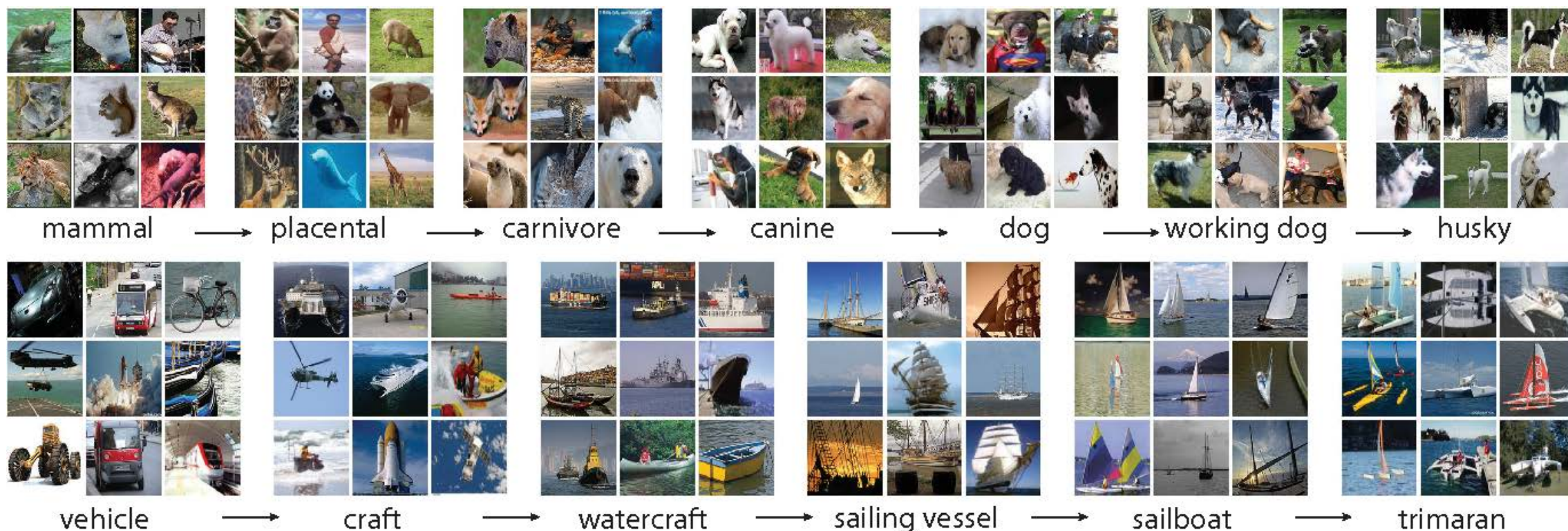


Learning to Generate 3D Structures

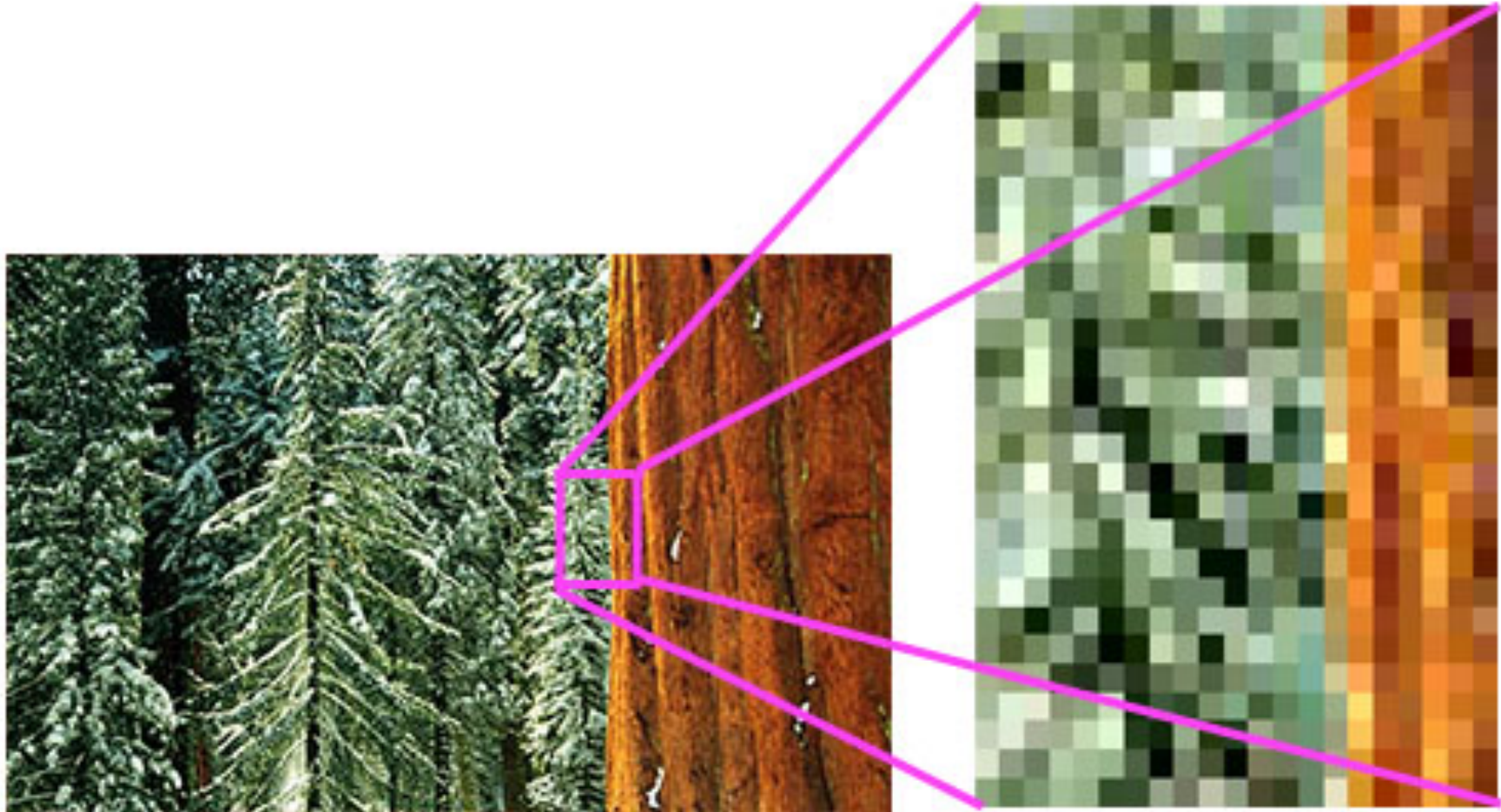
Siddhartha Chaudhuri
Adobe Research and IIT Bombay



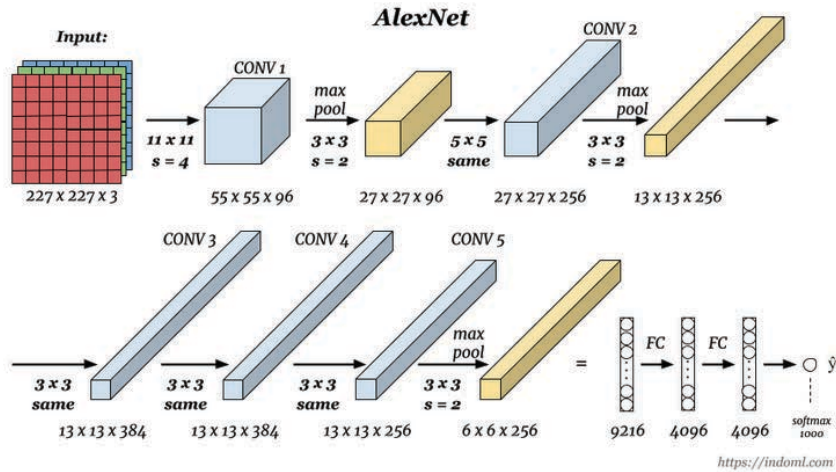
Visual machine learning is traditionally on 2D images



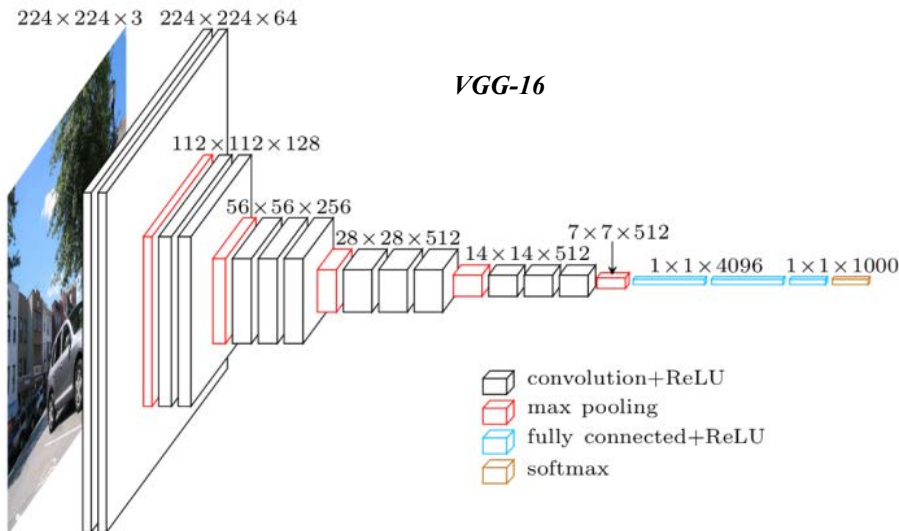
Images are 2D grids of pixels



2D deep networks do *convolution* on grids



- Convolution measures “weighted overlap” of f with another function g as it is (reversed and) shifted over f



2	3	1
0	5	1
1	0	8

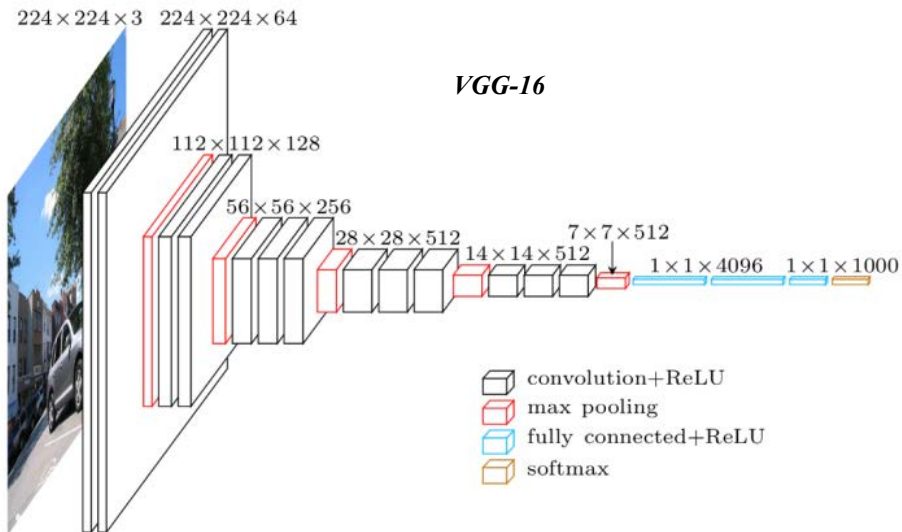
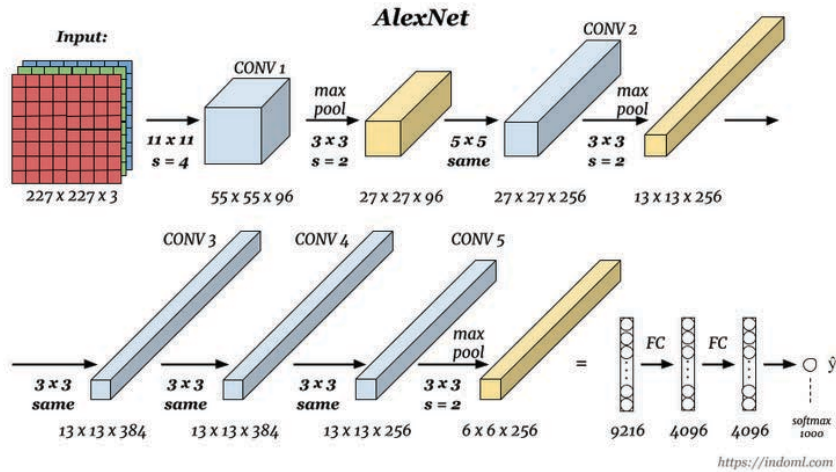
 $*$

0	-1	0
-1	5	-1
0	-1	0

 $= ?$

(we'll assume everything outside the 3×3 is zero)

2D deep networks do *convolution* on grids

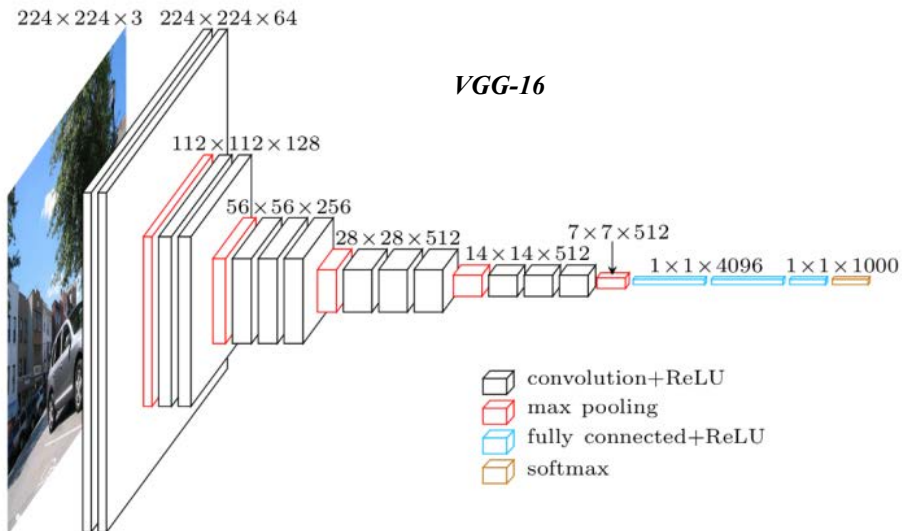
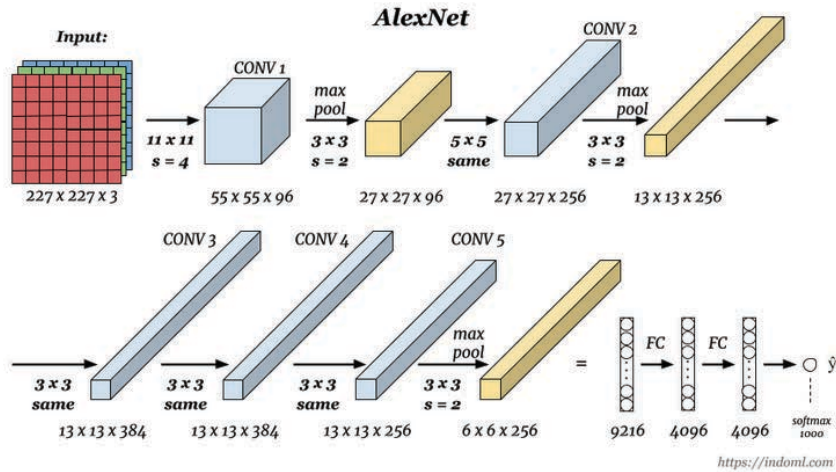


- Convolution measures “weighted overlap” of f with another function g as it is (reversed and) shifted over f

0	-1	0	
-1	25	31	1
0	01	50	1
	1	0	8

7		

2D deep networks do *convolution* on grids

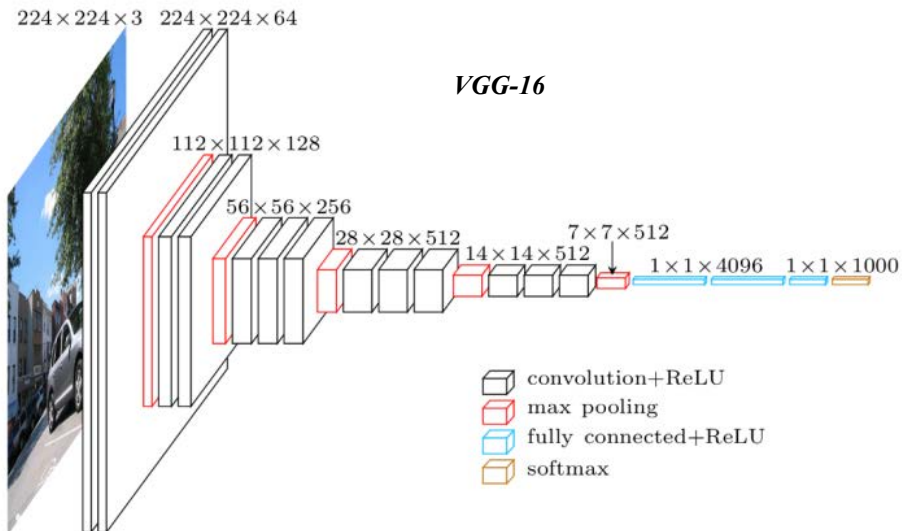
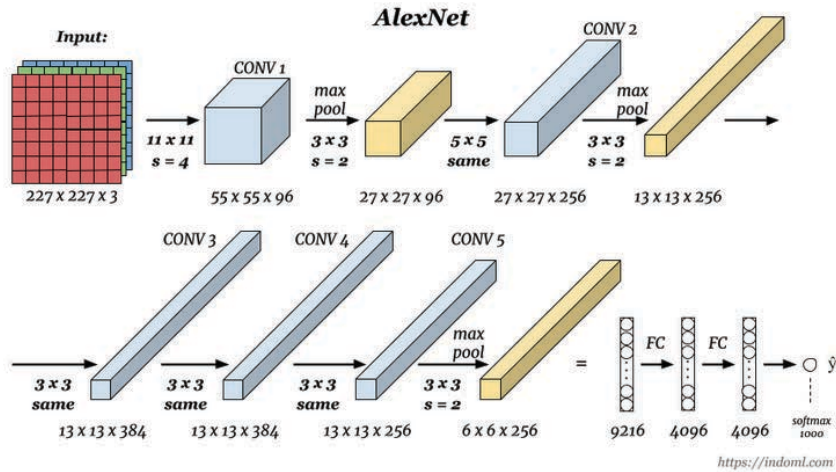


- Convolution measures “weighted overlap” of f with another function g as it is (reversed and) shifted over f

0	-1	0
21	35	11
00	51	10
1	0	8

7	7	

2D deep networks do *convolution* on grids

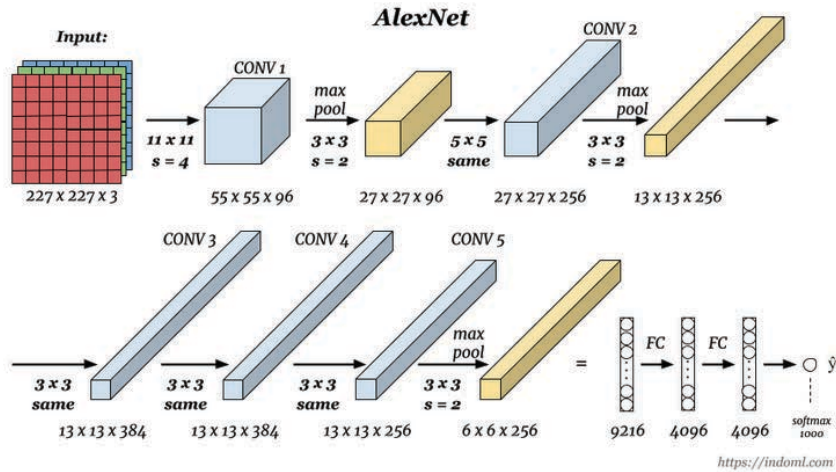


- Convolution measures “weighted overlap” of f with another function g as it is (reversed and) shifted over f

