

Neural networks for processing 3D objects

Martin Mirbauer

Survey and Evaluation of Neural 3D Shape Classification Approaches

Martin Mirbauer^{}, Miroslav Krabec, Jaroslav Křivánek^{}, Elena Šikudová^{}

Abstract—Classification of 3D objects – the selection of a category in which each object belongs – is of great interest in the field of machine learning. Numerous researchers use deep neural networks to address this problem, altering the network architecture and representation of the 3D shape used as an input. To investigate the effectiveness of their approaches, we conduct an extensive survey of existing methods and identify common ideas by which we categorize them into a taxonomy. Second, we evaluate 11 selected classification networks on two 3D object datasets, extending the evaluation to a larger dataset on which most of the selected approaches have not been tested yet. For this, we provide a framework for converting shapes from common 3D mesh formats into formats native to each network, and for training and evaluating different classification approaches on this data. Despite being partially unable to reach the accuracies reported in the original papers, we compare the relative performance of the approaches as well as their performance when changing datasets as the only variable to provide valuable insights into performance on different kinds of data. We make our code available to simplify running training experiments with multiple neural networks with different prerequisites.

Index Terms—3D shape analysis, classification algorithms, computer graphics, convolutional neural network, deep learning, image processing, machine learning, multi-layer neural network, neural networks, object recognition.



1 INTRODUCTION

Classification and generation of 3D shapes is one of the widely researched topics in the field of artificial intelligence. It is applied in a vast number of fields such as autonomous driving [1], analysis of medical data [2] as well as various fields of computer vision and graphics [3, 4]. Classification

feature extraction – one of the tasks in the broader context of machine *understanding* of shapes and scenes.

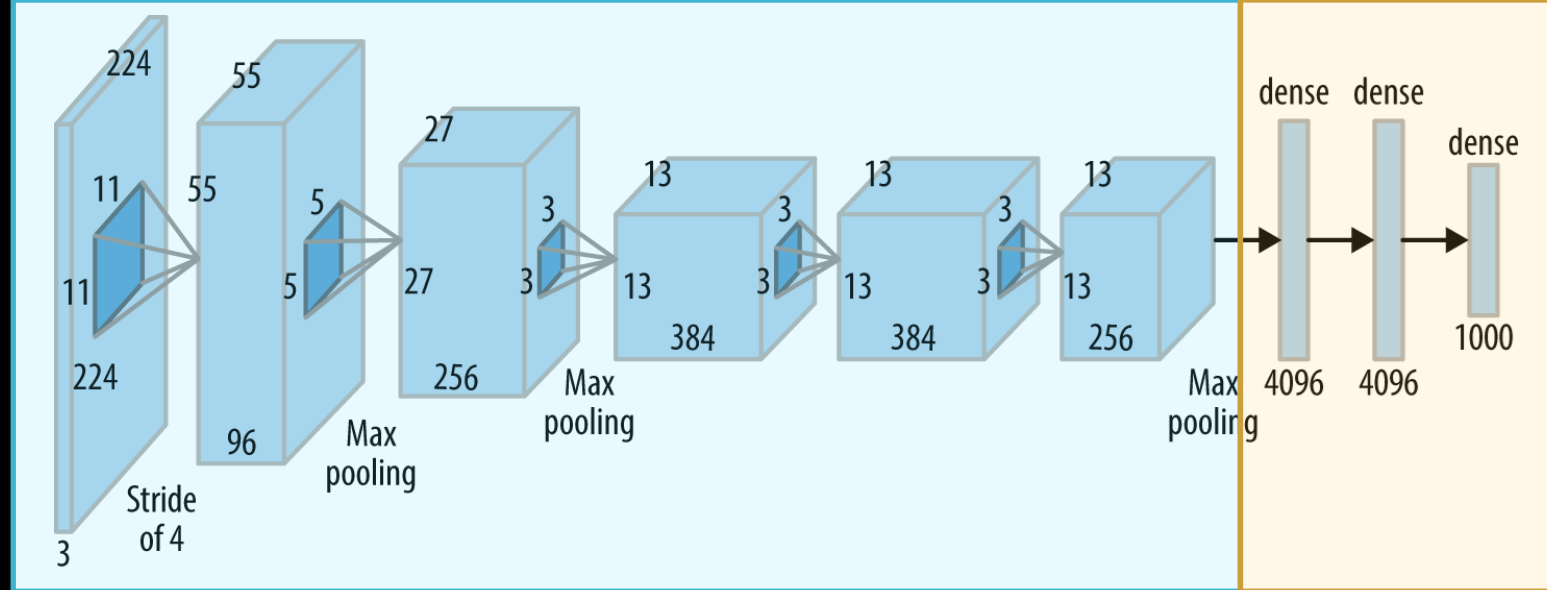
We define the classification task as follows: we are given a set of training examples $\{(x_1, y_1), \dots, (x_n, y_n)\}$, where x_i is a 3D shape representation and y_i is a numerical encoding of the corresponding label. Each shape belongs to exactly one class. A classification model is a parametric model

Goals

- Survey
 - Overview of approaches
 - *Taxonomy*
- Evaluation
 - Replicate
 - *Framework* for running the training experiments

Survey

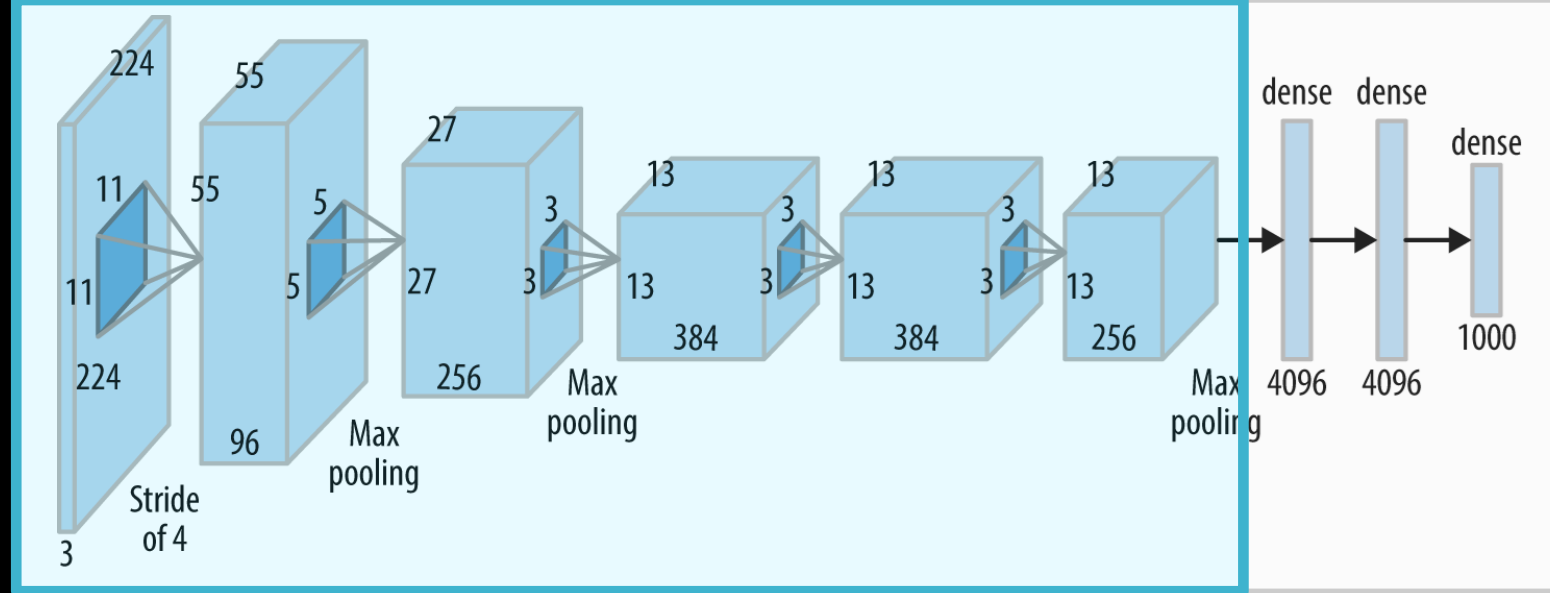
- Task: shape \rightarrow class
- Network structure:
 - **Feature extractor**
 - **Classifier**



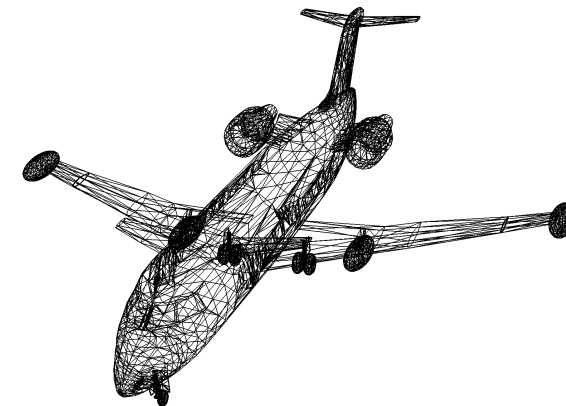
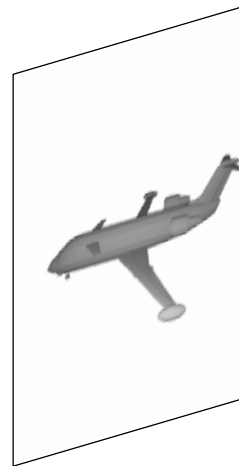
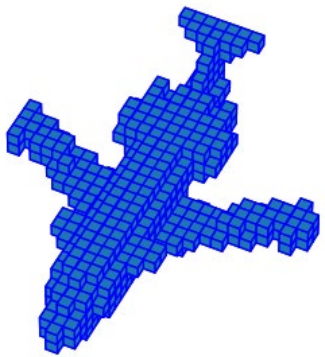
The AlexNet architecture
<https://www.oreilly.com/library/view/tensorflow-for-deep/9781491980446/ch01.html>

Survey

- Task: shape \rightarrow class
- Network structure:
 - **Feature extractor**
 - Classifier
- Representations:



The AlexNet architecture
<https://www.oreilly.com/library/view/tensorflow-for-deep/9781491980446/ch01.html>



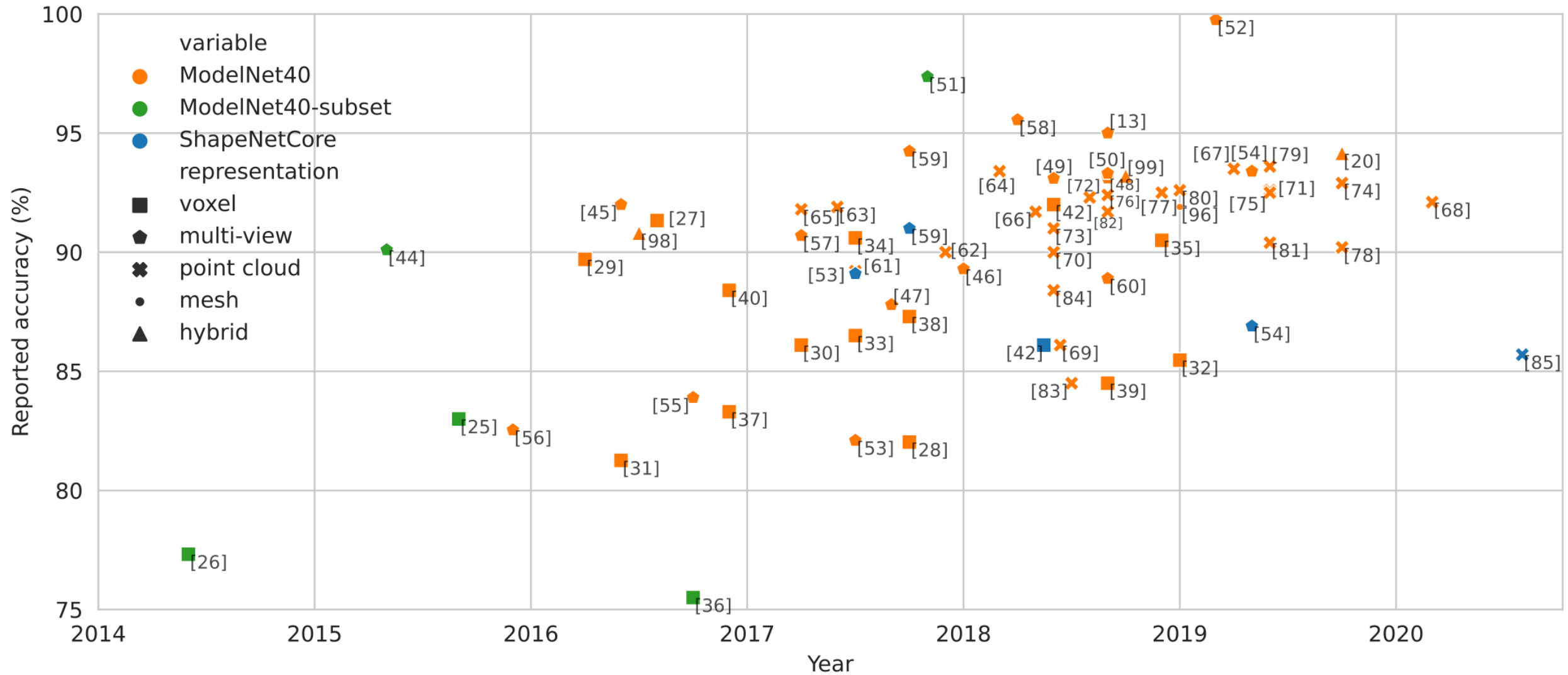


Fig. 2. Reported accuracies of the surveyed methods over time. Datasets and input representations are denoted by different colors and shapes.

