# LMSCNet: Lightweight Multiscale 3D Semantic Completion

Luis Roldão
Inria, AKKA Technologies
luis.roldao@inria.fr

Raoul de Charette
Inria
raoul.de-charette@inria.fr

Anne Verroust-Blondet
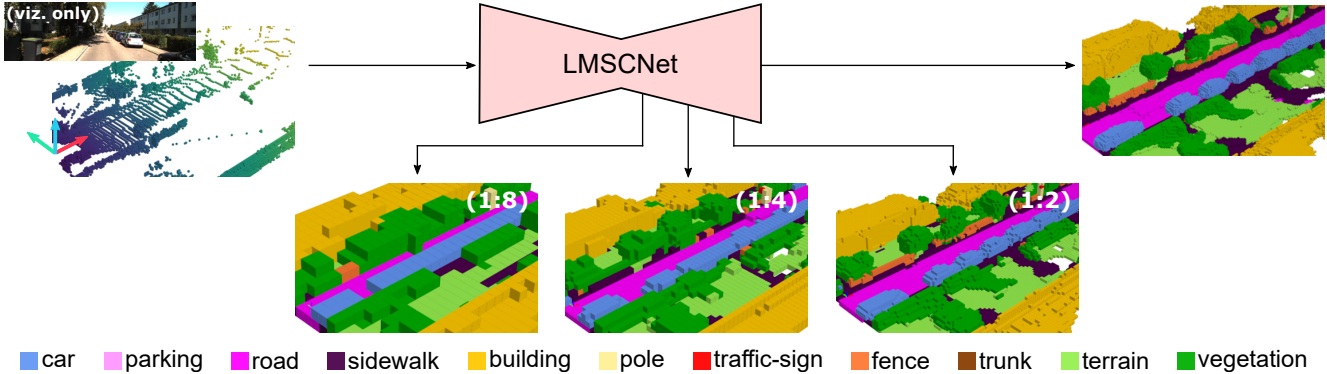Inria
anne.verroust@inria.fr

Figure 1: To prevent heavy computation overhead we use a mix of 2D/3D convolutions to infer multiscale 3D semantic scene completion from sparse voxelized input. Evaluation performed on the challenging SemanticKITTI [1] benchmark shows that our LMSCNet proposal reaches state-of-the-art performance at significantly faster computation speed.

## Abstract

*We introduce a new approach for multiscale 3Dsemantic scene completion from voxelized sparse 3D LiDAR scans. As opposed to the literature, we use a 2D UNet backbone with comprehensive multiscale skip connections to enhance feature flow, along with 3D segmentation heads. On the SemanticKITTI benchmark, our method performs on par on semantic completion and better on occupancy completion than all other published methods – while being significantly lighter and faster. As such it provides a great performance/speed trade-off for mobile-robotics applications. The ablation studies demonstrate our method is robust to lower density inputs, and that it enables very high speed semantic completion at the coarsest level. Our code is available at https://github.com/cv-rits/LMSCNet.*

## 1. Introduction

Understanding 3D surroundings is a natural ability for humans. While past experience allows us to reason about scene geometry and semantics of an entire scene, this proves difficult for computers given the inherently sparse nature of 3D sensors [2] (due to sparse sensing, limited field-of-view, and occlusions). Still, a comprehensive 3D sensing of the scene is crucial for applications like mobile robotics, and more especially for autonomous driving. Recently, semantic scene completion was proposed [36] as a new generative task, where both completion and semantic labels are inferred for the whole scene.

Unlike images, conveniently encoded as 2D tensors, 3D data causes representation challenges. It is thus common to encode the latter as voxel grids processed by 3D Convolutional Neural Networks (CNNs) [11, 36, 12, 26]. This shows good results but also requires heavy computation, as the memory requirement grows cubically with the input voxel resolution [28]. Consequently, most of the literature limits the predicted resolution and network depth, being incapable to perform the task at the same spatial resolution as the input [36, 15, 41]. This drawback has limited the deployment of such methods for real time applications – *i.e.* augmented and virtual reality [38], robotics perception and navigation [23], scene understanding [25], among others – that would greatly benefit of semantic scene completion from sparse LiDAR scan.

We tackle this problem, and propose a Lightweight Multiscale Semantic Completion, coined LMSCNet, where a 3D voxel grid is processed with considerably lighter 2D convolutions without need of additional modalities [1, 15].

1

This is achieved by convolving along one spatial axis (close in spirit to bird-eye view process [5]), while mapping to third dimension with 3D segmentation heads. In our proposal, multiscale completion is also possible given informative features map flow, preserving computation efficiency and enabling very fast inference at coarse levels. Fig. 1 shows the multiscale output of our LMSCNet on the SemanticKITTI dataset [1], using a single sparse LiDAR scan input encoded as voxel. While some works use progressive multiscale losses [25, 10, 12], the literature ignores the benefit of multiscale completion which we prove useful for reducing inference times. To summarize, the main contributions of our work are:

- a novel 3D semantic scene completion pipeline using an occupancy grid,

- a lightweight architecture with mix of 2D/3D convolutions leading significantly less parameters,

- a modular multiscale pipeline which allows coarser inference at very high speed,

- state of the art performance on SemanticKITTI [1] and better performance on completion.

## 2. Related Works

To process 3D data such as point-clouds, some use bird-eye-view [6] or 2D spherical [29] projection. Still, the common strategy relies either on point [32] or voxel [33] networks. The inherent limitation of voxel representation is the staggering memory requirement due to empty voxels, which led to optimized structures [33] or use of sparse convolutions [18] to prevent dilation of the data manifold. When it comes from real sensing, 3D data is inherently sparse (e.g. LiDAR scan, stereo, etc.) and its densification was initially framed as a *completion or reconstruction* task. Recent works though, also assign semantic labels to their output subsequently referring to this as *semantic completion*.

**3D completion & reconstruction.** First works approximate missing data as a set of primitive local surfaces either from the data structure [37, 35, 11, 37] or using truncated Signed Distance Functions (TSDFs) [8, 30, 34], while others use continuous energy minimization [21].

More recently, learning methods boosted the completion of occluded and unseen regions [14, 17, 43]. In [22], voxel labels are predicted using a Conditional Random Field (CRF), while [11] uses 3D-Encoder-Predictor to correlate observed scene with *a priori* known 3D shapes. A few works also benefit from Signed Distance Functions representation as it provides richer gradient flow [31, 33, 10]. For memory reason, [10] uses TSDFs input with sparse encoder and partially dense decoder to propagate features in unknown regions. While this effectively reduces memory, because TSDFs are denser than occupancy grids, we argue it would require a bigger and denser decoder, thus annihilating the benefit of any sparse encoding. Other end-to-end completion networks were also proposed in [9, 4]. Despite discretization, we preferred a voxel-based implementation due to size limitations of point-based networks, even if there are promising object completion results [40].

**3D Semantic Scene Completion.** SSCNet [36] was the first work to combine semantic segmentation and scene completion end-to-end with 3D CNNs. Further works also use additional RGB data by projecting or fusing semantic features from an image network [15, 26, 24, 27]. An alternative to 3D data only is to encode LiDAR scans as spherical projection [29], which enriches neighboring information [1, 19]. While this boosts performance, it also increases the network complexity and subsequently the inference time. Generative Adversarial Networks (GANs) have also been proposed to enforce realistic outputs [39, 7] but are harder to train. To lower memory consumption with the preferred voxelized representations, Spatial Group Convolutions (SGC) [41] divide input into groups for efficient processing at the cost of small performance drops.