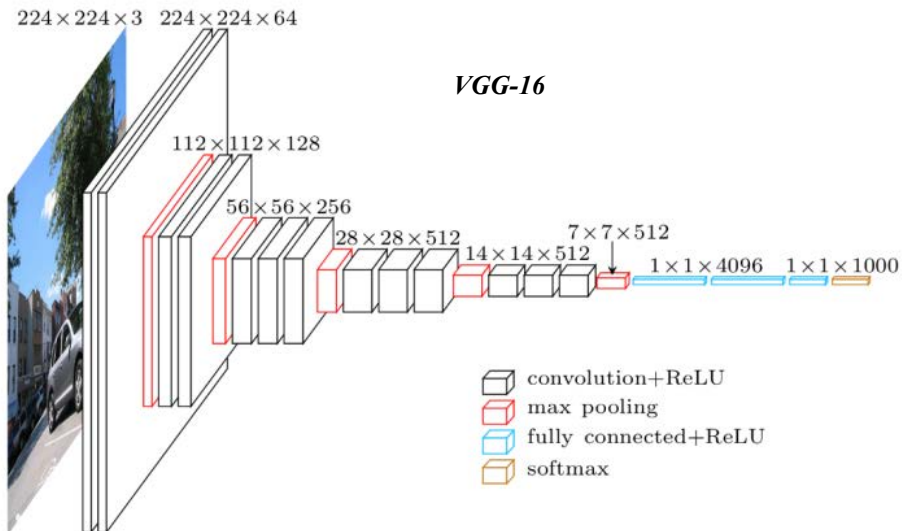
		0	-1	0
2	3	1	5	-1
0	5	0	1	0
1	0	8		

7	7	1

2D deep networks do *convolution* on grids



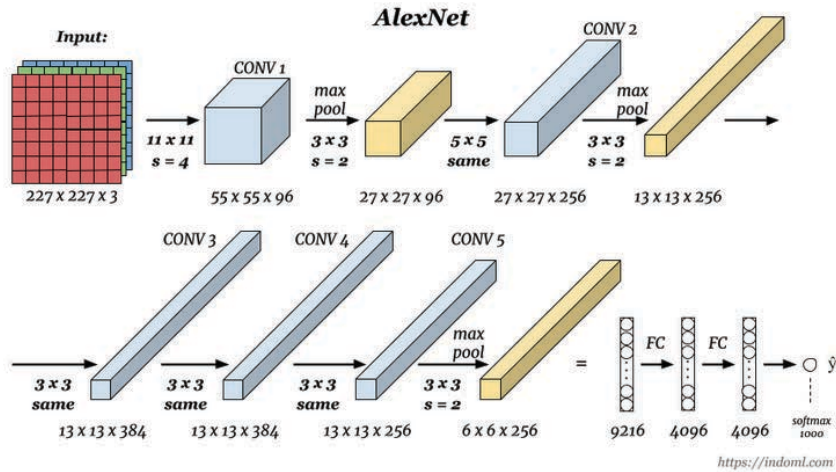
- Convolution measures “weighted overlap” of f with another function g as it is (reversed and) shifted over f



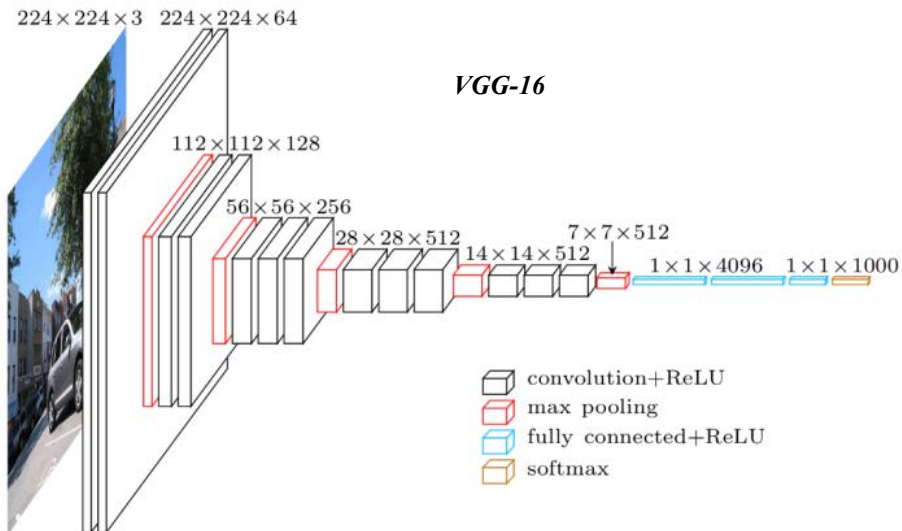
0	21	30	1
-1	05	51	1
0	11	00	8

7	7	1
-8		

2D deep networks do *convolution* on grids



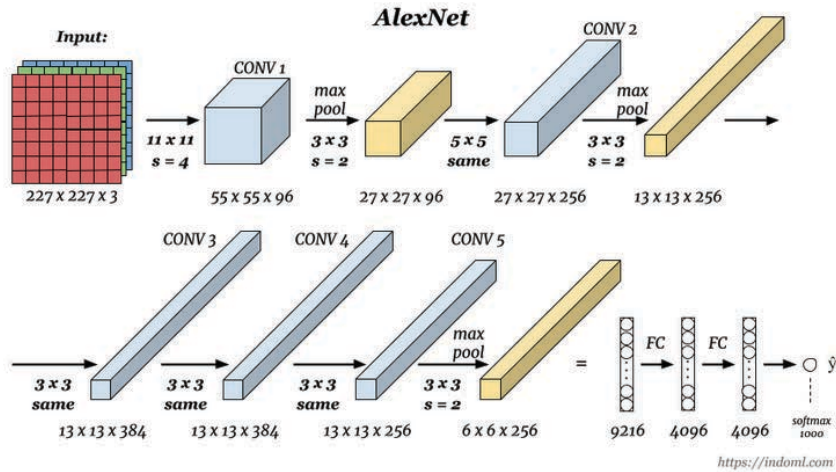
- Convolution measures “weighted overlap” of f with another function g as it is (reversed and) shifted over f



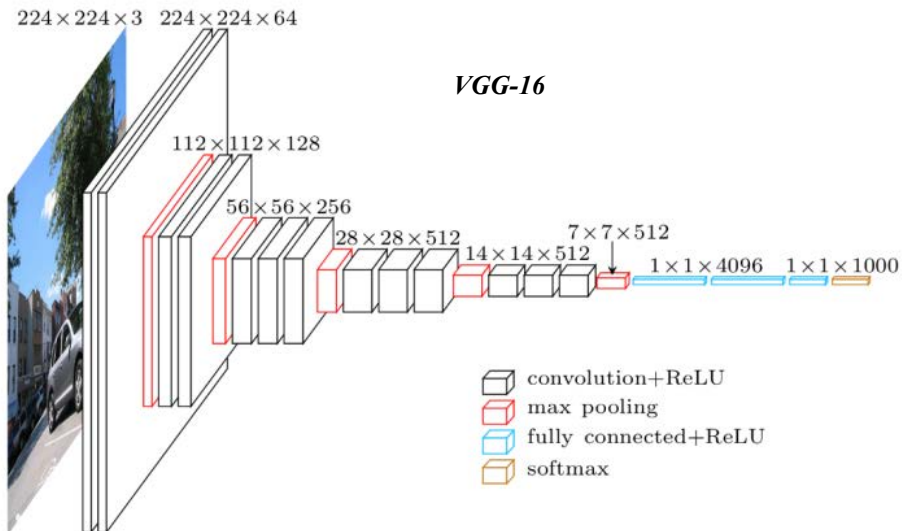
20	31	10
01	55	11
10	01	80

7	7	1
-8	21	

2D deep networks do *convolution* on grids



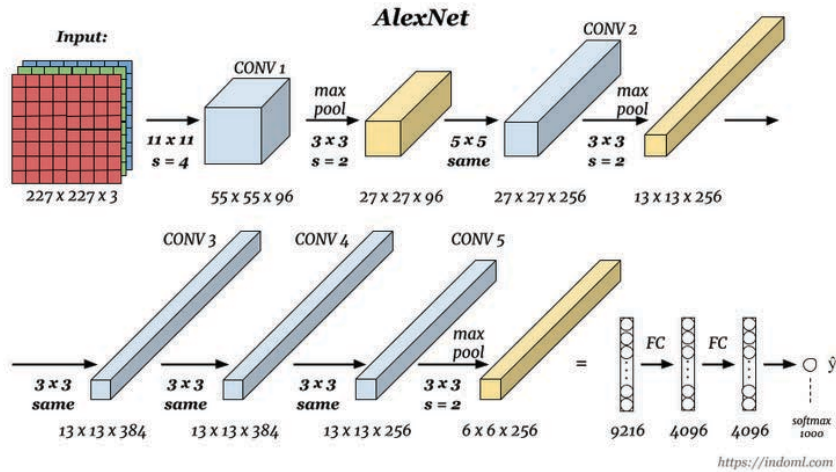
- Convolution measures “weighted overlap” of f with another function g as it is (reversed and) shifted over f



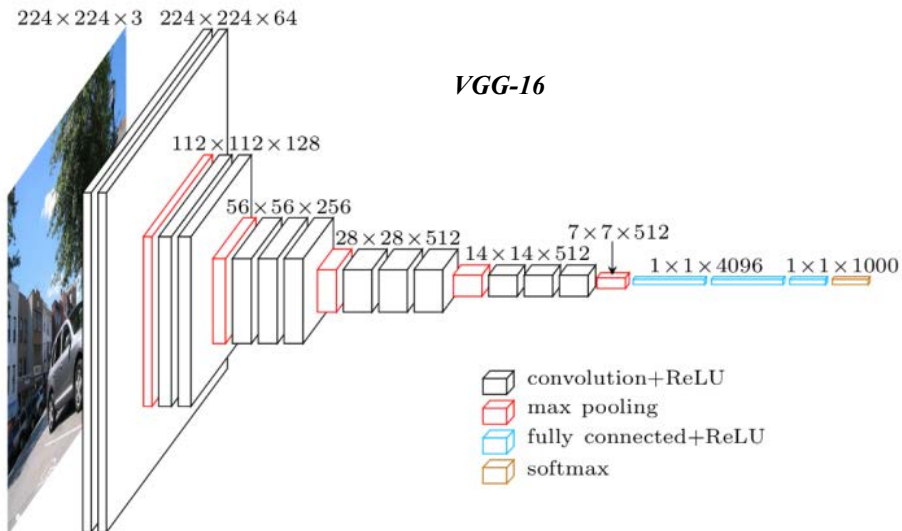
2	30	11	0
0	51	15	-1
1	00	81	0

7	7	1
-8	21	-9

2D deep networks do *convolution* on grids



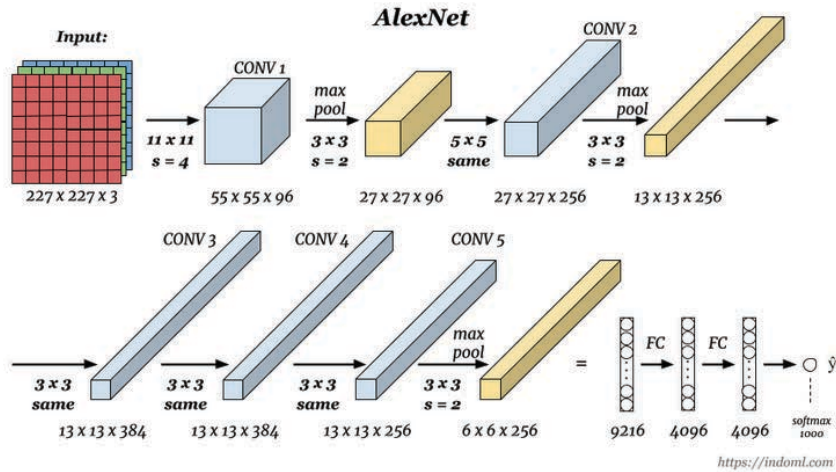
- Convolution measures “weighted overlap” of f with another function g as it is (reversed and) shifted over f



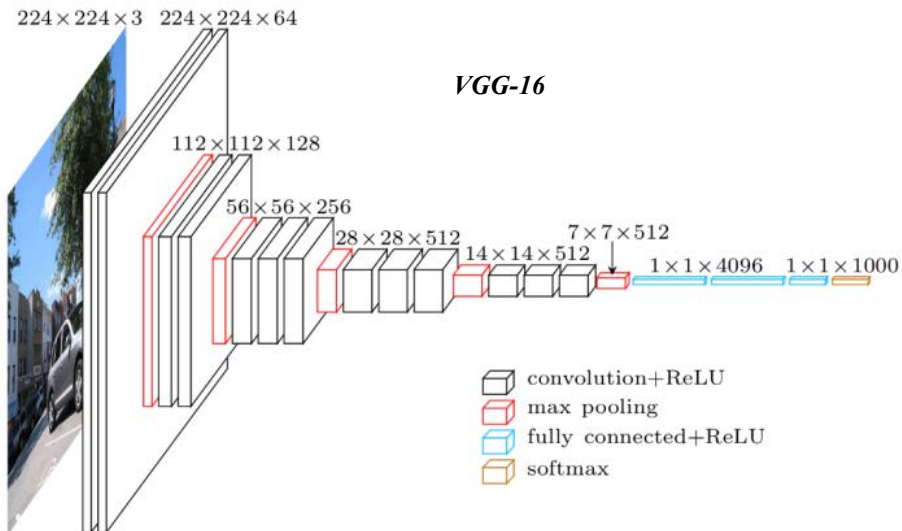
	2	3	1
0	0.1	5.0	1
-1	1.5	0.1	8
0	-1	0	

7	7	1
-8	21	-9
5		

2D deep networks do *convolution* on grids



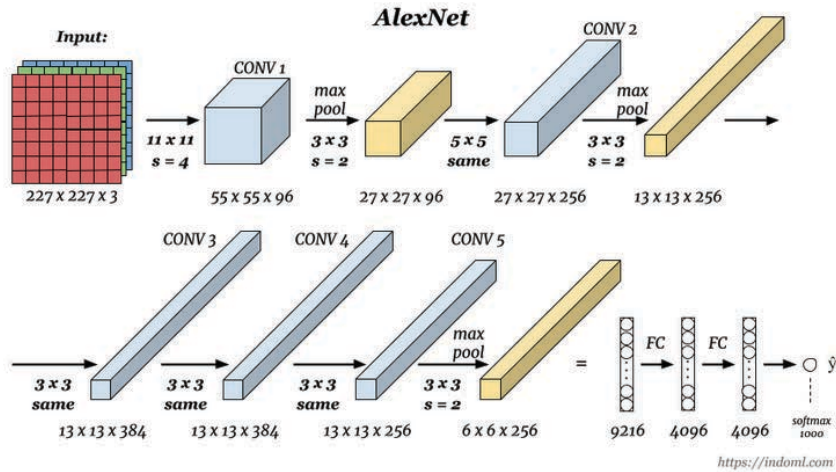
- Convolution measures “weighted overlap” of f with another function g as it is (reversed and) shifted over f



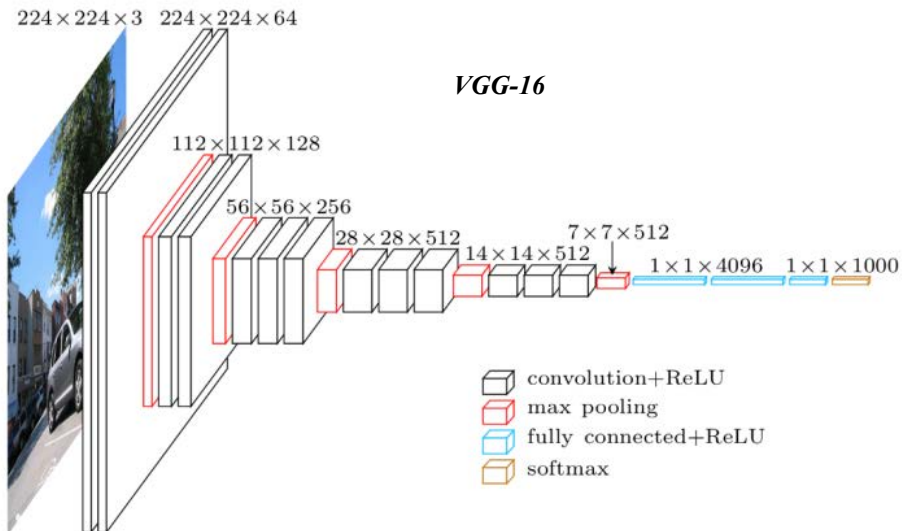
2	3	1
00	51	10
11	05	81
0	-1	0

7	7	1
-8	21	-9
5	-14	

2D deep networks do *convolution* on grids



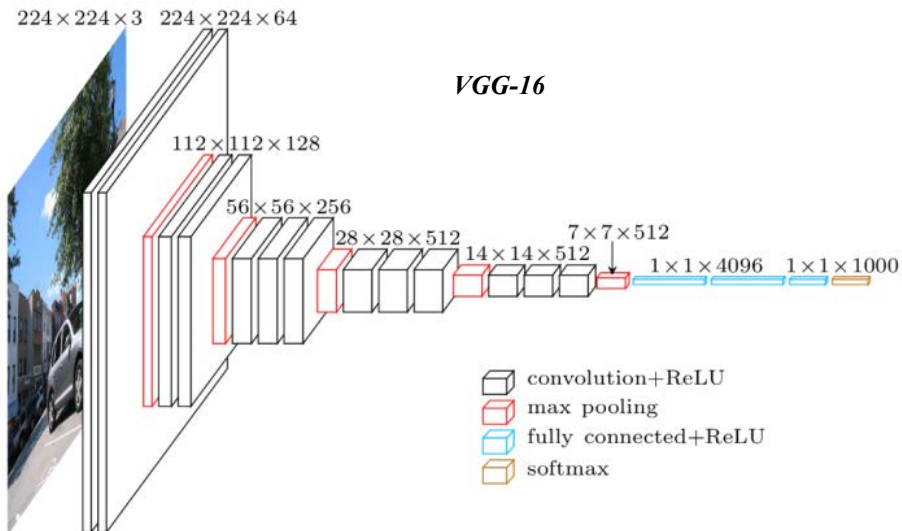
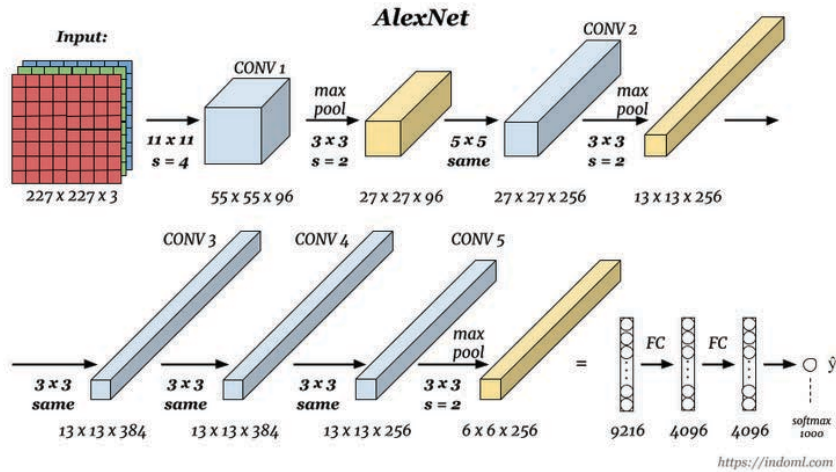
- Convolution measures “weighted overlap” of f with another function g as it is (reversed and) shifted over f



2	3	1	
0	50	11	0
1	01	85	-1
	0	-1	0

7	7	1
-8	21	-9
5	-14	39

2D deep networks do *convolution* on grids



- Convolution measures “weighted overlap” of f with another function g as it is (reversed and) shifted over f

