

TriCoLo: Trimodal Contrastive Loss for Text to Shape Retrieval

Yue Ruan^{1*} Han-Hung Lee^{1*} Yiming Zhang¹ Ke Zhang¹ Angel X. Chang^{1,2}
¹Simon Fraser University ²Alberta Machine Intelligence Institute (Amii)
 {yuer, hla300, yza440, ke_zhang_4, angelx}@sfu.ca
<https://3dlg-hcvc.github.io/tricolo/>

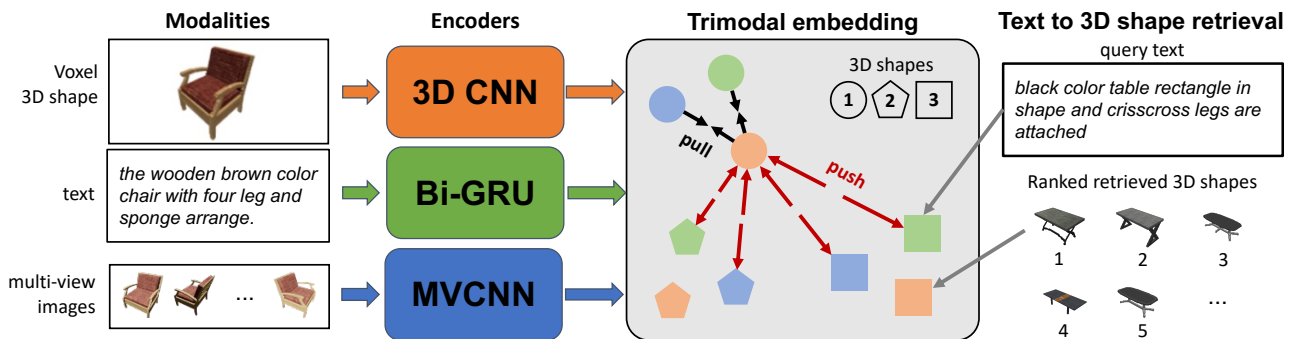


Figure 1. We introduce **TriCoLo**, a **trimodal contrastive loss** for text to 3D shape retrieval. We take objects represented by 3D colored voxels, text descriptions, and multi-view images and jointly use these three modalities to train a trimodal embedding space. This trimodal embedding allows us to perform fine-grained text to shape retrieval.

Abstract

Text-to-shape retrieval is an increasingly relevant problem with the growth of 3D shape data. Recent work on contrastive losses for learning joint embeddings over multimodal data [45] has been successful at tasks such as retrieval and classification. Thus far, work on joint representation learning for 3D shapes and text has focused on improving embeddings through modeling of complex attention between representations [53], or multi-task learning [25]. We propose a trimodal learning scheme over text, multi-view images and 3D shape voxels, and show that with large batch contrastive learning we achieve good performance on text-to-shape retrieval without complex attention mechanisms or losses. Our experiments serve as a foundation for follow-up work on building trimodal embeddings for text-image-shape.

1. Introduction

There has been a dramatic increase in the availability of 3D content in recent years. Improved scanning hardware and reconstruction algorithms are democratizing 3D content creation. The growth in virtual and augmented reality

applications has also driven demand for more synthetic (i.e. human-designed) 3D content. It is no wonder that operating systems now natively support viewing and editing 3D content (e.g., iOS/macOS and Windows). In addition to curated 3D object datasets for research [5, 14, 20, 46, 60], large repositories of 3D shapes provide both synthetic [49, 55, 56] and scanned objects [16, 43].

As 3D assets become more pervasive, we need techniques that allow users to easily and rapidly search through large 3D collections. In recent years, text to image search has seen renewed interest due to improved architectures [9, 34, 37, 45] and objectives [17, 32, 45, 63] for joint representation learning. In contrast, there has been little research on text-driven 3D content search.

Early work by Min et al. [39] compared the text query with text associated with the shape (essentially text-text retrieval). Chen et al. [7] were the first to jointly embed text and 3D shapes for text-to-shape retrieval. They learned the embedding space using triplet loss combined with learning by association [24]. Leveraging the ‘chairs and tables’ dataset [7], followup work investigated improved methods for text-to-shape retrieval [25, 53].

So far, prior work on text-to-shape retrieval has not provided a systematic investigation of: 1) whether 3D information is necessary for text-to-shape retrieval (or whether

single view images to represent a shape are sufficient); 2) whether there are benefits to incorporating information across three modalities; and 3) what contrastive learning setup and loss should be used for constructing joint text-shape embeddings. In our work, we present a systematic study of what is important for improved text-to-shape retrieval. We conduct experiments to examine the effect of input representation (single-view vs multi-view vs 3D voxels), loss function, batch size, and resolution. We show that recent contrastive learning algorithms [63] are sufficient to achieve good performance while avoiding more complex mechanisms, such as combining metric learning with learning by association [7] or training using part-based segmentation of the 3D shapes [53].

In addition, we propose a joint embedding that leverages the multiple modalities offered by 3D data. Specifically, we learn a joint embedding in a trimodal setting: voxel, multi-view images and text. Prior work on text-to-shape retrieval either learns a joint representation with voxels and text, or multi-view images and text, both of which are bimodal settings. We use all three modalities to learn the joint embedding space in an end-to-end fashion and show that trimodal works better than bimodal embedding for text-to-shape retrieval. In summary, our contributions are:

- We introduce a trimodal training scheme with contrastive loss that jointly embeds multi-view images, voxels, and language. We show the trimodal embedding is effective for text to 3D shape retrieval.
- We release ShapeNet c13, a dataset of paired shapes and captions for 13 object categories from ShapeNet [5].
- We present extensive experiments and analysis to provide guidelines on effective settings for applying contrastive loss for text-to-shape retrieval.
- We establish a high-performing baseline for text-to-shape retrieval. Our simple but effective approach outperforms more complex techniques from prior work. Since we introduced TriCoLo in 2022, several follow-up works [53, 57] have used it as a comparison baseline.

2. Related work

There has been growing interest in connecting language to 3D representations for several tasks: identifying 3D objects in scenes [2, 6, 28, 47, 62, 64], describing 3D objects [10, 26], using 3D augmentation in caption-driven image retrieval [59], generating [7] and disambiguating [1, 54] 3D shapes using natural language.

3D shape retrieval. Min et al. [39] was one of the first to address text to 3D shape retrieval by comparing the text query with textual information associated with the shape. Their approach was based purely on text, and relied on each shape having an associated description. Chen et al. [7] was the first work to create a joint embedding of text and 3D

shapes and use that for text-to-shape retrieval. The joint embedding was constructed using a CNN encoder on voxels and GRU encoders on text, with a combined triplet loss [51] and learning by association [24] to align the embedded representations. To improve retrieval, Han et al. [25] used a GRU to encode image features from multiple views to represent the shape, and use reconstruction losses (both intra and inter modalities) in addition to triplet loss and classification loss to train the joint embedding. In contrast, we use multi-view and voxel representation for the shape and do not rely on reconstruction losses. Tang et al. [53] incorporated part-level information, and used point cloud representations for the shapes. In their work, semantic part data was used to compute attention with words to model 3D part relationship with the descriptions. However, obtaining semantic part information can be difficult. Following our work, Wang et al. [57] shows improved performance for text-to-shape retrieval by better selection of positive and negative pairs for contrastive learning, even with just bimodal embeddings of text and multi-image views.

