



## SketchANIMAR: Sketch-Based 3D Animal Fine-Grained Retrieval

Trung-Nghia Le<sup>a,b</sup>, Tam V. Nguyen<sup>b,c</sup>, Minh-Quan Le<sup>b,a,b</sup>, Trong-Thuan Nguyen<sup>b,a,b</sup>, Viet-Tham Huynh<sup>b,a,b</sup>, Trong-Le Do<sup>b,a,b</sup>, Khanh-Duy Le<sup>b,a,b</sup>, Mai-Khiem Tran<sup>b,a,b</sup>, Nhat Hoang-Xuan<sup>b,a,b</sup>, Thang-Long Nguyen-Ho<sup>b,a,b</sup>, Vinh-Tiep Nguyen<sup>b,d,b</sup>, Nhat-Quynh Le-Pham<sup>a,b</sup>, Huu-Phuc Pham<sup>a,b</sup>, Trong-Vu Hoang<sup>a,b</sup>, Quang-Binh Nguyen<sup>a,b</sup>, Trong-Hieu Nguyen-Mau<sup>b,a,b</sup>, Tuan-Luc Huynh<sup>b,a,b</sup>, Thanh-Danh Le<sup>a,b</sup>, Ngoc-Linh Nguyen-Ha<sup>a,b</sup>, Tuong-Vy Truong-Thuy<sup>a,b</sup>, Truong Hoai Phong<sup>a,b</sup>, Tuong-Nghiem Diep<sup>a,b</sup>, Khanh-Duy Ho<sup>a,b</sup>, Xuan-Hieu Nguyen<sup>a,b</sup>, Thien-Phuc Tran<sup>a,b</sup>, Tuan-Anh Yang<sup>a,b</sup>, Kim-Phat Tran<sup>a,b</sup>, Nhu-Vinh Hoang<sup>a,b</sup>, Minh-Quang Nguyen<sup>a,b</sup>, Hoai-Danh Vo<sup>a,b</sup>, Minh-Hoa Doan<sup>a,b</sup>, Hai-Dang Nguyen<sup>b,a,b</sup>, Akihiro Sugimoto<sup>b,e</sup>, Minh-Triet Tran<sup>b,a,b,\*</sup>

<sup>a</sup>University of Science, VNU-HCM, Ho Chi Minh City, Vietnam

<sup>b</sup>Vietnam National University, Ho Chi Minh City, Vietnam

<sup>c</sup>University of Dayton, Ohio, U.S.

<sup>d</sup>University of Information Technology, VNU-HCM, Ho Chi Minh City, Vietnam

<sup>e</sup>National Institute of Informatics, Tokyo, Japan

### ARTICLE INFO

#### Article history:

Received August 10, 2023

3D object retrieval, fine-grained retrieval, and animal models.

### ABSTRACT

The retrieval of 3D objects has gained significant importance in recent years due to its broad range of applications in computer vision, computer graphics, virtual reality, and augmented reality. However, the retrieval of 3D objects presents significant challenges due to the intricate nature of 3D models, which can vary in shape, size, and texture, and have numerous polygons and vertices. To this end, we introduce a novel SHREC challenge track that focuses on retrieving relevant 3D animal models from a dataset using sketch queries and expedites accessing 3D models through available sketches. Furthermore, a new dataset named ANIMAR was constructed in this study, comprising a collection of 711 unique 3D animal models and 140 corresponding sketch queries. Our contest requires participants to retrieve 3D models based on complex and detailed sketches. We receive satisfactory results from eight teams and 204 runs. Although further improvement is necessary, the proposed task has the potential to incentivize additional research in the domain of 3D object retrieval, potentially yielding benefits for a wide range of applications. We also provide insights into potential areas of future research, such as improving techniques for feature extraction and matching and creating more diverse datasets to evaluate retrieval performance.

© 2023 Elsevier B.V. All rights reserved.

### 1. Introduction

The rapid development of 3D technologies has produced a remarkable number of 3D objects. Consequently, 3D object retrieval has garnered considerable attention and is beneficial in real-life applications [1, 2, 3], including but not limited to video games, artistic pursuits, cinematography, and virtual reality.

Sketch-based 3D object retrieval aims to retrieve 3D models from a user's hand-drawn 2D sketch. Due to the innate intuitive appeal of freehand drawings, sketch-based 3D object retrieval has drawn a significant amount of attention and is being utilized in numerous critical applications such as 3D scene reconstruction [4, 5, 6], 3D geometry video retrieval [7, 8, 9], and 3D augmented/virtual reality entertainment [10, 11]. However, sketch-based 3D object retrieval poses a formidable challenge in 3D object retrieval research, primarily due to the large discrepancy between the 2D and 3D modalities: non-realistic 2D

\*Corresponding author

e-mail: [tmtriet@fit.hcmus.edu.vn](mailto:tmtriet@fit.hcmus.edu.vn) (Minh-Triet Tran)

sketches differ significantly from their 3D counterparts and respective views.

Several SHREC challenge tracks [12, 13, 14, 15, 16, 17] have been organized to facilitate research on sketch-based 3D object retrieval. However, the existing datasets incorporated in these tracks primarily comprise generic objects with simplistic shapes and poses. To augment sketch-based 3D object retrieval research, we organize a new SHREC challenge track dedicated to *Sketch-based 3D ANIMAL model fine-grained Retrieval (SketchANIMAR)*<sup>1</sup>. This track aims to retrieve relevant 3D animal models from a dataset using sketch queries and expedites accessing 3D models through available sketches. Previous SHREC challenge tracks have focused on a limited number of general object categories, often lacking realism. Our challenge track for SHREC 2023 is significantly more challenging and can simulate real-life scenarios more effectively than its predecessors. After the challenge concluded, the dataset has been made publicly available for academic purposes.

First, conventional 3D object retrieval tasks consider only the object category, where the training and test samples are characterized by the same category settings. Consequently, features extracted from these methods are often optimized to fit the seen categories while lacking generalizability for unseen categories. Under such circumstances, the classification-based retrieval embedding learning methods become invalid in practice. Meanwhile, open-set 3D object retrieval can address this issue more effectively by dealing with unseen categories better. This technique involves training retrieval and representation models using seen-category 3D objects, with unseen-category 3D data subsequently used for retrieval. Nevertheless, our fine-grained retrieval task requires participants to conduct an accurate search to get 3D animal models whose shapes correspond to the query, necessitating consideration of unseen categories and poses (*cf.* Table 1). Compared to searching for 3D general objects of a given category, 3D animal model fine-grained retrieval poses a more significant challenge due to the substantial discrepancy in animal breeds and poses.

Second, participants in our track challenge must solve the considerable domain gap between sketches and 3D shapes when dealing with differently posed animals. Furthermore, human sketches on existing datasets tend to be semi-photorealistic and drawn by experts. In contrast, our dataset comprises more diverse sketches, including abstract sketches drawn by amateurs, semi-photorealistic sketches, and sketches in different styles (*cf.* Fig. 1). As such, this task proves significantly more challenging than conventional sketch-based object retrieval tasks. We anticipate that the sketch-based 3D animal fine-grained retrieval task can pave the way for a new research direction and exciting, practical applications.

The remainder of this paper is organized as follows. Section 2 provides an overview of related work. Section 3 presents the ANIMAR dataset and the evaluation measures used in this SHREC contest. Section 4 describes the participant statistics. In Section 5, the methods of the participating teams are presented. The evaluation results and an in-depth analysis of their

Table 1: SHREC challenge tracks for 3D object retrieval.

SHREC Challenge	Year	Query Type	Training Category	Testing Category
Hameed <i>et al.</i> [18]	2018	Image	Seen	Seen
Hameed <i>et al.</i> [19]	2019	Image	Seen	Seen
Li <i>et al.</i> [20]	2019	Image	Seen	Seen
Li <i>et al.</i> [21]	2020	Image	Seen	Seen
Feng <i>et al.</i> [22]	2022	Image	Seen	Unseen
Li <i>et al.</i> [12]	2012	Sketch	Seen	Seen
Li <i>et al.</i> [13]	2013	Sketch	Seen	Seen
Li <i>et al.</i> [14]	2014	Sketch	Seen	Seen
Juefei <i>et al.</i> [15]	2018	Sketch	Seen	Seen
Juefei <i>et al.</i> [16]	2019	Sketch	Seen	Seen
Qin <i>et al.</i> [17]	2022	Sketch	Seen	Seen
SketchANIMAR	2023	Sketch	Unseen	Unseen

performance are reported in Section 6. Finally, Section 7 concludes the paper and suggests directions for future work.

## 2. Related Work

To recover 3D objects from a database, content-based 3D object retrieval examines the visual contents of the objects, such as color, texture, form, and geometric aspects. Many tracks concentrating on similar problems have been held in previous SHREC competitions (see Table 1) to promote research on content-based 3D object retrieval.

Several SHREC tracks concentrate on retrieving 3D items in a database that resemble the 3D objects used as a query. The attractiveness of sketch-based 3D object retrieval, in particular, stems from the organic and intuitive quality of freehand sketches, and it has garnered much attention in recent years.