2	3	1
0	5	1
1	0	8

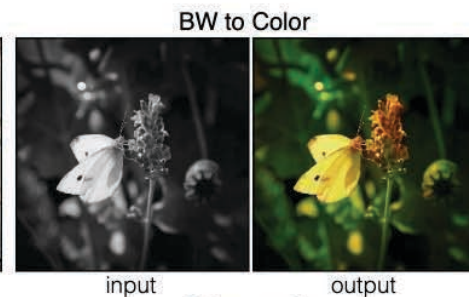
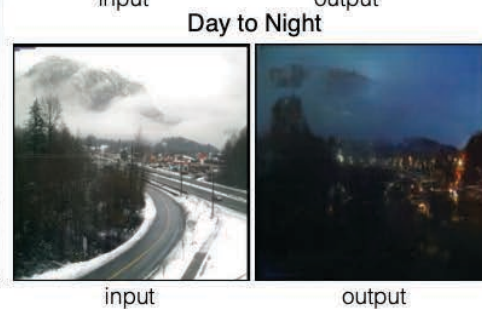
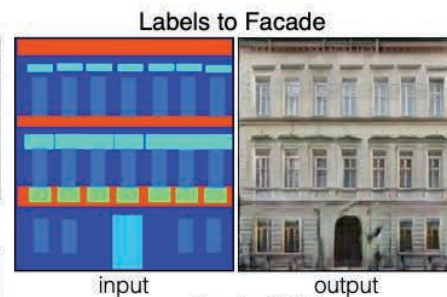
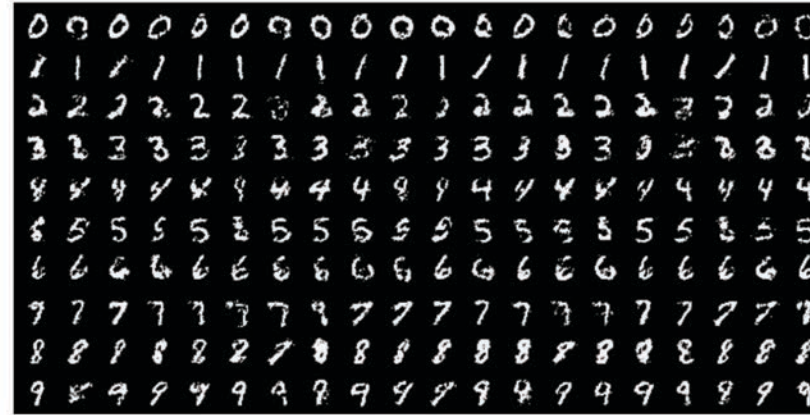
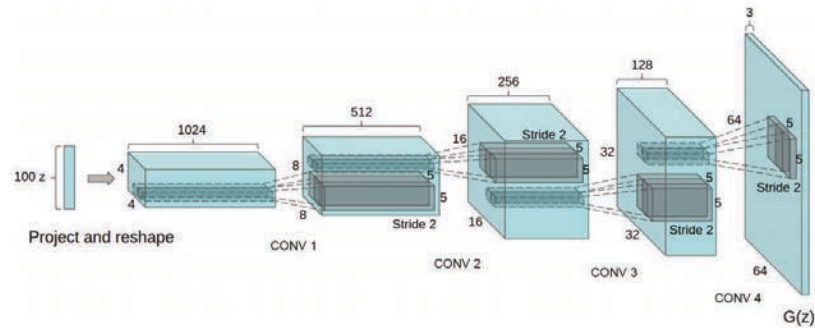
 $*$

0	-1	0
-1	5	-1
0	-1	0

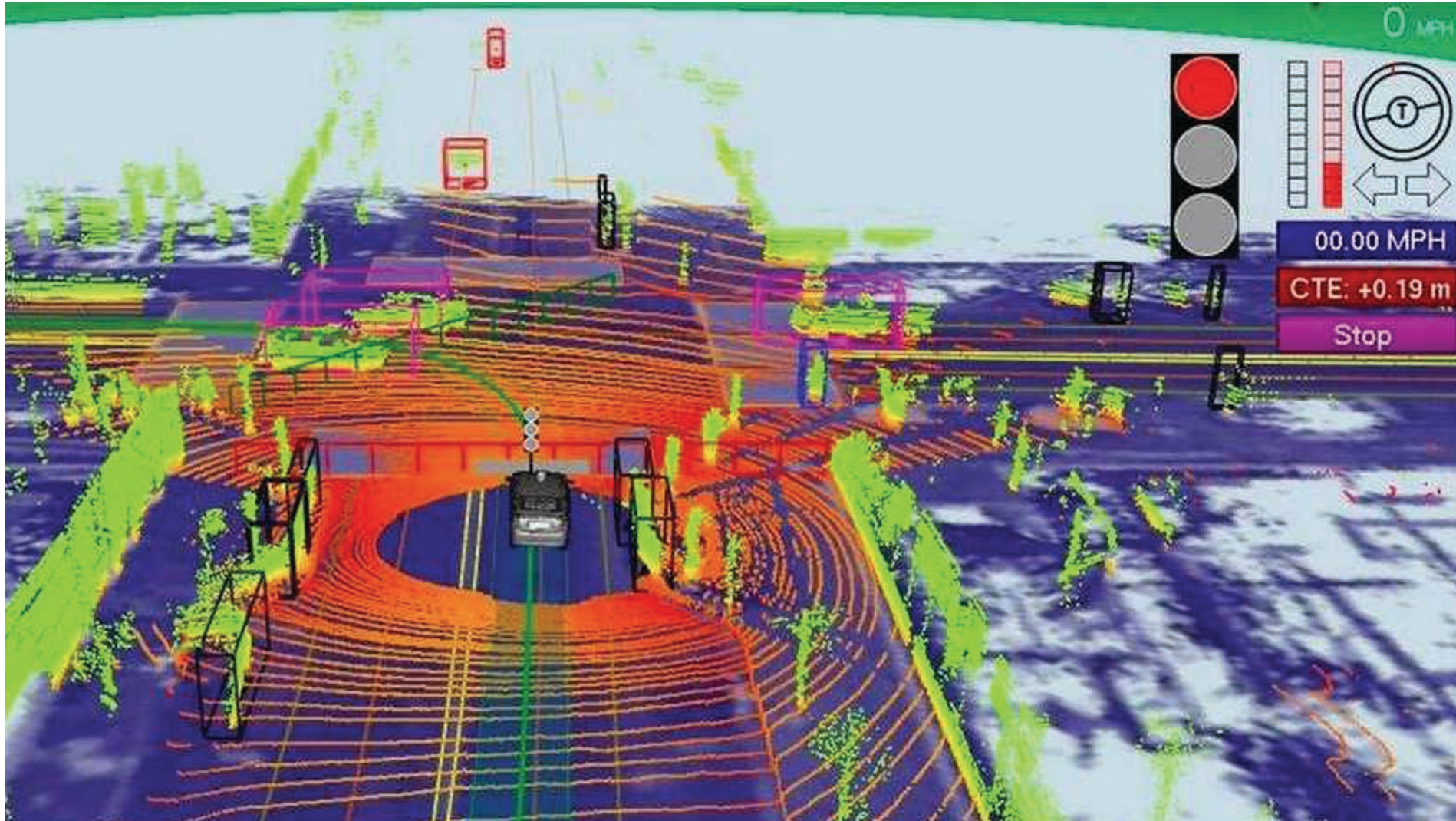
 $=$

7	7	1
-8	21	-9
5	-14	39

2D generative models use 2D convolution



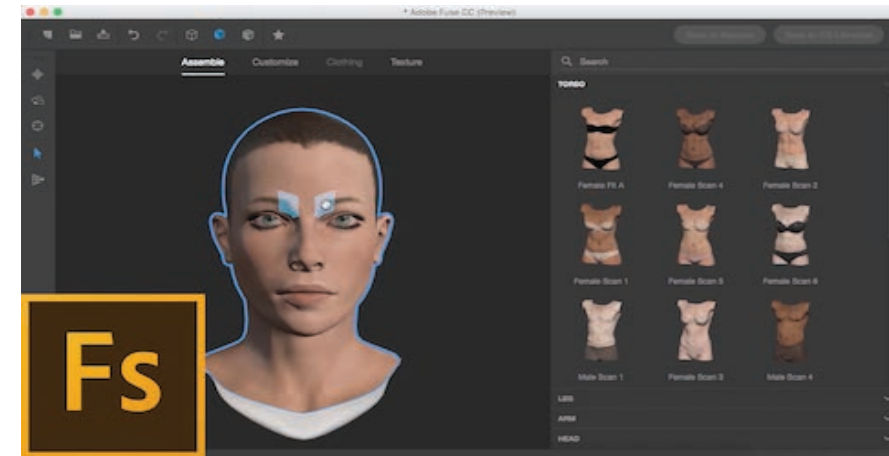
But the real world is (at least) 3D



Synthesizing 3D shapes is widely important



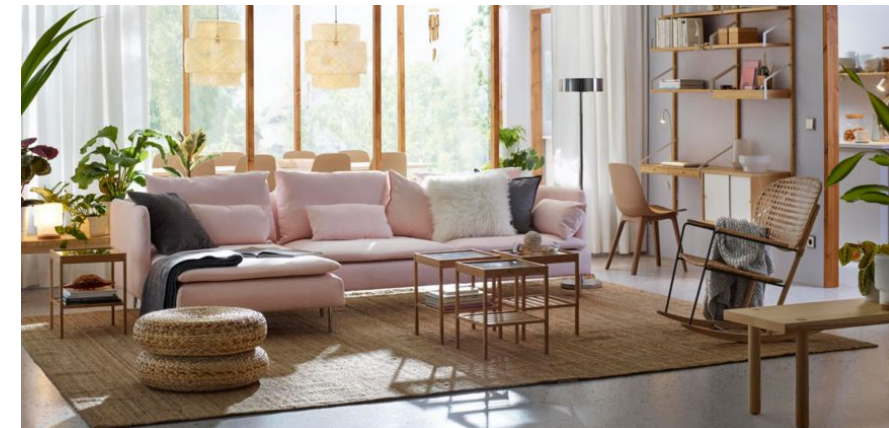
CAD/CAM



3D modeling for AR/VR/entertainment



Fabrication

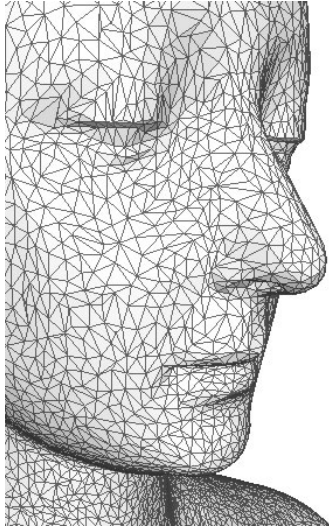


Architecture and interior design

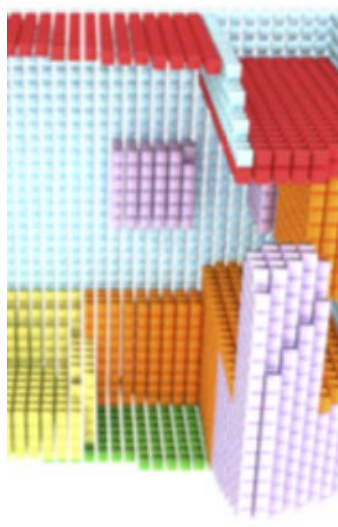
Many ways to represent 3D shapes



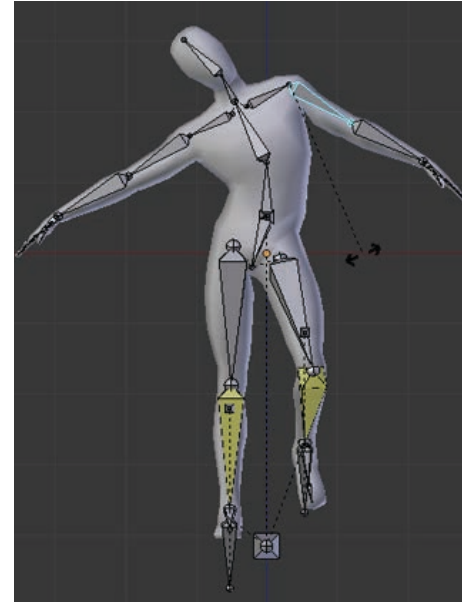
Point cloud



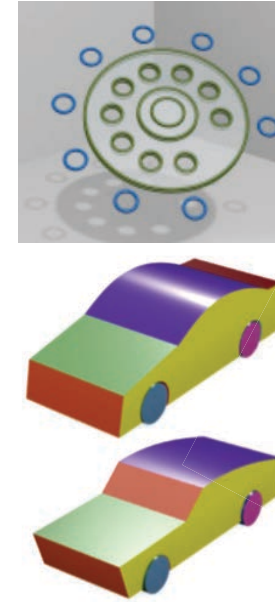
Polygon mesh



Voxel grid



Rigged skeleton



Wires and patches



Part and object assembly

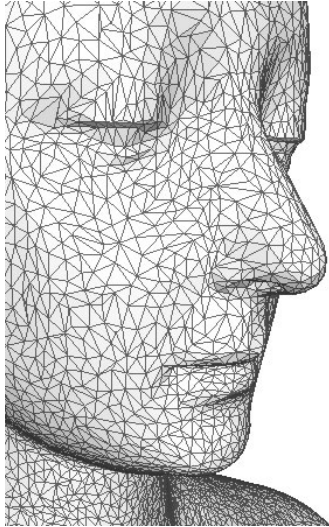


And 2D convolution doesn't extend directly to most of them 😞
(no grid structure with common parametrization)

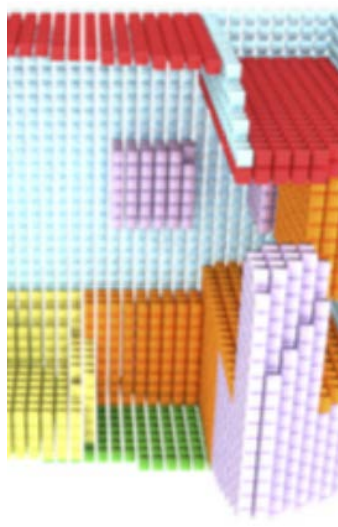
Some representations are *structure-aware*



Point cloud

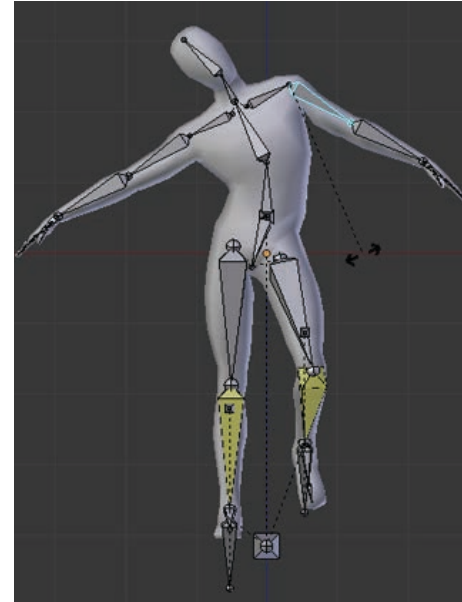


Polygon mesh

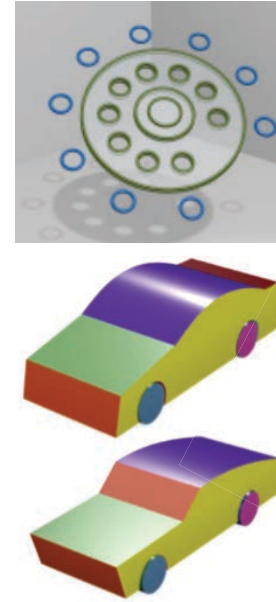


Voxel grid

Not
structure-aware



Rigged skeleton



Wires and
patches

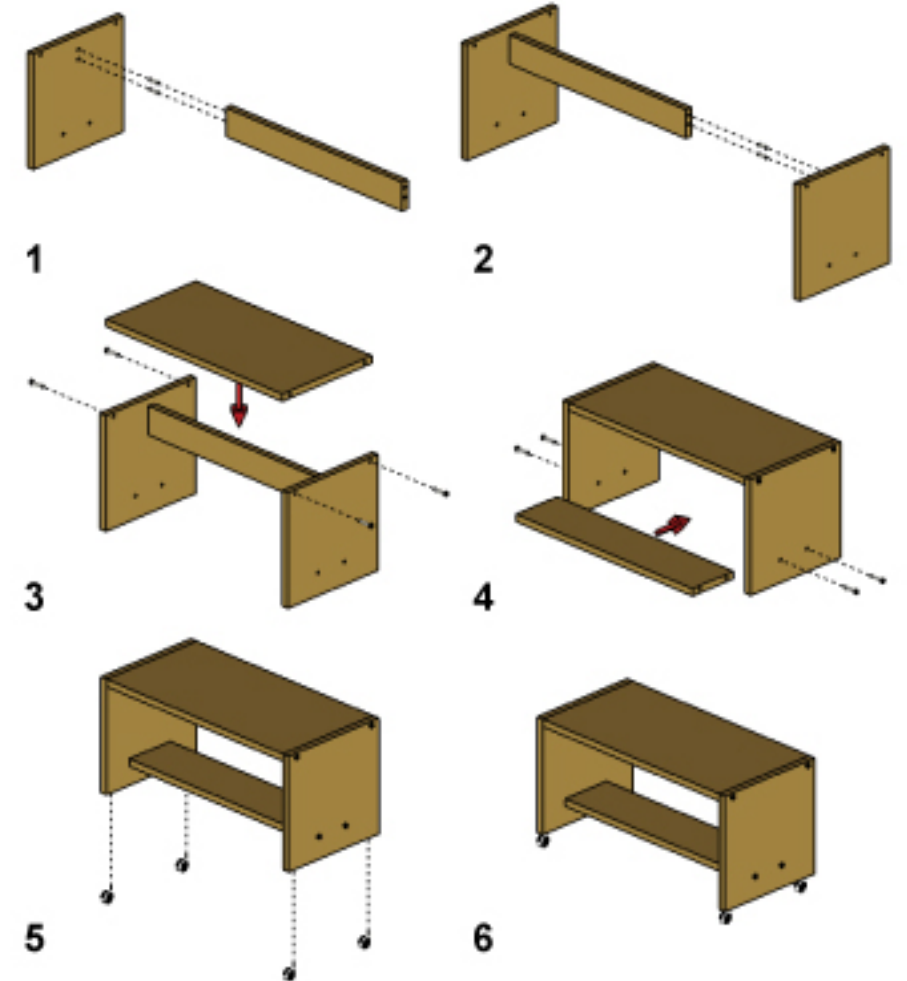


Part and object
assembly

Structure-aware

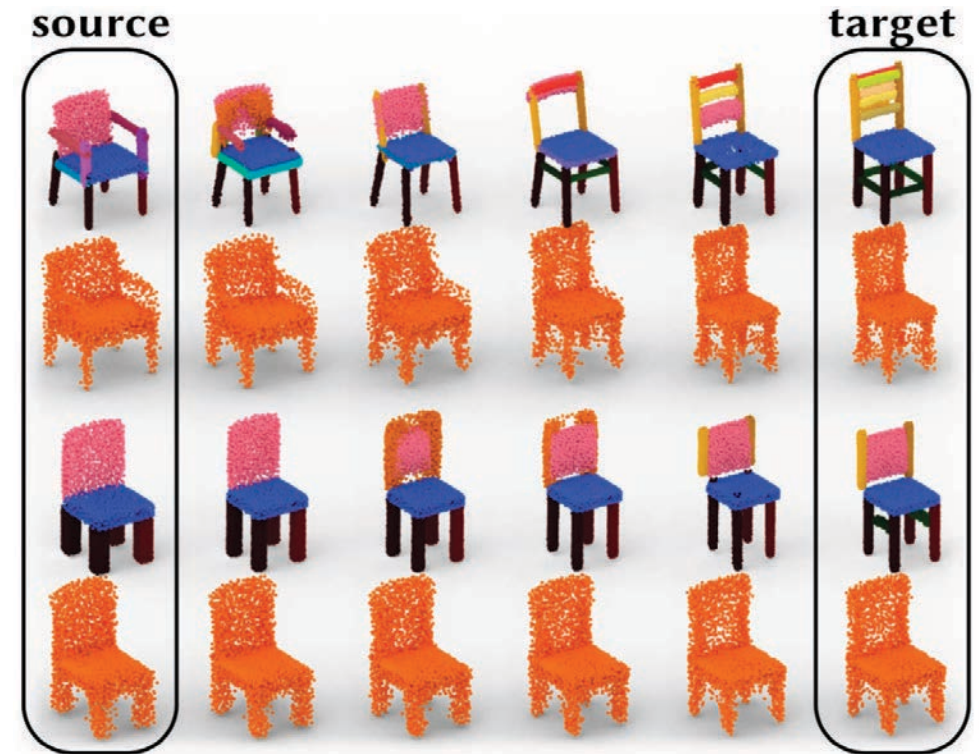
Shape structure: *what* and *why*?

