

3D Deep Learning

What is "3D"?



Parametric Surfaces



Implicit Surfaces







3D Tasks

Tasks: 3D Classification



Instance: 010.toilet_000000079.001 Predicted label: toilet True label: toilet

Class from 3D model (e.g., obtained with Kinect Scan)

[Maturana et al. 15] & [Qi et al. 16] 3D vs Multi-view

Tasks: 3D Semantic Segmentation



1500 densely annotated 3D scans; 2.5 mio RGB-D frames

[Dai et al. 17] ScanNet

Tasks: 3D Detection / Instance Segmentation





Prof. Niessner

Generative 3D Tasks



Parametric Surfaces



Implicit Surfaces





Truncated Signed Distance Field (TSDF)



Volumetric Grids

Volumetric Grids

Volumetric Data Structures

- Occupancy grids
- Ternary grids
- Distance Fields
- Signed Distance fields



(binary) Voxel Grid

Volumetric Grids

- 3D Convolutions
 - Direct extension of 2D convolutions

- Often: occupancy grid, signed/unsigned distance fields
- Operate on regular grid structures
 - Get neighborhood structures and spatial propagation
 - Well-defined pooling and unpooling for hierarchical processing
 - Cubic growth in memory

Object Classification /w 3DCNNs



ScanNet: Semantic Segmentation in 3D



Prof. Niessner

ScanNet: Sliding Window



Prof. Niessner

[Dai et al. 17] ScanNet

Generative Tasks: 3D Shape Completion



Prof. Niessner Works with 32 x 32 x 32 voxels...

[Dai et al. 17] CNNComplete

SurfaceNet: Stereo Reconstruction



Run on 32 x 32 x 32 blocks -> takes forever...

Prof. Niessner

[Ji et al. 17] Surface^{[N}et

ScanComplete: Fully Convolutional





[Dai et al. 18] ScanComplete

Dependent Predictions: Autoregressive Neural Networks



[Dai et al. 18] ScanComplete



[Dai et al. 18] ScanComplete

Prof. Niessner

ScanComplete: Fully Convolutional





[Dai et al. 18] ScanComplete

Conclusion so far

- Volumetric grids are easy and naturally extend 2D CNNs and other concepts
 - Encode free space
 - Encode distance fields
 - Need a lot of memory
 - Need a lot of processing time
 - But can be used sliding window or fully-conv.

Conclusion so far



Surface occupancy gets smaller with higher resolutions



Volumetric Hierarchies

Discriminative Tasks

Structure is known in advance!

State of the art is somewhere here...





(b) Accuracy (a) Layer 1: 32° (b) Layer 2: 16° (c) Layer 3: 8° <u>OctNet: Learning Deep 3D Representations at High Resolutions</u> (CVPR 2017)

Prof. Niess QrCNN: Octree-based Convolutional Neural Networks for 3D Shape Analysis (SIG17)23

Generative Tasks

Need to infer structure!



Octree Generating Networks: Efficient Convolutional Architectures for High-resolution Outputs OctNetFusion: Learning Depth Fusion from Data (that one not end to end) 24

Conclusion so far

- Hierarchies
 - are great for reducing memory and runtime
 - Comes at a performance hit
 - Easier for discriminative tasks when structure is known



Multi-view

Multiple Views: Classification

- RGB images from fixed views around object:
 - view pooling for classification (only RGB; no spatial corr.)



Prof. Niessner

Multi-view Convolutional Neural Networks for 3D Shape Recognition

Multiple Views: Segmentation



3D Shape Segmentation with Projective Convolutional Networks

This one is interesting in a sense that it does 3D shape segmentation (only on CAD models) But it uses multi-view and has a spatial correlation on top of the mesh surface Prof. Niessner

Fun thing...



#Views	Accuracy (class)	Accuracy (instance)
-	68.2	-
-	75.5	-
12	84.8	-
80	90.1	-
20	87.2	90.5
20	89.7	92.0
20	91.4	93.8
	#Views 12 80 20 20 20 20	#Views Accuracy (class) - 68.2 - 75.5 12 84.8 80 90.1 20 87.2 20 89.7 20 91.4

Table 3. Comparison of multi-view based methods. Numbers reported are classification accuracy (class average and instance average) on ModelNet40.

Figure 1. 3D shape representations.

