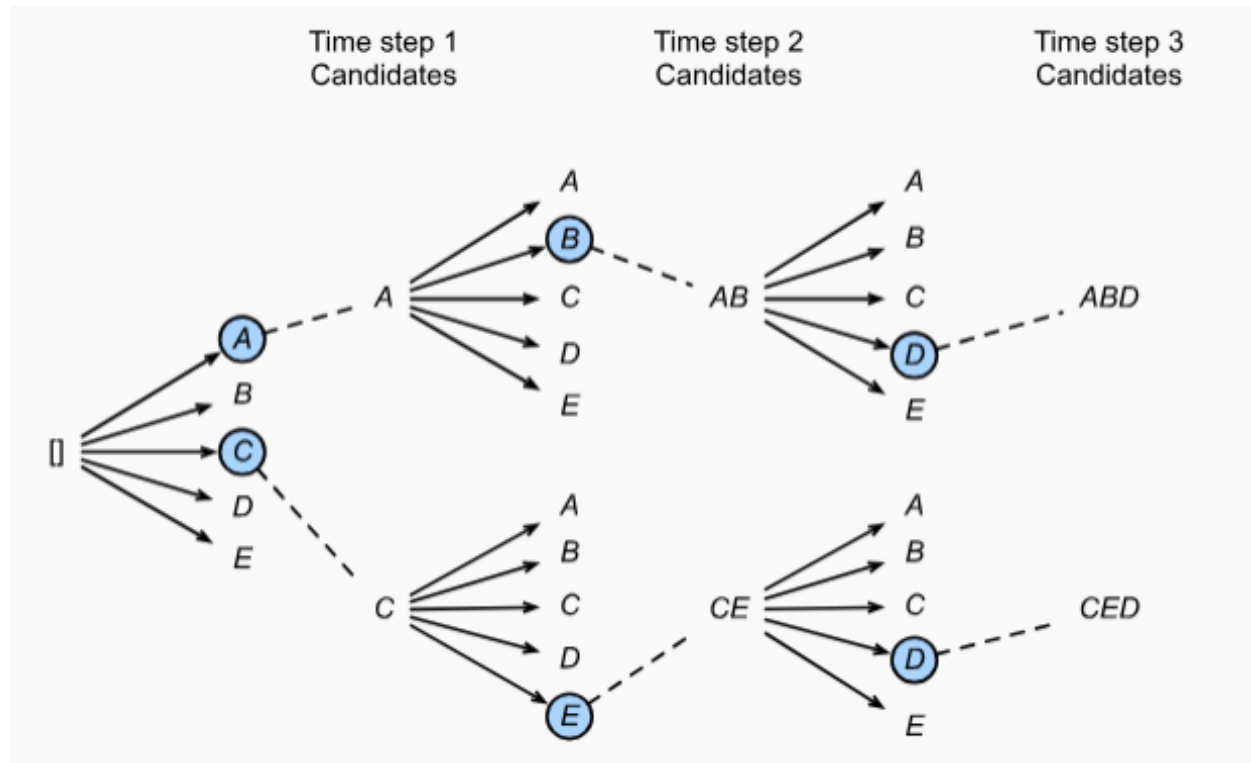
github.com/ShuYuMo2003/ArtFormer

Reconstruction pipeline:

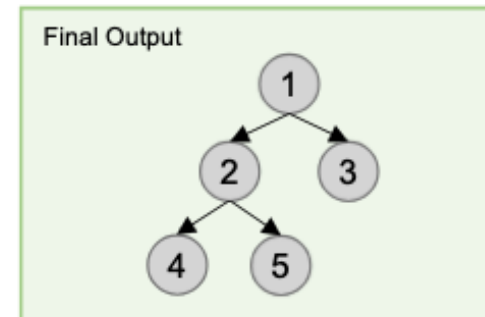
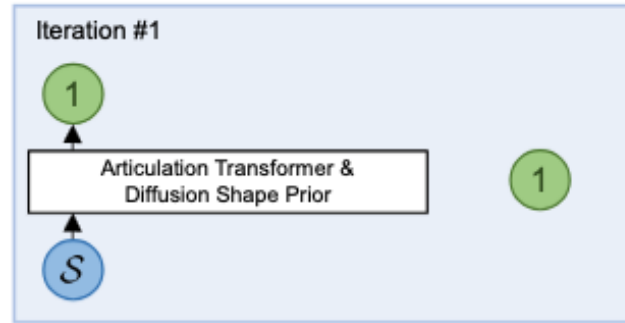
1. Infer the connectivity graph of parts from input image **by learning a transformer**
2. Infer the kinematics of connection (positions, directions and types of joints) **by learning a transformer**
3. Reconstruct the geometry and appearance of parts **by generating structured latents**