Different from the literature, we rely solely on 3D voxelized occupancy data avoiding any preprocessing (as for TSDFs, SDFs, etc.), and propose a lightweight architecture with additional multiscale capability.

## 3. Method

We tackle the problem of dense 3D semantic completion where the task is to assign a semantic label to each individual voxel. Given a sparse 3D voxel grid, the goal is to predict the 3D semantic scene representation, where each voxel is being assigned a semantic label $C = [c_0, c_1, \ldots, c_N]$, where $N$ is the number of semantic classes and $c_0$ stands for free voxels.

Our architecture, coined LMSCNet and shown in Fig. 2, uses a lightweight UNet style architecture to predict 3D semantic completion at multiple scales, allowing fast coarse inference, beneficial for mobile robotics applications. Instead of greedy 3D convolutions, we mostly employ 2D convolutions along the height axis; similar to a bird-eye view. In the following we detail our custom lightweight 2D/3D architecture (Sec. 3.1), the multiscale reconstruction (Sec. 3.1), and the overall training pipeline (Sec. 3.2).

### 3.1. Lightweight multiscale 2D/3D architecture

To infer a dense output from the sparse input voxel grid, we use a standard encoder-decoder UNet architecture with 4 levels, thus learning features at decreasing resolutions.
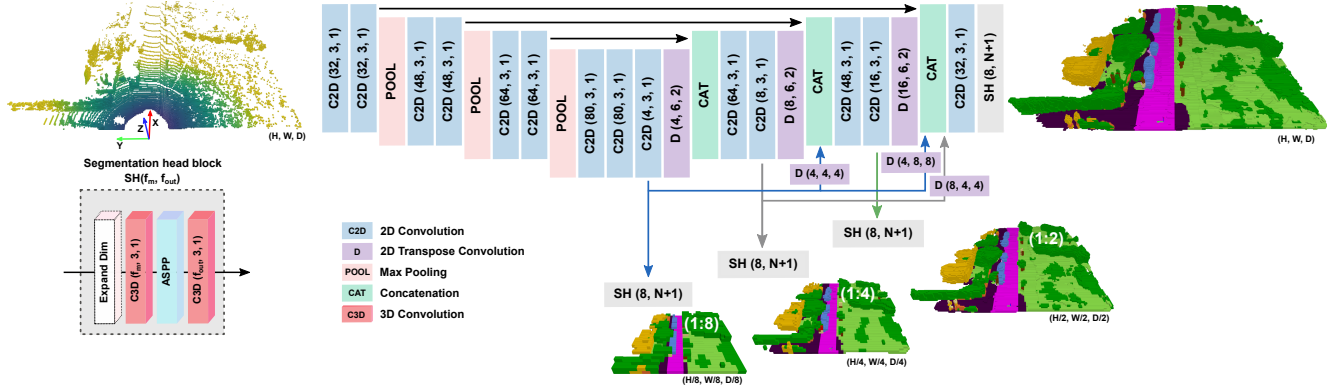
**Figure 2: LMSCNet: Lightweight Multiscale Semantic Completion Network.** Our pipeline uses a UNet architecture with 2D backbone convolutions (in blue) and 3D segmentation heads (in gray) to perform 3D semantic segmentation and completion at different scales, while preserving low complexity. Convolution parameters are shown as: (number of filters, kernel size and stride). Notice that we intentionally lower the 2D features dimension and use Atrous 3D convolutions (ASPP blocks from [26]) to preserve low inference complexity.

At each level, a series of convolution operations is applied followed by a pooling; downscaling the resolution size by 2. The reduction of spatial dimensions in UNets is beneficial for semantic tasks as it subsequently increases the kernels field-of-view at no cost. Note that dilated convolutions (a.k.a 'atrous') with increasing dilation rates cannot be used in the encoder due to the sparse input nature. Though dense convolutions in the encoder imply a dilation of the input manifold [18], we argue this is beneficial for 3D semantic completion, given the sparse↦dense nature of the task.

**2D backbone.** To preserve a lightweight architecture, we use 2D convolutions along the X,Y dimensions, thus turning the height dimension (Z) into a feature dimension. Notice that we directly process 3D data in contrast to other 2D/3D works that rely on 2.5D data (e.g. depth [19, 26], bird-eye view [6]). While using 2D convolutions implies loosing 3D spatial connexity, it also enables significantly lighter operations. To further reduce the memory requirements, we keep a minimum number of features in each convolution layer. Along with the standard skip connections, we also enhance information flow in the decoder by concatenating the output of each level to all lower levels. Technically, we upsample coarse feature maps learning ad-hoc deconvolution before concatenation to lower levels, which is shown with purple deconv blocks along blue and gray arrows in Fig. 2. Intuitively, this enables our network to use high level features from coarser resolutions, and thus enhancing the spatial contextual information.

**3D segmentation head.** Different from other works handling point cloud as bird-eye-view, the task of 3D semantic completion actually requires to retrieve the 3rd dimension "lost" with 2D convolutions. In other words, while 2D CNNs output 3D features maps, our decoder must output 4D tensor; the last dimension being the semantic class-wise probability distribution.

To address this, we introduce 3D segmentation heads depicted as gray blocks in Fig. 2. The heads use a series of dense and dilated convolutions. The latter, in the form of Atrous Spatial Pyramid Pooling (aka ASPP [5, 26]), is beneficial to fuse information from different receptive fields thanks to the convolutions with increasing dilation rates (here [1, 2 and 3]). Note that dilated convolutions, though light and powerful, are not appropriate for sparse inputs and, as such, cannot be used in the encoder. In our segmentation head, the benefit of preceding ASPP with dense 3D convolutions is dual: a) to further densify the feature maps, b) to ward off features from the segmentation heads and the backbone features. This last property is required to enable multiscale capacity, which we now describe.

**Multiscale completion.** In the same vein as [41, 10], we aim to output multiscale completion to enable both coarse scene representation and faster scene completion at lower resolution – beneficial for mobile robotics applications. We subsequently attach a 3D segmentation head after each level of the 2D UNet architecture, thus providing outputs at input relative scale of $\frac{1}{2^l} \ \forall \ l \in \{0, 1, 2, 3\}$. A sample output is shown in Fig. 3. As already mentioned, we noticed experimentally the importance of separating the segmentation features from the main features of the 2D backbone, which again justifies the additional 3D convolutions in the segmentation head. The main interests of our multiscale architecture is that it infers semantic scene completion at a desired scale as needed, reducing the computation and mem-

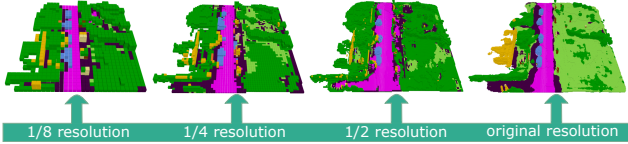| 1/8 resolution | 1/4 resolution | 1/2 resolution | original resolution |

Figure 3: Our pipeline enables multiscale reconstruction. To supervise coarser representation, we use majority vote pooling from the original resolution ground truth.

ory requirements. This is further analyzed in Sec. 4.1.1.

