

# SVL: Spike-based Vision-language Pretraining for Efficient 3D Open-world Understanding

Xuerui Qiu<sup>1,2</sup>, Peixi Wu<sup>3</sup>, Yaozhi Wen<sup>1,2</sup>, Shaowei Gu<sup>1,2</sup>, Yuqi Pan<sup>1</sup>,  
Xinhao Luo<sup>1</sup>, Bo Xu<sup>1</sup>, Guoqi Li<sup>1\*</sup>

<sup>1</sup>Institute of Automation, Chinese Academy of Sciences,

<sup>2</sup>School of Future Technology, University of Chinese Academy of Sciences,

<sup>3</sup> University of Science and Technology of China

## Abstract

Spiking Neural Networks (SNNs) provide an energy-efficient way to extract 3D spatio-temporal features. However, existing SNNs still exhibit a significant performance gap compared to Artificial Neural Networks (ANNs) due to inadequate pre-training strategies. These limitations manifest as restricted generalization ability, task specificity, and a lack of multimodal understanding, particularly in challenging tasks such as multimodal question answering and zero-shot 3D classification. To overcome these challenges, we propose a Spike-based Vision-Language (SVL) pretraining framework that empowers SNNs with open-world 3D understanding while maintaining spike-driven efficiency. SVL introduces two key components: (i) Multi-scale Triple Alignment (MTA) for label-free triplet-based contrastive learning across 3D, image, and text modalities, and (ii) Re-parameterizable Vision-Language Integration (Rep-VLI) to enable lightweight inference without relying on large text encoders. Extensive experiments show that SVL achieves a top-1 accuracy of 85.4% in zero-shot 3D classification, surpassing advanced ANN models, and consistently outperforms prior SNNs on downstream tasks, including 3D classification (+6.1%), DVS action recognition (+2.1%), 3D detection (+1.1%), and 3D segmentation (+2.1%) with remarkable efficiency. Moreover, SVL enables SNNs to perform open-world 3D question answering, sometimes outperforming ANNs. To the best of our knowledge, SVL represents the first scalable, generalizable, and hardware-friendly paradigm for 3D open-world understanding, effectively bridging the gap between SNNs and ANNs in complex open-world understanding tasks. Code is available [Here](#).

## 1 Introduction

Bio-inspired Spiking Neural Networks (SNNs) offer an efficient approach to learning superior representations from sparse 3D geometric data (e.g., event streams and point clouds) [1], owing to their distinctive spike-driven nature [2] and spatio-temporal processing capabilities [3]. For instance, the Speck [4] chip uses event-by-event sparse processing to handle 3D input data, with operational power consumption as low as 0.7 mW. However, existing SNNs [5; 6; 7] exhibit a significant performance gap compared to ANNs, and remain task-dependent, lacking both generalizable representations and the ability to achieve multimodal understanding in 3D open-world scenarios.

For instance, when deploying SNNs in real-world scenarios [4], they may struggle to generalize to input data from unseen categories not present in the training set. This highlights the critical need to

---

\*Corresponding author

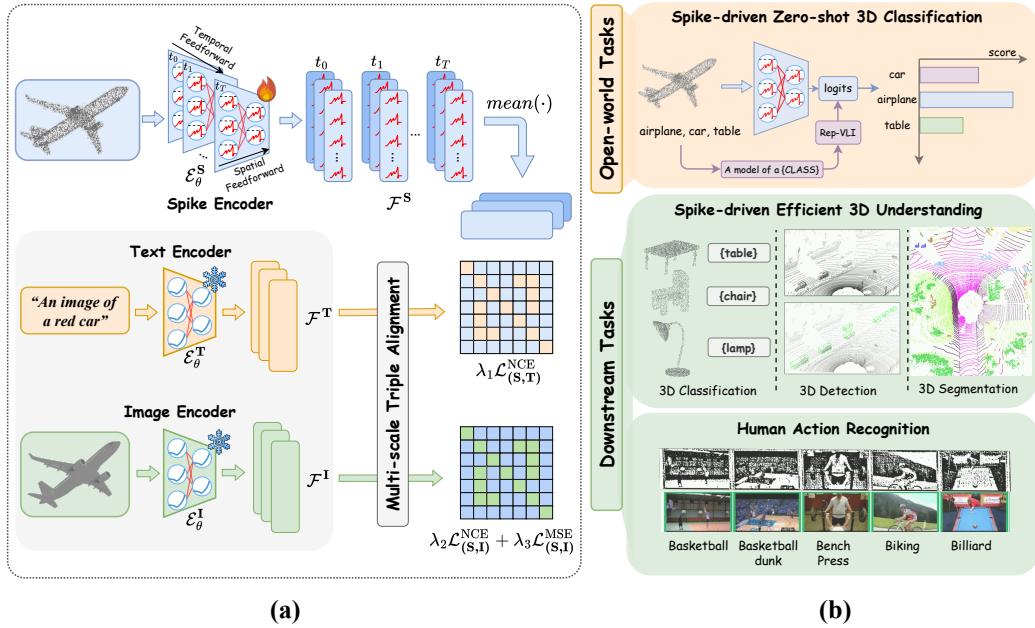


Figure 1: Overall architecture and applications of our SVL. (a) In pretraining, we proposed Multi-scale Triple Alignment (MTA) that jointly optimizes correlation alignment across text, image, and 3D inputs. (b) For downstream tasks, we propose Re-parameterizable Vision-Language Integration (Rep-VLI) to reparameterize the text embeddings generated by the text encoder into lightweight weights, enabling efficient spike-driven inference.

develop robust pretraining strategies to enhance the visual representation capabilities and adaptability of SNNs. Existing methods, such as STDP-based initialization [8] and knowledge distillation in SpikeBert and SpikeCLIP [9; 10; 11], refine spike-based representations, while SpikformerV2 and Spike-driven Transformer V3 [6; 7] employ masked image modeling to improve scalability. However, these approaches [8] either lose effectiveness as dataset complexity increases, demand substantial computational resources [7], which limits their feasibility for neuromorphic hardware deployment, or lack multimodal integration [9]. Moreover, pre-trained models often exhibit inadequate visual representation capabilities and limited transferability, restricting unified applicability to downstream tasks [6].

Another challenge is the limited availability of annotated 3D datasets, as the creation of such datasets is both labor-intensive and error-prone, rendering it often impractical for large-scale, real-world applications [12]. In response, Vision-Language Models (VLMs) [13] have been employed to explore the transfer of knowledge gleaned from extensive 2D datasets to facilitate open-world 3D understanding. However, most VLMs [14; 15; 16] depend heavily on large-scale text encoders during inference, which imposes substantial limitations on the practicality of hardware deployment.

In this paper, we introduce a universal Spike-based Vision-Language pretraining framework (SVL) that enhances SNNs' capability for open-world multimodal 3D understanding while maintaining efficient spike-driven inference. As shown in Fig. 1, our framework incorporates two key innovations: (i) Multi-scale Triple Alignment (MTA), which enables label-free triplet representation learning for capturing geometric properties of 3D data, and (ii) Re-parameterizable Vision-language Integration (Rep-VLI), which facilitates efficient deployment by addressing the computational overhead of text-driven inference. Moreover, our framework uses CLIP for its strong generalization. During pre-training, we freeze CLIP and train the 3D encoder by aligning 3D features with CLIP's textual and visual features via contrastive learning. The pre-trained 3D model can then be fine-tuned for downstream tasks. To demonstrate SVL's efficacy, we rigorously evaluated our models on a wide range of 3D tasks, including open-world understanding and downstream applications. We also highlight SVL's versatility through its potential in generative tasks like 3D object captioning and

open-world question answering, showcasing its broad applicability in 3D vision domains. Our main contribution can be summarized as:

- We introduce SVL, a universal spike-based vision-language pretraining framework that enables SNNs to understand open-world multimodal 3D information while maintaining efficient spike-driven inference.
- We propose two key technical innovations: (i) Multi-scale Triple Alignment (MTA), a label-free triplet learning mechanism for capturing geometric properties of 3D data across different scales, and (ii) Re-parameterizable Vision-language Integration (Rep-VLI), which achieves lightweight inference by reducing the computational overhead of text encoder.
- Extensive experiments on multiple benchmarks demonstrate the effectiveness of SVL, achieving state-of-the-art (SOTA) performance in both 3D open-world understanding and downstream tasks, as well as generative applications such as 3D object captioning and open-world question answering.

## 2 Related Works

**Pretraining Algorithms of SNNs** Numerous pretraining methods are proposed for spike-based representation learning. [8] utilized spike-timing-dependent plasticity [17] to initialize SNNs, enhancing the model’s robustness and training speed. While this approach has been successful on simple datasets with shallow networks, its effectiveness diminishes as the complexity of the datasets and networks increases. To address these issues, SpikeBert and SpikeCLIP [9; 10; 11] employ a two-stage knowledge distillation process from ANNs to enhance spike-based representations for complex downstream tasks. However, these methods rely on ANN weight initialization, limiting structural flexibility. Additionally, they use LayerNorm, which hinders neuromorphic hardware deployment. SpikformerV2 and Spike-driven Transformer V3 [6; 7; 18] apply a masked image modeling approach to address the performance degradation in SNNs as the model scales up. However, they require substantial storage and computational resources, and the lack of multimodal integration, particularly language guidance, limits their effectiveness in open-world understanding tasks.

**Vision-language Models (VLMs)** aim to align image and text embeddings for cross-modal transfer, with CLIP [13] being a seminal work that uses contrastive learning for zero-shot classification. Building on this foundation, subsequent methods have expanded cross-modal alignment to include other modalities. These approaches typically fall into two categories: dual-encoder and triple-encoder frameworks. Dual-encoder fine-tune both visual and textual encoders [11; 19; 20]. Triple-encoder frameworks incorporate additional modality-specific encoders [14; 15; 21], which combine triple models to achieve open-world understanding. This architecture is highly flexible, making it suitable for a variety of downstream tasks [12; 16]. However, triple-encoder frameworks still rely on large text encoders during inference, hindering hardware deployment.