Volumetric and Multi-View CNNs for Object Classification on 3D Data



Hybrid: Volumetric + Multi-view

2D + 3D Semantic Segmentation

	avg class accuracy
geometry only	54.4
geometry + voxel colors	55.9

Resolution Mismatch!

[Dai & Niessner 18] 3ÅMV

Prof. Niessner



Prof. Niessner

[Dai & Niessner 18] 3ĐMV





Prof. Niessner

[Dai & Niessner 18] 3ĐMV

	avg class accuracy
color only	58.2
geometry only	54.4

[Dai & Niessner 18] 3ĎMV

Prof. Niessner

	avg class accuracy
color only	58.2
geometry only	54.4
color+geometry	75.0


	avg class accuracy		
geometry only	54.4		
color+geometry (1 views)	70.1		
color+geometry (3 views)	73.0		
color+geometry (5 views)	75.0		

3D Volumetric + Multi-view

3D Volumetric + Multi-view

	wall	floor	cab		s	bath	other	avg
2d only (1 view)	37.1	39.1	26.7	_	2	36.3	20.4	27.1
2d only (3 views)	58.6	62.5	40.8		7	61.5	34.3	44.2
Ours (no geo input)	76.2	92.9	59.3		С	80.8	9.3	58.2
Ours (3d geo only)	60.4	95.0	54.4		3	87.0	20.6	54.4
Ours (3d geo+voxel color)	58.8	94.7	55.5		4	85.4	20.5	55.9
Ours (1 view, fixed 2d)	77.3	96.8	70.0	•••	3	87.0	58.5	69.1
Ours (1 view)	70.7	96.8	61.4		5	81.6	51.7	70.1
Ours (3 view, fixed 2d)	81.1	96.4	58.0		1	92.5	60.7	72.8
Ours (3 view)	75.2	97.1	66.4		1	89.9	57.2	73.0
Ours (5 view, fixed 2d)	77.3	95.7	68.9	1	7	93.5	59.6	74.5
Ours (5 view)	73.9	95.6	69.9		3	94.7	58.5	75.0

[Dai & Niessner 18] 3DMV

Conclusion so far

- Hybrid:
 - Nice way to combine color and geometry
 - Great performance (best so far for segmentation)
 - End-to-end helps less than we hoped for
 - Could be faster...



Point Clouds

Deep Learning on Point Clouds: PointNet



[Qi et al. 17] Point#Net

Prof. Niessner

Deep Learning on Point Clouds: PointNet

Classification Network mlp(64,64)mlp (64, 128, 1024) input feature max transform input points transform \ pool 1024 nx64 nx3 nx64 nx3 nx1024 shared shared global feature point features 64x64 3x3 T-Net T-Net \transform transform nx128 n x 1088 shared shared matrix matrix multiply multiply mlp (512,256,128) mlp (128,m)

Segmentation Network

[Qi et al. 17] PointNet

mlp

(512,256,k)

k

output scores

output scores.

nxm

Prof. Niessner

PointNet++

Main idea

- Learn hierarchical representation of point cloud
- Apply multiple (simplified) PointNets at different locations and scales
- Each Scale: Furthest-Point Sampling -> Query Ball Grouping -> PointNet
- Multi-scale or Multi-resolution grouping for sampling density robustness

Evaluations: Classification, Part-Segmentation, Scene-Segmentation





[Qi et al. 17] PointNet44

Point Convolutions

Main idea

- Transform points to continuous R3 representation (RBFs)
- Convolve in R3
- Restrict results to points

Uses Gaussian RBF representation.

Boils down to computing fixed weights for convolution.

Don't use real data as far as I know!

<u>Point Convolutional NN by Extension Operators</u> Matan Atzmon, Haggai Maron, Yaron Lipman (SIGGRAPH 2018)



Point Transformer



Figure 3. Point transformer networks for semantic segmentation (top) and classification (bottom).



https://arxiv.org/pdf/2012.09164

Figure 4. Detailed structure design for each module.

[Zhao et al. 21] Point Transformer 4

Conclusions so far

- PointNet variants:
 - Train super fast (also testing)
 - Can cover large spaces in one shot
 - Cannot represent free space
 - Performance (mostly) worse than pure volumetric
 - Still lots of ongoing research!