3D object disambiguation through language. The task of object disambiguation through language (also known as a reference game) is related to our text-to-shape retrieval. The main difference between the two tasks is a matter of scale. In shape retrieval, we retrieve all objects that match a textual query from a large set of candidate objects. In contrast, in 3D object disambiguation, there is a small set of objects (typically just two or three) from which we select the one that best matches the description. Reference games involving images and language have a long history [13, 18, 21, 30, 40], but there is significantly less work that takes advantage of the 3D nature of objects. Achlioptas et al. [1] used a speaker-listener model for selecting the correct object based on the text description from among three objects. They showed that combining 3D features (from point clouds) with 2D features (from images) is better than just using 3D or 2D features. More recently, Thomason et al. [54] showed that using multi-view images can improve the disambiguation power of a model. Unlike this line of prior work, we focus on text-to-3D shape retrieval and examine the benefit of combining multi-view images and colored 3D voxel representations. Note that the text-to-shape retrieval problem has different characteristics and challenges as it requires selecting from a large number of instances at inference time (vs disambiguating two or three).

Joint embedding. Joint embedding spaces for text and images [17, 19, 32, 45, 58, 63] have enabled retrieval and generation between text and 2D images. Most joint embedding approaches use contrastive learning, focusing on one modality such as images [8, 23, 48, 50], or two modalities [45]. Recently, an increasing number of works explore combinations of more modalities [3, 4, 36, 38], with prior work (typically in combining vision, audio, and lan-

guage) showing that multiple modalities can improve performance [3, 4, 38]. Liu et al. [36] introduce a data augmentation technique where modalities are disturbed to generate negative samples. These lines of prior work are orthogonal to our work as we investigate the use of trimodal contrastive loss on creating a joint embedding with 3D shape, language, and multiview images for text-to-shape retrieval. Since our work was introduced in 2022, other works started to investigate trimodal embeddings of text-image-shape [35, 61, 65], showing it is useful to train 3D encoders to align 3D embeddings against frozen CLIP text and vision embeddings. ULIP [61] and OpenShape [35] demonstrated that such aligned point cloud embeddings are useful by quantitatively evaluating on classification, and Zhao et al. [65] showed that using an aligned space results in more faithful text-to-shape generation. These works did not focus on evaluating how well the aligned space works for more fine-grained text-to-shape retrieval or comparing different encoders.

3. Problem statement

We tackle the problem of 3D shape retrieval given an input query sentence x_t . We use the Text2Shape [7] dataset which contains tables and chairs from ShapeNet [5] paired with several text descriptions for each object. The text descriptions provide fine-grained information about the appearance of the objects such as color, texture, shape, and whether the object has a certain part (e.g. armrest or a circular base). Accurate retrieval requires learning a good similarity measure between text description and 3D shape. To this end, we learn a shared latent space to facilitate the process of text-shape alignment.

4. Approach

Inspired by recent developments in multimodal contrastive learning [3, 4, 36, 38], we use 3D voxels, multi-view images and language to learn a shared embedding using contrastive learning. Fig. 2 shows how we encode different modalities with per-modality architectures. Embeddings of the same object are pulled closer, while those of different objects are pushed apart using contrastive loss.

Encoder models. We represent the input 3D voxels, text description and multi-view images as x_v, x_t and x_i respectively. For each modality $m \in (v, i, t)$, an encoder f_m takes the input x_m and outputs an encoding $u_m \in \mathbb{R}^d$. The text encoder f_t is a Bi-directional Gate Recurrent Unit (Bi-GRU) [12] which takes a text description $x_t \in \mathbb{R}^{L \times e_t}$ and outputs the embedding $u_t \in \mathbb{R}^d$, where L and e_t are the sentence and word embedding lengths respectively. For voxels, a 3D CNN model f_v takes a 3D input of $x_v \in \mathbb{R}^{r_v \times r_v \times r_v \times 4}$ and outputs $u_v \in \mathbb{R}^d$ where r_v is the voxel resolution. Finally, the image encoder takes M views of the object $x_i \in \mathbb{R}^{M \times r_i \times r_i \times 3}$ through an MVCNN [52] architecture

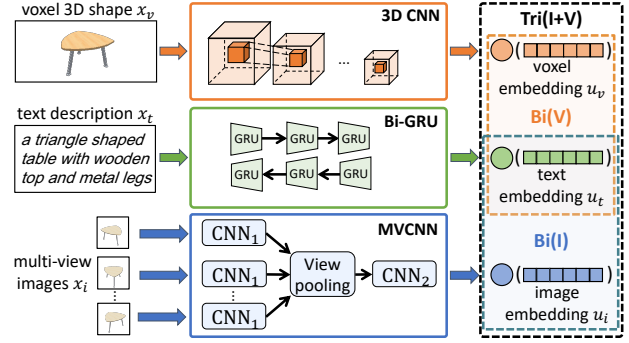


Figure 2. Given the voxel shapes x_v , input text description x_t and rendered images x_i , 3D CNN, Bi-GRU and MVCNN transform them to feature vectors u_v, u_t and u_i . We then minimize a bidirectional contrastive loss to learn effective shape, text, and image representations that are close to each other if they are from the same object.

with pretrained ResNet18 [27] backbone f_i to obtain the representation $u_i \in \mathbb{R}^d$ where r_i is the image resolution.

Loss function. We adopt the bimodal loss from ConVIRT [63]. Specifically for two modalities $m_1, m_2 \in (v, i, t)$ so that $m_1 \neq m_2$ and a batch size of N we construct N positive pairs (u_{m_1j}, u_{m_2j}) for embeddings of the same object and $N^2 - N$ negative pairs $(u_{m_1j}, u_{m_2k})_{j \neq k}$ for different objects. We then apply the symmetric NT-Xent contrastive loss from ConVIRT[63] and popularized by CLIP [45]:

$$l_j^{v \rightarrow t} = -\log \frac{\exp(\langle u_{v_j}, u_{t_j} \rangle / \tau)}{\sum_{k=1}^N \exp(\langle u_{v_j}, u_{t_k} \rangle / \tau)}, \quad (1)$$

$$l_j^{t \rightarrow v} = -\log \frac{\exp(\langle u_{t_j}, u_{v_j} \rangle / \tau)}{\sum_{k=1}^N \exp(\langle u_{t_j}, u_{v_k} \rangle / \tau)} \quad (2)$$

where $\tau \in \mathbb{R}^+$ is a temperature parameter that controls the concentration of the distribution and smoothness of softmax, and $\langle \cdot, \cdot \rangle$ is the cosine similarity. Finally we calculate a weighted sum of $l_j^{v \rightarrow t}$ and $l_j^{t \rightarrow v}$ and average over the mini-batch: $L(v, t) = \frac{1}{N} \sum_{j=1}^N (\alpha l_j^{v \rightarrow t} + (1 - \alpha) l_j^{t \rightarrow v})$, where $\alpha \in [0, 1]$ *Trimodal loss*: To extend the loss to three modalities we simply calculate the contrastive loss over all pair possibilities for the text, voxel and image representations. This gives the final loss: $L_{\text{tri}} = L(v, i) + L(v, t) + L(i, t)$.

