

# DREAM: Dual-Memory Reasoning for Anticipatory 3D Perception in Autonomous Driving

Anonymous CVPR submission

Paper ID 12

## Abstract

001 *Autonomous driving systems must not only interpret the cur-*  
002 *rent scene but also anticipate perceptual challenges such as*  
003 *occlusions, sparse observations, and recurring failure pat-*  
004 *terns. However, most existing LiDAR semantic segmenta-*  
005 *tion methods operate in a frame-centric manner, process-*  
006 *ing each scan independently and discarding past contextual*  
007 *experience, which limits robustness under dynamic scene*  
008 *changes. We propose DREAM, a Dual-memory REASONING*  
009 *framework for continual LiDAR seMantic segmentation that*  
010 *reframes perception as an experience-driven and anticipa-*  
011 *tory process. DREAM maintains two complementary mem-*  
012 *ory banks in a strictly online setting: a latent memory*  
013 *that stores compact semantic abstractions of previously ob-*  
014 *served scenes, and an error memory that records repre-*  
015 *sentations associated with uncertain predictions. At each*  
016 *timestep, relevant memory entries are retrieved via cosine*  
017 *similarity and integrated into the feature space through a*  
018 *lightweight modulation mechanism, enabling the model to*  
019 *reinforce consistent semantic patterns while suppressing re-*  
020 *curring failure modes. The backbone remains frozen and no*  
021 *past scans are replayed, ensuring computational efficiency*  
022 *and bounded memory growth. Extensive experiments on*  
023 *multiple large-scale LiDAR benchmarks demonstrate that*  
024 *DREAM achieves state-of-the-art performance, with con-*  
025 *sistent improvements on dynamic and small-scale objects,*  
026 *highlighting the effectiveness of persistent and error-aware*  
027 *memory for robust long-horizon perception.*

## 028 1. Introduction

029 Semantic scene understanding is a core capability for au-  
030 tonomous systems operating in complex environments such  
031 as autonomous driving. Among perception tasks, seman-  
032 tic segmentation provides dense structured representations  
033 of the surrounding scene that are critical for downstream  
034 planning and decision-making. However, reliable auton-  
035 omy requires perception systems that remain robust under

challenging sensing conditions such as occlusions, sparse 036  
measurements, and dynamic scene changes. Despite strong 037  
progress in deep learning-based perception, most LiDAR 038  
semantic segmentation methods remain fundamentally re- 039  
active, processing each scan independently or within short 040  
temporal windows [7]. This frame-centric formulation pre- 041  
vents perception systems from leveraging past experience 042  
to anticipate recurring perceptual challenges. 043

In real-world deployment, LiDAR observations arrive as 044  
a continuous stream where each scan provides only a par- 045  
tial and noisy view of the environment. Measurements are 046  
often affected by occlusions, sensor noise, and long-range 047  
sparsity. Under such conditions, treating frames indepen- 048  
dently [4, 13, 21, 27, 29, 32] leads to brittle predictions and 049  
temporal inconsistencies, particularly for distant or rarely 050  
observed structures. In contrast, human perception relies 051  
on accumulated experience: past observations inform fu- 052  
ture interpretation and help anticipate recurring patterns of 053  
uncertainty [14]. Enabling perception models to accumu- 054  
late and exploit such experience remains a key challenge 055  
for robust long-horizon scene understanding. 056

Three-dimensional representations provide a natural 057  
foundation for storing perceptual experience. Unlike im- 058  
age projections, 3D space preserves metric structure, allow- 059  
ing memory to maintain consistent object scales and spatial 060  
relationships across viewpoints. This enables stable seman- 061  
tic priors, more uniform representation across near and far 062  
regions, and improved reasoning about occlusions and mo- 063  
tion. Regions that become temporarily unobserved due to 064  
occlusions can therefore be retained in memory without in- 065  
terference from newly observed surfaces. However, main- 066  
taining spatial memory alone does not fully address a key 067  
limitation of deployed perception systems: the tendency to 068  
make similar errors under recurring sensing conditions re- 069  
peatedly. 070

Recent work has begun exploring memory-augmented 071  
perception for LiDAR segmentation. Earlier approaches, 072  
such as Cylinder3D [32] and Meta-RangeSeg [25], primar- 073  
ily focus on improving feature representations and spatial- 074  
temporal aggregation for semantic segmentation, but they 075

do not explicitly maintain persistent memory across observations. MemorySeg [12], for example, maintains a persistent voxel memory that aggregates information across streaming LiDAR scans and improves understanding of long-range and occluded regions. While effective, such approaches primarily accumulate spatial context and treat memory as a passive repository of past observations. In practice, perception systems frequently encounter recurring failure patterns, such as systematic uncertainty under severe sparsity or occlusion. Conventional architectures lack mechanisms to explicitly represent these failure modes and use them to guide future predictions.

To address this limitation, we propose **DREAM**, a dual-memory framework for continual LiDAR semantic segmentation that augments perception with both perceptual and error-aware memory. DREAM maintains two compact latent memory banks operating in a strictly online setting. A *latent perceptual memory* stores compressed semantic descriptors of previously observed scenes, while an *error memory* stores latent representations associated with uncertain predictions. At each timestep, a global latent descriptor extracted from the current scan retrieves relevant entries from both memory banks using cosine-similarity addressing. The retrieved latent and error memories are projected into the dense feature space and integrated through a lightweight modulation mechanism prior to classification. This allows the network to reinforce reliable semantic patterns while suppressing conditions that previously produced low-confidence predictions.

The proposed framework DREAM, processes one scan at a time, maintains bounded memory capacity, and avoids computationally expensive attention mechanisms or long-horizon backpropagation. The backbone encoder remains frozen, while lightweight memory modulation layers condition perception on retrieved experience. This design enables temporally aware semantic segmentation while maintaining efficiency suitable for real-time autonomous driving scenarios. The effectiveness of the proposed DREAM architecture is substantiated through comprehensive evaluations on three benchmark datasets: SemanticKITTI[2], nuScenes[3], and PandaSet[26].

In summary, the key contributions of this work are as follows:

- Proposal of DREAM, a Dual-Memory Reasoning for Anticipatory 3D Perception for Autonomous Driving.
- We introduce an error-aware memory mechanism that explicitly models recurring uncertainty patterns and helps mitigate future perceptual failures
- Proposal of a Dual Memory Modulation scheme for continual LiDAR semantic segmentation that integrates perceptual and error-aware memories to enable experience-conditioned perception, one of its first kinds.

## 2. Related Works

In this section, we review existing LiDAR semantic segmentation methods and discuss approaches that incorporate temporal reasoning for streaming LiDAR perception.

### 2.1. LiDAR Semantic Segmentation

LiDAR semantic segmentation methods can be broadly categorized based on the underlying data representation: point-based [11, 15, 16, 24], projection-based [6, 27, 31], voxel-based [4, 10, 32], and hybrid approaches that combine multiple representations [21, 28].