**Efficient 3D Recognition** is to learn better representations from sparse and irregular 3D geometric data (e.g., event and point cloud). Currently, there exist two main approaches. Point-based methods [22; 23] utilize the PointNet series to directly extract geometric features from raw 3D data and make predictions. Voxel-based methods [24] first convert the 3D data into regular voxels and then use 3D sparse convolutions for feature extraction. Nevertheless, the improved accuracy is often accompanied by increased computational costs, limiting its applicability in practical systems. Hence, numerous studies [25; 26; 27] in the SNN field combine spiking neurons with point-based methods like PointNet [22] for low energy consumption edge computing. These methods are too simplistic to apply to various tasks and datasets. It was not until E-3DSNN [5] used spike sparse convolution to achieve high performance in numerous 3D tasks with spike-driven nature. Our approach is orthogonal to the aforementioned 3D encoders. Pre-training with SVL can enhance the visual extraction capabilities of the 3D encoders, thereby improving their performance in 3D recognition tasks. SVL can also endow them with the ability to achieve open-world understanding while retaining their spike-driven nature.

## 3 Preliminaries

**Spiking Neurons** are inspired by the dynamics of biological neurons [3; 28], which are the fundamental units of Spiking Neural Networks (SNNs). Among these, the Leaky Integrate-and-

Fire (LIF) neuron is the most widely used due to its balance between biological plausibility and computational efficiency [3]. We begin by translating the LIF spiking neuron into an iterative expression using the Euler method [29], which is described as follows:

$$u_i^{(\ell)}[t+1] = h_i^{(\ell)}[t] + f(w^{(\ell)}, x_i^{(\ell-1)}[t]), \quad (1)$$

$$s_i^{(\ell)}[t] = \Theta(u_i^{(\ell)}[t+1] - \vartheta), \quad (2)$$

$$h_i^{(\ell)}[t+1] = \beta u_i^{(\ell)}[t+1](1 - s_i^{(\ell)}[t]), \quad (3)$$

Here,  $\beta$  is the time constant  $t$  and  $i$  represents the time step and the neuron index in the  $\ell$ -th layer, respectively. The weight matrix  $w$  defines the synaptic connections between adjacent layers, while  $f(\cdot)$  is a function that denotes operations such as convolution (Conv) or fully connected (FC). The input is represented by  $x$ , and  $\Theta(\cdot)$  denotes the Heaviside step function. When the membrane potential  $u$  exceeds the firing threshold  $\vartheta$ , the LIF neuron generates a spike,  $s$ . Additionally,  $h$  represents the membrane potential after the spike event, which is scaled by a constant factor  $\beta$ .

Directly training the above LIF-based SNNs requires the use of backpropagation through time (BPTT) [29], resulting in a time complexity of  $\mathcal{O}(LT)$ , where  $L$  and  $T$  are the number of layers and time steps. This significantly increases both the training time and memory requirements. To mitigate this issue, we use the Integer LIF Spiking Neuron.

**Integer LIF Spiking Neuron** is incorporated into our SVL to reduce the quantization error, training time, and memory [6; 30; 31], which allows us to rewrite Eq. (2) as:

$$s_i^{(\ell)}[t] = \lfloor \text{clip}\{u_i^{(\ell)}[t], 0, D^t\} \rfloor, \quad (4)$$

where  $\lfloor \cdot \rfloor$  denotes the rounding operator,  $\text{clip}\{x, a, b\}$  confines  $x$  within range  $[a, b]$ , and  $D^t$  is a hyperparameter indicating the maximum emitted integer value by I-LIF. Moreover, I-LIF will emit integer values while pretraining and convert them into binary spikes by expanding the virtual timestep to ensure that the inference is spike-driven with only sparse addition.

**CLIP** employs a contrastive learning framework to forge associations between images and textual descriptions [13]. It leverages a vast dataset comprising 400 million image-text pairs, training an image encoder based on either a ResNet [32] or Vision Transformer [33] architecture, and a text encoder utilizing the Transformer model [34]. These encoders work together to project both images and text into a unified embedding space. This training paradigm enables CLIP to perform zero-shot classification, allowing it to identify images solely based on textual descriptions, even without category-specific training. In this study, we aim to utilize CLIP’s semantically rich feature space to support spike-based encoders, facilitating energy-efficient 3D open-world understanding with a spike-driven nature.

## 4 Method

Our primary goal is to develop a spike-based encoder that accurately captures the geometric properties of 3D input data and efficiently achieves a unified representation for open-world 3D understanding with the spike-driven nature. To this end, we construct our triplet dataset  $\{(D_i^t, I_1^t, T_1^t), (D_2^t, I_2^t, T_2^t), \dots, (D_n^t, I_n^t, T_n^t)\}$ , which consists of a 3D input  $D_i^t$ , an image  $I_i$ , and a text description  $T_i$  at  $t$  time step.

### 4.1 3D Input Representation

In this part, we present the 3D input representation, such as point clouds and event streams. Event streams, in particular, require special handling. We define them as  $E_i = (x_i, y_i, t_i, p_i)$ . Using a sliding window technique [23; 25], we convert event streams into an event cloud, formulated as:

$$E_i = (x_i, y_i, z_i) \quad \text{where} \quad z_i = \frac{t_i - t_{\min}}{t_{\max} - t_{\min}},$$

By doing so, we treat event streams as a distinct kind of spatio-temporal point cloud. This allows us to consider both point clouds and event streams as collections of points, denoted by  $D^t = \{\mathcal{P}, \mathcal{I}\}$ . This includes voxel sets  $D_k^t = \{\mathcal{P}_k^t, \mathcal{F}_k^t\}$ , where  $\mathcal{P}_k^t \in \mathbb{R}^3$  represents the 3D coordinates and  $\mathcal{I}_k^t \in \mathbb{R}^D$

indicate the features across  $d$  channels at the time step  $t$ . Following this, we utilize our I-LIF spiking neuron to encode these 3D inputs into spatio-temporal spike trains, which are then transmitted to the spike encoder.

## 4.2 Multi-scale Triple Alignment

To develop a unified representation for open-world 3D understanding, we introduce a multi-scale triple alignment (MTA) framework that jointly optimizes correlation alignment across text, image, and 3D inputs. This framework integrates both semantic spike-text alignment and fine-grained spike-image alignment. Specifically, the overall architecture of SVL, illustrated in Fig. 1, comprises three encoders: (i) Text Encoder ( $\mathcal{E}_\theta^T$ ): embeds text into text features  $\mathcal{F}^T \in \mathbb{R}^C$ ; (ii) Spike-based Encoder ( $\mathcal{E}_\theta^S$ ): transforms spike inputs into spike trains  $\mathcal{F}^S \in \mathbb{R}^{T \times C}$ . (iii) Image Encoder ( $\mathcal{E}_\theta^I$ ): encodes images into image features  $\mathcal{F}^I \in \mathbb{R}^C$ . Here,  $C$  represents the embedding dimension. These encoders collaboratively embed the triplet texts, spikes, and images into their respective feature spaces, facilitating comprehensive and fine-grained alignment across different modalities.

**Semantic Spike-Text Alignment** To leverage the open-world recognition capabilities of the pre-trained CLIP model [13], we align the spike firing rate  $\mathcal{F}^S/T$  with the text embeddings  $\mathcal{F}^T$  obtained from CLIP, using a spike-text tuple  $\mathcal{B}_i = \{T_i^t, \mathcal{D}_i^t\}$  as input. The core idea is to bring the feature centroids of 3D instances and their corresponding text prompts closer together in the embedding space. To achieve this, we compute the InfoNCE loss [35] between the mean spike trains and the text features, as follows:

$$\mathcal{L}_{(S,T)}^{\text{NCE}} = -\frac{1}{2|\mathcal{B}|} \sum_i^{|B|} \log \frac{e^{\tau \mathbf{x}_i \cdot \mathbf{y}_i}}{\sum_j^{|B|} e^{\tau \mathbf{x}_i \cdot \mathbf{y}_j}} + \log \frac{e^{\tau \mathbf{x}_i \cdot \mathbf{y}_i}}{\sum_j^{|B|} e^{\tau \mathbf{x}_j \cdot \mathbf{y}_i}}, \quad (5)$$

where  $\mathbf{x}_i = \frac{\mathcal{F}^S/T}{\|\mathcal{F}^S/T\|_2}$  and  $\mathbf{y}_i = \frac{\mathcal{F}^T}{\|\mathcal{F}^T\|_2}$  represent the normalized spike and text features, respectively. The indices  $i$  and  $j$  are used for sampling, the dot product ( $\cdot$ ) denotes cosine similarity between vectors, and  $\tau$  is a learnable temperature parameter.