Retrieval. For the retrieval task, we are given an input text description and we have to return the matching object. To do this, we leverage the shared embedding space we have built between three modalities. We consider three strategies for matching the text and shape by calculating the cosine similarity between: 1) text and voxel embeddings, 2) text and image embeddings, and 3) text and sum of image and voxel embeddings.

5. Experiments

In the main paper, we present experiments on text-to-shape retrieval on the ‘chair and tables’ dataset from Text2Shape [7]. The ‘chairs and tables’ dataset consists of solid colored voxels of 6521 chairs and 8378 tables from ShapeNet [5], and text descriptions collected from humans (≈ 5 per shape). We follow the train/val/test split by Chen et al. [7] which ensures that shapes do not occur in the same split. We present additional results on shape-to-text retrieval, and results on the primitives dataset in the supplement. To show our method works beyond ‘chairs and tables’, we conduct text-to-shape retrieval on an extended set of 13 object categories from ShapeNet [5] (see supplement). Results across experiments consistently show that our trimodal model outperforms bimodal models.

5.1. Metrics

We follow prior work on text-to-shape retrieval [7, 25, 53] and use standard metrics of Recall Rate (RR@k) and Normalized Discounted Cumulative Gain (NDCG) [29]. RR@k deems a retrieval successful if the ground truth (GT) appears in the top k candidates (we set k to 1 and 5). NDCG compares the ranked retrieval results with optimal ranking. However, since we assume there is only one relevant shape for each query, we do not take full advantage of the NDCG metric. We also evaluate using Mean Reciprocal Rank (MRR), the average of the inverse of the rank of the GT.

We note that there are often multiple shapes that can match the text description. Since the text description can be underspecified, we also measure the similarity of the top k retrieved shapes to the GT shape. Following work in shape retrieval [33], we use a point-wise $F1^\tau$ with $\tau = 0.1$ to calculate shape similarity. $F1^\tau$ is the harmonic mean of the fraction of points from retrieved shapes within τ of a point from GT (point-wise precision), and the fraction of points from GT within τ of a point from retrieved shapes (point-wise recall). To compute $F1^\tau$, we sample 10K points uniformly on the mesh surface of GT and retrieved shapes. See the supplement for more details.

5.2. Implementation details

We use a one-layer bi-directional GRU [12] for the text encoder, and a 3D CNN architecture for the voxel encoder. We use the pretokenized and lemmatized text from Chen et al. [7], with a vocabulary consisting of 3587 unique words and 1 pad token. For the Bi-GRU, we use word embedding size of 256, and a hidden state size of 128. Word embeddings are randomly initialized from a standard Normal distribution. For the 3D CNN, we use 5 Conv3D layers of sparse convolutions from the `spconv` library.¹ For multi-

¹<https://github.com/traveller59/spconv>

view images we use the MVCNN [52] architecture with pre-trained ResNet18 [27] backbone. A fully-connected layer is added to ensure the output dimension for all encoders is 512. Unless otherwise specified, training uses batch size 128, voxel resolution 64^3 , image resolution 128^2 and 6 images for the MVCNN. In preprocessing, we normalize image and voxel values from 0-255 to 0-1. We implement our models using PyTorch [42] and train with the Adam optimizer [31]. We use a learning rate of 0.00035 for batch size 128, and adopt the linear scaling rule [22] to scale the learning rate for other batch sizes. We train for up to 20 epochs until convergence, and select the checkpoint with the best performance on the val set. All models are trained on an A40 GPU with each experiment taking about 1 hour. Our models are memory efficient with most models requiring less than 12GB (see supplement). We render the multiview images with Pyrender² from 12 camera positions elevated slightly above the object, pointing towards the object, and separated by 30 degrees. For multiview experiments using fewer images, we subsample so images are evenly spaced.

5.3. Models

Baselines. We compare to Text2Shape [7], Y2Seq2Seq [25] and Parts2Words [53]. Text2Shape [7] uses a triplet loss [51] combined with learning by association [24]. Y2Seq2Seq [25] uses a view-based model and a triplet constraint. For Parts2Words, we report results for a global model (no part modeling) and their full model that uses part information. Parts2Words [53] uses point clouds as input instead of voxels. The global model uses PointNet [44] as the feature encoder and a Bi-GRU as the text encoder and aggregates the point and text features. Parts2Words jointly embeds point clouds and text by aligning parts from shapes and words from sentences. Both the global and the part-based models use a semi-hard negative mining triplet ranking loss. We note that these methods use either more complex losses or require additional labelled data compared to our model. In addition to baselines from prior work, we use two random baselines: one computes the expected metric mathematically, and the other uses our architecture with random weights.

Our models. We train variants of our model with just two modalities (Bi) or all three modalities (Tri). For the bimodal models, we only consider text and image (**I**), or text and voxels (**V**). During retrieval, we compute the similarity of the text with image (**I**), or text with voxels (**V**). In the case of trimodal embedding, we also use a combination of the image and voxel when computing the similarity, with (**I+V**) denoting that the retrieval was done by calculating similarity with text and the sum of image and voxel representations.

²<https://github.com/mmatl/pyrender>
















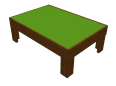

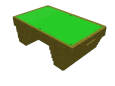


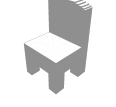









		top1	top2	top3	top4	top5
1	an L-shaped dark brown colored wooden table.	 17.31	 6.60	 20.34	 GT	 2.19
2	a luxurious gray leather modern concept plush chair with stainless steel frame feet	 GT	 2.89	 9.50	 1.78	 3.42
3	simple circular table with no leg and only one circular base.	 0.79	 5.77	 0.59	 GT	 6.32
4	This is greenish top wooden billiards table.	 15.79	 GT	 8.04	 19.59	 3.18
5	this is a boxy look gray chair. It appears to be made out of granite and is gray with 4 short legs and a high, arched back.	 GT	 4.64	 22.32	 12.97	 11.42
6	wooden armless dining room chair with open nine-square back.	 GT	 19.36	 13.31	 18.78	 12.38

Figure 3. Retrieval results on the *test* set with Tri(I+V). For each description, we use Tri(I+V) to retrieve the top-5 shapes. We show the $F1^{0.1}$ score (as a percentage) for each retrieved shape and mark the ground-truth shape (indicated by green GT). The expected F1 score for GT is 100. Shapes that are not a perfect match to the description are marked in dark orange (color mismatch), and gold (shape detail mismatch). Results show that our network has good language grounding ability overall. It can retrieve shapes that match *L-shaped* (row 1), *stainless steel frame feet* (row 2), *circular table* (row 3), *no leg* (row 3), *circular base* (row 3), *greenish top* (row 4), *wooden* (row 4), *boxy look* (row 5), *gray* (row 5), *armless* (row 6) and *nine-square back* (row 6). Though we assume one ground-truth shape, multiple shapes can match the query description (row 1, 3-5).

