

Survey on Modeling of Articulated Objects

Jiayi Liu Manolis Savva Ali Mahdavi-Amiri

Simon Fraser University

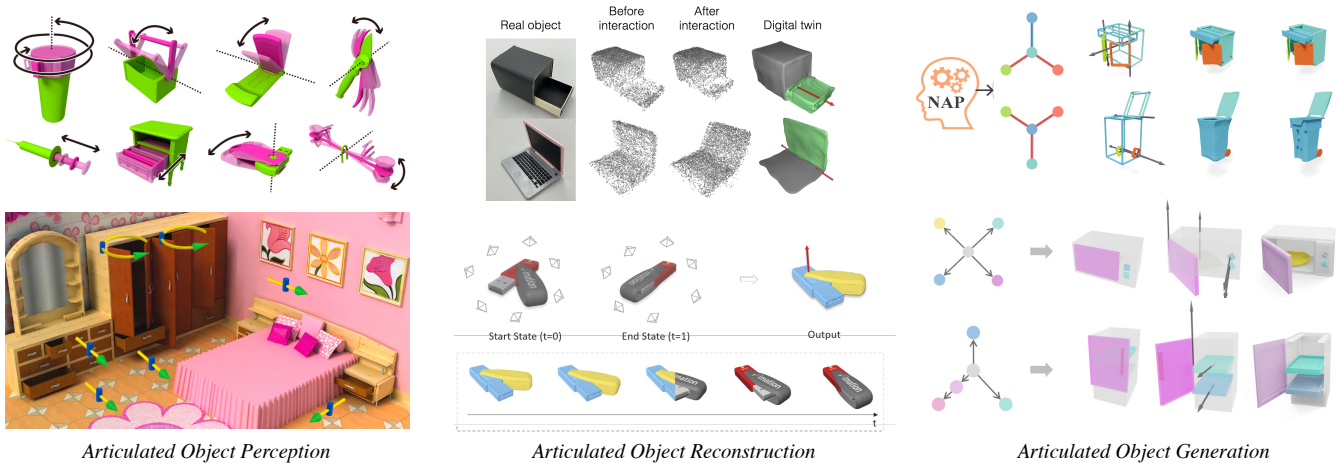


Figure 1: Example of works focusing on different tasks of 3D modeling of articulated objects. From left to right: 1) mobility perception from a single snapshot of an object [HLV*17] or a scene [SHL*14]; 2) articulated object reconstruction from point clouds [JHZ22] and multi-view images [LMS23]; 3) articulated object generation unconditionally [LDS*23] and constrained by a graph [LTMS23]. Figures reproduced from original papers [SHL*14; HLV*17; JHZ22; LMS23; LDS*23; LTMS23]

Abstract

3D modeling of articulated objects is a research problem within computer vision, graphics, and robotics. Its objective is to understand the shape and motion of the articulated components, represent the geometry and mobility of object parts, and create realistic models that reflect articulated objects in the real world. This survey provides a comprehensive overview of the current state-of-the-art in 3D modeling of articulated objects, with a specific focus on the task of articulated part perception and articulated object creation (reconstruction and generation). We systematically review and discuss the relevant literature from two perspectives: geometry processing and articulation modeling. Through this survey, we highlight the substantial progress made in these areas, outline the ongoing challenges, and identify gaps for future research. Our survey aims to serve as a foundational reference for researchers and practitioners in computer vision and graphics, offering insights into the complexities of articulated object modeling.

1. Introduction

Articulated objects, composed of multiple rigid parts connected by joints, are ubiquitous in our daily life, encompassing a wide range of objects, such as human bodies, mechanical assemblies, furniture pieces, etc. Modeling these objects is a research field intersecting computer vision, graphics, and robotics. It aims to comprehend and represent the shape and mobility of articulated components. This endeavor extends to creating 3D models that real-

istically reflect real-world articulated objects. Contributing to the dynamics and functionality in our physical world, articulated object modeling is essential for a wide range of applications, such as animation [YVN*22; CJS*23; QWM*23; LZWL23] and simulation [WCK19; YZW*22; XQM*20], robotic interaction and manipulation [HNOS15; GES21; MGM*21; QF23], embodied AI [KMH*17; PRB*18; GGS*19; SKM*19; PUS*23], etc.

The complexity of articulated object modeling stems from the

fact that articulated objects are not just represented by the geometry of their parts but also by their kinematic structure. Over the past decade, much research has focused on addressing these complexities in two main directions: articulated part perception, and articulated object reconstruction and generation. Articulated part perception focuses on the detection and segmentation of the articulated components of an object while estimating the articulation model that describes the kinematic structure of the object. Articulated object creation involves the reconstruction and synthesis of these parts into a compositional and hierarchical 3D model that reflects the real-world object that can be interacted with.

In this survey, we begin by laying the groundwork with an overview of the field, defining the scope and focus of our discussion in Section 2. Next, in Section 3, we compile the key datasets for articulated objects that have been collected and utilized within the research community, highlighting their critical role in advancing the 3D modeling of articulated objects. Our survey then delves into an in-depth analysis of the techniques and methodologies developed for 3D modeling of articulated objects. This analysis is organized around two pivotal axes: geometry processing and articulation modeling, which are the foundational components to tackle the complexities in this field. Geometry processing covers the representation and approaches employed to understand the shape of the articulated parts, whereas articulation modeling involves the articulation model and the methodologies used to estimate the part mobility and kinematic structure of the object. In the end, we conclude our survey by identifying the key challenges and potential research directions that can drive future progress in the field.

This survey provides a comprehensive overview of the current state-of-the-art in 3D modeling of articulated objects. Through this survey, we highlight the substantial progress made in these areas, outline the ongoing challenges, and identify the gaps for future research. Our survey aims to serve as a foundational reference for researchers and practitioners in computer vision and graphics, offering insights into the complexities of articulated object modeling and inspiring new research in this area.

Related surveys. The most related survey to ours is the recent survey by Pejić et al. [PŠDC22], which focuses on the manipulation of articulated objects in the realm of robotics and computer vision. In contrast, our focus is on the geometry and articulation modeling of articulated objects within the field of computer vision and graphics. This aspect represents a distinct yet critical dimension of articulated object research, orthogonal to the topic of manipulation. Additionally, another related survey by Hu et al. [HSvK18] explores the functionality of general objects within the scope of computer graphics literature. The functionality of an object is influenced not just by its mobility but also by its shape, material, and various other attributes that collectively determine its capability to perform a task. Our survey takes a more specialized focus on the articulated object and its mobility particularly, compared to the work by Hu et al. [HSvK18] which concentrates on functionality for general objects.

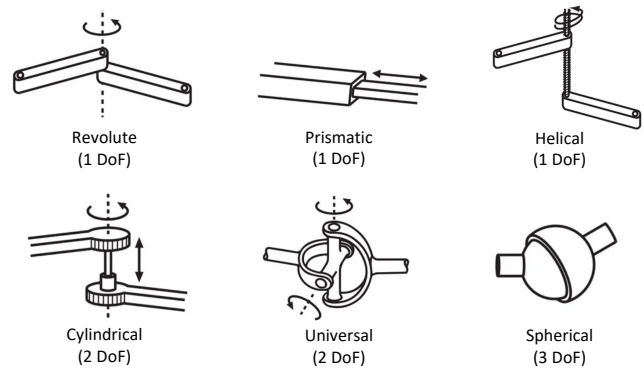


Figure 2: Common joint types. The first row shows the joints with linear motion which is the most common assumption considered in the literature. The second row shows more complex joints with more than 1 DoF. Figure reproduced from original paper [Mue19].

2. Background & Scope

2.1. Definition of Articulated Objects

An articulated object is composed of multiple rigid parts interconnected by joints. These joints serve as the pivotal points that allow for relative motion between the connected parts. The range and nature of the articulation of each part are constrained and defined by the type of joint it possesses. Among man-made objects, the most frequently encountered joint types facilitate linear motion, such as revolute, prismatic, and helical joints. For example, a laptop is an articulated object with a revolute joint connecting the screen to the keyboard. This joint allows the screen to rotate around a single axis relative to the keyboard, with the range of motion being restricted to a specific limit. In contrast, the human body is a more complex example of an articulated object, with a wider variety of joint types, such as ball-and-socket or spherical joints. These joints allow for motion beyond linear translation or rotation, enabling a more diverse range of movements. The illustration of several common joint types is shown in Figure 2.

Types of articulated objects. Articulated objects, with their distinctive characteristics and varied types, play a significant role in the dynamics of our living environments. In general, articulated objects can be classified into two types: *organic objects* and *human-made objects*.

- **Organic objects** refers to naturally occurring systems that are capable of initiating movement through articulated structures. The bodies of humans and other animals are common examples of organic articulated objects. They are crucial for various functions such as locomotion, manipulation, and interaction with the environment. A skeleton with bones and joints is usually used to represent the kinematic structure of these objects. The motion allowed for these objects is usually structurally complex with a relatively large number of joints varied in the topology of connections and degrees of freedom. However, the kinematic structure is relatively fixed and can be predefined for a specific species. How to extract effective skeletons, reconstruct 3D models, and fit dynamic motions for animals across various species

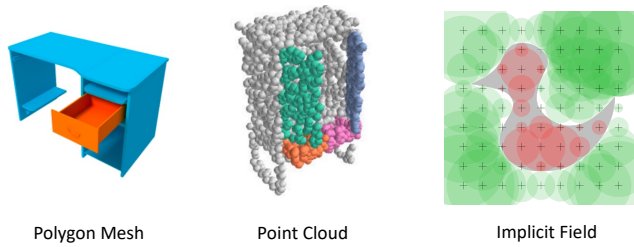


Figure 3: Common geometric representation for articulated object modeling. Figures reproduced from original papers [HLV*17; YHY*19].

is an active research topic [WLJ*23; AM23; LWP*23; LLL*24; YRH*24]. Leveraging the skeleton template of the human body, and coordination between the joints to achieve the desired animation with high fidelity is another line of research [SYZR21; BKY*22; SBR22; WSGT22; LZWL23; QWM*23]. In this survey, we leave the discussion of organic articulated objects out of the scope and only focus on human-made articulated objects.

- **Human-made objects** are engineered systems designed to mimic the natural articulation of organic objects or to serve specific functionalities that require movement. These objects are complex assemblies of rigid parts connected by joints, and their motion is often a direct result of interaction with the environment or driven by an active system. Common examples of human-made articulated objects include gadgets, furniture, vehicles, and other mechanical systems. These objects are ubiquitous in our daily lives and play an important role in our interaction with the physical world. Modeling human-made articulated objects presents several unique challenges in terms of physical and geometric modeling, articulation modeling, and environmental interaction. The modeling of these objects is the focus of this survey and will be discussed in detail in the following sections. We will only mention articulated objects for short in the rest of the survey to refer to human-made articulated objects.

2.2. Representation of Articulated Objects

Articulated object modeling is a multifaceted domain within computer vision, graphics, and robotics. It aims to understand the shape and interactability of the articulated components, represent the geometry and mobility of these objects, and create realistic models that reflect articulated objects in the real world. The inherent complexity of this task arises from the fact that the articulation of the object is not only determined by the geometry of the parts but also by the kinematic structure of the object. As a result, effective modeling of articulated objects requires the capturing of a dual representation of the geometry and articulation of the object. These two facets are deeply interconnected, each influencing and informing the other. The geometric decomposition of the parts lays the foundational groundwork for analyzing part mobility, while the process of articulation modeling also guides the shape understanding of the mobility parts. This symbiotic relationship between geometric and articulation modeling underscores the importance of a cohesive approach in articulated object modeling, where understanding and

representing both geometry and articulation are essential for capturing the full essence and functionality of complex entities. Over the last decade, there has been a notable surge in research within this field, addressing the inherent challenges predominantly from two perspectives: *geometry processing* and *articulation modeling*. We will further discuss these two axes in the following sections.