Li *et al.* [12, 13, 14] promoted this intriguing study by organizing SHREC tracks of sketch-based 3D shape retrieval. At that time, deep learning was not popular; thus, submitted solutions were based on hand-crafted features such as Scale-Invariant Feature Transform (SIFT), Histogram of Oriented Gradient (HOG), fourier descriptors, bag-of-features, and sparse coding. After that, Juefei *et al.* [15, 16] extended the task to 2D scene sketch-based 3D scene retrieval. Domain adaptation algorithms, such as two-stream CNN with triplet loss, adversarial training, and different data augmentation techniques, were used to resolve the disagreement between two domains (*i.e.*, sketch and 3D object). In addition, a competition for sketch-based 3D form retrieval in the wild was conducted by Qin *et al.* [17], further advancing the task. They used a variety of 3D forms, including models created by scanning genuine objects, as well as large-scale sketches created by amateur artists with a range of sketching abilities. Furthermore, technologies, such as point cloud and multi-view learning using various deep learning architectures, were created to emulate actual retrieval circumstances.

Sketch-based 3D object retrieval methods can be grouped into two categories: model-based and view-based approaches. Model-based methods commonly utilize 3D CNN to extract 3D shape features directly from the original 3D representations.

<sup>1</sup><https://aichallenge.hcmus.edu.vn/sketchanimar>

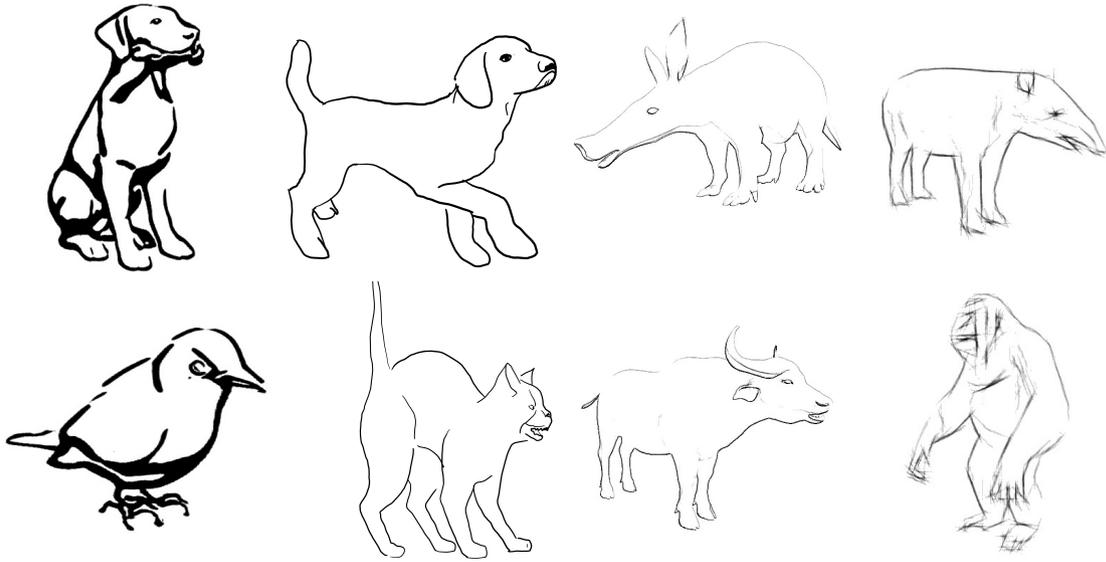


Fig. 1: Sketch ground-truth in ANIMAR dataset, including various sketch types and animal poses.

In the view-based approach, 2D convolutional neural networks (CNN) are frequently used to analyze shape features from a set of 2D view projections.

Regarding the model-based approach, Furuya *et al.* [23] propose Deep Local Feature Aggregation Network (DLAN), which extracts rotation-invariant 3D local features and aggregates them in a single deep architecture. More concretely, the DLAN uses a set of 3D geometric features invariant to local rotation to characterize local 3D regions of a 3D model. The DLAN then compiles the set of features into a (global) rotation-invariant and compact feature for each 3D model. Furthermore, an Octree-based Convolutional Neural Network (O-CNN) [24] is also proposed for 3D shape analysis. O-CNN executes 3D CNN operations on the octants filled by the 3D shape surface using the average normal vectors of a 3D model sampled in the smallest leaf octants as input.

Concerning the view-based approach, Wang *et al.* [25] propose two Siamese Convolutional Neural Networks for the views and the sketches. Moreover, the loss function is designed for within-domain and cross-domain similarities. Similarly, two deep CNNs are proposed by Xie *et al.* [26] for deep feature extraction of sketches and 2D projections of 3D shapes. Next, the authors compute the Wasserstein barycenters of deep features of multiple projections of 3D shapes to form a barycentric representation. Last but not least, Multi-view Convolutional Neural Network (MVCNN) [27] creates a single, compact shape descriptor from data from multiple views of a 3D shape, which improves recognition performance.

In recent competitions of 3D Shape Retrieval Contests (SHREC) [28, 29, 17], teams achieving high performances followed the view-based approach.

### 3. Dataset and Evaluation

#### 3.1. Dataset

In this competition, we constructed a new dataset, namely ANIMAR, which encompasses a corpus of 711 distinct 3D an-

imal models along with 140 sketch queries.

We collected an assemblage of 186 mesh models depicting over 50 diverse categories of animals. These models were diligently sourced from an array of publicly available online resources and video games, including the well-known Planet Zoo video game<sup>2</sup> [30]. The primary goal of our competition track was to simulate real-life scenarios in which users endeavor to identify and explore a diverse range of animal species. To achieve this, **we purposely concealed categorical information during both the training and retrieval stages.** Furthermore, we refined our model database by generating a series of watertight mesh models by reducing the number of faces by 25%, 50%, and 75%, yielding a total of 525 models. Following the work of Douze *et al.* [31], our 3D animal model database is employed for both the training and retrieval phases.

From 186 original mesh models, we randomly selected 60 models for sketch image creation. For each model, we rotated the model and generated 2-3 sketches from distinct viewpoints, thereby producing a total of 140 sketch images to describe the 3D animal models. Notably, we intentionally chose not to create sketches for all animal models in order to prevent participants from utilizing them to train retrieval solutions. Of the 140 sketches, 74 were aligned with their corresponding models in the database, yielding a set of 297 query-model pairs that were utilized for training purposes. The remaining 66 sketches were designated as queries, resulting in 265 query-model pairs employed during the retrieval phase. Unlike existing datasets, which primarily featured semi-photorealistic sketches drawn by experts, **our dataset comprises more diverse sketches, including abstract sketches drawn by amateurs, semi-photorealistic sketches, and sketches in different styles.** This diversity is exemplified in Fig. 1, where the varied nature of the sketches can be observed.

<sup>2</sup><https://www.planetzoogame.com>

### 3.2. Evaluation Metrics

We provide a comprehensive evaluation of the performance of different methods in this track. The following metrics are utilized:

- **Nearest Neighbor (NN)** evaluates top-1 retrieval accuracy.
- **Precision-at-10 (P@10)** is the ratio of relevant items in the top-10 returned results.
- **Normalized Discounted Cumulative Gain (NDCG)** is a measure of ranking quality defined as  $\sum_{i=1}^p \frac{rel_i}{\log_2(i+1)}$ , where  $p$  is the length of the returned rank list, and  $rel_i$  denotes the relevance of the  $i$ -th item.
- **Mean Average Precision (mAP)** is the area under the precision-recall curve. It measures the precision of methods at different levels and then takes the average. mAP is calculated as  $\frac{1}{r} \sum_{i=1}^r P(i)(R(i) - R(i - 1))$ , where  $r$  is the number of retrieved relevant items,  $P(i)$  and  $R(i)$  are the precision and recall at the position of the  $i$ <sup>th</sup> relevant item, respectively.
- **First Tier (FT)** denotes the recall of the top  $m$  retrieval results, where  $m$  is the number of relevant images in the whole database. It measures the accuracy of retrieving the most relevant images among all the possible matches. The FT score is calculated as:  $FT = (\text{number of relevant images retrieved in the top } m) / m$ .
- **Second Tier (ST)** denotes the recall of the top  $2m$  retrieval results, where  $m$  is the number of relevant images in the whole database. It measures the ability to retrieve relevant images within a broader set of results. The ST score is calculated as:  $ST = (\text{number of relevant images retrieved in the top } 2m) / m$ .
- **Fallout Rate (FR)** shows the ratio of non-relevant retrieved items in relation to the total number of non-relevant items available. It measures the system's ability to avoid retrieving non-relevant items. The FR score is calculated using the formula:  $FR = (\text{number of non-relevant items retrieved}) / (\text{number of total non-relevant items in the database})$

### 4. Participants

Eight groups participated in the SketchANIMAR challenge track. Each group was provided with three weeks to complete the challenge. Throughout the contest, a total of 204 runs were submitted. All participating groups were required to register and submit their results along with a detailed description of their methods. It is important to note that the organizers did not participate in the challenge. We remark that three teams opted not to disclose the methods they used in the competition against the SHREC spirit, which was born to compare the performance of algorithms on common data. Thus, they are not reported in this paper. The participant details are provided below (team members will be added upon acceptance):

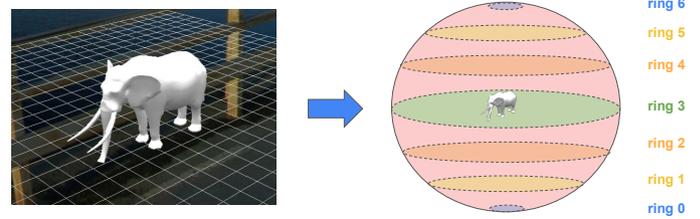


Fig. 2: 3D object represented as view sequences of 7 rings with 12 views on each ring. The chosen latitudes were 0 (the equator),  $\pm 90$  (the poles), and  $\pm 30, \pm 60$ .