## 3.2. Training strategy

We train our LMSCNet from scratch in a standard end-to-end fashion from pairs of sparse input voxel ($x$) and semi-dense semantically labeled voxel grid ($\hat{y}$). It is important to note that in a real setup, a dense ground truth is impractical for scene completion, due to occlusions and sensor field-of-view limitations. As such, the ground truth $\hat{y}$ is *also* sparse and encoded with N+2 classes ($N$ semantic classes, 1 free class, 1 unknown). Similar to others [36, 26, 15] we use a sparse loss strategy, back propagating the gradient only where ground truth is known.

For each scale $l$, we train with a cross-entropy loss defined as

$$\mathcal{L}_l = -\sum_{c=0}^{N} w_c \hat{y}_{i,c} \log \left( \frac{e^{y_{i,c}}}{\sum_{c'}^{N} e^{y_{i,c'}}} \right) , \qquad (1)$$

where $y$ is the network output, $i$ a voxel index, and $\hat{y}_{i,c}$ a one-hot vector (i.e. $\hat{y}_{i,c} = 1$ if voxel $i$ is labeled class $c$, otherwise $\hat{y}_{i,c} = 0$). Note that semantic tasks are by nature highly class-imbalanced problems. This is especially true in outdoor settings, which causes the prevalence of classes like road or vegetation. We account for the class-imbalance nature in Eq. 1 by weighting each class loss according to the inverse of the class-frequency $f_c$ as in [29], thus using $w_c = \frac{1}{\log(f_c+\epsilon)}$ (with $\epsilon \ll 1$). Finally, the complete network loss is a weighted sum of all level losses[1] and writes:

$$\mathcal{L} = \sum_{l=0}^{3} \alpha_l \mathcal{L}_l , \qquad (2)$$

where $\alpha_l$ is the per-level loss weight, written for generality, though we use $\alpha_l = 1, \forall l$ which works well and preserves multiscale capacity. Note that some of our choices were guided by faster training or inference speed. For example, unlike [41, 42, 13, 10, 36] we avoid using Truncated Signed Distance Function variants (TSDF) that require a greedy computation time and was found to be of little benefit [15, 1]. We also tried to encode input as N+2 classes, that is with *unknown* class, but we noticed little improvement –

---
[1]In Eq. 2, losses from heterogeneous resolutions can be summed due to the ad-hoc normalization in Eq. 1

if any – at the cost of a large pre-processing time for ray casting.

## 4. Experiments

We evaluate our LMSCNet method by training on the recent semantic scene completion benchmark SemanticKITTI [1] providing 3D voxel grids from semantically labeled scans of HDL-64E rotating LiDAR in outdoor urban scenes [16]. In [1], inputs are voxelized single scans, while the ground truth was obtained from the voxelized aggregation of successive registered scans. Grids are 256x256x32 with 0.2m voxel size, and it is important to note that input *and* ground truth are sparse, with average density of 6.7% and 65.8%, respectively. We use standard mIoU as a semantic completion metric, measuring the intersection over union averaged over all classes (20 semantic classes + *free*). Additionally, we consider completion metrics IoU, Precision, and Recall to provide a sense of the scene completion quality, regardless of the assigned semantic labels (i.e. considering the binary *free / occupied* setting). Note that completion is crucial for obstacle avoidance in mobile robotics.

**Implementation details.** We train using the original train/val splits with 3834/815 grids [1], adding x-y flipping augmentation for generalization. Adam optimizer is used ($\beta_1 = 0.9$, $\beta_2 = 0.999$) with learning rate of 0.001 scaled by $0.98^{\text{epoch}}$. Training fits in a single 11GB GPU with batch size 4, taking around 48 hours to converge ( 80 epochs).

## 4.1. Performance

In the following we report performance against four state-of-the-art methods: SSCNet [36], TS3D [15], TS3D+DNet [1], TS3D+DNet+SATNet [1]. Because SSCNet output is 4x downsampled, we also report performance using deconvolution to reach full input resolution, hereafter denoted SSCNet-full. We refer to the supplementary for details on the required architectures adjustments.

Hereafter, we denote our multiscale architecture as LMSCNet. We detail semantic completion performance and then demonstrate the speed and lightness of our architecture.

### 4.1.1 Semantic Scene Completion

Performance on the SemanticKITTI benchmark [1] is reported in Tab. 1 against all published methods and SSCNet-full. The evaluation was conducted on the official server (i.e. hidden test set) hence, with the full size ground truth.

Overall, we perform on par with the best methods, though 2nd on the semantic completion metric (mIoU). On the latter, TS3D+DNet+SATNet [1] is slightly better despite their significantly heavier and slower network. Note also

| Approach | scene completion | | | semantic scene completion | | | | | | | | | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | precision | recall | IoU | road (15.30%) | sidewalk (11.13%) | parking (1.12%) | other-ground (0.56%) | building (14.1%) | car (3.92%) | truck (0.16%) | bicycle (0.03%) | motorcycle (0.03%) | other-vehicle (0.20%) | vegetation (39.3%) | trunk (0.51%) | terrain (9.17%) | person (0.07%) | bicyclist (0.07%) | motorcyclist (0.05%) | fence (3.90%) | pole (0.29%) | traffic-sign (0.08%) | mIoU |
| SSCNet [36] | 31.71 | 83.40 | 29.83 | 27.55 | 16.99 | 15.60 | 6.04 | 20.88 | 10.35 | 1.79 | 0 | 0 | 0.11 | 25.77 | 11.88 | 18.16 | 0 | 0 | 0 | 14.40 | 7.90 | 3.67 | 9.53 |
| *SSCNet-full [36] | 59.64 | 75.52 | 49.98 | 51.15 | 30.76 | 27.12 | 6.44 | 34.53 | 24.26 | 1.18 | **0.54** | **0.78** | **4.34** | 35.25 | 18.17 | 29.01 | 0.25 | 0.25 | **0.03** | 19.87 | 13.10 | 6.73 | 16.14 |
| TS3D [15] | 31.58 | 84.18 | 29.81 | 28.00 | 16.98 | 15.65 | 4.86 | 23.19 | 10.72 | 2.39 | 0 | 0 | 0.19 | 24.73 | 12.46 | 18.32 | 0.03 | 0.05 | 0 | 13.23 | 6.98 | 3.52 | 9.54 |
| TS3D+DNet [1] | 25.85 | **88.25** | 24.99 | 27.53 | 18.51 | 18.89 | **6.58** | 22.05 | 8.04 | 2.19 | 0.08 | 0.02 | 3.96 | 19.48 | 12.85 | 20.22 | **2.33** | **0.61** | 0.01 | 15.79 | 7.57 | **6.99** | 10.19 |
| TS3D+DNet+SATNet [1] | 80.52 | 57.65 | 50.60 | 62.20 | 31.57 | 23.29 | 6.46 | 34.12 | 30.70 | **4.85** | 0 | 0 | 0.07 | 40.12 | **21.88** | **33.09** | 0 | 0 | 0 | **24.05** | **16.89** | 6.94 | **17.70** |
| LMSCNet (ours) | 77.11 | 66.19 | 55.32 | 64.04 | 33.12 | 24.91 | 3.22 | **38.67** | 29.48 | 2.53 | 0 | 0 | 0.11 | 40.53 | 18.97 | 30.77 | 0 | 0 | 0 | 20.52 | 15.72 | 0.54 | 17.01 |
| LMSCNet-singlescale (ours) | **81.55** | 65.07 | **56.72** | **64.80** | **34.68** | **29.02** | 4.62 | 38.08 | **30.89** | 1.47 | 0 | 0 | 0.81 | **41.31** | 19.89 | 32.05 | 0 | 0 | 0 | 21.32 | 15.01 | 0.84 | 17.62 |

* Own implementation.