Geometric representation. Unlike the geometric modeling of static objects, which typically concentrates solely on the object’s outer surface, the geometric representation of articulated objects demands a more nuanced approach. This involves capturing the spatial arrangement of the parts, modeling the part-level surfaces, and constructing the interior structure of the object. It is essential to accurately describe the shape, size, relative position, and orientation of each individual part that is structurally organized within the object. This comprehensive description is vital as it provides the foundational details needed to represent a complete object that can be interacted with. The data format chosen for representing the geometry of the objects varies, depending on data capture methods and the intended application. The most common access to the part geometry is through the 3D point cloud, which is a set of points sampled from the object surface that can be obtained from 3D scanning or depth projection. While the point cloud provides a raw, unstructured representation of the object’s geometry, it often undergoes further processing into more organized formats for detailed analysis. Another commonly employed representation is the 3D mesh, which consists of a collection of polygons that approximate the surface of each part. Due to its structured nature and ability to effectively model surface details, the mesh format is extensively utilized in applications involving simulation, rendering, and animation. In Section 4, we will discuss different choices of geometric representation in various related tasks in detail. Here we specify the geometric representations that are considered in the literature and visually illustrated in Figure 3.

- **Mesh:** is a collection of polygons that approximate the surface.
- **Point cloud (PC):** a set of 3D points sampled on object surface that can be obtained from 3D scanning or depth projection.
- **Implicit field:** is an implicit function that assigns a value to each point in the space to construct a field describing the object shape. Common forms of implicit fields include the signed distance field (SDF), occupancy field, density field, etc. The function parameterized by a neural network, known as a neural field, is a popular intermediate representation as it is inherently differentiable and easy to integrate with learning-based methods.
- **2D images:** some work approaches the 3D modeling problem from the 2D image domain, such as single-view RGB-D images and posed RGB images captured from multi-views or stereo cameras.

Articulation representation. The representation for articulation includes the description of the mobility of each part and the kinematic structure of the object. The kinematic structure describes how the object’s parts are connected and how they can move relative to each other. Typically, this kinematic structure is represented by a tree or graph, which is an effective way to model hierarchical relationships in articulated systems. The part mobility can be described by the parameters of each connected joint or the deformation field of the object. In Section 5, we will discuss ways of articulation modeling with different assumptions and representation in various related

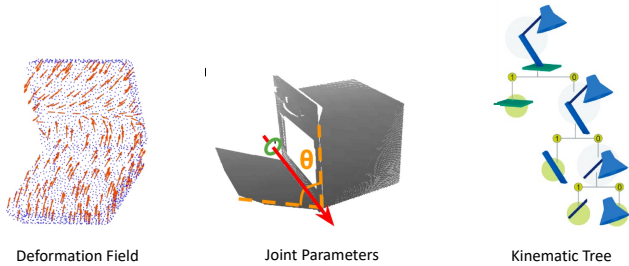


Figure 4: Common articulated motion representation. Figures reproduced from original papers [YHL*18; JLCN21; SHL*14].

tasks in detail. Here we specify the articulated motion representation as follows and visually illustrated in Figure 4.

- **Deformation field:** the deformation field is a general form of motion representation that describes the displacement of each point as a 3D vector. It can also be used to represent beyond the rigid motion of the shape in a continuous space, such as free-form deformation.
- **Joint parameters:** the parameters of each joint, including the *joint type*, location and orientation of the *joint axis*, the rotational angle or translational distance as *joint state*, and the left/right bound of joint as joint limit or *motion range*.
- **Kinematic tree:** the kinematic structure of the object represented by a tree or graph, where the nodes represent the articulated parts and the edges represent the joint connections.

3. Datasets

Recent advancements in deep learning have led to the creation of large-scale 3D datasets for various computer vision tasks, particularly in shape understanding. However, datasets for articulated objects are relatively scarce compared to those for static 3D objects. Gathering data for articulated objects is more laborious because it requires detailed modeling of each part’s geometry and careful annotation of articulation parameters. The challenges in creating datasets for articulated objects are multifaceted:

- The creation of synthetic datasets involves manually designing the geometry of each component. This labor-intensive process restricts the number of categories and instances in these datasets, often leading to a lack of diversity and complexity in the data samples. As a result, models may be less realistic or overly simplistic.
- Real-world datasets are typically gathered using sensors like RGB-D cameras, followed by post-processing to reconstruct shapes and annotate part attributes. One of the significant challenges in modeling articulated objects, as opposed to static surface models, is the difficulty in capturing the interior structure of each part from surface data alone. The data obtained from surface sensors is often incomplete and noisy, with the interior structure frequently occluded by the outer surface. This limitation impacts the geometric quality of the data and further reduces the diversity and complexity of data samples.

With the recent surge in interest in articulated objects, there has

been a growing effort to create datasets for this domain. We collect the existing datasets for articulated objects in both object-level and scene-level in Table 1.

3.1. Synthetic Data

The recent advancements in 3D modeling technology have led to a significant increase in the availability of 3D models with part-level structures accessible online. This progress has opened up new possibilities for detailed analysis and application in various fields, enhancing the resources for researchers and practitioners working with complex 3D structures. Articulated objects are one such class of complex 3D structures that have benefited from this progress. The availability of 3D datasets with part-level structures has enabled the creation of synthetic datasets for articulated objects.

To facilitate the prediction of part mobility via a learning-based approach, Hu et al. [HLV*17] collected a synthetic dataset for 3D articulated objects by sourcing shapes from ShapeNet [CFG*15] and SketchUp [Inc17b]. In this process, they manually segmented the shapes into parts and provided detailed annotations for the articulation parameters. Following this foundational work, RPM-Net [YHY*19] extends the scope of the dataset to include a greater variety of objects incorporating more complex motion, e.g. the opening of the umbrella cover. For each shape in the dataset, each pair of possible parts is labeled as a *reference part* and a *moving part*. And each pair is taken as a *mobility unit*, which is associated with annotated motion parameters. These parameters form a quadruple, including the transformation type (which can be translation, rotation, or the combination of the two), the position and direction of the transformation axis, and the range of the motion.

Shape2Motion [WZS*19] is another synthetic dataset in a larger scale that was introduced concurrently. The shapes are sourced from ShapeNet [CFG*15] and 3D Warehouse [Inc17a]. Each shape includes part segmentations and articulation parameters that are annotated in a similar way as RPM-Net [YHY*19] and Hu et al. [HLV*17]. However, the process is further refined with the use of a developed annotation tool. This tool enhances the efficiency of the annotation process by allowing the annotators to visually verify the correctness by animating the object with the annotated parameters. This allows for boosting the scale of the dataset with a much larger number of objects and movable parts in a wide range of categories.

The primary application of the datasets mentioned earlier involves sampling point clouds from the mesh surfaces to analyze part mobility. However, a notable limitation of these datasets is the absence of texture information. This lack of texture detail restricts their utility in applications such as simulation, the creation of digital twins, augmented reality, and other modeling tasks that require a more holistic visual representation. To address this gap, the PartNet-Mobility dataset, with a simulation environment SAPIEN [XQM*20], was introduced. This dataset is a subset and an extension from a part-level 3D shape dataset PartNet [MZC*19], enriched with articulation annotations organized in URDF files. The inclusion of diverse appearances and motions in PartNet-Mobility significantly enhances its value, inspiring increasingly more research in the field of articulated object modeling. Building upon the PartNet-Mobility dataset, OPDSynth [JMCS22] is introduced for

	Source	Level	Representation	Textured	# objects	# movable parts	# object categories
Hu et al. [HLV*17]	synthetic	object	mesh	✗	368	368	-
RPM-Net [YHY*19]	synthetic	object	mesh	✗	969	1,420	-
Shape2Motion [WZS*19]	synthetic	object	mesh	✗	2,440	6,762	45
RBO [MEB19]	synthetic, real	object	mesh, RGB-D	✓	14	21	14
PartNet-Mobility [XQM*20]	synthetic	object	mesh	✓	2,346	14,068	46
ReArt-48 [LXX*22]	real	object	mesh	✓	48	?	5
AKB-48 [LXF*22]	real	object	mesh	✓	2,037	?	48
OPDSynth [JMSC22]	synthetic	object	mesh, RGB	✓	683	1,343	11
OPDReal [JMSC22]	real	object	mesh, RGB-D	✓	284	875	8
MultiScan [MZJ*22]	real	scene	mesh, RGB-D	✓	10,957	5,129	20
GAPartNet [GXZ*23]	synthetic, real	object	mesh	✓	8,489	1,166	27

Table 1: Datasets for articulated objects with statistics on the source of collection, the level at which these datasets operate, data representation, whether the 3D model is texture, the number of objects, movable parts, and object categories. The '-' symbol indicates that the information is not available, and the '?' symbol indicates that the information is not reported in the original paper.

the task of *openable part detection* (OPD) by blending the rendering of the synthetic data with real-world RGB images. This dataset enables the mobility part detection from RGB images, aligning more closely with practical real-world applications.

3.2. Real-world Data

A primary challenge in working with synthetic data is the often unrealistic assumption that it is perfect, overly simplistic, complete, and free from noise. To fill this synthetic-real gap, more recent efforts have started to focus on collecting data from real-world sensors, such as RGB-D cameras.

RBO [MEB19] represents the pioneering dataset in this regard, being the first to collect data from RGB-D scans. It not only provides RGB-D recordings of human interaction with objects under varying experimental conditions, but also creates corresponding synthetic mesh models for each object. ReArt-48 [LXX*22] introduced a slightly larger dataset that reconstructs meshes from RGB-D scans and provides the articulation annotations across five object categories. Taking this further, AKB-48 [LXF*22] emerges as the first large-scale dataset for articulated objects based on real scans, where each object is described in a knowledge graph. To construct this dataset, a fast articulation knowledge modeling pipeline is presented, significantly reducing the cost and effort required for object modeling in the real world. This innovation enables the creation of 3D models on a scale comparable to PartNet-Mobility. Additionally, AKB-48 also annotates the physical properties of the objects, such as mass, to enhance the dataset’s applicability. This additional annotation is vital for bridging the generalization gap between simulation and real-world applications.

In parallel, the MultiScan [MZJ*22] dataset marks a groundbreaking advancement as the first large-scale scene-level dataset that documents multiple states of articulated objects in indoor settings. To compile this dataset, a scalable 3D environment acquisition pipeline is designed for processing raw RGB-D scans to produce 3D surface mesh reconstruction with texture and articulation annotations for each articulated object in the scene. MultiScan is an invaluable resource for advancing research and applications that

require an understanding of the dynamics and interactions of articulated objects in real-world environments at the scene level. Leveraging the MultiScan dataset, OPDMulti [SJSC23] is introduced for the OPD task of the multiple-parts version by extracting frames from the RGB-D scans and annotating the openable parts in each frame. It extends the OPD task to the scene level, which is more challenging and practical in real-life scenarios.