Idea: Auto-Regressive Models Naturally Generate Tree of Data

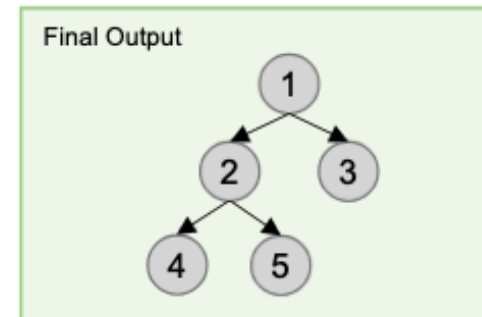
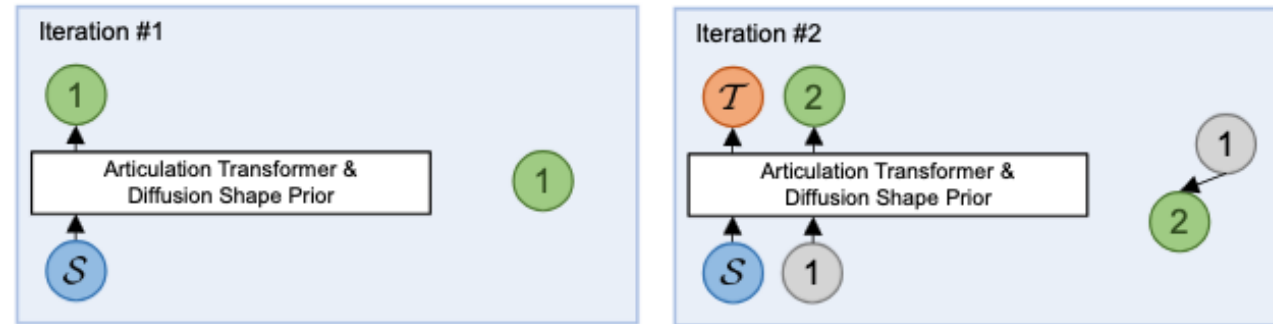
Example: Beam search with Large Language Models