**Fine-grained Spike-Image Alignment** A singular spike-text alignment fails to fully capture the semantic information embedded within both images and 3D data. To achieve a more comprehensive multimodal understanding, we further introduce an alignment between image and spike features. Specifically, we first employ the InfoNCE loss to align the image features, denoted as  $\mathcal{F}^I$ , with the average pulse signals, represented as  $\mathcal{F}^S/T$ . This alignment can be expressed as follows:

$$\mathcal{L}_{(S,I)}^{\text{NCE}} = -\frac{1}{2|\mathcal{C}|} \sum_i^{|C|} \log \frac{e^{\tau \mathbf{a}_i \cdot \mathbf{b}_i}}{\sum_j^{|C|} e^{\tau \mathbf{a}_i \cdot \mathbf{b}_j}} + \log \frac{e^{\tau \mathbf{a}_i \cdot \mathbf{b}_i}}{\sum_j^{|B|} e^{\tau \mathbf{a}_j \cdot \mathbf{b}_i}}, \quad (6)$$

where  $\mathcal{C}_i = \{I_i^t, \mathcal{D}_i^t\}$  a spike-image tuple,  $\mathbf{a}_i = \frac{\mathcal{F}^S/T}{\|\mathcal{F}^S/T\|_2}$  and  $\mathbf{b}_i = \frac{\mathcal{F}^I}{\|\mathcal{F}^I\|_2}$  represent the normalized spike and text features, respectively. However, this approach resulted in overly coarse alignment granularity, failing to account for the fine-grained and tightly coupled alignment between spikes and images. To address this, we incorporate the MSE loss on the basis of the InfoNCE loss to enhance the alignment granularity. The alignment objective between spike trains and images, which is formulated as follows:

$$\mathcal{L}_{(S,I)}^{\text{MSE}} = \sum_i^{|C|} \|\mathcal{F}_i^S - \mathcal{F}_i^I\|^2, \quad (7)$$

where  $\|\cdot\|^2$  is the  $\ell_2$  norm. Finally, we obtain the resultant total learning objective  $\mathcal{L}_{\text{total}}$  as the combination of  $\mathcal{L}_{(S,T)}$  and  $\mathcal{L}_{(S,I)}$ , where both alignments of semantic spike-text and fine-grained spike-image alignment are injected as:

$$\mathcal{L}_{\text{total}} = \lambda_1 \mathcal{L}_{(S,T)}^{\text{NCE}} + \lambda_2 \mathcal{L}_{(S,I)}^{\text{NCE}} + \lambda_3 \mathcal{L}_{(S,I)}^{\text{MSE}}, \quad (8)$$

where  $\lambda_1$ ,  $\lambda_2$ , and  $\lambda_3$  are hyperparameters that balance the influence of image features, text features, and spike trains. In this paper, the default setting of both is 1.

### 4.3 Re-parameterizable Vision-language Integration

After pretraining, the text prompt is applied to guide the learned spike-based representation to enable the following zero-shot transfer tasks. However, during inference, the traditional methods [14; 15] still require a high-parameter text encoder, which compromises the spike-driven nature and significantly impacts the deployment of neuromorphic hardware. To address this issue, we propose a novel structure called the re-parameterizable Vision-language module (Rep-VLI), which converts text embeddings into the weights of linear layers. It can be described as:

$$W_i^L = e^\tau \mathcal{E}_\theta^T(T_i^t), \quad (9)$$

where  $\mathcal{E}_\theta^T(\cdot)$  is the text encoder and  $W_i$  is the re-parameterized layer weights on the  $i$ -th input text prompt  $T_i^t$  at  $t$  time step. Then we concatenate the re-parameterized weights  $W_i$  by  $K$  text prompts to  $W^L \in \mathbb{R}^{K \times C}$ . And the classification logits are calculated with the spike trains  $\mathcal{F}^S$  and re-parameterized weights as:

$$\text{logits}_i = \text{softmax}\left[\frac{W^L}{T} \mathcal{E}_\theta^S(D_i^t)\right], \quad (10)$$

where  $\mathcal{E}_\theta^S(\cdot)$  is the spike-based encoder,  $T$  is the time step, and  $D_i^t$  is the input 3D data. Rep-VLI overcomes the limitations of traditional high-parameter models by removing the need for a text encoder during inference. It re-parameterizes text into a lightweight format aligned with the spike-driven framework, ensuring efficient deployment on neuromorphic hardware while supporting diverse tasks like 3D open-world detection and segmentation [36; 37] with spike-based encoders.

## 5 Experiments

To validate that our spike-based 3D encoder learns robust visual representations via SVL, we evaluate it on diverse 3D open-world tasks, including zero-shot classification and visual question answering. The pretrained encoder is also fine-tunable for downstream tasks like 3D classification, segmentation, detection, and action recognition. This section details the experimental setup, including backbones, datasets, and implementation, followed by quantitative results and ablation studies on modality count, time steps, and loss functions. For further details about implementation and the datasets, please refer to Appendix E and Appendix F, respectively.

**3D Backbone Networks** We evaluate three 3D backbones within our SVL framework: (i) Spike PointNet [26], a lightweight, energy-efficient spiking model for point clouds; (ii) E-3DSNN [5], a state-of-the-art model leveraging spike sparse convolution for sparse 3D inputs; and (iii) Spike PointFormer is our proposed Transformer-based SNN backbone for 3D data. Further details are provided in Appendix B. We use OpenCLIP’s [38] ViT-G/14 text encoder (frozen during pretraining) to define a shared feature space, into which the 3D modality is aligned. During inference, our Rep-VLI transforms text embeddings into compact weights, preserving the spike-driven nature of the above spike-based encoders.

### 5.1 3D Open-world Understanding

**Zero-shot classification** We evaluate the zero-shot classification performance of our models on the widely-used ModelNet40 [24] and the larger, more challenging Objaverse-LVIS [41]. Compared to other benchmarks, Objaverse-LVIS offers broader class coverage and a long-tailed distribution, providing a more realistic evaluation of open-world 3D understanding [16]. As shown in Tab. 1, our SVL-based E-3DSNN achieves 85.4% accuracy on ModelNet40 with only 17.7M parameters, outperforming both ANN and SNN baselines. This demonstrates SVL’s effectiveness in enhancing both accuracy and efficiency.

Specifically, compared to OpenShape and ULIP, our model achieves 85.4% accuracy (vs. 83.4% and 69.6%), consumes only 0.79 mJ of energy (vs. 161.8 mJ and 73.8 mJ), and uses fewer parameters (17.7M vs. 41.3M and 21.9M). It also delivers a 15.8% accuracy gain over ULIP-based PointBERT [14] while consuming just 11.4% of the energy. On Objaverse-LVIS, our model performs comparably to ULIP-2 [15] but with a  $204\times$  energy efficiency advantage. This is enabled by our Rep-VLI module, which reparameterizes text embeddings into compact, spike-compatible weights

Architecture	Model	Pre-train Method	Input	$T \times D$	Point+Text Param (M)	Energy (mJ)	Obj.	M40.
ANN	PointCLIP [19]	N/A	Image	N/A	25.5+57.3	24.9+20.3	1.9	20.2
	PointNet [22]	Openshape [16]	Point	N/A	3.47+202.5	20.1+71.7	24.4	74.9
	Point-Bert [39]		Point	N/A	21.9+202.5	83.2+71.7	43.2	82.8
	Sparseconv [40]		Voxel	N/A	5.3+202.5	0.61+71.7	31.7	78.8
			Voxel	N/A	41.3+202.5	2.13+71.7	43.4	83.4
SNN	Point-Bert [39]	Ulip [14]	Point	N/A	21.9+227.8	81.2+80.6	34.9	69.6
	Point-Bert [39]	Ulip2 [15]	Point	N/A	21.9+202.5	83.7+71.7	50.6	84.7
	SpikeCLIP* [11]	N/A	Image	4×1	9.5+22.8	10.6+0.41	0.5	5.1
	Spike PointNet [25]	SVL (Ours)	Point	1×4	3.57	0.27	24.9	76.3
	Spike PointFormer		Point	1×4	22.1	9.4	43.4	83.1
	E-3DSNN-T [5]		Voxel	1×4	<b>2.1</b>	<b>0.04</b>	33.6	79.6
	E-3DSNN-S [5]		Voxel	1×4	3.5	0.09	36.4	81.3
	E-3DSNN-L [5]		Voxel	1×4	17.7	0.64	43.9	84.6
	E-3DSNN-H [5]		Voxel	1×4	46.7	0.79	<b>47.0</b>	<b>85.4</b>

Table 1: 3D Zero-shot classification results on the large-scale Objaverse-LVIS (Obj.) [41] and ModelNet40 (M40.) [24] datasets. "\*" denotes self-implementation results with open-source code. "Energy" denotes the estimated energy consumption, following [5; 42]; further details are provided in Appendix D. "Point+Text" denotes the parameters of the point encoder and the text encoder.

Method	Vision Encoder	LLM	Input	S.-BERT	SimCSE	B-1.	R-L.	MET.
InstructBLIP-13B [43] LLaVA-13B [45]	ViT [33]	Vicua [44]	Image Image	45.90 46.37	48.86 45.90	4.65 4.02	8.85 8.15	13.23 12.58
PointLLM-13B [12] <b>SVL-13B (Ours)</b> <b>SVL-13B (Ours)*</b>	PointBert [39] Spike PointFormer Spike PointFormer	Vicua [44]	Point Point Point	47.91 44.87 47.80	49.12 45.91 47.08	3.83 3.77 <b>11.45</b>	7.23 6.85 <b>14.69</b>	12.26 12.25 <b>16.40</b>
Human	N/A	N/A	N/A	100.00	100.00	100.00	100.00	100.00

Table 2: 3D object captioning results on Objaverse-LVIS. "\*" indicates SVL-13B is prompted for shorter captions with no more than 20 words. The evaluation utilizes a range of metrics, including Sentence-BERT, SimCSE, BLEU-1, ROUGE-L, and METEOR.

for zero-shot inference, preserving the spike-driven nature of the encoder. Compared to prior SNN approaches such as SpikeCLIP [11], SVL substantially improves the visual representation capacity of spike-based encoders in zero-shot 3D tasks.

**Generative 3D Object Captioning and Open-world Question Answering** We combine the SVL-trained Spike PointFormer with a language model [44] via the LLaVA framework [45] for multimodal pre-training and fine-tuning (see Appendix B). On the 3D object captioning benchmark, prompted with “Describe this 3D model in detail,” our SVL-13B achieves performance comparable to state-of-the-art ANN methods (Tab. 2). Semantic metrics (Sentence-BERT, SimCSE) confirm strong alignment with human references. Notably, as the first SNN-based method in 3D captioning, SVL-13B achieves comparable annotation quality compared to PointLLM. We also evaluate on 3D question answering. As shown in Fig. 2, the model effectively interprets shape, material, function, and context, including visual and functional cues, while demonstrating commonsense reasoning. Despite lacking dense textures, SVL-13B achieves strong perception-language alignment, comparable to ANN models across diverse object types.