Point-based methods [11, 15, 16, 24] operate directly on raw point clouds and learn point-wise features using permutation-invariant networks. While effective, many of these approaches rely on aggressive downsampling to maintain computational efficiency. For example, KPConv [24] introduces deformable point convolutions that significantly improve segmentation performance, but at the cost of increased computational complexity.

Projection-based methods [6, 27, 31] project 3D point clouds into structured 2D representations such as range-view or bird’s-eye-view images. These approaches benefit from efficient 2D convolutional networks and often achieve real-time inference. However, the projection process introduces geometric distortions and information loss, which can limit segmentation accuracy.

Voxel-based approaches [4, 10, 32] convert point clouds into structured 3D grids and process them using sparse convolutional networks [5, 22]. These methods have demonstrated strong performance by preserving spatial structure while enabling efficient computation. Nevertheless, discretization into voxels may lead to loss of fine geometric details when voxel resolution is reduced.

To leverage the strengths of different representations, recent works combine multiple modalities [21, 28–30]. For instance, SPVNAS [21] integrates point-based and voxel-based features within a unified architecture, while RPVNet [28] employs a tri-branch network that jointly processes range-view, point-wise, and voxel features. These hybrid approaches improve segmentation accuracy by exploiting complementary geometric representations.

### 2.2. Temporal LiDAR Reasoning

To incorporate temporal context, several works attempt to aggregate information from multiple LiDAR frames. Early approaches [18, 20, 25] integrate a small number of past frames to enhance segmentation accuracy. For example, SpSequenceNet [20] introduces cross-frame attention and interpolation mechanisms to combine features from consecutive scans, while MetaRangeSeg [25] incorporates residual depth information from previous frames. Although these methods improve short-term consistency, they oper-

ate within limited temporal windows and discard historical context beyond a few frames.

Other works explore recurrent formulations for temporal state maintenance. Duerr et al. [8] propose a recurrent model that aligns range-view features across streaming scans, while MemorySeg [12] introduces a persistent latent voxel memory for improved segmentation in long-range and occluded regions. StrObe [9] similarly maintains a multi-scale BEV memory for incremental object detection. However, range-view and BEV representations are prone to projection artifacts and occlusion conflicts, and none of these approaches explicitly model recurring perceptual failures.

In contrast, DREAM augments LiDAR segmentation with persistent latent memory designed for long-horizon perception. Rather than storing raw observations or geometric maps, our framework maintains compact representations capturing both semantic experience and historical uncertainty. Retrieving relevant past contexts while suppressing recurring failure patterns enables robust, anticipatory perception, making DREAM one of the first approaches to handle LiDAR segmentation entirely within the latent space.

### 3. Methodology

Most LiDAR segmentation networks process each scan independently, ignoring temporal context critical for handling occlusions, sparsity, and recurring scene configurations. To address this, we introduce **DREAM** (Fig. 1), a dual-memory reasoning framework for continual LiDAR semantic segmentation. DREAM treats segmentation as an experience-driven process where structured latent memories evolve, enabling anticipatory reasoning over streaming data in a strictly online regime.

#### 3.1. Problem Formulation

Refer to Fig. 1 for the detailed architecture diagram of the proposed DREAM framework. We model segmentation as a recurrent reasoning process. At each timestep  $t$ , the model receives a LiDAR scan  $X_t$  and produces semantic predictions  $\hat{Y}_t$  conditioned on both the current input and an accumulated memory state  $\mathcal{M}_{t-1}$ . The system operates online: past scans are never revisited, and only compact latent representations are retained.

Given a temporally ordered sequence  $\{X_t\}_{t=1}^T$ , the model defines the mapping

$$\hat{Y}_t = f(X_t, \mathcal{M}_{t-1}; \theta), \quad (1)$$

where  $\theta$  denotes learnable parameters. After prediction, the memory is updated as

$$\mathcal{M}_t = \mathcal{U}(\mathcal{M}_{t-1}, X_t, \hat{Y}_t). \quad (2)$$

The central challenge is to design a compact yet expressive representation for  $\mathcal{M}_t$ , along with mechanisms for retrieval and integration into perception.

#### 3.2. Perceptual Encoding

Each raw LiDAR scan  $X_t$  is projected into a range-view encoder  $I_t \in \mathbb{R}^{5 \times H \times W}$  containing range, spatial coordinates  $(x, y, z)$ , and intensity. A convolutional backbone  $g(\cdot)$  extracts hierarchical features:

$$F_t, Z_t = g(I_t), \quad (3)$$

where  $F_t \in \mathbb{R}^{C \times H' \times W'}$  are dense spatial features for segmentation and  $Z_t \in \mathbb{R}^{D \times h \times w}$  are higher-level latent features. The backbone is frozen to stabilize representations and decouple long-term reasoning from low-level feature drift.

We summarize the latent representation via global average pooling:

$$\bar{z}_t = \text{GAP}(Z_t) \in \mathbb{R}^D, \quad (4)$$

which serves as the interface between perception and memory.

#### 3.3. Latent Memory Representation

LiDAR observations are inherently partial, with objects often sparsely sampled, occluded, or incomplete across individual scans. Relying solely on the current frame leads to unstable predictions, particularly for distant or rarely observed objects. To address this, we introduce a latent perceptual memory that stores compressed semantic descriptors from past observations, enabling the model to retrieve relevant contextual information during inference.

Formally, we maintain a fixed-capacity latent memory bank

$$\mathcal{M}_t^L = \{m_i^L \in \mathbb{R}^D\}_{i=1}^{K_L}, \quad (5)$$

where each entry encodes a compressed semantic abstraction of a previously observed scene. Given the global descriptor  $\bar{z}_t$  extracted from the current scan, the system retrieves relevant past contexts through content-based addressing using cosine similarity:

$$s_i = \frac{\bar{z}_t^\top m_i^L}{\|\bar{z}_t\| \|m_i^L\|}. \quad (6)$$

The top- $k$  most similar entries are selected to form the retrieved latent context

$$M_t^L \subset \mathcal{M}_{t-1}^L. \quad (7)$$

By retrieving semantically related past descriptors, the model can incorporate previously observed contextual information when interpreting the current scan. This mechanism provides a compact form of long-term contextual memory that improves robustness to sparse observations, partial visibility, and recurring scene configurations.