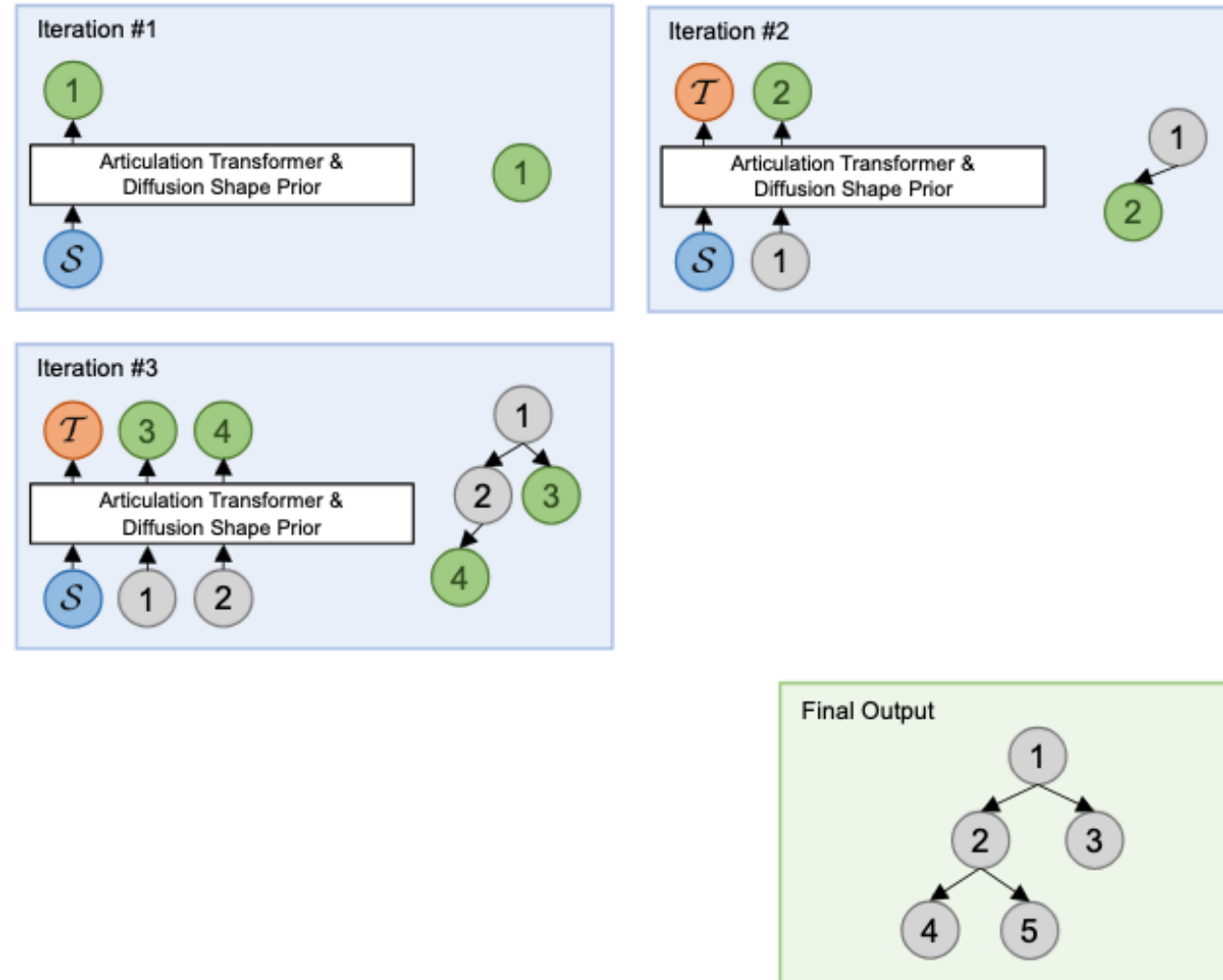
Idea: Sample Tree of Data with Transformers