- TikTorch team submitted by Nhat-Quynh Le-Pham, Huu-Phuc Pham, Trong-Vu Hoang, Quang-Binh Nguyen, and Hai-Dang Nguyen (see Section 5.2).
- THP team submitted by Truong Hoai Phong (see Section 5.3).
- Etinifni team submitted by Tuong-Nghiem Diep, Khanh-Duy Ho, Xuan-Hieu Nguyen, Thien-Phuc Tran, Tuan-Anh Yang, Kim-Phat Tran, Nhu-Vinh Hoang, and Minh-Quang Nguyen (see Section 5.4).
- V1olet team submitted by Trong-Hieu Nguyen-Mau, Tuan-Luc Huynh, Thanh-Danh Le, Ngoc-Linh Nguyen-Ha, and Tuong-Vy Truong-Thuy (see Section 5.5).
- DH team submitted by Hoai-Danh Vo and Minh-Hoa Doan (see Section 5.6).

## 5. Methods

### 5.1. Overview of Submitted Solutions

All submissions to our track are built upon the foundation of view-based learning. This approach captures the essence of each 3D object by presenting it as a sequence of ring images, as illustrated in Fig. 2. These images are acquired by strategically maneuvering a camera around the object along a predefined path, with each ring consisting of a series of images. In particular, when the camera's trajectory aligns parallel to the ground plane relative to the object, the multi-view method demonstrates remarkable effectiveness, generating valuable images that greatly assist in extracting features for representing three-dimensional objects.

The view-based learning shows more advantages than directly learning the point clouds. This matches well with the settings of our challenge. The 3D objects in our dataset have high-density point clouds with a large number of points. This detail can make it difficult for feature extraction models on point cloud such as PointNet [32] and PointMLP [33] when these models usually randomly sample a specific number of points (1024, for example) in the point cloud of 3D objects. This approach proves particularly useful in scenarios where 3D models are not readily available for querying, but sketches of the objects are, which is often the case in real-world applications.

To facilitate the retrieval task, TikTorch, THP, and Etinifni teams considered the problem as contrastive learning (as shown in Sections 5.2, 5.3, 5.4). Meanwhile, V1olet team formulated

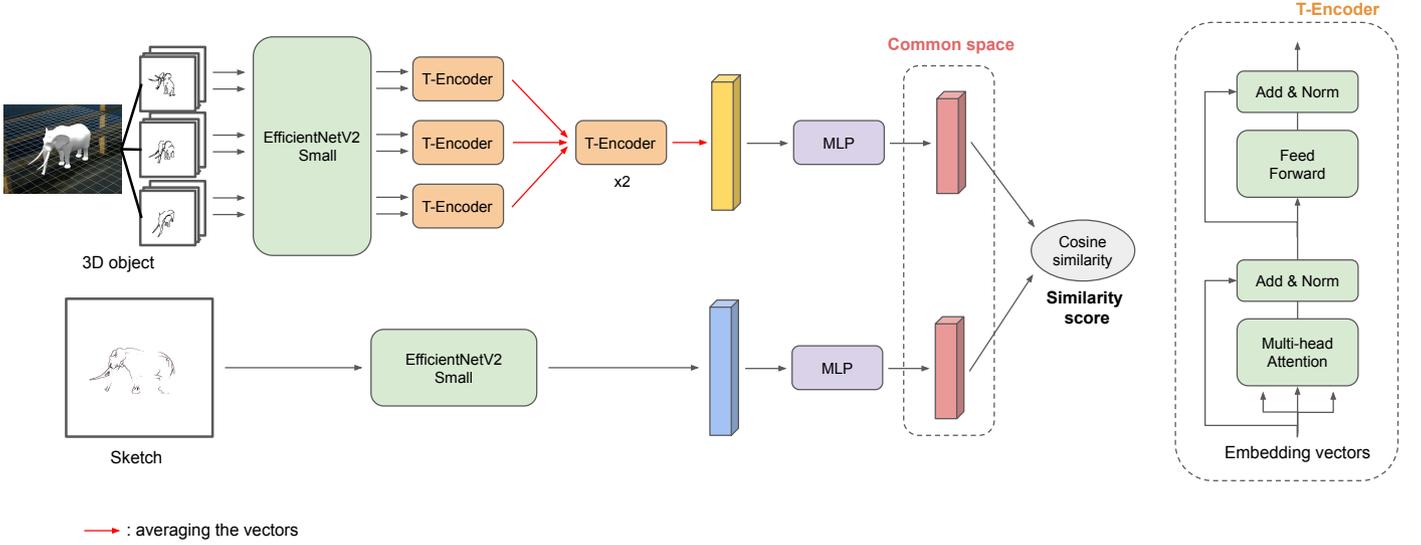


Fig. 3: Proposed framework of TikTorch team.

the task as a classical classification problem (as depicted in Section 5.5). On the other hand, DH team directly extracted and compared non-deep learning features between sketch queries and generated sketches from 3D objects (see Section 5.6).

## 5.2. TikTorch Team

### 5.2.1. Proposed Contrastive Learning Solution

To retrieve 3D objects from sketch queries, they propose a contrastive learning framework where embedding vectors of 3D objects and 2D sketches are learned. The embedding vectors of similar objects and sketches should be closer to each other and vice versa.

The overall architecture of their method is presented in Fig. 3, containing two separate feature extractors for 3D objects and sketch images. The extracted feature vectors are then embedded in the common vector space by two Multi-layer Perceptron (MLP) networks. The contrastive loss used for simultaneous learning of the parameters for models is a customized version of Normalized Temperature-scaled Cross Entropy Loss (NT-Xent) [34].

**Sketch feature extractor.** To extract the features of sketch images, they fine-tune EfficientNetV2-Small [35] pretrained on ImageNet dataset [36]. The models in the EfficientNetV2 family reduce the parameter size significantly while maintaining competitive accuracy on many datasets, which is desirable for simple images, especially sketch images.

**3D object feature extractor.** Each 3D object is represented as a set of 3 rings, and each ring contains 12 images. The 3D object feature extractor has two main phases: extracting the features of each ring (ring extractor) and combining the features of 3 rings to obtain the features of the object.

In the ring extractor, they also fine-tune EfficientNetV2-Small [35], similar to the sketch feature extractor module, to extract the features of 12 images of each ring. These 12 feature vectors then go through an encoder block called T-Encoder [37] to learn the relationship between images in the same ring to decide which image is essential in the current ring and which is

not. After that, they combine these vectors by simply calculating their average to get a single feature vector for each ring.

When they obtain the feature vectors of 3 rings, these vectors are passed into 2 T-encoder blocks to know which ring is useful for the model to learn the features of the current object. Then, the vectors are averaged to get the feature vector of the 3D object.

**Embedding into common space.** To compute the similarity between objects and sketches, their feature vectors must be embedded into a shared space. Since the feature vectors of 3D objects and sketches may have different dimensions, two MLP networks with two layers are utilized. The output layer of each network has the same number of units, ensuring that the feature vectors are transformed into the same vector space. In addition, a Dropout layer [38] is added to each network to prevent overfitting. Once the feature vectors are embedded in the common space, the similarity between two embedding vectors,  $\mathbf{u}$  and  $\mathbf{v}$ , can be computed using the cosine similarity metric.

$$\text{sim}(\mathbf{u}, \mathbf{v}) = \frac{\mathbf{u} \cdot \mathbf{v}}{\|\mathbf{u}\| \|\mathbf{v}\|} \quad (1)$$

**Loss function.** The contrastive loss function used is a customized version of Normalized Temperature-scaled Cross Entropy Loss (NT-Xent) [34]. Given a mini batch of  $2N$  samples  $\{\mathbf{x}_i\}_{i=1}^{2N}$  containing  $N$  objects and  $N$  sketches. They denote  $\mathbf{z}_i$  as the embedding vector of the sample  $\mathbf{x}_i$  in the common space. Let  $P_i$  be the set of indices of samples that are similar to  $\mathbf{x}_i$  in the current mini-batch exclusive of  $i$ , i.e.,  $(\mathbf{x}_i, \mathbf{x}_j)$  is a positive pair for  $j \in P_i$ . Here,  $\mathbf{x}_i$  can belong to many positive pairs, such as two 3D objects that are similar to the same sketch. The loss function for a positive pair  $(\mathbf{x}_i, \mathbf{x}_j)$  is defined as:

$$l_{i,j} = -\log \frac{\exp(\text{sim}(\mathbf{z}_i, \mathbf{z}_j) / \tau)}{\sum_{k=1}^{2N} \mathbb{I}_{[k \neq i, k \notin P_i]} \exp(\text{sim}(\mathbf{z}_i, \mathbf{z}_k) / \tau)}, \quad (2)$$

where  $\mathbb{I}_{[k \neq i, k \notin P_i]} \in \{0, 1\}$  is an indicator function evaluating to 1 if and only if  $k \neq i$  and  $k \notin P_i$ ,  $\tau$  is a temperature parameter.

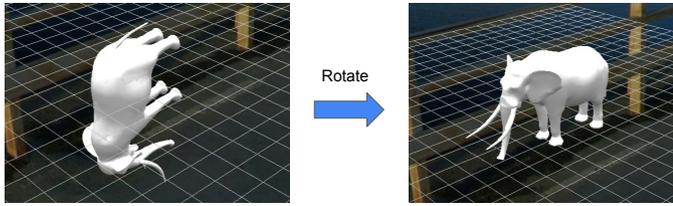


Fig. 4: Example of rotating a 3D object whose axis is not aligned with the majority of objects.

