# Learning Visual and Motion Aware 3D Model Representations for Semantic Retrieval

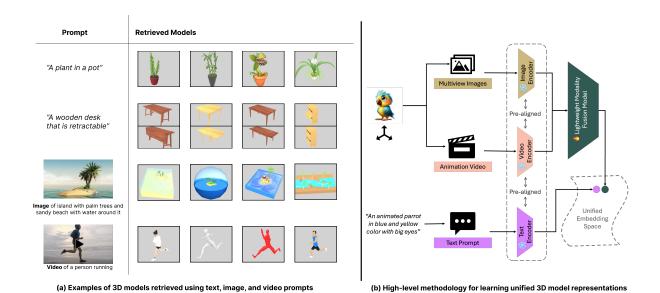


Figure 1: **3D model retrieval examples and methodology.** (a) Examples of **3D models retrieved for different input prompts.** Prompts can be textual (e.g., descriptions of visual or dynamic features such as a "retractable table"), or based on sample images and videos that capture visual appearance or motion. (b) High-level methodology for learning **3D model representations.** Each **3D model** is represented using multiview images to capture appearance and videos to capture dynamics. A trainable multimodal fusion model combines the outputs of pretrained encoders into a unified embedding space, enabling cross-modal retrieval.

## **A**BSTRACT

VR content creators today face lengthy trial-and-error cycles when searching for 3D models, which relies heavily on manually annotated metadata. However, this approach struggles to capture the rich visual and dynamic details inherent to 3D models. Prior work has attempted to address this by rendering models as sets of images, but has largely remained restricted to category-level search and static geometry. In VR environments where asset behavior is often as important as appearance, we present a visual and motion aware 3D retrieval method that represents each model through multiview renders and short animations, extending search beyond static geometry to include dynamic behaviors critical for VR applications. Pretrained image-text and video-text encoders extract features that are combined into a unified embedding by training a lightweight multimodal fusion model. This enables 3D model retrieval from text, images, and videos without requiring abundant and detailed ground truth data, supporting more nuanced and semantically meaningful search. On text queries referencing both appearance and motion, our method achieves a 23-percentage-point improvement in top-5 retrieval accuracy on the Objaverse dataset [4], while also attaining competitive accuracy on the ModelNet40 dataset [35] and maintaining robustness as the model library size scales. A user study further shows that participants strongly preferred our retrievals compared to the one used by Sketchfab, one of today's most widely used 3D

repositories. To support future work, we release multi-view renders, animation clips, and 500 manually annotated visual and motion aware captions for a subset of Objaverse.

**Index Terms:** Multimodal Machine Learning, Virtual Reality

#### 1 Introduction

The rapid growth of virtual reality (VR) and immersive applications has created a strong demand for effective search and retrieval of 3D models. Current solutions like online repositories (e.g. Sketchfab [30]) or asset libraries in 3D authoring tools (like Unity[32] and Unreal[5]) largely rely on metadata-based search, where results are organized by manually assigned tags or broad category labels. While effective for coarse filtering, these approaches fall short in supporting natural language queries and in capturing the semantics of how 3D models look and move. As a result, creators often face long trial-and-error cycles when sourcing assets for interactive experiences. In addition, with the growing use of Large Language Models (LLMs) in authoring 3D scenes, being able to semantically retrieve 3D models is essential.

Recent advances in Vision Language Models (VLMs) have transformed retrieval in 2D domains, enabling robust alignment between text, images, and video. However, their potential for 3D content remains underexplored. Prior work primarily focuses on static shape descriptors or learned embeddings of geometry, with limited consideration of motion. This neglects a critical aspect of 3D models in VR, where many models are designed to be animated, and their dynamic behavior is central to their meaning (e.g., a "running dog" vs. a "sleeping dog").