5.4. Results

We present qualitative retrieval examples and quantitative evaluations comparing our method to prior work. We examine the choice of different loss functions and hyperparameters (see supplement for experiments on image and voxel resolution, and experiments with different backbones). We train models with different seeds and report mean and standard error across 3 runs.

Example retrievals Fig. 3 shows successful retrievals of shapes using Tri(I+V), our best-performing model (see supplement for more examples). Our model successfully grounds language describing shape (*L-shaped*, *boxy*), color (*brown*, *greenish*), and texture (*wooden*). It can also handle negation (*armless*). Note that many shapes match the de-

scription despite not being the ground-truth shape, indicating that there are indeed many matching shapes for a given description. For example, in row 3 the retrieved shapes all match the description, but the four of the five shapes would be negatives in our training and retrieval metrics.

Comparison with prior work. We report the text-to-shape retrieval results in Tab. 1. The full Parts2Words [53] model assumes prior part segmentation knowledge to compute attention with the word embeddings and trains using the triplet loss with negative sampling. In contrast, we do not leverage any part prior knowledge, or attention mechanisms. For comparison, we include Parts2Words (global) with average pooling, which does not use any part information. Tab. 1 shows that our method performs better on all retrieval metrics, and we can achieve slightly better per-

	RR@1	RR@5	NDCG@5
Random (expected)	0.06	0.30	0.20
Random (weights)	0.08	0.32	0.20
Text2Shape [7]	0.40	2.37	1.35
Y2Seq2Seq [25]	2.93	9.23	6.05
Parts2Words (global) [53]	8.60	24.82	16.83
Parts2Words (full) [53]	12.72	32.98	23.13
Bi(I) (our)	11.09	30.78	21.10
Bi(V) (our)	8.93	26.71	17.98
Tri(I+V) (our)	12.22	32.23	22.46

Table 1. Text to shape retrieval comparison against prior work on the *test* set. We report the recall rate (RR@1, RR@5) and NDCG@5 as percentages. We train with a batch size of 128, 64^3 voxels, and 6 multi-view images at a resolution of 128^2 each. Our bimodal joint embedding (Bi(I), Bi(V)) trained using the NT-Xent loss outperforms prior work with global matching. Our trimodal embedding (Tri(I+V)) further improves performance, and is close to Parts2Words [53] which uses part annotations for training.

	RR@1(↑)	RR@5(↑)	NDCG@5(↑)	MRR(↑)	$F1^{0.1}$ (↑)
Bi(I)	11.61 ± 0.20	30.65 ± 0.19	21.36 ± 0.23	21.46 ± 0.25	16.69 ± 0.50
Tri(I)	12.19 ± 0.45	32.33 ± 0.60	22.54 ± 0.54	22.62 ± 0.49	17.39 ± 0.41
Bi(V)	9.59 ± 0.27	27.14 ± 0.48	18.54 ± 0.13	19.03 ± 0.08	14.96 ± 0.20
Tri(V)	9.83 ± 0.21	27.75 ± 0.35	18.97 ± 0.21	19.32 ± 0.20	15.08 ± 0.23
Tri(I+V)	12.52 ± 0.28	32.67 ± 0.61	22.87 ± 0.46	22.68 ± 0.32	17.45 ± 0.30

Table 2. Comparison of bimodal and trimodal models for text-to-shape retrieval on the *val* set. Trimodal embeddings (Tri(I),Tri(V)) give better performance than bimodal embeddings (Bi(I),Bi(V)). By summing the image and voxel embeddings from the trimodal model (Tri(I+V)), we further improve retrieval performance.

formance than even the full Parts2Words model. Note that there are several differences in the prior work compared to our own: the network architectures and specifics of the loss functions, as well as different input representations. Chen et al. [7] used 32^3 colored voxels, while Y2Seq2Seq [25] used multi-view images, and Parts2Words [53] used colored point clouds.

Bimodal vs Trimodal. We compare the trimodal joint embedding with bimodal ones (see Tab. 2). The modalities in the parentheses indicate which representation was used to retrieve the 3D shapes with respect to the text embeddings. We see that the trimodal embedding improves retrieval performance across all metrics when retrieving by both images and voxels. We obtain the best result when we sum the image and voxel embeddings, indicating that the information in the voxels is complementary to the multi-view images. Fig. 4 shows a comparison of retrieved shapes for an example description from the validation set. The retrieved shapes using Tri(I+V) conform to the description more closely than Bi(I) or Bi(V) shapes. See supplement for more examples.

Can pretrained CLIP encoders help? We investigate whether using pretrained vision-language encoders can help improve performance. Specifically, we experiment with

	Text	Image	Voxels	RR@1	RR@5	NDCG@5
ZS*	CLIP	CLIP	-	5.43	16.57	11.27
Bi(I)	CLIP	MVCNN	-	5.79	16.90	11.49
Bi(I)	GRU	CLIP	-	7.29	22.57	14.85
Bi(I)	CLIP	CLIP	-	5.76	19.13	12.41
Tri(I+V)	CLIP	CLIP	3D-CNN	6.72	21.95	14.52
Bi(I)	GRU	MVCNN	-	11.43	30.07	20.92
Bi(V)	GRU	-	3D-CNN	8.98	26.76	17.99
Tri(I+V)	GRU	MVCNN	3D-CNN	12.11	32.39	22.42

Table 3. Comparison of text to shape retrieval performance using CLIP-based models on the *val* set. We report the recall rate (RR@1, RR@5) and NDCG@5 as percentages. It can be seen that zero-shot CLIP [45] has relatively good performance considering that it has not been trained on the Text2Shape [7] ‘chairs and tables’ dataset. Training an MLP to project the CLIP embeddings (CLIP-MLP) drastically improves the retrieval performance, but still underperforms our Tri(I+V) model.

