

# VMNet: Voxel-Mesh Network for Geodesic-Aware 3D Semantic Segmentation

<https://arxiv.org/pdf/2107.13824.pdf>

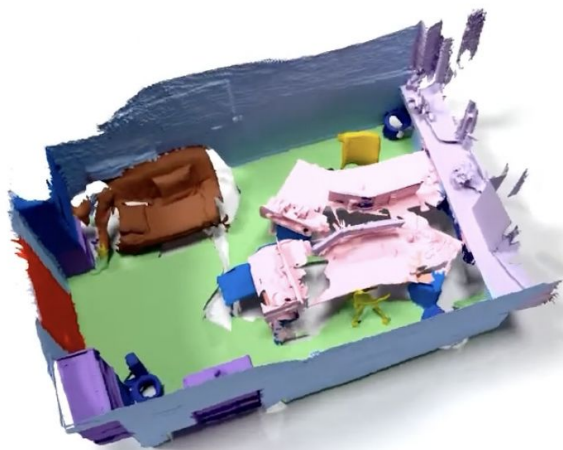
# Abstract

VoxelMesh-ConvNet over 3d meshes which combine geodesic and euclidean information from convolutions performed on actual 3d meshes and its corresponding voxel representation respectively.

They perform semantic segmentation over 3d scene from scannet dataset.



VoxelMesh-Net

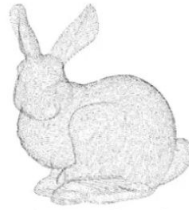


# 3D data representation



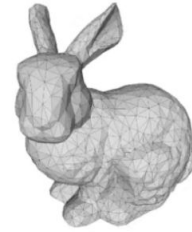
Voxelized

discrete Euclidean neighborhood



Point Cloud

continuous Euclidean neighborhood



Mesh

adjacent mesh neighbors

voxelization process omits all mesh edges and only retains Euclidean positions of mesh vertices.

Point cloud in essence are only points embedded in 3d euclidean space.

Mesh are those point cloud (extracted using lidar, etc) and then connected to form several triangles to represent the surfaces.

# Euclidean Information

voxel- based methods suffer from ambiguous features on spatially close objects and struggle with handling complex and irregular geometries due to the lack of geodesic information.

Compared to Point Cloud-

1. they lose resolution compared to point clouds, since several distinct points representing intricate structures will be binned into one voxel if they're close together.
2. voxel grids can lead to unnecessarily high memory usage compared to point clouds in sparse environments, since they actively consume memory to represent free and unknown space whereas point clouds contain only known points.

# Voxels

Imperfect approximation of an underlying 3D structure, hence useful in multi-view problems where earlier 2d images were used. Can't make enough sense to use it.

# Intention of using Voxels, as I understood it

Being different from DCM-Net. No obvious reason was stated in the paper.

Though they show better performance by using voxels but that can be attributed to the use some attention modules which work on both domains of input (voxels and mesh).



# Why both voxels and meshes

1. CNN generate similar features for voxels that are close in the Euclidean domain, even though these voxels may belong to different objects and are distant in the geodesic domain.
2. Without the geodesic information about shape surfaces, these Euclidean convolutions may struggle with learning specific object shapes. As

# Novel Idea

1. Aggregate intra-domain features
2. Fuse inter-domain features.

Hence two key components of VMNet:

- Intra- domain Attentive Aggregation Module (work on only on mesh)
- Inter-domain Attentive Fusion Module. (work on fusing voxel and mesh features)

Also adding attention to the network was also novel, earlier work like DCM-Net didn't use that.

# VM-Net Arch - Like U-Net

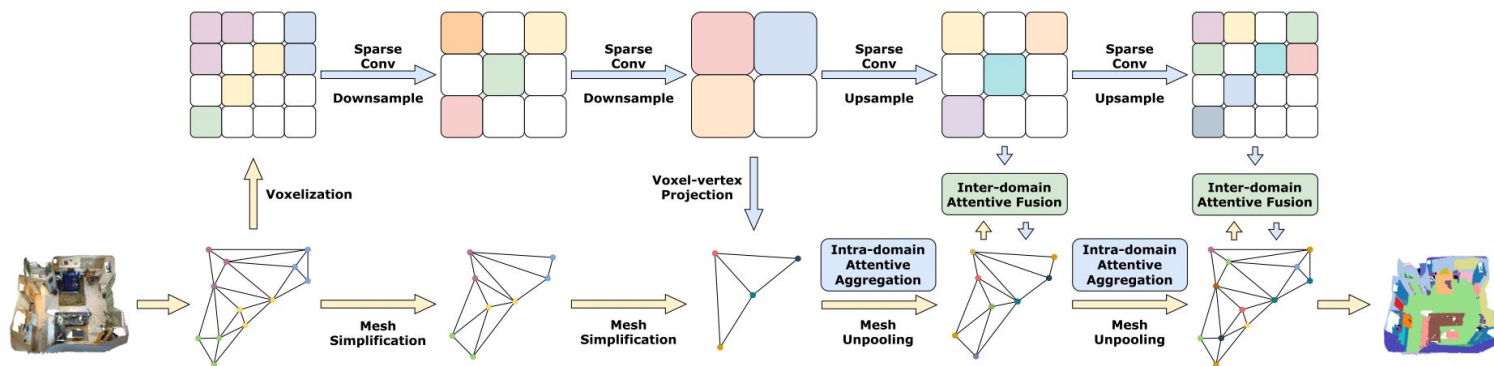


Figure 3. **Overview of Voxel-Mesh Network (VMNet).** Taking a colored mesh as input, we first rasterize it and apply voxel-based sparse convolutions to extract contextual information in the Euclidean domain. These features are then projected from voxels to vertices, and are further aggregated and fused in the geodesic domain producing distinctive per-vertex features. For simplicity, skip connections between the encoder and decoder are neglected here and only three levels of hierarchical voxel downsampling and mesh simplification are shown. The detailed network structure can be found in **Supplementary Section A**.