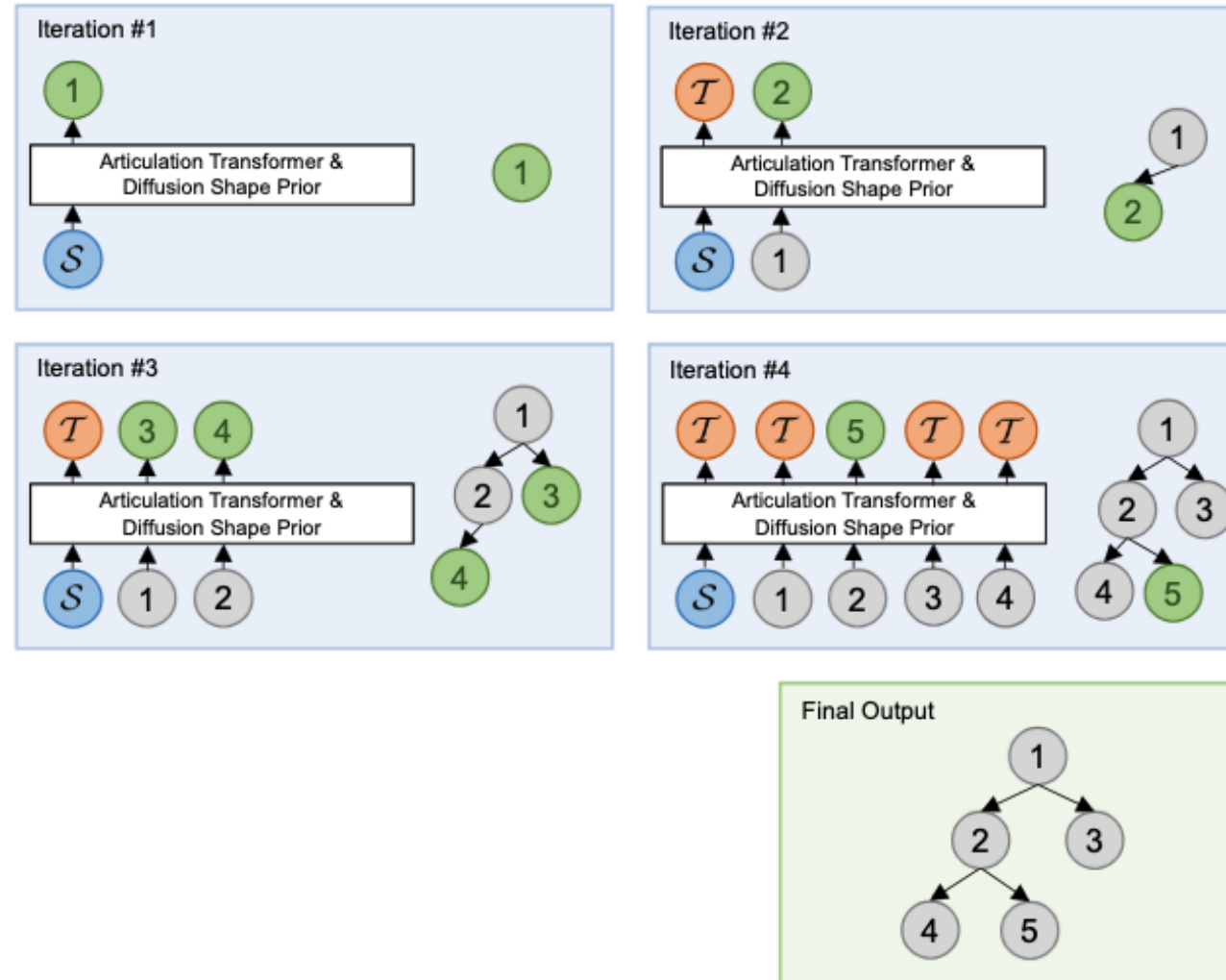
Idea: Sample Tree of Data with Transformers