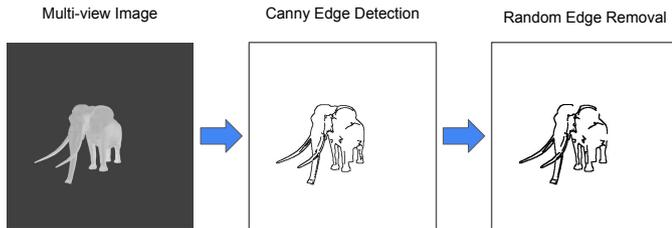


Fig. 5: Multi-view image processing steps.

**Training phase.** During the training process, the optimizer used for training was AdamW [39], along with the StepLR algorithm to reduce the learning. They also applied the  $k$ -fold cross-validation technique with  $k = 5$ .

**Retrieval phase.** They ensemble the results of models trained on  $k$ -fold by max-voting. The similarity between a 3D object and a sketch image is the largest value of the similarity score computed by the five models.

### 5.2.2. Data Augmentation

**Generation of multi-view images for 3D objects.** Before generating batch images from 3D objects, it is essential to ensure axial synchronization of the objects so that the resulting multi-view images with their corresponding camera angles are consistent. To achieve this, they carefully examine the available dataset and identify several objects rotated at a 90-degree angle along the  $O_x$  axis. Then, they apply a consistent rotation to align these objects with the majority of the dataset as in Fig. 4.

Among seven rings, as in Fig. 2, they find that the most informative views are captured from rings 2, 3, and 4, which provide a 360-degree perspective around the object. Hence, they focus on processing the images from these rings to extract the relevant features and information.

To enhance the sketch-like appearance of the multi-view images, they utilize the Canny edge detector [40] to extract edge information. They also add some noises and variations to the edge information by randomly removing edges in the image using a traversal algorithm while preserving the underlying structure and content of the image (see Fig. 5). Figure 6 illustrates the outcome of generating multi-view images for a 3D object.

**Generation of training sketch query images.** Firstly, they cluster similar 3D objects together by some algorithms to check the similarity between the distribution of points in the point clouds and also the manual checking as post-processing. After that, they identify the best-quality object in each cluster, which usually is the most fine-grained object (*i.e.*, the highest number of points in the point cloud). Then, when they generate a sketch-line image for each object in this cluster and use it for

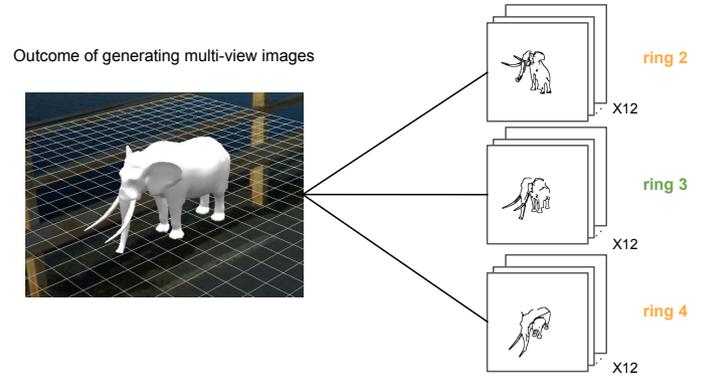


Fig. 6: Outcome of generating multi-view images for 3D objects. Each 3D object is represented by a set of 3 rings, and each ring is a collection of 12 images.

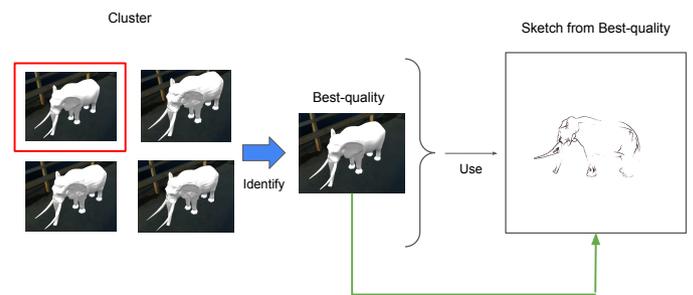


Fig. 7: Optimizing object sketches through clustering: the sketch of the best-quality object for all objects within a group of similar 3D objects is used.

contrastive learning, they pick the image of the best-quality object as the representative sketch (depicted in Fig. 7).

To expand training samples for contrastive learning, they develop a method to generate three queries per object. Each query is randomly chosen from rings 2, 3, and 4 (see Fig. 2) with probabilities of 0.2, 0.6, and 0.2, respectively, as they observe that the majority of informative queries are in ring 3. Once a ring is selected, they randomly choose an image within that ring from the cluster this object belongs to and apply random Canny edge [40] or Artline [41] techniques, along with image horizontal flip and rotation transformations. By implementing this process, they can significantly increase the number of our training samples from about 100 to 2500 while maintaining a high level of quality.

## 5.3. THP Team

### 5.3.1. Architecture of Proposed Network

Fig. 8 illustrates an overview of the proposed network. To evaluate the similarity of two given sketch images, they compare the distance of global features and local features and then combine the results.

To extract global features, they use the pre-trained CLIP model [42]. The input of the CLIP model is the dilated Canny edge extracted from multi-view images. After that, CLIP feature vectors of view images and sketch queries are matched using cosine similarity. The final score is calculated as the maximum of scores of 4 views, in which each view score is calculated by the sum of the six highest similarities.

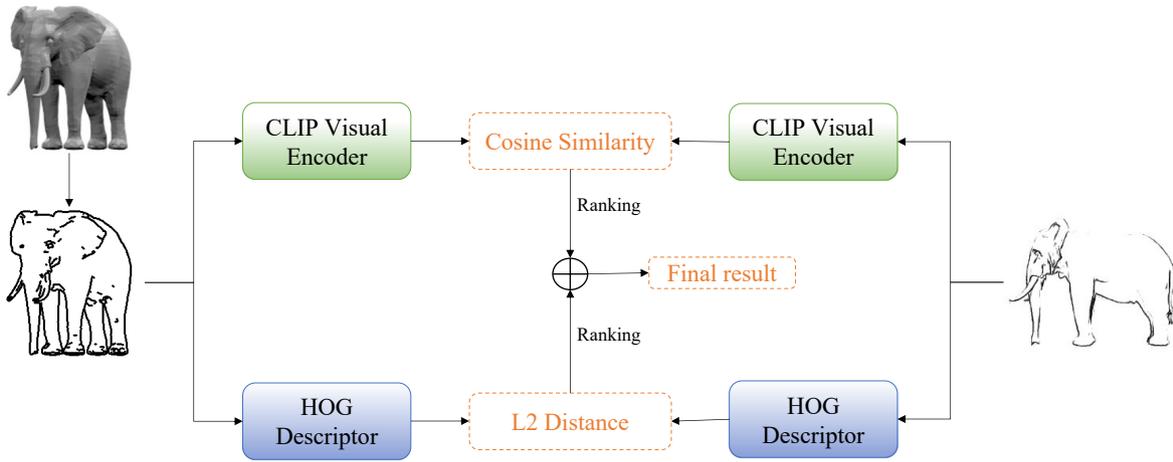


Fig. 8: Proposed framework of THP team.

To increase local information awareness, they use the HOG descriptor [43] on both sketch and multi-view images. The HOG vectors are then matched using the L2 distance. They are also ranked like CLIP feature vectors.

Finally, they combine CLIP and HOG similarity scores as follows:  $Score = \alpha * CLIP\ score + (1 - \alpha) * HOG\ score$ , where  $\alpha = 0.7$ .

### 5.3.2. 2D Shape Projection

They use four camera setups to take multiple views of 3D objects:

- For the first camera setup, assuming the 3D object is initially aligned along the  $z$ -axis, the camera is aligned on the  $Oxy$  plane and looks at the center of the object. The camera is moved around the subject to create 12 views from a distance of 30 degrees each time.
- For the second camera setup, the camera is raised to 30 degrees above the  $Oxy$  plane and moved around to create the next 12 views.
- For the third camera setup, the camera is placed on the  $Oyz$  plane and looks at the object's center. The camera is moved like the first setup to create the next 12 views.
- Camera for the last setup is raised from the third setup to 30 degrees to create the next 12 views.

There are a total of 48 views for each object. Note that it is possible to create images from other directions, but according to their tests, from these 48 views, they can observe the characteristics of objects.

### 5.3.3. Sketch Pre-processing

To reduce the domain gap, they apply Canny edge algorithm [40] for each 2D projection image to create an image similar to the sketch image because the sketch is also a special type of edge. Both sketch and Canny edge images are then cropped and resized to  $224 \times 224$  with the padding of 5 pixels. After that, edges are further clarified using the dilation morphology algorithm. The sketch pre-processing pipeline is shown in Fig. 9.

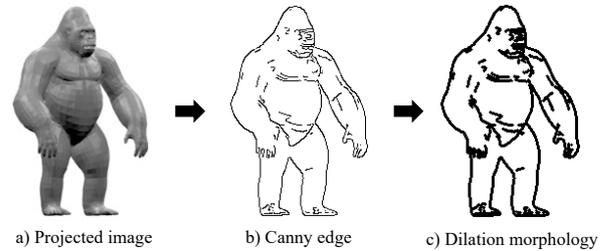


Fig. 9: Pipeline of sketch pre-processing.