- **What:** Low-dimensional semantic/functional abstraction (parts, layouts, deformations...)



Shape structure: *what* and *why*?

- **What:** Low-dimensional semantic/functional abstraction (parts, layouts, deformations...)
- **Why** is structure useful for shape generation?
 - Low-dimensional: less training data, more captured variability
 - Explicit factorization assists learning, model architectures can be decoupled
 - High-quality output: e.g. synthesize novel structure, reuse fine-grained part geometry
 - Output pre-rigged for high-level modeling

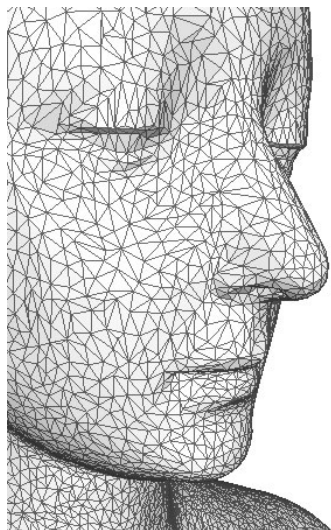


VAE latent space interpolation, **with** and **without** structure-aware model

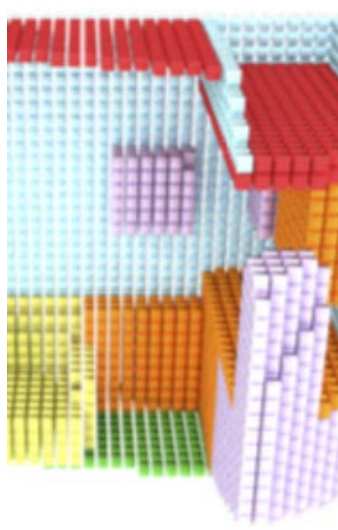
Many ways to generate 3D shapes



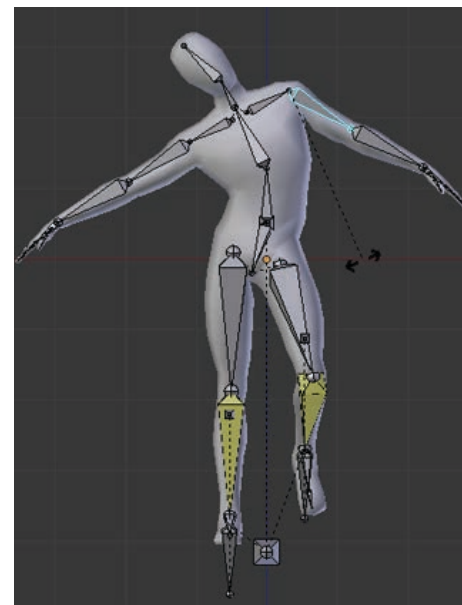
Point cloud



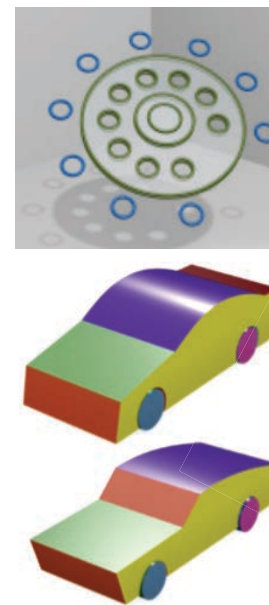
Polygon mesh



Voxel grid



Rigged skeleton



Wires and patches



Part and object assembly

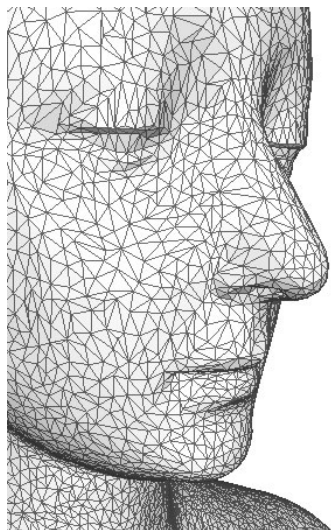


Achlioptas et al. 2017

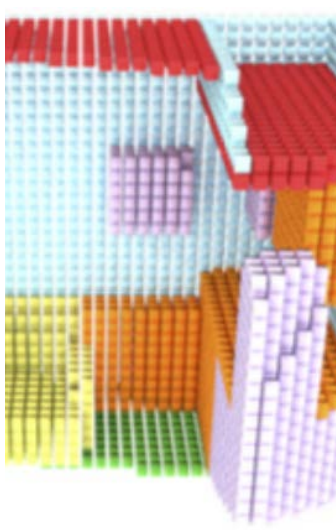
Many ways to generate 3D shapes



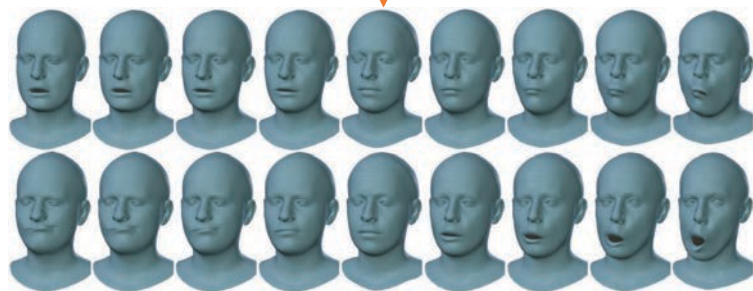
Point cloud



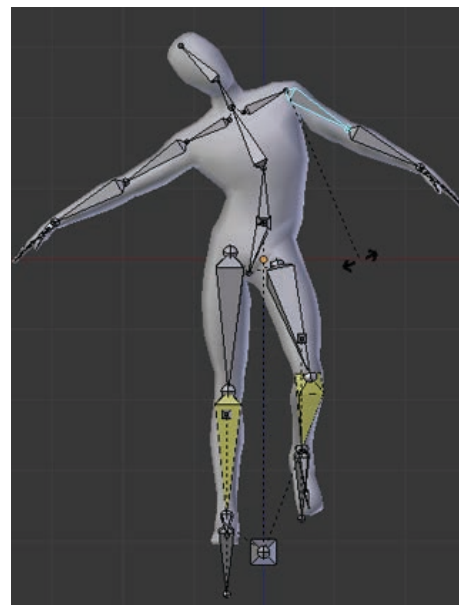
Polygon mesh



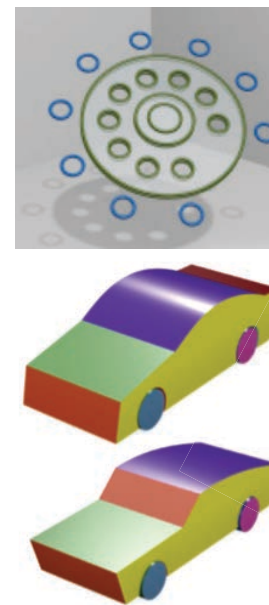
Voxel grid



Ranjan et al. 2018



Rigged skeleton



Wires and patches

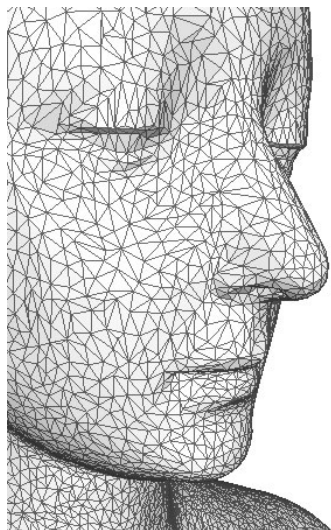


Part and object assembly

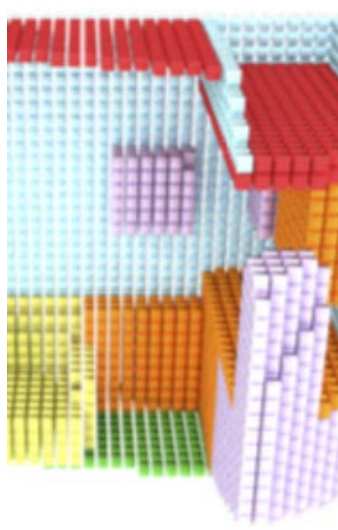
Many ways to generate 3D shapes



Point cloud



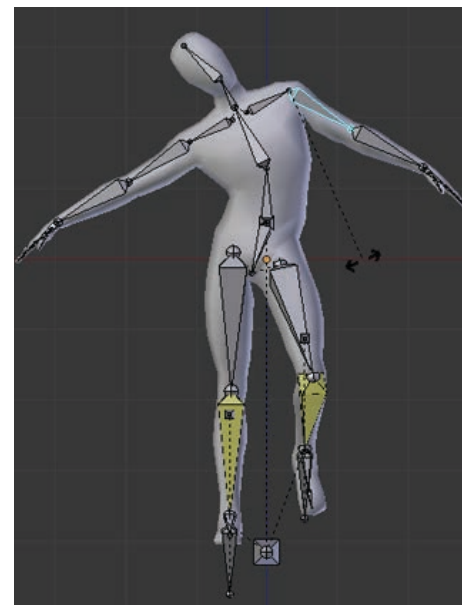
Polygon mesh



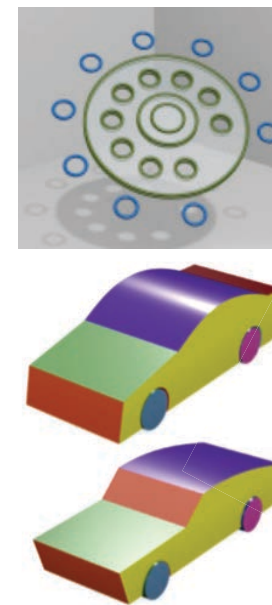
Voxel grid



Wu et al. 2015



Rigged skeleton



Wires and patches

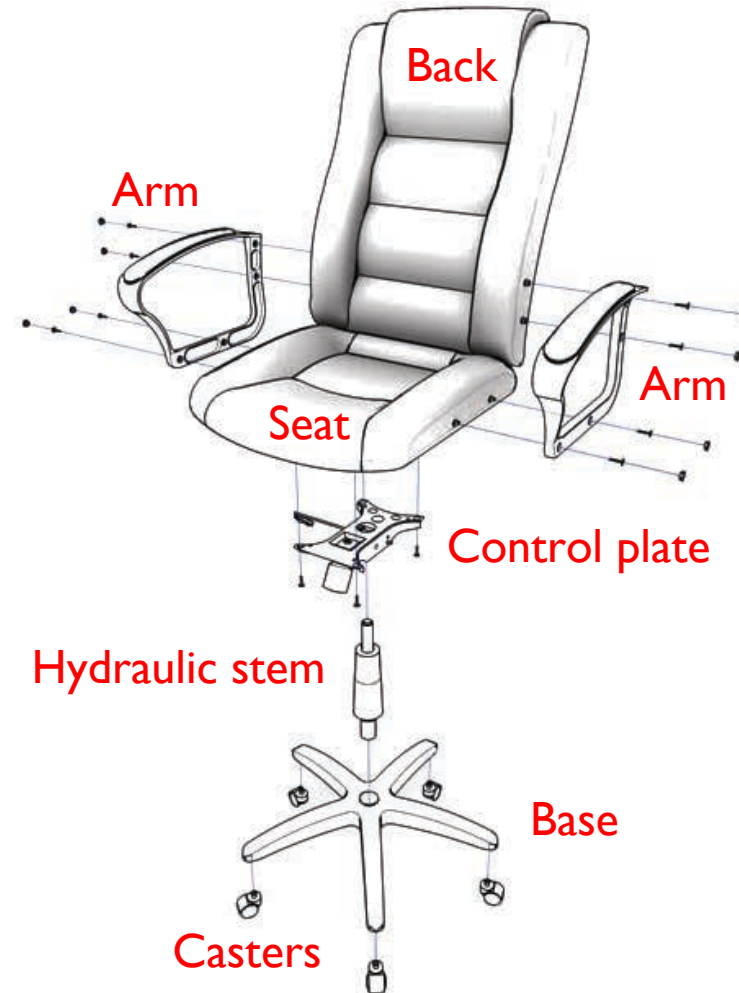


Part and object assembly

Focus of this talk

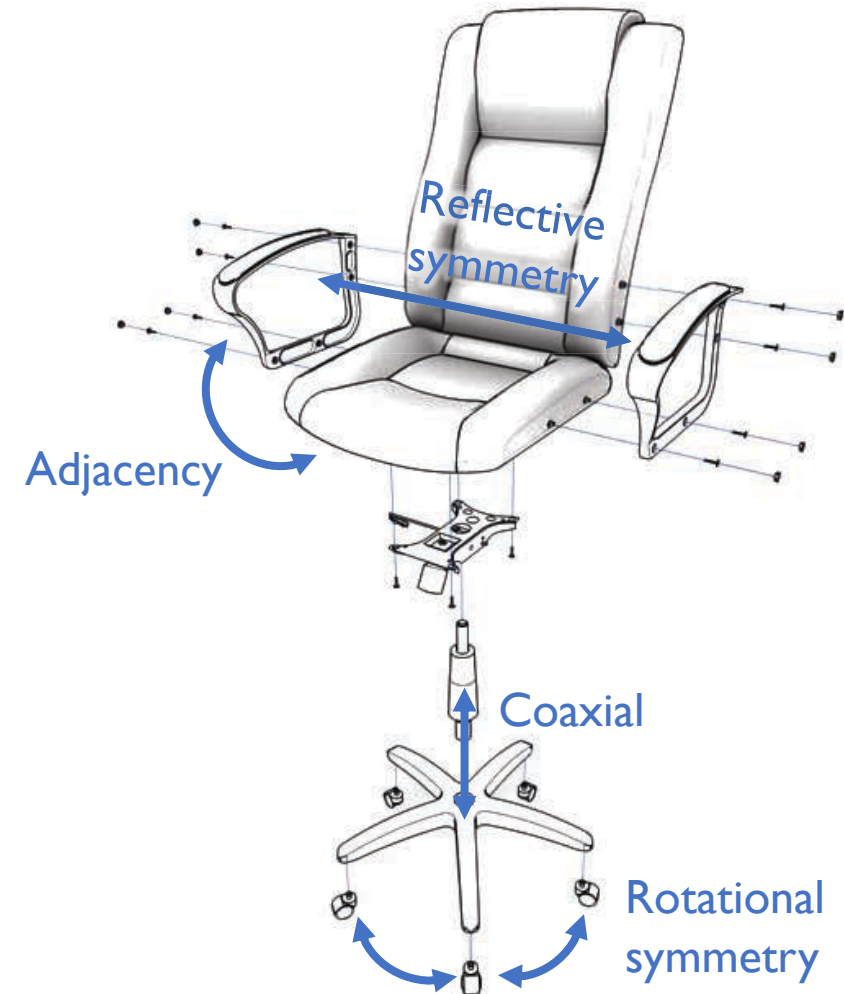
Structure is commonly described by (hyper-)graphs

- Parts (*vertices*)
 - Discrete and continuous variation



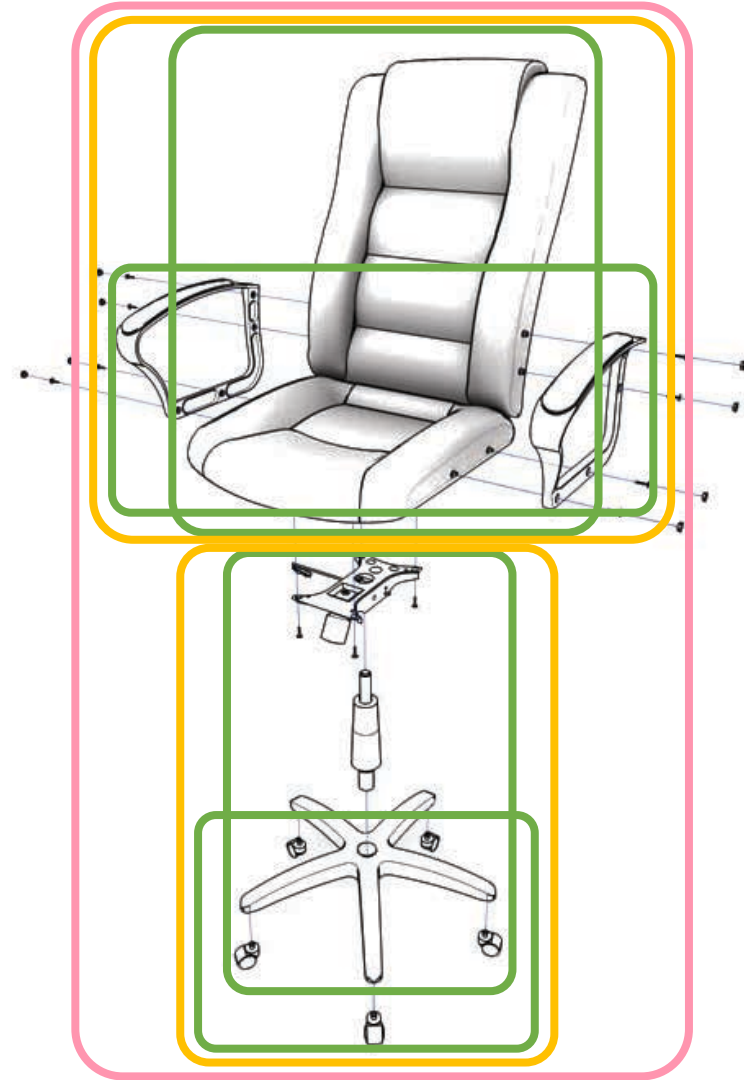
Structure is commonly described by (hyper-)graphs

- Parts (*vertices*)
 - Discrete and continuous variation
- Layout relationships (*edges*)