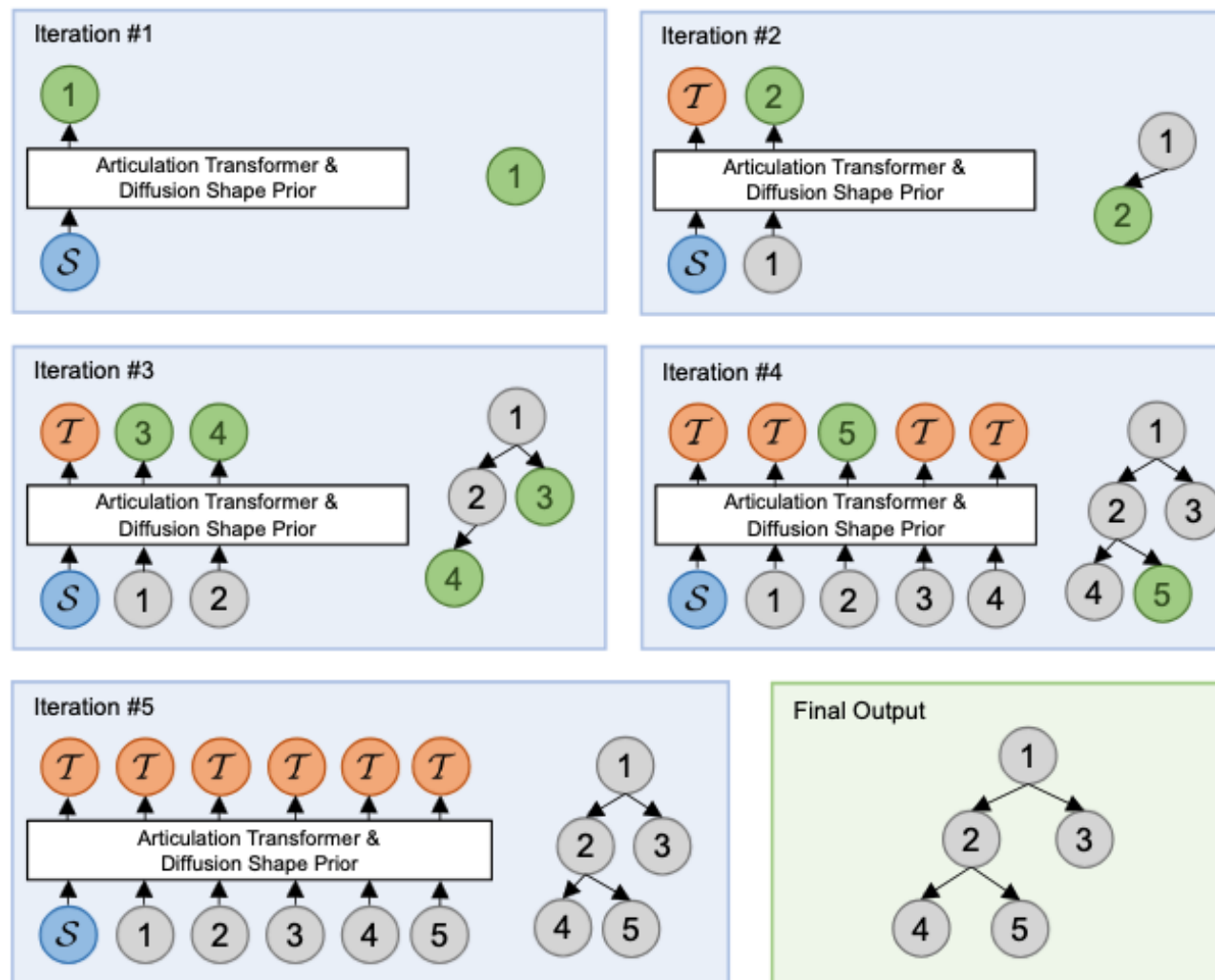
Idea: Sample Tree of Data with Transformers