Point Sets (local)

RBF

Point Convolutional NN by Extension Operators (SIGGRAPH 2018) Tangent Convolutions for Dense Prediction in 3D (CVPR 2018)

Nearest point neighborhoods

Dynamic edge-conditioned filters in convolutional neural networks on graphs (CVPR17) 3D Graph Neural Networks for RGBD Semantic Segmentation (ICCV17) PPFNet: Global context aware local features for robust 3d point matching (CVPR18) FeaStNet: Feature-Steered Graph Convolutions for 3D Shape Analysis (CVPR18)

Very interesting combination where convolutions are essentially over line segments in 3D, and where both locations and are being optimized <u>https://arxiv.org/abs/1605.06240</u>

Idea is great, performance could be a bit better (probably hard to optimize)



Sparse Convolutions

Regular, dense 3x3 Convolution -> set of actives (non-zeros) grows rapidly -> need a lot of memory -> takes a long time for feature prop.









Prof. Niessner

Submanifold Sparse Conv: -> set of active sites is unchanged -> active sites look at active neighbors (green) -> non-active sites (red) have no overhead



Submanifold Sparse Conv:

- -> disconnected components do not communicate at first
- -> although they will merge due to effect of stride, pooling, convs, etc.



from left: (i) an active point is highlighted; a convolution with stride 2 sees the green active sites (ii) and produces output (iii), 'children' of highlighted active point from (i) are highlighted; a submanifold sparse convolution sees the green active sites (iv) and produces output (v); a deconvolution operation sees the green active sites (vi) and produces output (vii).





Conv7

https://github.com/NVIDIA/MinkowskiEngine

Dimension	Name in 'torch.nn'	Use cases
1	Conv1d	Text, audio
2	Conv2d	Lines in 2D space, e.g. handwriting
3	Conv3d	Lines and surfaces in 3D space or (2+1)D space-time
4	_	Lines, etc, in (3+1)D space-time

https://github.com/facebookresearch/SparseConvNet

SparseConv Generators

SparseConv Generators



SG-NN: Sparse Convs + Self-Supervised

Input Scan

SG-NN (Ours)





Depth Frames



Depth Frames



Target Scan

Depth Frames



Target Scan

Depth Frames

Target Scan Input Scan

Depth Frames Input Scan Target Scan Self-Supervision **Unobserved Space**

[Meng et al. 24]: Latent 3D Scene Diffusion





[Meng et al. 24]: Latent 3D Scene Diffusion



[Meng et al. 24]: Latent 3D Scene Diffusion

Conclusions so far

- Spares (volumetric) Convs:
 - Implemented with spatial hash function
 - Features only around "surface"
 - Require significantly less memory
 - Allow for much higher resolutions
 - It's slower, but much higher accuracy



3D Meshes

3DMeshes

- Collection of vertices, edges, and faces
- Defines surface geometry



- Interpret meshes as graphs
- How to define a learnable operators / convs over a graph

3DMeshes

- Graph Networks
 - Message Passing
 - Graph Convolutions
 - Transformers?

MeshCNN, Spectral Graph Convolutions, Geodesic
CNNs

3D Meshes: MeshCNN

- A CNN for triangle meshes
- Define convolution and pooling operators on mesh edges
- Each edge has a feature
- Each edge has four edge neighbors (from 2 incident faces)
- Convolution filters applied to ea feature and its 4 neighbors
- Pooling by edge-collapse





Machine Learning for 3D Geometry

MeshCNN [Hanocka et al. '19]

Scan2Mesh: From Unstructured Range Scans to 3D Meshes



CVPR'19 [Dai and Niessner]: Scan2Mesh
Scan2Mesh: From Unstructured Range Scans to 3D Meshes



CVPR'19 [Dai and Niessner]: Scan2Mesh

Scan2Mesh: From Unstructured Range Scans to 3D



CVPR'19 [Dai and Niessner]: Scan2Mesh

PolyDiff: Generating 3D Meshes with Diffusion



arXiv'23 [Alliegro et al.]: PolyDiff

PolyDiff: Generating 3D Meshes with Diffusion



arXiv'23 [Alliegro et al.]: PolyDiff

Mesh Generation with Transformers



arXiv'20 [Nash et al.]: PolyGen: An Autoregressive Generative Model of 3D Meshes

Mesh Generation with Transformers



arXiv'20 [Nash et al.]: PolyGen: An Autoregressive Generative Model of 3D Meshes

GPTs for Mesh Generation



GPTs for Mesh Generation



GPTs for Mesh Generation



MeshGPT: Inferencing the Transformer



MeshGPT: Inferencing the Transformer





MeshGPT: Generating Triangle Meshes with Decoder-Only Transformers

Generative 3D Mesh Approaches

- ScanzMesh (Graph Network)
- PolyGen (Diffusion Model)
- PolyDiff (Transformer + PointerNetwork)
- MeshGPT (GPT-style Transformer)