## 5.4. Etinifni Team

### 5.4.1. Overview of Proposed Solution

Due to the nature of the retrieval tasks, they build a deep learning framework using the CLIP model [42] as the backbone (See Fig. 10). The framework performs the following steps:

1. Pre-process the dataset to re-direct the 3D objects into one single vertical orientation.
2. Extract multi-views (*i.e.*, 12 random views around the objects in uniform angles) from the 3D objects.
3. Encode the 2D sketches and the view images of 3D objects using the AutoEncoder built upon the ResNet50.
4. Reduce the size of the feature vectors from 512 to 128 through the projection head.
5. Compare feature vectors of sketches and the 3D objects using the cosine similarity function. After that, they can identify the matching pairs of sketches and 3D objects.

### 5.4.2. Data Pre-Processing

**Resize objects:** They resize the 3D objects to fit the 3D object inside an imaginary box of  $2 \times 2 \times 2$  by re-scaling the dimensions  $(x, y, z)$  of the 3D objects into the new dimension  $(x', y', z')$  not greater than 2 using the following formula:

$$(x', y', z') = \left( \frac{2x}{\max(x, y, z)}, \frac{2y}{\max(x, y, z)}, \frac{2z}{\max(x, y, z)} \right) \quad (3)$$

**Re-orientate objects:** They re-direct the orientation of the 3D objects manually so that the 3D objects are standing (*i.e.*, the objects are in an erect position).

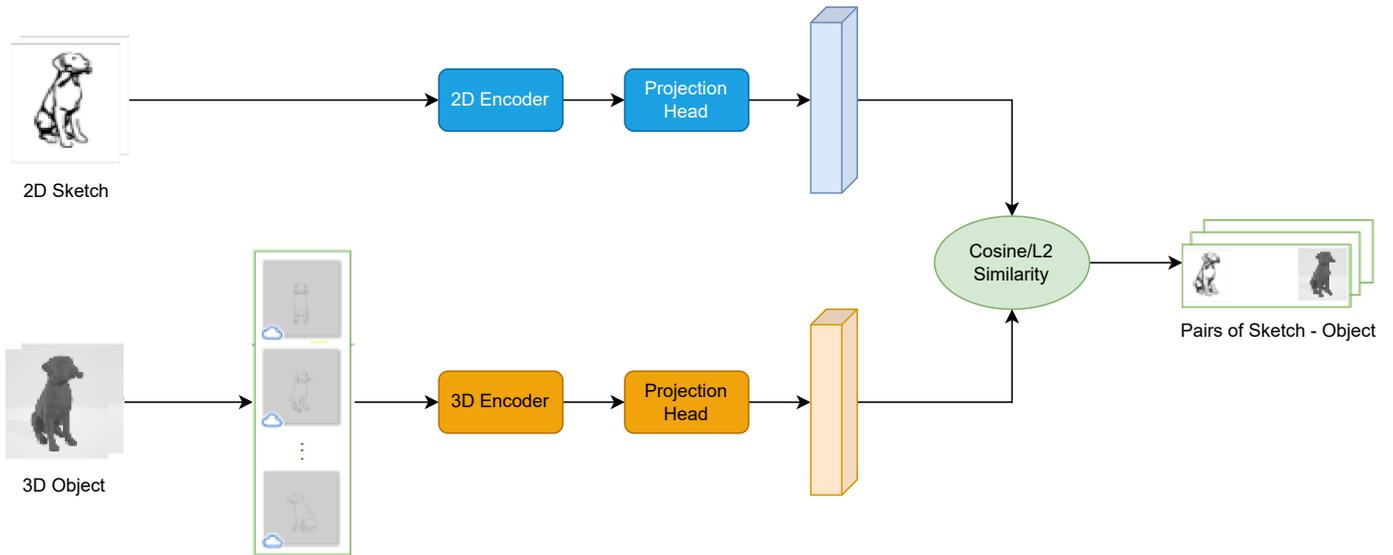


Fig. 10: Proposed framework of Etinifni team.

#### 5.4.3. Multi-view Generation

They first rotate the camera as in Fig. 2. They apply 3D image rendering at each position to provide a more accurate and detailed view of objects. This enables the model to extract the necessary features for retrieval with greater precision, even where subtle variations in shape and texture can be crucial in determining an object's identity. Models then extract the most detailed and accurate information possible, leading to more robust and reliable results.

#### 5.4.4. Data Augmentation

Due to the limited available training data, data augmentation techniques are applied to increase the number of training data:

##### Outline modifications:

- **Thickness:** is a significant attribute of 2D objects as it plays a crucial role in defining their physical properties and functionality. By adjusting the thickness, they generate new sketches with varying line lengths, thereby introducing diversity in the visual depiction of the objects.
- **Geometry:** explores the properties and relationships of shapes, sizes, positions, and dimensions of objects in space. By modifying the stroke of these edges, they can generate variations in the appearance of the zigzag patterns. This allows for a deeper investigation into the geometric characteristics of the objects and the effects of stroke adjustments on their overall geometry.

##### Image processing:

- **Random deletion:** To ensure the robustness of the sketch, they partition the 2D images into multiple blocks. Each block is equal in size and contains a subset of pixels from the original image. They randomly select a portion of these blocks and remove them from the image to simulate missing strokes on the data set.

- **Image compression:** Because the provided query images are of low quality, a compressor is utilized to reduce the image quality to generate the sketches.

##### View extraction:

- The camera is set up at a suitable distance, height, and orientation to ensure comprehensive coverage of the object's surface. The camera views are then randomly selected to capture diverse angles for robust 3D reconstruction.

#### 5.5. Violet Team

##### 5.5.1. Proposed Classification Approach

They formulate the retrieval task as a classical classification problem. Notably, several potential CNNs, including EfficientNet [44], EfficientNetV2 [35], and ConvNeXt [45] are employed to recognize whether the sketch query and 3D objects are the same class. These networks are renowned for their remarkable feature extraction capabilities and are considered state-of-the-art in image recognition. The used models also are lightweight and suitable for real-life applications while achieving considerable performance.

**Ensemble Solution.** Figure 11 illustrates the proposed ensemble approach by averaging the predictions of each model. This approach effortlessly helps mitigate individual models' potential shortcomings, resulting in improved performance and robust generalization to varying data distributions.

To further improve the accuracy of models, they also utilize Test Time Augmentation (TTA), which involves applying a range of transformations such as rotations, flips, and translations to test images and averaging the results to obtain the final prediction. Specifically, they utilize horizontal flipping to provide additional perspectives of the original images. This technique not only enhances generalization but also enables

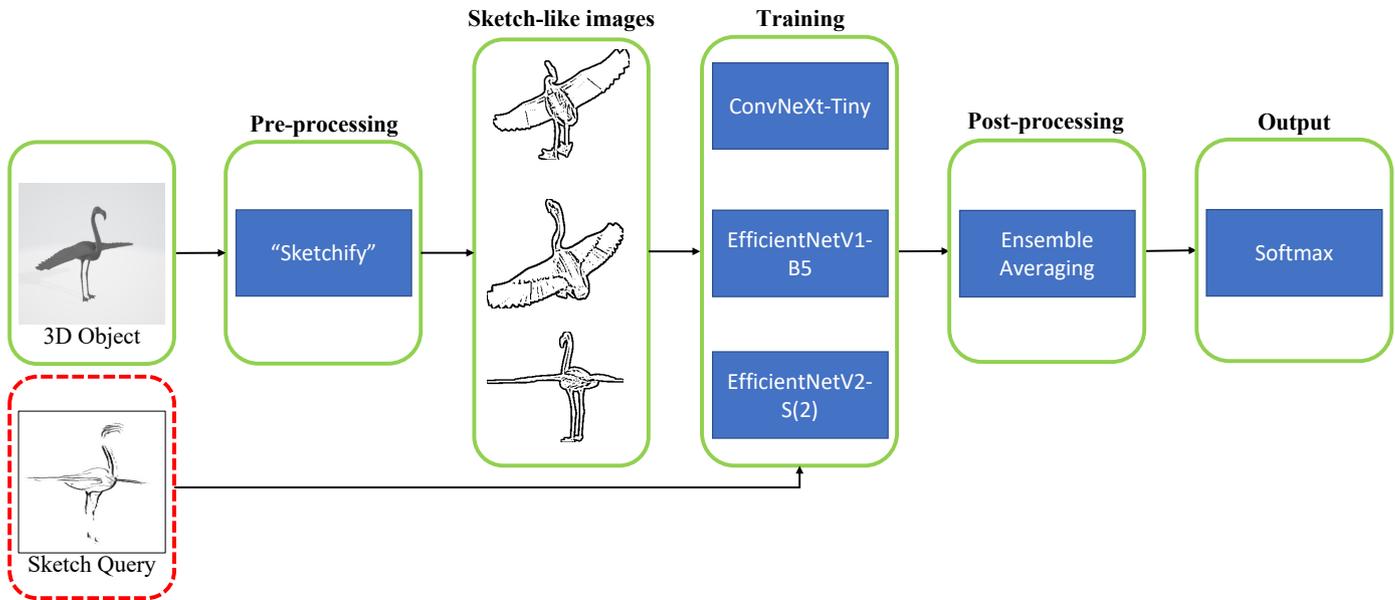


Fig. 11: Proposed framework of V1olet team.

the model to recognize objects that may be oriented differently from those in the training set.

**Training Phase.** To evaluate the performance of used models, they created a validation set by randomly leaving out 10% of samples from each class in the training data. Pre-trained models on ImageNet were fine-tuned using the remaining training sets. They also employed the cross-entropy loss with label smoothing of 0.1 to prevent overfitting and improve the generalizability of models. Sketchified multi-view images were jointly trained with the original sketch queries to enforce networks to recognize them as the same class. During both training and inference, an image size of  $384 \times 384$  was utilized. All networks were trained for 20 epochs using the Adam optimizer [46] with a learning rate of 0.0001. Finally, they selected the models with the best validation accuracy for ensembling.