Idea: Sample Tree of Data with Transformers

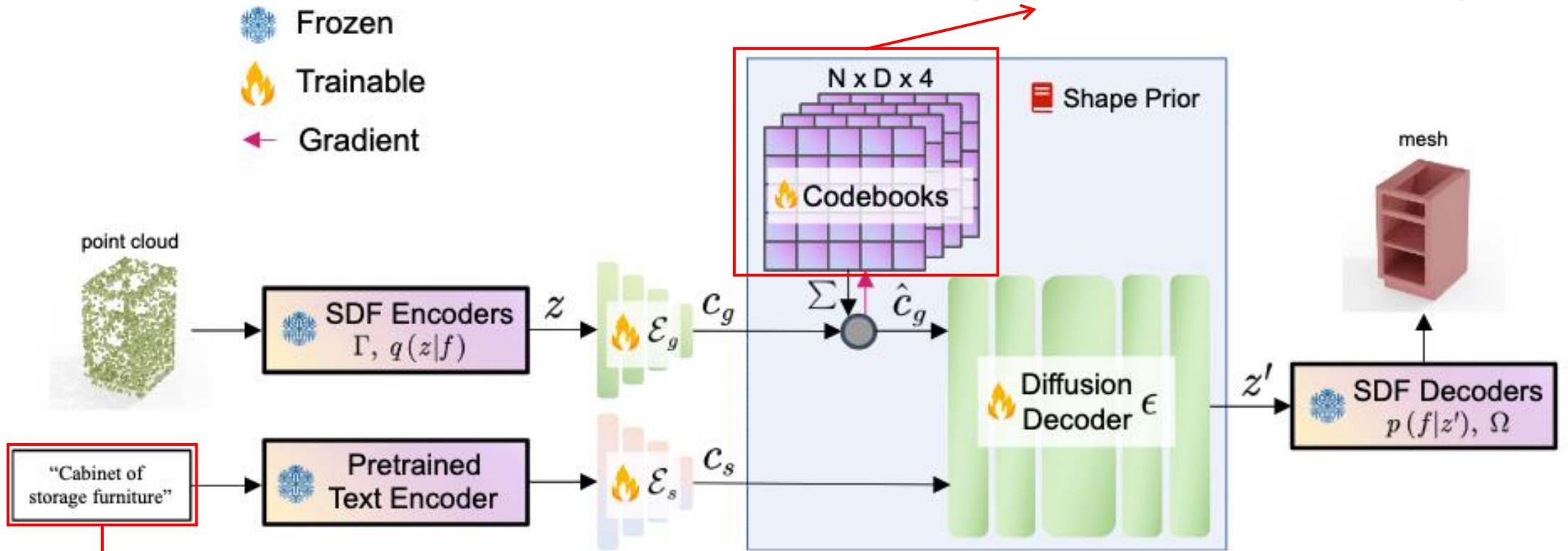


Idea: Sample Tree of Data with Transformers