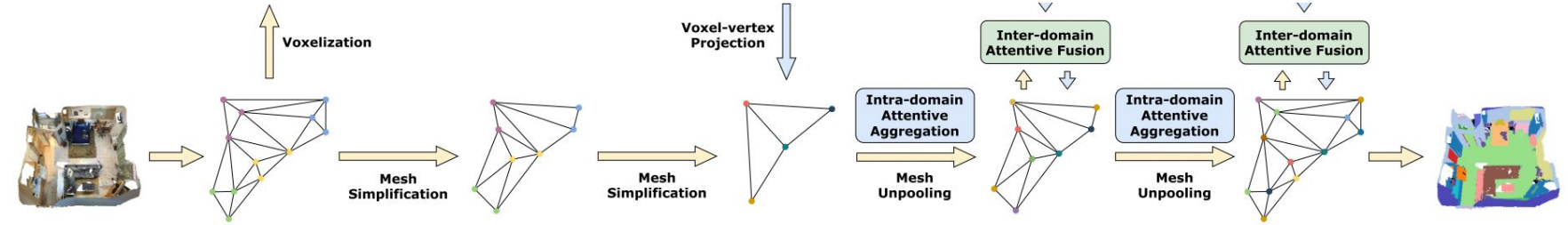
# Input to the model

Taking a mesh as input, the colored vertices are first voxelized and then fed to the Euclidean branch.

Euclidean Branch uses sparse voxel-based convolutions.

MLP in Attention layers are Graph Convs.

# Geodesic branch



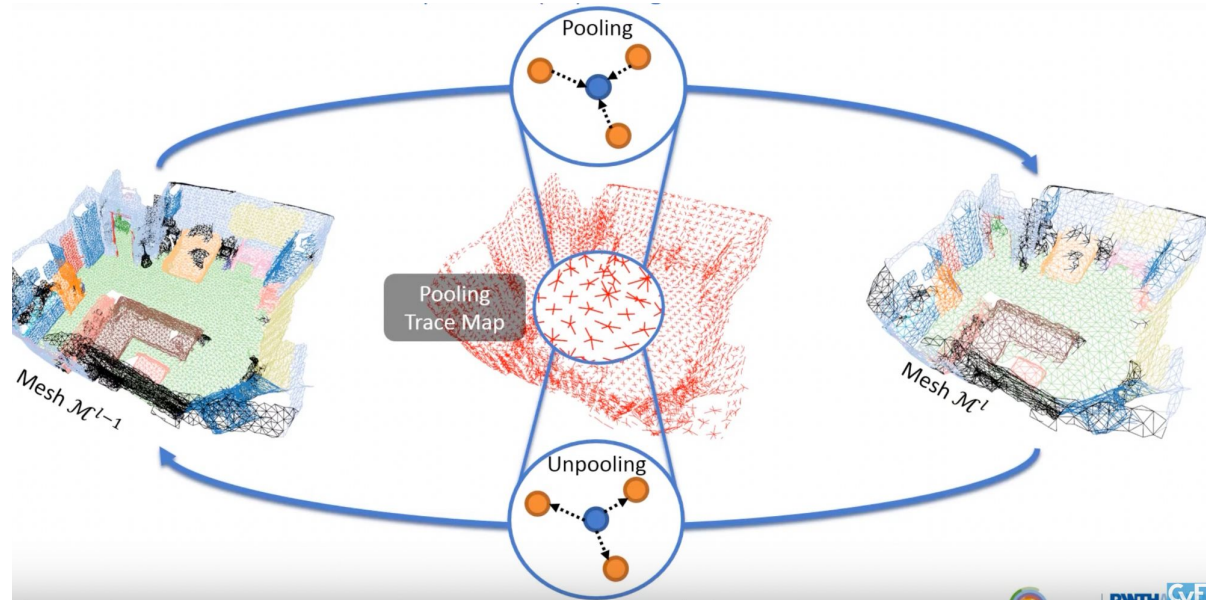
prepare hierarchy of simplified meshes ( $M_0, \dots, M_I, \dots, M_L$ )

each level of simplified mesh  $M_I$  corresponds to a the simplified mesh. They are saved for unpooling downsampling level of sparse voxels  $S_I$ . Trace maps of the simplified mesh are saved for unpooling operations between mesh levels.

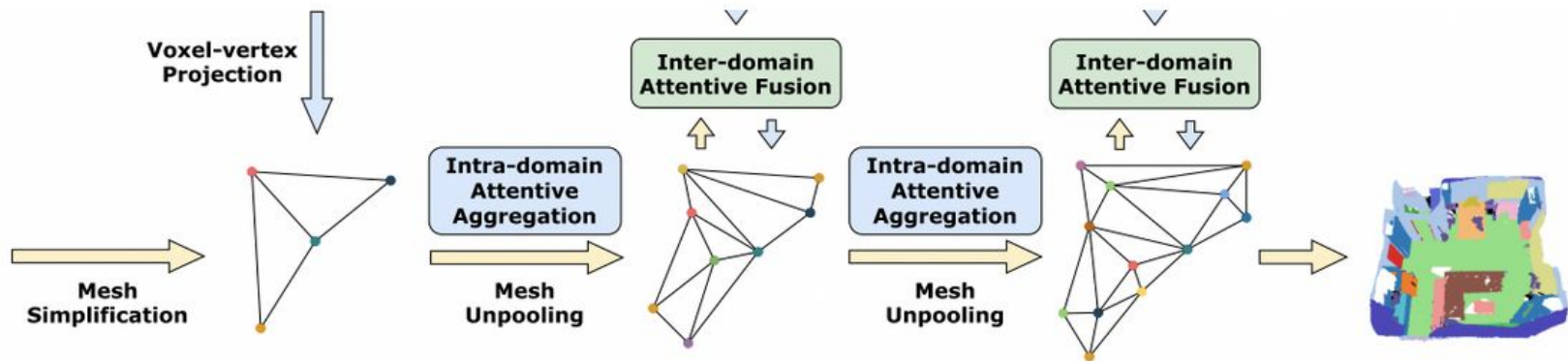
Mesh Simplification methods used: Vertex Clustering (VC) and Quadric Error Metrics (QEM)

# Trace maps

It defines which vertices are collapsed together. It is a permutation invariant function. Mesh structure is preserved in all pooling levels using these maps. And when unpooling these maps helps in developing accurate feature maps over mesh.