Table 1: Comparison of published methods on the official SemanticKITTI [1] benchmark. Despite light mixed 2D/3D reasoning, our network performs 2nd on the semantic metrics (mIoU), outdistanced by the more complex TS3D+DNet+SATNet also twice slower than us. On the completion metrics (IoU), we perform 1st with a comfortable margin. The last two rows show that LMSCNet is better in its single scale version (*LMSCNet-singlescale*), though this comes at the cost of loosing multiscale capacity. Except for SSCNet-full, all results originate from [1].

that TS3D uses additional RGB input, and all TS3D+DNet use also LiDAR refraction intensity. Conversely, LMSCNet is more versatile as it uses only occupancy grid input. Notice that the highly imbalanced class frequencies (shown in parenthesis in Tab. 1) also illustrate the task complexity. Specifically, we outperform others on the largest four classes but performs on par or lower on the others, which advocates for some improvement in our balancing strategy. On the completion metrics (IoU) our method outperforms all others by a comfortable margin. Again, completion is of high importance for practical mobile robotics applications.

In addition to the multiscale proposal (LMSCNet), we also report LMSCNet-singlescale – a variation of LMSCNet where we train with $\mathcal{L} = \mathcal{L}_0$ –, which logically performs a little better at full size though at the cost of loosing crucial multiscale capacity.

**Qualitative performance.** We compare qualitatively full size outputs of our LMSCNet and SSCNet-full in Fig. 4, with views pairs from 4 scenes of the SemanticKITTI validation set[2]. At the rightmost, ground truth visualization also illustrates the sparse supervision complexity since holes are still visible. Our method produces visually smoother semantic labels, easily noticeable in rows 5-8, and is able to reconstruct thin structures, like trees or cars (rows 6 or 7). For comprehensive analysis, we further test the same model (trained on 64-layer LiDAR) on the popular nuScenes dataset [3], which has been registered using a 32-layers lidar sensor. Fig. 5 shows that our network better adjusts to the change of density and maintains the smoothness

---

[2]Note that SemanticKITTI benchmark (i.e. test set) does not provide any visual results. Hence, we omit TS3D baselines due to retraining complexities and their use of additional modalities (RGB or LiDAR intensity).

| LMSCNet scale | IoU | mIoU |
|---|---|---|
| 1:1 (full size) | 54.22 | 16.78 |
| 1:2 | 56.27 | 16.78 |
| 1:4 | 59.36 | 17.19 |
| 1:8 | 65.45 | 17.37 |

Table 2: LMSCNet multiscale semantic completion performance on SemanticKITTI validation set. We reach good performance at all levels, even better at the coarsest levels .

in the reconstruction.

**Multiscale performance.** Tab. 2 shows multiscale performance of our method on the SemanticKITTI validation set, where the scale is relative to the full size resolution (level 0). From Sec. 3.1, scale at level $l$ is $\frac{1}{2^l}$. Ground truths at lower resolution were obtained from majority vote pooling of the full size ground truth. From the above table, our architecture maintains a good performance in all resolutions, with best performance logically reached at the lowest resolution (highest level). Qualitative multiscale completion is visible in Fig. 3. We argue that our architecture reaches multiscale capacity thanks to the disentanglement of segmentation features with our custom head. Additionally, at coarser resolution our network reaches very fast inference, which will be described in details in the following section.

### 4.1.2 Architectures comparison.

Tab. 3 reports networks statistics for our architecture and all above mentioned baselines. From the latter, even at full size LMSCNet has significantly less parameters (0.35M) and lower computational cost for inference (72.6G FLOPs). Compared to any TS3D baselines it is at least an order of magnitude faster. However, SSCNet (original or full) is

| | Input | SSCNet-full [36] | LMSCNet (ours) | Ground Truth |
|---|---|---|---|---|

bicycle  car  motorcycle  truck  other-vehicle  person  bicyclist  motorcyclist  parking  road
sidewalk  other-ground  building  pole  traffic-sign  fence  trunk  terrain  vegetation
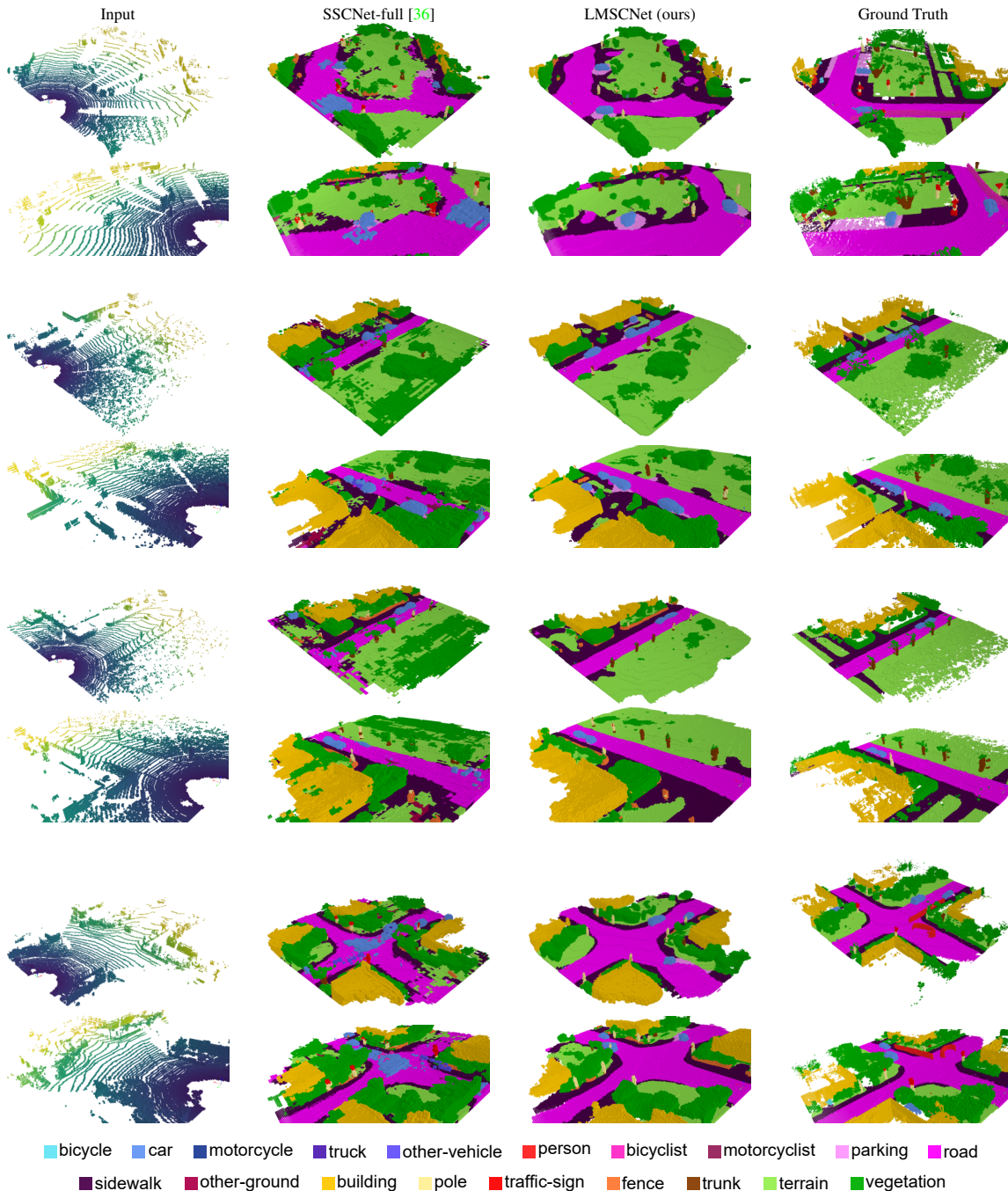
Figure 4: Qualitative 3D semantic completion at full size on the SemanticKITTI [1] validation set. Each pairs of rows show a single scene with different viewpoints. Compared to SSCNet-full [36], our LMSCNet provides smoother semantics labels and is capable of retrieving finer details. This is evident when looking at the cars (rows 7-8) or the trees (rows 5-6).