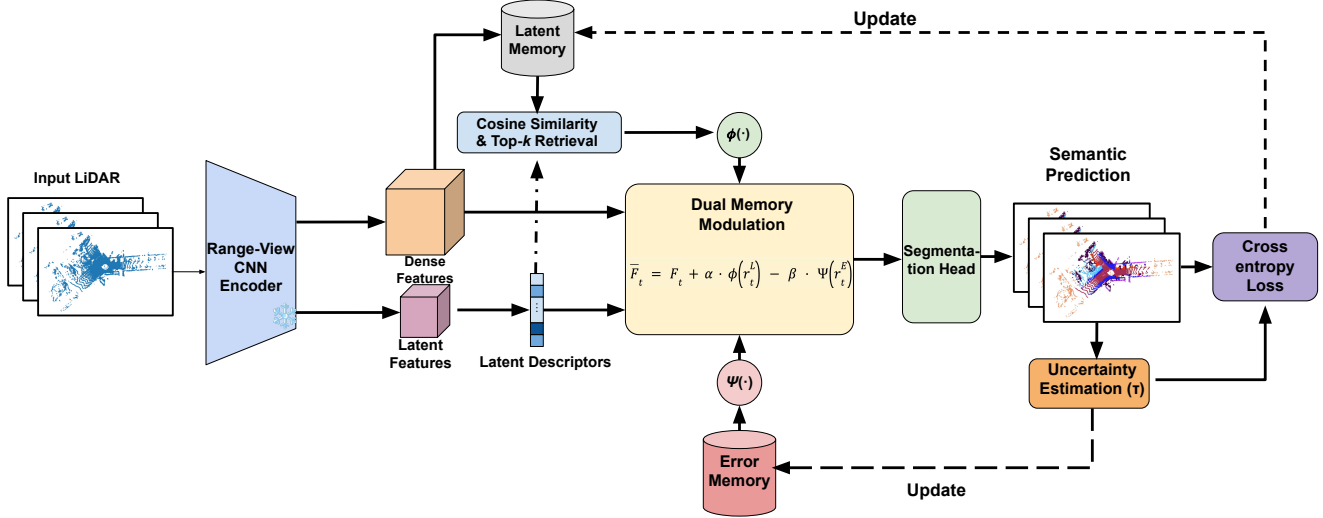


Figure 1. Overview of the proposed DREAM architecture. The network processes each LiDAR scan through hierarchical feature encoding and dual-memory reasoning modules that leverage past experience to guide current perception.

### 270 3.4. Error Memory Representation

271 Complementary to perceptual memory, we introduce an er-  
272 ror memory bank

$$273 \mathcal{M}_t^E = \{m_j^E \in \mathbb{R}^D\}_{j=1}^{K_E}, \quad (8)$$

274 which explicitly encodes latent representations associated  
275 with unreliable predictions. Let  $p_t(c|u)$  denote the pre-  
276 dicted class probability at location  $u$ . Uncertain regions are  
277 identified using a confidence threshold  $\tau$ :

$$278 \mathcal{E}_t = \{Z_t(u) \mid \max_c p_t(c|u) < \tau\}. \quad (9)$$

279 Latent vectors from  $\mathcal{E}_t$  are aggregated using mean pooling  
280 to produce a compact descriptor  $\bar{e}_t \in \mathbb{R}^D$ , which is then  
281 inserted into  $\mathcal{M}_t^E$ . While  $\mathcal{M}_t^L$  captures stable semantic  
282 abstractions,  $\mathcal{M}_t^E$  accumulates representations of recurring  
283 perceptual failures. By explicitly modeling these patterns,  
284 the system can anticipate similar uncertainty conditions in  
285 future observations and adjust its predictions accordingly.

### 286 3.5. Dual Memory Modulation

287 The retrieved latent memories  $M_t^L$  and error memories  $M_t^E$   
288 are integrated into the dense feature space prior to classifi-  
289 cation. Error memory retrieval follows the same content-  
290 based addressing mechanism as the latent memory, where  
291 cosine similarity between the global descriptor  $\bar{z}_t$  and stored  
292 error memory entries is computed and the most similar en-  
293 try is selected using top-1 retrieval. Let  $\phi(\cdot)$  and  $\psi(\cdot)$  denote  
294 learnable projection functions mapping latent vectors to the  
295 dense feature space. We define the modulated features as

$$296 \tilde{F}_t = F_t + \alpha \cdot \phi(M_t^L) - \beta \cdot \psi(M_t^E), \quad (10)$$

297 where  $\alpha$  and  $\beta$  control the strength of positive reinforcement  
298 and error suppression, respectively. In practice,  $\alpha$  and  $\beta$   
299 are treated as fixed hyperparameters selected via validation to  
300 balance contextual enhancement against over-suppression.  
301 Empirically, moderate values ensure stable training and  
302 consistent gains across datasets. The additive formulation  
303 preserves spatial alignment while injecting global semantic  
304 priors derived from past experience. Crucially, modulation  
305 occurs before classification, allowing retrieved memories to  
306 influence perception before predictions are produced.

### 307 3.6. Segmentation and Optimization

308 The modulated features  $\tilde{F}_t$  are fed to a lightweight segmen-  
309 tation head  $h(\cdot)$ :

$$310 \hat{Y}_t = h(\tilde{F}_t). \quad (11)$$

311 Training minimizes a class-weighted cross-entropy loss

$$312 \mathcal{L}_t = \text{CE}(\hat{Y}_t, Y_t), \quad (12)$$

313 where only the segmentation head and modulation layers  
314 are optimized, while the backbone remains frozen.

### 315 3.7. Online Memory Update

316 After prediction, the global descriptor  $\bar{z}_t$  is inserted into  
317  $\mathcal{M}_t^L$ , and the aggregated error descriptor  $\bar{e}_t$  derived from  
318 uncertain latents  $\mathcal{E}_t$  is inserted into  $\mathcal{M}_t^E$ . Both memories  
319 operate with fixed capacity using a replacement strategy  
320 (e.g., FIFO), ensuring bounded memory growth while re-  
321 taining long-term experience in compressed form.

### 322 3.8. Memory-Conditioned Perception

323 DREAM formulates LiDAR semantic segmentation as a  
324 perception problem conditioned on accumulated context

325 from streaming observations. Rather than processing each  
326 scan independently, the model maintains two latent mem-  
327 ory banks,  $\mathcal{M}_t^L$  and  $\mathcal{M}_t^E$ , which store compact representa-  
328 tions of previously observed semantic structures and his-  
329 torical low-confidence predictions. At each timestep, a  
330 global latent descriptor retrieves relevant entries through  
331 content-based addressing, and the retrieved memories mod-  
332 ulate the dense feature representation prior to classification.  
333 This mechanism allows the model to incorporate contex-  
334 tual information from past observations while maintaining  
335 a feed-forward inference pipeline. As a result, DREAM  
336 integrates long-horizon contextual cues into LiDAR per-  
337 ception in a computationally efficient manner, enabling  
338 temporally aware semantic segmentation in streaming au-  
339 tonomous driving scenarios.