At the first level of the decoding process (level L), the features are projected from intra-domain attentive aggregation. The resulting geodesic features of ML are unpooled to the next level ML-1.

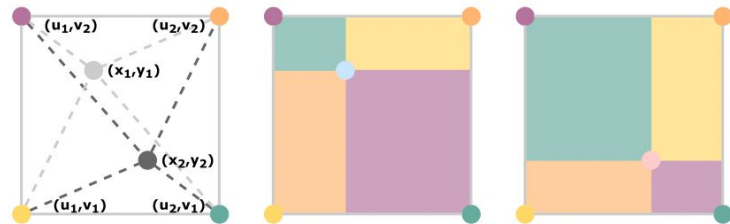


# Voxel-vertex projection

Aggregation of euclidean features present in layer I. Now transform the features of voxels SI back to vertices MI for further processing in the geodesic domain.

Trilinear projection is used. Trilinear interpolation involves only the immediate voxel neighborhood of a sample point, which are the voxels of the enclosing volume cell.

In fig. Voxels (u,v), points (x,y)



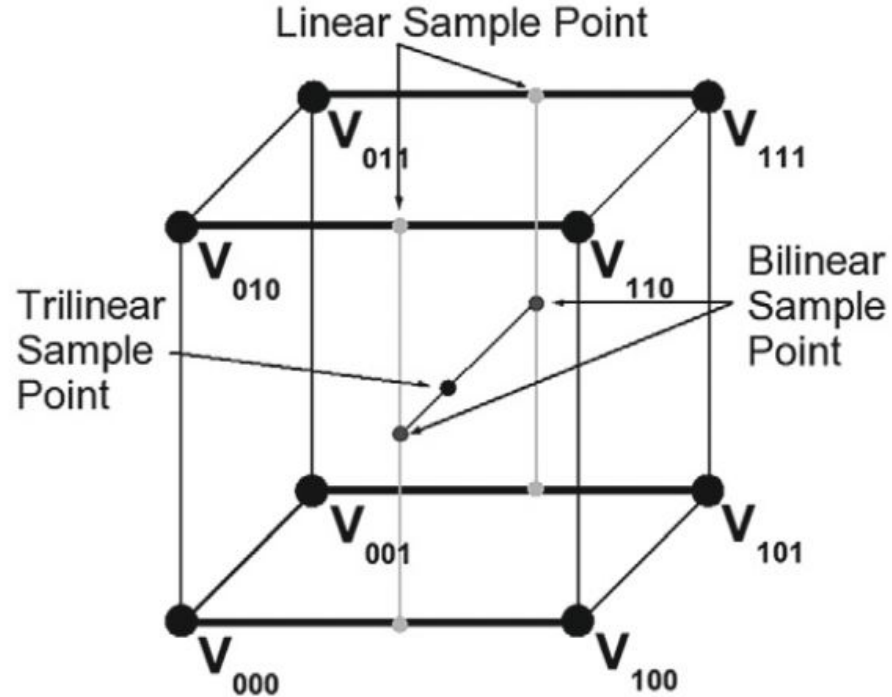
$$L(\alpha) = V_0 * (1 - \alpha) + V_1 * \alpha | 0 \leq \alpha \leq 1$$

$$\begin{aligned} B(\alpha_1, \alpha_2) &= L_0(\alpha_1) * (1 - \alpha_2) + L_1(\alpha_1) * \alpha_2 \\ &= (V_{00} * (1 - \alpha) + V_{10} * \alpha) * (1 - \alpha_2) \\ &\quad + (V_{01} * (1 - \alpha) + V_{11} * \alpha) * \alpha_2 \end{aligned}$$

$$\begin{aligned} T(x, y, z) &= B_0(x, y) * (1 - z) + B_1(x, y) * z \\ &= ((L_0(x) * (1 - y) + L_1(x) * y)) * (1 - z) \\ &\quad + ((L_2(x) * (1 - y) + L_3(x) * y)) * z \\ &= (((V_{000} * (1 - x) + V_{100} * x)) * (1 - y) \\ &\quad + ((V_{010} * (1 - x) + V_{110} * x)) * y) * (1 - z) \\ &\quad + (((V_{001} * (1 - x) + V_{101} * x)) * (1 - y) \\ &\quad + ((V_{011} * (1 - x) + V_{111} * x)) * y) * z \end{aligned}$$



<https://sci-hub.se/https://doi.org/10.1016/B978-0-12-415873-3.00002-X>

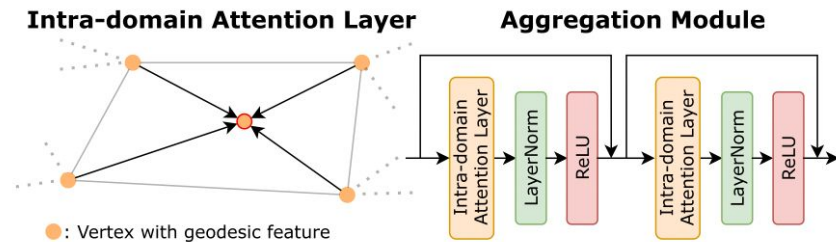


# Intra-domain Attention Layer

$N_i$  is the one-ring neighborhood of vertex  $V_i$ . The functions  $\rho^{intra}$ ,  $\alpha^{intra}$ ,  $\phi^{intra}$ , and  $\psi^{intra}$  are vertex-wise feature transformations implemented by MLP,  $\omega_{ij}$  is the attention coefficient, and  $d$  is the size of output feature channels.