twice faster than LMSCNet, though with more parameters and worse performance (cf. Tab. 1). Since lighter models does not *always* run faster due to the sequentiality of some operations on GPU, we conjecture the higher speed of SSCNet is caused by the lower number of convolutional operations compared to LMSCNet full scale (16 vs. 25).

In last rows of Tab. 3, we report statistics for coarser completion, removing *unnecessary* parts of our network at inference. Lower resolution inference allows significant speedups in the processing, reaching 372 FPS at the high-
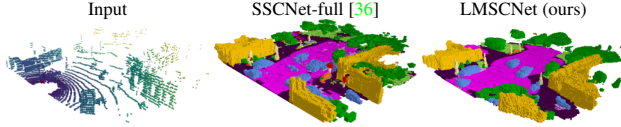
Figure 5: Inference results on nuScenes [3] with 32-layers LiDAR, while being trained on 64-layers SemanticKITTI. Our method performs well with sharp scene labeling, despite the change of input density.

| Method | Params (M) | FLOPs (G) | FPS |
|--------|-----------|-----------|-----|
| *SSCNet [36] | 0.93 | 82.5 | 56.90 |
| *SSCNet-full [36] | 1.09 | 769.6 | 45.94 |
| *TS3D [15] | 43.77 | 2016.7 | 9.79 |
| *TS3D+DNet [1] | 51.31 | 847.1 | 8.72 |
| *TS3D+DNet+SATNet [1] | 50.57 | 905.2 | 1.27 |
| LMSCNet | 0.35 | 72.6 | 21.28 |
| LMSCNet (1:2) | 0.32 | 13.7 | 126.38 |
| LMSCNet (1:4) | 0.28 | 5.7 | 323.46 |
| LMSCNet (1:8) | 0.24 | 4.4 | 372.24 |

\* Own implementation to compute network statistics.

Table 3: Network statistics. Even at full resolution LMSCNet (ours) has significantly less parameters with lower FLOPs. On a speed basis, we are twice slower than SSCNet-full [36] which performs worse than us (see Tab. 1). Still, our multiscale versions – denoted *LMSCNet (1:x)* – enable very fast inference.
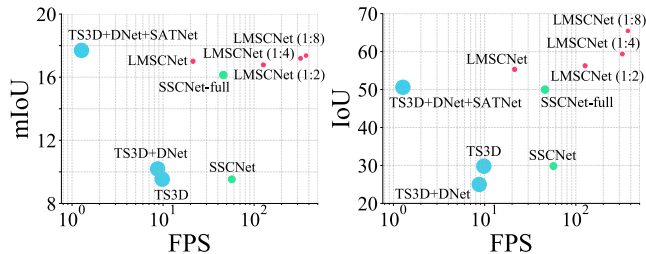


Figure 6: Architectures performance versus speed (markers are scaled with # of parameters). Notice that TS3D+DNet+SATNet is the only better method on semantics (+0.69 mIoU) though less time performant (x17 slower) and worse on completion (-4.72 IoU).

est scale – 6x faster than SSCNet and 300x faster than TS3D+DNet+SATNet –. Fig. 6 illustrates the architectures performance versus speed. Notice that *even at full scale* we provide a better speed-performance balance. Because semantic completion is an application of high interest for mobile robotics, like autonomous driving, our lighter architecture is beneficial for embedded GPUs and enables coarse scene analysis at high speed.
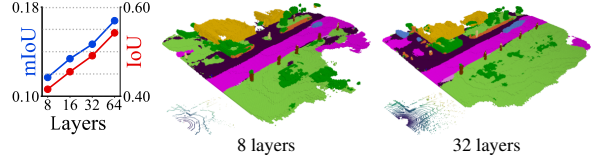


Figure 7: Semantic scene completion results from simulated lower resolution LiDAR sensors (downsampled from 64 layers input). Even with only 8 layers input our LMSC-Net correctly predicts the scene outline.

| Method | IoU | mIoU |
|--------|-----|------|
| LMSCNet (ours) | 54.22 | 16.78 |
| w/o Deconv | 52.79 | 15.64 |
| w/o ASPP | 53.81 | 16.21 |
| w/o Multiscale UNet | 53.54 | 16.22 |

Table 4: Ablation study of our model design choices on the SemanticKITTI [1] validation set.

## 4.2. Ablation studies

To study the benefit of our design choices, we conduct a series of ablation studies on SemanticKITTI validation set. This is done by modifying important blocks of our architecture and evaluating its performance.

**Influence of input resolution.** We evaluate our robustness, by retrieving the original 64-layers KITTI scans used in SemanticKITTI and simulating 8/16/32 layers LiDARs with layers subsampling[3], as in [20].

Fig. 7 shows quantitative and qualitative performance using simulated and original LiDAR. As expected, lower layers input deteriorate the performance, especially in areas far from the sensor location, but our network still performs reasonably well on semantics (mIoU) and completion (IoU). This is visible in the middle image, as 8 layers input (2.10% density) is sufficient to retrieve the general outline of the scene.

**Deconv versus Upsampling.** As we aimed to preserve a lightweight architecture, we tried to remove the parameters-greedy deconv layers from our network (cf. Sec. 2), replacing them with up-sampling layers. From Tab. 4, performance *without deconv* introduces a 1.43% and 1.14% performance drop for completion and semantic completion respectively, with only 3% less parameters.