Taxonomy of approaches

Volumetric grid

Multi-view (images)

Point cloud

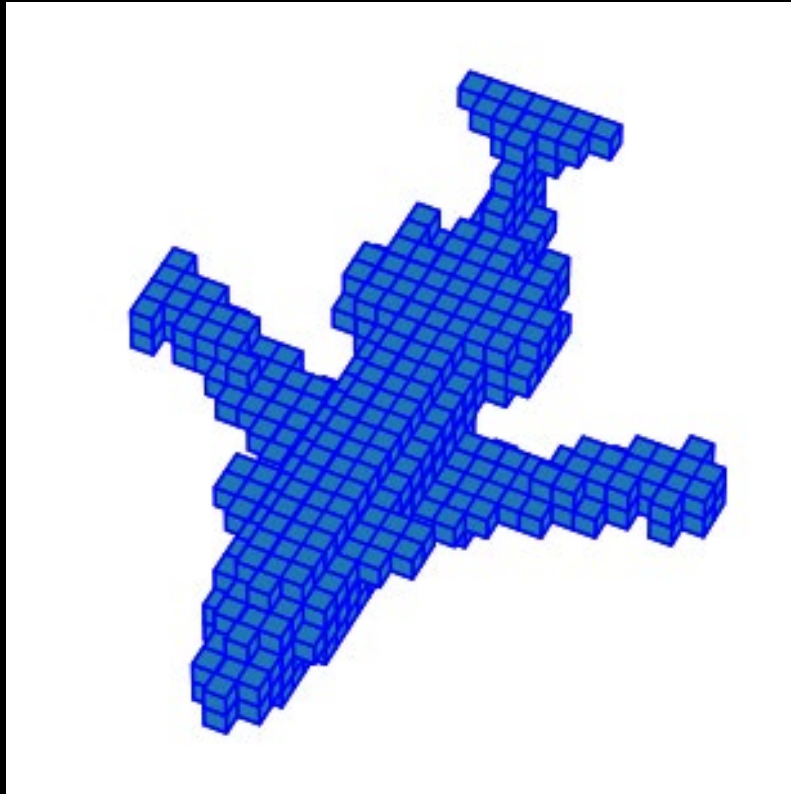
Surface shape

Hybrid

Manifold-based

Volumetric grid-based approaches

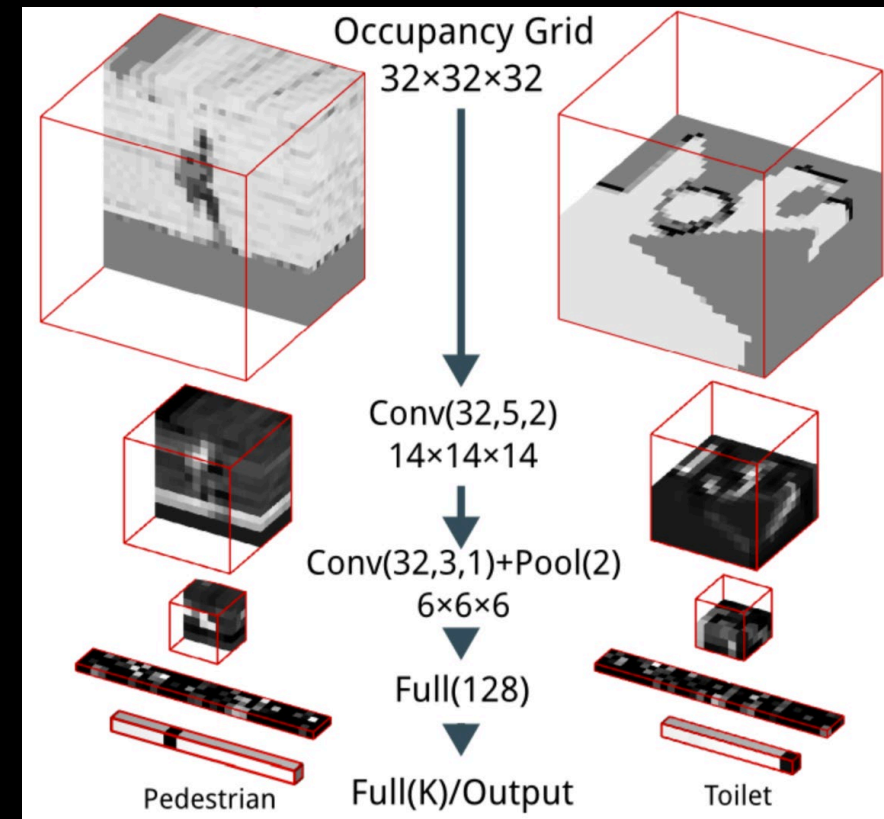
Volumetric grid
Basic Architectures
Voxel CNN with Residual Connections
Auxiliary Task
Network Architecture Optimization
Octree-represented Voxel Grid
Unsupervised Representation Learning
Non-convolutional Approaches
Conversion from a Point Cloud



Volumetric grid-based approaches

Volumetric grid
Basic Architectures
Voxel CNN with Residual Connections
Auxiliary Task
Network Architecture Optimization
Octree-represented Voxel Grid
Unsupervised Representation Learning
Non-convolutional Approaches
Conversion from a Point Cloud

- **VoxNet** [Maturana and Scherer 2015]
- **3D ShapeNets** [Wu et al. 2014]
 - ModelNet dataset



Volumetric grid-based approaches

Volumetric grid
Basic Architectures
Voxel CNN with Residual Connections
Auxiliary Task
Network Architecture Optimization
Octree-represented Voxel Grid
Unsupervised Representation Learning
Non-convolutional Approaches
Conversion from a Point Cloud

- „**Voxception-ResNet**“ in *Generative and discriminative voxel modeling with convolutional neural networks* [Brock et al. 2016] (“vrn”)

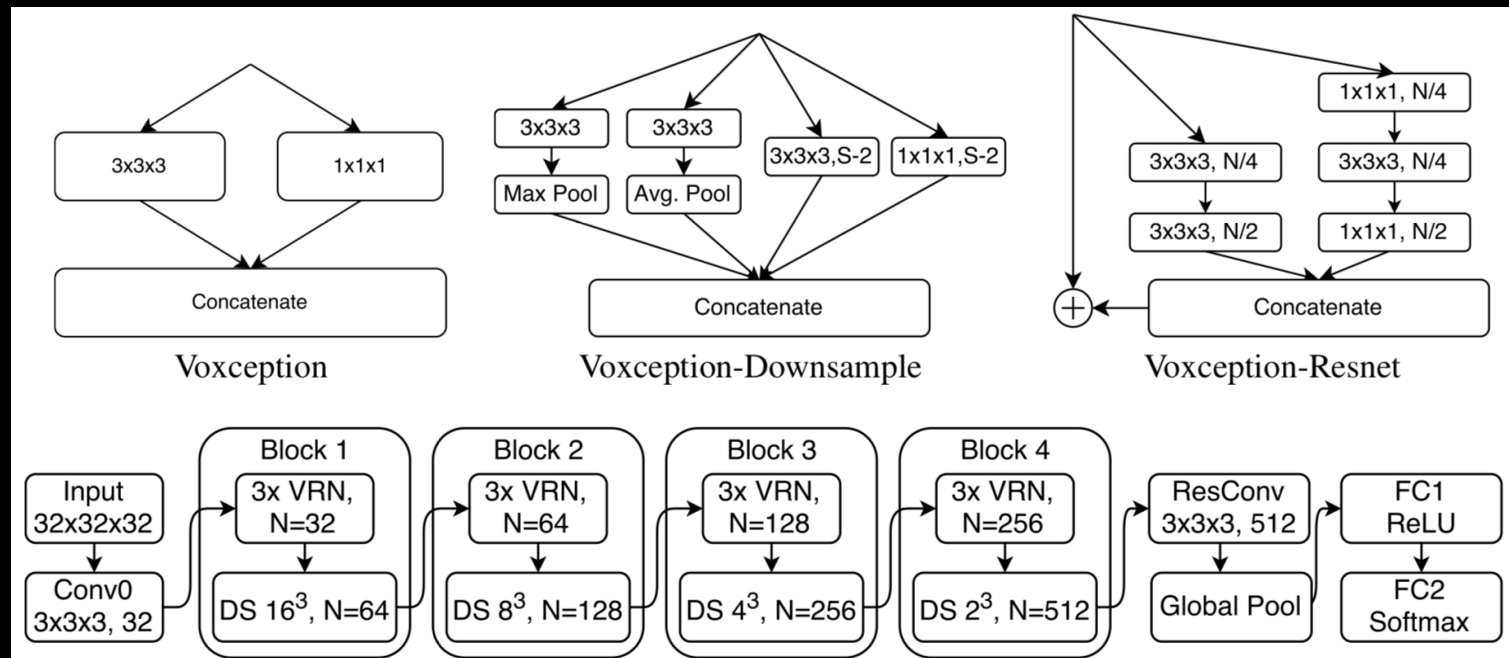
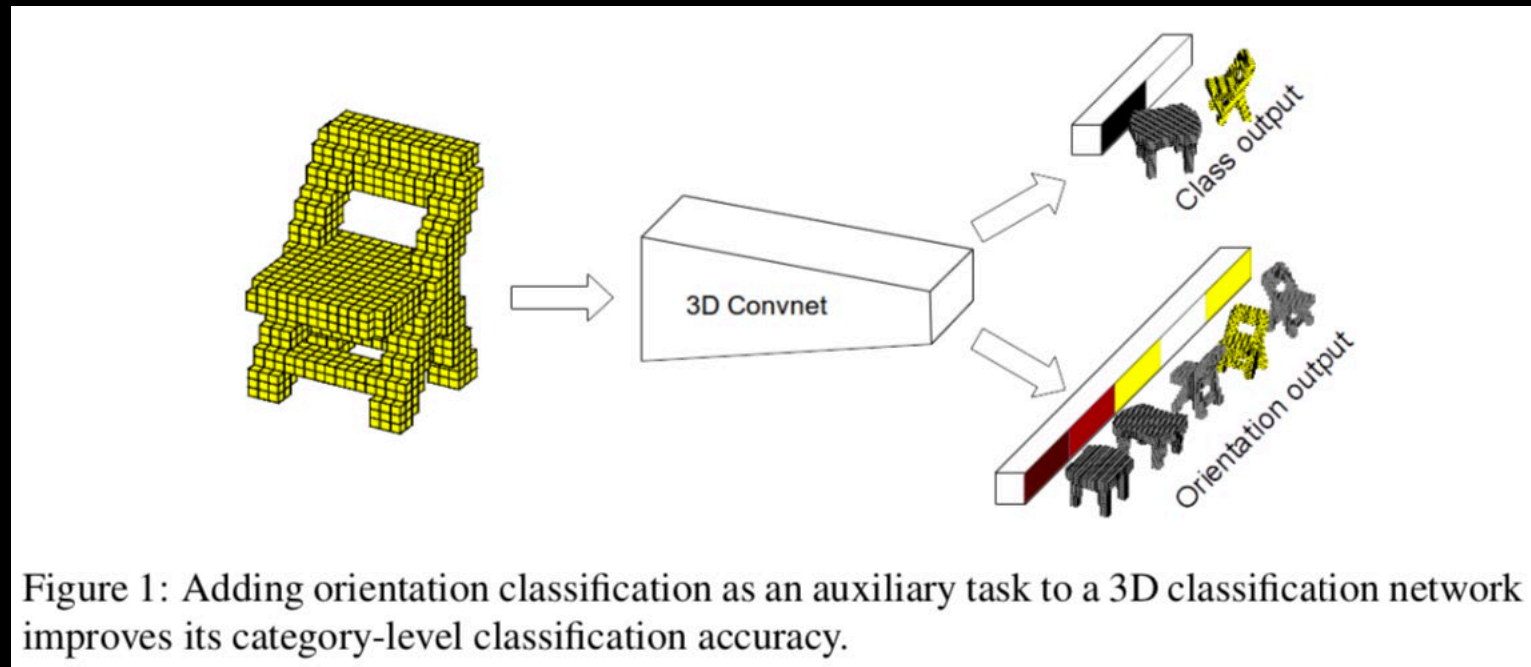