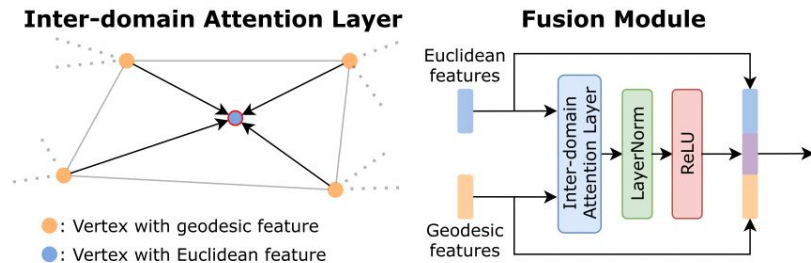
$$f_i'^{geo} = \rho_k^{intra}(f_i^{geo}) + \sum_{j \in N_i} \omega_{ij} \alpha_k^{intra}(f_j^{geo}),$$

$$\omega_{ij} = \text{softmax}\left(\frac{\phi_k^{intra}(f_i^{geo})^T \psi_k^{intra}(f_j^{geo})}{\sqrt{d}}\right),$$



# Inter-domain Attentive Fusion Module

where  $N_i$  is the same one-ring neighborhood of vertex  $V_i$  as the one used for intra-domain aggregation



$$f_i^{fuse} = \rho_k^{inter}(f_i^{euc}) + \sum_{j \in N_i} \omega_{ij} \alpha_k^{inter}(f_j^{geo}),$$

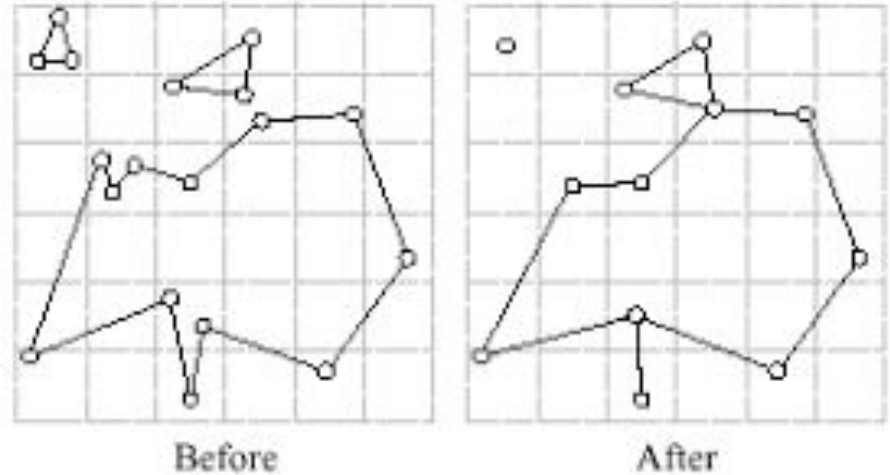
$$\omega_{ij} = \text{softmax}\left(\frac{\varphi_k^{inter}(f_i^{euc})^T \psi_k^{inter}(f_j^{geo})}{\sqrt{d}}\right),$$

# Mesh Simplification

Each level of simplified mesh corresponds to a level of downsampled 3D sparse voxels.

# Vertex Clustering

- (1) construct a bounding box for all vertices on the triangle mesh;
  - (2) regularly subdivide the bounding box into smaller cells (cell size depends on the approximation error);
  - (3) associate each vertex with a single cell which enclose it and form all the vertices inside a specific cell a cluster;
  - (4) collapse all vertices in all clusters into a representative vertices;
  - (5) remove triangles having more than one vertex in a cluster and return the new mesh;
- 



# Quadratic Error Metrics

determines how far a vertex is from an ideal spot after collapse, thus has explicit control over mesh topology.

The plane equation is  $ax + by + cz + d = 0$  where  $a^2 + b^2 + c^2 = 1$ .  $K_p$  is quadric. It approximate error around vertex  $p$

$$K_p = pp^T = \begin{bmatrix} a^2 & ab & ac & ad \\ ab & b^2 & bc & bd \\ ac & bc & c^2 & cd \\ ad & bd & cd & d^2 \end{bmatrix}$$

# QEM Algorithm

1. Calculate this  $K$  for all vertices  $V$ .  $K$  is loss

Now For a given contraction  $(v_1, v_2) \rightarrow v'$ ,  $\text{loss}(v')$  should be minimum wrt  $v_1, v_2$ .

Since the loss function is quadratic, finding its minimum is a linear problem. Thus, we find  $v'$  by solving  $\partial l / \partial x = \partial l / \partial y = \partial l / \partial z = 0$ .

2. Select all valid pairs.
3. Compute the optimal contraction target  $\bar{v}$  for each valid pair  $(v_1, v_2)$ . The error  $v'^T (Q_1 + Q_2) v'$  of this target vertex becomes the cost of contracting that pair.
4. Place all the pairs in a heap keyed on cost with the minimum cost pair at the top.
5. Iteratively remove the pair  $(v_1, v_2)$  of least cost from the heap, contract this pair, and update the costs of all valid pairs involving  $v_1$ .

# Note

Directly applying the QEM method on the original meshes results in high-frequency signals in noisy areas.

VM-Net apply the VC method on the original mesh for the first two mesh levels and then apply the QEM method for the remaining mesh levels.



# Result

Metric: mean IoU

Method	mIoU(%)	Conv Category
TangentConv [57]	43.8	2D-3D
SurfaceConvPF [68]	44.2	
3DMV [9]	48.3	
TextureNet [23]	56.6	
JPBNet [6]	63.4	
MVPNet [25]	64.1	
V-MVFusion [29]	74.6	
BPNet* [20]	<b>74.9</b>	
PointNet++ [45]	33.9	PointConv
FCPN [46]	44.7	
PointCNN [33]	45.8	
DPC [13]	59.2	
MCCN [19]	63.3	
PointConv [65]	66.6	
KPConv [58]	68.4	
JSENet [21]	69.9	
SparseConvNet [17]	72.5	SparseConv
MinkowskiNet [7]	73.6	
SPH3D-GCN [32]	61.0	GraphConv
HPEIN [26]	61.8	
DCM-Net [51]	65.8	
<b>VMNet (Ours)</b>	<b>74.6</b>	Sparse+Graph Conv

Table 1. **Mean intersection over union scores on ScanNet Test [8].** Detailed results can be found on the ScanNet benchmarking website<sup>2</sup>. \* indicates a concurrent work.