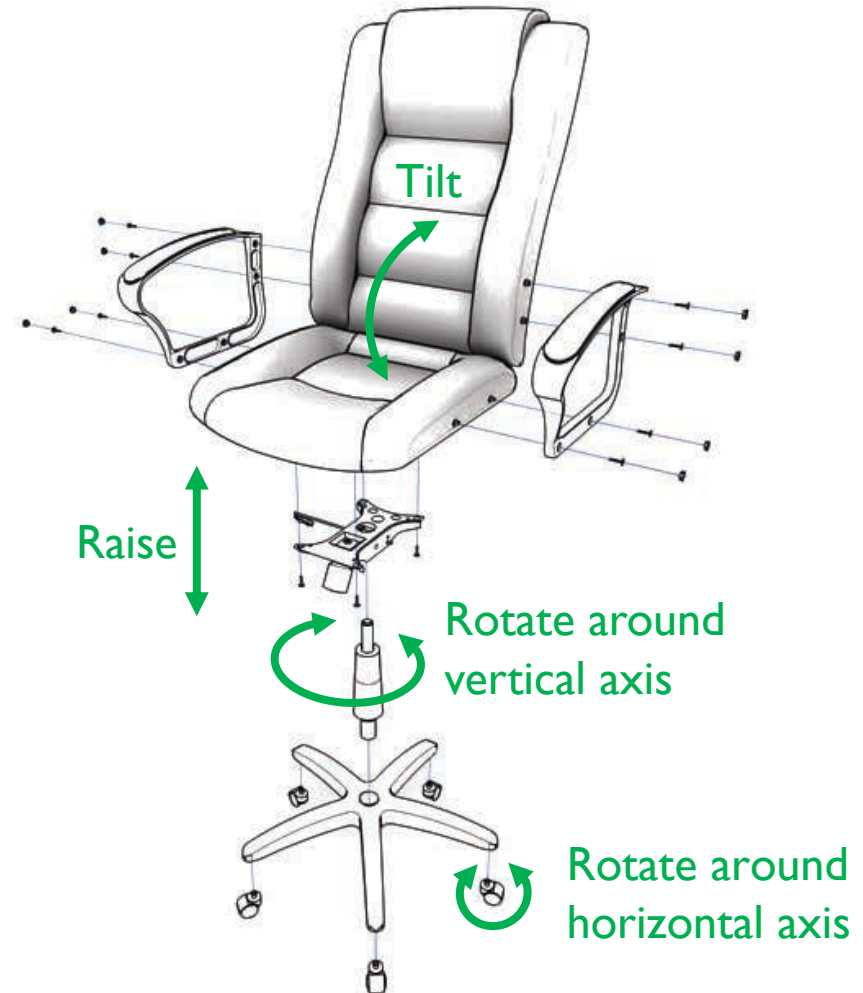
Structure is commonly described by (hyper-)graphs

- Parts (*vertices*)
 - Discrete and continuous variation
- Layout relationships (*edges*)
 - Hierarchical



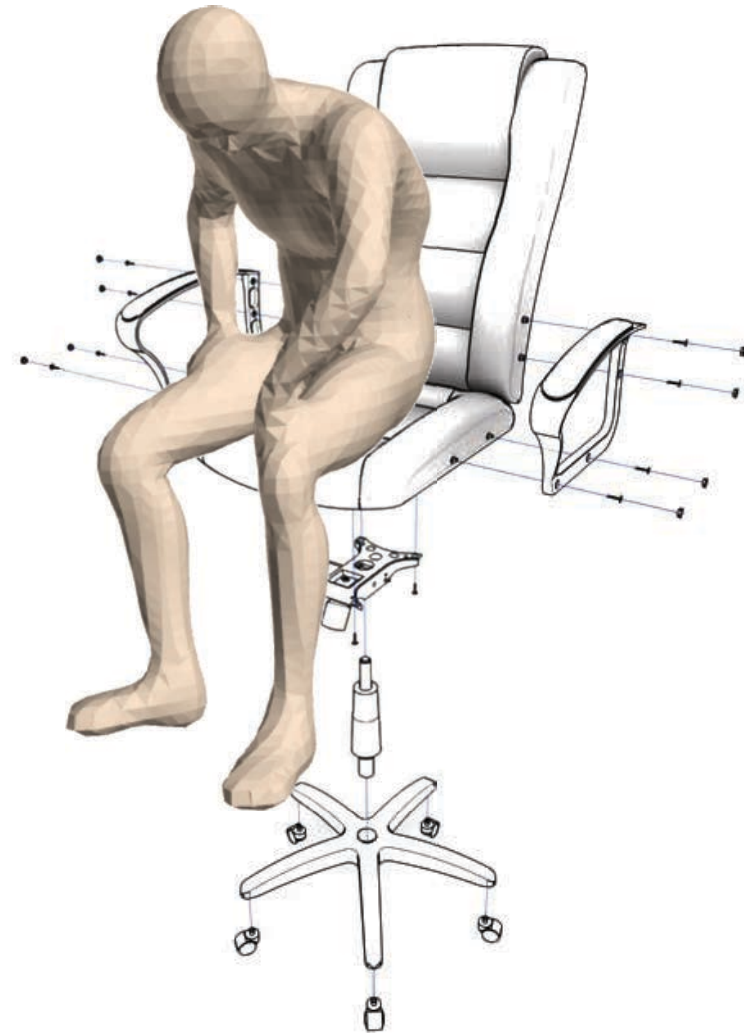
Structure is commonly described by (hyper-)graphs

- Parts (*vertices*)
 - Discrete and continuous variation
- Layout relationships (*edges*)
 - Hierarchical
- Motion



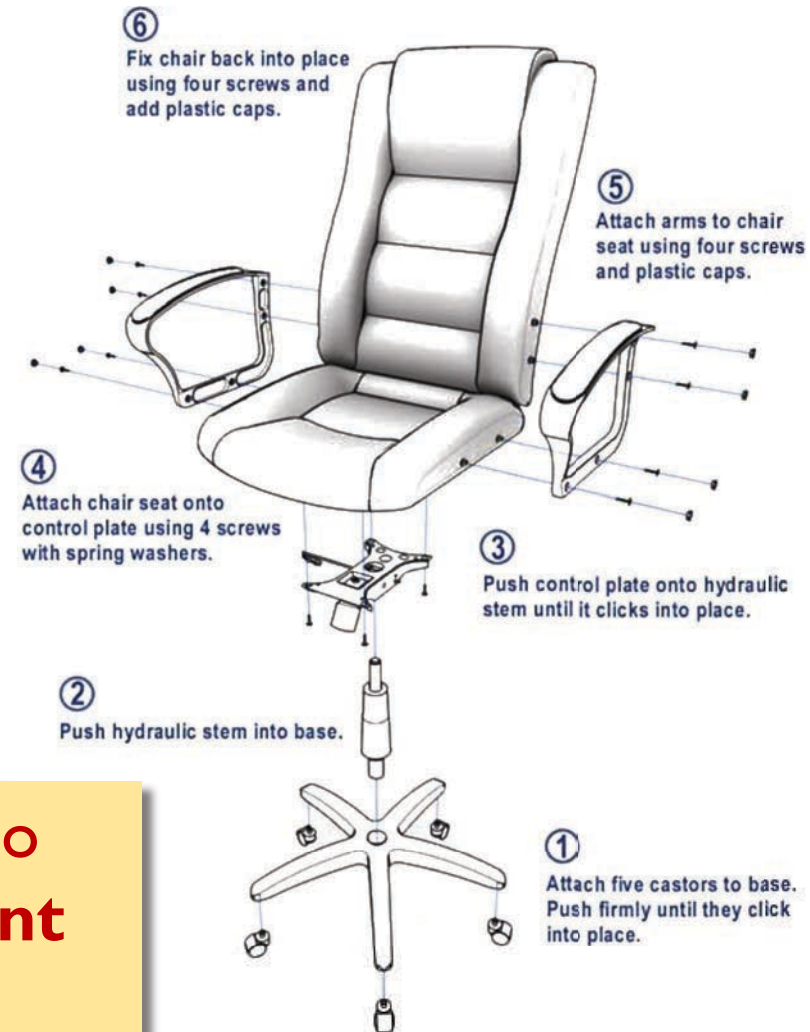
Structure is commonly described by (hyper-)graphs

- Parts (*vertices*)
 - Discrete and continuous variation
 - Layout relationships (*edges*)
 - Hierarchical
 - Motion
 - Function
- } *attributes*



Structure is commonly described by (hyper-)graphs

- Parts (*vertices*)
 - Discrete and continuous variation
 - Layout relationships (*edges*)
 - Hierarchical
 - Motion
 - Function
 - Text and other annotations
- } *attributes*

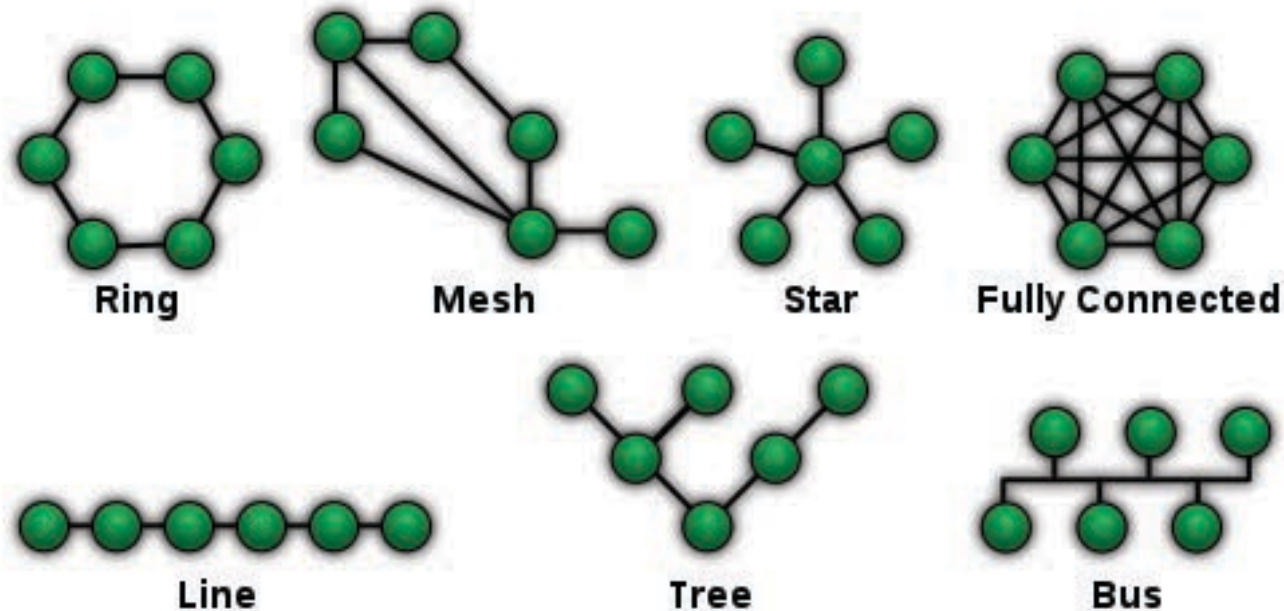


The web of interdependencies is typically too complex to specify by hand. Hence, it is **learnt** from a repository of exemplar structures

Learning generative models of **shape structure**



Learning generative models over **spaces of graphs**



- Each point in the space is a graph
- Can't assume fixed topology

What is a **generative** model?

- Let \mathbf{x} be observed data (e.g. shape) and \mathbf{y} be target data (e.g. label)
- The distribution $P(\mathbf{y} \mid \mathbf{x})$ is a **discriminative model**
 - No way to sample \mathbf{x} itself
 - Great for vision applications where we need to predict class etc (\mathbf{y}) for an observed object (\mathbf{x})
- The distribution $P(\mathbf{x}, \mathbf{y})$, or just $P(\mathbf{x})$, is a **generative model**
 - Samplers like MCMC or GAN/VAE can sample \mathbf{x} according to the distribution
 - Great for graphics applications where we need to **synthesize new objects!**
 - $P(\mathbf{x} \mid \mathbf{y})$ is a *conditional* generative model, e.g. \mathbf{y} is a sketch and \mathbf{x} is the corresponding 3D shape

Focus of this talk!

Probabilistic graphical model

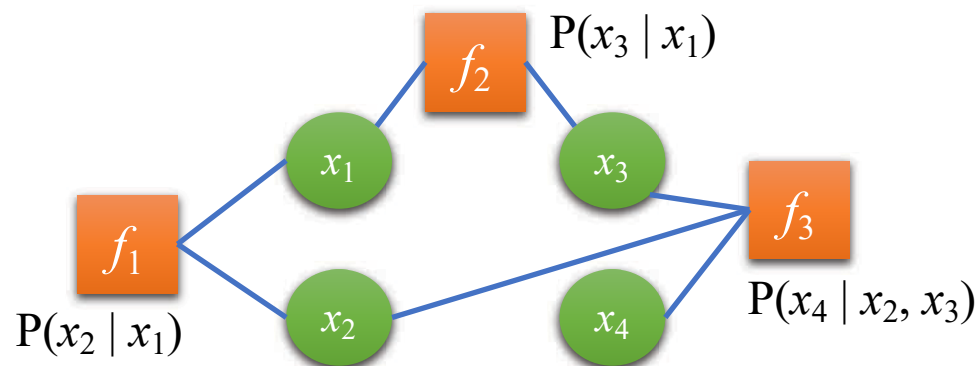
Factor graph, Bayesian network, Markov random field...

- If \mathbf{x} is high-dimensional, $P(\mathbf{x})$ can be very complicated
- Often, it can be written as a product of smaller, independent **factors**

$$P(\mathbf{x}) = P(\mathbf{x}_1 \mid \dot{\mathbf{x}}_1) \times P(\mathbf{x}_2 \mid \dot{\mathbf{x}}_2) \times \dots \times P(\mathbf{x}_k \mid \dot{\mathbf{x}}_k)$$

where $\mathbf{x}_1 \dots \mathbf{x}_k$ and (potentially empty) $\dot{\mathbf{x}}_1 \dots \dot{\mathbf{x}}_k$ are subsets of the variables of \mathbf{x}

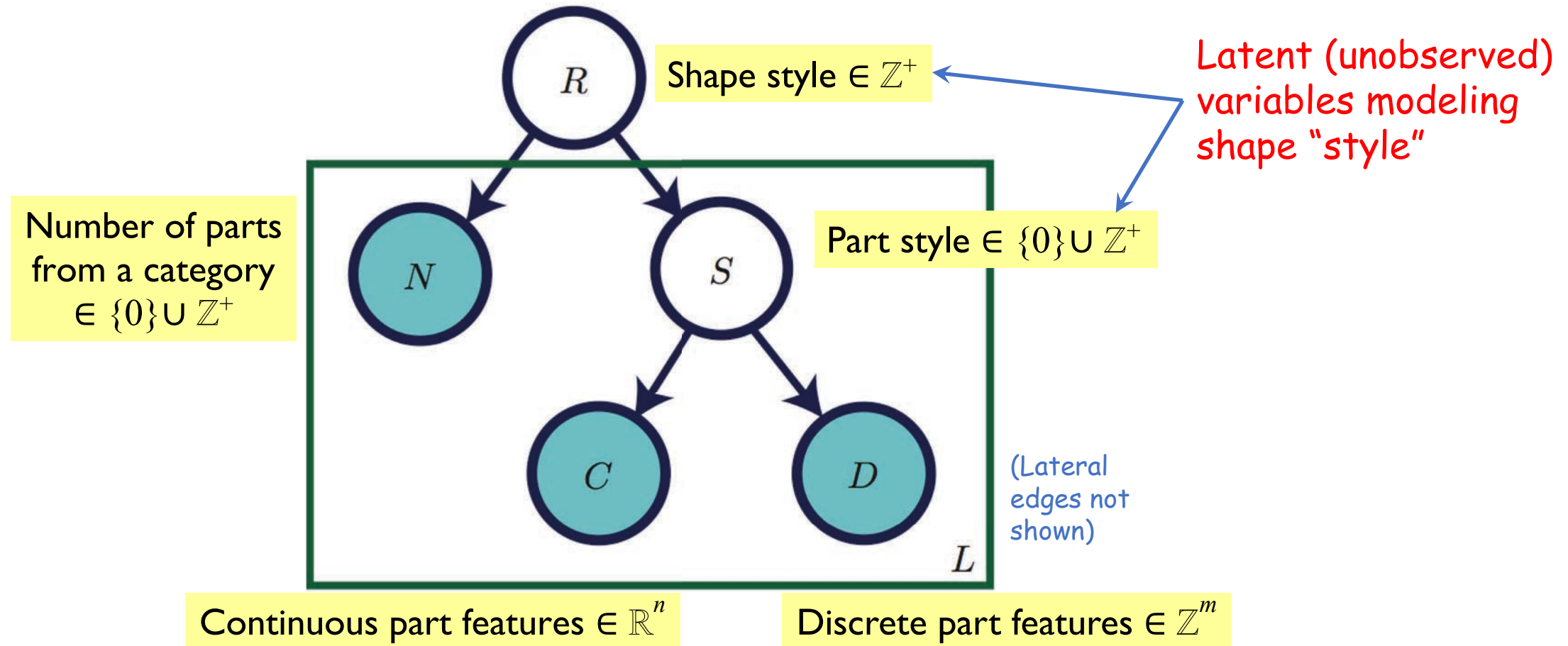
- A **graphical model** visually interprets this factorization as a graph



$$\begin{aligned} P(\mathbf{x}) = & P(x_4 \mid x_2, x_3) \\ & \times P(x_3 \mid x_1) \\ & \times P(x_2 \mid x_1) \\ & \times P(x_1) \end{aligned}$$

A graphical model for assembling existing parts

Part graph encoded with fixed number of variables, layout is deterministic



$$P(\mathbf{X}) = P(R) \prod_{l \in \mathcal{L}} [P(S_l | R) P(N_l | R, \pi(N_l)) P(\mathbf{C}_l | S_l, \pi(\mathbf{C}_l)) P(\mathbf{D}_l | S_l, \pi(\mathbf{D}_l))]$$

Sampling the learned model



Sampling the learned model



Learned latent shape “styles”



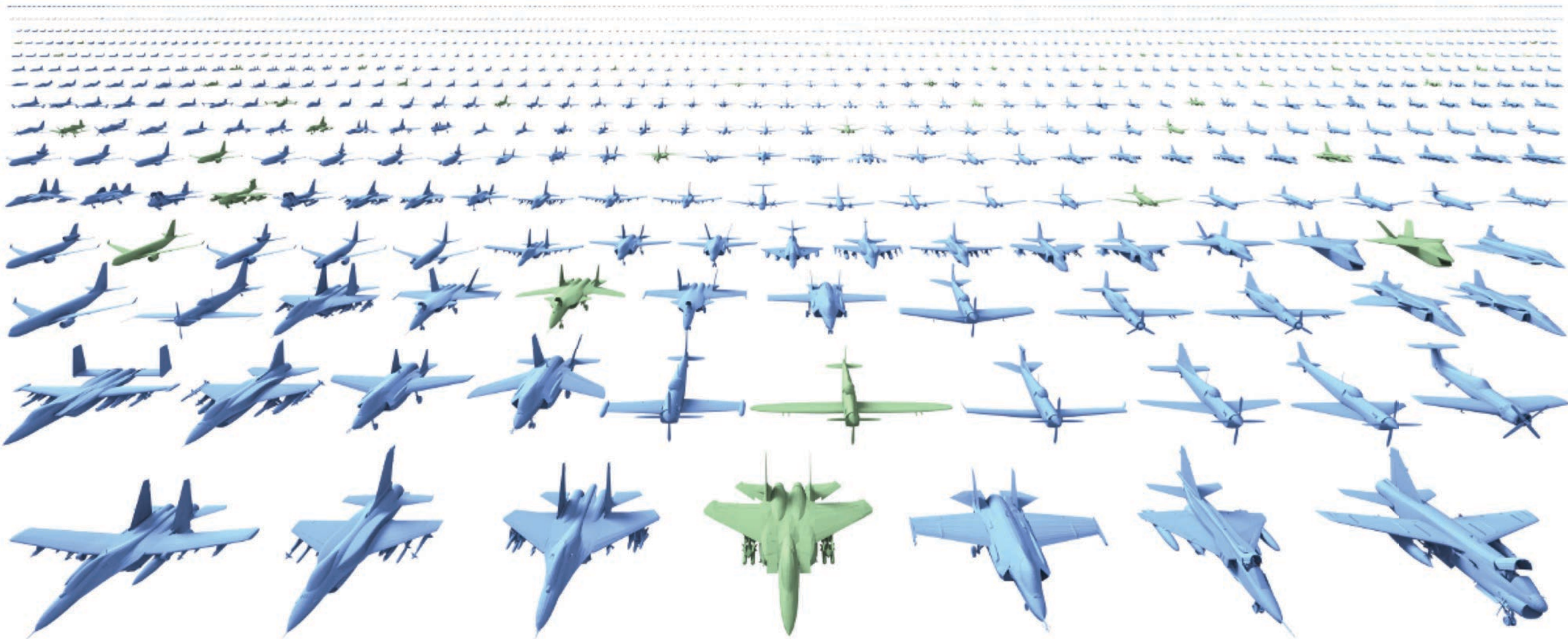
Learned latent component “styles”