Unlike images and videos on the Internet, there is a massive scarcity of richly annotated 3D models, particularly those with animations or motion information. Public repositories typically provide only minimal textual descriptions or user-assigned tags, which rarely capture fine-grained visual or motion semantics such as "a person wearing a black baseball cap walking", or "a three ducks of different sizes swimming across a pond". This lack of detailed captions makes it difficult to train dedicated 3D-text encoders at scale, in stark contrast to the abundance of data available for image-text or video-text learning. The problem is further amplified in VR contexts, where asset behavior is often as important as asset appearance. To overcome this, we adopt a proxy strategy: rather than learning directly from scarce 3D annotations, we generate 2D renderings and animated sequences of each model, enabling the direct reuse of pretrained image and video retrieval models.

In our work, we present a visual and motion aware 3D model representation learning and retrieval strategy that treats assets as multimodal entities composed of both visual features and dynamics. Each model is represented through a combination of static multiview images and short rendered animation clips, which are aligned with natural language queries using state of the art pretrained image—text and video—text retrieval models. Surprisingly, we find that lightweight fusion strategies, such as simple score averaging across modalities, consistently outperform more complex joint embedding schemes. These results highlight motion as an essential retrieval dimension while also showing that simple integration methods with lower resource requirements can be highly effective.

We validate our approach through both benchmarks and user studies. Quantitatively, we compare our retrieval accuracy against prior state-of-the-art methods and ablated variants of our model. Our method achieves a 23-percentage-point improvement in top-5 retrieval accuracy (87% vs. 63.8% for the next best prior work) on the Objaverse database [4] for prompts involving both visual and motion cues, and competitive performance (85%) on the ModelNet40 database[35]. The gains are most pronounced on larger datasets, highlighting that our approach learns more discriminative features and remains robust as the corpus scales. Qualitatively, we conduct a user study comparing our system to Sketchfab, the industry standard for 3D model search. Participants consistently preferred our results, rating our retrievals higher (3.88 vs. 3.07 out of 5) and scoring our animations substantially better (3.65 vs. 2.05) for alignment between prompts and retrieved models. These findings demonstrate the practical benefits of motion-aware search in real-world VR workflows. Finally, to support future research, we release a dataset of rendered images and video clips for Objaverse, together with 500 manually annotated motion-aware captions, providing a compact yet valuable benchmark for the community.

Our key contributions are:

- A 3D representation learning strategy that leverages multiview images and short animation clips, enabling the reuse of powerful pretrained image and video language models.
- Identification of motion as an essential retrieval dimension, showing that dynamic behavior is critical for aligning language with 3D model semantics.
- A benchmark and dataset for text-to-3D retrieval, including 26k multi-view images and animation clips, as well as 500 detailed motion-aware captions for Objaverse.
- Quantitative benchmarks and user studies showing the effectiveness of multimodal fusion and the value of semantic retrieval over Sketchfab.

## 2 RELATED WORK

## 2.1 3D Representation Learning

Learning effective 3D representations has been a long-standing goal, with approaches spanning geometry-based, view-based, and multimodal paradigms. Early 3D encoders learn directly from geometry, capturing invariances to permutation, viewpoint, and rigid transforms. Representative models include PointNet/PointNet++ [25, 26], DGCNN [34], and multi-view CNNs such as MVCNN [31] and MVTN [9]. Subsequent work explored self-supervised and transformer-based pre-training [36, 40, 23], improving robustness but remaining limited to category-level semantics and underrepresenting fine-grained or open-vocabulary attributes, primarily due to the lack of richly annotated training data.

The rise of large-scale pretraining has shifted focus toward generalizable semantic encoders aligned across modalities. Image-text models like CLIP [28] and ALIGN [13], video-text encoders such as VideoCLIP [37] and InternVideo [33], and multimodal spaces like LanguageBind [45] and ImageBind [8] enable zero-shot comparison across images, videos, and text. Although training analogous encoders directly on 3D models is not yet feasible due to limited large-scale 3D-text data, these image- and video-text models provide transferable priors that we leverage as building blocks.