## Result with and without attention modules

Information	mIoU(%)	Baseline	Intra	Inter	mIoU(%)
Geo Only	58.1	✓			70.2
Euc Only	71.0	✓	✓		72.1
VMNet(Geo+Euc)	<b>73.3</b>	✓	✓	✓	<b>73.3</b>

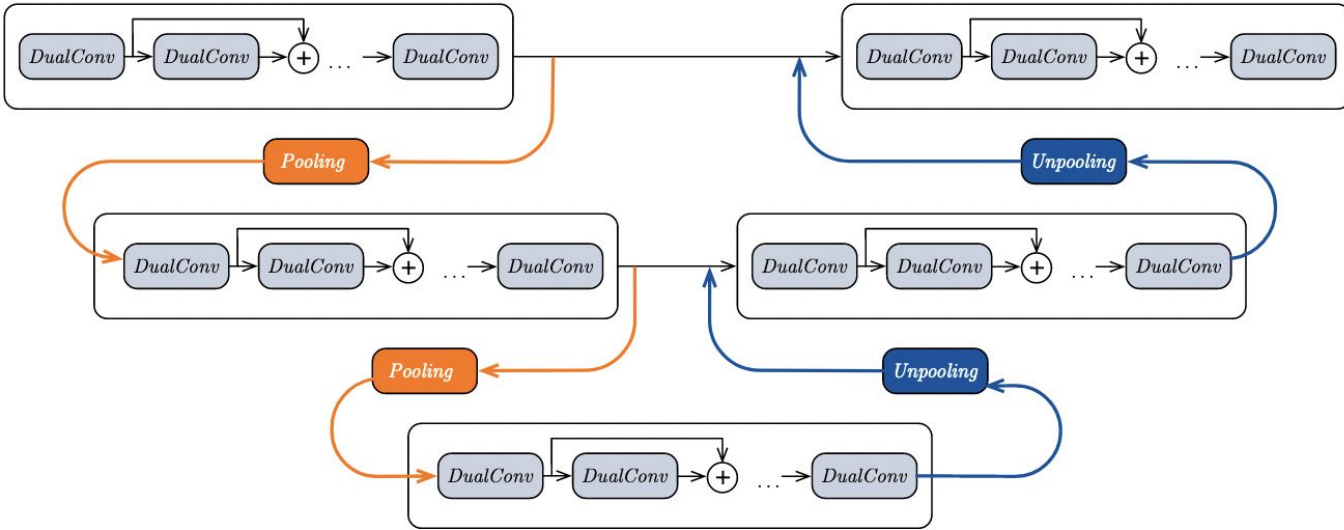
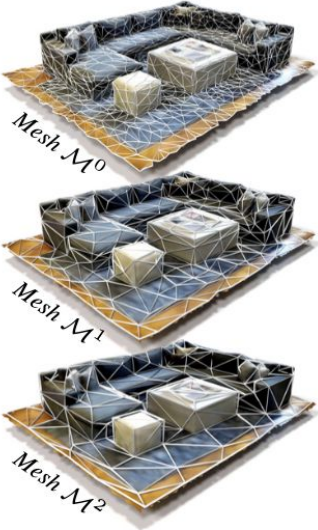
Table 4. **Ablation study: (Left)** Euclidean and geodesic information; **(Right)** Network components.

## Mesh Simplification and Attention operators

Operator	mIoU(%)	Method	mIoU(%)
Vector Attention	72.3	VC only	72.3
EdgeConv	72.6	QEM only	72.9
Scalar Attention	<b>73.3</b>	VC + QEM	<b>73.3</b>

Table 5. **Ablation study:** (**Left**) Attentive operators; (**Right**) Mesh simplification.

# DCM-Net Arch



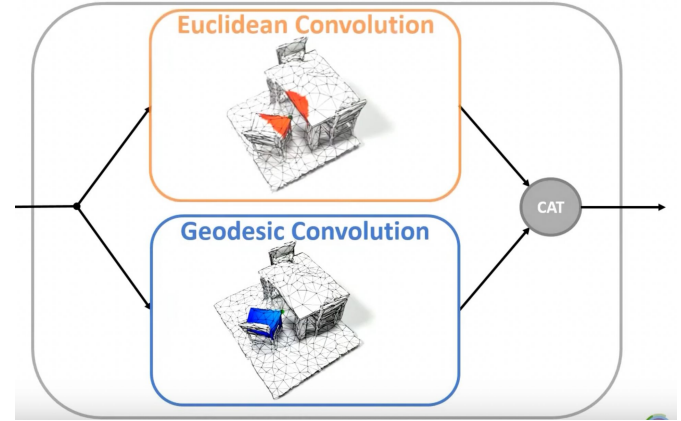
# Difference from DCM-Net

Only difference between VM-Net and DCN-Net, is the type of conv operators.

DCM-Net - GraphConv

VM-Net - SparseConv + GraphConv

Then concat of 2 features is done DCM-Net they are not agg as in VM-Net using attention modules.

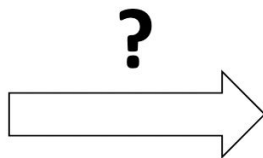
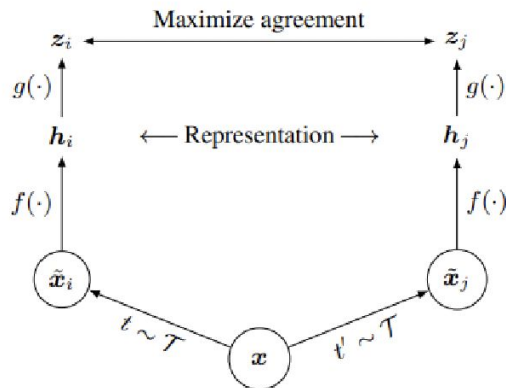


# Self-Damaging Contrastive Learning

<https://arxiv.org/pdf/2106.02990.pdf>

# Abstract

- Real world data is imbalance and sparse in some classes.
- When large networks are pruned, it forgets difficult-to-memorize samples.  
(generally they are long-tail samples)



# Novel idea

create a dynamic **self-competitor model** to contrast with the **target model**, which is a pruned version of the latter.