More recently, the GAPartNet [GXZ*23] dataset is introduced to capture finer-level part details that have been previously overlooked in existing datasets. This dataset is tailored to enhance the generalization of object perception and manipulation tasks across various object categories by focusing on the concept of *Generalizable and Actionable Parts*. The underlying premise is that functional parts, such as buttons and handles, are fundamental elements whose identification and extraction can significantly improve generalizability within and across object categories. They argue that the functional parts such as buttons and handles are more elementary, and the extraction of these parts can improve the generalizability in an intra-category manner. To achieve a comprehensive and versatile dataset, GAPartNet selectively compiles data from both the PartNet-Mobility [XQM*20] and AKB-48 [LXF*22] datasets to encompass both synthetic and real data. It provides a rich resource for robust algorithms that are applicable across scenarios in the realm of articulated object modeling and manipulation.

Slightly different from the previous datasets, the OPDReal [JMSC22] dataset is introduced for the OPD task by collecting RGB-D scans from real-world scenes. By providing a set of real RGB images capturing various articulation states of the objects, along with reconstructed mesh models, OPDReal presents new opportunities for research in the field of articulated object modeling from 2D visuals. OPDMulti [SJSC23] follows this trend and extends the OPD task to multiple objects in the scene by collecting RGB-D frames, which is more challenging and practical in real-life scenarios.

4. Geometry Processing

In this section, we discuss the recent works on articulated object modeling from the perspective of geometry processing. Based on

	Input		Input Assumptions				Methodology			Output	
	geo rep.	# states	part seg	# parts	partial/noisy	aligned	intermediate rep.	strategy	supervision	geo rep.	part seg
Articulated Part Perception											
Xu et al. [XWY*09]	mesh	1	✗	✗	✗	-	surface patch	handcrafted	-	mesh	✓
Mitra et al. [MY*10]	mesh	1	✓	✓	✗	-	surface patch	handcrafted	-	mesh	✗
Sharf et al. [SHL*14]	mesh	multi	✗	✗	✓	✗	surface patch	handcrafted	-	mesh	✓
Yuan et al. [YLX*16]	PC	multi	✗	✗	✓	✗	3D trajectory	handcrafted	-	PC	✓
Li et al. [LWL*16]	RGB-D	multi	✗	✗	✓	✗	3D trajectory	handcrafted	-	PC	✓
Hu et al. [HLV*17]	mesh	1	✓	✓	✗	-	surface patch	-	-	-	✗
Yi et al. [YHL*18]	PC	2	✗	✗	✓	✗	PC feat.	supervised	labeled PC	PC	✓
Shape2Motion [WZS*19]	PC	1	✗	✗	✗	-	PC feat.	supervised	labeled PC	PC	✓
RPM-Net [YHY*19]	PC	1	✗	✗	✓	-	PC feat.	supervised	labeled PC	PC	✓
Abbatematteo et al. [ATK19]	RGB-D	1	✗	✓	✓	-	image feat.	supervised	labeled PC	PC	✓
Li et al. [LWY*20]	PC	1	✗	✓	✓	-	NAOCS, NPCS	supervised	labeled PC	PC	✓
Shi et al. [SCZ21]	PC	multi	✗	✗	✓	✓	3D trajectory	-	-	PC	✓
ScrewNet [JLCN21]	PC	multi	✗	✓	✓	✓	PC feat.	-	-	-	✗
Abdul-Rashid et al. [AFA*22]	PC	1	✗	✗	✓	-	PC feat.	supervised	labeled PC	PC	✓
Qian et al. [QJR*22]	RGB	multi	✗	✗	✓	✓	image feat.	supervised	labeled PC	3D plane	✓
OPD [JMSC22]	RGB	1	✗	✓	-	-	image feat.	supervised	labeled PC	2D mask	✓
OPDMulti [SJSC23]	RGB	1	✗	✗	-	-	image feat.	supervised	2D mask	2D mask	✓
Liu et al. [LXX*22]	RGB-D	1	✗	✗	✓	-	NPCS	supervised	labeled PC	PC	✓
GAPartNet [GXZ*23]	RGB-D	1	✗	✓	✓	-	NPCS	supervised	labeled PC	PC	✓
Liu et al. [LZH*23]	PC	1	✗	✓	✓	-	PC feat.	self-supervised	PC corr.	PC	✓
Liu et al. [LSH*23]	PC	1	✓*	✗	✗	-	PC feat.	supervised	labeled PC	PC	✓
Banana [DLS*23]	PC	1	✗	✓	✗	-	PC feat.	supervised	labeled PC	PC	✓
Articulated Object Creation											
Pekelny and Gotsman [PG08]	PC	multi	✓*	✓	✓	✓	PC	handcrafted	-	mesh	✓
A-SDF [MQK*21]	SDF	1	✗	✓	✗	✓	neural field	supervised	SDF	mesh	✗
Ditto [JHZ22]	PC	2	✗	✓	✓	✓	neural field	supervised	labeled PC	mesh	✓
CLA-NeRF [TLYS22]	MV RGB	1	✗	✓	-	✓	neural field	supervised	2D mask	mesh	✓
Wei et al. [WCM*22]	MV RGB	1	✗	✓	-	✓	neural field	supervised	MV RGB	mesh	✗
CARTO [HIZ*23]	MV RGB	1	✗	✓	-	-	neural field	supervised	SDF	mesh	✗
PARIS [LMS23]	MV RGB	2	✗	✓	-	✓	neural field	self-supervised	MV RGB	mesh	✓
NAP [LDS*23]	-	-	-	✗	-	✓	neural field	generative	-	mesh	✓
CAGE [LTMS23]	-	-	-	✗	-	✓	bbox	generative	-	mesh	✓

Table 2: Summary of the related works in the aspect of geometry processing. The table provides information about the geometric representation and the number of articulated states observed in the input, the assumptions made on the input data, the methodology used in the geometry processing, and the output data representation. Please refer to section 4.2.1 and fig. 6 for the explanation of the intermediate representations.

the goal of the tasks, we categorize the works into two groups: articulated part perception and articulated object creation. The former group of works focuses on understanding the articulated part structure of the objects from either 3D geometry or 2D observations. The latter group focuses on reconstructing or generating the 3D geometry of the articulated objects with either articulation parameters estimated from input observations or the ability to animate the objects to different articulation states. In Table 2, we summarize these works in terms of the choice of geometric representation, the assumptions that each work made for the input, and the methodology each work adopted for shape analysis. We specify the meaning of each column in Table 2 as follows:

- **geo rep.:** the geometric representation as input, output, or intermediate representation during processing (under *methodology*).
- **# states:** the number of articulation states observed in the input.
- **part seg:** whether the input shape of the object is pre-segmented into articulated parts with labels.
- **aligned:** whether the object is aligned across different states when multiple articulated states are observed in the input; whether the objects are aligned in a canonical space when there are multiple objects in the input; for generative models, objects are assumed to be canonicalized by default.
- **# parts:** whether the number of articulated parts is known.
- **partial/noisy:** whether the input geometry can be partially observed (e.g. point cloud projected from the single view depth) or noisy (e.g. raw scans with sensory noise).
- **intermediate rep.:** the intermediate representation used during the processing of the input shape or to represent the shape.
- **strategy:** the methodology adopted for shape analysis, including handcrafted methods, supervised learning, and self-supervised learning. *Handcrafted* refers to the methods that rely on human-designed descriptors, heuristic rules, or mathematical algorithms to analyze the shape and it is not learning-based. *Supervised* refers to the data-driven methods that require ground truth labeled data for training, such as 2D masks, segmentation labels, displacement for each point, etc. *Self-supervised* refers to the data-driven methods that only rely on the input data for training, such as leveraging the geometric consistency between different states of the object.
- **supervision:** the source of supervision used for the learning-based methods, including distance function (*distance func.*) typically used for metric learning, segmentation labels on the point cloud (*labeled PC*), segmentation mask on images (*2D mask*).

4.1. Geometric Representation for Task Setting

4.1.1. Articulated Part Perception

To ultimately analyze the part mobility of articulated objects, the recognition of the part segmentation from an input object shape is usually necessary as a first step. One exception to this is the early work from Mitra et al. [MY*10] and Hu et al. [HLV*17] that made the assumption that part segmentation is given as input to simplify the problem. ScrewNet [JLCN21] analyzes the part mobility from a point cloud sequence without knowing the part structure, but they directly predict the articulation parameters without explicitly assigning the points to different parts. Apart from these, other works for articulated part perception address the segmentation problem jointly with the part mobility analysis.

Perception from meshes. As the most widely used geometric representation in the industry, polygonal meshes are naturally chosen by the early works [XWY*09; MY*10; SHL*14; HLV*17] as the initial input representation to work with. These meshes can be obtained from CAD models that are carefully designed by the artists with part-level structure and sharp edges. It implies that the input geometry in these works is assumed to be complete and noise-free, which simplifies the problem by observing the object in a geometrically perfect condition in 3D. Mitra et al. [MY*10] and Hu et al. [HLV*17] further simplify the problem of part mobility analysis by assuming that the partition of the object into articulated parts is accessible as input to the system. Sharf et al. [SHL*14] pioneered addressing the segmentation problem for objects and articulated parts from an indoor scene as a mesh. Their segmentation relies on the difference between the object and part pose being observed in multiple states. They did not make any assumptions on the mesh properties, such as self-intersection or non-manifoldness. It loosens the requirement of the input geometry to be perfect, but it still assumes that the object is observed in a complete form with part structure available. It is expensive to obtain such a delicate mesh model from the real world, especially when we try to replicate a real-world object with a complex shape and part structure.

Perception from point clouds. Point clouds of the objects can be obtained from 3D scans or depth sensors, which are more accessible and easier to obtain than the polygonal meshes. This accessibility makes the point clouds a more practical choice for the input representation as it can be directly obtained from the real world. It encourages more works later that consume point clouds as input for articulated part perception. Leveraging the ease of acquisition, a line of work [LWL*16; YLX*16; SCZ21; JLCN21] take point clouds that observe a sequence of articulated motion of the object as input. These point cloud sequences are assumed to be incomplete and noisy, which makes the robustness of the motion fitting an additional critical requirement for the system. Another implied assumption is that the object is aligned in the same coordinate system across different states, which is a prerequisite for the motion fitting process. However, observing a sequence of articulated motion of an object is unnatural in the real-world setting for passively articulated objects as they cannot move on their own. Consequently, the motion-capturing process typically demands extensive human intervention, making it costly and labor-intensive. Instead of observing motion sequences, Yi et al. [YHL*18] propose to take a pair of point clouds of different articulation states as input (illustrated in).

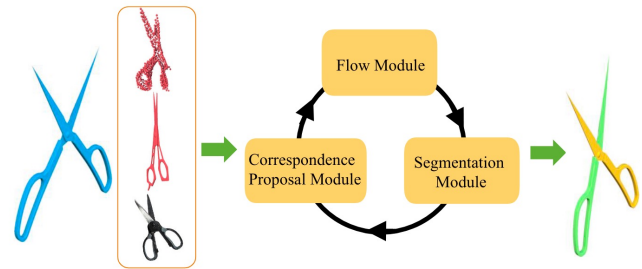


Figure 5: Method proposed by Yi et al. [YHL*18] to analyze the articulated part structure from a pair of point clouds from functionally related objects in arbitrary articulation states. Figure reproduced from original paper[YHL*18]

Their work also relaxes several assumptions made by the previous works that the object can be misaligned in two states and the input only needs to be from objects that are functionally related but not necessarily from the same instance. More recent works [WZS*19; YHY*19; LWY*20; LZH*23; DLS*23] tackle the problem from a single state of the object at inference time to make it more accessible on the input. Abdul-Rashid et al. [AFA*22], Liu et al. [LXX*22], and GPartNet [GXZ*23] further extend the input to colored point clouds, which can leverage more visual information than the raw point clouds. One note for the work from Liu et al. [LSH*23] is that they conduct the articulated part analysis from an over-segmented point cloud, which is a different problem setting from the other works in this category.