To transfer these priors to 3D, view-based methods render objects and aggregate 2D features. PointCLIP [42], CLIP2Point [11], and CLIP-Goes-3D [10] pool CLIP features across multiviews, while captioning pipelines such as CAP3D [19] leverage renderings to produce natural-language supervision. More recently, language-anchored pretraining approaches, like ULIP [38] and OpenShape [17], jointly align 3D, image, and text encoders to support open-vocabulary retrieval and zero-shot tasks [43, 15, 16].

Despite these advances, most methods operate on static geometry or image sets, limiting their ability to distinguish models with similar appearance but different dynamics. Moreover, reliance on synthetic captions may introduce noise through hallucinated attributes. Our work targets this gap by combining pre-aligned image and video encoders with pooled multiview renderings and short turntable clips to learn motion-aware 3D embeddings without requiring explicit 3D captions.

# 2.2 3D Search and Retrieval

3D search systems span a spectrum of approaches, from geometric similarity and structured metadata to context-aware and language-driven retrieval. Early methods focused on shape descriptors and sketch-based search [7, 22, 20, 14], later extended by scene context and relational reasoning [6]. Metadata-driven systems leveraged tags, parts, and annotations, which are mostly manually annotated, for structured queries [44, 21, 1].

Recent work has shifted toward embedding-based retrieval using multimodal encoders aligned across text, image, and 3D modalities [38, 17, 10]. These models support open-vocabulary and zero-shot retrieval, and have been integrated into authoring tools for VR scene generation and layout composition [41, 39, 12]. In immersive settings, such systems allow users to retrieve, place, and manipulate 3D models using natural language or high-level intent, bypassing traditional menu-based interfaces. For example, VRCopilot and Holodeck support compositional scene control through voice or text prompts, enabling iterative scene editing and spatial reasoning in real time. Other frameworks like SceneSuggest and ATISS [29, 24] incorporate learned spatial priors to recommend plausible object placements, while diffusion-based models [18, 27] generate realistic room layouts conditioned on functional context. These tools reflect a broader trend toward multimodal, context-aware VR authoring workflows that unify retrieval, layout, and interaction into a single semantic interface. Together, these techniques reflect the growing convergence of 3D understanding, semantic alignment, and interactive content creation.

<sup>&</sup>lt;sup>1</sup>Data and Code will be released with the camera ready paper

#### 3 OUR APPROACH

Our goal is to enable text-driven retrieval of 3D models by learning multimodal embeddings that capture both *appearance* and *motion*. Unlike images or videos, 3D models rarely come with high-quality natural language descriptions, making it difficult to train a dedicated 3D-text encoder. To address this, we adopt a proxy strategy: we render 3D models as **multiview images** and **short animation clips**, and leverage powerful pretrained encoders that jointly aligns text, image, and video in a shared representation space. In our implementation we use the LanguageBind [45] pretrained encoders, but the method is encoder-agnostic: any model that provides a unified image–video–text space can be swapped in without changing the pipeline. This design supports text, image, and video based retrieval of 3D models, enabling more expressive search (Fig. 1b). This section details how we learn these multimodal 3D representations from the individual modalities (Fig. 2).

## 3.1 Multiview Image Representation

For each 3D model, we generate M = 6 canonical renderings from fixed orientations  $(\pm x, \pm y, \pm z)$ . We chose these views because they give uniform coverage of all principal sides with minimal redundancy, whereas fewer views miss at least one side. Each view  $I_m$  is encoded using the LanguageBind image encoder  $f_I$ :

$$z_m = f_I(I_m), \quad z_m \in \mathbb{R}^d, \quad m = 1, \dots, 6. \tag{1}$$

While each  $z_m$  captures a valid perspective of the model, they contain overlapping information and may emphasize different features. To form a single representation  $z_{\rm img}$ , we train lightweight fusion models that learn to combine the six views. We experiment with:

- MLP fusion (order-invariant): a lightweight fusion model that scores each view with a *shared* per-view MLP, softmaxnormalizes the scores, and forms a symmetric weighted sum before a final projection. This parameter-efficient design is robust on smaller datasets.
- Attention fusion (content-adaptive): a self-attention layer over the six view embeddings that models inter-view interactions and assigns content-aware weights, emphasizing views with unique or discriminative cues; the extra flexibility typically helps on larger datasets at the cost of more compute.