## 340 4. Experimental Setup

341 We evaluate the proposed framework on three large-  
342 scale autonomous driving datasets: SemanticKITTI [2],  
343 nuScenes [3], and PandaSet [26]. These datasets are col-  
344 lected using different LiDAR sensors and capture diverse  
345 urban environments, making them suitable for evaluating  
346 perception systems operating in streaming conditions. Our  
347 experiments focus on assessing the effectiveness of persis-  
348 tent memory for long-horizon perception and its ability to  
349 improve robustness under sparse observations, occlusions,  
350 and dynamic scene changes. In addition to comparisons  
351 with existing methods, we perform ablation studies to an-  
352 alyze the contribution of each component of the proposed  
353 dual-memory design.

### 354 4.1. Datasets

355 We evaluate on three large-scale autonomous driving  
356 datasets. **SemanticKITTI** [2] comprises 22 sequences cap-  
357 tured with a Velodyne HDL-64E in Germany, with se-  
358 quences 0–10 for training and 11–21 for testing. La-  
359 bels are mapped to 19 classes (single-scan) and 25 classes  
360 (multi-scan), where we primarily report multi-scan results  
361 as our method models temporal context. **nuScenes** [3]  
362 contains 1000 driving scenes from Boston and Singapore  
363 (700/150/150 train/val/test split), captured with a Velodyne  
364 HDL-32E at 2 Hz, with 31 classes mapped to 16 categories.  
365 **PandaSet** [26] provides temporally continuous sequences  
366 recorded in Silicon Valley using a Hesai Pandar64, with  
367 103 sequences and fine-grained classes grouped into 14 cat-  
368 egories following [8].

### 369 4.2. Evaluation Protocol

370 During evaluation, each sequence is processed in temporal  
371 order from the first frame to the last. The memory state  
372 is updated online during inference, allowing the model to  
373 accumulate experience across the sequence. Performance  
374 is evaluated using Intersection-over-Union (IoU) [17], with

mean IoU (mIoU) across all semantic classes reported as  
the primary metric. In addition, we report seen-class mIoU,  
which averages IoU only over classes present in the evalu-  
ated sequence, as well as class-wise IoU scores for detailed  
analysis.

### 4.3. Implementation Details

380 Training proceeds in two stages: the segmentation head is  
381 first trained for 50 epochs in a single-scan setting, after  
382 which the dual-memory mechanism and modulation lay-  
383 ers are fine-tuned for 20 additional epochs in an online  
384 recurrent manner. The backbone remains frozen through-  
385 out. Each LiDAR scan is projected into a  $64 \times 1024$   
386 range-view with five channels (range,  $x$ ,  $y$ ,  $z$ , intensity).  
387 A pre-trained ResNet-18 extracts dense features from the  
388 third residual stage and a 512-dimensional global descrip-  
389 tor from the final stage for memory retrieval. The latent  
390 and error memory banks store  $K_L=64$  and  $K_E=32$  vec-  
391 tors ( $D=512$ ) respectively, using cosine similarity for top-1  
392 retrieval and FIFO replacement. The error memory is up-  
393 dated via low-confidence regions. Optimization uses Adam  
394 ( $\text{lr} = 5 \times 10^{-4}$ ) with class-weighted cross-entropy and  
395 mixed precision training. Augmentations include global  
396 scaling, random translation, and vertical-axis rotation. The  
397 model is trained on 8 NVIDIA A100 GPUs (80GB) using  
398 distributed data parallelism. 399

### 4.4. Quantitative Analysis

400 We evaluate DREAM on SemanticKITTI, nuScenes, and  
401 PandaSet, with results summarized in Tables 1, 2, and 3. 402

403 On SemanticKITTI (Tab. 1), DREAM achieves 61.7%  
404 mIoU, outperforming MemorySeg (58.3%) and Meta-  
405 RangeSeg (49.7%). Gains are most pronounced in dy-  
406 namic categories (moving cars, bicyclists, pedestrians, mo-  
407 torcyclists) and structural classes (building, trunk, terrain),  
408 reflecting the benefit of temporal reasoning via the dual-  
409 memory design.

410 On nuScenes (Tab. 2), DREAM achieves **81.8%** mIoU,  
411 with notable improvements in barrier, pedestrian, and traffic  
412 cone categories, suggesting the memory mechanism effec-  
413 tively addresses sparsity and scale challenges.

414 On PandaSet (Tab. 3), DREAM achieves **71.0%** mIoU,  
415 performing consistently across foreground classes (truck,  
416 person, traffic sign) and background classes (sidewalk, ter-  
417 rain), confirming generalization across sensor configura-  
418 tions and urban environments.

419 Across all benchmarks, persistent latent memory enables  
420 robust segmentation under sparse observations, occlusions,  
421 and dynamic scene changes.

### 4.5. Comparison Against State-of-the-Art

422 We compare DREAM with prior LiDAR-only methods  
423 across all three benchmarks (Tab. 1–3). 424

Method	mIoU	Car	Bicycle	Motorcycle	Truck	Other Vehicle	Person	Bicyclist	Motorcyclist	Road	Parking	Sidewalk	Other Ground	Building	Fence	Vegetation	Trunk	Terrain	Pole	Traffic Sign	Car (m)	Bicyclist (m)	Person (m)	Motorcyclist (m)	Other Vehicle (m)	Truck (m)
TangentConv [23]	34.1	84.9	2.0	18.2	21.1	18.5	1.6	0.0	0.0	83.9	38.3	64.0	15.3	85.8	49.1	79.5	43.2	56.7	36.4	31.2	40.3	1.1	6.4	1.9	30.1	42.2
DarkNet53Seg [1]	41.6	84.1	30.4	32.9	20.2	20.7	7.5	0.0	0.0	91.6	64.9	75.3	27.5	85.2	56.5	78.4	50.7	64.8	38.1	53.3	61.5	14.1	15.2	0.2	28.9	37.8
KPCConv [24]	51.2	93.7	44.9	47.2	42.5	38.6	21.6	0.0	0.0	86.5	58.4	70.5	26.7	90.8	64.5	84.6	70.3	66.0	57.0	53.9	69.4	67.4	67.5	47.2	4.7	5.8
Cylinder3D [32]	52.5	94.6	67.6	63.8	41.3	38.8	12.5	1.7	0.2	90.7	<b>65.0</b>	74.5	<b>32.3</b>	92.6	66.0	<b>85.8</b>	72.0	68.9	63.1	61.4	74.9	68.3	65.7	11.9	0.1	0.0
SpSequenceNet [20]	43.1	88.5	24.0	26.2	29.2	22.7	6.3	0.0	0.0	90.1	57.6	73.9	27.1	91.2	66.8	84.0	66.0	65.7	50.8	48.7	53.2	41.2	26.2	36.2	2.3	0.1
TemporallLidarSeg [8]	47.0	92.1	47.7	40.9	39.2	35.0	14.4	0.0	0.0	91.8	59.6	<b>75.8</b>	23.2	89.8	63.8	82.3	62.5	64.7	52.6	60.4	68.2	42.8	40.4	12.9	12.4	2.1
TemporallLatticeNet [19]	47.1	91.6	35.4	36.1	26.9	23.0	9.4	0.0	0.0	91.5	59.3	75.3	27.5	89.6	65.3	84.6	66.7	70.4	57.2	60.4	59.7	41.7	51.0	48.8	5.9	0.0
Meta-RangeSeg [25]	49.7	90.8	50.0	49.5	29.5	34.8	16.6	0.0	0.0	90.8	62.9	74.8	26.5	89.8	62.1	82.8	65.7	66.5	56.2	64.5	69.0	60.4	57.9	22.0	16.6	2.6
MemorySeg [12]	58.3	94.0	68.3	68.8	51.3	40.9	27.0	0.3	2.8	89.9	64.3	74.8	29.2	92.2	69.3	84.8	75.1	70.1	65.5	<b>68.5</b>	71.7	<b>74.4</b>	<b>71.7</b>	73.9	15.1	13.6
<b>DREAM (Ours)</b>	<b>61.7</b>	<b>97.6</b>	<b>69.1</b>	<b>68.9</b>	<b>55.7</b>	<b>43.1</b>	<b>29.0</b>	<b>1.7</b>	<b>3.2</b>	<b>93.5</b>	62.6	74.8	29.2	92.2	69.3	84.8	<b>77.3</b>	<b>71.0</b>	<b>65.6</b>	68.4	71.5	70.3	71.5	<b>74.3</b>	<b>37.6</b>	<b>49.7</b>