Figure 4: Voxception-ResNet 45 Layer Architecture. DS are Voxception-Downsample blocks.

Volumetric grid-based approaches

Volumetric grid
Basic Architectures
Voxel CNN with Residual Connections
Auxiliary Task
Network Architecture Optimization
Octree-represented Voxel Grid
Unsupervised Representation Learning
Non-convolutional Approaches
Conversion from a Point Cloud

- “**ORION**” in **O**rientation-boosted **v**oxel **n**ets for 3D object recognition [Sedaghat et al. 2016]



Volumetric grid-based approaches

Volumetric grid
Basic Architectures
Voxel CNN with Residual Connections
Auxiliary Task
Network Architecture Optimization
Octree-represented Voxel Grid
Unsupervised Representation Learning
Non-convolutional Approaches
Conversion from a Point Cloud

- **LightNet**: A lightweight 3D convolutional neural network for real-time 3D object recognition [Zhi et al. 2017]
- Beam search for learning a deep convolutional neural network of 3D shapes [Xu and Todorovic 2016]

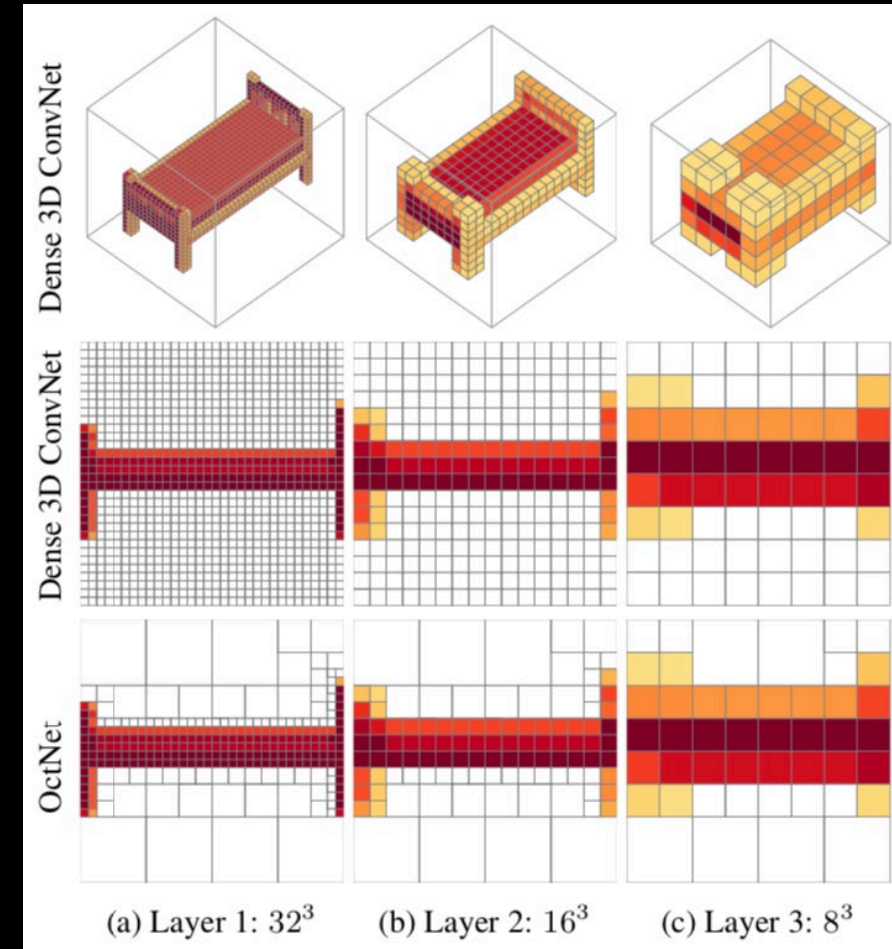
Volumetric grid-based approaches

Volumetric grid
Basic Architectures
Voxel CNN with Residual Connections
Auxiliary Task
Network Architecture Optimization
Octree-represented Voxel Grid
Unsupervised Representation Learning
Non-convolutional Approaches
Conversion from a Point Cloud

- **OctNet** [Riegler et al. 2017]
- **O-CNN** [Wang et al. 2017] (“octree”)
- **Adaptive O-CNN** [Wang et al. 2018] (“octree-adaptive”)

Network	without voting	with voting
O-CNN(3)	85.5%	87.1%
O-CNN(4)	88.3%	89.3%
O-CNN(5)	89.6%	90.4%
O-CNN(6)	89.9%	90.6%
O-CNN(7)	89.5%	90.1%
O-CNN(8)	89.6%	90.2%

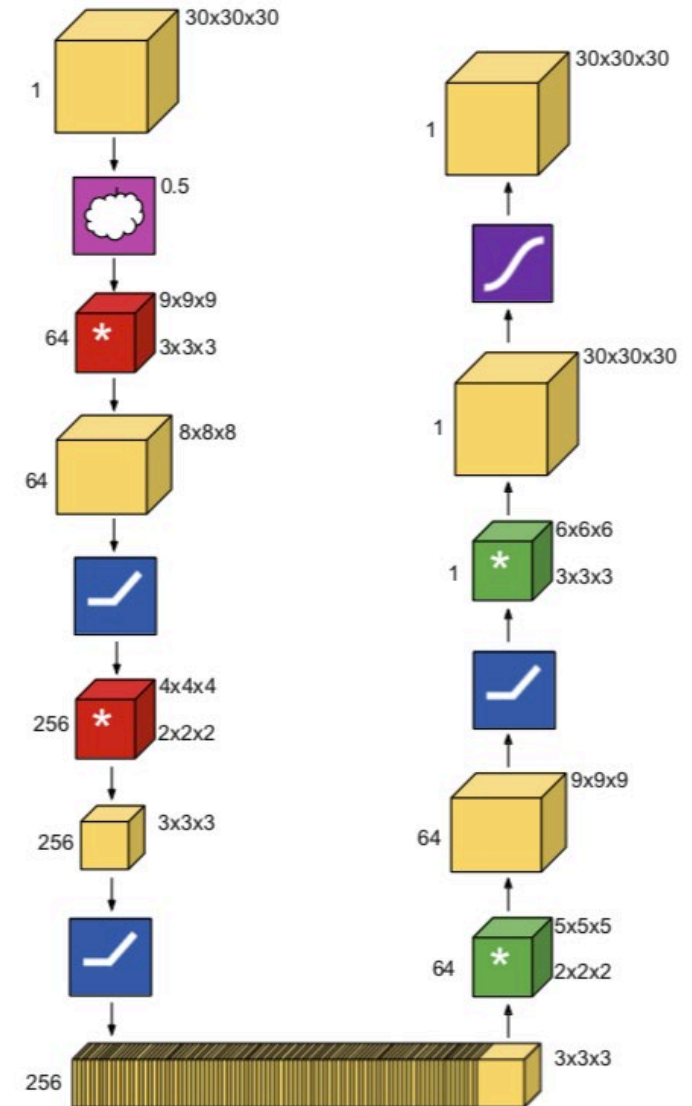
Table 1. Object classification results on ModelNet40 dataset. The number shown in the table is the accuracy of object recognition. The second and third columns show the results of the network with and without voting. Numbers in parentheses are the resolutions of voxels. A number in boldface emphasizes the best result.



Volumetric grid-based approaches

- **VConv-DAE** [Sharma et al. 2016]

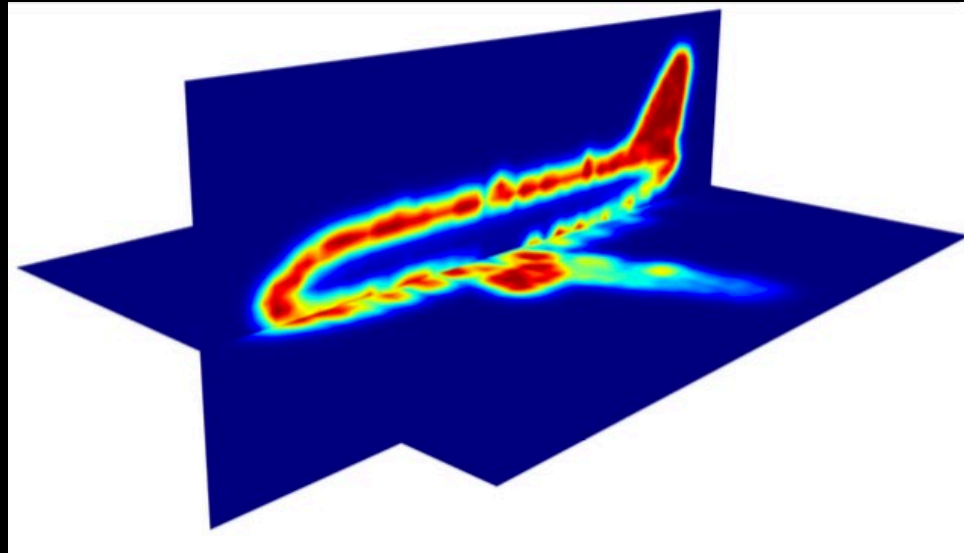
Volumetric grid
Basic Architectures
Voxel CNN with Residual Connections
Auxiliary Task
Network Architecture Optimization
Octree-represented Voxel Grid
Unsupervised Representation Learning
Non-convolutional Approaches
Conversion from a Point Cloud



Volumetric grid-based approaches

Volumetric grid
Basic Architectures
Voxel CNN with Residual Connections
Auxiliary Task
Network Architecture Optimization
Octree-represented Voxel Grid
Unsupervised Representation Learning
Non-convolutional Approaches
Conversion from a Point Cloud