Perception from monocular RGB images. Even more accessible than point clouds, monocular RGB(-D) images have also served as input for articulated part perception recently. Qian et al. [QJR*22] proposed to extract 3D planes to represent articulated parts from RGB videos recording the human-object interaction. OPD [JMCS22] first introduces the task of detecting openable parts of objects from a single-view image. OPDMulti [SJSC23] extends the scenario to multiple objects in the scene as a follow-up. The articulation analysis in 3D from a single-view RGB image is a challenging but useful task, especially for the robotics system equipped with RGB cameras.

4.1.2. Articulated Object Creation

Articulated object reconstruction. The goal of articulated object reconstruction is to recover the 3D geometry of the articulated objects with either articulation parameters estimated from input observations or the ability to animate the objects to different articulation states. Since triangle meshes are the most widely used geometric representation for 3D objects, it is naturally chosen as the output representation for the works in this category. Pekelny and Gotsman [PG08] presents the earliest work that reconstructs the surface from a sequence of point clouds, where they assume the part segmentation and skeleton of the rigid pieces are given in the first frame. Once all the other frames are aligned with the first one by tracing the skeleton, the mesh surface of each part is accumulated and recovered using point cloud meshing algorithms [KBH06; ACTD07]. A-SDF [MQK*21] is the first work that reconstructs articulated

objects from signed distance fields (SDFs) of a single snapshot in a way that the objects can be deformed to any other articulation states. Wei et al. [WCM*22] share the same goal with A-SDF but take multi-view RGB images as input for the reconstruction. CARTO [HIZ*23] further extends the input to a single stereo RGB image. One limitation of these three works is that they both reconstruct the articulated object as a whole into a single surface. The lack of part-level structure in the output geometry makes it difficult to actually interact with the object in a meaningful way. CLANeRF [TLYS22], Ditto [JHZ22], and PARIS [LMS23] address this limitation by reconstructing the articulated object in the part level. CLANeRF and PARIS both take multi-view RGB images as input, while Ditto consumes point clouds. Ditto and PARIS share the configuration that the object is observed from two articulated states and the two states are assumed to be pre-aligned. However, research in this category typically assumes a predetermined number of articulated parts, concentrating on modeling single moving parts or relying on the object’s category to infer its structure. Such assumptions highlight a notable gap in the literature: the challenge of reconstructing objects with an arbitrary number of articulated parts in varying structures remains insufficiently addressed.

Articulated object generation. This category of work aims to synthesize the 3D geometry of articulated objects with articulation parameters compatible with each part. This line of work falls into the category of generative models, which typically learn the distribution of the input data and generate new samples from it. NAP [LDS*23] is a pioneering work that generates a full description of the articulated object, including part geometry, articulation graph, and joint parameters from random noises unconditionally. CAGE [LTMS23] extends the task in a conditional setting, where the object category and a part connectivity graph as given as the user-driven input. The output of these works targets the mesh representation composited in varied structures. The input for these works is not explicitly specified, but it is assumed that the input is aligned and noise-free, which is a common assumption for generation tasks.

4.2. Methodology

4.2.1. Intermediate Representation

In this section, we discuss the different choices of intermediate geometric representations used for part mobility analysis or object reconstruction, some of which are visualized in Figure 6.

Articulated part perception. For the work focusing on articulated part perception, the intermediate representation plays an important role in the methodology to facilitate the analysis of the part mobility by providing a more informative and expressive representation of the input geometry. Below we discuss each representation used for this category.

- **Surface patch.** The methods taking mesh as input commonly leverage the surface patches as the proxy for analysis. The local patches can be detected from slippable analysis as kinematic surfaces [GG04], which are used to associate with rigid motions for each part [XWY*09]. The surface patches can also be useful for revealing the self-similarity and symmetry on the mechanical part to characterize the motion configuration [MY*10]. The patches at the intersection area between components can

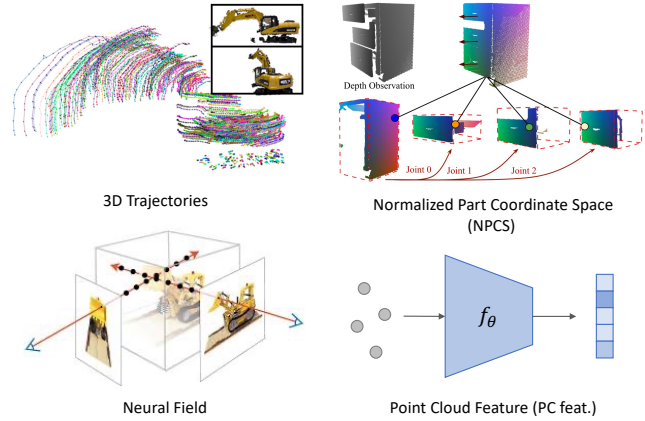


Figure 6: A selection of the intermediate representation used for articulated part perception. Figures reproduced from original papers [LWY*20; LWL*16; MST*21].

be leveraged to determine the supporting-supported relationship [SHL*14], which are informative for the part decomposition and structure analysis. Similarly, the interaction pattern described at the connecting region between surface patches is leveraged in Hu et al. [HLV*17] to measure the similarity between snapshots of the object for motion classification.

- **3D trajectory.** The 3D trajectory is a spatial-temporal representation that records the position of the points in the 3D space over time. For works capturing motion with a sequence of point clouds as input, extracting the 3D trajectory is a common choice for the intermediate representation [PG08; YLX*16; SCZ21]. As the points are initially unordered and untracked, establishing the correspondence between the points across different frames is a critical step for motion fitting. Once the trajectory is robustly estimated, the rigid pieces in the object can be grouped based on the transformational consistency of the points.
- **NAOCS, NPCS.** Articulation-aware Normalized Coordinate Space Hierarchy (ANCSH) consists of two coordinate spaces: Normalized Articulated Object Coordinate Space (NAOCS) and Normalized Part Coordinate Space (NPCS). It is a two-level hierarchical representation that was first proposed by Li et al. [LWY*20]. It is designed manually for articulated objects to align the instances and parts to a canonical space. NAOCS provides a canonical space on the object level that binds with a specific object category, while NPCS provides a canonical reference for each articulated part from a specific object instance. Knowing the object category as input, mapping the input shape to the NAOCS will produce part segmentation as it is defined by the structure template for each object category. The object category is not only a semantic label but also a structural label that defines the spatial arrangement of the parts. For example, a cabinet stacked with two drawers is classified separately from a cabinet with three drawers. Mapping the segmented parts to the NPCS can be useful for estimating the articulated state or pose for each part. Inspired by this idea, Liu et al. [LXX*22] and

GAPartNet [GXZ*23] propose to use NPCS as the intermediate representation for the articulated part analysis.

- **Point cloud feature (PC feat.).** Point cloud latent features are commonly used as the intermediate representation in the methods that take point clouds as input. The learning-based methods typically extract useful features in 3D from the input point cloud and use them to assign the part label and to regress the articulation parameters. PointNet [QSMG16], PointNet++ [QYSG17], and Sparse U-Net [GEvdM18] are the common feature extractors backbone in the existing works.
- **Image feature (image feat.).** Image latent features are used as the intermediate representation in the works that take RGB or RGB-D images as input. Similar to the point cloud features, the image features are extracted from pixels for later use in mobility part analysis. Some works follow this scheme to extract 2D features from depth images instead of per-point features from point clouds. ResNet [HZRS16] is the common backbone used for the feature extraction in the existing works.

Articulated object creation. For the work focusing on articulated object creation, the neural field is the most commonly used intermediate representation. A neural field is an implicit function parameterized by a neural network to represent the 3D geometry of the object in a continuous space. The network itself can be regarded as an SDF, occupancy field, neural density field, or neural radiance field, depending on the way the network is trained. Once constructed, the neural field can be queried to reconstruct or generate the objects and further be extracted to the mesh representation. To be more specific, A-SDF [MQK*21], Wei et al. [WCM*22], CARTO [HIZ*23] leverage neural SDF to represent the 3D geometry of the object. Ditto [JHZ22] constructs a neural occupancy field to indicate the occupancy status for each articulated part. CLA-NeRF [TLYS22] and PARIS [LMS23] both use a neural radiance field to represent both the shape and appearance of each part of the object.

4.2.2. Strategy and Supervision Signal

In this section, we discuss the different choices of algorithms or learning strategies adopted by the works when analyzing the articulated part structure.

Handcrafted methods. To infer part structure from the input geometry, early works rely on heuristic rules, human-designed descriptors, or mathematical algorithms to analyze the shape.

For the works that take a single-state mesh as input, the approaches based on the *slippage analysis*, *geometric properties*, and *relationship descriptors* are adopted to effectively segment the parts associated with the articulated motion. The *slippage analysis* is proposed by Gelfand and Guibas [GG04] to analyze the shape by discovering the slippable motions. It can segment the geometry into slippable portions which they call kinematic surfaces to associate with rigid motions for each part. Leveraging this approach, Xu et al. [XWY*09] detected primitive surfaces from the input object and link to motion joints that connect the articulated parts. Geometric attributes on the mesh surfaces can be characteristic for inferring different joint configurations. Based on this insight, Mitra et al. [MY*10] proposed to characterize the part motion by detecting the *self-similarity* and *symmetry* on the mesh surface for each mechanical part in a mathematical way. Spatial relationships

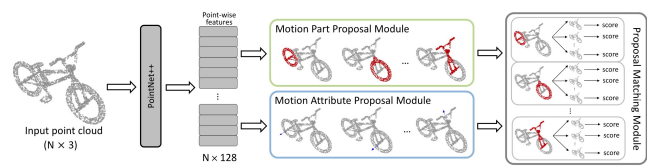


Figure 7: Shape2Motion [WZS*19] method pipeline that jointly segments the articulated parts and estimates motion attributes from a single point cloud as input. Figure reproduced from original paper [WZS*19].

between the parts are also informative for segmenting the components. Building on this observation, Sharf et al. [SHL*14] proposed to leverage the *supporting-supported* relationship between surface patches to decompose a scene into a set of objects and their articulated parts. Similarly, Hu et al. [HLV*17] proposed to leverage the *interaction pattern* between parts described by a human-designed interaction descriptor [HZvK*15]. Based on these geometric features, a distance function is designed to measure the similarity between snapshots of the object.

For the works that take point clouds observing a motion sequence as input, mathematical algorithms, and handcrafted descriptors are the two common choices to find the correspondence between adjacent frames. Once the correspondence is established, the part segmentation can be inferred from the 3D trajectory by clustering the points that are moving together. Given a known skeleton of the object in the first frame, Pekelny and Gotsman [PG08] proposed to trace the points that are skinned on each bone in the other frames using the *Iterated Closest Point (ICP)* algorithm [BM92]. Consuming an unlabelled point cloud sequence, Yuan et al. [YLX*16] and Shi et al. [SCZ21] proposed to leverage a deformable 3D shape registration algorithm [PB11] to estimate the 3D trajectory. Taking RGB-D sequence with additional color information as input, Li et al. [LWL*16] proposed to use *scene flow* [JSGC15] to produce a raw dense correspondence, from which *SIFT features* [Low04] are extracted from RGB images to refine and prune the initial proposals in the dense trajectory.