Table 1. Comparison with state-of-the-art LiDAR semantic segmentation methods on the **SemanticKITTI** multi-scan test benchmark. (m) denotes moving classes. Metrics are reported in mIoU (%).

Method	mIoU	FW-mIoU	Barrier	Bicycle	Bus	Car	Construction	Motorcycle	Pedestrian	Traffic Cone	Trailer	Truck	Drivable	Other Flat	Sidewalk	Terrain	Manmade	Vegetation
PolarNet [31]	69.4	87.4	72.2	16.8	77.0	86.5	51.1	69.7	64.8	54.1	69.7	63.5	96.6	67.1	77.7	72.1	87.1	84.5
Cylinder3D [32]	77.2	89.9	82.8	29.8	84.3	89.4	63.0	<b>79.3</b>	77.2	73.4	84.6	69.1	97.7	<b>70.2</b>	80.3	75.5	90.4	87.6
SPVCNN [21]	77.4	89.7	80.0	30.0	91.9	90.8	64.7	79.0	75.6	70.9	81.0	74.6	97.4	69.2	80.0	76.1	89.3	87.1
(AF) <sup>2</sup> -S3Net [4]	78.3	88.5	78.9	<b>52.2</b>	89.9	84.2	<b>77.4</b>	74.3	77.3	72.0	83.9	73.8	97.1	66.5	77.5	74.0	87.7	86.8
MemorySeg [12]	80.6	91.4	84.9	40.2	91.2	92.4	71.2	73.5	<b>85.9</b>	77.8	88.0	76.4	<b>97.9</b>	69.0	81.2	77.6	92.6	89.7
<b>DREAM (Ours)</b>	<b>81.8</b>	<b>91.4</b>	<b>86.3</b>	40.2	<b>92.3</b>	<b>94.5</b>	<b>79.4</b>	73.9	84.7	<b>78.0</b>	<b>88.4</b>	<b>79.1</b>	96.2	68.3	<b>81.9</b>	76.3	<b>93.5</b>	<b>90.4</b>

Table 2. Comparison with state-of-the-art LiDAR-only semantic segmentation methods on the **nuScenes** test set [3]. Metrics are reported in mIoU (%). Best results are highlighted in **bold**.

425 On SemanticKITTI (Tab. 1), DREAM outperforms  
 426 Cylinder3D (52.5%), Meta-RangeSeg (49.7%), and Mem-  
 427 orySeg (58.3%), with the largest gains in dynamic  
 428 classes (moving cars, bicyclists, pedestrians, motorcy-  
 429 clists), demonstrating effective temporal modeling via dual  
 430 memory.

431 On nuScenes (Tab. 2), DREAM sets a new state-of-the-  
 432 art among LiDAR-only methods, surpassing MemorySeg  
 433 (80.6%) and (AF)<sup>2</sup>-S3Net (78.3%), with notable improve-  
 434 ments in small-object categories such as barrier, pedestrian,  
 435 and traffic cone.

436 On PandaSet (Tab. 3), DREAM achieves the best over-  
 437 all performance over MemorySeg (70.3%) and SPVCNN  
 438 (64.7%), with consistent gains across foreground and back-  
 439 ground categories including truck, person, and sidewalk.

440 Across all datasets, persistent latent memory reasoning  
 441 consistently outperforms both single-frame and short-term  
 442 temporal fusion approaches.

## 443 4.6. Qualitative Analysis

### 444 4.6.1. Temporal Consistency Analysis

445 Fig. 3 presents qualitative results across four consecutive  
 446 LiDAR frames ( $t-2, t-1, t, t+1$ ). The top row shows the  
 447 input LiDAR intensity visualization, while the middle and

bottom rows illustrate the ground-truth semantic labels and  
 predictions from DREAM, respectively.

450 Despite the sparse and partially observed nature of indi-  
 451 vidual LiDAR scans, the predicted semantic structure re-  
 452 mains stable across time. In particular, large structural  
 453 classes such as road and vegetation are consistently seg-  
 454 mented across frames. This temporal stability demonstrates  
 455 the benefit of the proposed memory-based reasoning mech-  
 456 anism, which allows the model to leverage historical context  
 457 when interpreting the current observation.

### 458 4.6.2. Occlusion Recovery

459 Fig. 4 illustrates the behavior of the DREAM framework  
 460 under partial occlusions. The top row shows the input Li-  
 461 DAR intensity visualization, where the highlighted region  
 462 exhibits varying point density across consecutive frames  
 463 due to occlusion and sparse observations.