	RR@1(↑)	RR@5(↑)	NDCG@5(↑)	MRR(↑)
Bi(I)	6.44 ± 0.36	21.6 ± 0.61	14.08 ± 0.51	14.99 ± 0.44
Bi(V)	6.19 ± 0.14	20.85 ± 1.02	13.58 ± 0.46	14.51 ± 0.30
Tri(I+V)	8.12 ± 0.20	26.39 ± 0.58	17.96 ± 0.55	18.55 ± 0.42

Table 4. Text-to-shape retrieval on the *val* set using triplet loss with semi-hard negative mining. Performance is lower than NT-Xent (Tab. 2).

CLIP [45], a popular text-image embedding that was trained using contrastive learning with the NT-Xent loss on a large corpus of 400M image-text pairs. We took the pretrained CLIP with ViT-L/14 [15] backbone, and aggregated the embeddings from 6 multi-view images.

We compared the performance of using CLIP in a zero-shot (ZS) manner (without training any weights) and using the CLIP image and text encoders in our models. For zero-shot retrieval, we use the average of the multi-view image embeddings as our overall image embedding, and match that against the text embedding by taking the dot product. To incorporate CLIP into our models, we project frozen CLIP embeddings using a two-layer MLP (see supplement). We present results for zero-shot CLIP and variations of using the text or image CLIP encoders in Table 3. We find that zero-shot CLIP does not perform well, likely due to the domain gap between the rendered images and the CLIP training data. Nevertheless, it can beat the baseline method from the original Text2Shape [7] without being trained on the dataset. Incorporating CLIP into our models and training the MLP results in higher performance, showing the value of training a task-specific MLP on the Text2Shape [7] ‘chairs and tables’ data. We find that our models, which are trained from scratch, are able to outperform the CLIP variants. We also experiment on the ShapeNet c13 data (see supplement), where the CLIP-based models perform much better. We believe that is likely due to less training data for the other categories.

NT-Xent vs triplet loss. To validate the choice of NT-Xent

	# of images	RR@1(↑)	RR@5(↑)	NDCG@5(↑)	MRR(↑)
Bi(I)	1	9.15 ± 0.11	26.34 ± 0.32	17.84 ± 0.20	18.36 ± 0.19
	3	11.19 ± 0.19	30.22 ± 0.25	20.97 ± 0.05	21.23 ± 0.12
	6	11.61 ± 0.20	30.65 ± 0.19	21.36 ± 0.23	21.46 ± 0.25
	12	11.23 ± 0.20	31.13 ± 0.10	21.43 ± 0.09	21.50 ± 0.15

Table 5. Comparison of number of images on shape retrieval for Bi(I) on the *val* set. We find that having multiple views is important for improved performance, but increasing the number of images beyond 6 causes a slight decrease in RR@1. We believe that 6 views is sufficient to capture the necessary information, and increasing it further increases the number of parameters and requires more compute.

	batch size	RR@1(↑)	RR@5(↑)	NDCG@5(↑)	MRR(↑)
Tri(I+V)	32	10.62 ± 0.27	30.19 ± 0.61	20.60 ± 0.16	20.86 ± 0.11
	64	11.48 ± 0.34	31.40 ± 0.55	21.67 ± 0.37	21.80 ± 0.32
	128	12.52 ± 0.28	32.67 ± 0.61	22.87 ± 0.46	22.68 ± 0.32
	256	12.43 ± 0.35	32.25 ± 0.58	22.53 ± 0.49	22.65 ± 0.40

Table 6. Comparison of batch-size on shape retrieval for Tri(I+V) on the *val* set. We find that increasing the batch size increases the performance. However, the performance decreased for the largest batch size we tried (256). This could be due to overfitting on the limited amount of negatives, or the presence of more noisy negatives in the large batch.

as our loss function, we compare the performance of our model using a hinge-based triplet loss [48] instead of NT-Xent. We use semi-hard negative mining with a margin of 0.025. Semi-hard negatives have been shown to improve performance for contrastive losses [8]. Specifically Tang et al. [53] showed it worked better than either triplet-loss by itself or hard negatives for retrieval with the Text2Shape dataset. Tab. 4 shows that the retrieval performance with triplet loss is significantly lower than with NT-Xent. Overall, our findings are consistent with prior work [11]. Note that our model outperforms Y2Seq2Seq [25] even with just triplet loss. We find that with NT-Xent loss, our bimodal models surpass the performance of Parts2Words [53].

5.5. Hyper-parameter analysis

We compare the performance of bimodal models on the validation set with different numbers of input images and batch sizes. We use the bimodal models as they are faster to train and require less memory than the trimodal model.

Do we need multi-view images? For Bi(I), we experiment with number of images ranging from 1 to 12 and find that performance increases as we increase the number of images to 6, after which there are diminishing returns and even a small drop in performance (see Tab. 5). The results indicate that multi-view images provide a benefit over a single view.

Does larger batch size always help? We also compare batch sizes of 32, 64, 128 for Tri(I+V) and find that performance increases with increasing batch size from 32 to 128 (see Tab. 6). This is consistent with findings from prior work on contrastive learning [8, 41]. However, the perfor-

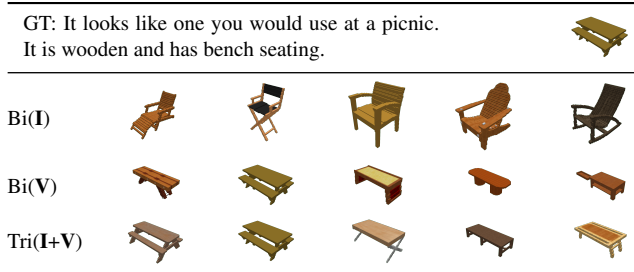


Figure 4. Top 5 retrieved shapes from the *val* set using Bi(I), Bi(V), and Tri(I+V). We see that Bi(I) understands abstract concepts such as *picnic* poorly. Tri(I+V) retrieves the most shapes consistent with the description.

mance drops when the batch size increases to 256 for Bi(I). For Bi(V), increasing the batch size to 256 makes little difference. We hypothesize this is due to more false negatives in the batch since the text description may apply to multiple shapes. Another reason may be that since our dataset size is small compared to image datasets used in prior work [8, 45], having a big batch size might overfit our model. We also note that variance is quite high between runs, which we again attribute to false negatives in the batch and randomness introduced when sampling batches. However, more investigation is warranted.

Impact of other parameters. We also conduct additional experiments (see supplement) examining the effect of different resolutions for image and voxels (higher resolution is better), sparse convolutions (similar performance with less memory, but more time to train), image backbone (smaller backbone works better), zero-shot performance of CLIP (works but not as good as model trained on the data) as well as CLIP embeddings projected using a trained MLP (similar performance as Bi(I) on ‘chairs and tables’ but slightly better for ShapeNet c13).

5.6. Error analysis

Manual analysis. We conduct a manual analysis of the top 5 results returned for 50 text queries from the validation set for Bi(I), Bi(V), and Tri(I+V). We count the number of query results (shapes) that match the description exactly, and categorize the error into color mismatch, large shape mismatch, shape detail mismatch, and missing part (see supplement for examples and Tab. 7 for analysis summary). As expected from the quantitative results, Tri(I+V) has the most shapes that match the text. With the limited number of queries we examined, all models have similar performance on color and missing parts. The Bi(I) model had difficulty getting small shape details correct, and Tri(I+V) obtained the best performance on matching the overall shape.