Training objectives. Because high-quality text annotations for 3D models are rare, we train the fusion model without paired captions using two self-consistency losses.

**Cosine Similarity (intra-model).** We encourage the fused representation to stay close to its constituent views:

$$\mathcal{L}_{\cos} = 1 - \max_{m \in \{1,\dots,6\}} \cos(z_{\text{img}}, z_m). \tag{2}$$

**Batch-Wise Contrastive Loss (inter-model).** Let the batch size be B. For each model  $i \in \{1, \dots, B\}$  we sample one positive view index  $m_i^+$  uniformly from  $\{1, \dots, 6\}$  and take the corresponding normalized embedding  $v_i = \frac{z_{i,m_i^+}}{\|z_{i,m_i^+}\|}$ . We also L2-normalize the fused

embedding  $o_i = rac{z_{ ext{img},i}}{\|z_{ ext{img},i}\|}$  . We form the logits matrix

$$\ell_{ij} = \frac{o_i^\top v_j}{\tau}, \quad i, j \in \{1, \dots, B\},\tag{3}$$

and apply cross-entropy with identity labels (diagonal as positives, off-diagonals as negatives):

$$\mathcal{L}_{\text{ctr}} = -\frac{1}{B} \sum_{i=1}^{B} \log \frac{\exp(\ell_{ii})}{\sum_{j=1}^{B} \exp(\ell_{ij})}.$$
 (4)

where  $\tau$  is a temperature hyperparameter.

For each fused embedding, the randomly chosen view of the *same* model is treated as the positive, while the corresponding views from all other models in the batch serve as negatives.

**Total objective.**, where  $\lambda_{cos}$  and  $\lambda_{ctr}$  are the weights representing how much each loss contributes to the total loss.

$$\mathcal{L}_{img} = \lambda_{cos} \mathcal{L}_{cos} + \lambda_{ctr} \mathcal{L}_{ctr}. \tag{5}$$

# 3.2 Video Representation

To capture dynamics, we render a short animation clip V for each 3D model and encode it with the frozen LanguageBind video encoder  $f_V$ :

$$z_{\text{vid}} = f_V(V), \quad z_{\text{vid}} \in \mathbb{R}^d.$$
 (6)

This representation captures temporal semantics (e.g., "walking," "flying," "collapsing") that static images cannot convey. We chose a 6 seconds video clip given our dataset, as the animations were shorted than that time frame, although this length can be tuned given the dataset and the encoder's context window.

#### 3.3 Multimodal Fusion

Finally, we construct a unified multimodal embedding  $z_{3D}$  that integrates both  $z_{\rm img}$  and  $z_{\rm vid}$ . Although we experimented with trainable fusion layers (results in Table 1), similar to the multiview image fusion model, we observed that a simple average provided consistently higher accuracy:

$$z_{3D} = \frac{1}{2} (z_{\text{img}} + z_{\text{vid}}). \tag{7}$$

This result suggests that pretrained encoders are already well aligned across modalities, and that computationally lightweight fusion is both effective and robust. Moreover, because  $z_{3D}$  lies in the same embedding space as the pretrained text encoder, retrieval can be performed directly.