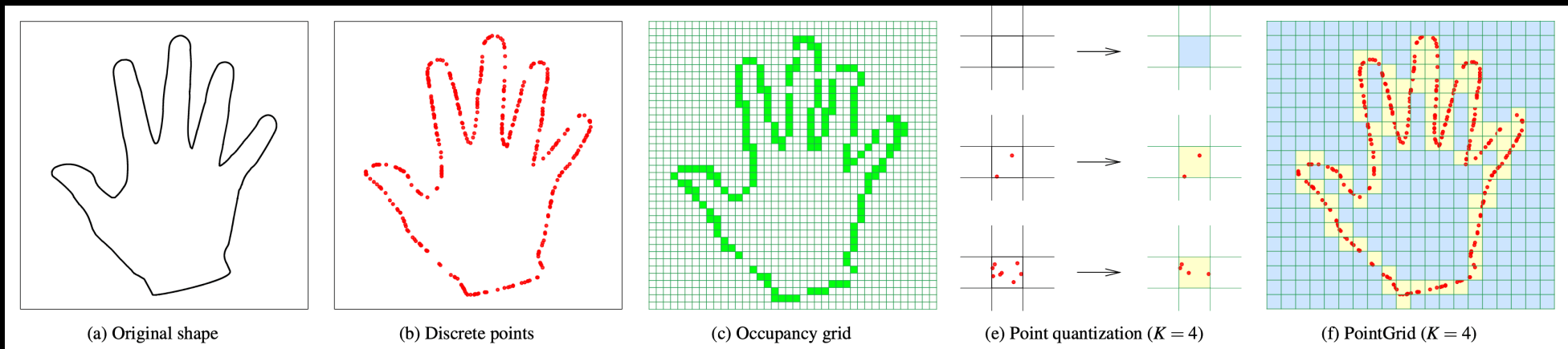
- **FPNN**: Field probing neural networks for 3D data [Li et al. 2016]



Volumetric grid-based approaches

Volumetric grid
Basic Architectures
Voxel CNN with Residual Connections
Auxiliary Task
Network Architecture Optimization
Octree-represented Voxel Grid
Unsupervised Representation Learning
Non-convolutional Approaches
Conversion from a Point Cloud

- **VoxNet** [Maturana and Scherer 2015]
- **PointGrid** [Le and Duan 2018]



Volumetric grid-based approaches

- 3D convolutions
- Skip connections, residual blocks
- Data augmentation (rotation), architecture optimizations
- Low resolution, slow evaluation (except for octree representation)

Multi-view image-based approaches

Multi-view (images)	
Basic Architectures	
Multiple Modalities	
Axis-aligned Views	
Learned View Grouping	
Unsupervised Viewpoints Assignment	
Unsupervised Representation Learning	
Using Auxiliary Data	
Special Projections	Geometry Images
	Panorama
	Spherical



Multi-view image-based approaches

Multi-view (images)

Basic Architectures

Multiple Modalities

Axis-aligned Views

Learned View Grouping

Unsupervised Viewpoints Assignment

Unsupervised Representation Learning

Using Auxiliary Data

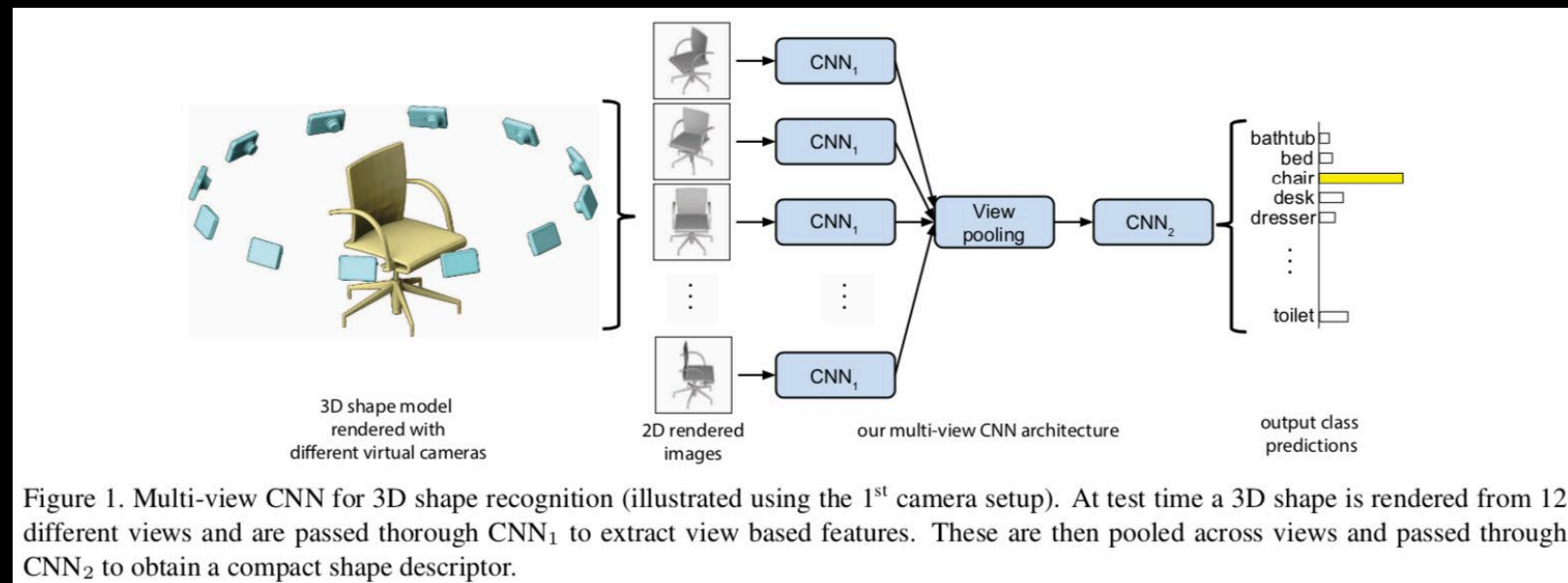
Special Projections

Geometry Images

Panorama

Spherical

- Learning methods for generic object recognition with invariance to pose and lighting [Lecun and Huang **2004**]
- „MVCNN“ Multi-view convolutional neural networks for 3D shape recognition [Su et al. 2015]



- A deeper look at 3D shape classifiers [Su et al. 2018] (“*mvcnn2*”)

Multi-view image-based approaches

Multi-view (images)

Basic Architectures

Multiple Modalities

Axis-aligned Views

Learned View Grouping

Unsupervised Viewpoints Assignment

Unsupervised Representation Learning

Using Auxiliary Data

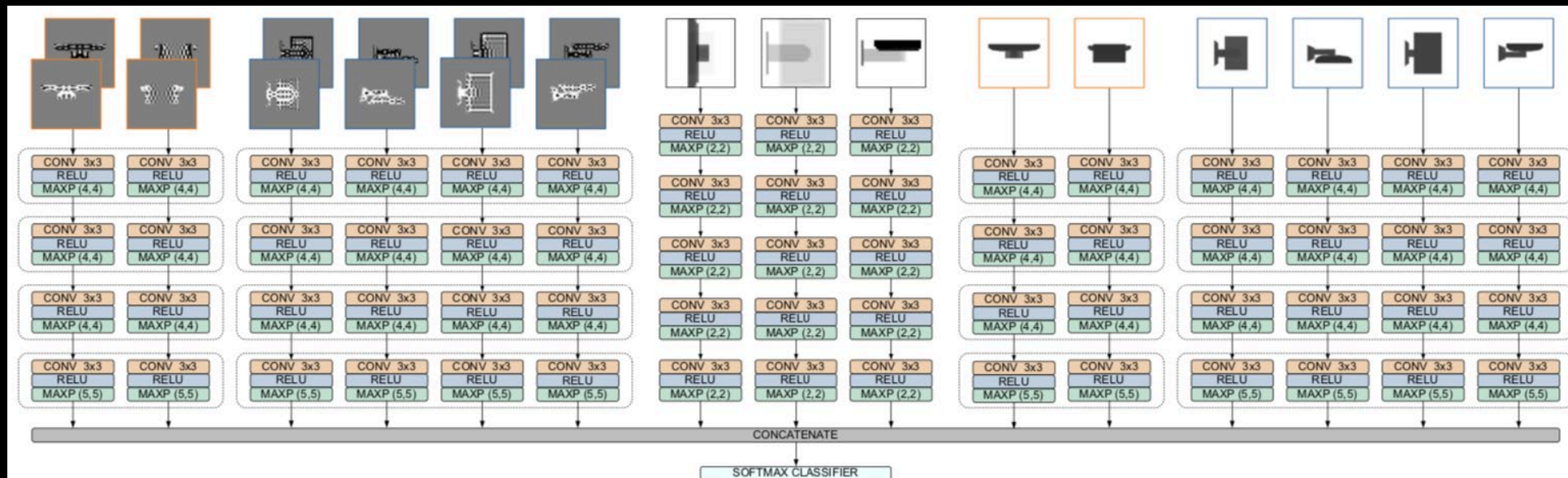
Special Projections

Geometry Images

Panorama

Spherical

- Pairwise decomposition of image sequences for active multi-view recognition [Johns et al. 2016]
- Deep learning for 3D shape classification based on volumetric density and surface approximation clues [Minto et al. 2018]



Multi-view image-based approaches

Multi-view (images)

Basic Architectures

Multiple Modalities

Axis-aligned Views

Learned View Grouping

Unsupervised Viewpoints Assignment

Unsupervised Representation Learning

Using Auxiliary Data

Special Projections

Geometry Images

Panorama

Spherical

- Deep learning for 3D shape classification from multiple depth maps [Zanuttigh and Minto 2017]
- Learning 3D shapes as multi-layered height-maps using 2D convolutional networks [Sarkar et al. 2018]

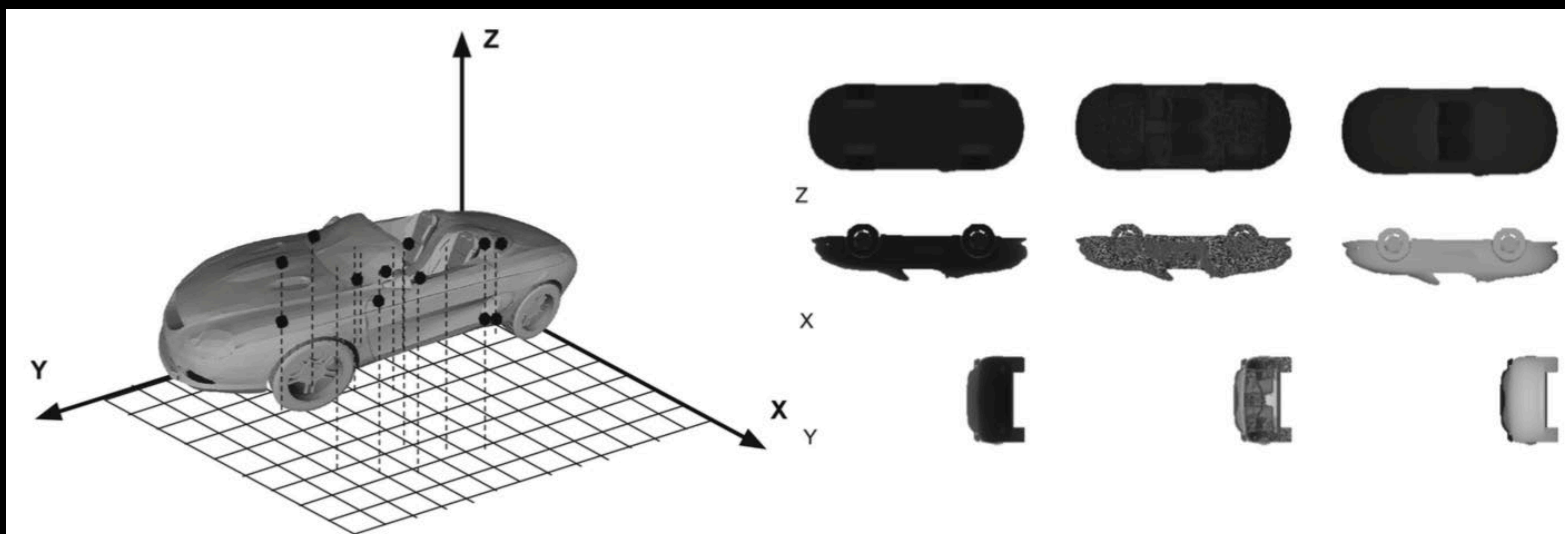


Fig. 1. (Left) Multi-layered height-map descriptors for a shape with the view along Z. (Right) Visualization of the corresponding descriptor with $k = 3$ from 3 different views of X, Y and Z.