Handcrafted methods are excellent in terms of interpretability and computational efficiency, but their performance can be affected by the noisiness and incompleteness of the input data, initialization of the parameters in the algorithms, etc. With the advent of machine learning and deep learning, handcrafted methods are gradually replaced by learning-based methods, which can automatically learn the underlying features from the data and adapt to the input conditions. Later works after 2016 resort to data-driven methods to address the articulated object modeling problem, where the common practice converges into two main streams: supervised learning and self-supervised learning.

Supervised learning methods. For the task of articulated part perception, most of the existing works choose to learn the part assignment from a single state of the object in a supervised manner. Under this setting, the part segmentation is essentially formulated as a classification problem, where the goal is to predict the part label for each point in the input shape by transferring

the knowledge learned from training data. Most of the work under this group takes the point cloud optionally with color as input, and the classification model is trained with supervision from signals in 3D that are available from the human-annotated datasets. Shape2Motion [WZS*19] is the pioneering work that contributes both the dataset and the learning model for the task. Their network is trained to propose motion parts and attributes followed by a matching module to refine the final prediction. Figure 7 illustrates their method pipeline. With more 3D data with part annotations available, more works [YHY*19; ATK19; LWY*20; AFA*22; LXX*22; GXZ*23] follow the trend of joint prediction for segmentation and motion with varied network designs. Note that the part semantics are not always aligned with the mobility, but the data with semantic segmentation labels can be an initial source to refine the articulated part segmentation. Inspired by this idea, Liu et al. [LSH*23] contribute to extending a semantic segmentation dataset PartNet [MZC*19] with more articulation annotation by transferring the knowledge trained on PartNet-Mobility [XQM*20] dataset in a supervised way. To make the segmentation robust to arbitrary articulated states of the object, a common practice is to augment the training objects with different articulation states. It requires the data to be annotated with the articulation parameters, which is expensive and not at scale. In an effort to reduce reliance on articulation annotation available, Banana [DLS*23] design a network to be aware of the inter-part equivariance so that the knowledge acquired from objects only in resting states can be generalized to any other articulated states. This strategy allows the datasets compiled from static objects with segmentation labels to be usable for articulated part perception. Another line of work in this category attempts to perceive the part structure from RGB(-D) images. Given a pair of segmented parts in a single-view RGB-D image, Zeng et al. [ZLLK21] proposed FormNet to estimate connectedness and motion flow between the parts in a supervised way. OPD [JMCS22] and OPDMulti [SJSC23] proposed to learn 2D segmentations from single-view RGB images in a supervised manner, where the ground truth masks are available from the datasets they collected. Qian et al. [QJR*22] first proposed to infer mobility from an RGB video observing the human-object interaction. They train a network to regress 3D plane parameters to represent the articulated parts by leveraging the supervision from the 3D datasets without articulation annotation.

For the task of articulated object creation, supervised learning methods are also widely adopted. A-SDF [MQK*21], Wei et al. [WCM*22], and CARTO [HIZ*23] is a line of work for the object-level reconstruction of a shape that is deformable to a specific state. The model training is supervised by either ground truth SDF or multi-view RGB images. These supervision signals guide the model to learn an underlying distribution of the geometry that is associated with the articulation states. To further achieve part-level reconstruction, Ditto [JHZ22] and CLA-NeRF [TLYS22] require additional supervision from the part segmentation labels either on the point cloud or the multi-view images.

Self-supervised learning methods. Since the articulated part segmentation and articulation annotations are expensive to obtain, self-supervised learning methods have been proposed recently to reduce the reliance on these supervision signals. Taking a point cloud with a known number of parts as input, Liu et al. [LZH*23] proposed

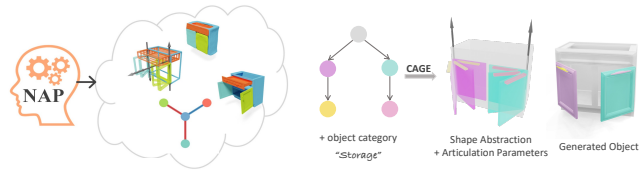


Figure 8: Generative models for articulated object synthesis. NAP [LDS*23] unconditionally synthesizes articulated objects with a full description, and CAGE [LTMS23] learns to generate objects with part structure conforming to the conditional input graph. Figures reproduced from original papers [LDS*23; LTMS23].

to learn the segmentation by extracting part-level equivariant features and factorizing the shape into canonical states. Through the iterative processing of canonicalization and then reconstruction using the factorized parameters, part decomposition can be learned by minimizing the reconstruction error on the input point cloud. PARIS [LMS23] reconstructs separate parts from a pair of multi-view RGB images that observe the object in two different states by leveraging differential rendering. By identifying the static and moving components through iterative optimization, the part-level geometry can be reconstructed by minimizing the photometric error between the input images and the rendered ones from the reconstructed surface.

Generative models. Leveraging the recent advance of denoising diffusion models [HJA20] in 3D generation, a line of work NAP [LDS*23] and CAGE [LTMS23] train generative models to synthesize the articulated object with part structure in a compositional way. The process of data distribution learning requires extensive data with complete annotations for articulated objects, such as a kinematic tree, part geometry represented as SDF in NAP [LDS*23] or as bounding primitives in CAGE [LTMS23]. See Figure 8 for the overview of the NAP and CAGE work. The acquisition of high-quality synthetic and real data with these annotations is still expensive and not at scale, which limits the generation ability of these models beyond the training data distribution.

5. Articulation Modeling

The articulation model of an object is a representation for describing the part mobility and the relationship between the parts. Understanding the kinematic structure of an object is a shared goal for both tasks of articulated part perception and articulated object creation. One of the key challenges in articulation modeling is how to deal with objects with different kinematic structures, such as the number of joints, the type of joints, the degrees of freedom (DoF) of the joints, and what hierarchy the joints are organized. In table 3, we summarize the related works from the perspective of articulation modeling in terms of the assumptions made on the articulation model, the representation of articulated motion, and the methodology used for articulation modeling. In the following, we will discuss each aspect in detail.

	Assumption			Articulated Motion Representation						Methodology	
	# joints	# DoF	obj. cat.	joint type	joint axis	joint state	motion range	kine. tree	deform. flow	strategy	supervision
Articulated Part Perception											
Xu et al. [XWY*09]	X	3	X	✓	✓	X	✓	X	X	handcrafted	-
Mitra et al. [MY*10]	✓	3	X	✓	✓	X	✓	✓	X	handcrafted	-
Sharf et al. [SHL*14]	X	1	X	✓	✓	X	X	✓	X	handcrafted	-
Li et al. [LWL*16]	X	3	X	✓	✓	X	X	✓	X	handcrafted	-
Hu et al. [HLV*17]	✓	1	X	✓	✓	X	✓	X	X	supervised	distance func.
Yi et al. [YHL*18]	X	3	X	X	X	X	X	X	✓	supervised	displacement
Wang et al. [WZS*19]	X	1	X	✓	✓	X	X	X	X	supervised	joint params.
RPM-Net [YHY*19]	X	arbitrary	X	✓	✓	X	X	✓	✓	supervised	displacement
Abbatematteo et al. [ATK19]	✓	1	✓	X	✓	✓	X	X	X	supervised	joint params.
Li et al. [LWY*20]	✓	1	✓	✓	✓	✓	✓	X	X	supervised	joint params.
Shi et al. [SCZ21]	✓	1	X	✓	✓	X	✓	X	X	self-supervised	point corr.
ScrewNet [JLCN21]	✓	1	X	✓	✓	✓	X	X	X	supervised	joint params.
Abdul-Rashid et al. [AFA*22]	X	1	X	✓	✓	X	X	✓	X	supervised	joint params.
Qian et al. [QJR*22]	X	1	X	✓	✓	X	X	X	X	supervised	joint params.
OPD [JMSC22]	✓	1	X	✓	✓	X	X	X	X	supervised	joint params.
OPDMulti [SJS23]	X	1	X	✓	✓	X	X	X	X	supervised	joint params.
Liu et al. [LXX*22]	X	1	X	✓	✓	✓	X	✓	X	supervised	joint params.
Liu et al. [LZH*23]	✓	1	✓	✓	✓	✓	X	X	X	self-supervised	point corr.
GAPartNet [GXZ*23]	X	1	X	✓	✓	✓	X	X	X	supervised	NPCS maps
Liu et al. [LSH*23]	X	1	X	✓	✓	X	X	✓	X	supervised	joint params.
Articulated Object Creation											
A-SDF [MQK*21]	✓	1	✓	X	X	X	X	X	✓	supervised	joint params.
Wei et al. [WCM*22]	✓	1	✓	X	X	✓	X	X	✓	supervised	joint params.
Ditto [JHZ22]	✓	1	X	✓	✓	✓	X	X	X	supervised	joint params.
CLA-NeRF [TLYS22]	✓	1	✓	X	✓	X	X	X	X	handcrafted	-
CARTO [HIZ*23]	✓	1	X	✓	X	✓	X	X	X	supervised	joint params.
PARIS [LMS23]	✓	1	X	✓	✓	✓	X	X	X	self-supervised	multi-view RGB
NAP [LDS*23]	X	1	X	✓	✓	✓	✓	✓	X	generative	-
CAGE [LTMS23]	X	1	X	✓	✓	X	✓	✓	X	generative	-

Table 3: Summary of the related works from the perspective of articulation modeling in terms of the assumptions about articulation structure each work makes, the representation of articulated motion, and the methodology used in the motion estimation process.

5.1. Representations of Articulated Motion

Deformation flow. Deformation flow is a general form of motion representation that can describe not only the rigid motion but also the motion as free-form deformation. The deformation flow describes the motion as a vector field, where each vector denotes the spatial displacement for each point in space. Yi et al. [YHL*18] and RPM-Net [YHY*19] both choose the deformation flow to represent the articulation motion as the output of their networks. They both leverage recurrent neural networks, such as LSTM, to directly predict the deformation flow from the input point cloud. An example of the deformation field prediction network is illustrated in fig. 10. And the training is supervised by the ground truth deformation flow. One of the advantages of using deformation flow is that it can handle the motion of the parts in a free-form manner. This flexibility is useful when the number of parts is unknown as the deformation on each point can be later grouped to form articulated parts in arbitrary numbers. The deformation flow can be particularly useful when modeling the non-trivial motion. For example, RPM-Net can model the opening and closing motion of the cover of an umbrella, which is a non-linear motion that cannot be represented by a simple rigid transformation. Similarly, A-SDF [MQK*21] and Wei et al. [WCM*22] also leverage the deformation field to represent articulation for the objects in the task of articulated object creation. They both leverage the neural field as a coordinate-based network to implicitly represent the deformation by conditioning on the articulation state of the object. Once the network is trained, the reconstructed object can be deformed to different states by querying

the neural field. However, one of the limitations of using deformation flow or field is that it is less constrained than a strict form of rigid transformation when we are dealing with rigidly moving parts. To alleviate this issue, Yi et al. [YHL*18] tried to optimize the output deformation flow with an As-rigid-as-possible (ARAP) objective [SA07] to preserve the local rigidity of the parts. RPM-Net [YHY*19] also tried to enforce rigidity by adding a regularization term to the training process.