## 3.4 Multimodal Retrieval

At query time, a natural language description T is encoded with the LanguageBind text encoder  $f_T$  to obtain  $z_{\text{text}} = f_T(T)$ . Each model's multimodal embedding  $z_{3D}$  is then compared with  $z_{\text{text}}$  in the shared space using cosine similarity:

$$sim(z_{3D}, z_{text}) = \frac{z_{3D} \cdot z_{text}}{\|z_{3D}\| \|z_{text}\|}.$$
 (8)

Since the LanguageBind framework aligns text, image, and video encoders in a common space, our system naturally supports retrieval from multiple query modalities: text, images, or videos. For scalability, all  $z_{3D}$  embeddings are stored in a vector database, specifically ChromaDB [2], which allows an efficient approximate search for nearest neighbor using cosine distance, and hence can be used to retrieve similar 3D model embeddings. This allows interactive, low-latency retrieval across thousands of models, making the system practical for VR content creation applications.

#### 4 DATASET GENERATION

To train and evaluate our method, we construct a multimodal dataset based on Objaverse [4], supplemented with ModelNet40 [35] for additional baselines.

Objaverse We generate 6 canonical multiview renderings for each 3D model and record 6-second video clips to capture motion and animation dynamics. This results in a dataset of over 26,000 models. While prior work has released static renderings, our contribution is the first large-scale release of rendered *videos* for 3D models. We make all data and rendering scripts (for Blender [3]) publicly available to facilitate reproducibility and future research.

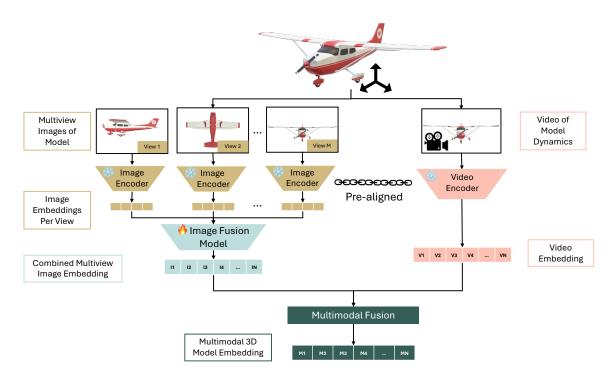


Figure 2: **Multimodal 3D model representation learning for text-to-3D retrieval.** Each 3D model is rendered from M canonical views. A frozen, pretrained image encoder produces per-view embeddings, which a lightweight, trainable image-fusion model aggregates into a single multiview image embedding. In parallel, a frozen video encoder computes a video embedding from a short dynamics clip. The final 3D model representation is obtained by a non-trainable multimodal fusion step implemented as the mean of the multiview image and video embeddings. The resulting 3D embedding is aligned in the same embedding space as the text encoder, enabling text-to-3D retrieval.

Captions We initially used the automatically generated Cap3D captions [19] as a starting point. However, these captions often contain errors as they are generated using large language models and lack descriptions of dynamic behaviors. They also have high variability in terms of length and descriptiveness. As such, we do not use them as our ground truth retrieval. Instead, we manually curate a benchmark set of 500 models with detailed captions describing both appearance and motion. This provides a small but valuable benchmark for motion-aware 3D retrieval.

ModelNet40 For comparison, we also generate multiview images and videos for ModelNet40. Because these models are static, the resulting videos reduce to extended renderings of fixed geometry without meaningful motion. Although they lack dynamics, this provides a useful baseline for evaluating retrieval performance when only appearance cues are available.

Overall, our dataset contributions are: (i) the first large-scale release of rendered videos for Objaverse models, (ii) a curated benchmark of 500 motion-aware captions, and (iii) reproducible Blender pipelines for extending the dataset.

#### 5 EVALUATION

# 5.1 Text-to-3D Retrieval Accuracy

To evaluate how well our 3D representation captures both visual and motion semantics, we measure text-based retrieval accuracy on two benchmarks. On the Objaverse dataset, which consists of publicly available models from Sketchfab [30], we use our manually curated set of 500 detailed captions that describe both appearance and dynamics. On ModelNet40, we use the standard category-level captions, which are more representative of prior baseline setups. We compare against three representative baselines: ULIP[38], OpenShape [17], and CLIP-Goes-3D [10]. We

also compare against metadata-based search, GPT-4o-generated descriptions, and single modality image-only and video-only based retrieval. In addition, we ablate different multiview image fusion strategies, including simple averaging, max pooling (closest image match), and trainable fusion modules (MLP and attention), as described in Section 3. We also experiment with a trainable multimodal fusion model, following the same training strategy as the multiview image fusion, and compare it to a simple mean of the multiview image and video embeddings (our solution).