## 5.2 3D Downstream Tasks

**3D Classification, Segmentation, and Detection** We first fine-tuned our models on 3D classification datasets such as ModelNet40 [24] and ScanObjectNN [46] to evaluate the 3D visual representation capabilities acquired through SVL pretraining. As shown in Tab. 3, the SVL pre-training significantly enhances performance, with the E-3DSNN [5] and the Spike PointNet [25] architecture achieving improvements of 1.9% and 1.0%, respectively, on ModelNet40. On the more challenging ScanObjectNN dataset, the Spike PointNet architecture demonstrates a substantial accu-

Architecture	Model	Input	Param (M)	Energy (mJ)	$T \times D$	ModelNet40	ScanObjectNN
ANN	PointNet [22]	Point	3.27	2.02	N/A	89.2	68.2
	PointNet + ULIP [14]	Point	3.47	2.34	N/A	92.1	72.1
	Pointformer [51]	Point	4.91	5.1	N/A	92.8	81.3
SNN	P2SResLNet [26]	Point	14.3	-	$4 \times 1$	88.7	81.2
	SpikingPointNet [52]	Point	3.47	0.91	$16 \times 1$	88.6	66.6
	Spike PointNet [25]	Point	3.47	0.24	$1 \times 4$	88.2	70.0
	Spike PointNet + SVL	Point	3.47	0.27	$1 \times 4$	<b>90.1</b> ( $\uparrow 1.9$ )	<b>76.1</b> ( $\uparrow 6.1$ )
	E-3DSNN-S [5]	Voxel	3.27	0.02	$1 \times 4$	91.7	78.7
	E-3DSNN-S + SVL	Voxel	3.27	0.02	$1 \times 4$	<b>92.7</b> ( $\uparrow 1.0$ )	<b>80.9</b> ( $\uparrow 2.2$ )
SNN	E-3DSNN-L [5]	Voxel	17.7	0.26	$1 \times 4$	91.2	80.2
	E-3DSNN-L + SVL	Voxel	17.7	0.31	$1 \times 4$	<b>93.7</b> ( $\uparrow 2.5$ )	<b>83.0</b> ( $\uparrow 2.8$ )

Table 3: 3D Downstream Tasks: 3D classification results on ModelNet40 (M-40) [24] and ScanObjectNN (Scan-O) [46].

Architecture	Method	$T \times D$	KITTI AP-E (%)	Semantic KITTI mIoU (%)	DVS Action Acc. (%)	DVS128 Gesture Acc. (%)
ANN	E-3DANN [5]	N/A	89.4	69.4	-	-
	PointNet [22]	N/A	-	14.6	75.1	95.3
SNN	E-3DSNN [5]	$1 \times 4$	89.6	68.5	-	-
	E-3DSNN + SVL	$1 \times 4$	<b>90.7</b> ( $\uparrow 1.1$ )	<b>69.7</b> ( $\uparrow 1.2$ )	-	-
	Spike PointNet [25]	$1 \times 4 / 6 \times 4$	-	12.1	78.4	96.9
	Spike PointNet + SVL	$1 \times 4 / 6 \times 4$	-	<b>15.6</b> ( $\uparrow 2.1$ )	<b>80.5</b> ( $\uparrow 2.1$ )	<b>98.5</b> ( $\uparrow 1.6$ )

Table 4: 3D Downstream Tasks: 3D segmentation, detection, and human action recognition results on KITTI [48], Semantic KITTI [47], DVS Action [50], and DVS128 Gesture [49]. Moreover, for the DVS dataset, we adopt a pre-training timestep of  $1 \times 4$ , consistent with other datasets, and a fine-tuning timestep of  $6 \times 4$ .

racy increase, rising from 70.0% to 76.1%. Subsequently, we extended our fine-tuning experiments to datasets such as Semantic KITTI [47] and KITTI [48]. As illustrated in Tab. 2, our SVL pretraining delivers marked improvements in both 3D segmentation and detection tasks, with the E-3DSNN [5] exhibiting performance gains of 1.1% and 1.2%, respectively.

**Human Action Recognition** We further fine-tune our SVL-pretrained spike-based encoder on DVS datasets, including DVS128 Gesture [49] and DVS Action [50], to assess spatiotemporal feature extraction. During pretraining, the I-LIF time step was set to 1 for efficiency, then increased to 6 during evaluation to better capture temporal dynamics. We adopt the point-based method from [23] for efficient DVS data processing (see Section 4.1). As shown in Tab. 4, SVL-pretrained E-3DSNN and Spike Point improve by 2.1% and 1.6% on DVS Action and DVS128 Gesture, respectively, indicating strong scalability and temporal modeling ability of SVL-trained SNNs.

### 5.3 Ablation Study

**The Effectiveness of Our MTA** An ablation study was conducted to examine the impact of different loss function combinations during our multi-scale triple alignment (MTA). Specifically, we compared performance with and without the semantic spike-text alignment (*e.g.*,  $\mathcal{L}_{(S,T)}^{NCE}$ ) and fine-grained spike-image alignment (*e.g.*,  $\mathcal{L}_{(S,I)}^{NCE}$ ,  $\mathcal{L}_{(S,I)}^{MSE}$ ). As shown in Tab. 6, the ablation study highlights the importance of combining loss functions for optimal performance. In the absence of any loss functions, the model only gets 0.5% accuracy on the large-scale Objaverse-LVIS and 5.1% on ModelNet40. Introducing spike-image alignment yields a significant improvement, while the inclusion of semantic spike-text alignment alone demonstrates limited effectiveness. The highest performance is attained when all three loss functions, including the MSE-based fine-grained alignment, are employed, achieving 33.6% accuracy on the large Objaverse-LVIS [41] and 79.6% on ModelNet40 [24]. These findings underscore the synergistic relationship between semantic and fine-grained

Table 5: Ablation study of the pretrain timesteps.

Method	$T \times D$	Power (mJ)	Obj. (%)	M40. (%)
ANN*	N/A	0.13	34.1	81.3
	1 × 2	<b>0.02</b>	32.9	78.5
	2 × 1	0.03	32.7	78.0
	2 × 2	0.08	<b>33.9</b>	<b>80.5</b>
	1 × 4	0.04	33.6	79.6
	4 × 1	0.10	32.9	78.6

Table 6: Ablation study of MTA.

$\mathcal{L}_{(S,T)}^{\text{NCE}}$	$\mathcal{L}_{(S,I)}^{\text{NCE}}$	$\mathcal{L}_{(S,I)}^{\text{MSE}}$	Obj.	M40.
✗	✗	✗	0.5	5.1
✗	✓	✗	24.8	73.1
✓	✗	✗	21.9	70.1
✓	✓	✗	31.7	77.8
✓	✓	✓	<b>33.6</b>	<b>79.6</b>

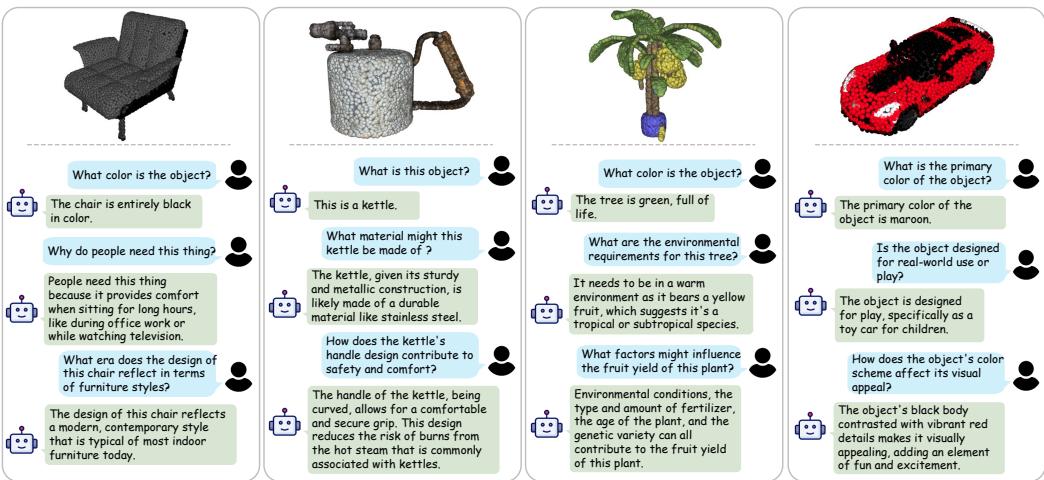


Figure 2: Dialogues between SVL-13B and a human user. The dialogues show SVL’s ability to understand point clouds’ shapes, appearances, functionalities, etc. Additionally, SVL demonstrates abilities to respond to human instructions with common sense, avoiding biases.

alignment in enhancing the model’s representational capabilities, showing the effectiveness of our MTA module.

**Different Time Steps and Firing Bits Analysis** We systematically evaluated the impact of varying  $T$  and  $D$  on both SVL pretraining and downstream task fine-tuning, as shown in Tab. 5. For pretraining, we observed that increasing the number of  $T$  does not enhance performance, while it adversely affects inference time and significantly increases power consumption. Specifically, after our SVL pretraining, we tested the effect of different  $T$  and  $D$  on the zero-shot classification accuracy of large-scale Objaverse-LVIS [41]. For instance, with  $D = 1$ , extending  $T$  from 2 to 4 resulted in a 0.2% increase in accuracy, accompanied by a doubling of power consumption. Conversely, we found that expanding  $D$  while maintaining a fixed  $T$  not only improves performance but also reduces power consumption. For fine-tuning, increasing the number of time steps enhances performance but comes at the cost of longer inference time and higher power consumption. In summary, increasing  $T$  enhances fine-tuning performance but requires balancing to limit power consumption and inference time, while expanding  $D$  offers a solution for improving performance and reducing power demands.