**Dilated convolutions.** We evaluate the benefit of dilated convolutions in the decoder by ablating ASPP blocks from

---

[3]Every 2nd, 4th and 8th layer are subsampled to simulate 32, 16 and 8 layer LiDARs, respectively. Unlike [20], note that data SemanticKITTI uses KITTI odometry set in which data is already untwisted.

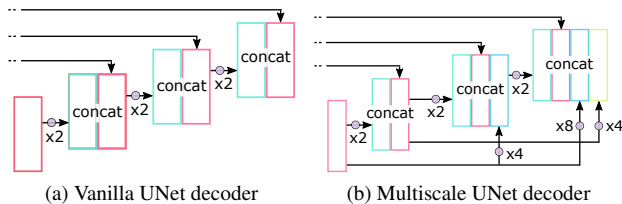(a) Vanilla UNet decoder      (b) Multiscale UNet decoder

Figure 8: Decoders comparison. While Vanilla UNet decoder only considers features from the previous level (a), we instead use Multiscale UNet where all coarser levels enhance spatial contextual information (b). Circles show intermediary operations to reach required feature maps size.

the segmentation head (see Fig. 2). Tab. 4 indicates that mIoU drops by 0.41% without ASPP. We conjecture that the boost of ASSP results come from the increasing receptive fields of the inner dilated convolutions, providing richer features.

**Multiscale UNet decoder.** As illustrated in Fig. 8a, unlike vanilla UNet decoder we concatenate the features at the end of each decoder level to all other levels. This is intended to aggregate multiscale features and should intuitively help considering coarser semantic features for fine resolutions.

We assess the benefit of our *multiscale UNet* by evaluating *Vanilla UNet* in the last row of Tab. 4, which shows that our proposal boosts completion by 0.68% and semantic completion by 0.56%.

## 5. Conclusion

We proposed a novel method, coined LMSCNet, for 3D Semantic Scene Completion, which benefits from mixing 2D/3D convolutions to preserve lightweight architecture, while enabling inference at multiple scales. On the challenging SemanticKITTI benchmark, we perform on par with other methods on semantic completion with a much lighter architecture and at faster inference speed. For completion, we outperform the state of the art. Results show that loss of 3D spatial connexity caused by the 2D convolutions *does not* impair performance. We attribute this to the uniform dimensional variance in used application (i.e. constant sensor viewpoint, urban outdoor scenes). We conjecture that data with higher variance in all directions would cause a higher impact. Also of interest for mobile robotics, our proposal is robust to much lower input density and our multiscale capacity enables scene completion for lower resolution at very high speed.

## References

[1] J. Behley, M. Garbade, A. Milioto, J. Quenzel, S. Behnke, C. Stachniss, and J. Gall. SemanticKITTI: A dataset for semantic scene understanding of LiDAR sequences. *International Conference on Computer Vision (ICCV)*, pages 9296–9306, 2019. 1, 2, 4, 5, 6, 7, 10, 11

[2] M. Berger, A. Tagliasacchi, L. M. Seversky, P. Alliez, G. Guennebaud, J. A. Levine, A. Sharf, and C. T. Silva. A survey of surface reconstruction from point clouds. *Comput. Graph. Forum*, 36:301–329, 2017. 1

[3] H. Caesar, V. Bankiti, A. H. Lang, S. Vora, V. E. Liong, Q. Xu, A. Krishnan, Y. Pan, G. Baldan, and O. Beijbom. nuScenes: A multimodal dataset for autonomous driving. *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020. 5, 7, 10, 11

[4] A. X. Chang, A. Dai, T. A. Funkhouser, M. Halber, M. Nießner, M. Savva, S. Song, A. Zeng, and Y. Zhang. Matterport3D: Learning from RGB-D data in indoor environments. *International Conference on 3D Vision (3DV)*, pages 667–676, 2017. 2

[5] L.-C. Chen, G. Papandreou, I. Kokkinos, K. Murphy, and A. L. Yuille. DeepLab: Semantic image segmentation with Deep Convolutional Nets, Atrous Convolution, and Fully Connected CRFs. *Transactions on Pattern Analysis and Machine Intelligence (PAMI)*, 40:834–848, 2018. 2, 3

[6] X. Chen, H. Ma, J. Wan, B. Li, and T. Xia. Multi-view 3D object detection network for autonomous driving. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 1907–1915, 2017. 2, 3

[7] Y. Chen, M. Garbade, and J. Gall. 3D semantic scene completion from a single depth image using adversarial training. In *International Conference on Image Processing (ICIP)*, pages 1835–1839, 2019. 2

[8] B. Curless and M. Levoy. A volumetric method for building complex models from range images. In *SIGGRAPH*, 1996. 2

[9] A. Dai, A. X. Chang, M. Savva, M. Halber, T. A. Funkhouser, and M. Nießner. ScanNet: Richly-annotated 3D reconstructions of indoor scenes. *Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 2432–2443, 2017. 2

[10] A. Dai, C. Diller, and M. Nießner. SG-NN: Sparse generative neural networks for self-supervised scene completion of RGB-D scans. *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020. 2, 3, 4

[11] A. Dai, C. R. Qi, and M. Nießner. Shape completion using 3D-encoder-predictor CNNs and shape synthesis. *Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 6545–6554, 2017. 1, 2

[12] A. Dai, D. Ritchie, M. Bokeloh, S. Reed, J. Sturm, and M. Nießner. ScanComplete: Large-scale scene completion and semantic segmentation for 3D scans. *Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 4578–4587, 2018. 1, 2

[13] A. Dourado, T. E. de Campos, H. S. Kim, and A. Hilton. EdgeNet: Semantic scene completion from RGB-D images. *ArXiv*, abs/1908.02893, 2019. 4

[14] M. Firman, O. M. Aodha, S. J. Julier, and G. J. Brostow. Structured prediction of unobserved voxels from a single depth image. *Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 5431–5440, 2016. 2

[15] M. Garbade, J. Sawatzky, A. Richard, and J. Gall. Two stream 3D semantic scene completion. *Conference on Computer Vision and Pattern Recognition Workshops (CVPRW)*, pages 416–425, 2019. 1, 2, 4, 5, 7, 11

[16] A. Geiger, P. Lenz, and R. Urtasun. Are we ready for autonomous driving? The KIITI vision benchmark suite. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 3354–3361, 2012. 4

[17] A. Geiger and C. Wang. Joint 3D object and layout inference from a single RGB-D image. In *German Conference on Pattern Recognition (GCPR)*, 2015. 2

[18] B. Graham, M. Engelcke, and L. van der Maaten. 3D semantic segmentation with submanifold sparse convolutional networks. *Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 9224–9232, 2018. 2, 3

[19] Y.-X. Guo and X. Tong. View-volume network for semantic scene completion from a single depth image. In *International Joint Conference on Artificial Intelligence (IJCAI)*, 2018. 2, 3

[20] M. Jaritz, R. de Charette, É. Wirbel, X. Perrotton, and F. Nashashibi. Sparse and dense data with CNNs: Depth completion and semantic segmentation. In *International Conference on 3D Vision (3DV)*, pages 52–60, 2018. 7

[21] M. M. Kazhdan, M. Bolitho, and H. Hoppe. Poisson surface reconstruction. In *Symposium on Geometry Processing (SGP)*, 2006. 2