During training, contrasting the two models will lead to adaptive online mining of the most easily forgotten samples for the current target model, and implicitly emphasize them more in the contrastive loss.

Extensive experiments across multiple datasets and imbalance settings show that SDCLR significantly improves not only overall accuracies but also balancedness, in terms of linear evaluation on the full-shot and few-shot settings.



# What do compressed deep neural networks forget?

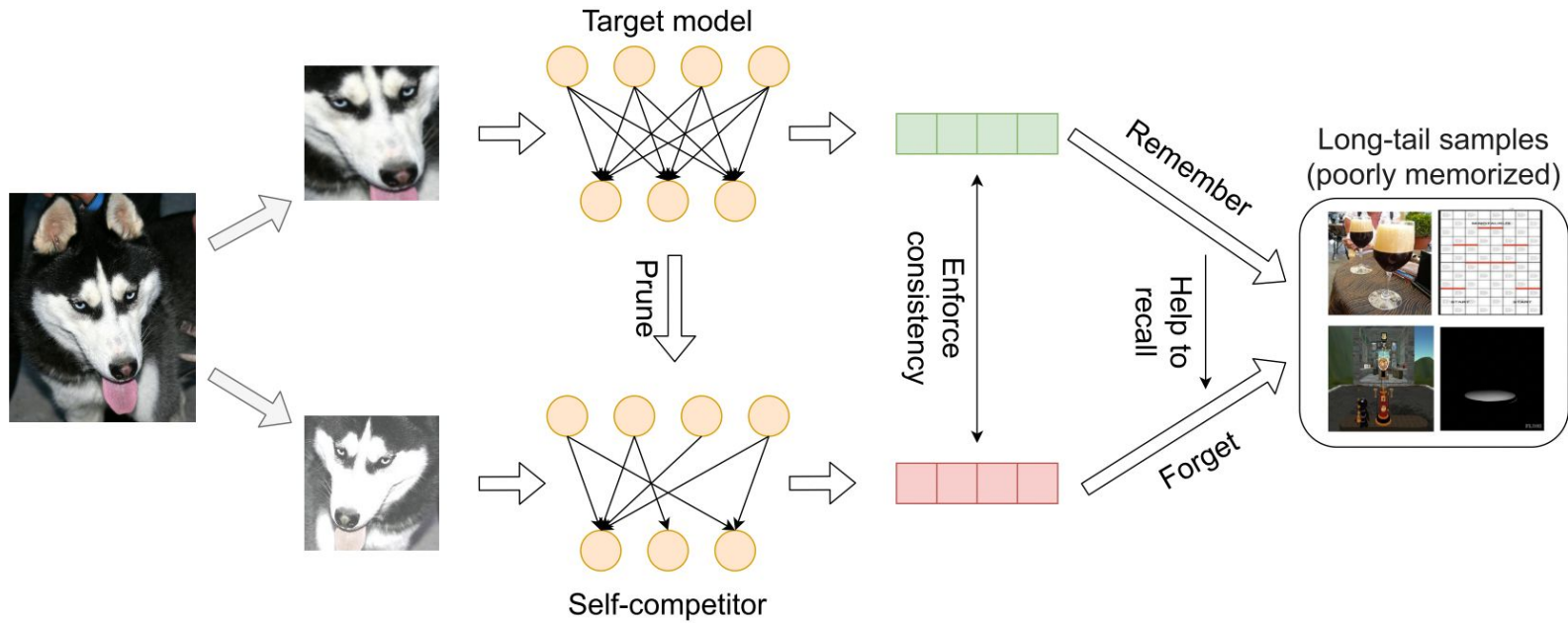
<https://arxiv.org/pdf/1911.05248.pdf>

## Findings:

1. top-1 or top-5 test-set accuracy hide critical details in the ways that pruning impacts model generalization.
2. The examples most impacted by pruning, Pruning Identified Exemplars (PIEs), are more challenging for both models and humans to classify (Human study done here). Compression impairs the model's ability to predict accurately on the long-tail of less frequent instances.
3. Pruned networks are more sensitive to natural adversarial images and corruptions. This sensitivity is amplified at higher levels of compression.

# PIE Samples

Pruning Identified Exemplars (PIEs) are images where there is a high level of disagreement between the predictions of pruned and non-pruned models. PIEs can be thought as images that are forgotten after pruning.



# Self-Competitor Network

Self competitor network is pruned version of target network. **It happens at training step.**

Since the self-competitor is always obtained and updated from the latest target model, the two branches will co-evolve during training. Their contrasting will implicitly give more weights on long-tail samples.

This also makes it different from SimCLR.

# How to Measure Representation Balancedness

The balancedness of a feature space can be reflected by the linear separability w.r.t. all classes.

1. learn the visual representation  $f_v$  on the training dataset with LCL.
2. training a linear classifier layer  $L$  on the top of  $f_v$  with a labeled balanced dataset (by default, the full dataset where the imbalanced subset is sampled from).
3. evaluating the accuracy of the linear classifier  $L$  on the testing set.

Here such accuracy measure is called **linear separability performance**.

# Few Shot performance measure

In place of step 2.

we use only 1% samples of the full dataset from which the pre-training imbalanced dataset is sampled.

Such a metric is few shot performance measure.

# Experiment

Divide each dataset to three disjoint groups: {Many, Medium, Few}.

In subsets of CIFAR10/CIFAR100, Many and Few each include the largest and smallest 1/3 classes, respectively.