464 The middle row shows the ground-truth semantic la-  
 465 bels, while the bottom row presents predictions produced  
 466 by DREAM. Despite incomplete observations in individual  
 467 LiDAR scans, DREAM maintains temporally consistent se-  
 468 mantic predictions with 61.7% mIoU and 81.8% mIoU on  
 469 SemanticKITTI and nuScenes datasets, respectively, across  
 470 frames, demonstrating the proposed memory mechanism’s

Method	mIoU	Car	Bicycle	Motorcycle	Truck	Other Vehicle	Person	Road	Road Barriers	Sidewalk	Building	Vegetation	Terrain	Background	Traffic Sign
SqueezeSegv3 [27]	55.7	92.8	24.1	18.0	36.5	54.3	63.0	91.1	11.9	71.3	86.2	85.0	61.3	63.2	20.6
SalsaNext [6]	57.8	92.1	40.7	31.7	28.7	56.2	69.0	90.0	22.6	67.1	85.6	83.4	58.5	63.3	20.6
TemporalLidarSeg [8]	60.0	93.7	33.6	38.0	37.1	59.9	72.0	91.1	14.6	70.6	88.2	88.4	63.8	68.4	20.7
SPVCNN [21]	64.7	95.8	38.1	46.3	44.0	74.1	78.6	91.2	<b>28.3</b>	70.3	87.2	87.5	61.5	67.5	35.6
MemorySeg [12]	70.3	97.2	60.2	58.4	62.9	74.3	82.6	92.1	27.7	74.1	89.4	90.7	<b>64.9</b>	72.8	36.4
<b>DREAM (Ours)</b>	<b>71.0</b>	<b>97.8</b>	<b>60.5</b>	<b>59.1</b>	<b>64.0</b>	<b>75.4</b>	<b>82.9</b>	<b>93.6</b>	28.1	<b>75.4</b>	<b>90.7</b>	<b>89.2</b>	64.5	<b>74.6</b>	<b>37.1</b>

Table 3. Comparison with state-of-the-art LiDAR semantic segmentation methods on the PandaSet test set [26]. Metrics are reported in mIoU (%). Best results are highlighted in **bold**.

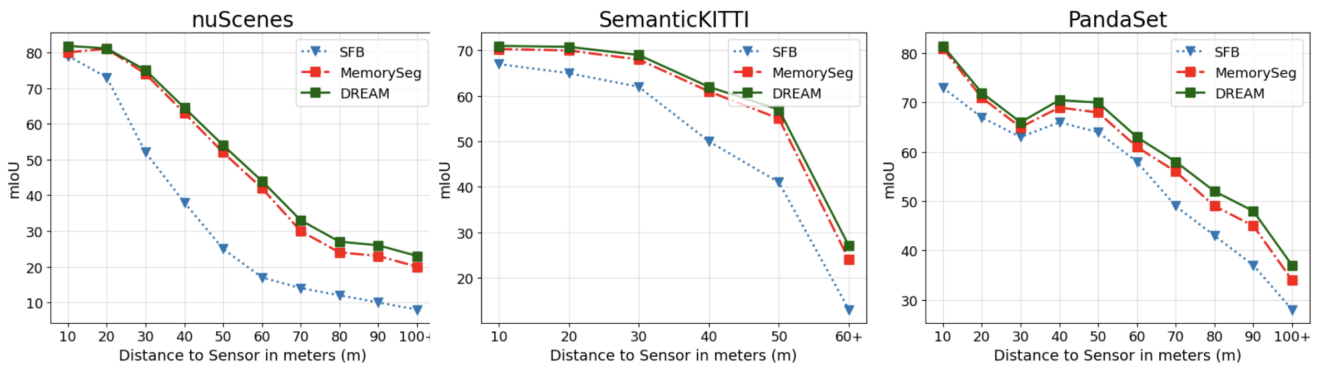


Figure 2. Distance-based performance comparison of SFB (dotted blue), MemorySeg (solid red), and DREAM (solid green) across nuScenes, SemanticKITTI, and PandaSet. DREAM demonstrates consistently improved robustness at mid- and long-range distances, with gains increasing as point sparsity and occlusion become more severe.

Method	Params (M)	FPS
SFB	14.2	29
DREAM (Full)	16.8	25

Table 4. Parameter count and inference speed comparison on SemanticKITTI. DREAM introduces a modest increase in the number of parameters and runtime while maintaining real-time performance.

471 ability to recover semantic structure in occluded or sparsely  
472 observed regions.

## 473 4.7. Ablation Studies

### 474 4.7.1. Importance of the Memory

475 To evaluate the impact of persistent memory, we compare  
476 DREAM against a single-frame baseline (SFB) that  
477 removes both latent perceptual memory and error mem-  
478 ory components. The comparison is performed using the  
479 distance-based evaluation shown in Fig. 2. Results indicate

480 that incorporating memory significantly improves segmen-  
481 tation performance, particularly at larger distances where  
482 LiDAR measurements are sparse and noisy. This demon-  
483 strates that maintaining a persistent latent representation of  
484 past observations helps the model recover semantic context  
485 that may be missing in individual frames.

### 486 4.7.2. Distance-based Analysis

487 Fig. 2 highlights the benefit of persistent dual memory at  
488 longer ranges, where LiDAR point clouds become increas-  
489 ingly sparse. We compare DREAM against a single-frame  
490 baseline (SFB) that removes both memory modules.

491 DREAM consistently outperforms the SFB across all  
492 distance bins, with the most pronounced gains at 80m:  
493 26.47% mIoU on nuScenes and 52.33% mIoU on PandaSet,  
494 where point density is lowest. The progressive improve-  
495 ment with distance confirms that persistent memory be-  
496 comes more beneficial as perceptual difficulty increases,  
497 enhancing semantic consistency through compact scene rep-  
498 resentations and explicit uncertainty modeling.

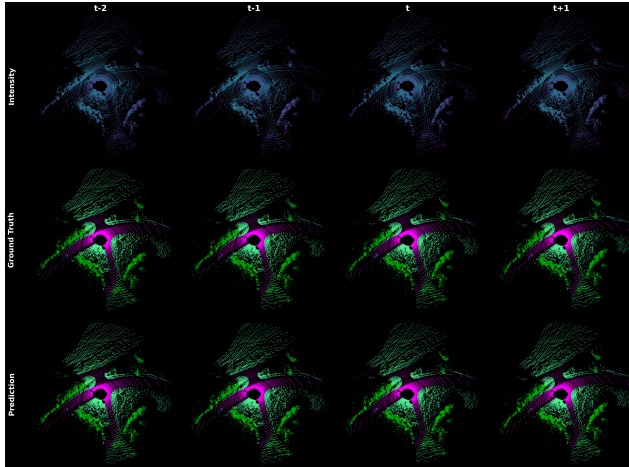


Figure 3. Temporal qualitative results across consecutive LiDAR frames. Columns correspond to four consecutive frames ( $t - 2$ ,  $t - 1$ ,  $t$ ,  $t + 1$ ). The top row shows the input LiDAR intensity visualization, the middle row shows the ground-truth semantic labels, and the bottom row shows the predictions produced by the proposed DREAM framework.

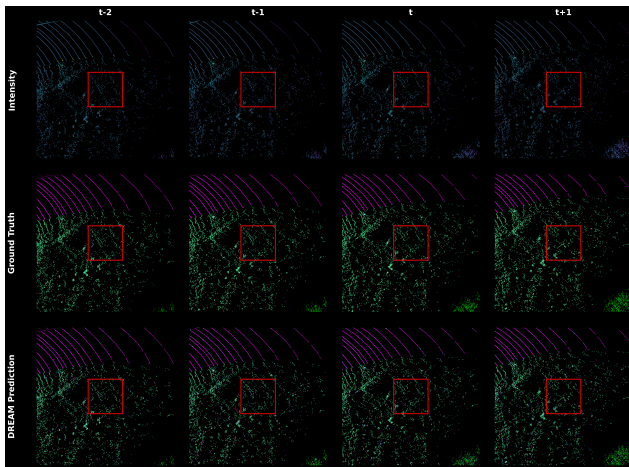


Figure 4. Qualitative example illustrating occlusion recovery across consecutive frames. As objects become partially occluded, DREAM maintains stable semantic predictions by leveraging temporal context stored in latent memory.