**Retrieval Phase.** Given a sketch query, the ensemble averaging of CNN models produces a set of softmax probabilities. These probabilities are then used to identify whether the sketch query and sketch-like images generated from the 3D object belong to the same class. The softmax probabilities also serve as a ranking metric, allowing for sorting the retrieved 3D objects by relevance.

### 5.5.2. Data Pre-processing

Figure 12 demonstrates an overview of the proposed data pre-processing pipeline. In general, it can be divided into three steps: Ringview extraction, color reversal, and sketchify.

**Ringview extraction.** Extracting multiple views of an object can be highly advantageous for various applications, including 3D object retrieval. They extract multiple views from 7 rings with 12 views on each ring, like Fig. 2. By providing different perspectives of the 3D object, these multiple views can extract more robust and detailed features to train models with greater accuracy for 3D object retrieval. Thus, ringview processing is a valuable technique that can improve the accuracy of models for 3D object retrieval.

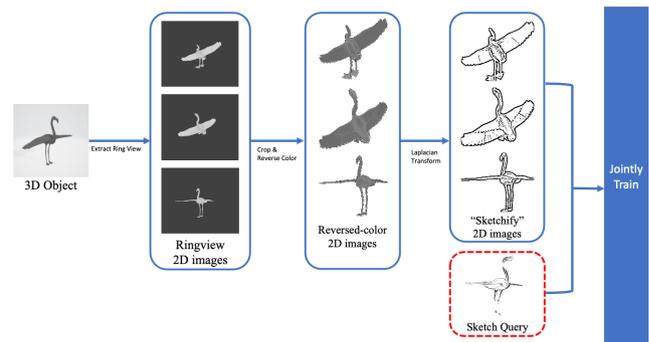


Fig. 12: The proposed data pre-processing pipeline

**Color reversal.** It is also crucial to consider the background of the images when matching 3D objects with the querying sketches. Typically, these sketches usually have a white background, while the multiple 2D images obtained from the ringview extraction step have a grey background. To solve this problem, they merely flip the color of the ringview image, making the backdrop translucent to match the background of the target sketch queries. Therefore, the resulting images better resemble the sketch queries and support the further "sketchify" procedure.

**Sketchify.** Laplacian Transform is utilized to produce images that are more similar to sketches. Particularly, the Laplacian Transform, as a linear operator, is applied to the grayscale image to generate a second-order derivative image that enhances the edges and transitions of the image [47]. The operator produces a sketch-like version of an image by thresholding the Laplacian image to obtain a binary edge map, which is used to synthesize a sketch-like representation of the image. This transformation enables more efficient comparison of the views with the sketch query as the critical matching features become more pronounced, facilitating accurate matching.

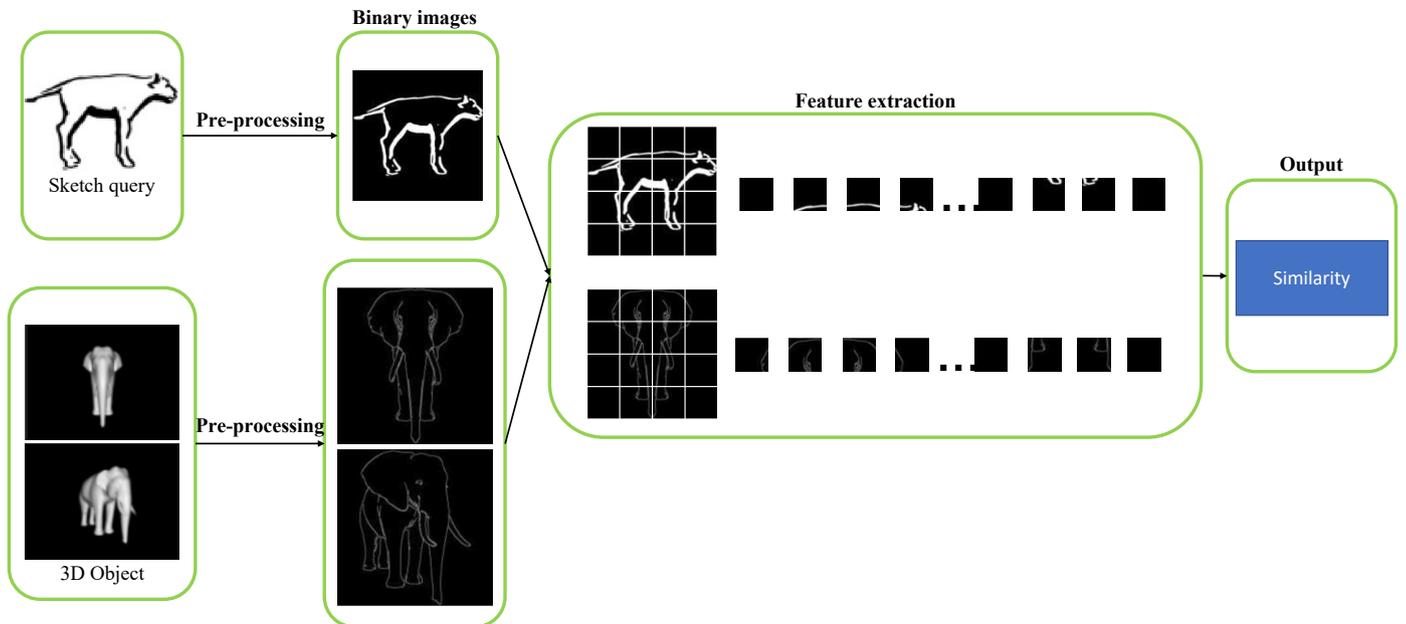


Fig. 13: Proposed framework of DH team.

## 5.6. DH Team

### 5.6.1. Proposed Method

They propose a simple method to measure the distance between 2 images with fast execution speed due to the small number of computations. Based on the observation that 3D objects with the same shape will have the same distribution of pixels in 2D projection space, they propose dividing the image into small parts and then comparing each area to measure the similarity between the two images. Fig. 13 shows the proposed framework, including four main modules: 2D sketch processing, 3D model processing, image feature extraction, and feature matching.

**2D sketch processing.** They crop images intending to keep only the part containing the animals and remove the background. Since sketches in the dataset only contain the animal and no other extraneous details in the background, they define the animal’s bounding box as follows: First, sketches are converted to binary images and then inverted. Then they find the top and bottom image lines with the value 255. These two lines correspond to the top and bottom edges of the bounding box. The same method is applied to the left and right edges. After that, they crop the part containing the animal into a square with the size of a maximum bounding-box height and width and then resize the image to  $224 \times 224$ .

**3D model processing.** They first rotate the 3D objects using the Open3D library and capture the 3D model from 21 perspectives. After that, view images are cropped similarly to 2D sketches.

**Image feature extraction.** The images are divided into  $4 \times 4$  squares, and then the ratio of total pixel values in each square to total pixel values of the whole image is considered as the score of the square. At the end of this step, a 16-dimensional feature vector represents an image, including the sketch and view images.

**Feature matching.** Given a sketch query image  $Q$  and a 3D model  $R$ , features are extracted to obtain  $f^Q$  representing the sketch image and  $f_i^R, i = 1..21$  representing the 21 view images corresponding to the 3D model. The similarity score between the sketch query  $Q$  and the 3D model  $R$  is defined as follows:

$$D(Q, R) = \min_i \|f^Q - f_i^R\|_2. \quad (4)$$

## 6. Results and Discussions

In our SketchANIMAR track, the submitted runs are evaluated on two subsets: the public and private tests. The private test comprised 66 sketch queries, resulting in 265 query-model pairs. To ensure fairness and prevent potential cheating, approximately half of the private test (30 sketch queries) was randomly selected and designated as the public test subset. The leaderboard for the private test was revealed after the challenge concluded.

Tables 2 and 3 display the leaderboard outcomes for the public and the private test, reporting only the best-performing runs submitted by each team. However, to ensure a fair comparison, our analysis focuses solely on the results from the private test, which evaluated all the submitted sketch queries.

As seen in Table 3, the TikTorch and Violet team’s presented methods repeatedly stood out as the top effective strategies. On 6 out of 7 performance metrics (P@10, NDCG, mAP, FT, ST, and FR), TikTorch outperformed rival teams by a wide margin. In terms of NN metric, the Violet team secured the top 1 position, indicating that this team focuses on the best search instead of neighboring search. Meanwhile, THP took up the third position on public and private leaderboards. Most teams achieving high performances apply the view-base approach, which analyzes shape features from 2D view projections. For the public test, similar findings are also shown in Table 2, excepting that the top-1 team (TikTorch) only takes the fourth position on

Table 2: Leaderboard results of SketchANIMAR competition. Best run results on the public test.

Team	NN	P@10	NDCG	mAP	FT	ST	FR
TikTorch	<b>0.533 (1)</b>	<b>0.280 (1)</b>	<b>0.708 (1)</b>	<b>0.570 (1)</b>	0.192 (4)	0.333 (4)	<b>0.0102 (1)</b>
V1olet	0.467 (2)	0.213 (2)	0.613 (2)	0.411 (2)	0.317 (2)	0.492 (2)	0.0111 (2)
THP	0.433 (3)	0.207 (3)	0.601 (3)	0.399 (3)	0.300 (3)	0.450 (3)	0.0112 (3)
Etinifni	0.200 (4)	0.147 (4)	0.489 (4)	0.303 (4)	<b>0.475 (1)</b>	<b>0.650 (1)</b>	0.0121 (4)
DH	0.100 (5)	0.080 (5)	0.361 (5)	0.140 (5)	0.133 (5)	0.192 (5)	0.0130 (5)

Table 3: Leaderboard results of SketchANIMAR competition. Best run results on the private test.