We evaluate retrieval on a model library of 500 models and report Top-1, Top-5, and Top-10 accuracy for both datasets. We also compute the Mean Reciprocal Rank (MRR), defined as MRR =  $\frac{1}{N}\sum_{i=1}^{N}\frac{1}{\mathrm{rank}_{i}}$ , where  $\mathrm{rank}_{i}$  is the position of the first correct retrieval for query i and N is the number of retrieval queries made, which captures how highly the correct result is ranked on average across all queries. Having an MRR of 1 indicates that the top-1 retrieval is the expected model. For Objaverse, each caption is instance-specific and corresponds to a single target model. In contrast, ModelNet40 uses category-level captions, so multiple models may match a query, and success is defined as retrieving any model of the correct category within the top-k results. The results can be seen in Table 1.

Table 1 shows that our multimodal approach outperforms all baselines on Objaverse, achieving 87.3% Top-5 accuracy and an MRR of 0.806. Among baselines, GPT-40-based text descriptions are the strongest on Objaverse, while OpenShape achieves the best results on ModelNet40, which primarily reflects category-level appearance cues. Our ablations highlight several trends. Multi-view image features consistently improve over single-image retrieval, and attention-based fusion yields the strongest image-only results. Video-only retrieval demonstrates the value of motion, and com-

Table 1: Zero-shot text to 3D retrieval on Objaverse and ModelNet40. Top-k results are reported as percentages. MRR is shown as mean.

Method	Objaverse				ModelNet40			
	Top1	Top5	Top10	MRR	Top1	Top5	Top10	MRR
Baselines								
ULIP	40.5	59.0	66.0	0.483	37.5	55.0	60.0	0.453
OpenShape	36.8	63.8	72.3	0.481	80.0	95.0	95.0	0.868
Clip-Goes-3D	1.8	9.5	14.3	0.052	12.5	25.0	30.0	0.174
Metadata	25.8	41.5	47.8	0.321			_	_
Text Description (GPT-4o)	55.0	74.5	78.8	0.634	60.0	85.0	87.5	0.694
Ours (ablations)								
Single Image	40.0	58.8	63.3	0.481	22.5	35.0	40.0	0.283
Average Across All View Images	65.3	76.8	79.8	0.704	60.0	80.0	82.5	0.700
Maximum Similarity Across All View Images	35.3	60.5	70.5	0.461	65.0	82.5	<u>85.0</u>	0.740
Video Only	62.5	77.8	80.8	0.692	55.0	65.0	70.0	0.585
Multiview Image Fusion using MLP	68.5	82.3	84.8	0.743	<u>70.0</u>	77.5	82.5	0.745
Multiview Image Fusion using Attention	70.5	83.3	85.8	0.760	<u>70.0</u>	80.0	<u>85.0</u>	0.740
Ours: Multimodal Retrieval (Trained Fusion Model)	60.5	83.0	85.7	0.699			_	_
Ours: Multimodal Retrieval (Mean of Multiview Image + Video)	75.5	87.3	89.5	0.806	<u>70.0</u>	<u>85.0</u>	<u>85.0</u>	0.721

bining images with videos further boosts performance. Notably, simple averaging of image and video embeddings surpasses more complex trained fusion, suggesting that pretrained encoders are already well aligned across modalities. Together, these results confirm that motion-aware multimodal embeddings provide the most robust retrieval performance.