Joint parameters. The most common way to represent the motion for the articulated part is to use the joint parameters. This representation is standard in the robotics community, where the joint parameters can be directly recorded in the URDF file for simulation and control. The complete set of joint parameters includes the joint type, the joint axis, the joint state, and joint limits. The joint types that are most commonly seen in passively articulated objects are the 1 DoF joints, such as the revolute joint and the prismatic joint. The joint with more than 1 DoF, such as the ball joint, can also be seen in more complex objects, e.g. mechanical parts. An overview of the joint types that are covered in the literature is shown in fig. 2. The joint axis is represented by a direction vector as the axis orientation and a position vector as the location of the joint pivot. The joint state denotes the current articulation state of the joint, which is usually represented by the rotation angle for the revolute joint and the translational distance for the prismatic joint. The joint limits denote the range of motion that is allowed for the joint at two ends. An illustration of the joint parameters estimated from the input mesh is shown in fig. 9. Most of the existing work formulates the mobil-

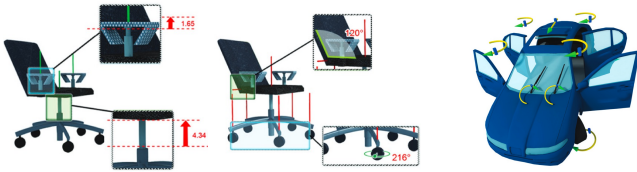


Figure 9: Examples of joint parameters output from Sharf et al. [SHL*14] by decomposing an object into a mobility tree and estimating articulation with heuristic rules. Figure reproduced from original paper [SHL*14].

ity analysis as a joint parameter estimation problem from various inputs, which can be converted into a combination of classification and regression problems. Usually, only a subset of these parameters are estimated as the output, depending on the specific task. While predicting the joint type with its axis, some works focus on estimating the current joint state of the object being observed [ATK19; LWY*20; LXX*22; LZH*23; GXZ*23] whose task is also termed as the part pose estimation problem. And other works aim to understand the joint limit [XWY*09; MYY*10; HLV*17; LWY*20; SCZ21], which can be more challenging as it requires the network to understand the physical constraints of the object.

Kinematic tree. The kinematic tree serves as a hierarchical representation of the joints and parts of an articulated object. This structure is particularly crucial for depicting objects with multiple parts and joints, as it clearly outlines the dependencies between them. An example of the kinematic tree is shown in fig. 4. Extracting the kinematic tree is a challenging problem since the kinematic structure varies from object to object with different numbers of parts and is structured in different topologies. Some existing works bypass this step by assuming the kinematic tree is predefined by the semantic object class [ATK19; LWY*20; LZH*23; MQK*21; WCM*22] or by simplifying the problem to a single part with a single joint [HLV*17; JLCN21; JHZ22; TLYS22; HIZ*23; LMS23]. And a large portion of the works dealing with multiple parts and joints also do not explicitly represent the kinematic tree in their model [YHL*18; WZS*19; YHY*19; SCZ21; SJSC23; GXZ*23]. There are only a few works that explicitly model the kinematic tree. The early work by Mitra et al. [MY*10] and Sharf et al. [SHL*14] extract the kinematic tree from the input mesh. They share a similar idea to leverage the intersection of the part surfaces to identify the connection between the parts. Mitra et al. [MY*10] detected the contact area between the parts to determine the topology of the interaction graph, and Sharf et al. [SHL*14] applied a segment-cluster process to compute the support-tree hierarchy from decomposed parts. Later works by Abdul-Rashid et al. [AFA*22] and Liu et al. [LXX*22] attempt to predict the kinematic tree from the input point cloud. Abdul-Rashid et al. [AFA*22] use a graph neural network to predict the labeling for each segmented part. Liu et al. [LXX*22] leverage an encoder-decoder architecture backbone by PointNet++ [QYSG17] to predict the connectivity between each pair of parts.

5.2. Assumption on the Articulation Model

There are several axes of assumptions made on the articulation model by the existing works. These assumptions usually are made to simplify the problem to make it more tractable. We summarize these common assumptions to show the research progress below.

The number of degrees of freedom (# DoF). Most of the existing works assume to only deal with the 1 DoF joints, including the revolute joint, the prismatic joint, and the helical joint. These joints are the most common types of joints that are seen in the daily articulated objects, such as cabinets, dishwashers, refrigerators, etc. And their motion is constrained to be linear and easily represented by a combination of rotation and translation constrained by a single joint axis. When dealing with more complex objects, such as delicate mechanical systems, the joint with more than 1 DoF, such as the ball joint, can also be seen. A few methods [XWY*09; MYY*10; LWL*16; YHL*18] have been proposed to deal with joint up to 3 DoF, which is achieved by either leveraging the deformation field to represent the motion or by combining with the handcrafted features on the part geometry to classify the joint type.

Known object category (Obj. cat.) As mentioned above, one of the main challenges in articulation modeling is the diversity of the kinematic structure of the objects. One way to generalize the knowledge learned from the training data to the unseen objects is to build a category-level model, which is based on the assumption that the objects in the same category share the same structures. The structure implied by this line of work [ATK19; LWY*20; LZH*23] is not only about the kinematic structure but also the geometric structure. It means that the number of parts, how the parts are arranged spatially, and the way each part is articulated are all shared among the objects in the same category. This presents a strong prior for the model to learn from. In some cases, this level of category can be aligned with the semantic class of the objects whose kinematic structure is relatively simple and fixed, such as eyeglasses, laptops, scissors, etc. However, this assumption is highly constrained since this is not the case for more objects with more complex structures. For example, a refrigerator with two doors stacked vertically and a refrigerator with two doors stacked horizontally are considered to be from the same semantic class but with different structures. A separate model should be trained for each kind of refrigerator under this assumption. It makes the model less scalable to real-world applications with a diversity of objects.

Known number of joints (# Joints). Whether the number of joints is known is another assumption made by the existing works. The previous assumption of knowing the object category inherently implies that the number of joints is known. There are also a few works [HLV*17; JLCN21; JHZ22; TLYS22; HIZ*23; LMS23] assume to estimate only one joint at a time but not knowing the motion the joint exhibits. It can be useful to generalize the method across objects with varied kinematic ways for simple objects. However, it is not scalable to more complex objects with more than one joint, highlighting a limitation in this line of work.

5.3. Methodologies for Articulation Modeling

Handcrafted methods. Early works leverage non-learning-based methods to infer the articulation model. Xu et al. [XWY*09] ap-

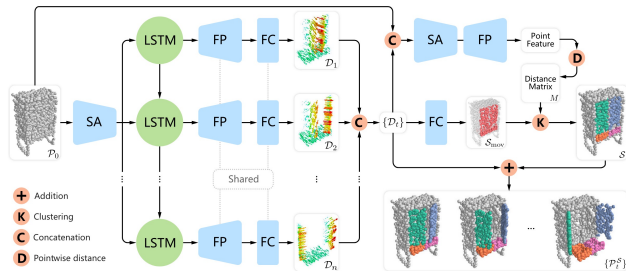


Figure 10: An illustration of using a recurrent neural network to predict deformation field for each articulated part from RPM-Net [YHY*19]. Figure reproduced from original paper [YHY*19].

plied slippage analysis to compute a set of joint parameters from a mesh input. The output includes the number of rotational and translational degrees of freedom which is thresholded to be up to three, and the corresponding axes for each slippable motion. To determine the motion range for each joint, they design a trial-and-error bisection process to decide the feasibility. They estimate the joint limit by probing the possible motions until the motion results in an unrealistic geometry, e.g. penetration occurs between the parts. Mitra et al. [MY*10] defines a list of heuristic rules to determine the joint type, axis, and limit for each mechanical part. Once detecting the supporting-supported relationship between parts, Sharf et al. [SHL*14] extract this connectivity to form a kinematic tree. By leveraging the recurrence of the objects and parts in the scene with different poses, they design rules to determine the motion type and axes by using the PCA axes of the part geometry. Taking a point cloud sequence as input, Li et al. [LWL*16] propose a random sampling consensus method in 4D to fit motion on 3D trajectories. Then these motions can be clustered to convert to the joint parameters mathematically and the joints can be organized in a mobility graph.

Supervised learning. Most of the later works leverage supervised learning to predict the articulation model as shown in table 3. Depending on the articulation representation chosen, the training process can be supervised by the ground truth joint parameters or deformation flows. Different from these works that explicitly articulation annotation, Hu et al. [HLV*17] propose to classify the motion type from a mesh input through a metric learning model. The training is supervised by a distance function with several constraints to pull the object pairs with the same motion type closer and push the ones with different motion types away. Illustrated in fig. 11, Li et al. [LWY*20] proposed a method to canonicalize an articulated object using coordinate spaces at object-level as Normalized Articulated Object Coordinate Space (NAOCS) and in part-level as Normalized Part Coordinate Space (NPCS). The canonical space is defined for each object category. Once the object is canonicalized, the joint parameters can be estimated by fitting the attributes on the template object. GPartNet [GXZ*23] follows the definition of the above coordinate space for each part to supervise the training with the ground truth NPCS maps. After learning the NPCS mapping from the input state to a canonical state, they apply the Umeyama algorithm [Ume91] to estimate the rigid transformation from the mapping to obtain the joint parameters.

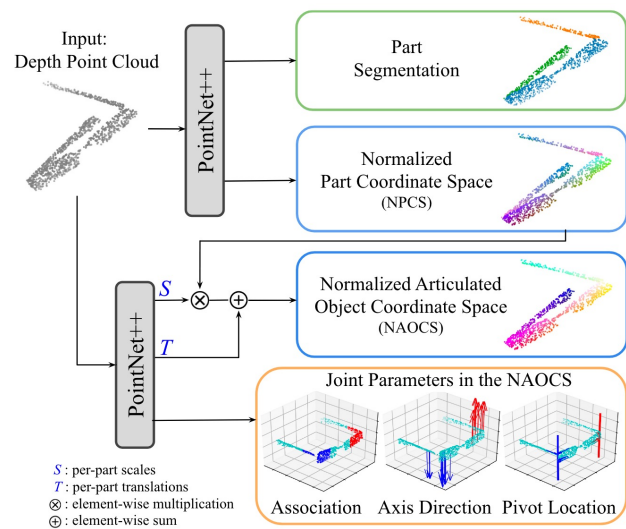


Figure 11: Overview of the ANCSH method proposed by Li et al. [LWY*20] that maps the object at an arbitrary articulated state to a hierarchical normalized coordinate space. Figure reproduced from original paper [LWY*20].

Self-supervised learning. Although supervised learning is effective in learning the articulation model, it requires a large amount of labeled data to train the model. To alleviate this issue, a few works propose to leverage self-supervised learning to estimate the articulation model. Shi et al. [SCZ21] leverage a network to learn the point correspondence across different frames of the point cloud sequence from the input. Once the 3D trajectories are constructed by learning the consistent point correspondence, the joint parameters can be estimated by fitting the motion on the trajectories. One limitation of this work is that it requires observing a motion sequence as input, which usually involves human intervention while capturing the data. This tedious motion-capturing process makes the method less scalable to real-world applications. To alleviate this issue, Liu et al. [LZH*23] propose to learn the joint parameters from a single point cloud in a self-supervised manner. By factorizing each articulated part to its own normalized canonical space and then reconstructing it back to the input state using an optimized transformation, the reconstruction consistency in this process can be used to supervise the training. However, this work is limited by the assumption that the model is category-specific, meaning the geometric and kinematic structure of the object is known in advance. In the task of articulated object reconstruction, PARIS [LMS23] proposes to learn the joint parameters from a pair of multi-view RGB images self-supervisedly. The estimation of the joint parameters is supervised by the consistency of the photometric loss across two articulated states given as input. Although this method can be generalizable across object categories and does not rely on any 3D or articulation supervision, it assumes only one part is moving in the input observation.