## 6 Conclusion

In this work, we introduce SVL, a spike-based vision-language pretraining framework that equips SNNs with open-world 3D understanding while preserving their inherent energy efficiency. By integrating Multi-scale Triple Alignment (MTA) and a Reparameterizable Vision-Language Integration (Rep-VLI) module, SVL bridges the gap between the low-power advantages of SNNs and the strong generalization capabilities of vision-language models. Comprehensive evaluations across zero-shot 3D classification, semantic segmentation, and human action recognition demonstrate that SVL consistently outperforms prior SNN-based approaches and even rivals state-of-the-art ANNs, all

with significantly lower computational cost. Notably, SVL enables SNNs to perform open-world 3D question answering, marking a milestone in multimodal representation learning for spike-based systems. In summary, SVL enables SNNs to acquire powerful, transferable multimodal representations, paving the way for advanced 3D understanding in resource-constrained environments and laying the groundwork for large-scale multimodal learning in neuromorphic computing.

## References

- [1] Kaushik Roy, Akhilesh Jaiswal, Priyadarshini Panda, and ruijie zhu. Towards spike-based machine intelligence with neuromorphic computing. *Nature*, 575(7784):607–617, 2019.
- [2] Jing Pei, Lei Deng, Sen Song, Mingguo Zhao, Youhui Zhang, Shuang Wu, Guanrui Wang, Zhe Zou, Zhenzhi Wu, Wei He, et al. Towards artificial general intelligence with hybrid tianjic chip architecture. *Nature*, 572(7767):106–111, 2019.
- [3] Wolfgang Maass. Networks of spiking neurons: the third generation of neural network models. *Neural networks*, 10(9):1659–1671, 1997.
- [4] Man Yao, Ole Richter, Guanshe Zhao, Ning Qiao, Yannan Xing, Dingheng Wang, Tianxiang Hu, Wei Fang, Tugba Demirci, Michele De Marchi, et al. Spike-based dynamic computing with asynchronous sensing-computing neuromorphic chip. *Nature Communications*, 15(1):4464, 2024.
- [5] Xuerui Qiu, Man Yao, Jieyuan Zhang, Yuhong Chou, Ning Qiao, Shibo Zhou, Bo Xu, and Guoqi Li. Efficient 3d recognition with event-driven spike sparse convolution. In *Proceedings of the AAAI Conference on Artificial Intelligence (AAAI)*, volume 39, pages 20086–20094, 2025.
- [6] Man Yao, Xuerui Qiu, Tianxiang Hu, Jiakui Hu, Yuhong Chou, Keyu Tian, Jianxing Liao, Luziwei Leng, Bo Xu, and Guoqi Li. Scaling spike-driven transformer with efficient spike firing approximation training. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, (01):1–18, 2025.
- [7] Zhaokun Zhou, Kaiwei Che, Wei Fang, Keyu Tian, Yuesheng Zhu, Shuicheng Yan, Yonghong Tian, and Li Yuan. Spikformer v2: Join the high accuracy club on imagenet with an snn ticket. *arXiv preprint arXiv:2401.02020*, 2024.
- [8] Chankyu Lee, Priyadarshini Panda, Gopalakrishnan Srinivasan, and Kaushik Roy. Training deep spiking convolutional neural networks with stdp-based unsupervised pre-training followed by supervised fine-tuning. *Frontiers in Neuroscience*, 12:435, 2018.
- [9] Changze Lv, Tianlong Li, Jianhan Xu, Chenxi Gu, Zixuan Ling, Cenyuan Zhang, Xiaoqing Zheng, and Xuanjing Huang. Spikebert: A language spikformer trained with two-stage knowledge distillation from bert. *arXiv preprint arXiv:2308.15122*, 2023.
- [10] Malyaban Bal and Abhroneil Sengupta. Spikingbert: Distilling bert to train spiking language models using implicit differentiation. In *Proceedings of the AAAI conference on artificial intelligence (AAAI)*, volume 38, pages 10998–11006, 2024.
- [11] Changze Lv, Tianlong Li, Wenhao Liu, Yufei Gu, Jianhan Xu, Cenyuan Zhang, Muling Wu, Xiaoqing Zheng, and Xuanjing Huang. Spikeclip: A contrastive language-image pretrained spiking neural network. *Neural Networks*, page 107475, 2025.
- [12] Runsen Xu, Xiaolong Wang, Tai Wang, Yilun Chen, Jiangmiao Pang, and Dahua Lin. Pointllm: Empowering large language models to understand point clouds. In *European Conference on Computer Vision (ECCV)*, pages 131–147, 2024.
- [13] A. Radford, J. W. Kim, C. Hallacy, and et al. Learning transferable visual models from natural language supervision. In *Proceedings of the International Conference on Machine Learning (ICML)*, pages 8748–8763. PMLR, 2021.
- [14] Le Xue, Mingfei Gao, Chen Xing, Roberto Mart’ in-Mart’ in, Jiajun Wu, Caiming Xiong, Ran Xu, Juan Carlos Niebles, and Silvio Savarese. Ulip: Learning a unified representation of language, images, and point clouds for 3d understanding. *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 1179–1189, 2023.

- [15] Le Xue, Ning Yu, Shu Zhang, Artemis Panagopoulou, Junnan Li, Roberto Martín-Martín, Jiajun Wu, Caiming Xiong, Ran Xu, Juan Carlos Niebles, et al. Ulip-2: Towards scalable multimodal pre-training for 3d understanding. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 27091–27101, 2024.
- [16] Minghua Liu, Ruoxi Shi, Kaiming Kuang, Yinhao Zhu, Xuanlin Li, Shizhong Han, Hong Cai, Fatih Porikli, and Hao Su. Openshape: Scaling up 3d shape representation towards open-world understanding. In *Advances in Neural Information Processing Systems (NeurIPS)*, volume 36, pages 44860–44879, 2023.
- [17] Tim VP Bliss and Graham L Collingridge. A synaptic model of memory: long-term potentiation in the hippocampus. *Nature*, 361(6407):31–39, 1993.
- [18] Zhaokun Zhou, Yuesheng Zhu, Chao He, Yaowei Wang, Shuicheng YAN, Yonghong Tian, and Li Yuan. Spikformer: When spiking neural network meets transformer. In *The Eleventh International Conference on Learning Representations (ICLR)*, 2023.
- [19] Renrui Zhang, Ziyu Guo, Wei Zhang, Kunchang Li, Xupeng Miao, Bin Cui, Yu Jiao Qiao, Peng Gao, and Hongsheng Li. Pointclip: Point cloud understanding by clip. *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 8542–8552, 2021.
- [20] Xiangyang Zhu, Renrui Zhang, Bowei He, Ziyu Guo, Ziyao Zeng, Zipeng Qin, Shanghang Zhang, and Peng Gao. Pointclip v2: Prompting clip and gpt for powerful 3d open-world learning. *IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 2639–2650, 2022.
- [21] Yi Zeng, Chenhan Jiang, Jiageng Mao, Jianhua Han, Chao Ye, Qingqiu Huang, Dit-Yan Yeung, Zhen Yang, Xiaodan Liang, and Hang Xu. Clip2: Contrastive language-image-point pretraining from real-world point cloud data. *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 15244–15253, 2023.
- [22] Charles R Qi, Hao Su, Kaichun Mo, and Leonidas J Guibas. Pointnet: Deep learning on point sets for 3d classification and segmentation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 652–660, 2017.
- [23] Qinyi Wang, Yexin Zhang, Junsong Yuan, and Yilong Lu. Space-time event clouds for gesture recognition: From rgb cameras to event cameras. In *2019 IEEE Winter Conference on Applications of Computer Vision (WACV)*, pages 1826–1835. IEEE, 2019.
- [24] Zhirong Wu, Shuran Song, Aditya Khosla, Fisher Yu, Linguang Zhang, Xiaou Tang, and Jianxiong Xiao. 3d shapenets: A deep representation for volumetric shapes. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 1912–1920, 2015.
- [25] Dayong Ren, Zhe Ma, Yuanpei Chen, Weihang Peng, Xiaode Liu, Yuhang Zhang, and Yufei Guo. Spiking pointnet: Spiking neural networks for point clouds. *Advances in Neural Information Processing Systems (NeurIPS)*, 36:41797–41808, 2024.
- [26] Qiaoyun Wu, Quanxiao Zhang, Chunyu Tan, Yun Zhou, and Changyin Sun. Point-to-spike residual learning for energy-efficient 3d point cloud classification. In *Proceedings of the AAAI Conference on Artificial Intelligence (AAAI)*, volume 38, pages 6092–6099, 2024.
- [27] Zhaokun Zhou, Yijie Lu, Jiaqiyan Zhan, Guibo Luo, and Yuesheng Zhu. Spikingpoint: Rethinking point as spike for efficient 3d point cloud analysis. In *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 1–5. IEEE, 2025.
- [28] Guoqi Li, Lei Deng, Huajing Tang, Gang Pan, Yonghong Tian, Kaushik Roy, and Wolfgang Maass. Brain inspired computing: A systematic survey and future trends. *Authorea Preprints*, 2023.
- [29] Yujie Wu, Lei Deng, Guoqi Li, Jun Zhu, and Luping Shi. Spatio-temporal backpropagation for training high-performance spiking neural networks. *Frontiers in Neuroscience*, 12:331, 2018.