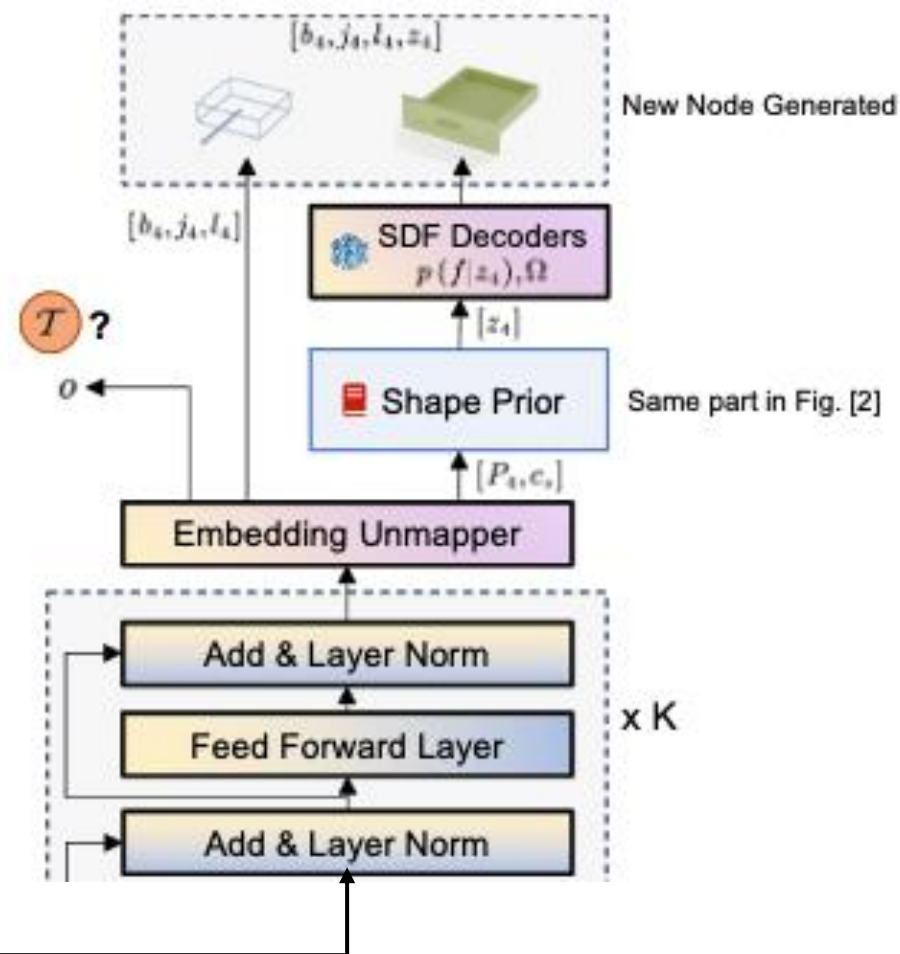
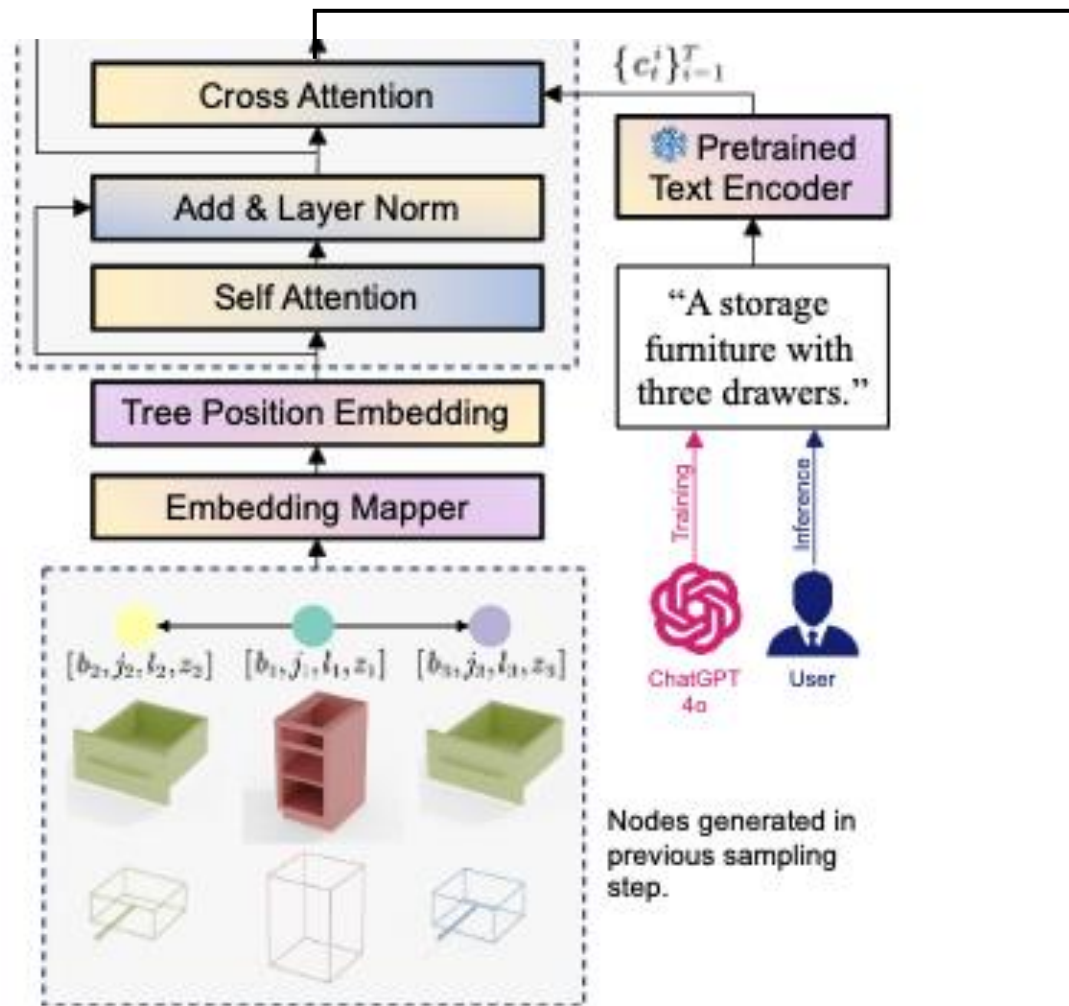
Learn Structured Latents with VAE