On ModelNet40, our approach achieves competitive performance, though OpenShape remains the strongest baseline. This gap is expected, as ModelNet40's 3D models are static meshes without textures or motion, limiting the benefit of our multimodal design. In this setting, the video modality effectively is another rendering of static geometry, providing little additional signal beyond the multiview images. Consequently, our multimodal fusion results are very close to those of the image-only models, consistent with the absence of dynamic cues in the dataset.

# 5.2 Retrieval Robustness to Model Library Size

While our model performs well on smaller libraries (e.g., 500 models), practical model libraries can be much larger. We therefore study how retrieval accuracy scales with library size. To quantify robustness, we measure the proportion of retrieval accuracy retained as the library grows, where higher values indicate greater resilience to scaling. Specifically, let Top-5@N denote the probability that the groundtruth model appears among the top five results retrieved from an N-item library. We report the retained accuracy as

$$\mathrm{retained}(N) = \frac{\mathrm{Top-5@N}}{\mathrm{Top-5@500}} \times 100,$$

where 100% indicates no degradation relative to the 500-model setting. We evaluate the retained retrieval accuracy at library sizes of 500, 1k, 10k, and 26k models.

From Figure 3, we observe that retrieval performance decreases for all models as the library size grows, but the rate of degradation differs significantly. Baseline methods such as OpenShape and ULIP drop sharply, while CLIP-Goes-3D struggles even at moderate scales. In contrast, our model retains 57% of its top-5 accuracy even at 26k candidates, demonstrating substantially greater robustness to scaling. This indicates that our multimodal embeddings capture more discriminative and transferable features, enabling reliable retrieval performance even in large-scale model libraries.

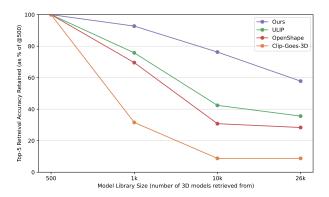


Figure 3: **Top-5 retrieval accuracy retained as the model library scales.** Each line shows the top-5 accuracy retained for our model and baselines as model library scales from 500 to 1k, 10k and 26k.

#### 5.3 Qualitative Retrieval Accuracy User Study

To complement our quantitative evaluation, we conduct a user study comparing our retrieval system with Sketchfab. We recruit 14 participants, eight with prior experience searching for 3D models and six with no previous exposure to 3D modeling.

Each participant is initially asked to use the same four fixed prompts (two text, one image, one video) when searching for 3D models on both the retrieval systems (ours and sketchfab). The participant is then asked to come with 4 more prompts on their own. For each prompt on each system, they inspect the top-5 retrievals and answer three main questions on a 5-point Likert scale (higher is better): (1) overall relevance to the prompt, (2) relevance of visual appearance, and (3) relevance of motion and animation. The motion score is ignored in our data analysis when the prompt does not specify any kind of animation.

The average scores for all prompts are shown in Fig. 4. Our retrieval system outperforms Sketchfab on every criterion: overall relevance (3.88 vs. 3.07), visual match (3.46 vs. 2.87), and motion/animation match (3.65 vs. 2.05). The largest improvement is in motion scores, indicating better alignment with dynamic seman-

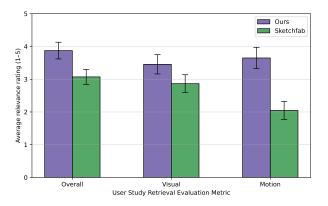


Figure 4: User study (n=14) comparing our retrieval system and Sketchfab. Bars show mean Likert ratings (1–5; higher is better) for Overall, Visual, and Motion relevance of the top-5 retrieved models across all prompts and error bars denote SEM. Our retrieval system scores higher than Sketchfab on all three metrics, with the largest margin on Motion.