[22] B.-S. Kim, P. Kohli, and S. Savarese. 3D scene understanding by voxel-CRF. *International Conference on Computer Vision (ICCV)*, pages 1425–1432, 2013. 2

[23] J. Kober, J. A. Bagnell, and J. Peters. Reinforcement learning in robotics: A survey. *International Journal of Robotics Research (IJRR)*, 32:1238 – 1274, 2013. 1

[24] J. Li, Y. Liu, D. Gong, Q. Shi, X. Yuan, C. Zhao, and I. D. Reid. RGBD based dimensional decomposition residual network for 3D semantic scene completion. *Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 7685–7694, 2019. 2

[25] L.-J. Li, R. Socher, and L. Fei-Fei. Towards total scene understanding: Classification, annotation and segmentation in an automatic framework. *Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 2036–2043, 2009. 1, 2

[26] S. Liu, Y. Hu, Y. Zeng, Q. Tang, B. Jin, Y. Han, and X. Li. See and think: Disentangling semantic scene completion. In *NeurIPS*, 2018. 1, 2, 3, 4, 11

[27] Y. W. Liu, J. Li, Q. Yan, X. Yuan, C.-X. Zhao, I. Reid, and C. Cadena. 3D gated recurrent fusion for semantic scene completion. *ArXiv*, abs/2002.07269, 2020. 2

[28] Z. Liu, H. Tang, Y. Lin, and S. Han. Point-voxel CNN for efficient 3D deep learning. In *NeurIPS*, 2019. 1

[29] A. Milioto, I. Vizzo, J. Behley, and C. Stachniss. RangeNet ++: Fast and accurate LiDAR semantic segmentation. *International Conference on Intelligent Robots and Systems (IROS)*, pages 4213–4220, 2019. 2, 4, 11

[30] R. A. Newcombe, S. Izadi, O. Hilliges, D. Molyneaux, D. Kim, A. J. Davison, P. Kohli, J. Shotton, S. Hodges, and A. W. Fitzgibbon. KinectFusion: Real-time dense surface mapping and tracking. *International Symposium on Mixed and Augmented Reality*, pages 127–136, 2011. 2

[31] J. J. Park, P. Florence, J. Straub, R. A. Newcombe, and S. Lovegrove. DeepSDF: Learning continuous signed distance functions for shape representation. *Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 165–174, 2019. 2

[32] C. R. Qi, L. Yi, H. Su, and L. J. Guibas. Pointnet++: Deep hierarchical feature learning on point sets in a metric space. In *NeurIPS*, 2017. 2

[33] G. Riegler, A. O. Ulusoy, H. Bischof, and A. Geiger. OctNetFusion: Learning depth fusion from data. In *International Conference on 3D Vision (3DV)*, pages 57–66, 2017. 2

[34] L. Roldão, R. de Charette, and A. Verroust-Blondet. 3D surface reconstruction from voxel-based LiDAR data. *Intelligent Transportation Systems Conference (ITSC)*, pages 2681–2686, 2019. 2

[35] S. Satkin and M. Hebert. 3DNN: Viewpoint invariant 3D geometry matching for scene understanding. *International Conference on Computer Vision (ICCV)*, pages 1873–1880, 2013. 2

[36] S. Song, F. Yu, A. Zeng, A. X. Chang, M. Savva, and T. A. Funkhouser. Semantic scene completion from a single depth image. *Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 190–198, 2017. 1, 2, 4, 5, 6, 7, 10, 11

[37] S. Thrun and B. Wegbreit. Shape from symmetry. *International Conference on Computer Vision (ICCV)*, 2:1824–1831 Vol. 2, 2005. 2

[38] D. W. F. van Krevelen and R. Poelman. A survey of augmented reality technologies, applications and limitations. *International Journal of Virtual Reality (IJVR)*, 9:1–20, 2010. 1

[39] Y. Wang, D. J. Tan, N. Navab, and F. Tombari. Adversarial semantic scene completion from a single depth image. In *International Conference on 3D Vision (3DV)*, pages 426–434, 2018. 2

[40] W. Yuan, T. Khot, D. Held, C. Mertz, and M. Hebert. Pcn: Point completion network. *2018 International Conference on 3D Vision (3DV)*, pages 728–737, 2018. 2

[41] J. Zhang, H. Zhao, A. Yao, Y. Chen, L. Zhang, and H. Liao. Efficient semantic scene completion network with spatial group convolution. In *European Conference on Computer Vision (ECCV)*, 2018. 1, 2, 3, 4

[42] P. Zhang, W. Liu, Y. Lei, H. Lu, and X. Yang. Cascaded context pyramid for full-resolution 3D semantic scene completion. *International Conference on Computer Vision (ICCV)*, pages 7800–7809, 2019. 4

[43] K. Zimmermann, T. Petrícek, V. Salanský, and T. Svoboda. Learning for active 3D mapping. *International Conference on Computer Vision (ICCV)*, pages 1548–1556, 2017. 2
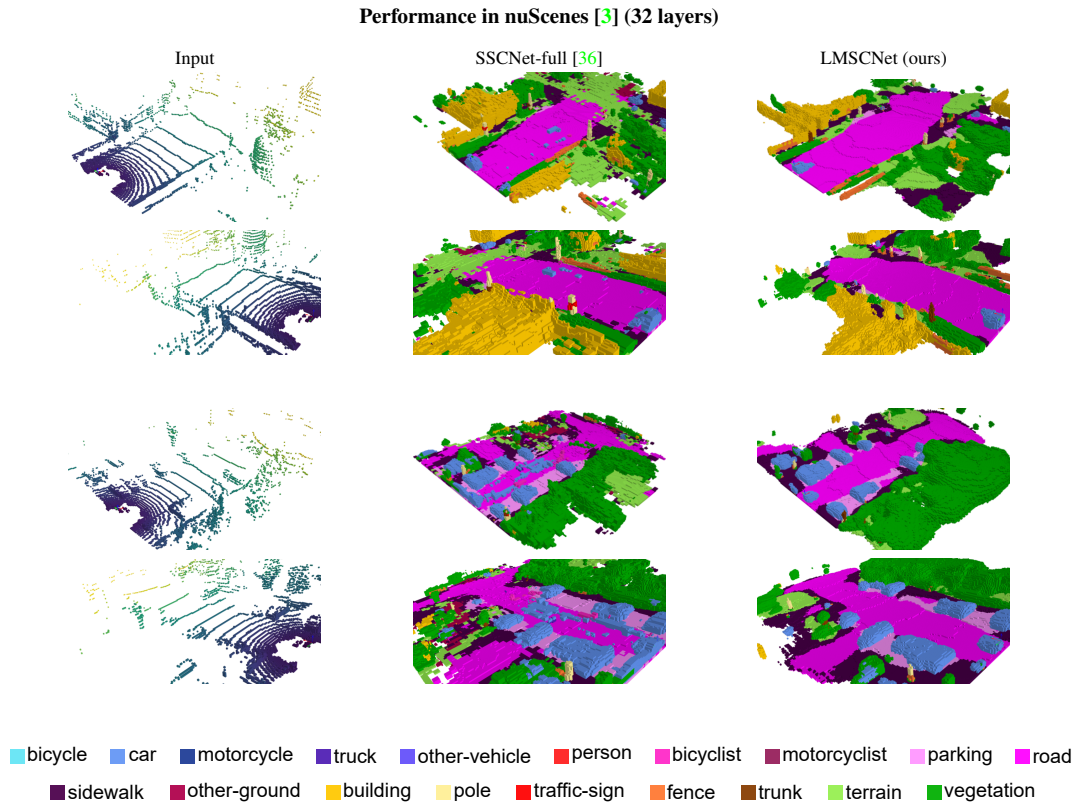
**Performance in SemanticKITTI [1] (64 layers)**