Multi-view image-based approaches

Multi-view (images)

Basic Architectures

Multiple Modalities

Axis-aligned Views

Learned View Grouping

Unsupervised Viewpoints Assignment

Unsupervised Representation Learning

Using Auxiliary Data

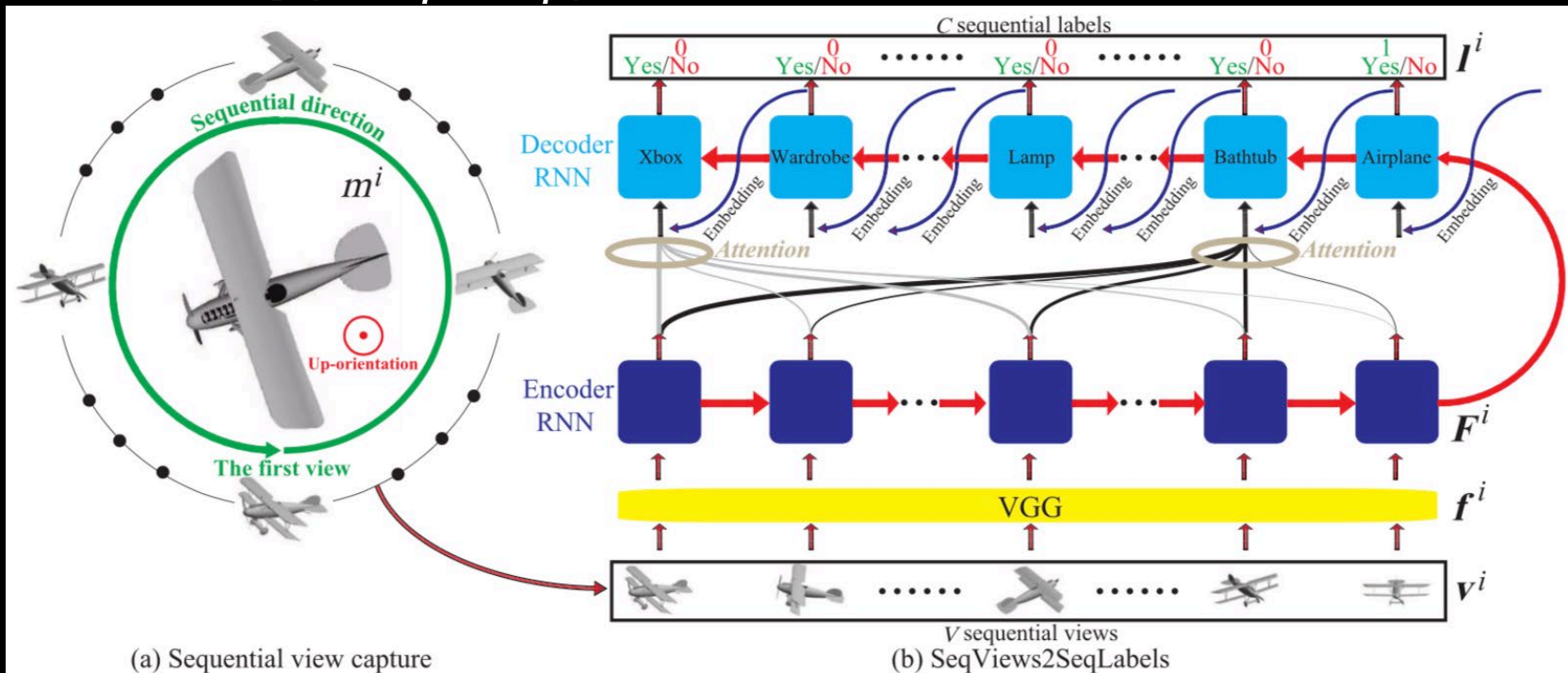
Special Projections

Geometry Images

Panorama

Spherical

- GVCNN: Group-view convolutional neural networks for 3D shape recognition [Feng et al. 2018]
- SeqViews2SeqLabels: Learning 3D global features via aggregating sequential views by RNN with attention [Zhizhong et al. 2018] (“seq2seq”)



Multi-view image-based approaches

Multi-view (images)

Basic Architectures

Multiple Modalities

Axis-aligned Views

Learned View Grouping

Unsupervised Viewpoints Assignment

Unsupervised Representation Learning

Using Auxiliary Data

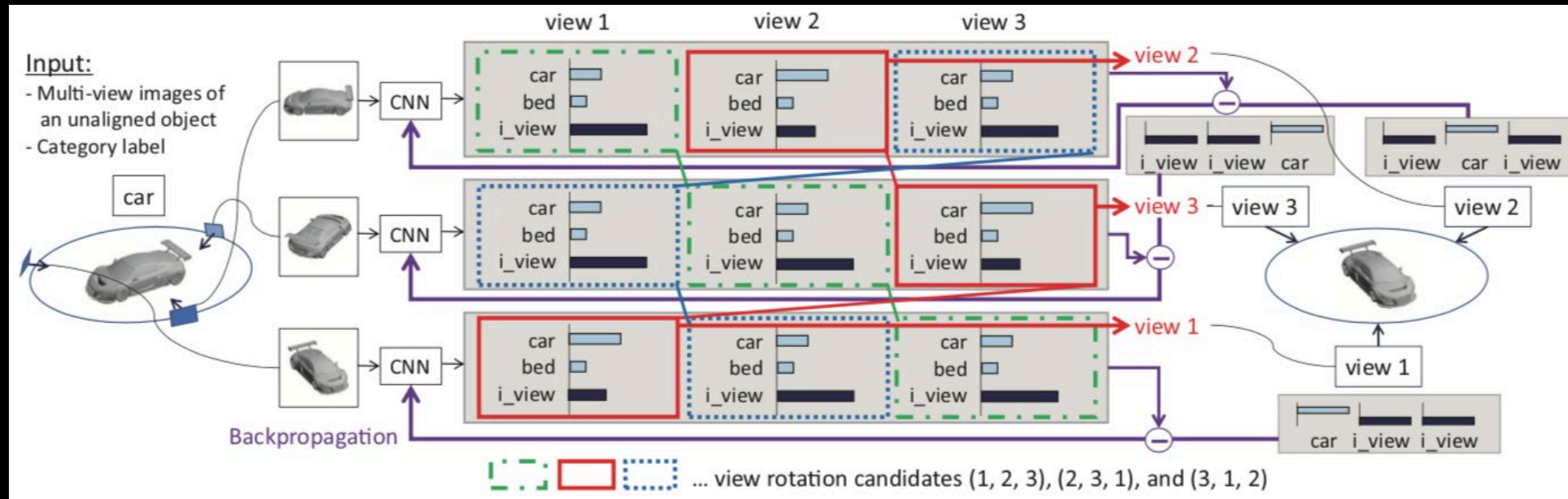
Special Projections

Geometry Images

Panorama

Spherical

- RotationNet: Joint object categorization and pose estimation using multiviews from unsupervised viewpoints [Kanezaki et al. 2018] (*"rotnet"*)



Multi-view image-based approaches

Multi-view (images)	
Basic Architectures	
Multiple Modalities	
Axis-aligned Views	
Learned View Grouping	
Unsupervised Viewpoints Assignment	
Unsupervised Representation Learning	
Using Auxiliary Data	
Special Projections	Geometry Images
	Panorama
	Spherical

Multi-view image-based approaches

Multi-view (images)	
Basic Architectures	
Multiple Modalities	
Axis-aligned Views	
Learned View Grouping	
Unsupervised Viewpoints Assignment	
Unsupervised Representation Learning	
Using Auxiliary Data	
Special Projections	Geometry Images
	Panorama
	Spherical

Multi-view image-based approaches

Multi-view (images)

Basic Architectures

Multiple Modalities

Axis-aligned Views

Learned View Grouping

Unsupervised
Viewpoints Assignment

Unsupervised
Representation
Learning

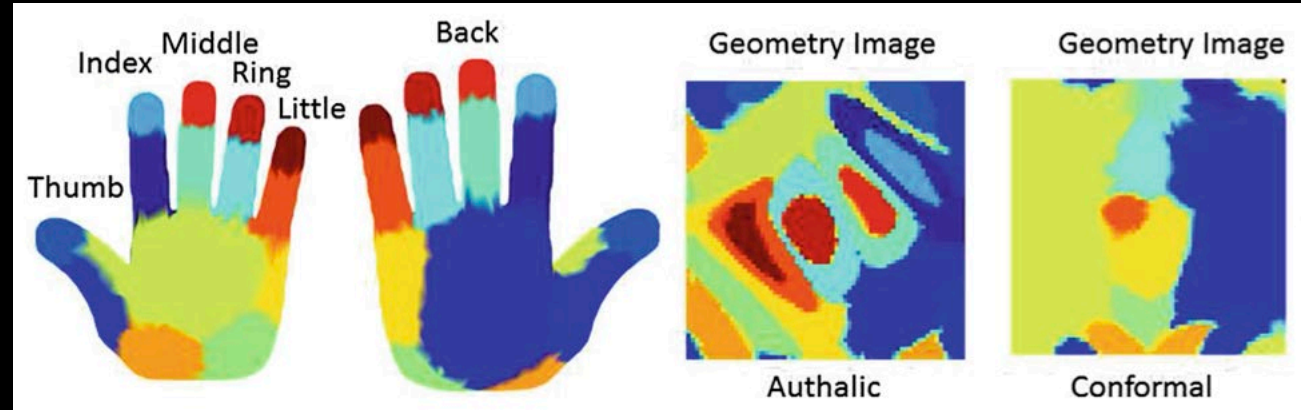
Using Auxiliary Data

Special
Projections

Geometry
Images

Panorama

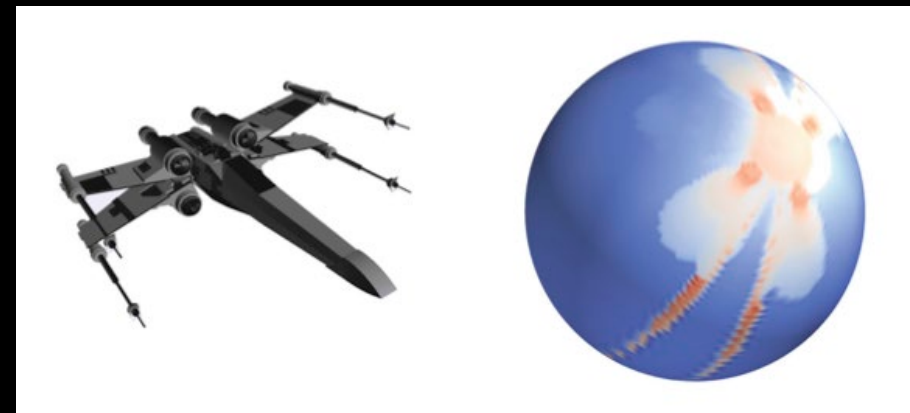
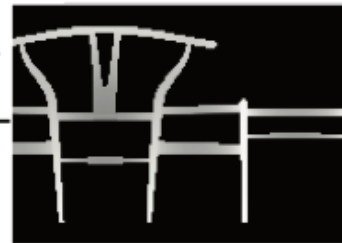
Spherical