Failure cases. Fig. 5 shows failure cases for the Tri(I+V) model. For the top row, while our model did not retrieve the ground-truth (GT) shape in the top 5, all top 5 retrieved













	GT	top1	top2	top3	top4	top5
round surface with interconnected leg		 0.62	 25.19	 5.06	 0.73	 22.28
taupe one seater sofa . it has a light brown wooden frame with four leg support . it has two <u>wide</u> arm rest		 5.53	 7.09	 8.88	 6.19	 0.47

Figure 5. Failed retrievals on the *val* set with Tri(I+V) (unmatched text is underlined). The ground truth (GT) is shown in the first column, followed by the retrieved results with the $F1^{0.1}$ score for each. We see that some descriptions are general and our retrieved shapes match the description but is not GT shape (first row), and retrieval of shapes with part details (*wide arm rest*) is hard (second row).

	match	color mismatch	big shape error	small shape error	missing part
Bi(I)	106	65	22	85	5
Bi(V)	103	67	26	76	5
Tri(I+V)	113	64	17	74	5

Table 7. Manual analysis of the top 5 results returned for 50 text queries. We group the results into whether they perfectly match the description, or whether there is a mismatch in color or shape. We confirm that Tri(I+V) has the best overall performance with the most perfect matches and the least number of shape mismatches.

shapes match the input description (28/100 samples that we manually inspected belonged to this error case). This illustrates that the evaluation data and protocol should be improved beyond what was established in Text2Shape [7]. Despite this, our manual analysis indicates the metrics do a good job of ranking the models. The second row shows the challenge of retrieving shapes with fine detail. While our model can retrieve *taupe* colored armchairs, the retrieved objects did not have *wide armrest*. Our model is good at the overall shape and color, but can miss fine details and infrequent terms such as *foldable*. More data and modeling part-to-text correspondence (as in Parts2Words [53]) can help to reduce these types of failures.

5.7. Limitations

We investigated a trimodal loss for text-to-shape retrieval and found that with careful tuning we outperformed the SoTA as of early 2022, when this work was initially performed. We restricted our study to voxel-based 3D representations which often do not capture geometric details and fine-grained surface textures. It would be interesting to consider other modalities such as point clouds, depth images, and textured 3D polygonal meshes which may alleviate these limitations. One big challenge of incorporating additional modalities is the memory cost. In addition, we focused on a specific type of contrastive loss. Other

contrastive losses, data augmentation, as well as other loss terms such as captioning loss and reconstruction loss are promising directions for further improvement. Our dataset is limited in the style of the text and the coverage of shapes. The evaluation also assumes that there is only one correct shape but as we have noted, multiple shapes can match a description. Thus, a significant challenge is to handle false negative pairs in a mini-batch due to the descriptions being ambiguous. These limitations suggest opportunities for future work. We believe our work can serve as a good foundation for follow-up work in text-to-shape retrieval.

6. Conclusion

We carry out a systematic study of contrastive losses for text-to-shape retrieval. We show that using simple contrastive losses can achieve comparable results to text-to-shape retrieval methods relying on extra annotation and complex losses. Our experiments show that incorporating 3D information either via voxels or multi-view images is helpful for the task. We identify important challenges to solve for the development of useful text-to-retrieval models. In addition, we propose a trimodal contrastive loss which further improves performance by considering both 2D and 3D representations. We hope our systematic study will serve as a foundation encouraging more work on text-to-shape retrieval, which is an increasingly important task as there are more and more 3D data repositories.

Acknowledgements. This work is funded by the Canada CIFAR AI Chair program, an NSERC Discovery Grant, and a TUM-IAS Hans Fischer Fellowship (Focus Group Visual Computing). This research was enabled in part by support provided by [WestGrid](#) and [Compute Canada](#). We thank Dave Zhenyu Chen for collecting the text descriptions for ShapeNet c13. We also thank the anonymous reviewers for their feedback, and Manolis Savva for proofreading and editing suggestions.