3D Datasets

(mostly indoor)

3D Shapes (Synthetic)

- ShapeNet
 - Main dataset, everyone uses it
 - 55 classes (51.3k shapes); relatively uniform distribution but also mostly chairs; mediocre textures
- Objaverse (links only)
 - 800k 3D shapes
- Objaverse-XL (links only)
 - 10mio 3D shapes (pretty heterogeneous)

3D Scenes (Synthetic)

- 3D-FRONT: 3D Furnished Rooms with Layouts and Semantics
 - 18,797 rooms (mostly procedural geometry)
 - 7,302 furniture objects





- ScanNet
 - Kinect-style reconstructions; 1500 scenes; 2.5 mio views; semantic + instance annotations



3D Scenes

Method	Info	avg iou	bathtub	bed	bookshelf	cabinet	chair	counter	curtain	desk	door	floor	otherfurniture	picture	refrigerator	shower curtain
		•	~	~	~	~	~	~	~	~	~	~	~	~	~	~
PTv3 ScanNet		0.794 1	0.941 s	0.813 17	0.851 7	0.782 5	0.890 2	0.597 1	0.916 2	0.696 7	0.713 s	0.979 1	0.635 1	0.384 2	0.793 2	0.907 7
Xiaoyang Wu, Li Jiang, Peng-	Shuai Wang	, Zhijian Liu, Xihu	i Liu, Yu Qiao,	Wanli Ouyang,	Tong He, Heng	shuang Zhao:	Point Transform	ner V3: Simple	r, Faster, Stron	ger. OVPR 202	24					
PonderV2		0.785 2	0.978 1	0.800 25	0.833 21	0.788 s	0.853 15	0.545 16	0.910 5	0.713 1	0.705 4	0.979 1	0.596 s	0.390 1	0.769 11	0.832 40
Haoyi Zhu, Honghui Yang, Xia	oyang Wu, I	Di Huang, Sha Z	hang, Xianglor	ng He, Tong He	, Hengshuang Z	hao, Chunhua	a Shen, Yu Qia	o, Wanli Ouyan	g: PonderV2: I	Pave the Way f	for 3D Foundat	aion Model w	ith A Universal Pre-tra	ining Paradigm	h.	
Mix3D	Р	0.781 s	0.964 2	0.855 1	0.843 15	0.781 s	0.858 11	0.575 s	0.831 s1	0.685 13	0.714 2	0.979 1	0.594 7	0.310 26	0.801 1	0.892 15
Alexey Nekrasov, Jonas Schu	it, Or Litany,	Bastian Leibe, P	Francis Engelm	ann: Mix3D: O	ut-of-Context Da	ata Augmentati	ion for 3D Scer	nes. 3DV 2021	(Oral)							
Swin3D	Р	0.779 4	0.861 20	0.818 13	0.836 18	0.790 2	0.875 4	0.576 •	0.905 s	0.704 4	0.739 1	0.969 10	0.611 2	0.349 10	0.758 20	0.958 1 0
TTT-KD		0.773 a	0.646 as	0.818 13	0.809 33	0.774 a	0.878 s	0.581 2	0.943 1	0.687 11	0.704 •	0.978 4	0.607 5	0.336 15	0.775 a	0.912 :
Jisa Weijler, Muhammad Jeha	inzeb Mirza,	Leon Sick, Can	Ekkazan, Pedr	o Hermosilla: T	TT-KD: Test-Tim	e Training for 3	3D Semantic S	egmentation th	rough Knowled	dge Distillation	from Foundatio	n Models.				
ResLFE_HDS		0.772 s	0.939 4	0.824 s	0.854 s	0.771 🛛	0.840 29	0.564 10	0.900 s	0.686 12	0.677 11	0.961 16	0.537 29	0.348 11	0.769 11	0.903 • 0
OctFormer	Р	0.766 7	0.925 7	0.808 21	0.849 9	0.786 4	0.846 25	0.566 9	0.876 14	0.690 9	0.674 13	0.960 17	0.576 16	0.226 📾	0.753 22	0.904 s C
Peng-Shuai Wang: OctForme	r: Octree-ba	ased Transformer	s for 3D Point	Clouds, SIGGF	RAPH 2023											