Input 3D shape



Panoramic view
construction

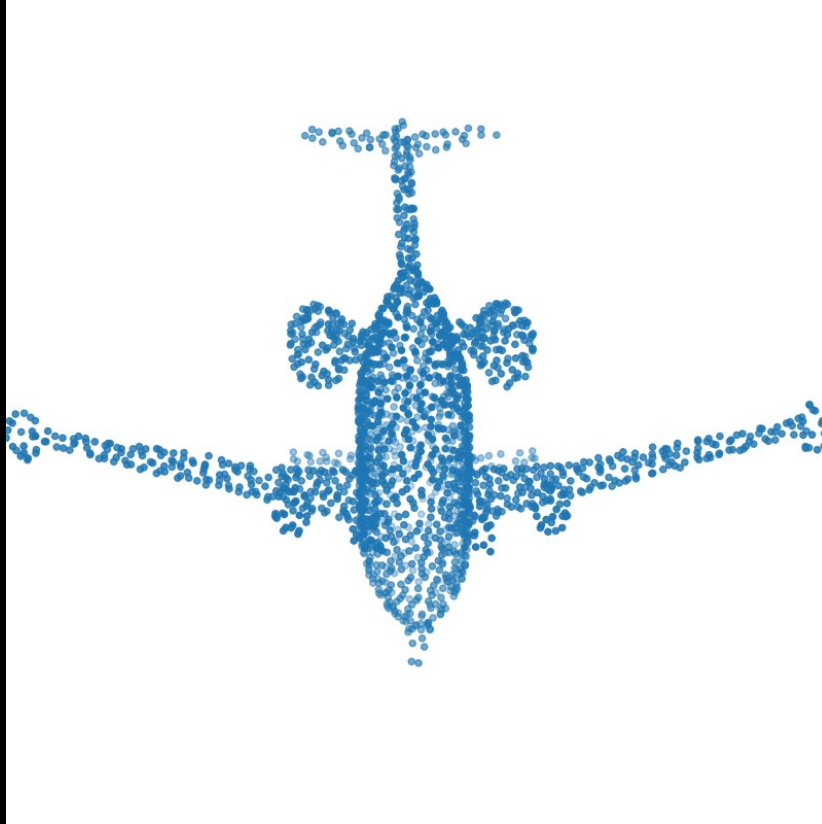


Multi-view image-based approaches

- Using 2D CNNs
- Possible re-use of
 - *architecture* ideas from 2D CNNs (residual blocks)
 - pre-trained weights for feature extractors (transfer learning)
- View aggregation: max./avg. pooling, RNN

Point cloud-based approaches

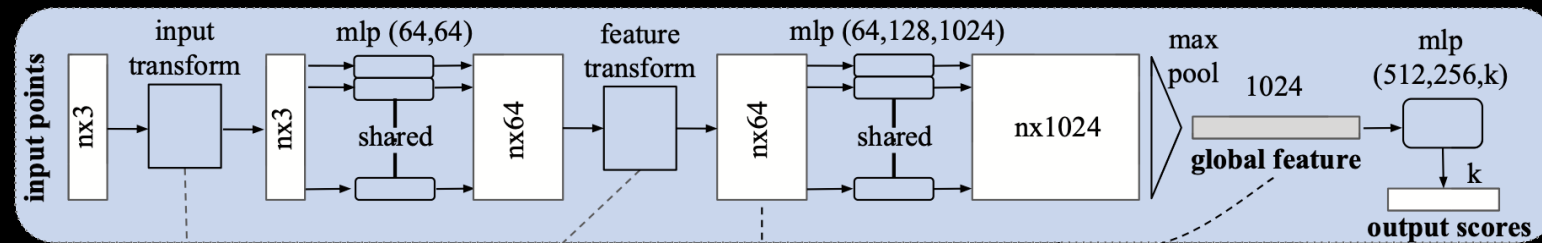
Point cloud	
Symmetric Operation on Points	
Hierarchical Feature Extraction	
Convolution on Neighborhood Graph	
Defining Convolution on Points	"Grid" Around the Query Point
	Continuous Convolution
Sequential Point Cloud Processing	Using Attention Mechanism
	Encoding Locality into the Order of Points
Unsupervised Learning of Shapes	



Point cloud-based approaches

Point cloud	
Symmetric Operation on Points	
Hierarchical Feature Extraction	
Convolution on Neighborhood Graph	
Defining Convolution on Points	"Grid" Around the Query Point
	Continuous Convolution
Sequential Point Cloud Processing	Using Attention Mechanism
	Encoding Locality into the Order of Points
Unsupervised Learning of Shapes	

- **PointNet:** Deep learning on point sets for 3D classification and segmentation [Qi et al. 2017] (*"pointnet"*)

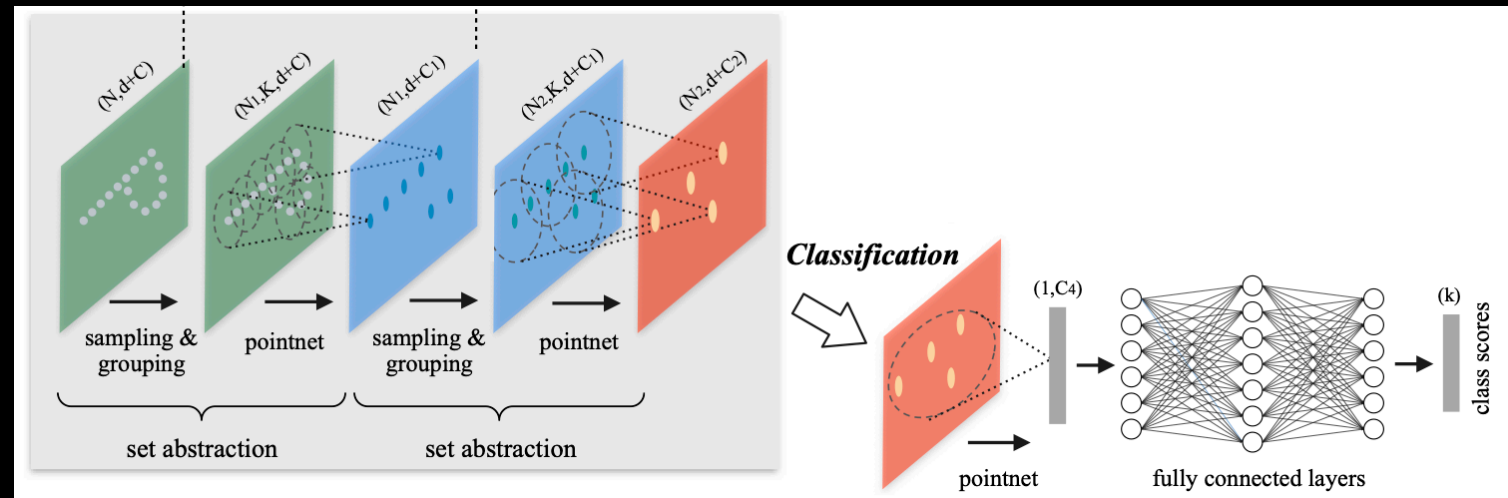


- Deep Sets [Zaheer et al. 2017]

Point cloud-based approaches

Point cloud	
Symmetric Operation on Points	
Hierarchical Feature Extraction	
Convolution on Neighborhood Graph	
Defining Convolution on Points	"Grid" Around the Query Point
	Continuous Convolution
Sequential Point Cloud Processing	Using Attention Mechanism
	Encoding Locality into the Order of Points
Unsupervised Learning of Shapes	

- **PointNet++** [Qi et al. 2017] ("*pointnet2*")



- **Escape from cells: Deep Kd-Networks** for the recognition of 3D point cloud models [Klokov and Lempitsky 2017] ("*kdnet*")

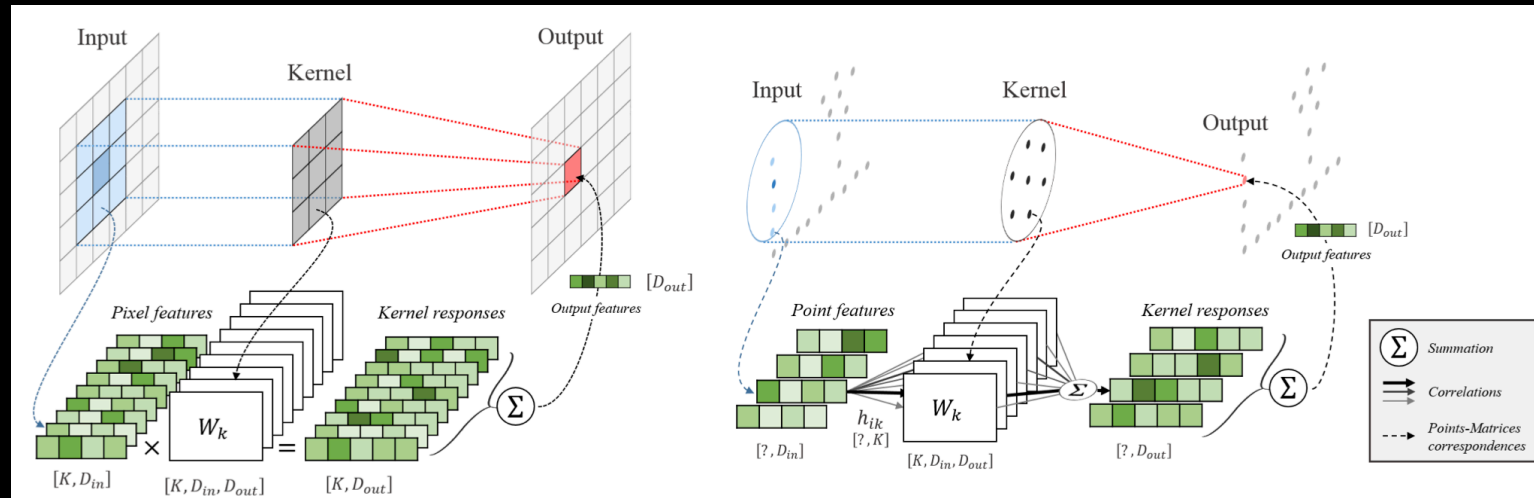
Point cloud-based approaches

Point cloud	
Symmetric Operation on Points	
Hierarchical Feature Extraction	
Convolution on Neighborhood Graph	
Defining Convolution on Points	"Grid" Around the Query Point
	Continuous Convolution
Sequential Point Cloud Processing	Using Attention Mechanism
	Encoding Locality into the Order of Points
Unsupervised Learning of Shapes	

Point cloud-based approaches

Point cloud	
Symmetric Operation on Points	
Hierarchical Feature Extraction	
Convolution on Neighborhood Graph	
Defining Convolution on Points	"Grid" Around the Query Point
	Continuous Convolution
Sequential Point Cloud Processing	Using Attention Mechanism
	Encoding Locality into the Order of Points
Unsupervised Learning of Shapes	

- Mining point cloud local structures by kernel correlation and graph pooling [Shen et al. 2018]
- KPConv [Thomas et al. 2019]



Point cloud-based approaches

Point cloud	
Symmetric Operation on Points	
Hierarchical Feature Extraction	
Convolution on Neighborhood Graph	
Defining Convolution on Points	"Grid" Around the Query Point
	Continuous Convolution
Sequential Point Cloud Processing	Using Attention Mechanism
	Encoding Locality into the Order of Points
Unsupervised Learning of Shapes	

- Multiresolution tree networks for 3D point cloud processing [Gadelha et al. 2018]

Point cloud-based approaches

Point cloud	
Symmetric Operation on Points	
Hierarchical Feature Extraction	
Convolution on Neighborhood Graph	
Defining Convolution on Points	"Grid" Around the Query Point
	Continuous Convolution
Sequential Point Cloud Processing	Using Attention Mechanism
	Encoding Locality into the Order of Points
Unsupervised Learning of Shapes	

Point cloud-based approaches

- Processing points:
 - aggregate using a symmetric function
 - *re-invent* convolution and pooling
 - group/cluster points → hierarchy
 - (nearest neighbor graph)
 - reorder points and use 1D convolution or attention

Surface shape-based approaches

Surface shape

Manifold-based
Convolution

Graph-based
Convolution

Native Mesh-
based Approaches



Surface shape-based approaches

Surface shape

Manifold-based
Convolution

Graph-based
Convolution

Native Mesh-
based Approaches

- **MeshNet**: mesh neural network for 3D shape representation [Feng et al. 2019]
- **MeshCNN**: a network with an edge [Hanocka et al. 2019]

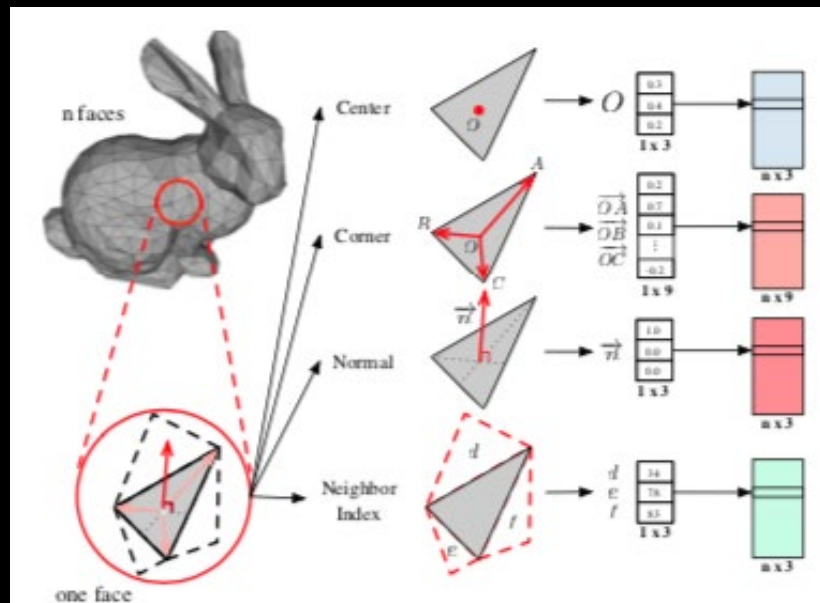


Figure 3: Initial values of each face. There are four types of initial values, divided into two parts: center, corner and normal are the face information, and neighbor index is the neighbor information.