#### 499 4.7.3. Effect of Error Memory

500 We evaluate the contribution of the proposed error memory  
 501 on the SemanticKITTI multi-scan validation set (Tab. 5).  
 502 Compared to the single-frame baseline (SFB), introducing  
 503 latent perceptual memory improves mIoU from 78.6% to  
 504 81.0%, highlighting the benefit of persistent semantic ag-  
 505 gregation. Adding the error memory further increases per-  
 506 formance to 81.8%, with more pronounced gains in long-  
 507 range regions, where mIoU improves from 24.8% to 27.3%.  
 508 These results indicate that explicitly modeling recurrent

Method	mIoU (%)
SFB	58.4
DREAM w/o Error Memory	60.9
DREAM (Full)	<b>61.5</b>

Table 5. Ablation study on the SemanticKITTI multi-scan validation set. Latent perceptual memory accounts for the majority of the improvement over the single-frame baseline (SFB). The proposed error memory further enhances performance, particularly in long-range regions (60m+), where point sparsity and occlusions are more severe.

low-confidence patterns complements spatial memory and enhances robustness under sparse and occluded observations.

#### 4.7.4. Parameter Count.

DREAM introduces only a modest computational overhead compared to the single-frame baseline (SFB). As shown in Tab. 4, the proposed dual-memory design increases the parameter count from 14.2M to 16.8M (approximately 18%) and reduces inference speed from 29 FPS to 25 FPS (approximately 15%), while maintaining real-time performance.

## 5. Conclusion

We presented DREAM, a dual-memory reasoning framework for continual LiDAR semantic segmentation that reformulates perception as a temporally grounded process. DREAM maintains two complementary memory banks: a perceptual memory that accumulates stable semantic abstractions and an error memory that records recurring failure patterns. Retrieved memories modulate dense features before classification, reinforcing consistent semantics while suppressing previously observed errors. The framework operates strictly online, avoids scan replay, and keeps the backbone frozen, yielding a bounded-memory design suitable for real-world deployment. Experiments on SemanticKITTI, nuScenes, and PandaSet demonstrate state-of-the-art performance, with the most significant gains in dynamic and small-object categories. DREAM advances LiDAR segmentation from stateless prediction to memory-augmented reasoning, offering a scalable foundation for long-horizon autonomous perception.

## References

- [1] Jens Behley, Martin Garbade, Andres Milioto, Jan Quenzel, Sven Behnke, C. Stachniss, and Juergen Gall. A dataset for semantic segmentation of point cloud sequences. *ArXiv*, abs/1904.01416, 2019. 6
- [2] 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. In *Proc.*