References

- [1] Panos Achlioptas, Judy Fan, Robert Hawkins, Noah Goodman, and Leonidas J Guibas. ShapeGlot: Learning language for shape differentiation. In *Proc. of International Conference on Computer Vision (ICCV)*, 2019. 2
- [2] Panos Achlioptas, Ahmed Abdelreheem, Fei Xia, Mohamed Elhoseiny, and Leonidas Guibas. ReferIt3D: Neural listeners for fine-grained 3D object identification in real-world scenes. In *Proc. of European Conference on Computer Vision (ECCV)*, 2020. 2
- [3] Hassan Akbari, Linagzhe Yuan, Rui Qian, Wei-Hong Chuang, Shih-Fu Chang, Yin Cui, and Boqing Gong. VATT: Transformers for multimodal self-supervised learning from raw video, audio and text. *arXiv preprint arXiv:2104.11178*, 2021. 2, 3
- [4] Jean-Baptiste Alayrac, Adria Recasens, Rosalia Schneider, Relja Arandjelovic, Jason Ramapuram, Jeffrey De Fauw, Lucas Smaira, Sander Dieleman, and Andrew Zisserman. Self-supervised multimodal versatile networks. *Advances in neural information processing systems*, 2(6):7, 2020. 2, 3
- [5] Angel X. Chang, Thomas Funkhouser, Leonidas Guibas, Pat Hanrahan, Qixing Huang, Zimo Li, Silvio Savarese, Manolis Savva, Shuran Song, Hao Su, Jianxiong Xiao, Li Yi, and Fisher Yu. ShapeNet: An information-rich 3D model repository. *arXiv preprint arXiv:1512.03012*, 2015. 1, 2, 3, 4
- [6] Dave Zhenyu Chen, Angel X Chang, and Matthias Nießner. ScanRefer: 3D object localization in RGB-D scans using natural language. In *Proc. of European Conference on Computer Vision (ECCV)*, 2020. 2
- [7] Kevin Chen, Christopher B Choy, Manolis Savva, Angel X Chang, Thomas Funkhouser, and Silvio Savarese. Text2shape: Generating shapes from natural language by learning joint embeddings. In *Proc. of Asian Conference on Computer Vision (ACCV)*, 2018. 1, 2, 3, 4, 6, 8
- [8] Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework for contrastive learning of visual representations. In *International Conference on Machine Learning (ICML)*, 2020. 2, 7
- [9] Yen-Chun Chen, Linjie Li, Licheng Yu, Ahmed El Kholy, Faisal Ahmed, Zhe Gan, Yu Cheng, and Jingjing Liu. Uniter: Learning universal image-text representations. In *Proc. of European Conference on Computer Vision (ECCV)*, 2020. 1
- [10] Zhenyu Chen, Ali Gholami, Matthias Nießner, and Angel X Chang. Scan2Cap: Context-aware dense captioning in RGB-D scans. In *Proc. of Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021. 2
- [11] Zhiqin Chen and Hao Zhang. Learning implicit fields for generative shape modeling. In *Proc. of Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018. 7
- [12] Kyunghyun Cho, Bart van Merriënboer, Çağlar Gülçehre, Fethi Bougares, Holger Schwenk, and Yoshua Bengio. Learning phrase representations using RNN encoder-decoder for statistical machine translation. In *Proceedings of the conference on empirical methods in natural language processing (EMNLP)*, 2014. 3, 4
- [13] Herbert H Clark and Deanna Wilkes-Gibbs. Referring as a collaborative process. *Cognition*, 22(1):1–39, 1986. 2
- [14] Jasmine Collins, Shubham Goel, Achleshwar Luthra, Leon Xu, Kenan Deng, Xi Zhang, Tomas F Yago Vicente, Himanshu Arora, Thomas Dideriksen, Matthieu Guillaumin, and Jitendra Malik. ABO: Dataset and benchmarks for real-world 3D object understanding. *arXiv preprint arXiv:2110.06199*, 2021. 1
- [15] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. In *Proc. of International Conference on Learning Representations (ICLR)*, 2020. 6
- [16] Laura Downs, Anthony Francis, Nate Koenig, Brandon Kinman, Ryan Hickman, Krista Reymann, Thomas B McHugh, and Vincent Vanhoucke. Google scanned objects: A high-quality dataset of 3D scanned household items. *arXiv preprint arXiv:2204.11918*, 2022. 1
- [17] Fartash Faghri, David J Fleet, Jamie Ryan Kiros, and Sanja Fidler. VSE++: Improving visual-semantic embeddings with hard negatives. *arXiv preprint arXiv:1707.05612*, 2017. 1, 2
- [18] Michael C Frank and Noah D Goodman. Predicting pragmatic reasoning in language games. *Science*, 336(6084):998–998, 2012. 2
- [19] Andrea Frome, Greg Corrado, Jonathon Shlens, Samy Bengio, Jeffrey Dean, Marc’Aurelio Ranzato, and Tomas Mikolov. Devise: A deep visual-semantic embedding model. In *Advances in neural information processing systems*, 2013. 2
- [20] Huan Fu, Rongfei Jia, Lin Gao, Mingming Gong, Binqiang Zhao, Steve Maybank, and Dacheng Tao. 3D-FUTURE: 3D Furniture shape with TextURE. *arXiv preprint arXiv:2009.09633*, 2020. 1
- [21] Dave Golland, Percy Liang, and Dan Klein. A game-theoretic approach to generating spatial descriptions. In *Proceedings of the 2010 conference on empirical methods in natural language processing*, pages 410–419, 2010. 2
- [22] Priya Goyal, Piotr Dollár, Ross B. Girshick, Pieter Noordhuis, Lukasz Wesolowski, Aapo Kyrola, Andrew Tulloch, Yangqing Jia, and Kaiming He. Accurate, large minibatch SGD: Training ImageNet in 1 hour. *ArXiv*, abs/1706.02677, 2017. 4
- [23] Raia Hadsell, Sumit Chopra, and Yann LeCun. Dimensionality reduction by learning an invariant mapping. In *Proc. of Conference on Computer Vision and Pattern Recognition (CVPR)*, volume 2, pages 1735–1742. IEEE, 2006. 2
- [24] Philip Haeusser, Alexander Mordvintsev, and Daniel Cremers. Learning by association—a versatile semi-supervised training method for neural networks. In *Proc. of Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 89–98, 2017. 1, 2, 4
- [25] Zhizhong Han, Mingyang Shang, Xiyang Wang, Yu-Shen Liu, and Matthias Zwicker. Y2Seq2Seq: Cross-modal representation learning for 3D shape and text by joint reconstruction and prediction of view and word sequences. In *Proceedings of the AAAI Conference on Artificial Intelligence*, 2019. 1, 2, 4, 6, 7
- [26] Zhizhong Han, Chao Chen, Yu-Shen Liu, and Matthias