Team	NN	P@10	NDCG	mAP	FT	ST	FR
TikTorch	0.470 (2)	<b>0.255 (1)</b>	<b>0.669 (1)</b>	<b>0.522 (1)</b>	<b>0.424 (1)</b>	<b>0.583 (1)</b>	<b>0.0105 (1)</b>
V1olet	<b>0.500 (1)</b>	0.232 (2)	0.640 (2)	0.453 (2)	0.379 (2)	0.534 (2)	0.0109 (2)
THP	0.409 (3)	0.226 (3)	0.608 (3)	0.421 (3)	0.333 (3)	0.504 (3)	0.0110 (3)
Etinifni	0.136 (4)	0.158 (4)	0.473 (4)	0.274 (4)	0.174 (4)	0.345 (4)	0.0119 (4)
DH	0.136 (4)	0.088 (5)	0.372 (5)	0.158 (5)	0.133 (5)	0.205 (5)	0.0129 (5)

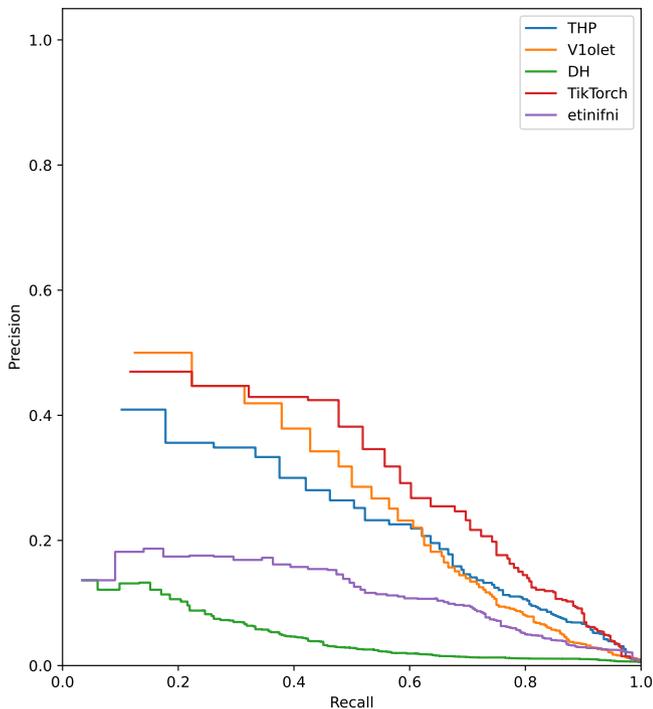


Fig. 14: The visualization of precision-recall curves of submissions on the private test of teams. It can be seen that the TikTorch team achieves the best average performance with the highest area under the curve, while V1olet obtains the highest precision at low recall (recall  $\leq 0.2$ ) among the eight teams. These insights are compatible with results in Table 3, TikTorch and V1olet secure the top 1 position regarding mAP and NN, respectively.

FT and ST metrics. The first two teams, TikTorch and V1olet, achieve the top two ranks. In contrast, the performance of the other teams (*i.e.*, THP, Etinifni, and DH) is constant across test sets, both private and public. Summary, the findings show the challenges of our ANIMAR dataset when participants did not achieve excellent performance (the best NN result is only around 0.5, and P@10 results are smaller than 0.3). It also suggests that there is room to improve the performance of this research direction.

Figure 14 illustrates the precision-recall curves of submis-

sions on the private test of teams. It is clear that among the teams, the TikTorch team has the best average performance with the biggest area under the curve, while V1olet has the maximum precision at low recall (recall  $\leq 0.2$ ). The main difference between TikTorch’s and V1olet’s methods is that TikTorch follows a contrastive learning approach while V1olet leverages a classification-based one with a softmax layer. The softmax score just implies whether an object and a sketched query belong to the same category and can not indicate much similarity between them. In general, V1olet’s method works with objects in the same category as sketched queries and is less robust to negative samples than the contrastive learning approach. These conclusions are consistent with Table 3’s findings, which show that TikTorch and V1olet rank first in terms of mAP and NN. Furthermore, the precision-recall curves of three teams, TikTorch, V1olet, and THP, share a similar shape, which shows the effectiveness at low recall threshold (recall  $\leq 0.4$ ) but drops significantly at high one (recall  $\geq 0.8$ ).

In conclusion, view-based learning methods have proven effective in achieving high performance. The difficulty of feature extraction models [32, 48] can be attributed to the high point cloud density of the 3D objects in the ANIMAR dataset. It is important to remember that these models often randomly sample a certain number of pointclouds (*e.g.*, 1024). Contrarily, by utilizing the semantic data and representation of 3D objects, the use of view pictures obtained by moving the trajectory camera, as demonstrated in Fig. 2, enhances feature learning. This further demonstrates the effectiveness of the view-based learning strategy for retrieving 3D objects.

## 7. Conclusion

This paper introduces a novel track for sketch-based retrieval of fine-grained 3D animal models along with a newly constructed ANIMAR dataset. Our SHREC 2023 challenge track is designed to simulate real-life scenarios and requires participants to retrieve 3D animal models based on complex and detailed sketches. The challenge received submissions from eight teams; however, the evaluated results in this paper include

methods from five of the eight teams. These submissions resulted in a total of 204 runs with different approaches. The evaluated results of this track were satisfactory but also revealed the difficulties of the task at hand.

In future research, we aim to expand the dataset by collecting a more diverse range of 3D animal models that encompass a wider variety of species, environmental contexts, and postures. This can enhance the generalization capability of potential solutions and improve performance on unseen 3D animal models. Additionally, we intend to generate synthetic data and texture to augment 3D animal models with different postures, backgrounds, and patterns to train more effective and robust representation models. We believe that by exploring these research avenues, we can advance the state-of-the-art in 3D object retrieval.

### CRedit authorship contribution statement

**Trung-Nghia Le:** Conceptualization, Writing – review & editing, Project administration, Supervision. **Tam V. Nguyen:** Conceptualization, Writing – review & editing. **Minh-Quan Le:** Software, Writing – review & editing. **Trong-Thuan Nguyen:** Software, Writing – review & editing. **Viet-Tham Huynh:** Data curation. **Trong-Le Do:** Software, Investigation. **Khanh-Duy Le:** Visualization. **Mai-Khiem Tran:** Data curation. **Nhat Hoang-Xuan:** Software. **Thang-Long Nguyen-Ho:** Software. **Vinh-Tiep Nguyen:** Conceptualization. **Nhat-Quynh Le-Pham:** Methodology, Writing – original draft. **Huu-Phuc Pham:** Methodology, Writing – original draft. **Trong-Vu Hoang:** Methodology, Writing – original draft. **Quang-Binh Nguyen:** Methodology, Writing – original draft. **Hai-Dang Nguyen:** Methodology, Writing – original draft. **Trong-Hieu Nguyen-Mau:** Methodology, Writing – original draft. **Tuan-Luc Huynh:** Methodology, Writing – original draft. **Thanh-Danh Le:** Methodology, Writing – original draft. **Ngoc-Linh Nguyen-Ha:** Methodology, Writing – original draft. **Tuong-Vy Truong-Thuy:** Methodology, Writing – original draft. **Truong Hoai Phong:** Methodology, Writing – original draft. **Tuong-Nghiem Diep:** Methodology, Writing – original draft. **Khanh-Duy Ho:** Methodology, Writing – original draft. **Xuan-Hieu Nguyen:** Methodology, Writing – original draft. **Thien-Phuc Tran:** Methodology, Writing – original draft. **Tuan-Anh Yang:** Methodology, Writing – original draft. **Kim-Phat Tran:** Methodology, Writing – original draft. **Nhu-Vinh Hoang:** Methodology, Writing – original draft. **Minh-Quang Nguyen:** Methodology, Writing – original draft. **Hoai-Danh Vo:** Methodology, Writing – original draft. **Minh-Hoa Doan:** Methodology, Writing – original draft. **Akihiro Sugimoto:** Conceptualization. **Minh-Triet Tran:** Conceptualization, Supervision, Funding acquisition, Writing - Review & Editing.

### Data availability

After the challenge concluded, the dataset has been made publicly available for academic purposes.

### Declaration of competing interest

The authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.

### Acknowledgments

This work was funded by the Vingroup Innovation Foundation (VINIF.2019.DA19) and National Science Foundation Grant (NSF#2025234).