- 548 of the *IEEE/CVF International Conf. on Computer Vision*  
549 (*ICCV*), 2019. 2, 5
- 550 [3] Holger Caesar, Varun Bankiti, Alex H. Lang, Sourabh Vora,  
551 Venice Erin Liong, Qiang Xu, Anush Krishnan, Yuxin Pan,  
552 Giancarlo Baldan, and Oscar Beijbom. nuscenes: A mul-  
553 timodal dataset for autonomous driving. *2020 IEEE/CVF*  
554 *Conference on Computer Vision and Pattern Recognition*  
555 (*CVPR*), pages 11618–11628, 2019. 2, 5, 6
- 556 [4] Ran Cheng, Ryan Razani, Ehsan Moeen Taghavi, Enxu Li,  
557 and Bingbing Liu. (af)<sup>2</sup>-s3net: Attentive feature fusion with  
558 adaptive feature selection for sparse semantic segmentation  
559 network. *2021 IEEE/CVF Conference on Computer Vision*  
560 *and Pattern Recognition (CVPR)*, pages 12542–12551, 2021.  
561 1, 2, 6
- 562 [5] Christopher Bongsoo Choy, JunYoung Gwak, and Silvio  
563 Savarese. 4d spatio-temporal convnets: Minkowski convo-  
564 lutional neural networks. *2019 IEEE/CVF Conference on*  
565 *Computer Vision and Pattern Recognition (CVPR)*, pages  
566 3070–3079, 2019. 2
- 567 [6] Tiago Cortinhal, George Tzelepis, and Eren Erdal Aksoy.  
568 Salsanext: Fast semantic segmentation of lidar point clouds  
569 for autonomous driving. *ArXiv*, abs/2003.03653, 2020. 2, 7
- 570 [7] Ayush Dewan and Wolfram Burgard. Deeptemporalseg:  
571 Temporally consistent semantic segmentation of 3d lidar  
572 scans. *2020 IEEE International Conference on Robotics and*  
573 *Automation (ICRA)*, pages 2624–2630, 2019. 1
- 574 [8] Fabian Duerr, Mario Pfaller, Hendrik Weigel, and Jürgen  
575 Beyerer. Lidar-based recurrent 3d semantic segmentation  
576 with temporal memory alignment. *2020 International Con-*  
577 *ference on 3D Vision (3DV)*, pages 781–790, 2020. 3, 5, 6,  
578 7
- 579 [9] Davi Frossard, Simon Suo, Sergio Casas, James Tu, Rui Hu,  
580 and Raquel Urtasun. Strobe: Streaming object detection  
581 from lidar packets. *ArXiv*, abs/2011.06425, 2020. 3
- 582 [10] Yuenan Hou, Xinge Zhu, Yuexin Ma, Chen Change Loy, and  
583 Yikang Li. Point-to-voxel knowledge distillation for lidar se-  
584 mantic segmentation. *2022 IEEE/CVF Conference on Com-*  
585 *puter Vision and Pattern Recognition (CVPR)*, pages 8469–  
586 8478, 2022. 2
- 587 [11] Qingyong Hu, Bo Yang, Linhai Xie, Stefano Rosa, Yulan  
588 Guo, Zhihua Wang, Agathoniki Trigoni, and A. Markham.  
589 Randla-net: Efficient semantic segmentation of large-scale  
590 point clouds. *2020 IEEE/CVF Conference on Computer Vi-*  
591 *sion and Pattern Recognition (CVPR)*, pages 11105–11114,  
592 2019. 2
- 593 [12] Enxu Li, Sergio Casas, and Raquel Urtasun. Memoryseg:  
594 Online lidar semantic segmentation with a latent memory.  
595 *2023 IEEE/CVF International Conference on Computer Vi-*  
596 *sion (ICCV)*, pages 745–754, 2023. 2, 3, 6, 7
- 597 [13] Andres Milioto, Ignacio Vizzo, Jens Behley, and C. Stach-  
598 niss. Rangenet ++: Fast and accurate lidar semantic segmen-  
599 tation. *2019 IEEE/RSJ International Conference on Intel-*  
600 *ligent Robots and Systems (IROS)*, pages 4213–4220, 2019.  
601 1
- 602 [14] Alireza Modirshanechi, Sophia Becker, Johanni Brea, and  
603 Wulfram Gerstner. Surprise and novelty in the brain. *Current*  
604 *Opinion in Neurobiology*, 82:102758, 2023. 1
- [15] C. Qi, Hao Su, Kaichun Mo, and Leonidas J. Guibas. Point-  
net: Deep learning on point sets for 3d classification and seg-  
mentation. *2017 IEEE Conference on Computer Vision and*  
*Pattern Recognition (CVPR)*, pages 77–85, 2016. 2
- [16] C. Qi, L. Yi, Hao Su, and Leonidas J. Guibas. Pointnet++:  
Deep hierarchical feature learning on point sets in a metric  
space. *ArXiv*, abs/1706.02413, 2017. 2
- [17] Seyed Hamid Rezaatofghi, Nathan Tsoi, JunYoung Gwak,  
Amir Sadeghian, Ian D. Reid, and Silvio Savarese. General-  
ized intersection over union: A metric and a loss for bound-  
ing box regression. *2019 IEEE/CVF Conference on Com-*  
*puter Vision and Pattern Recognition (CVPR)*, pages 658–  
666, 2019. 5
- [18] Peer Schütt, Radu Alexandru Rosu, and Sven Behnke.  
Abstract flow for temporal semantic segmentation on the  
permutohedral lattice. *2022 International Conference on*  
*Robotics and Automation (ICRA)*, pages 5139–5145, 2022.  
2
- [19] Peer Schutt, Radu Alexandru Rosu, and Sven Behnke. Ab-  
stract flow for temporal semantic segmentation on the per-  
mutohedral lattice. In *2022 International Conference on*  
*Robotics and Automation (ICRA)*, page 5139–5145. IEEE  
Press, 2022. 6
- [20] Hanyu Shi, Guosheng Lin, Hao Wang, Tzu-Yi Hung, and  
Zhenhua Wang. Spsequencenet: Semantic segmentation net-  
work on 4d point clouds. In *Proceedings of the IEEE/CVF*  
*Conference on Computer Vision and Pattern Recognition*  
*(CVPR)*, 2020. 2, 6
- [21] Haotian Tang, Zhijian Liu, Shengyu Zhao, Yujun Lin, Ji  
Lin, Hanrui Wang, and Song Han. Searching efficient 3d  
architectures with sparse point-voxel convolution. *ArXiv*,  
abs/2007.16100, 2020. 1, 2, 6, 7
- [22] Haotian Tang, Zhijian Liu, Xiuyu Li, Yujun Lin, and Song  
Han. Torchsparse: Efficient point cloud inference engine.  
*ArXiv*, abs/2204.10319, 2022. 2
- [23] Maxim Tatarchenko, Jaesik Park, Vladlen Koltun, and Qian-  
Yi Zhou. Tangent convolutions for dense prediction in 3d.  
*2018 IEEE/CVF Conference on Computer Vision and Pat-*  
*tern Recognition*, pages 3887–3896, 2018. 6
- [24] Hugues Thomas, C. Qi, Jean-Emmanuel Deschaud, Beatriz  
Marcotegui, François Goulette, and Leonidas J. Guibas. Kp-  
conv: Flexible and deformable convolution for point clouds.  
*2019 IEEE/CVF International Conference on Computer Vi-*  
*sion (ICCV)*, pages 6410–6419, 2019. 2, 6
- [25] Song Wang, Jianke Zhu, and Ruixiang Zhang. Meta-  
rangeseq: Lidar sequence semantic segmentation using mul-  
tiple feature aggregation. *IEEE Robotics and Automation*  
*Letters*, 7:9739–9746, 2022. 1, 2, 6
- [26] Pengchuan Xiao, Zhenlei Shao, Steven Hao, Zishuo Zhang,  
Xiaolin Chai, Judy Jiao, Zesong Li, Jian Wu, Kai Sun, Kun  
Jiang, Yunlong Wang, and Diange Yang. Pandaset: Ad-  
vanced sensor suite dataset for autonomous driving. *2021*  
*IEEE International Intelligent Transportation Systems Con-*  
*ference (ITSC)*, pages 3095–3101, 2021. 2, 5, 7
- [27] Chenfeng Xu, Bichen Wu, Zining Wang, Wei Zhan, Péter  
Vajda, Kurt Keutzer, and Masayoshi Tomizuka. Squeeze-  
segv3: Spatially-adaptive convolution for efficient point-

- 662 cloud segmentation. In *European Conference on Computer*  
663 *Vision*, 2020. 1, 2, 7
- 664 [28] Jianyun Xu, Ruixiang Zhang, Jian Dou, Yushi Zhu, Jie Sun,  
665 and Shiliang Pu. Rpvnet: A deep and efficient range-point-  
666 voxel fusion network for lidar point cloud segmentation.  
667 *2021 IEEE/CVF International Conference on Computer Vi-*  
668 *sion (ICCV)*, pages 16004–16013, 2021. 2
- 669 [29] Xu Yan, Jiantao Gao, Chaoda Zheng, Chaoda Zheng,  
670 Ruimao Zhang, Shenghui Cui, and Zhen Li. 2dpass: 2d pri-  
671 ors assisted semantic segmentation on lidar point clouds. In  
672 *European Conference on Computer Vision*, 2022. 1
- 673 [30] Maosheng Ye, Shuangjie Xu, Tongyi Cao, and Qifeng Chen.  
674 Drinet: A dual-representation iterative learning network for  
675 point cloud segmentation. *2021 IEEE/CVF International*  
676 *Conference on Computer Vision (ICCV)*, pages 7427–7436,  
677 2021. 2
- 678 [31] Yang Zhang, Zixiang Zhou, Philip David, Xiangyu Yue, Ze-  
679 rong Xi, and Hassan Foroosh. Polarnet: An improved grid  
680 representation for online lidar point clouds semantic segmen-  
681 tation. *2020 IEEE/CVF Conference on Computer Vision and*  
682 *Pattern Recognition (CVPR)*, pages 9598–9607, 2020. 2, 6
- 683 [32] Xinge Zhu, Hui Zhou, Tai Wang, Fangzhou Hong, Wei Li,  
684 Yuexin Ma, Hongsheng Li, Ruigang Yang, and Da Lin.  
685 Cylindrical and asymmetrical 3d convolution networks for  
686 lidar-based perception. *IEEE Transactions on Pattern Anal-*  
687 *ysis and Machine Intelligence*, 44:6807–6822, 2021. 1, 2,  
688 6