For instance

- CIFAR-100: the classes with [500-106, 105-20, 19-5] samples belong to [Many (34 classes), Medium (33 classes), Few (33 classes)] categories, respectively.
- ImageNet, we follow OLTR (Liu et al., 2019) to define Many as classes each with over training 100 samples, Medium as classes each with 20-100 training samples and Few as classes under 20 training samples.

# Loss Function

**Normalized Temperature-scaled Cross Entropy Loss** as in SimCLR.

$$\mathcal{L}_{\text{CL}} = \frac{1}{N} \sum_{i=1}^N -\log \frac{s(v_i, v_i^+, \tau)}{s(v_i, v_i^+, \tau) + \sum_{v_i^- \in V^-} s(v_i, v_i^-, \tau)}$$

$$s(v_i, v_i^+, \tau) = \exp(v_i \cdot v_i^+ / \tau)$$



# Result

Table 1. Comparing the *linear separability performance* for models learned on balanced subset  $D_b$  and long-tail subset  $D_i$  of CIFAR10 and CIFAR100. *Many*, *Medium* and *Few* are split based on class distribution of the corresponding  $D_i$ .

Dataset	Subset	<i>Many</i>	<i>Medium</i>	<i>Few</i>	All
CIFAR10	$D_b$	$82.93 \pm 2.71$	$81.53 \pm 5.13$	$77.49 \pm 5.09$	$80.88 \pm 0.16$
	$D_i$	$78.18 \pm 4.18$	$76.23 \pm 5.33$	$71.37 \pm 7.07$	$75.55 \pm 0.66$
CIFAR100	$D_b$	$46.83 \pm 2.31$	$46.92 \pm 1.82$	$46.32 \pm 1.22$	$46.69 \pm 0.63$
	$D_i$	$50.10 \pm 1.70$	$47.78 \pm 1.46$	$43.36 \pm 1.64$	$47.11 \pm 0.34$

# Result

Table 2. Comparing the *few-shot performance* for models learned on balanced subset  $D_b$  and long-tail subset  $D_i$  of CIFAR10 and CIFAR100. *Many*, *Medium* and *Few* are split according to class distribution of the corresponding  $D_i$ .

Dataset	Subset	<i>Many</i>	<i>Medium</i>	<i>Few</i>	All
CIFAR10	$D_b$	$77.14 \pm 4.64$	$74.25 \pm 6.54$	$71.47 \pm 7.55$	$74.57 \pm 0.65$
	$D_i$	$76.07 \pm 3.88$	$67.97 \pm 5.84$	$54.21 \pm 10.24$	$67.08 \pm 2.15$
CIFAR100	$D_b$	$25.48 \pm 1.74$	$25.16 \pm 3.07$	$24.01 \pm 1.23$	$24.89 \pm 0.99$
	$D_i$	$30.72 \pm 2.01$	$21.93 \pm 2.61$	$15.99 \pm 1.51$	$22.96 \pm 0.43$

# Result

Table 3. Comparing the *linear separability performance* and *few-shot performance* for models learned on balanced subset  $D_b$  and long-tail subset  $D_i$  of ImageNet and ImageNet-100. We consider two long tail distributions for ImageNet: Pareto and Exp, which corresponds to ImageNet-LT and Imagenet-LT-exp, respectively. *Many*, *Medium* and *Few* are split according to class distribution of the corresponding  $D_i$ .

Dataset	Long tail type	Split type	<i>linear separability</i>				<i>few-shot</i>			
			<i>Many</i>	<i>Medium</i>	<i>Few</i>	All	<i>Many</i>	<i>Medium</i>	<i>Few</i>	All
ImageNet	<i>Pareto</i>	$D_b$	58.03	56.02	56.71	56.89	29.26	26.97	27.82	27.97
		$D_i$	58.56	55.71	56.66	56.93	31.36	26.21	27.21	28.33
ImageNet	<i>Exp</i>	$D_b$	57.46	57.70	57.02	57.42	32.31	32.91	32.17	32.45
		$D_i$	58.37	56.97	56.27	57.43	35.98	29.56	28.02	32.12
ImageNet-100	<i>Pareto</i>	$D_b$	68.87	66.33	61.85	66.74	48.82	44.71	41.08	45.84
		$D_i$	69.54	63.71	59.69	65.46	48.36	39.00	35.23	42.16

# Comparison with SimCLR

Table 4. Compare the proposed SDCLR with SimCLR in terms of the *linear separability performance*.  $\uparrow$  means the metric the higher the better and  $\downarrow$  means the metric is the lower the better.

Dataset	Framework	<i>Many</i> $\uparrow$	<i>Medium</i> $\uparrow$	<i>Few</i> $\uparrow$	<i>Std</i> $\downarrow$	<i>All</i> $\uparrow$
CIFAR10-LT	SimCLR	$78.18 \pm 4.18$	$76.23 \pm 5.33$	$71.37 \pm 7.07$	$5.13 \pm 3.66$	$75.55 \pm 0.66$
	SDCLR	$86.44 \pm 3.12$	$81.84 \pm 4.78$	$76.23 \pm 6.29$	$5.06 \pm 3.91$	$82.00 \pm 0.68$
CIFAR100-LT	SimCLR	$50.10 \pm 1.70$	$47.78 \pm 1.46$	$43.36 \pm 1.64$	$3.09 \pm 0.85$	$47.11 \pm 0.34$
	SDCLR	$58.54 \pm 0.82$	$55.70 \pm 1.44$	$52.10 \pm 1.72$	$2.86 \pm 0.69$	$55.48 \pm 0.62$
ImageNet-100-LT	SimCLR	69.54	63.71	59.69	4.04	65.46
	SDCLR	70.10	65.04	60.92	3.75	66.48

# Comparison with SimCLR

Table 5. Compare the proposed SDCLR with SimCLR in terms of the *few-shot performance*.  $\uparrow$  means the metric the higher the better and  $\downarrow$  means the metric is the lower the better.