PT-SpUNet-Joint		0.766 7	0.932 •	0.794 aı	0.829 23	0.751 21	0.854 13	0.540 20	0.903 7	0.630 32	0.672 14	0.963 14	0.565 20	0.357 a	0.788 s	0.900 11
Kaoyang Wu, Zhuotao Tian, >	(in Wen, Bo	hao Peng, Xihui l	Liu, Kaichengʻ	Yu, Hengshuan	g Zhao: Towards	s Large-scale (3D Representa	tion Learning v	vith Multi-datas	set Point Promp	ot Training, CVP	PR 2024				
ocuSeg+Semantic		0.764 s	0.758 57	0.796 29	0.839 17	0.746 23	0.907 1	0.562 11	0.850 23	0.680 15	0.672 14	0.978 4	0.610 s	0.335 17	0.777 6	0.819 43
CU-Hybrid Net		0.764 9	0.924 s	0.819 11	0.840 16	0.757 16	0.853 15	0.580 s	0.848 24	0.709 s	0.643 22	0.958 20	0.587 11	0.295 32	0.753 22	0.884 19 0
O-CNN	Р	0.762 11	0.924 s	0.823 7	0.844 14	0.770 10	0.852 17	0.577 4	0.847 26	0.711 2	0.640 26	0.958 20	0.592 s	0.217 71	0.762 18	0.888 16 0
eng-Shuai Wang, Yang Liu, '	Yu-Xiao Guo), Chun-Yu Sun, 3	Xin Tong: O-C	NN: Octree-ba	aed Convolution	al Neural Netv	vorks for 3D Sh	ape Analysis.	SIGGRAPH 20	17						
OA-CNN- L_ScanNet20		0.758 12	0.783 43	0.826 •	0.858 4	0.776 7	0.837 32	0.548 15	0.896 11	0.649 24	0.675 12	0.962 15	0.586 12	0.335 17	0.771 10	0.802 47 0
ConDaFormer		0.755 ₁₃	0.927 6	0.822 s	0.836 18	0.801 1	0.849 20	0.516 30	0.864 20	0.651 23	0.680 10	0.958 20	0.584 14	0.282 40	0.759 18	0.855 so 0
unhao Duan, Shanshan Zha	o, Nan Xue,	Mingming Gong	, Guisong Xia,	Dacheng Tao:	ConDaFormer :	Disassembled	d Transformer v	ith Local Struc	ture Enhancer	nent for 3D Poi	int Cloud Unde	rstanding. Ne	euripa, 2023			
PNE		0.755 13	0.786 41	0.835 4	0.834 20	0.758 14	0.849 20	0.570 <mark>s</mark>	0.836 30	0.648 25	0.668 16	0.978 4	0.581 18	0.367 s	0.683 33	0.856 28
P. Hermosilla: Point Neighbort	rood Embed	ddings.														
PointTransformerV2		0.752 15	0.742 65	0.809 20	0.872 1	0.758 14	0.860 10	0.552 13	0.891 12	0.610 39	0.687 s	0.960 17	http:	/3/41	/WW	.sear

Evaluation on Hidden Test Set on ScanNet

Xiaoyang Wu, Yixing Lao, Li Jiang, Xihui Liu, Hengshuang Zhao; Point Transformer V2: Grouped Vector Attention and Partition-based Pooling. NeurIPS 2022

3D Scenes

- SceneNN
 - ScanNetStyle (100 scenes)
- Matterport3D
 - 90 buildings
- ARKitScenes
 - 1661 scenes
 - Faro + iPhone (no labels)
- ScanNet++
 - 450 3D indoor scenes
 - Faro Scans
 - DSLR + iPhone
 - Instance Labels

Many more....

Textured 3D Mesh





Panoramas



DSLR Image



iPhone RGB-D

Important Concepts

- Regular Structures vs Irregular structures
 - Convolutions, MLPs, Transformers, etc.
- Spatial data structures:
 - Octress, Hashes, Hierarchies
- Pooling:
 - n->1 mappings (points, views, etc.)
- Architecture vs Loss formulation
 - Transformer, Diffusion

Next Lectures

Neural Scene Representations

• Neural Fields

• Neural Radiance Fields (NeRF)



Thanks and Questions