Input　　　　　SSCNet-full [36]　　　　LMSCNet (ours)　　　　Ground Truth



**Performance in nuScenes [3] (32 layers)**

Input　　　　　SSCNet-full [36]　　　　LMSCNet (ours)



■ bicycle　■ car　■ motorcycle　■ truck　■ other-vehicle　■ person　■ bicyclist　■ motorcyclist　■ parking　■ road
■ sidewalk　■ other-ground　■ building　■ pole　■ traffic-sign　■ fence　■ trunk　■ terrain　■ vegetation
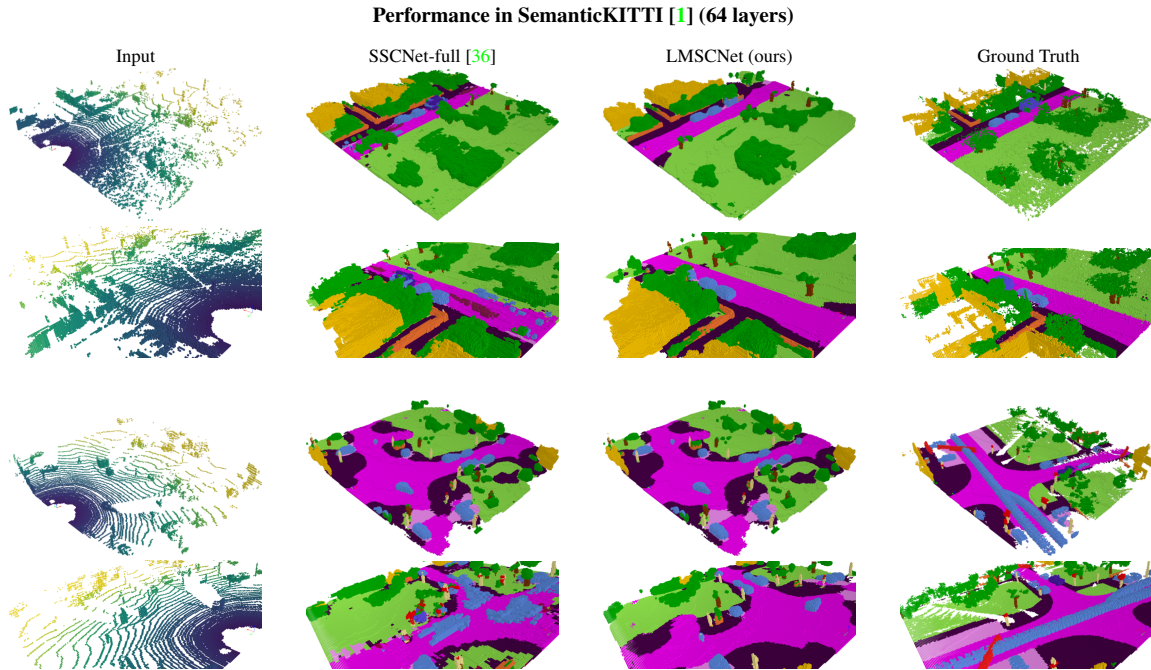
Figure 9: Additional qualitative 3D semantic completion results of our method in both SemanticKITTI [1] (top rows) and nuScenes [3] (bottom rows). Each pair of rows shows a single scene with different viewpoints. Groundtruth images are not shown for nuScenes due to the absence of point-wise semantic labels.

## A. Technical details

### A.1. Baselines implementations

Hereafter, we provide additional details on the implementation of the baselines listed in main article Tab. 1.

**TS3D baselines.** We compare our method with 3 variants of the Two Stream 3D (TS3D) network as reported in [1], which originate from [15]. As in their original work, TS3D uses an additional RGB modality, TS3D+DNet and TS3D+DNet+RangeNet use instead LiDAR intensity. The network is modified in two ways: first, by directly using projected semantic labels to the input grid obtained by a LiDAR-based semantic segmentation network [29] (TS3D+DNet); and secondly, by exchanging the 3D-CNN backbone by SATNet [26] (TS3D+DNet+SATNet). The semantic labels obtained from the 2D branch on the 3 variants are one-hot encoded and lifted to the 3D grid resulting on a $(N+1)\times H\times W\times D$ input tensor.

**SSCNet baselines.** Following the practice in [1], we use SSCNet [36] without the flipped TSDF as input encoding. However, while [1] only compares with SSCNet, which predictions are $1/4$ of the original input resolution, we also propose SSCNet-full, which outputs fullsize predictions. This is done by applying a $4\times4\times4$ transpose convolution to the last layer of the network to retrieve original dimensions. For data balancing, we use their strategy by randomly subsampling occluded free space, conserving a 2:1 free-occupied ratio.

### A.2. Architecture comparison

Fig. 10 provides additional insight about the benefit of each architecture, where the top rightmost corresponds to the best speed-performance trade-off. Even though SSCNet-full achieves faster inference than our method at the original scale, the performance is slightly lower and more noisy as observed in Fig. 9. TS3D+DNet+RangeNet achieves slightly higher performance but the inference time and the number of parameters are considerably higher as seen in main article Tab. 3. Considering this, our network keeps the best speed-performance trade-off. The interest of the coarser scale inferences of our method can be highlighted by the considerably lower inference times and high performances.

## B. Qualitative results

In Fig. 9 further qualitative results of our method in both SemanticKITTI [1] and nuScenes [3] datasets are provided. Notice our network performs smoother and less
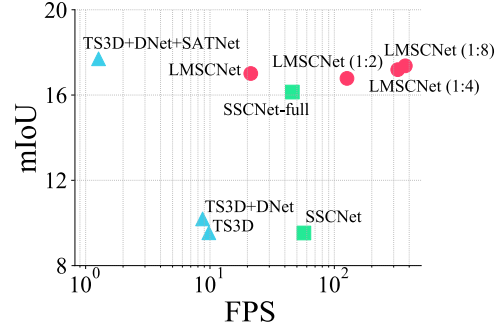


Figure 10: Network speed (FPS) and performance (mIoU). FPS are shown in log-scale. Notice that our method keeps good performance and fast inference time, this is specially noticeable for our coarse scale versions LMSCNet (1:x).

noisy reconstructions. Even though the ground-truth in SemanticKITTI accumulates scans of dynamic objects as seen in rows 3-4, our network reconstructs the vehicles correctly. This can be explained by the abundance of parked vehicles in the dataset. Performance in nuScenes can be seen in rows 5 to 8. It can be observed again the smoothness of the reconstruction when compared to SSCNet-full, with less noisy objects in the middle of the road. Notice that the network has been trained on SemanticKITTI, which explains the high vegetation predictions in nuScenes. We refer the reader to the supplementary video for more qualitative insights.