### References

- [1] Stotko, P, Krumpen, S, Hullin, MB, Weinmann, M, Klein, R. Slamcast: Large-scale, real-time 3D reconstruction and streaming for immersive multi-client live telepresence. *IEEE transactions on visualization and computer graphics* 2019;25(5):2102–2112.
- [2] Liu, X, Kofman, J. Real-time 3D surface-shape measurement using background-modulated modified fourier transform profilometry with geometry-constraint. *Optics and Lasers in Engineering* 2019;115:217–224.
- [3] Wang, J, Mueller, F, Bernard, F, Sorli, S, Sotnychenko, O, Qian, N, et al. RGB2hands: real-time tracking of 3D hand interactions from monocular RGB video. *ACM Transactions on Graphics (ToG)* 2020;39(6):1–16.
- [4] Guo, H, Peng, S, Lin, H, Wang, Q, Zhang, G, Bao, H, et al. Neural 3D scene reconstruction with the manhattan-world assumption. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 2022, p. 5511–5520.
- [5] Yookwan, W, Chinnasarn, K, So-In, C, Horkaew, P. Multimodal fusion of deeply inferred point clouds for 3D scene reconstruction using cross-entropy icp. *IEEE Access* 2022;10:77123–77136.
- [6] Li, J, Gao, W, Wu, Y, Liu, Y, Shen, Y. High-quality indoor scene 3D reconstruction with RGB-D cameras: A brief review. *Computational Visual Media* 2022;8(3):369–393.
- [7] Gümeli, C, Dai, A, Nießner, M. Roca: robust CAD model retrieval and alignment from a single image. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 2022, p. 4022–4031.
- [8] Manda, B, Kendre, PP, Dey, S, Muthuganapathy, R. Sketchcleanet—a deep learning approach to the enhancement and correction of query sketches for a 3D CAD model retrieval system. *Computers and Graphics* 2022;107:73–83.
- [9] Salihi, D, Steinbach, E. SGPCR: Spherical gaussian point cloud representation and its application to object registration and retrieval. In: *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*. 2023, p. 572–581.
- [10] Koca, BA, Çubukçu, B, Yüzgeç, U. Augmented reality application for preschool children with unity 3D platform. In: *2019 3rd International Symposium on Multidisciplinary Studies and Innovative Technologies (ISMSIT)*. IEEE; 2019, p. 1–4.
- [11] Guo, C, Jiang, T, Chen, X, Song, J, Hilliges, O. Vid2avatar: 3D avatar reconstruction from videos in the wild via self-supervised scene decomposition. *arXiv preprint arXiv:230211566* 2023;.
- [12] Li, B, Schreck, T, Godil, A, Alexa, M, Boubekur, T, Bustos, B, et al. SHREC'12 track: Sketch-based 3D shape retrieval. In: *3DOR@Eurographics*. 2012, p. 109–118.
- [13] Li, B, Lu, Y, Godil, A, Schreck, T, Aono, M, Johan, H, et al. SHREC'13 track: Large scale sketch-based 3D shape retrieval. 2013.
- [14] Li, B, Lu, Y, Li, C, Godil, A, Schreck, T, Aono, M, et al. Shrec'14 track: Extended large scale sketch-based 3D shape retrieval. In: *Eurographics workshop on 3D object retrieval*; vol. 2014. 2014, p. 121–130.
- [15] Yuan, J, Li, B, Lu, Y, Bai, S, Bai, X, Bui, NM, et al. Shrec'18 track: 2d scene sketch-based 3D scene retrieval. *Eurographics Workshop on 3D Object Retrieval* 2018;18:70.
- [16] Yuan, J, Abdul-Rashid, H, Li, B, Lu, Y, Schreck, T, Bui, NM, et al. Shrec'19 track: Extended 2d scene sketch-based 3D scene retrieval. *Eurographics Workshop on 3D Object Retrieval* 2019;18:70.

- [17] Qin, J, Yuan, S, Chen, J, Ben Amor, B, Fang, Y, Hoang-Xuan, N, et al. Shrec'22 track: Sketch-based 3D shape retrieval in the wild. *Computers and Graphics* 2022;.
- [18] Abdul-Rashid, H, Yuan, J, Li, B, Lu, Y, Bai, S, Bai, X, et al. 2D Image-Based 3D Scene Retrieval. In: Telea, A, Theoharis, T, Veltkamp, R, editors. *Eurographics Workshop on 3D Object Retrieval*. 2018;.
- [19] Abdul-Rashid, H, Yuan, J, Li, B, Lu, Y, Schreck, T, Bui, NM, et al. SHREC'19 track: Extended 2d scene image-based 3D scene retrieval. *Eurographics Workshop on 3D Object Retrieval* 2019;700:70.
- [20] Li, W, Liu, A, Nie, W, Song, D, Li, Y, Wang, W, et al. SHREC 2019-monocular image based 3D model retrieval. In: *Eurographics Workshop 3D Object Retrieval*. 2019, p. 1–8.
- [21] Li, W, Song, D, Liu, A, Nie, W, Zhang, T, Zhao, X, et al. SHREC 2020 track: extended monocular image based 3d model retrieval. In: *Eurographics Workshop 3D Object Retrieval*. 2020;.
- [22] Feng, Y, Gao, Y, Zhao, X, Guo, Y, Bagewadi, N, Bui, NT, et al. SHREC'22 track: Open-set 3D object retrieval. *Computers & Graphics* 2022;107:231–240.
- [23] Furuya, T, Ohbuchi, R. Deep aggregation of local 3D geometric features for 3D model retrieval. In: *Proceedings of the British Machine Vision Conference (BMVC)*. 2016;.
- [24] Wang, PS, Liu, Y, Guo, YX, Sun, CY, Tong, X. O-CNN: Octree-based convolutional neural networks for 3D shape analysis. *ACM Trans Graph* 2017;.
- [25] Wang, W, Shen, J, Porikli, F. Saliency-aware geodesic video object segmentation. In: *CVPR*. 2015, p. 3395–3402.
- [26] Xie, S, Girshick, R, Dollár, P, Tu, Z, He, K. Aggregated residual transformations for deep neural networks. In: *CVPR*. 2017, p. 1492–1500.
- [27] Su, H, Maji, S, Kalogerakis, E, Learned-Miller, EG. Multi-view convolutional neural networks for 3D shape recognition. In: *ICCV*. 2015;.
- [28] Savva, M, Yu, F, Su, H, Kanezaki, A, Furuya, T, Ohbuchi, R, et al. Shrec'17 track large-scale 3d shape retrieval from shapenet core55. In: *Proceedings of the Workshop on 3D Object Retrieval*. 2017;.
- [29] Moscoso Thompson, E, Biasotti, S, Giachetti, A, Tortorici, C, Werghe, N, Obeid, AS, et al. Shrec 2020: Retrieval of digital surfaces with similar geometric reliefs. *Computers and Graphics* 2020;.
- [30] Wu\*, Y, Chen\*, Z, Liu, S, Ren, Z, Wang, S. CASA: Category-agnostic skeletal animal reconstruction. In: *Neural Information Processing Systems*. 2022;.
- [31] Douze, M, Tolias, G, Pizzi, E, Papakipos, Z, Chausson, L, Radenovic, F, et al. The 2021 image similarity dataset and challenge. *arXiv preprint arXiv:210609672* 2021;.
- [32] Qi, CR, Su, H, Mo, K, Guibas, LJ. Pointnet: Deep learning on point sets for 3D classification and segmentation. In: *Conference on Computer Vision and Pattern Recognition*. 2017, p. 652–660.
- [33] Ma, X, Qin, C, You, H, Ran, H, Fu, Y. Rethinking network design and local geometry in point cloud: A simple residual mlp framework. *arXiv preprint arXiv:220207123* 2022;.
- [34] Chen, T, Kornblith, S, Norouzi, M, Hinton, G. A simple framework for contrastive learning of visual representations. In: *International conference on machine learning*. 2020, p. 1597–1607.
- [35] Tan, M, Le, Q. Efficientnetv2: Smaller models and faster training. In: *International conference on machine learning*. 2021, p. 10096–10106.
- [36] Russakovsky, O, Deng, J, Su, H, Krause, J, Satheesh, S, Ma, S, et al. ImageNet Large Scale Visual Recognition Challenge. *International Journal of Computer Vision* 2015;115(3):211–252.
- [37] Vaswani, A, Shazeer, N, Parmar, N, Uszkoreit, J, Jones, L, Gomez, AN, et al. Attention is all you need. *Advances in neural information processing systems* 2017;30.
- [38] Hinton, GE, Srivastava, N, Krizhevsky, A, Sutskever, I, Salakhutdinov, RR. Improving neural networks by preventing co-adaptation of feature detectors. *arXiv preprint arXiv:12070580* 2012;.
- [39] Loshchilov, I, Hutter, F. Decoupled weight decay regularization. *arXiv preprint arXiv:171105101* 2017;.
- [40] Canny, J. A computational approach to edge detection. *IEEE T-PAMI* 1986;(6):679–698.
- [41] Madhavan, V. Artline. <https://github.com/vijishmadhavan/ArtLine>; 2021. [Online; accessed 15-March-2023].
- [42] Radford, A, Kim, JW, Hallacy, C, Ramesh, A, Goh, G, Agarwal, S, et al. Learning transferable visual models from natural language supervision. In: *International conference on machine learning*. 2021, p. 8748–8763.
- [43] Dalal, N, Triggs, B. Histograms of oriented gradients for human detection. In: *CVPR*. 2005, p. 886–893.
- [44] Tan, M, Le, Q. Efficientnet: Rethinking model scaling for convolutional neural networks. In: *International conference on machine learning*. 2019, p. 6105–6114.
- [45] Liu, Z, Mao, H, Wu, CY, Feichtenhofer, C, Darrell, T, Xie, S. A convnet for the 2020s. In: *Conference on Computer Vision and Pattern Recognition*. 2022, p. 11976–11986.
- [46] Kingma, DP, Ba, J. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980* 2014;.
- [47] Nealen, A, Sorkine, O, Alexa, M, Cohen-Or, D. A sketch-based interface for detail-preserving mesh editing. In: *ACM SIGGRAPH*. 2005, p. 1142–1147.
- [48] Muzahid, A, Wan, W, Sohel, F, Wu, L, Hou, L. Curvenet: Curvature-based multitask learning deep networks for 3D object recognition. *IEEE/CAA Journal of Automatica Sinica* 2020;8(6):1177–1187.