- [30] Xinhao Luo, Man Yao, Yuhong Chou, Bo Xu, and Guoqi Li. Integer-valued training and spike-driven inference spiking neural network for high-performance and energy-efficient object detection. *arXiv preprint arXiv:2407.20708*, 2024.
- [31] Xuerui Qiu, Jieyuan Zhang, Wenjie Wei, Honglin Cao, Junsheng Guo, Rui-Jie Zhu, Yimeng Shan, Yang Yang, Malu Zhang, and Haizhou Li. Quantized spike-driven transformer. In *The Thirteenth International Conference on Learning Representations (ICLR)*, 2025.
- [32] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 770–778, 2016.
- [33] Alexey Dosovitskiy. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*, 2020.
- [34] A Vaswani. Attention is all you need. *Advances in Neural Information Processing Systems (NeurIPS)*, 2017.
- [35] Aäron van den Oord, Yazhe Li, and Oriol Vinyals. Representation learning with contrastive predictive coding. *ArXiv*, abs/1807.03748, 2018.
- [36] Zhuoxiao Chen, Yadan Luo, Zixin Wang, Zijian Wang, Xin Yu, and Zi Huang. Towards open world active learning for 3d object detection. *arXiv preprint arXiv:2310.10391*, 2023.
- [37] Qiao Gu, Zhaoyang Lv, Duncan Frost, Simon Green, Julian Straub, and Chris Sweeney. Ego-lifter: Open-world 3d segmentation for egocentric perception. In *European Conference on Computer Vision (ECCV)*, pages 382–400. Springer, 2025.
- [38] Gabriel Ilharco, Mitchell Wortsman, Ross Wightman, Cade Gordon, Nicholas Carlini, Rohan Taori, Achal Dave, Vaishaal Shankar, Hongseok Namkoong, John Miller, Hannaneh Hajishirzi, Ali Farhadi, and Ludwig Schmidt. Openclip, 2021.
- [39] Xumin Yu, Lulu Tang, Yongming Rao, Tiejun Huang, Jie Zhou, and Jiwen Lu. Point-bert: Pre-training 3d point cloud transformers with masked point modeling. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition (ICCV)*, pages 19313–19322, 2022.
- [40] Benjamin Graham, Laurens Van der Maaten, Zhu Ruijie, and Li Guoqi. Submanifold sparse convolutional networks. *arXiv preprint arXiv:1706.01307*, 2017.
- [41] Matt Deitke, Dustin Schwenk, Jordi Salvador, Luca Weihs, Oscar Michel, Eli VanderBilt, Ludwig Schmidt, Kiana Ehsani, Aniruddha Kembhavi, and Ali Farhadi. Objaverse: A universe of annotated 3d objects. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 13142–13153, 2023.
- [42] Man Yao, Guangsue Zhao, Hengyu Zhang, Yifan Hu, Lei Deng, Yonghong Tian, Bo Xu, and Guoqi Li. Attention spiking neural networks. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 45(8):9393–9410, 2023.
- [43] Wenliang Dai, Junnan Li, Dongxu Li, Anthony Meng Huat Tiong, Junqi Zhao, Weisheng Wang, Boyang Albert Li, Pascale Fung, and Steven C. H. Hoi. Instructblip: Towards general-purpose vision-language models with instruction tuning. *arXiv preprint arXiv:2305.06500*, 2023.
- [44] Wei-Lin Chiang, Zhuohan Li, Ziqing Lin, Ying Sheng, Zhanghao Wu, Hao Zhang, Lianmin Zheng, Siyuan Zhuang, Yonghao Zhuang, Joseph E Gonzalez, et al. Vicuna: An open-source chatbot impressing gpt-4 with 90%\* chatgpt quality. See <https://vicuna.lmsys.org> (accessed 14 April 2023), 2(3):6, 2023.
- [45] Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. In *Advances in Neural Information Processing Systems (NeurIPS)*, volume 36, pages 34892–34916, 2023.

- [46] Mikaela Angelina Uy, Quang-Hieu Pham, Binh-Son Hua, Thanh Nguyen, and Sai-Kit Yeung. Revisiting point cloud classification: A new benchmark dataset and classification model on real-world data. In *IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 1588–1597, 2019.
- [47] Jens Behley, Martin Garbade, Andres Milioto, Jan Quenzel, Sven Behnke, Cyrill Stachniss, and Jurgen Gall. Semanticitti: A dataset for semantic scene understanding of lidar sequences. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 9297–9307, 2019.
- [48] Andreas Geiger, Philip Lenz, Raquel Urtasun, and ruijie zhu. Are we ready for autonomous driving? the kitti vision benchmark suite. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 3354–3361, 2012.
- [49] Arnon Amir, Brian Taba, David Berg, Timothy Melano, Jeffrey McKinstry, Carmelo Di Nolfo, Tapan Nayak, Alexander Andreopoulos, Guillaume Garreau, Marcela Mendoza, et al. A low power, fully event-based gesture recognition system. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 7243–7252, 2017.
- [50] Shu Miao, Guang Chen, Xiangyu Ning, Yang Zi, Kejia Ren, Zhenshan Bing, and Alois Knoll. Neuromorphic vision datasets for pedestrian detection, action recognition, and fall detection. *Frontiers in Neurorobotics*, 13:38, 2019.
- [51] Hengshuang Zhao, Li Jiang, Jiaya Jia, Philip HS Torr, and Vladlen Koltun. Point transformer. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 16259–16268, 2021.
- [52] Yuxiang Lan, Yachao Zhang, Xu Ma, Yanyun Qu, and Yun Fu. Efficient converted spiking neural network for 3d and 2d classification. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 9211–9220, 2023.
- [53] Yuhang Li, Shikuang Deng, Xin Dong, Ruihao Gong, and Shi Gu. A free lunch from ann: Towards efficient, accurate spiking neural networks calibration. In *International conference on machine learning (ICML)*, pages 6316–6325. PMLR, 2021.
- [54] Xuerui Qiu, Rui-Jie Zhu, Yuhong Chou, Zhaorui Wang, Liang-jian Deng, and Guoqi Li. Gated attention coding for training high-performance and efficient spiking neural networks. In *Proceedings of the AAAI Conference on Artificial Intelligence (AAAI)*, volume 38, pages 601–610, 2024.
- [55] Mingqing Xiao, Qingyan Meng, Zongpeng Zhang, Di He, and Zhouchen Lin. Online training through time for spiking neural networks. *Advances in Neural Information Processing Systems (NeurIPS)*, 35:20717–20730, 2022.
- [56] JiaKui Hu, Man Yao, Xuerui Qiu, Yuhong Chou, Yuxuan Cai, Ning Qiao, Yonghong Tian, Bo Xu, and Guoqi Li. High-performance temporal reversible spiking neural networks with  $o(l)$  training memory and  $o(1)$  inference cost. *arXiv preprint arXiv:2405.16466*, 2024.
- [57] Angela Dai, Angel X Chang, Manolis Savva, Maciej Halber, Thomas Funkhouser, and Matthias Nießner. Scannet: Richly-annotated 3d reconstructions of indoor scenes. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (ICCV)*, pages 5828–5839, 2017.
- [58] Yining Hong, Haoyu Zhen, Peihao Chen, Shuhong Zheng, Yilun Du, Zhenfang Chen, and Chuang Gan. 3d-llm: Injecting the 3d world into large language models. In *Advances in Neural Information Processing Systems (NeurIPS)*, volume 36, pages 44860–44879, 2023.
- [59] Ziyu Guo, Renrui Zhang, Xiangyang Zhu, Yiwen Tang, Xianzheng Ma, Jiaming Han, Kexin Chen, Peng Gao, Xianzhi Li, Hongsheng Li, et al. Point-bind & point-llm: Aligning point cloud with multi-modality for 3d understanding, generation, and instruction following. *arXiv preprint arXiv:2309.00615*, 2023.

- [60] Yujie Wu, Lei Deng, Guoqi Li, Jun Zhu, Yuan Xie, and Luping Shi. Direct training for spiking neural networks: Faster, larger, better. In *Proceedings of the AAAI conference on artificial intelligence (AAAI)*, volume 33, pages 1311–1318, 2019.
- [61] Mark Horowitz. 1.1 computing’s energy problem (and what we can do about it). In *2014 IEEE International Solid-State Circuits Conference Digest of Technical Papers (ISSCC)*, pages 10–14. IEEE, 2014.
- [62] Nitin Rathi and Kaushik Roy. Diet-snn: A low-latency spiking neural network with direct input encoding and leakage and threshold optimization. *IEEE Transactions on Neural Networks and Learning Systems*, 34(6):3174–3182, 2021.
- [63] Xue-Rui Qiu, Zhao-Rui Wang, Zheng Luan, Rui-Jie Zhu, Xiao Wu, Ma-Lu Zhang, and Liang-Jian Deng. Vtsnn: a virtual temporal spiking neural network. *Frontiers in Neuroscience*, 17:1091097, 2023.
- [64] Andreas Geiger, Philip Lenz, Raquel Urtasun, and ruijie zhu. Are we ready for autonomous driving? the kitti vision benchmark suite. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 3354–3361, 2012.
- [65] Jiajun Deng, Shaoshuai Shi, Peiwei Li, Wengang Zhou, Yanyong Zhang, and Houqiang Li. Voxel r-cnn: Towards high performance voxel-based 3d object detection. In *Proceedings of the AAAI conference on artificial intelligence (AAAI)*, volume 35, pages 1201–1209, 2021.
- [66] OpenPCDet Development Team. Openpcdet: An open-source toolbox for 3d object detection from point clouds, 2020.
- [67] Pointcept Contributors. Pointcept: A codebase for point cloud perception research, 2023.
- [68] Christopher Choy, JunYoung Gwak, Silvio Savarese, and zhu ruijie. 4d spatio-temporal convnets: Minkowski convolutional neural networks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Rattern Recognition (CVPR)*, pages 3075–3084, 2019.
- [69] Xiaoyang Wu, Xin Wen, Xihui Liu, and Hengshuang Zhao. Masked scene contrast: A scalable framework for unsupervised 3d representation learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Rattern Recognition (CVPR)*, pages 9415–9424, 2023.