- Zwicker. ShapeCaptioner: Generative caption network for 3D shapes by learning a mapping from parts detected in multiple views to sentences. In *Proceedings of the 28th ACM International Conference on Multimedia*, 2020. 2
- [27] Kaiming He, X. Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proc. of Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 770–778, 2016. 3, 4
- [28] Pin-Hao Huang, Han-Hung Lee, Hwann-Tzong Chen, and Tyng-Luh Liu. Text-guided graph neural networks for referring 3D instance segmentation. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 35, pages 1610–1618, 2021. 2
- [29] Kalervo Järvelin and Jaana Kekäläinen. Cumulated gain-based evaluation of IR techniques. *ACM Trans. Inf. Syst.*, 20(4):422–446, 2002. 4
- [30] Sahar Kazemzadeh, Vicente Ordonez, Mark Matten, and Tamara Berg. Referitgame: Referring to objects in photographs of natural scenes. In *Proceedings of the conference on empirical methods in natural language processing (EMNLP)*, pages 787–798, 2014. 2
- [31] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014. 4
- [32] Ryan Kiros, Ruslan Salakhutdinov, and Richard S Zemel. Unifying visual-semantic embeddings with multimodal neural language models. *arXiv preprint arXiv:1411.2539*, 2014. 1, 2
- [33] Weicheng Kuo, Anelia Angelova, Tsung-Yi Lin, and Angela Dai. Mask2CAD: 3D shape prediction by learning to segment and retrieve. In *Proc. of European Conference on Computer Vision (ECCV)*, pages 260–277. Springer, 2020. 4
- [34] Xiujun Li, Xi Yin, Chunyuan Li, Pengchuan Zhang, Xiaowei Hu, Lei Zhang, Lijuan Wang, Houdong Hu, Li Dong, Furu Wei, et al. Oscar: Object-semantic aligned pre-training for vision-language tasks. In *Proc. of European Conference on Computer Vision (ECCV)*, pages 121–137. Springer, 2020. 1
- [35] Minghua Liu, Ruoxi Shi, Kaiming Kuang, Yin hao Zhu, Xuanlin Li, Shizhong Han, Hong Cai, Fatih Porikli, and Hao Su. OpenShape: Scaling up 3D shape representation towards open-world understanding. *arXiv preprint arXiv:2305.10764*, 2023. 3
- [36] Yunze Liu, Qingnan Fan, Shanghang Zhang, Hao Dong, Thomas Funkhouser, and Li Yi. Contrastive multimodal fusion with TupleInfoNCE. In *Proc. of International Conference on Computer Vision (ICCV)*, pages 754–763, 2021. 2, 3
- [37] Jiasen Lu, Dhruv Batra, Devi Parikh, and Stefan Lee. VILBERT: Pretraining task-agnostic visiolinguistic representations for vision-and-language tasks. In *Advances in neural information processing systems*, 2019. 1
- [38] Sijie Mai, Ying Zeng, Shuangjia Zheng, and Haifeng Hu. Hybrid contrastive learning of tri-modal representation for multimodal sentiment analysis. *ArXiv*, abs/2109.01797, 2021. 2, 3
- [39] Patrick Min, Michael Kazhdan, and Thomas Funkhouser. A comparison of text and shape matching for retrieval of online 3D models. In *International Conference on Theory and Practice of Digital Libraries*, pages 209–220. Springer, 2004. 1, 2
- [40] Will Monroe, Robert XD Hawkins, Noah D Goodman, and Christopher Potts. Colors in context: A pragmatic neural model for grounded language understanding. *Transactions of the Association for Computational Linguistics*, 5:325–338, 2017. 2
- [41] Aaron van den Oord, Yazhe Li, and Oriol Vinyals. Representation learning with contrastive predictive coding. *arXiv preprint arXiv:1807.03748*, 2018. 7
- [42] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Kopf, Edward Yang, Zachary DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. Pytorch: An imperative style, high-performance deep learning library. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d’Alché-Buc, E. Fox, and R. Garnett, editors, *Advances in Neural Information Processing Systems 32*, pages 8024–8035. Curran Associates, Inc., 2019. 4
- [43] Polycam. Polycam. <https://poly.cam/explore>, 2022. Accessed: 2022-03-07. 1
- [44] Charles R Qi, Hao Su, Kaichun Mo, and Leonidas J Guibas. PointNet: Deep learning on point sets for 3D classification and segmentation. In *Proc. of Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 652–660, 2017. 4
- [45] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever. Learning transferable visual models from natural language supervision. In *International Conference on Machine Learning (ICML)*, 2021. 1, 2, 3, 6, 7
- [46] Jeremy Reizenstein, Roman Shapovalov, Philipp Henzler, Luca Sbordone, Patrick Labatut, and David Novotny. Common objects in 3D: Large-scale learning and evaluation of real-life 3D category reconstruction. In *Proc. of International Conference on Computer Vision (ICCV)*, pages 10901–10911, 2021. 1
- [47] Junha Roh, Karthik Desingh, Ali Farhadi, and Dieter Fox. LanguageRefer: Spatial-language model for 3D visual grounding. In *Proc. of Conference on Robot Learning (CoRL)*, 2021. 2
- [48] Florian Schroff, Dmitry Kalenichenko, and James Philbin. Facenet: A unified embedding for face recognition and clustering. In *Proc. of Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 815–823, 2015. 2, 7
- [49] SketchFab. Sketchfab. <https://sketchfab.com/features/free-3d-models>, 2021. Accessed: 2021-10-30. 1
- [50] Kihyuk Sohn. Improved deep metric learning with multi-class n-pair loss objective. In *Advances in neural information processing systems*, 2016. 2
- [51] Hyun Oh Song, Yu Xiang, Stefanie Jegelka, and Silvio Savarese. Deep metric learning via lifted structured feature embedding. In *Proc. of Conference on Computer Vision and*

- Pattern Recognition (CVPR)*, pages 4004–4012, 2016. 2, 4
- [52] Hang Su, Subhransu Maji, Evangelos Kalogerakis, and Erik G. Learned-Miller. Multi-view convolutional neural networks for 3D shape recognition. In *Proc. of International Conference on Computer Vision (ICCV)*, pages 945–953, 2015. 3, 4
- [53] Chuan Tang, Xi Yang, Bojian Wu, Zhizhong Han, and Yi Chang. Parts2words: Learning joint embedding of point clouds and texts by bidirectional matching between parts and words. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6884–6893, 2023. 1, 2, 4, 5, 6, 7, 8
- [54] Jesse Thomason, Mohit Shridhar, Yonatan Bisk, Chris Paxton, and Luke Zettlemoyer. Language grounding with 3D objects. In *Proc. of Conference on Robot Learning (CoRL)*, 2021. 2
- [55] Trimble. Sketchup 3D warehouse. <https://3dwarehouse.sketchup.com>, 2021. Accessed: 2021-10-30. 1
- [56] TurboSquid. Turbosquid. <https://www.turbosquid.com/Search/3D-Models>, 2021. Accessed: 2021-10-30. 1
- [57] Ye Wang, Bowei Jiang, Changqing Zou, and Rui Ma. MXM-CLR: A unified framework for contrastive learning of multifold cross-modal representations. *arXiv preprint arXiv:2303.10839*, 2023. 2
- [58] Jason Weston, Samy Bengio, and Nicolas Usunier. Wsabie: Scaling up to large vocabulary image annotation. In *Twenty-Second International Joint Conference on Artificial Intelligence*, 2011. 2
- [59] Xiaoshi Wu, Hadar Averbuch-Elor, Jin Sun, and Noah Snavely. Towers of babel: Combining images, language, and 3D geometry for learning multimodal vision. In *Proc. of International Conference on Computer Vision (ICCV)*, 2021. 2
- [60] Zhirong Wu, Shuran Song, Aditya Khosla, Fisher Yu, Linguang Zhang, Xiaoou Tang, and Jianxiong Xiao. 3D shapenets: A deep representation for volumetric shapes. In *Proc. of Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 1912–1920, 2015. 1
- [61] Le Xue, Mingfei Gao, Chen Xing, Roberto Martín-Martín, Jiajun Wu, Caiming Xiong, Ran Xu, Juan Carlos Nieves, and Silvio Savarese. ULIP: Learning a unified representation of language, images, and point clouds for 3D understanding. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1179–1189, 2023. 3
- [62] Zhihao Yuan, Xu Yan, Yinghong Liao, Ruimao Zhang, Sheng Wang, Zhen Li, and Shuguang Cui. InstanceRefer: Cooperative holistic understanding for visual grounding on point clouds through instance multi-level contextual referring. In *Proc. of International Conference on Computer Vision (ICCV)*, pages 1791–1800, 2021. 2
- [63] Yuhao Zhang, Hang Jiang, Yasuhide Miura, Christopher D Manning, and Curtis Langlotz. Contrastive learning of medical visual representations from paired images and text. *arXiv preprint arXiv:2010.00747*, 2021. 1, 2, 3
- [64] Lichen Zhao, Daigang Cai, Lu Sheng, and Dong Xu. 3DVG-transformer: Relation modeling for visual grounding on point clouds. In *Proc. of International Conference on Computer Vision (ICCV)*, pages 2928–2937, 2021. 2
- [65] Zibo Zhao, Wen Liu, Xin Chen, Xianfang Zeng, Rui Wang, Pei Cheng, Bin Fu, Tao Chen, Gang Yu, and Shenghua Gao. Michelangelo: Conditional 3D shape generation based on shape-image-text aligned latent representation. *arXiv preprint arXiv:2306.17115*, 2023. 3