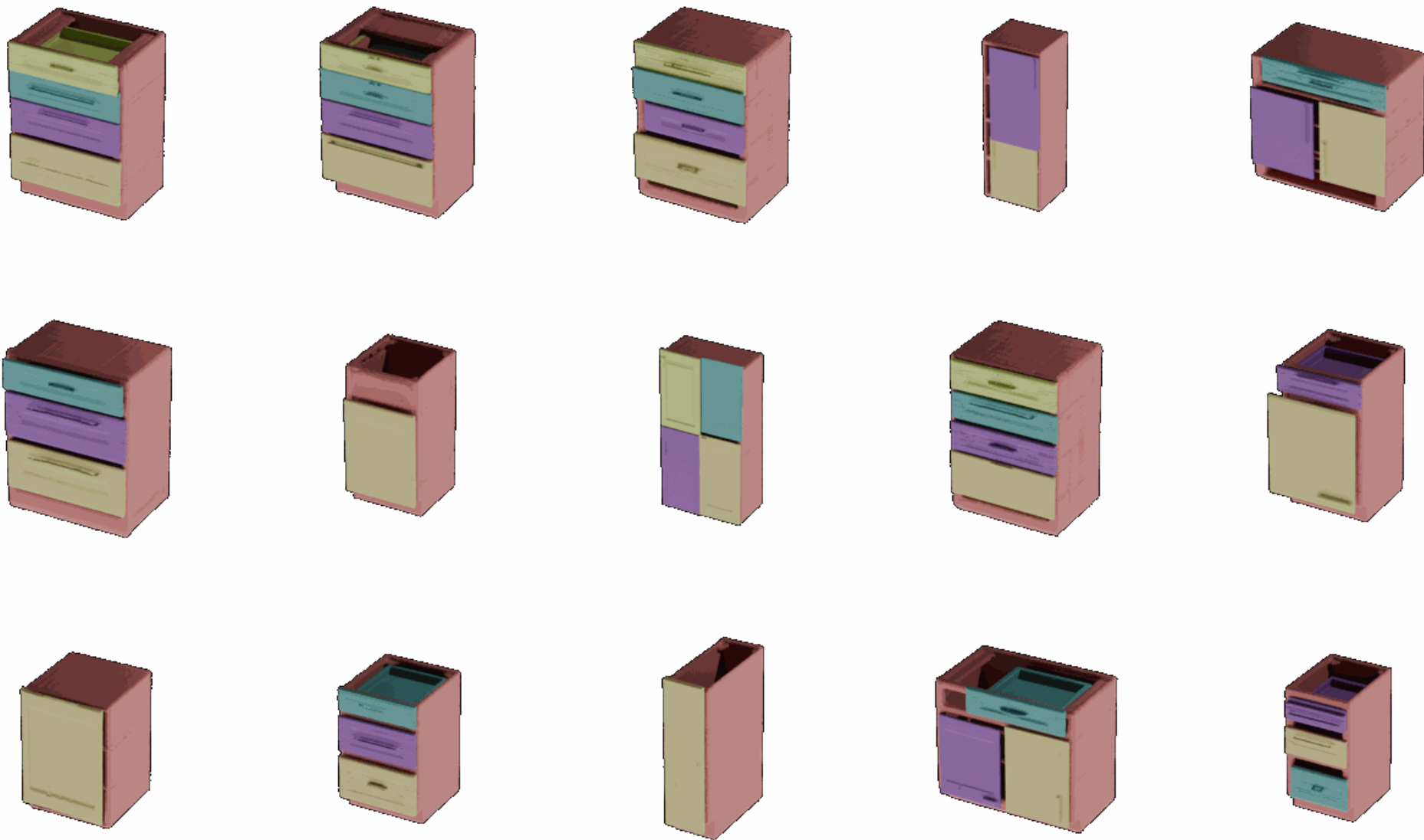
Condition on codebooks of structured latents extracted from training dataset to enhance diversity



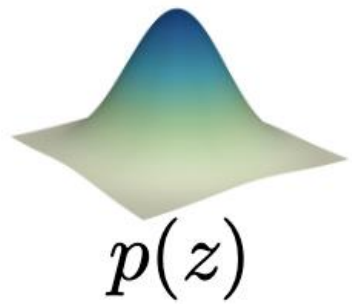
Condition on name of the object part

Put Everything Together





A Harder Problem: Generate 3D Articulated Objects from Scratch



generate



We need to generate:

1. The connectivity graph of parts
2. The geometry and appearance of parts
3. The kinematics of connection (positions, directions and types of joints)

Resources

- You can find additional information from
 - [Physics Simulation in Visual Computing \(with Javascript tutorial\)](#)
 - [Physics-Based Animation Site](#) by Cristopher's Batty

What We Will Cover the Next Week

- Physical modeling
 - Time Integration
 - Mass-Spring System
 - Position-Based Dynamics
 - The Finite Element Method