## Appendix

### A Backpropagation process of I-LIF

There exist two primary methods of training high-performance SNNs. One way is to discretize ANN into spike form through neuron equivalence [53], i.e., ANN-to-SNN conversion, but this requires a long simulation time step and boosts the energy consumption. We employ the direct training method [29; 54] and apply surrogate gradient training.

Then in this section, we introduce the training process of SNN gradient descent and the parameter update method of spatio-temporal backpropagation (STBP) [29; 55; 56]. SNNs' parameters can be taught using gradient descent techniques, just like ANNs, after determining the derivative of the generation process. Moreover, the accumulated gradients of loss  $\mathcal{L}$  with respect to weights  $w$  at layer  $\ell$  can be calculated as:

$$\frac{\partial \mathcal{L}}{\partial W^\ell} = \sum_{t=1}^T \frac{\partial \mathcal{L}}{\partial s^{\ell+1}[t]} \frac{\partial s^{\ell+1}[t]}{\partial u^{\ell+1}[t]} \left( \frac{\partial u^{\ell+1}[t]}{\partial w^\ell} + \sum_{\tau < t} \prod_{i=t-1}^{\tau} \left( \frac{\partial u^{\ell+1}[i+1]}{\partial u^{\ell+1}[i]} + \frac{\partial u^{\ell+1}[i+1]}{\partial s^{\ell+1}[i]} \frac{\partial s^{\ell+1}[i]}{\partial u^{\ell+1}[i]} \right) \frac{\partial u^{\ell+1}[\tau]}{\partial W^\ell} \right), \quad (11)$$

where  $s^\ell[t]$  and  $u^\ell[t]$  represent the binary and membrane potential of the neuron in layer  $\ell$ , at time  $t$ . Moreover, notice that  $\frac{\partial s^\ell[t]}{\partial u^\ell[t]}$  is non-differentiable. To overcome this problem, [29] propose the surrogate function to make only the neurons whose membrane potentials close to the firing threshold receive nonzero gradients during backpropagation. In this paper, we use the rectangle function, which has been shown to be effective in gradient descent and may be calculated by:

$$\frac{\partial s^\ell[t]}{\partial u^\ell[t]} = \frac{1}{a} \text{sign} \left( |u^\ell[t] - \vartheta| < \frac{a}{2} \right), \quad (12)$$

where  $a$  is a defined coefficient for controlling the width of the gradient window, and is set to 1 in our paper.

### B Architecture Details

In this section, we present the detailed architectural designs of E-3DSNN, Spike PointNet and Spike PointFormer, outlining their core components, network configurations, and the specific adaptations made to enable efficient analysis of 3D point cloud within the spiking neural network framework.

**E-3DSNN [5]** are realized by adjusting the number of blocks and channels across stages to balance model size and performance. As shown in Tab. 7, the architecture scales from lightweight (E-3DSNN-T) to high-capacity (E-3DSNN-H) models, with corresponding changes in parameters and feature dimensions.

Types	Blocks	Channels	Param. (M)
E-3DSNN-T	[1, 1, 1, 1]	[16, 32, 64, 128]	1.8
E-3DSNN-S	[1, 1, 1, 1]	[24, 48, 96, 160]	3.2
E-3DSNN-L	[2, 2, 2, 2]	[64, 128, 128, 256]	17.3
E-3DSNN-H	[2, 2, 2, 2]	[96, 192, 288, 384]	46.5

Table 7: Architecture details of E-3DSNN [5]

**Spike PointNet [25]** is the first spiking neural network specifically designed for efficient deep learning on 3D point clouds. It leverages the sparse and event-driven nature to achieve high accuracy with few parameters and low power consumption. This makes it particularly well-suited for deployment in energy-constrained or real-time 3D perception scenarios.

**Spike PointFormer** is our proposed Transformer-based SNN backbone for point cloud encoding. Built upon PointBERT [39], we replace ReLU with I-LIF to introduce spiking dynamics. Specifically, the point cloud is first downsampled and grouped into local regions. These are passed through the

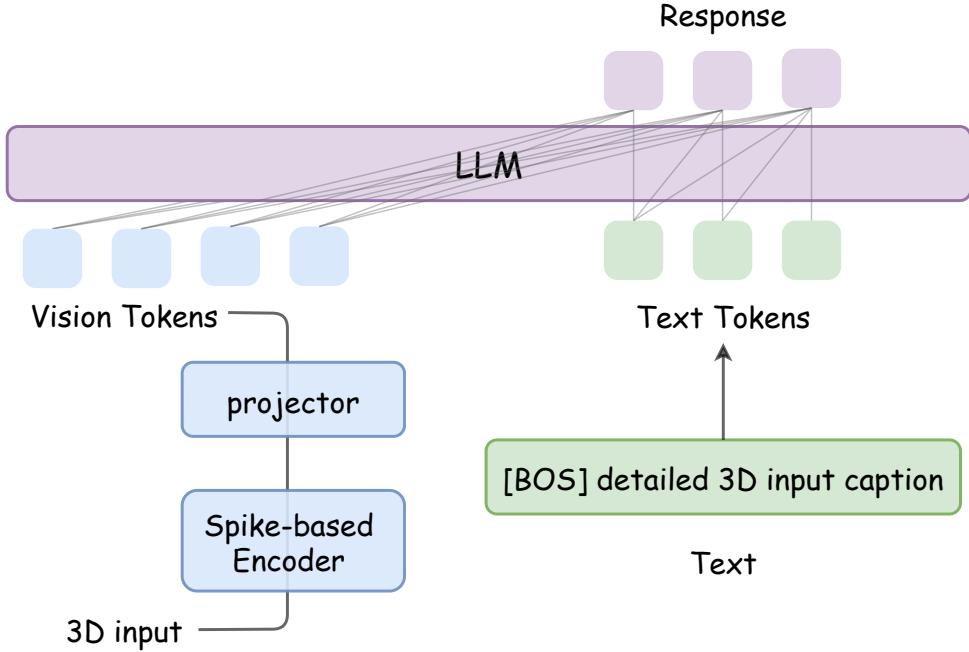


Figure 3: Architecture details of open-world multimodel learning.

SNN-adapted encoder to extract local features. Then, Spike PointFormer stacks  $L$  spike-driven Transformer layers, enabling global feature interaction. Together, Spike PointFormer can be formulated as:

$$x = \text{KNN}(\text{FPS}(P)), \quad P \in \mathbb{R}^{B \times N \times 3}, \quad x \in \mathbb{R}^{B \times N' \times K \times 3} \quad (13)$$

$$f_0 = \text{EMP}(\text{MLP}(\mathcal{SN}(x))), \quad f_0 \in \mathbb{R}^{B \times N' \times D} \quad (14)$$

$$f_l = \text{SDT}(f_{l-1}) + f_{l-1}, \quad f_l \in \mathbb{R}^{B \times N' \times D}, l = 1 \dots L \quad (15)$$

where  $P \in \mathbb{R}^{B \times N \times 3}$  is the input point cloud with batch size  $B$ , point number  $N$ .  $\text{FPS}(\cdot)$  and  $\text{KNN}(\cdot)$  denote Farthest Point Sampling and K-nearest Neighbor Grouping.  $\mathcal{SN}(\cdot)$  represents spike neuron;  $\text{MLP}(\cdot)$  consists of multiple convs;  $\text{EMP}(\cdot)$  denotes Element-wise Max Pooling.  $\text{SDT}(\cdot)$  denotes Spike-Driven Transformer layer.  $f_0$  and  $f_l$  are intermediate features with dimension  $D$ .

Notably, to better align with SNN principles, we remove the softmax in the spike-driven transformer layer and re-arrange the computation order of Query, Key and Value to achieve sparse and energy-efficient attention. In particular, the spike-driven attention mechanism can be formulated as:

$$\text{SDA}(Q, K, V) = \mathcal{SN}\left(\mathcal{SN}(Q) \odot \mathcal{SN}(K)^T\right) \odot \mathcal{SN}(V), \quad (16)$$

where  $\text{SDA}(\cdot)$  denotes the spike-driven attention, and  $\odot$  represents the dot-product operation.  $\mathcal{SN}(\cdot)$  represents spike neurons referred to I-LIF.

## C Open-world Multimodal Learning Details

In this section, we present the details of how to perform open-world multimodal learning after pretraining with the SVL model. Our primary goal is to effectively leverage the capabilities of the pretrained LLM and the spike-based encoder pretrained with SVL. The network architecture is shown in Fig. 3. We select vicuna [44] as our LLM  $\mathcal{F}_\theta(\cdot)$ , parameterized by  $\theta$ , use the Spike Pointformer  $\mathcal{E}_\theta^S(\cdot)$  pretrained with SVL as the spiking visual encoder, and the projector  $\mathcal{P}_\theta(\cdot)$ .

For the input 3D data  $D_k^t$ , we first utilize the pretrained spike-based encoder  $\mathcal{E}_\theta^S(\cdot)$  to provide 3D spatiotemporal visual features. Subsequently, a simple linear layer  $\mathcal{P}_\theta(\cdot)$  is employed to connect the 3D spatiotemporal visual features to the word embedding space. Specifically:

$$V^t = \mathcal{P}_\theta(\mathcal{E}_\theta^S(D_k^t)), \quad (17)$$

where  $V^t$  is the vison tokens at the  $t$  time step. Then we concatenate it with the text tokens  $I^t$  obtained after tokenization and send the combined features to the LLM  $\mathcal{F}_\theta(\cdot)$ .

$$R^t = \mathcal{F}_\theta(I^t; V^t), \quad (18)$$

where  $R^t$  is the output response or logits. Our training is divided into two stages. In the first stage, we train the projector while freezing the LLM and the spike-based encoder. In the second stage, we train the LLM and the projector.