Sampling the learned model



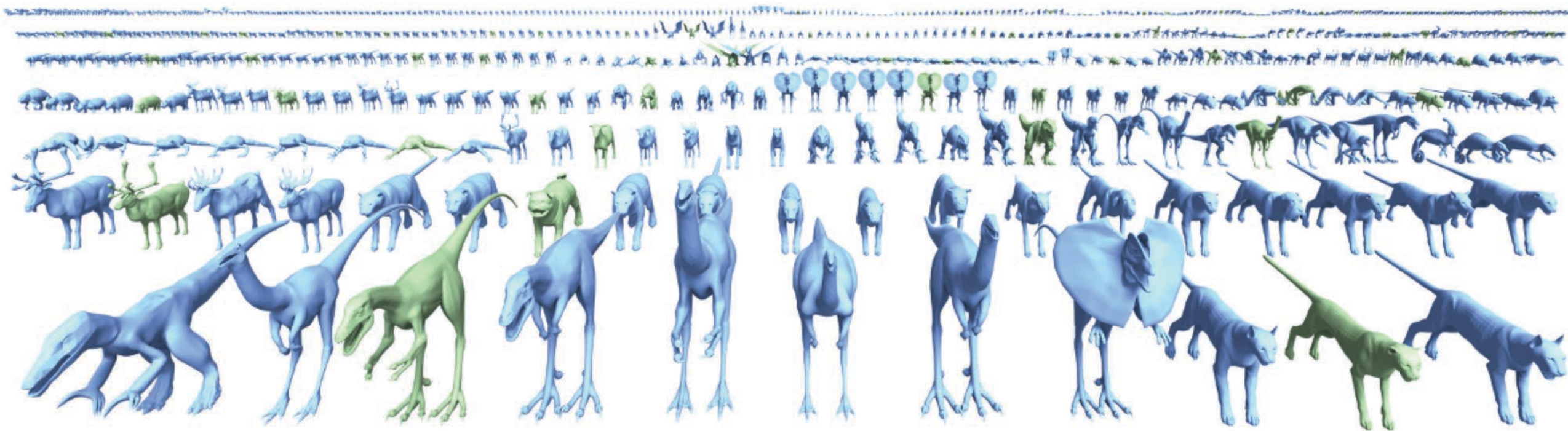
Sampling the learned model



Sampling the learned model

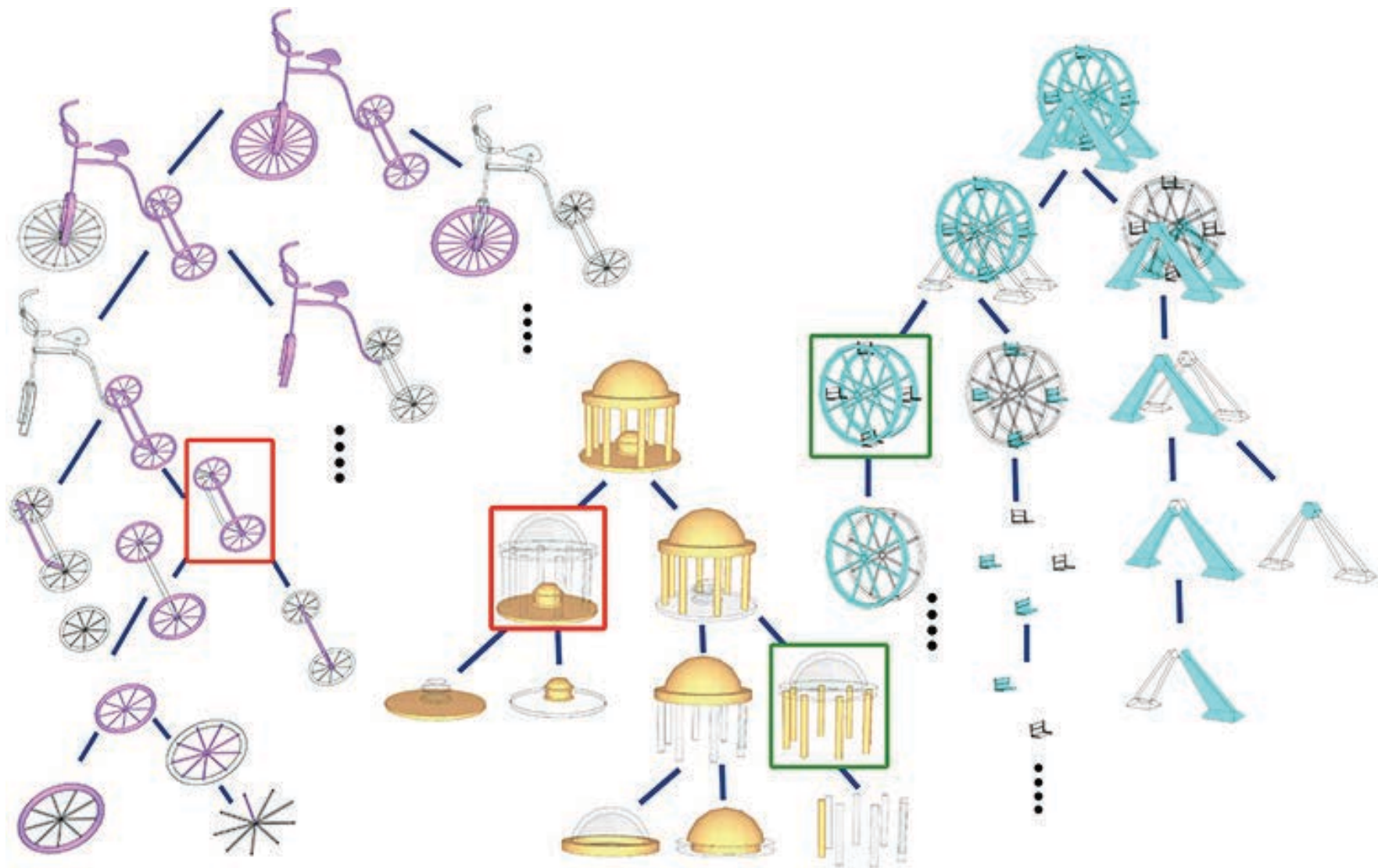


Sampling the learned model



Shapes have natural hierarchical structure

Hierarchical
grouping driven
by cognitive or
functional
imperatives



Probabilistic context-free grammars

- PCFGs are “dynamic” graphical models
- Set of expansion rules, each with a probability

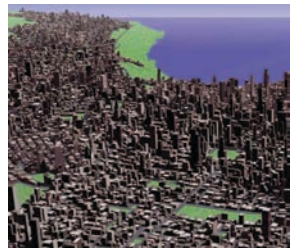
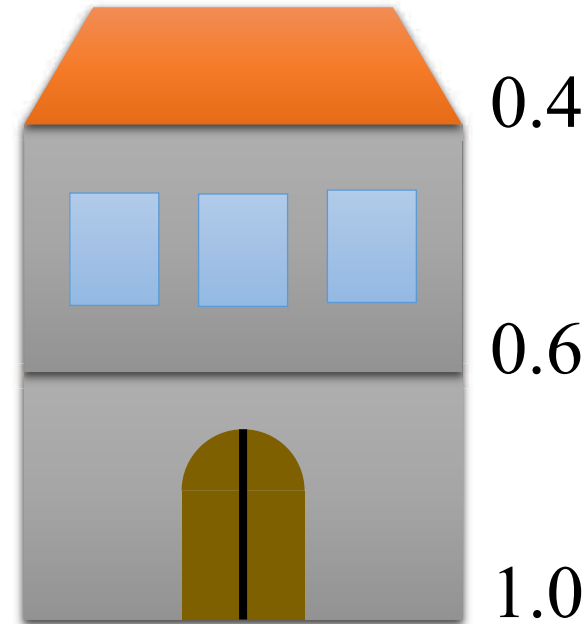
START : *G*

G \rightarrow [Floor Door] *F* (1.0)

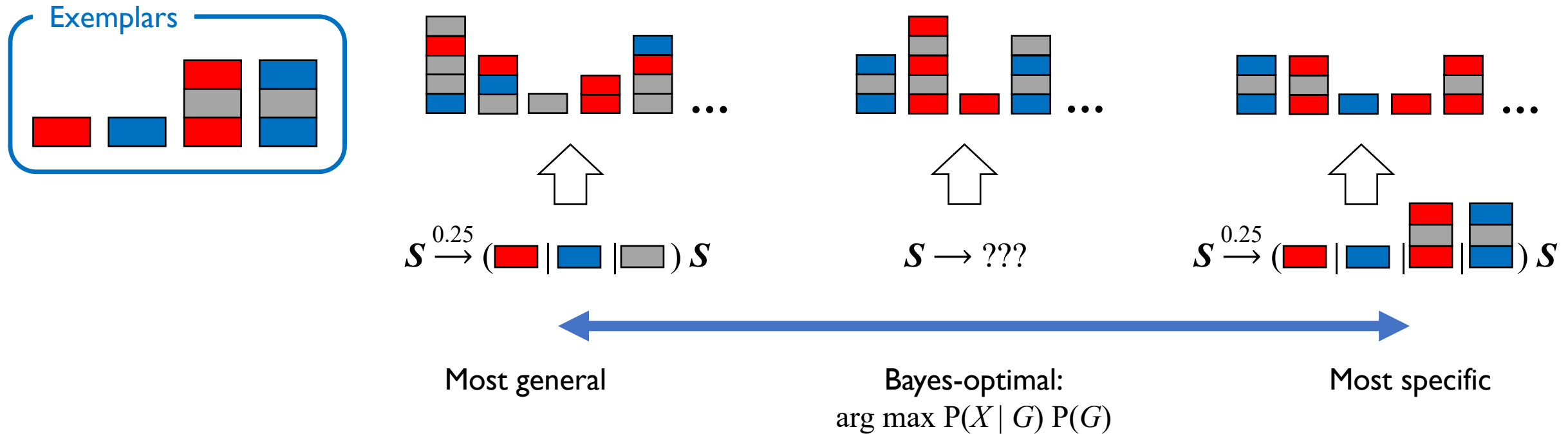
F \rightarrow [Floor Windows] *F* (0.6)

F \rightarrow [Roof] (0.4)

$$P(\mathbf{x}) = 1.0 \times 0.6 \times 0.4$$

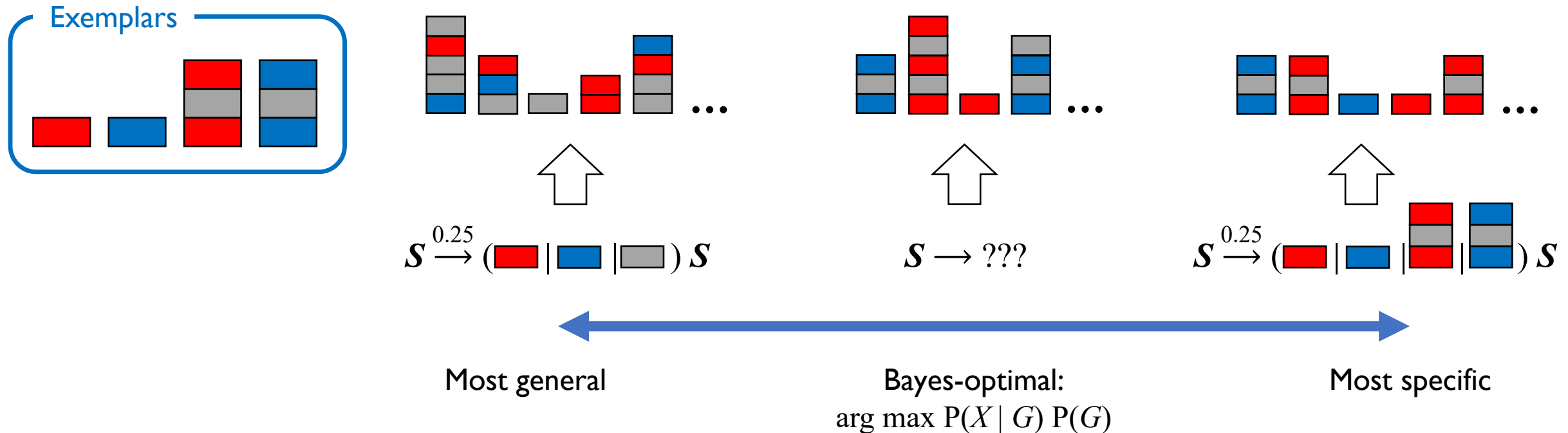


Learning a PCFG from exemplar hierarchies



$P(X | G)$ favors grammars which are likely to produce the exemplars
 $P(G)$ favors grammars which are compact (model prior)

Learning a PCFG from exemplar hierarchies



Optimal grammar G^* is learned by **Bayesian model merging**

- Start with most specific grammar
- Explore more general (compact) grammars by merging and splitting nonterminals to improve $P(X \mid G) P(G)$ via MCMC

Shapes sampled from learned grammars

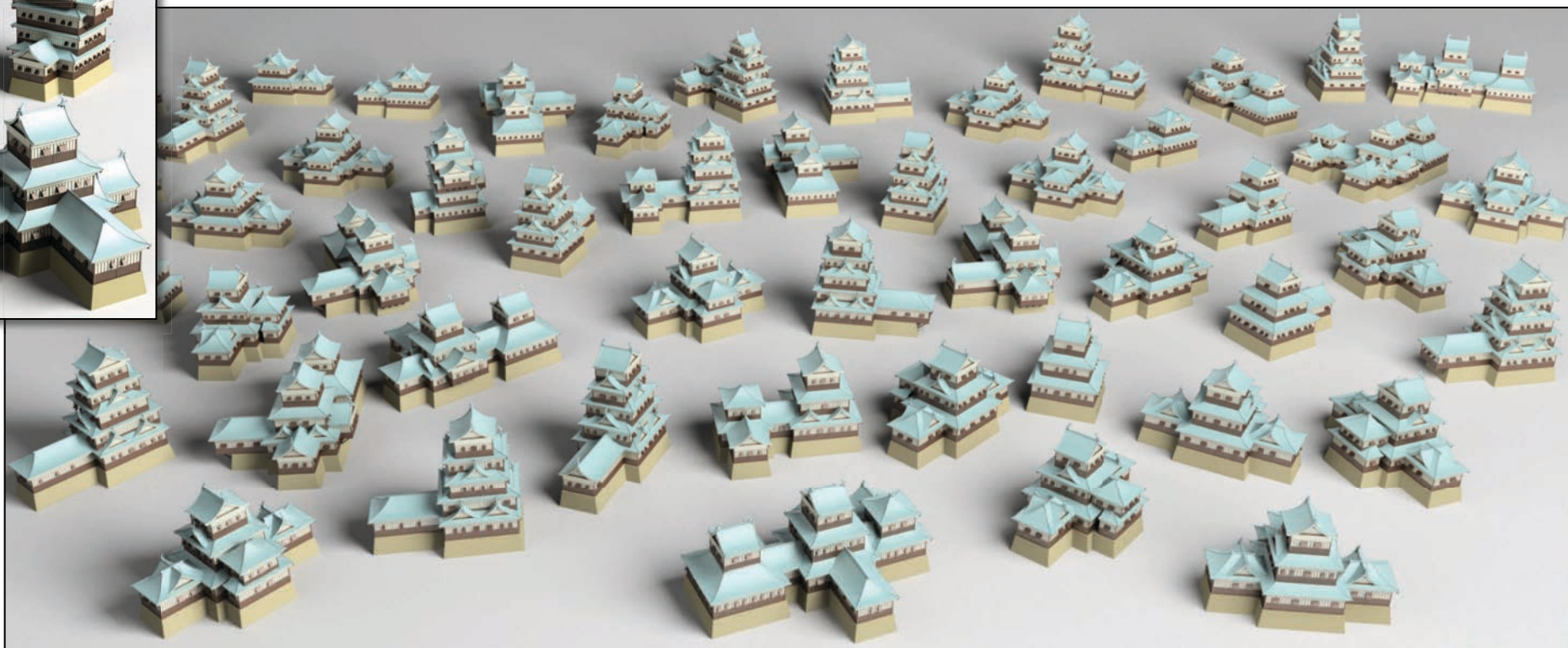
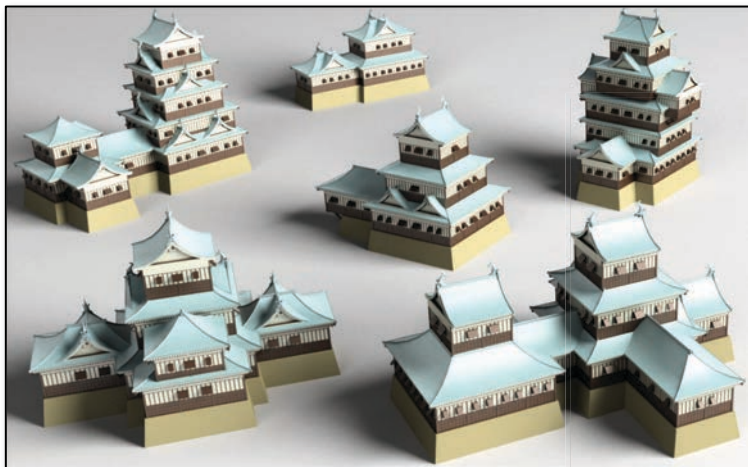
Exemplars



Sampled

Shapes sampled from learned grammars

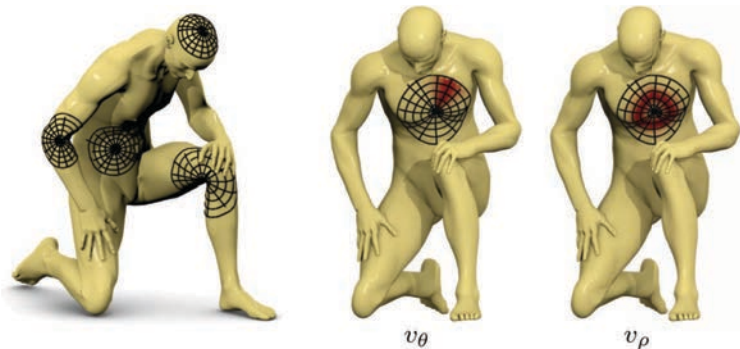
Exemplars



Sampled

Deep generative models of 3D structure

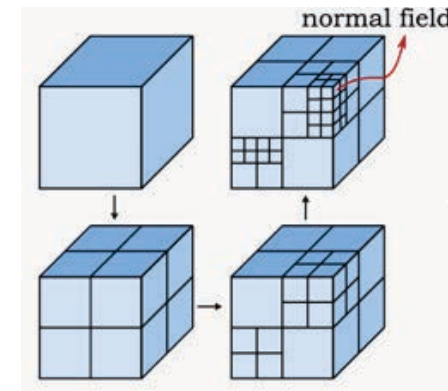
- Classical graphical models are hard to scale to complex, high-dimensional spaces with significant structural variation
- Deep neural networks are powerful high-dimensional models, but do not lend themselves naturally to non-grid-structured (“irregular”) domains such as part layouts and scene graphs
 - Elegant adaptations for low-level, non-structure-aware 3D representations have been developed