6. Future Discussion and Conclusion

As the field of 3D modeling for articulated objects continues to evolve, several promising research directions have emerged from our survey. The key areas for future exploration include:

Generalization across objects with varying structures. Handling the diversity of articulated objects with large variations in the geometric and kinematic structures is one of the major challenges in the field. Existing methods are often designed for specific object categories or assume a certain level of prior knowledge about the object is available to the system. Given limited training data and the high cost of data collection, it is important to develop methods that can generalize to novel objects in a data-efficient manner.

Improved accuracy and robustness in the perception models. Although significant progress has been made in the development of articulated part perception models, there is still a need for more accurate and robust methods that can handle occlusions, clutter, and noise in the input data. This is particularly important for applications such as robotics and autonomous driving, where the models need to be able to handle real-world scenarios.

Incorporating physical constraints into the modeling process. Articulated objects are subject to physical constraints that are often not considered in the existing modeling process. Examining the physical constraints of the objects can help to improve the accuracy of the articulation modeling and enable the generation of physically plausible animations and simulations. This consideration is also useful in articulated object creation for applications such as virtual reality, where the goal is to create realistic and physically plausible models.

Detailed interior and part geometry modeling. Modeling interior and part-level geometry in detail for articulated objects is an important aspect of geometric modeling that has not been fully explored. For the object reconstruction task, the quality of the models is often limited by the quality of the input observation which is often partial due to view occlusion and sensor noises. Developing methods that can handle these challenges and generate detailed and accurate 3D models is an important direction for future research.

Generative AI for articulated objects. The recent advances in generative models and large language models have shown great potential in various computer vision and graphics tasks. Applying these techniques to the domain of articulated object modeling can be a promising direction to provide a better understanding and accessibility of the complex 3D articulated structures.

Conclusion. This survey systematically reviewed the progress in 3D modeling for articulated objects, with a focus on articulated part perception and articulated object creation. By examining the development of articulated object modeling across the two axes of geometric processing and articulation modeling, we have identified significant advancements and ongoing challenges in the field. We also highlighted the potential research directions that can drive future progress in the field. We hope that this survey serves as a foundational reference for researchers and practitioners in computer vision and graphics, offering insights into the complexities of articulated object modeling and inspiring new research in this area.

Acknowledgements. This work was funded in part by a Canada Research Chair and an NSERC Discovery grant. We thank Angel X. Chang and Hao (Richard) Zhang for helpful discussions and feedback on preliminary drafts.

References

- [ACTD07] ALLIEZ, PIERRE, COHEN-STEINER, DAVID, TONG, YIYING, and DESBRUN, MATHIEU. “Voronoi-based variational reconstruction of unoriented point sets”. *Proceedings of the Eurographics Symposium on Geometry Processing*. Vol. 7. 2007, 39–48 7.
- [AFA*22] ABDUL-RASHID, HAMEED, FREEMAN, MILES, ABBATEMATTEO, BEN, et al. “Learning to infer kinematic hierarchies for novel object instances”. *Proceedings of the International Conference on Robotics and Automation (ICRA)*. 2022, 8461–8467 6, 7, 10–12.
- [AM23] AYGÜN, MEHMET and MAC AODHA, OISIN. “SAOR: Single-View Articulated Object Reconstruction”. *arXiv preprint arXiv:2303.13514* (2023) 3.
- [ATK19] ABBATEMATTEO, BEN, TELLEX, STEFANIE, and KONIDARIS, GEORGE. “Learning to generalize kinematic models to novel objects”. *Proceedings of the Conference on Robot Learning (CoRL)*. 2019, 1289–1299 6, 10–12.
- [BKY*22] BERGMAN, ALEXANDER, KELLNHOFFER, PETR, YIFAN, WANG, et al. “Generative neural articulated radiance fields”. *Advances in neural information processing systems (NeurIPS)* 35 (2022), 19900–19916 3.
- [BM92] BESL, PAUL J and MCKAY, NEIL D. “Method for registration of 3-D shapes”. *Sensor fusion IV: control paradigms and data structures*. Vol. 1611. Spie. 1992, 586–606 9.
- [CFG*15] CHANG, ANGEL X, FUNKHOUSER, THOMAS, GUIBAS, LEONIDAS, et al. “ShapeNet: An information-rich 3D model repository”. *arXiv preprint arXiv:1512.03012* (2015) 4.
- [CJS*23] CHEN, XU, JIANG, TIANJIAN, SONG, JIE, et al. “Fast-SNARF: A fast deformer for articulated neural fields”. *IEEE Transactions on Pattern Analysis and Machine Intelligence* (2023) 1.
- [DLS*23] DENG, CONGYUE, LEI, JIAHUI, SHEN, BOKUI, et al. “Banana: Banach fixed-Point network for pointcloud segmentation with inter-part equivariance”. *arXiv preprint arXiv:2305.16314* (2023) 6, 7, 10.
- [GES21] GADRE, SAMIR YITZHAK, EHSANI, KIANA, and SONG, SHURAN. “Act the part: Learning interaction strategies for articulated object part discovery”. *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*. 2021, 15752–15761 1.
- [GEvdM18] GRAHAM, BENJAMIN, ENGELCKE, MARTIN, and van der MAATEN, LAURENS. “3D Semantic Segmentation with Submanifold Sparse Convolutional Networks”. *CVPR* (2018), 9224–9232 9.
- [GG04] GELFAND, NATASHA and GUIBAS, LEONIDAS J. “Shape segmentation using local slippage analysis”. *Proceedings of the Eurographics Symposium on Geometry Processing*. 2004, 214–223 8, 9.
- [GGS*19] GAO, XIAOFENG, GONG, RAN, SHU, TIANMIN, et al. “Vrkitthen: an interactive 3d virtual environment for task-oriented learning”. *arXiv preprint arXiv:1903.05757* (2019) 1.
- [GXZ*23] GENG, HAORAN, XU, HELIN, ZHAO, CHENGYANG, et al. “GAPartNet: Cross-category domain-generalizable object perception and manipulation via generalizable and actionable parts”. *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. 2023, 7081–7091 5–7, 9–13.
- [HIZ*23] HEPPERT, NICK, IRSHAD, MUHAMMAD ZUBAIR, ZAKHAROV, SERGEY, et al. “CARTO: Category and Joint Agnostic Reconstruction of ARTiculated Objects”. *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. 2023, 21201–21210 6, 8–12.
- [HJA20] HO, JONATHAN, JAIN, AJAY, and ABBEEL, PIETER. “Denoising diffusion probabilistic models”. *Advances in neural information processing systems (NeurIPS)* 33 (2020), 6840–6851 10.

- [HLV*17] HU, RUIZHEN, LI, WENCHAO, VAN KAICK, OLIVER, et al. “Learning to predict part mobility from a single static snapshot”. *ACM Transactions on Graphics (TOG)* 36.6 (2017), 1–13 1, 3–9, 11–13.
- [HNOS15] HAUSMAN, KAROL, NIEKUM, SCOTT, OSENTOSKI, SARAH, and SUKHATME, GAURAV S. “Active articulation model estimation through interactive perception”. *Proceedings of the International Conference on Robotics and Automation (ICRA)*. 2015, 3305–3312 1.
- [HSvK18] HU, RUIZHEN, SAVVA, MANOLIS, and van KAICK, OLIVER. “Functionality representations and applications for shape analysis”. *Computer Graphics Forum*. Vol. 37. 2. 2018, 603–624 2.
- [HZRS16] HE, KAIMING, ZHANG, XIANGYU, REN, SHAOQING, and SUN, JIAN. “Deep residual learning for image recognition”. *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. 2016, 770–778 9.
- [HSvK*15] HU, RUIZHEN, ZHU, CHENYANG, van KAICK, OLIVER, et al. “Interaction context (ICON) towards a geometric functionality descriptor”. *ACM Transactions on Graphics (TOG)* 34.4 (2015), 1–12 9.
- [Inc17a] INC., TRIMBLE. *3D Warehouse*. Accessed: 2024-01-22. 2017. URL: <https://3dwarehouse.sketchup.com/>.
- [Inc17b] INC., TRIMBLE. *SketchUp*. Accessed: 2024-01-22. 2017. URL: <https://www.sketchup.com/>.
- [JHZ22] JIANG, ZHENYU, HSU, CHENG-CHUN, and ZHU, YUKE. “Ditto: Building digital twins of articulated objects from interaction”. *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. 2022, 5616–5626 1, 6, 8–12.
- [JLCN21] JAIN, AJINKYA, LIOUTIKOV, RUDOLF, CHUCK, CALEB, and NIEKUM, SCOTT. “ScrewNet: Category-Independent Articulation Model Estimation From Depth Images Using Screw Theory”. *Proceedings of the International Conference on Robotics and Automation (ICRA)*. 2021, 13670–13677 4, 6, 7, 11, 12.
- [JMSc22] JIANG, HANXIAO, MAO, YONGSEN, SAVVA, MANOLIS, and CHANG, ANGEL X. “OPD: Single-view 3D openable part detection”. *Proceedings of the European Conference on Computer Vision (ECCV)*. 2022, 410–426 4–7, 10, 11.
- [JSGC15] JAIMEZ, MARIANO, SOUIAI, MOHAMED, GONZALEZ-JIMENEZ, JAVIER, and CREMERS, DANIEL. “A primal-dual framework for real-time dense RGB-D scene flow”. *Proceedings of the International Conference on Robotics and Automation (ICRA)*. 2015, 98–104 9.
- [KBH06] KAZHDAN, MICHAEL, BOLITHO, MATTHEW, and HOPPE, HUGUES. “Poisson surface reconstruction”. *Proceedings of the Eurographics Symposium on Geometry Processing*. Vol. 7. 2006 7.
- [KMH*17] KOLVE, ERIC, MOTTAGHI, ROOZBEH, HAN, WINSON, et al. “AI2-THOR: An interactive 3d environment for visual AI”. *arXiv preprint arXiv:1712.05474* (2017) 1.
- [LDS*23] LEI, JIAHUI, DENG, CONGYUE, SHEN, BOKUI, et al. “NAP: Neural 3D articulation prior”. *arXiv preprint arXiv:2305.16315* (2023) 1, 6, 8, 10, 11.
- [LLL*24] LI, ZIZHANG, LITVAK, DOR, LI, RUINING, et al. “Learning the 3D Fauna of the Web”. *arXiv preprint arXiv:2401.02400* (2024) 3.
- [LMS23] LIU, JIAYI, MAHDAVI-AMIRI, ALI, and SAVVA, MANOLIS. “PARIS: Part-level reconstruction and motion analysis for articulated objects”. *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 2023, 352–363 1, 6, 8–13.
- [Low04] LOWE, DAVID G. “Distinctive image features from scale-invariant keypoints”. *International journal of computer vision* 60 (2004), 91–110 9.
- [LSH*23] LIU, GENGXIN, SUN, QIAN, HUANG, HAIBIN, et al. “Semi-weakly supervised object kinematic motion prediction”. *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. 2023, 21726–21735 6, 7, 10, 11.
- [LTMS23] LIU, JIAYI, TAM, HOU IN IVAN, MAHDAVI-AMIRI, ALI, and SAVVA, MANOLIS. “CAGE: Controllable Articulation GEneration”. *arXiv preprint arXiv:2312.09570* (2023) 1, 6, 8, 10, 11.
- [LWL*16] LI, HAO, WAN, GUOWEI, LI, HONGHUA, et al. “Mobility fitting using 4D RANSAC”. *Computer Graphics Forum*. Vol. 35. 5. 2016, 79–88 6–9, 11–13.
- [LWP*23] LEI, JIAHUI, WANG, YUFU, PAVLAKOS, GEORGIOS, et al. “GART: Gaussian articulated template models”. *arXiv preprint arXiv:2311.16099* (2023) 3.
- [LWY*20] LI, XIAOLONG, WANG, HE, YI, LI, et al. “Category-level articulated object pose estimation”. *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. 2020, 3706–3715 6–8, 10–13.
- [LXF*22] LIU, LIU, XU, WENQIANG, FU, HAOYUAN, et al. “AKB-48: A real-world articulated object knowledge base”. *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. 2022, 14809–14818 5.
- [LXX*22] LIU, LIU, XUE, HAN, XU, WENQIANG, et al. “Toward real-world category-level articulation pose estimation”. *IEEE Transactions on Image Processing* 31 (2022), 1072–1083 5–8, 10–12.
- [LZH*23] LIU, XUEYI, ZHANG, JI, HU, RUIZHEN, et al. “Self-supervised category-Level articulated object pose estimation with part-Level SE(3) Equivariance”. *Proceedings of the International Conference on Learning Representations (ICLR)*. 2023 6, 7, 10–13.
- [LZWL23] LI, ZHE, ZHENG, ZERONG, WANG, LIZHEN, and LIU, YEBIN. “Animatable gaussians: Learning pose-dependent gaussian maps for high-fidelity human avatar modeling”. *arXiv preprint arXiv:2311.16096* (2023) 1, 3.
- [MEB19] MARTÍN-MARTÍN, ROBERTO, EPPNER, CLEMENS, and BROCK, OLIVER. “The RBO dataset of articulated objects and interactions”. *The International Journal of Robotics Research* 38.9 (2019), 1013–1019 5.
- [MGM*21] MO, KAICHUN, GUIBAS, LEONIDAS J, MUKADAM, MUSTAFA, et al. “Where2act: From pixels to actions for articulated 3D objects”. *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*. 2021, 6813–6823 1.
- [MQK*21] MU, JITENG, QIU, WEICHAO, KORTYLEWSKI, ADAM, et al. “A-SDF: Learning Disentangled Signed Distance Functions for Articulated Shape Representation”. *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*. 2021, 12981–12991 6, 7, 9–12.
- [MST*21] MILDENHALL, BEN, SRINIVASAN, PRATUL P, TANCIK, MATTHEW, et al. “NeRF: Representing scenes as neural radiance fields for view synthesis”. *Communications of the ACM* 65.1 (2021), 99–106 8.
- [Mue19] MUELLER, ANDREAS. “Modern robotics: Mechanics, planning, and control [bookshelf]”. *IEEE Control Systems Magazine* 39.6 (2019), 100–102 2.
- [MY*10] MITRA, NILOY J, YANG, YONG-LIANG, YAN, DONG-MING, et al. “Illustrating how mechanical assemblies work”. *ACM Transactions on Graphics (TOG)* 29.4 (2010), 58 6–9, 11–13.
- [MZC*19] MO, KAICHUN, ZHU, SHILIN, CHANG, ANGEL X, et al. “PartNet: A large-scale benchmark for fine-grained and hierarchical part-level 3D object understanding”. *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. 2019, 909–918 4, 10.
- [MZJ*22] MAO, YONGSEN, ZHANG, YIMING, JIANG, HANXIAO, et al. “MultiScan: Scalable RGBD scanning for 3D environments with articulated objects”. *Advances in neural information processing systems (NeurIPS)* 35 (2022), 9058–9071 5.
- [PB11] PAPAJOV, CHAVDAR and BURSCHEKA, DARIUS. “Deformable 3D shape registration based on local similarity transforms”. *Computer Graphics Forum*. Vol. 30. 5. 2011, 1493–1502 9.
- [PG08] PEKELNY, YURI and GOTSMAN, CRAIG. “Articulated object reconstruction and markerless motion capture from depth video”. *Computer Graphics Forum*. Vol. 27. 2. 2008, 399–408 6–9.
- [PRB*18] PUIG, XAVIER, RA, KEVIN, BOBEN, MARKO, et al. “Virtual-home: Simulating household activities via programs”. *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. 2018, 8494–8502 1.