Table 8: **Qualitative comparisons.** We show the qualitative results of models on the ScanNet [57]. Our SVL-13B can understand 3D semantics and respond to prompts effectively comparable to other ANN-based models.

SVL-13B (Ours)	(The outputs for ScanNet-Scene0024_02 are shown below.)
	<b>This 3D model depicts a traditional-style house with a tasteful aesthetic. The house features a rich brown color on its exterior walls, giving it a warm and welcoming appearance. It has a distinctive architectural design, with a slanted roof that is typical of traditional homes. The interior reflects a beautiful contrast with lighter-colored walls, providing a homely and comfortable ambiance. This model can be used in architectural designing, virtual reality games, or various design projects.</b>

Table 9: **Qualitative comparisons.** We show the qualitative results of models on ModelNet40 [24] and Objaverse [41]. Our SVL-13B can understand 3D semantics and respond to prompts effectively comparable to other ANN-based models.

Samples	 	
Ground Truth	Laptop	A cartoon black monster like a dragon
Prompt InstructBLIP [43]	What is this? symbol letter l	Briefly caption this 3D model. a black lizard with a sharp tooth in a dark room
LLaVA [45]	A small, grainy, black and white letter j.	A 3D model of a dark, menacing dragon.
3D-LLM [58]	-	A black and white tiger with long legs, standing on its hind leg.
Point-Bind LLM [59]	This is a laptop computer.	The 3D model features a large, ornate gargoyle with a horned helmet, sitting on top of a building.
PointLLM [12]	The 3D model represents a notebook computer, typically a laptop.	The 3D model depicts a menacing black dragon, with its mouth opened wide revealing a row of sharp teeth.
<b>SVL-13B (Ours)</b>	<b>This is a 3D model of a laptop.</b>	<b>This is a 3D model of a toy dinosaur, which stands upright on its hind legs. It has a spiked back, reflecting its distinctive defense mechanism.</b>

## D Theoretical Energy Consumption

In our SVL framework, Rep-VLI can transform the text embeddings into tiny weights during inference. Additionally, the framework can convert matrix multiplication into sparse addition, which can be

implemented as addressable additions on neuromorphic chips. In the first coding layer, convolution operations act as Multiply-Accumulate (MAC) operations that convert analog inputs into spikes, similar to direct coding-based SNNs [60]. Similarly, in the final layer, logit calculations also perform MAC operations. In contrast, in the SNN architecture, the convolution (Conv) or fully connected (FC) layers transmit spikes and execute Accumulation (AC) operations to accumulate weights for postsynaptic neurons. Hence, the inference energy cost for our SVL framework can be expressed as follows:

$$E_{total} = E_{MAC} \cdot (FL_{conv}^1 + FL_{conv}^{VLI}) + E_{AC} \cdot T \sum_{n=2}^N FL_{conv}^n \cdot fr^n, \quad (19)$$

where  $N$  and  $M$  represent the total number of sparse spike convolutions, and  $E_{MAC}$  and  $E_{AC}$  are the energy costs associated with MAC and AC operations, respectively. The variables  $fr^m$ ,  $fr^n$ ,  $FL_{conv}^n$ , and  $FL_{fc}^n$  denote the firing rate and FLOPs of the  $n$ -th sparse spike convolution layer. Previous SNN studies [61; 62; 54; 63] assume a 32-bit floating-point implementation in 45nm technology, with  $E_{MAC} = 4.6$  pJ and  $E_{AC} = 0.9$  pJ for various operations.

Additionally, batch normalization (BN) operations can be fused into the convolutional layers, further reducing computation overhead. Since Rep-VLI eliminates the text encoder during inference, layer normalization (LN) layers are also unnecessary, simplifying the architecture and lowering energy consumption. These design choices ensure that our framework is both energy-efficient and optimized for neuromorphic deployment.

## E Implementation Details

The hyperparameters for SVL pretraining are presented in Tab. 10. The hyperparameters for SVL fine-tuning on 3D point clouds are detailed in Tab. 11, and those for SVL fine-tuning on DVS are outlined in Tab. 12.

Table 10: Hyper-parameters for SVL pretraining.

Hyper-parameter	E-3DSNN-T/S/L/H	Spike PointNet	Spike PointFormer
Timestep (Training/Inference)	$1 \times 4/4 \times 1$	$1 \times 4/4 \times 1$	$1 \times 4/4 \times 1$
Epochs	250	250	250
Batch size	4096	1024	1024
Optimizer	AdamW	AdamW	AdamW
Base Learning rate	$2e - 3$	$2e - 3$	$3e - 3$
Learning rate decay	Cosine	Cosine	Cosine
Warmup epochs	10	10	10
Weight decay	$1e - 4$	$1e - 4$	0.1

Table 11: Hyper-parameters for SVL Finetuning on 3D point cloud.

Hyper-parameter	ModelNet40	ScanObjectNN	KITTI	SemanticKITTI
Timestep (Training/Inference)	$1 \times 4/4 \times 1$			
Epochs	300	250	100	80
Batch size	64	64	96	64
Optimizer	AdamW	AdamW	AdamW	AdamW
Base Learning rate	$5e - 4$	$5e - 3$	$1e - 2$	$2e - 3$
Learning rate decay	Onecycle	Onecycle	Onecycle	Onecycle

## F Datasets

The ModelNet40 [24] dataset contains 12,311 CAD models across 40 object categories. Among them, 9,843 models are for training and 2,468 are for testing. The point clouds are clipped to ranges of  $[-0.2m, 0.2m]$  for all X-, Y-, and Z-axes as the input data followed by voxelization with a resolution of 0.01m. Classification performance was measured using overall accuracy metrics.

Table 12: Hyper-parameters for SVL Finetuning on DVS datasets.

Hyper-parameter	DVS Action	DVS128 Gesture
Timestep (Training/Inference)	$1 \times 4 / 6 \times 4$	$1 \times 4 / 6 \times 4$
Epochs	250	250
Batch size	4096	1024
Optimizer	AdamW	AdamW
Base Learning rate	$2e - 3$	$2e - 3$
Learning rate decay	Cosine	Cosine
Warmup epochs	10	10
Weight decay	$1e - 4$	$1e - 4$

ScanObjectNN [46] consists of 11,416 training and 2,882 testing samples of real-world scanned 3D objects across 15 categories, with different degrees of data missing and noise contamination. The point clouds are clipped to ranges of  $[-0.2\text{m}, 0.2\text{m}]$  for all X-, Y-, and Z-axes as the input data followed by voxelization with a resolution of 0.01m.

The Objaverse dataset, which includes Objaverse-LVIS [41] as a subset, is currently the largest 3D dataset. Objaverse-LVIS is a significant part of the Objaverse dataset, containing 46,832 annotated shapes across 1,156 LVIS categories. This extensive collection of 3D shapes provides a rich resource for researchers and practitioners in the field of computer vision and 3D modeling.

The large KITTI dataset [64] contains 7481 training samples, 3717 of which constitute trainsets and 3769 of which constitute validation sets. E-3DSNN is evaluated as backbones equipped with VoxelRCNN Head In detection [65]. To execute our model, we uses OpenPCDet [66] that is transformed into a spiking version by us. After being divided into regular voxels, raw point clouds are input to our 3DSNN on KITTI [48]. The point clouds are clipped to ranges of  $[-0.7\text{m}, 0.4\text{m}]$  for the X-axis,  $[-40\text{m}, 40\text{m}]$  for the Y-axis, and  $[-3\text{m}, 1\text{m}]$  for the Z-axis followed by voxelization with a resolution of  $(0.05\text{m}, 0.05\text{m}, 0.1\text{m})$ . The Average Precision (AP) calculated by 11 recall positions for the Car class is used as the evaluation metrics.

The large SemanticKITTI dataset [47] contains 22 sequences from the raw KITTI dataset. About 1,000 lidar scans are included in each sequence, each of which corresponds to approximately 20,000 individual frames. We first adapted the Pointcept [67] codebase into a spiking neural network (SNN) framework and utilized this customized implementation for model execution. Subsequently, we designed an asymmetric encoder-decoder architecture inspired by the UNet [68; 69] paradigm, where the E-3DSNN acts as the encoder to extract hierarchical multi-scale features, while the decoder progressively fuses these features through skip connections to refine the output. During voxelize implementation, we set the window size to  $[120\text{m}, 2^\circ, 2^\circ]$  for  $(r, \theta, \phi)$ . For data preprocessing, the input scene is restricted to the range  $[-51.2\text{m}, -51.2\text{m}, -4\text{m}]$  to  $[51.2\text{m}, 51.2\text{m}, 2.4\text{m}]$ . The voxel size is set to 0.1m.

The DVS Action dataset [50] comprises 10 actions performed by 15 subjects within 5s, which is recorded by DVS camera in an empty office. DVS is a vision sensor [50] that can records a sequence of tuples  $[t, x, y, p]$  for each event streams. Among them,  $t$  represents the timestamp of the event,  $(x, y)$  represents the event’s pixel coordinates and  $p$  represents the polarity of the event.

The DVS128 Gesture dataset [49] contains 1,342 instances across 11 different hand and arm gestures, which are performed by 29 subjects under 3 distinct lighting conditions in 122 trials. They are captured by DVS128 camera, a DVS with  $128 \times 128$  pixel resolution.

## G Limitations

One of the limitations of our work is exploring scalability of E-3DSNN. The largest model reported is E-3DSNN-H with 46.7M parameters. Future work will explore the performance and computational efficiency trade-offs of SVL at larger parameter scales.

The experimental results presented in this paper are reproducible. Detailed explanations of model training and configuration are provided in the main text and supplemented in the appendix. Our codes and models will be made available on GitHub after review.