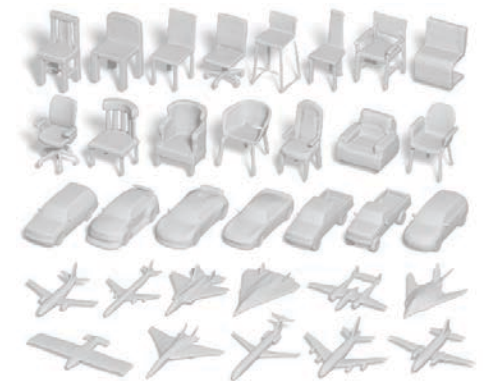
Convolution on curved manifolds



Permutation-invariant NNs



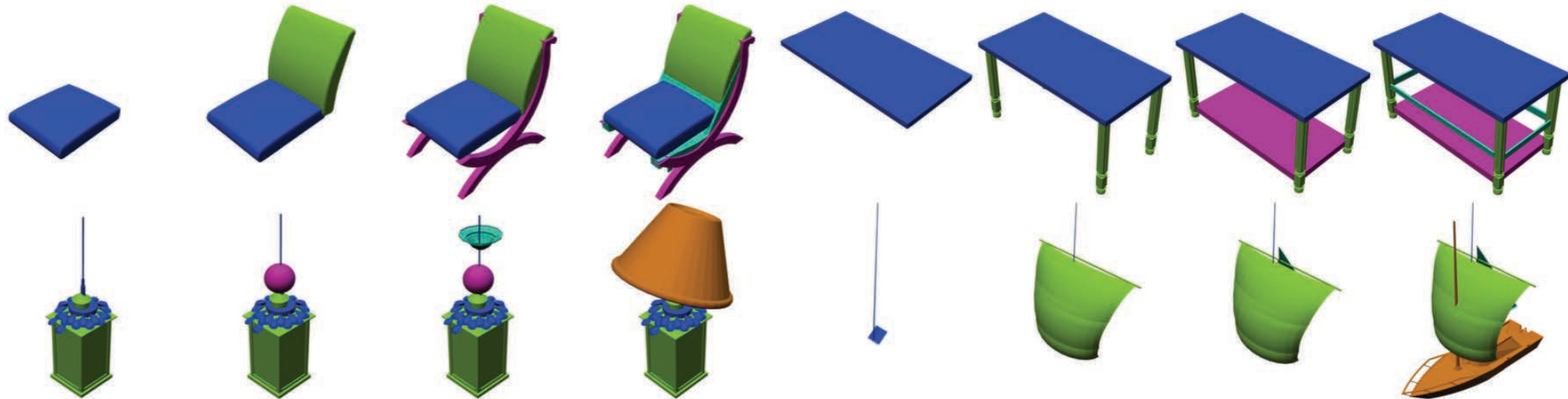
Convolution on octrees



Implicit field NNs

Incremental synthesis from existing parts

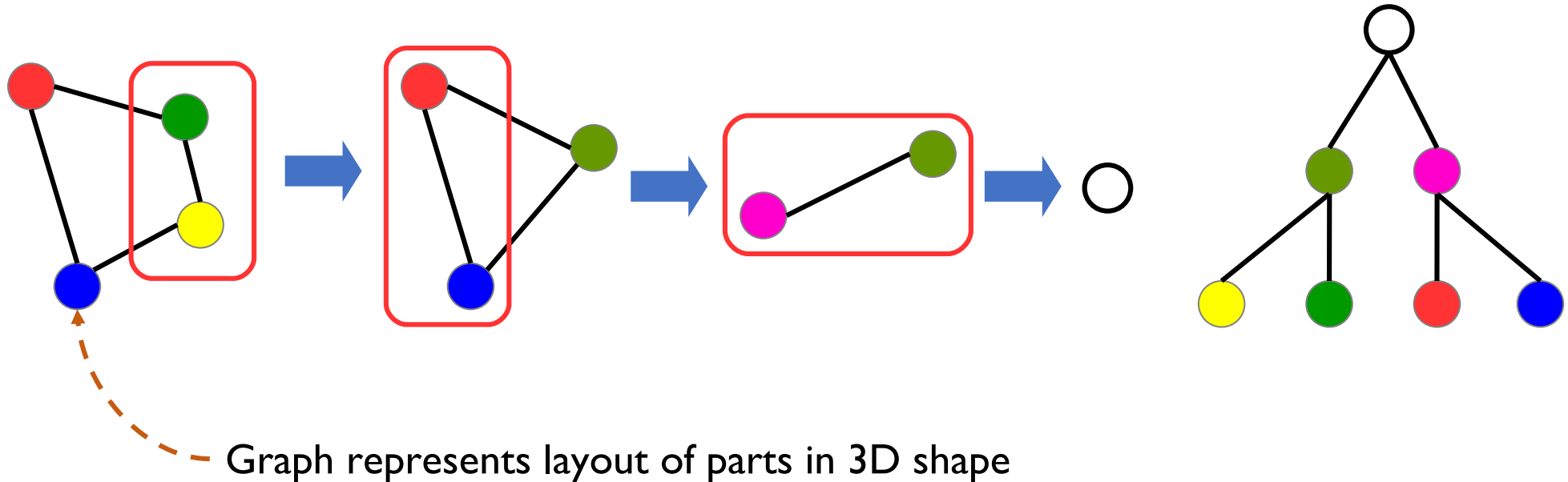
- “Recurrently” construct the shape one part at a time



- Each new part is chosen and placed (by trained networks) to complement the partially constructed shape
- Greedy strategy, avoids joint optimization over the whole structure

Encoding graphs as trees

- Edges of a graph can be collapsed sequentially to form a tree
- Each node of the tree has an associated n-D feature vector, computed recursively
- The process can be reversed to reconstruct the graph (modulo cycles)

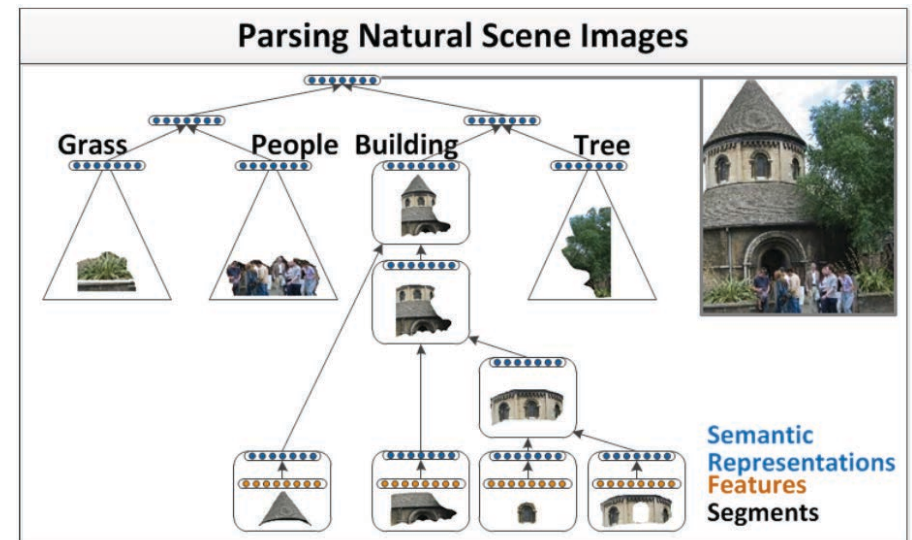
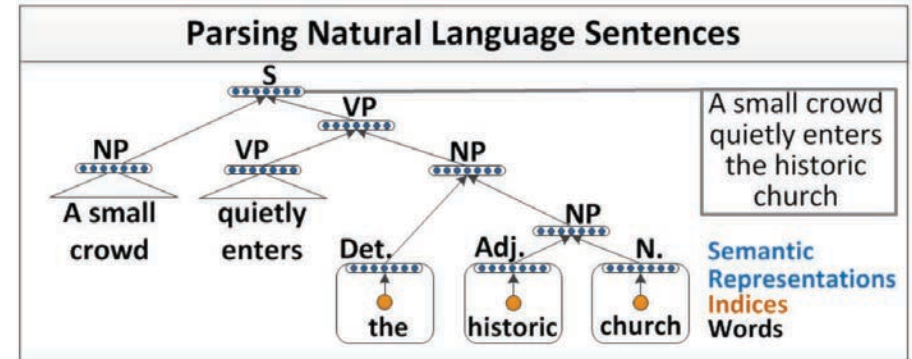
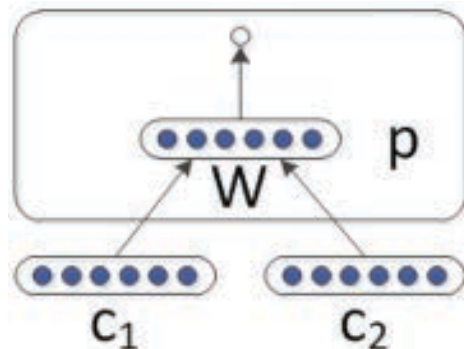


Recursive neural networks (RvNN)

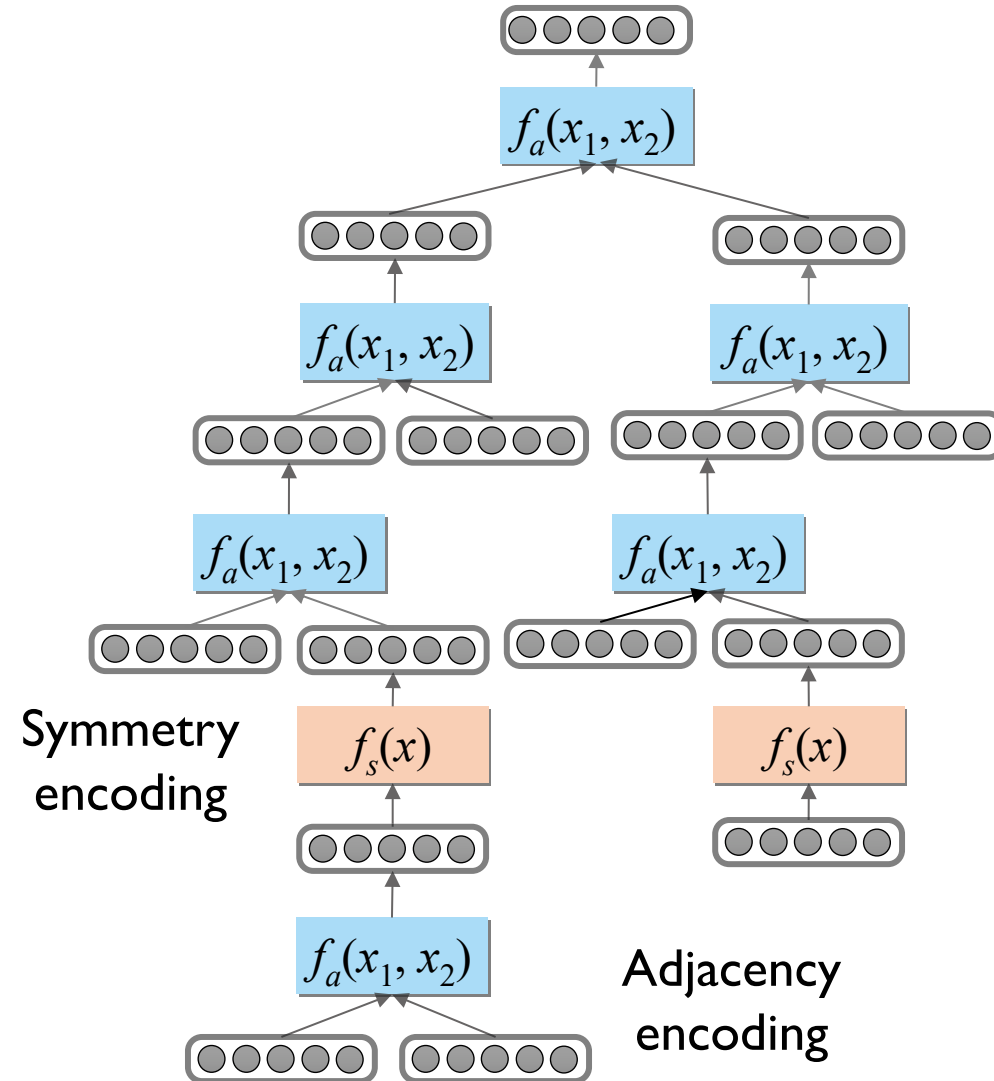
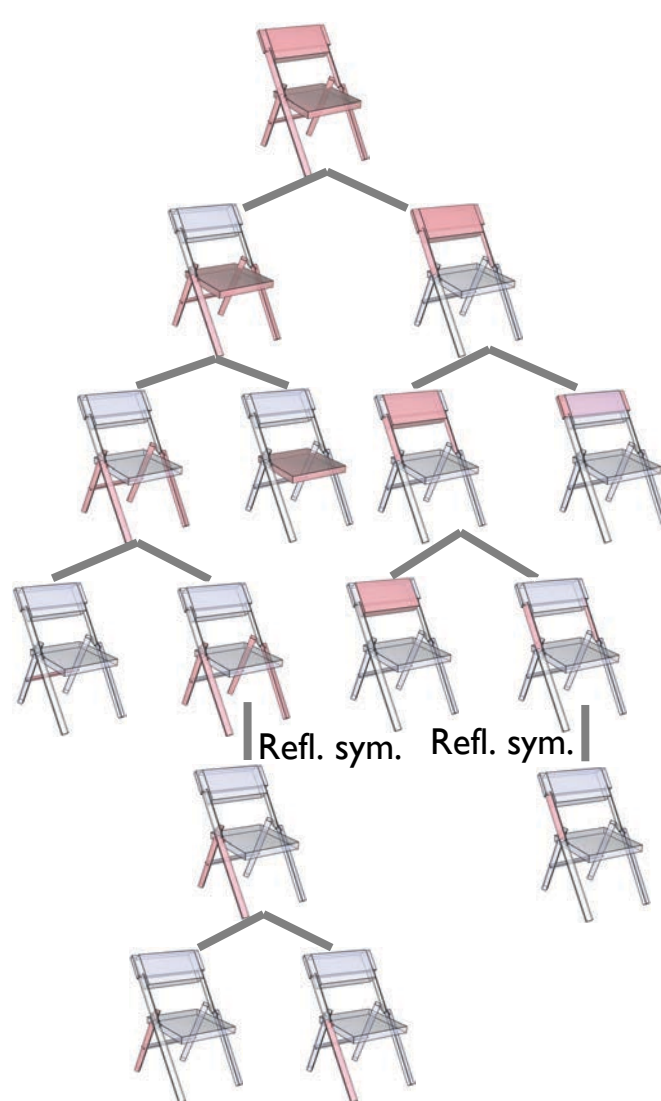
- Repeatedly merge two nodes into one
- Each node has an n-D feature vector, computed recursively

$$p = f(W[c_1; c_2] + b)$$

- End up with a fixed n-D vector for the root node, encoding the whole tree

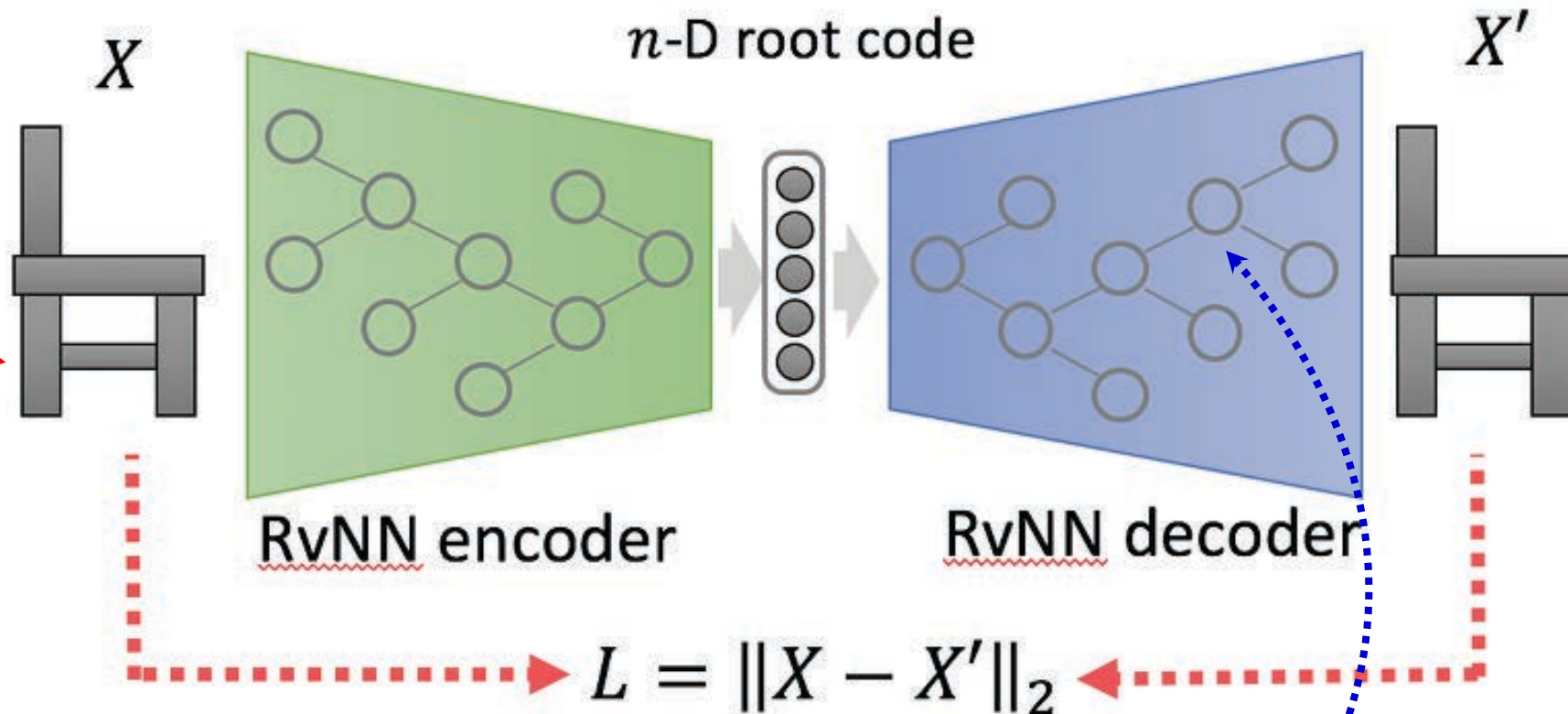


GRASS: Encoding 3D part graphs with RvNNs



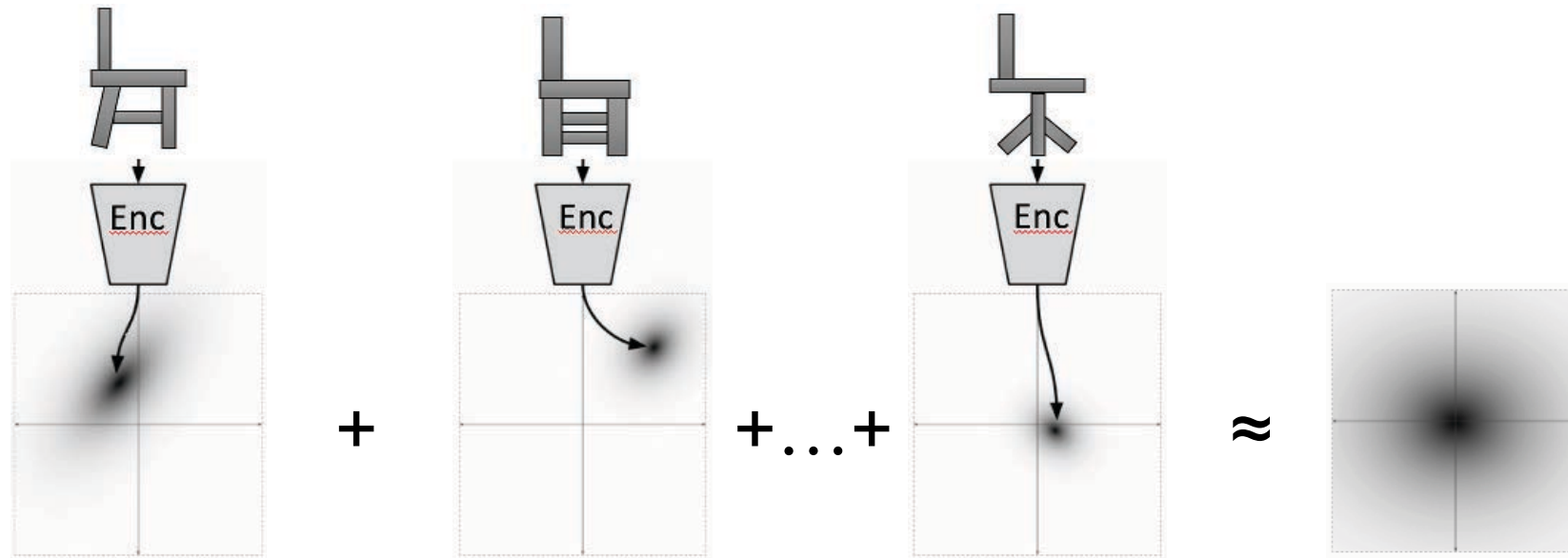
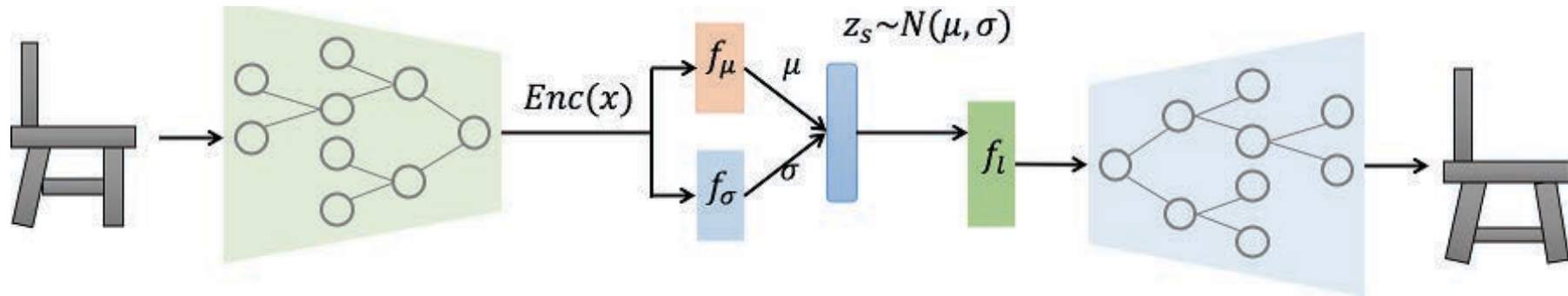
A recursive autoencoder