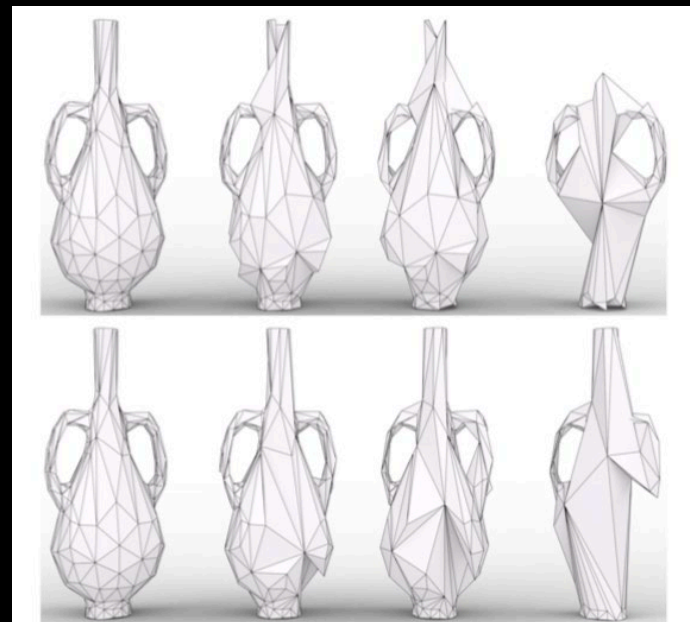
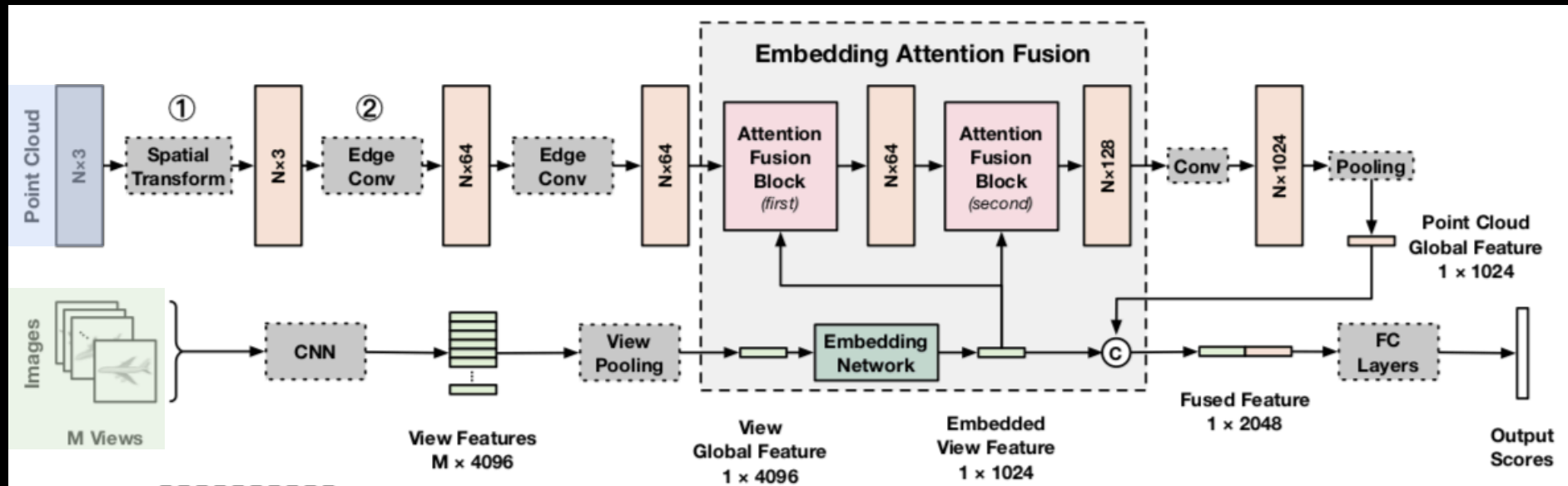


Fig. 1. Mesh pooling operates on irregular structures and adapts spatially to the task. Unlike geometric simplification (removes edges with a minimal geometric distortion), mesh pooling delegates which edges to collapse to the network. Top row: MeshCNN trained to classify whether a vase has a handle, bottom row: trained on whether there is a neck (top-piece).

Hybrid approaches

Hybrid
Ensembling
Descriptor Merging

- **FusionNet** [Hegde and Zadeh 2016]
- **PVNet**: A joint convolutional network of point cloud and multi-view for 3D shape recognition [You et al. 2018]



Survey conclusion

- Different representations, each with many possible approaches
- Taxonomy

Volumetric grid	Multi-view (images)	Point cloud		Surface shape	Hybrid
Basic Architectures	Basic Architectures	Symmetric Operation on Points		Manifold-based Convolution	Ensembling
Voxel CNN with Residual Connections	Multiple Modalities	Hierarchical Feature Extraction		Graph-based Convolution	Descriptor Merging
Auxiliary Task	Axis-aligned Views	Convolution on Neighborhood Graph		Native Mesh-based Approaches	
Network Architecture Optimization	Learned View Grouping	Defining Convolution on Points	"Grid" Around the Query Point		
Octree-represented Voxel Grid	Unsupervised Viewpoints Assignment		Continuous Convolution		
Unsupervised Representation Learning	Unsupervised Representation Learning	Sequential Point Cloud Processing	Using Attention Mechanism		
Non-convolutional Approaches	Using Auxiliary Data		Encoding Locality into the Order of Points		
Conversion from a Point Cloud	Special Projections	Unsupervised Learning of Shapes			
	Geometry Images				
	Panorama				
	Spherical				

- Paper:
 - One paragraph describing the architecture and results for each surveyed approach
 - Evaluation of selected approaches

Evaluation

- Data conversion: ModelNet/ShapeNet to network-native (voxel grid, multi-view images, point cloud)
- Runtime: **Train and evaluate**
 - Requirements:
 - x86_64 CPU, NVIDIA GPU
 - Linux + Docker + NVIDIA Container Toolkit
 - SW dependencies provided in Docker containers
- Visualization:
 - Per-network accuracy and loss on train/test set
 - Accuracy comparison plot

Results

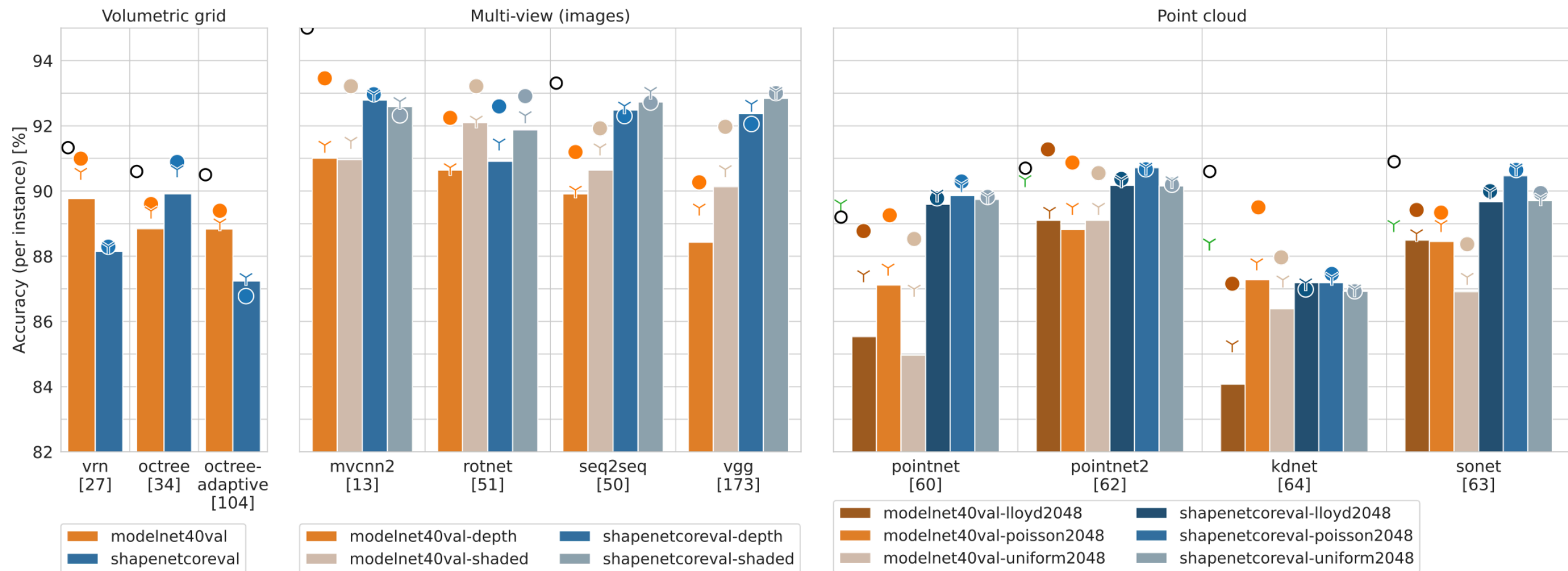


Fig. 6. The measured accuracies on different datasets. The bars show the test set accuracy at the epoch with the best validation accuracy and same-colored Υ and \bullet mark the highest achieved accuracy on the test and validation subsets on the same dataset, respectively. \circ marks accuracies reported on ModelNet40 and Υ marks the replicated test set accuracy on the point clouds provided by Qi et al. [61] using 1024 points (*qi1024*).