- [PŠDC22] PEJIĆ, PETRA, ŠIMUNDIĆ, VALENTIN, DŽIJAN, MATEJ, and CUPEC, ROBERT. “Articulated Objects: From Detection to Manipulation—Survey”. *International Conference on Intelligent Autonomous Systems*. 2022, 495–508 2.
- [PUS*23] PUIG, XAVIER, UNDERSANDER, ERIC, SZOT, ANDREW, et al. “Habitat 3.0: A co-habitat for humans, avatars and robots”. *arXiv preprint arXiv:2310.13724* (2023) 1.
- [QF23] QIAN, SHENGYI and FOUHEY, DAVID F. “Understanding 3d object interaction from a single image”. *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*. 2023, 21753–21763 1.
- [QJR*22] QIAN, SHENGYI, JIN, LINYI, ROCKWELL, CHRIS, et al. “Understanding 3D object articulation in internet videos”. *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. 2022, 1599–1609 6, 7, 10, 11.
- [QSMG16] QI, C, SU, H, MO, K, and GUIBAS, LJ. “PointNet: deep learning on point sets for 3D classification and segmentation. 2017”. *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. 2016, 77–85 9.
- [QWM*23] QIAN, ZHIYIN, WANG, SHAOFEI, MIHAJLOVIC, MARKO, et al. “3dgs-avatar: Animatable avatars via deformable 3d gaussian splatting”. *arXiv preprint arXiv:2312.09228* (2023) 1, 3.
- [QYSG17] QI, CHARLES RUIZHONGTAI, YI, LI, SU, HAO, and GUIBAS, LEONIDAS J. “Pointnet++: Deep hierarchical feature learning on point sets in a metric space”. *Advances in neural information processing systems (NeurIPS)* 30 (2017) 9, 12.
- [SA07] SORKINE, OLGA and ALEXA, MARC. “As-rigid-as-possible surface modeling”. *Proceedings of the Eurographics Symposium on Geometry Processing*. Vol. 4. 2007, 109–116 11.
- [SBR22] SU, SHIH-YANG, BAGAUTDINOV, TIMUR, and RHODIN, HELGE. “Danbo: Disentangled articulated neural body representations via graph neural networks”. *European Conference on Computer Vision*. Springer. 2022, 107–124 3.
- [SCZ21] SHI, YAHAO, CAO, XINYU, and ZHOU, BIN. “Self-Supervised Learning of Part Mobility from Point Cloud Sequence”. *Computer Graphics Forum*. Vol. 40. 6. 2021, 104–116 6–9, 11–13.
- [SHL*14] SHARF, ANDREI, HUANG, HUI, LIANG, CHENG, et al. “Mobility-trees for indoor scenes manipulation”. *Computer Graphics Forum*. Vol. 33. 1. 2014, 2–14 1, 4, 6–9, 11–13.
- [SJSC23] SUN, XIAOHAO, JIANG, HANXIAO, SAVVA, MANOLIS, and CHANG, ANGEL XUAN. “OPDMulti: Openable Part Detection for Multiple Objects”. *arXiv preprint arXiv:2303.14087* (2023) 5–7, 10–12.
- [SKM*19] SAVVA, MANOLIS, KADIAN, ABHISHEK, MAKSYMETS, OLEKSANDR, et al. “Habitat: A platform for embodied ai research”. *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*. 2019, 9339–9347 1.
- [SYZR21] SU, SHIH-YANG, YU, FRANK, ZOLLHÖFER, MICHAEL, and RHODIN, HELGE. “A-NeRF: Articulated neural radiance fields for learning human shape, appearance, and pose”. *Advances in neural information processing systems (NeurIPS)* 34 (2021), 12278–12291 3.
- [TLYS22] TSENG, WEI-CHENG, LIAO, HUNG-JU, YEN-CHEN, LIN, and SUN, MIN. “CLA-NeRF: Category-level articulated neural radiance field”. *Proceedings of the International Conference on Robotics and Automation (ICRA)*. 2022, 8454–8460 6, 8–12.
- [Ume91] UMEYAMA, SHINJI. “Least-squares estimation of transformation parameters between two point patterns”. *IEEE Transactions on Pattern Analysis & Machine Intelligence* 13.04 (1991), 376–380 13.
- [WCK19] WENG, CHUNG-YI, CURLESS, BRIAN, and KEMELMACHER-SHLIZERMAN, IRA. “Photo wake-up: 3D character animation from a single photo”. *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. 2019, 5908–5917 1.
- [WCM*22] WEI, FANGYIN, CHABRA, ROHAN, MA, LINGNI, et al. “Self-supervised neural articulated shape and appearance models”. *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 2022, 15816–15826 6, 8–12.
- [WLJ*23] WU, SHANGZHE, LI, RUINING, JAKAB, TOMAS, et al. “Magicpony: Learning articulated 3D animals in the wild”. *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 2023, 8792–8802 3.
- [WSGT22] WANG, SHAOFEI, SCHWARZ, KATJA, GEIGER, ANDREAS, and TANG, SIYU. “Arah: Animatable volume rendering of articulated human sdfs”. *European conference on computer vision*. Springer. 2022, 1–19 3.
- [WZS*19] WANG, XIAOGANG, ZHOU, BIN, SHI, YAHAO, et al. “Shape2Motion: Joint analysis of motion parts and attributes from 3D shapes”. *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. 2019, 8876–8884 4–7, 9–12.
- [XQM*20] XIANG, FANBO, QIN, YUZE, MO, KAICHUN, et al. “SAPIEN: A SimulATED Part-based Interactive ENvironment”. *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. 2020, 11097–11107 1, 4, 5, 10.
- [XWY*09] XU, WEIWEI, WANG, JUN, YIN, KANGKANG, et al. “Joint-aware manipulation of deformable models”. *ACM Transactions on Graphics (TOG)* 28.3 (2009), 1–9 6–9, 11, 12.
- [YHL*18] YI, LI, HUANG, HAIBIN, LIU, DIFAN, et al. “Deep part induction from articulated object pairs”. *ACM Transactions on Graphics (TOG)* 37.6 (2018), 1–15 4, 6, 7, 11, 12.
- [YHY*19] YAN, ZIHAO, HU, RUIZHEN, YAN, XINGGUANG, et al. “RPM-Net: recurrent prediction of motion and parts from point cloud”. *ACM Transactions on Graphics (TOG)* 38.6 (2019), 1–15 3–7, 10–13.
- [Y LX*16] YUAN, QING, LI, GUIQING, XU, KAI, et al. “Space-time co-segmentation of articulated point cloud sequences”. *Computer Graphics Forum*. Vol. 35. 2. 2016, 419–429 6–9.
- [YRH*24] YAO, CHUN-HAN, RAJ, AMIT, HUNG, WEI-CHIH, et al. “AR-TIC3D: Learning robust articulated 3d shapes from noisy web image collections”. *Advances in neural information processing systems (NeurIPS)* 36 (2024) 3.
- [YVN*22] YANG, GENGSHAN, VO, MINH, NEVEROVA, NATALIA, et al. “Banmo: Building animatable 3d neural models from many casual videos”. *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. 2022, 2863–2873 1.
- [YZW*22] YANG, JI, ZUO, XINXIN, WANG, SEN, et al. “Object Wake-up: 3D Object Rigging from a Single Image”. *Proceedings of the European Conference on Computer Vision (ECCV)*. Springer. 2022, 311–327 1.
- [ZLLK21] ZENG, VICKY, LEE, TABITHA EDITH, LIANG, JACKY, and KROEMER, OLIVER. “Visual identification of articulated object parts”. *IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. 2021, 2443–2450 10.