Input: collection of
unlabeled parts



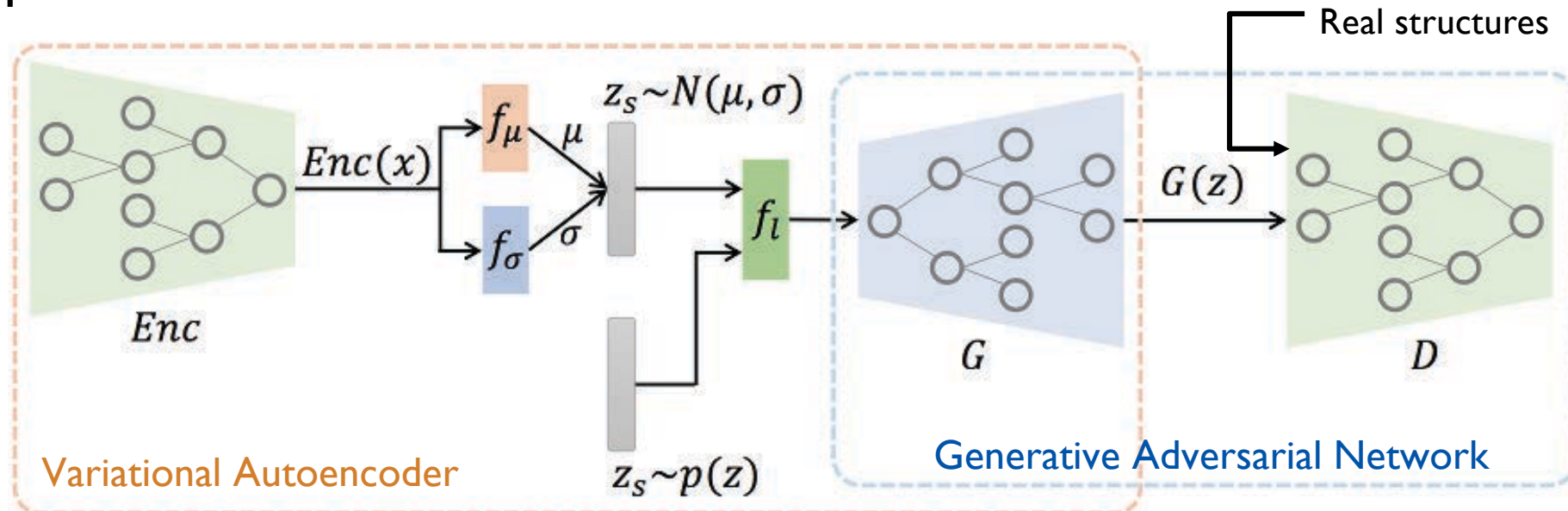
A classifier trained in parallel tells us whether to split a node into two adjacent parts, or into a symmetric group

Making the model generative: Variational Autoencoder (VAE)



Making the VAE robust

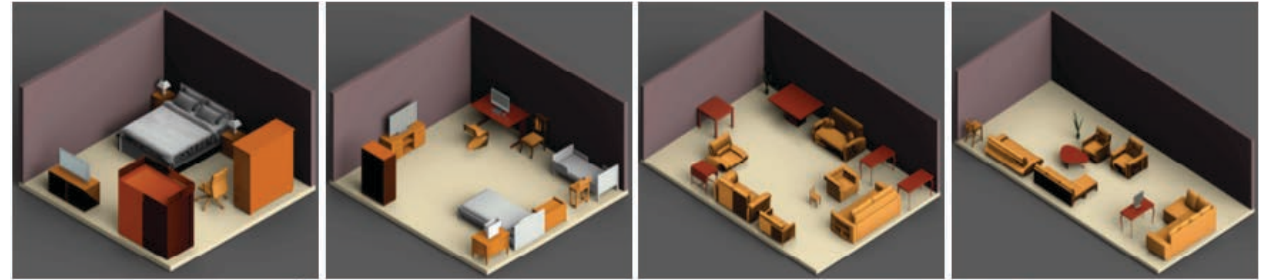
- Adversarial training (VAE-GAN):
 - **Discriminator** tries to tell plausible structure (training hierarchy) from implausible one
 - **Generator** tries to fool discriminator by mapping random codes to plausible structures



Some applications of shape RvNNs



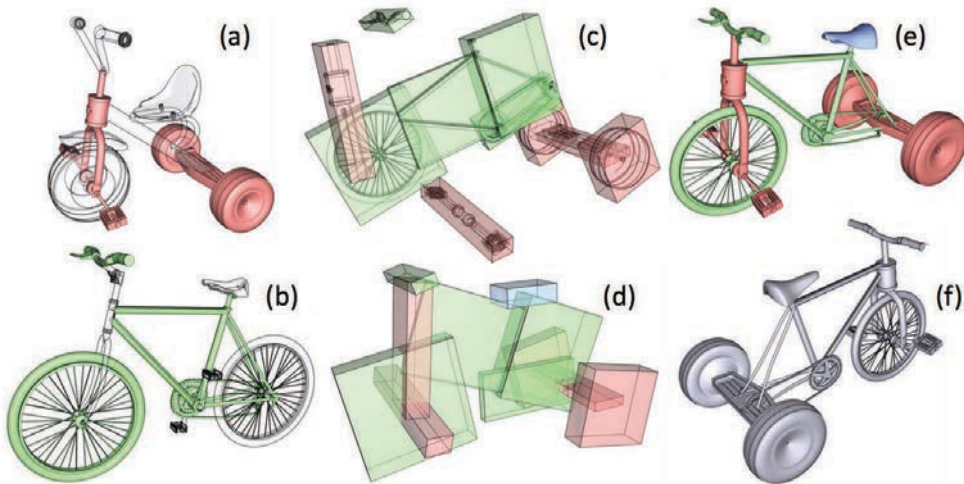
Shape interpolation and synthesis



Bedrooms

Living rooms

Scene synthesis



Shape composition



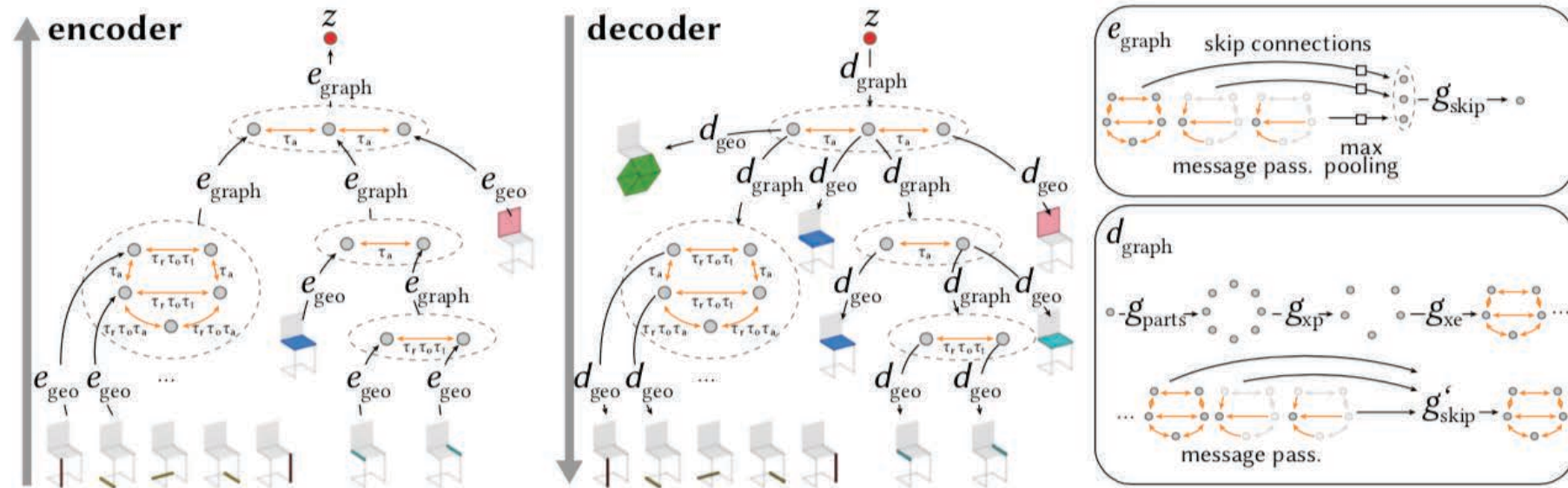
Evolving shape collections



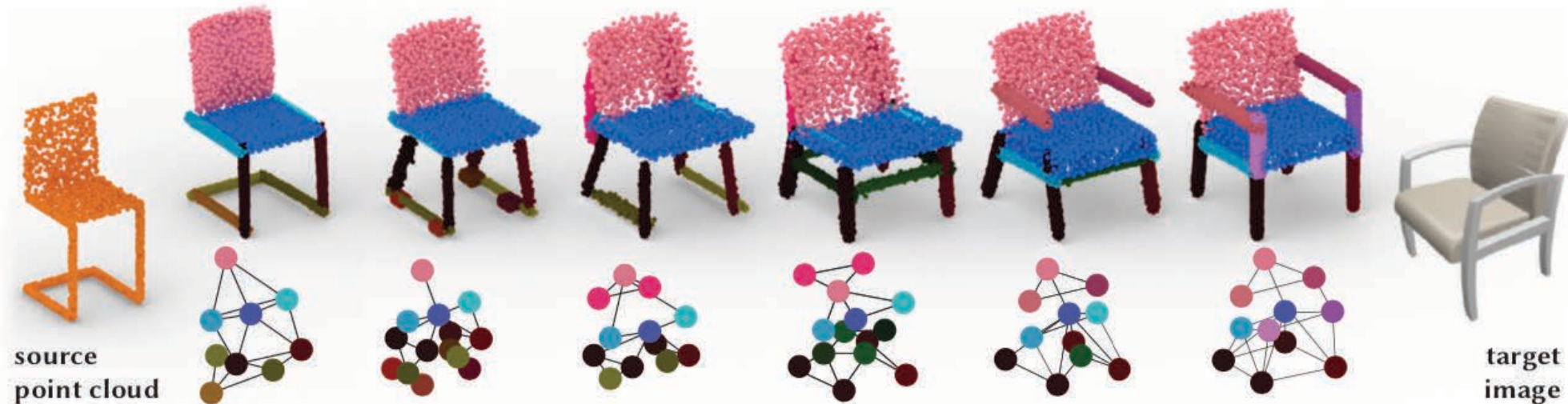
Scan reconstruction

Extending RvNNs with n -ary merges

StructureNet

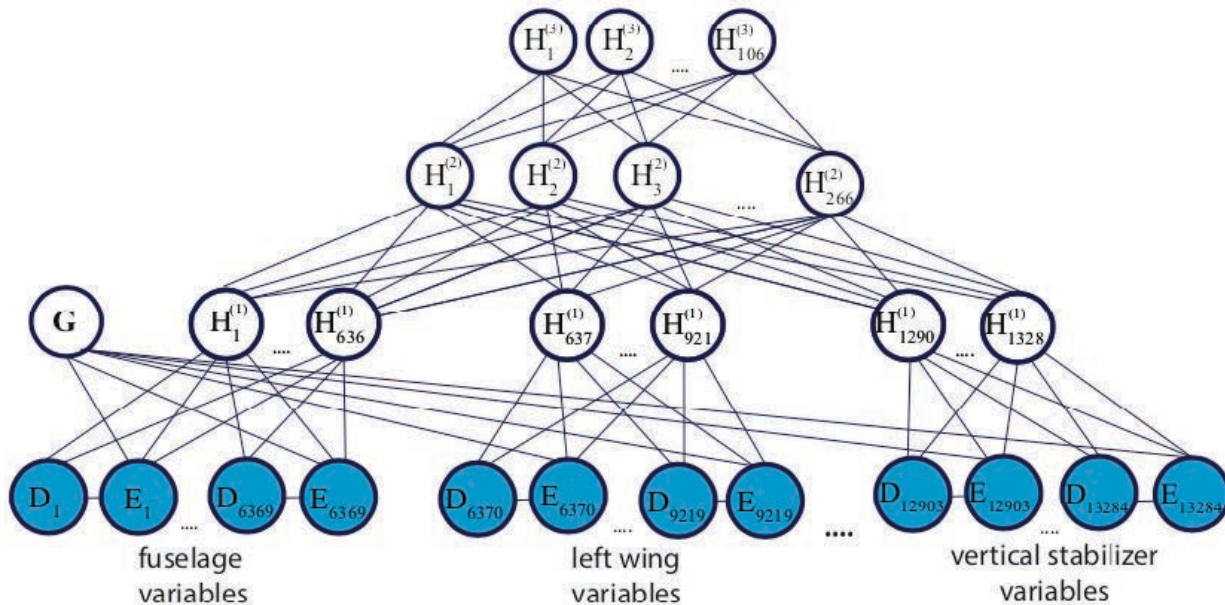


Graph convnets
encode/decode
variable-degree
nodes



Modeling fine-grained geometry

- RvNNs are powerful models of layout variability, but not (yet) of fine-grained local geometry
- Huang et al. [2015] developed a structure-aware surface generator, using a Deep Boltzmann Machine



Modeling fine-grained geometry

- SDM-Net [Gao et al. 2019] models parts as deformed boxes, and chains a part VAE with a simple structure VAE

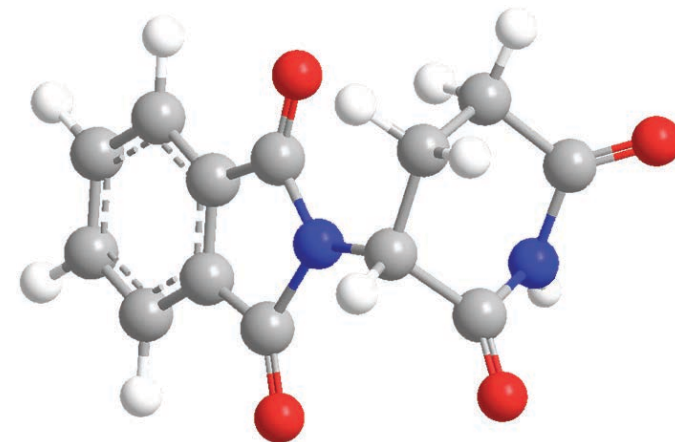


Takeaways

- Structure is an efficient low-dimensional factorization that captures rich shape variation
- Modeling structural variation requires new (or newly borrowed) learning architectures
- Structure-aware generation produces annotated, high-level-editable output “for free”
- Models of semantic and functional intent can layer on top of structural shape priors
[UIST '13, SIGGRAPH '15]



Research Directions



- **Richer structural models**

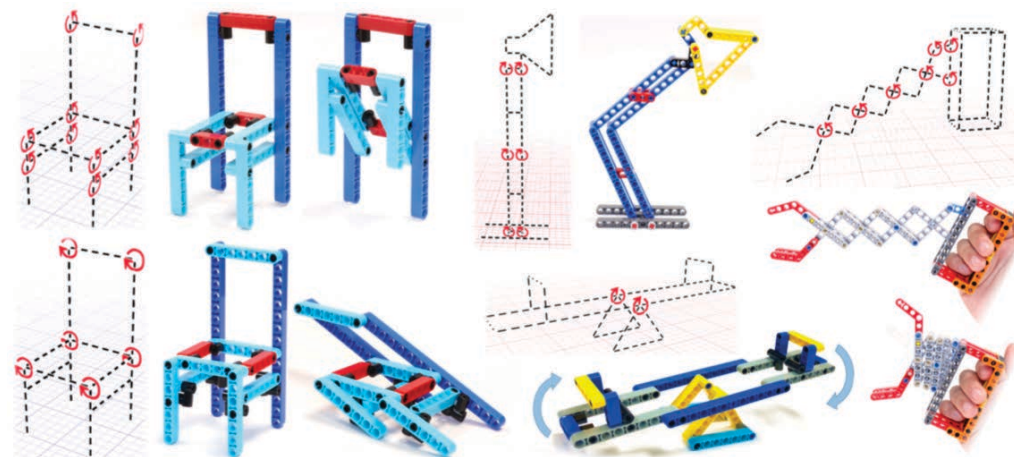
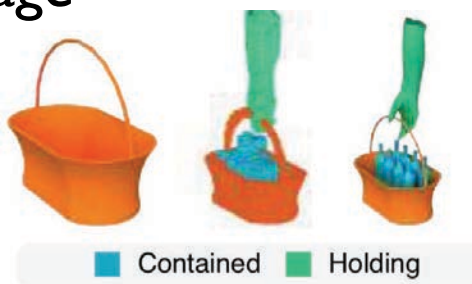
- Domain-specific architectures
- Probabilistic program induction

- **Easier training**

- Weakly/semi-supervised training (infer structure without training structures)
- Few-shot generalization with domain knowledge

- **Richer objectives**

- Linking structure with physical materials and fabrication constraints
- Function, motion, language
 - furniture
 - mechanical assemblies
 - drug design ...



For much more, see
our forthcoming
Eurographics
State-of-the-Art
Report (STAR)

<http://www.cse.iitb.ac.in/~sidch>

Learning Generative Models of 3D Structures

Siddhartha Chaudhuri^{1,2}

Daniel Ritchie³

Jiajun Wu⁴

Kai Xu⁵

Hao Zhang⁶

¹Adobe Research

²IIT Bombay

³Brown University

⁴Stanford University

⁵National University of Defense Technology

⁶Simon Fraser University