Results overview

- Best achieved accuracy (average over all testsets)
 1. multi-view images:
 - *mv cnn2* 91.83% \pm 0.98 pp
 - *rotnet* 91.38% \pm 0.71 pp (faster, smaller)
 2. point cloud:
 - *pointnet2* 89.67% \pm 0.76 pp
 3. voxel grid:
 - *octree* 89.37% \pm 0.75 pp

Speed, memory consumption

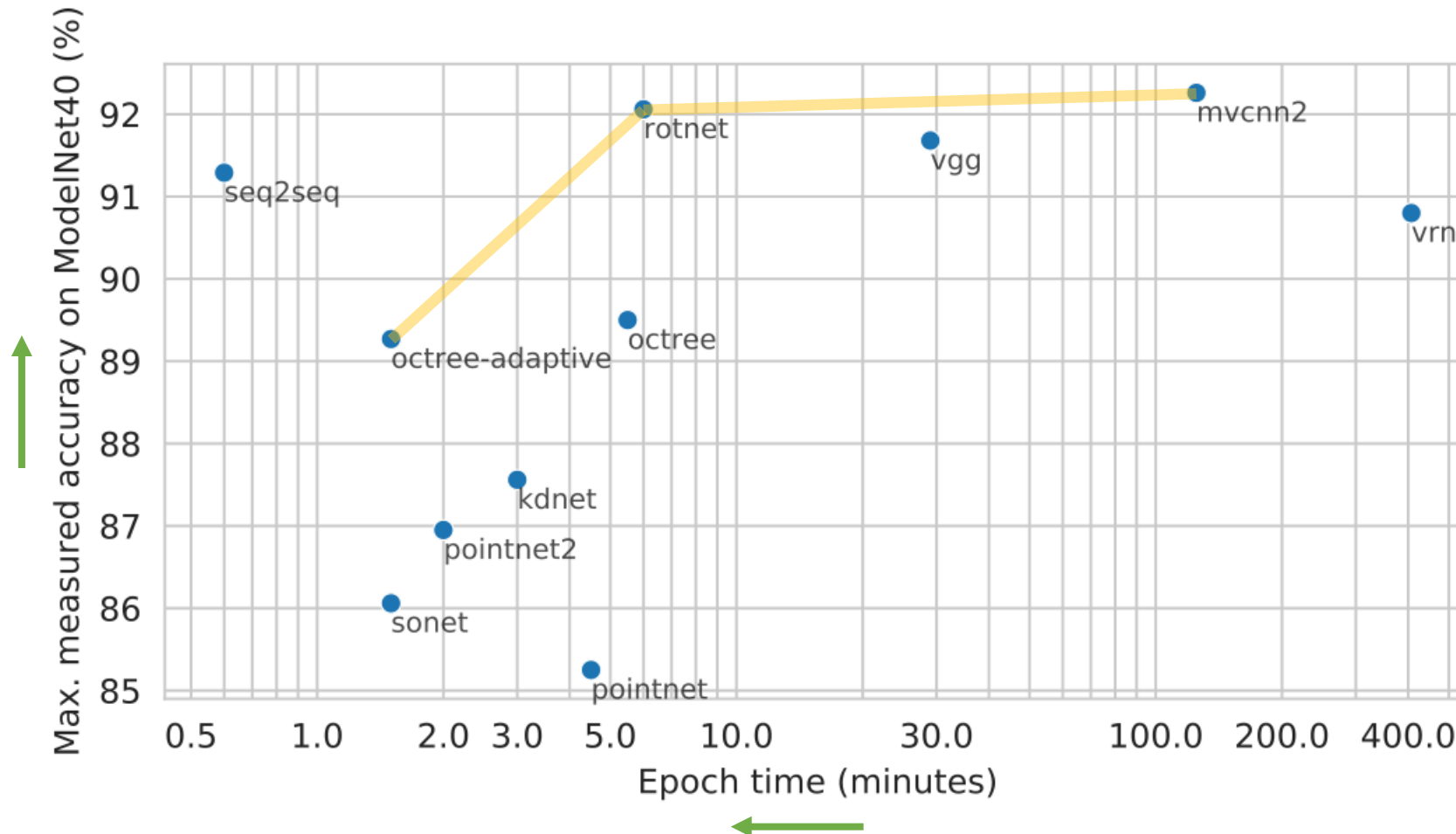


Fig. 4. Relationship of the achieved accuracy and training time.

Speed, **memory consumption**

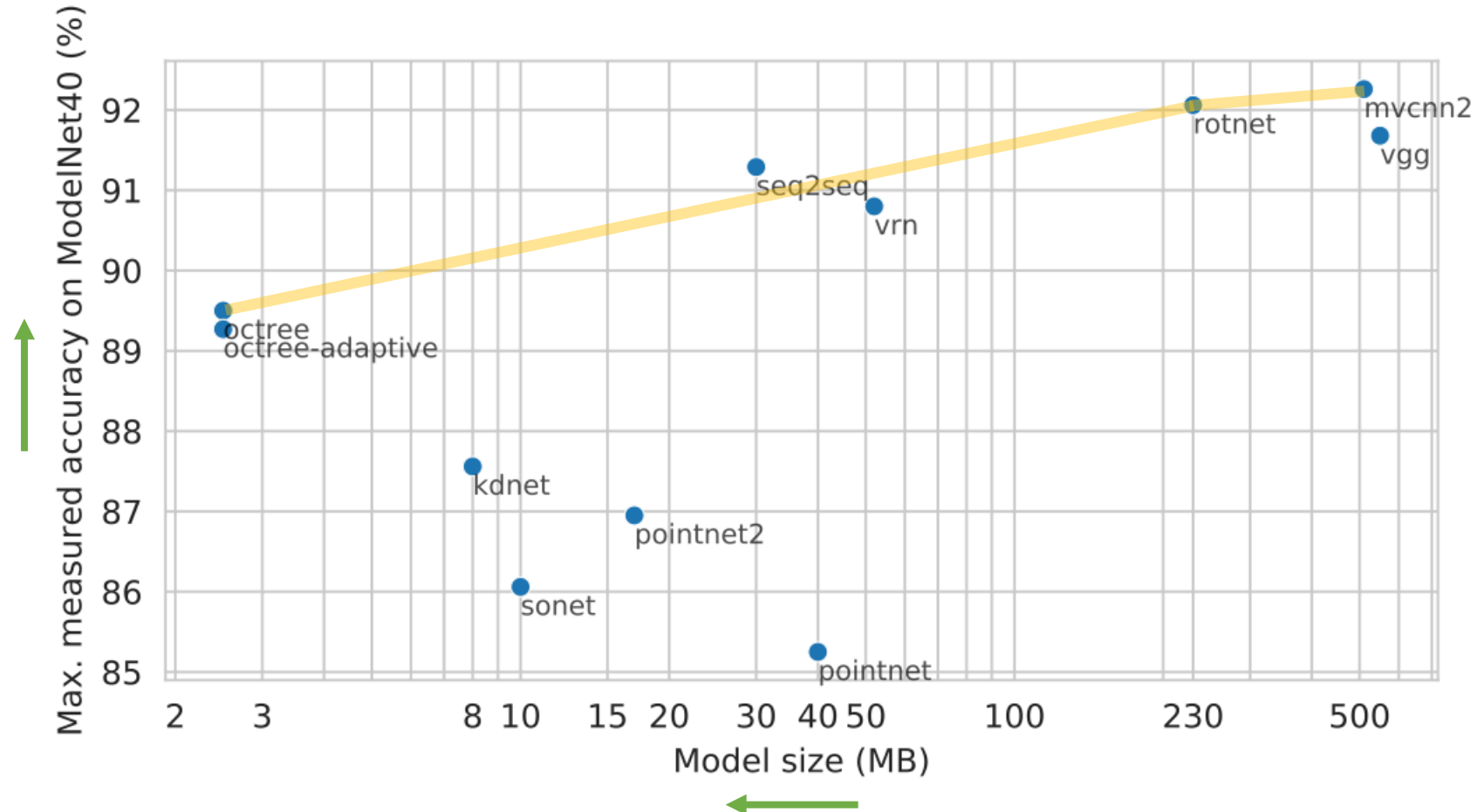


Fig. 5. Relationship of the achieved accuracy and model size.

Conclusions

- Multi-view approaches: best accuracy at the cost of size and speed
- Octree-based *volumes*: less demanding, good accuracy
- Good *price-to-performance* ratio:
 - Volumetric: O-CNN or Adaptive O-CNN
 - Multi-view images: RotationNet
 - Point clouds: Kd-network for classifying

Survey and Evaluation of Neural 3D Shape Classification Approaches

Martin Mirbauer, Miroslav Krabec, Jaroslav Křivánek, Elena Šikudová

Survey and Evaluation of Neural 3D Shape Classification Approaches

Martin Mirbauer[✉], Miroslav Krabec, Jaroslav Křivánek[✉], Elena Šikudová[✉]

Abstract—Classification of 3D objects – the selection of a category in which each object belongs – is of great interest in the field of machine learning. Numerous researchers use deep neural networks to address this problem, altering the network architecture and representation of the 3D shape used as an input. To investigate the effectiveness of their approaches, we conduct an extensive survey of existing methods and identify common ideas by which we categorize them into a taxonomy. Second, we evaluate 11 selected classification networks on two 3D object datasets, extending the evaluation to a larger dataset on which most of the selected approaches have not been tested yet. For this, we provide a framework for converting shapes from common 3D mesh formats into formats native to each network, and for training and evaluating different classification approaches on this data. Despite being partially unable to reach the accuracies reported in the original papers, we compare the relative performance of the approaches as well as their performance when changing datasets as the only variable to provide valuable insights into performance on different kinds of data. We make our code available to simplify running training experiments with multiple neural networks with different prerequisites.

Index Terms—3D shape analysis, classification algorithms, computer graphics, convolutional neural network, deep learning, image processing, machine learning, multi-layer neural network, neural networks, object recognition.

1 INTRODUCTION

Classification and generation of 3D shapes is one of the widely researched topics in the field of artificial intelligence. It is applied in a vast number of fields such as autonomous driving [1], analysis of medical data [2] as well as various fields of computer vision and graphics [3–4]. Classification of objects in 2D images has been revolutionized by deep convolutional neural networks [5, 6] and has been shown to achieve super-human accuracy [7]. This is not yet the case for 3D shapes, perhaps because of the lack of a representation that is both expressive and easy to process by a neural network.

Numerous network architectures working with different 3D shape representations have been designed, and new ones are still being developed. However, their relative performance needs further evaluation and comparison.

As the number of published approaches increases, understanding existing approaches, finding the proper representation and approach for a given application, and following new ones becomes more difficult. Categorizing them into a taxonomy and comparing the methods which use different representations is essential to simplify orientation in the landscape of approaches.

In this work, we focus on supervised learning, specifically the classification task, which is closely related to global

feature extraction – one of the tasks in the broader context of machine understanding of shapes and scenes.

We define the classification task as follows: we are given a set of training examples $\{(x_1, y_1), \dots, (x_n, y_n)\}$, where x_i is a 3D shape representation and y_i is a numerical encoding of the corresponding label. Each shape belongs to exactly one class. A classification model is a parametric model $P(\theta) : X \rightarrow Y$, where X is a space of 3D shapes, Y is a space of labels, and θ are trainable parameters. With θ optimized to minimize a prediction error metric, $P(\theta)$ should predict the correct class label for each 3D shape from X .

The contributions of this work are:

First, we extensively survey deep learning-based 3D shape classification approaches published before October 2019 and categorize them based on common approach ideas, which provides researchers with an overview of approaches suitable for processing 3D shapes.

Second, we select several existing techniques of 3D shape classification to replicate their reported results, compare and evaluate them on publicly available CAD datasets. We provide a pipeline which simplifies evaluating quality of new classifiers and methods for converting between different shape representations. The code is available on the project's website.

1.1 Related Work

There are existing works surveying the machine learning methods which process 3D shapes, however Zelner [8], Ioannidou et al. [9] and Carvalho and von Wangenheim [10] survey publications before the year 2016, which we extend until the end of year 2020 thus including current state-of-the-art methods. Other works by Arnold et al. [11], Griffiths and Boehm [12] focus on processing scanned data (RGB-D

[✉] Manuscript received [TODO when] ... The work was supported by the Charles University Grant Agency project GAKR 866119. This work was supported by the Charles University grant SVV-260588. (Corresponding author: Martin Mirbauer)
[✉] The authors are with Charles University, Faculty of Mathematics and Physics, Prague, Czech Republic.
E-mail: {martin.mirbauer, miroslav.krabec, jaroslav.krivanecek, elena.sikudova}@mff.cuni.cz
J. Křivánek, also with Czech Academy of Sciences, Prague, Czech Republic.

1 <https://cg.mff.cuni.cz/~martin/papers/2021-survey-eval>



Webpage of the paper:
<https://cg.mff.cuni.cz/publications/survey-and-evaluation-of-neural-3d-shape-classification-approaches/>