Dataset	Framework	<i>Many</i> $\uparrow$	<i>Medium</i> $\uparrow$	<i>Few</i> $\uparrow$	<i>Std</i> $\downarrow$	All $\uparrow$
CIFAR10	SimCLR	$76.07 \pm 3.88$	$67.97 \pm 5.84$	$54.21 \pm 10.24$	$9.80 \pm 5.45$	$67.08 \pm 2.15$
	SDCLR	$76.57 \pm 4.90$	$70.01 \pm 7.88$	$62.79 \pm 7.37$	$6.99 \pm 5.20$	$70.47 \pm 1.38$
CIFAR100	SimCLR	$30.72 \pm 2.01$	$21.93 \pm 2.61$	$15.99 \pm 1.51$	$6.27 \pm 1.20$	$22.96 \pm 0.43$
	SDCLR	$29.72 \pm 1.52$	$25.41 \pm 1.91$	$20.55 \pm 2.10$	$3.98 \pm 0.98$	$25.27 \pm 0.83$
Imagenet-100-LT	SimCLR	48.36	39.00	35.23	5.52	42.16
	SDCLR	48.31	39.17	36.46	5.07	42.38

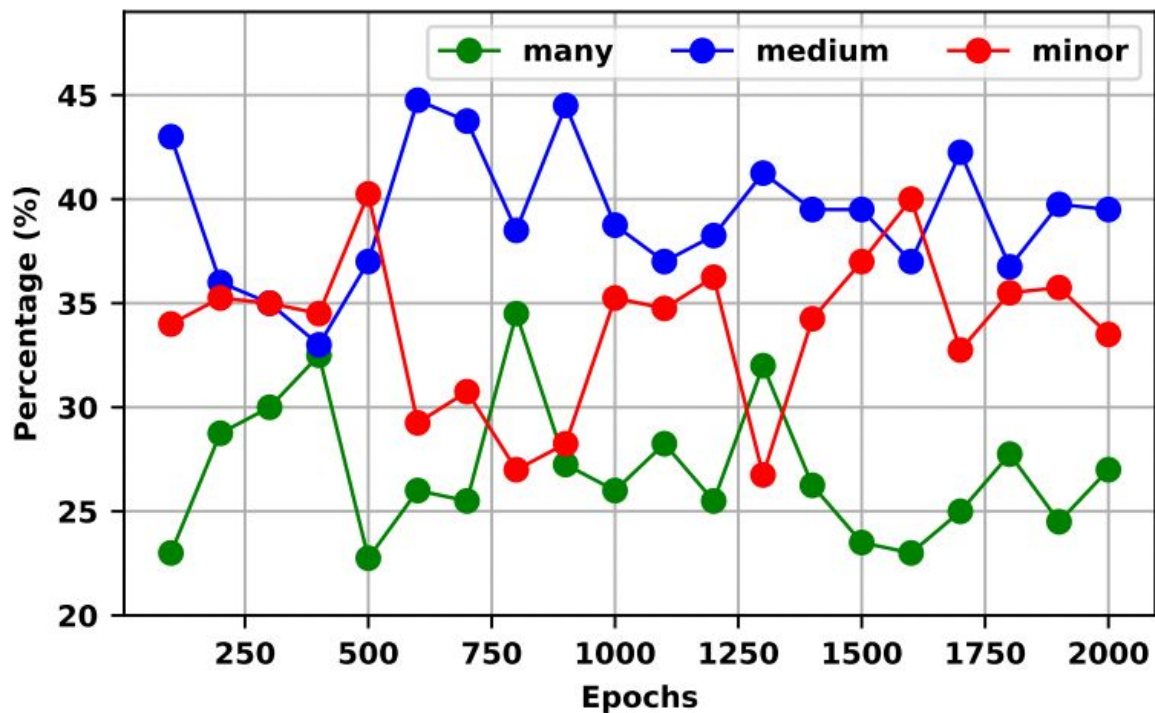
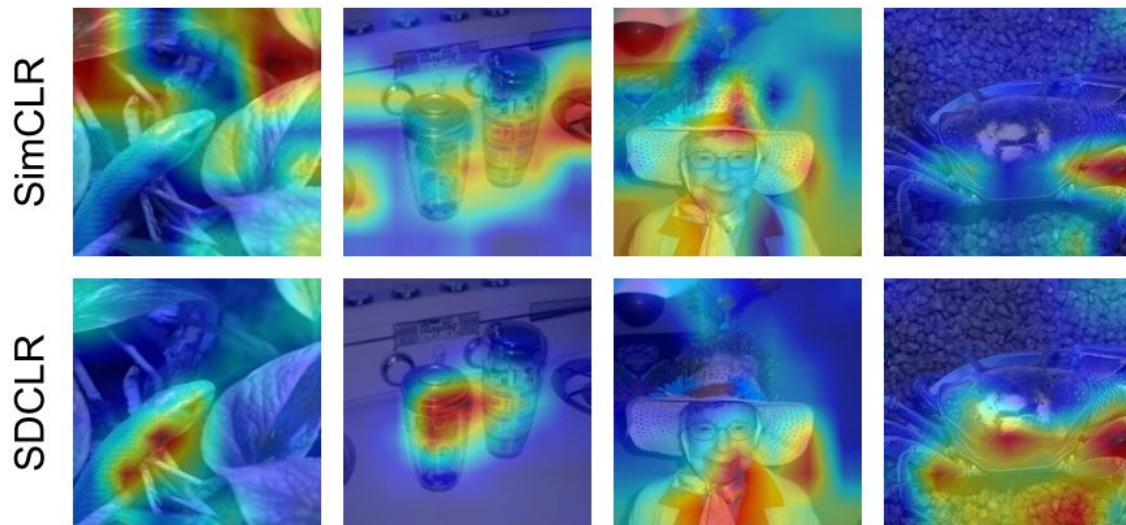


Figure 2. Pre-training on imbalance splits of CIFAR100, The percentage of *many* (●), *medium* (●) and *few* (●) in 1% most easily forgotten data under different training epochs.

# Visualization



*Figure 5. Visualization of attention on tail class images with Grad-CAM (Selvaraju et al., 2017). The first and second row corresponds to SimCLR and SDCLR, respectively.*