tics, though we also observe general improvement in retrieval quality, particularly for detailed prompts. When asked to select their preferred platform for each prompt, participants chose our retrieval system 83.9% of the time for the four fixed prompts that we provided. However, for the prompts the participants created themselves, our system was preferred only 57.1% of the time. We attribute this discrepancy primarily to the limited scope of our model library (26k Objaverse models) compared to Sketchfab's substantially larger, open-ended database. Consequently, specific models that users searched for (e.g., "library with a worm resting on one of the books" or "a hacker typing at a keyboard") were often unavailable in our collection, while Sketchfab's broader catalog increased the likelihood of finding relevant matches. Additionally, since Sketchfab does not natively support image- or video-based queries, participants had to reformulate such prompts as text descriptions. Given these limitations, we report descriptive statistics without claims of statistical significance.

## 5.4 Performance Overhead

On a Linux desktop with an RTX 4090, the **per-model time to add** a **new 3D model (capture**  $\rightarrow$  **embedding) is**  $\approx$  31.4 s in total: image capture (6 views) 1.84 s, video capture (6 s clip) 19.30 s, image embedding (6 images) 8.38 s (median), and video embedding 1.90 s (median). Multimodal fusion (multiview fusion over six images, averaged with video) and adding to the database takes 1.8 ms per model, which is quite small compared to the previous steps.

For retrieval on a 26k-item database, the database lookup itself is fast: top-10 retrieval latency is 1.5 ms (p50) and 52 ms (p95). **Endto-end retrieval (text embedding** + **retrieval) is 2.20** s (p50) and 2.25 s (p95), dominated by getting the text embedding at  $\approx$  2.19 s.

The most memory-consuming phase is image capture (GPU  $\sim 8.38\,\mathrm{GB}$  VRAM; CPU  $\sim 7.52\,\mathrm{GB}$ ), followed by video capture (GPU  $\sim 7.52\,\mathrm{GB}$ ; CPU  $\sim 8.59\,\mathrm{GB}$ ). The embedding generation steps use less GPU memory ( $\sim 3.54$ – $4.08\,\mathrm{GB}$ ) but show higher CPU peaks ( $\sim 9.13\,\mathrm{GB}$ ).

Although processing and indexing new models require the bulk of computation and time, these steps occur only once at ingestion. Once a model has been captured, embedded, and added to the database, subsequent searches incur negligible overhead. This design makes the system practical and scalable as the one-time ingestion cost is amortized over all future queries, while retrieval remains lightweight and responsive even at large corpus sizes. this All figures are rough estimates intended to convey operational cost.

The absolute values vary with software and hardware infrastructure, scene geometry and textures, and corpus scale.

#### 6 Discussion

Our results highlight the promise and current limitations of multimodal 3D retrieval. By jointly leveraging multiview images and rendered motion clips, our approach outperforms existing baselines on Objaverse, especially for prompts that reference dynamic behaviors. The ability to align text with visual appearance and motion cues makes retrieval more expressive and closer to how creators naturally think about models when building immersive environments.

At the same time, our experiments reveal caveats. On Model-Net40, with static, textureless geometry, the benefits of multimodality are diminished, as the video stream provides little additional signal beyond the rendered views. While our results are competitive, this suggests that motion-aware retrieval is most valuable for models designed to be animated, rather than purely geometric datasets. Similarly, while our system demonstrates robustness as the retrieval library scales, accuracy inevitably declines with very large collections, pointing to an open challenge for scalable 3D search.

Another notable finding is that simple averaging of image and video embeddings often outperforms more complex trainable fusion schemes. This suggests that pretrained encoders are already well aligned across modalities and lightweight integration can suffice. However, this may also reflect the limited scale of current training data. With larger and more diverse multimodal datasets, learned fusion strategies could potentially yield greater gains.

Although our current framework uses six canonical views and a single rendered video per model, there is substantial room to explore richer input modalities. Extending to multiview video could provide more comprehensive coverage of both geometry and motion, particularly for models with complex dynamics. Beyond simply adding more views, an important open question is how to identify which views or motion segments are the most informative. This could improve robustness and efficiency by focusing attention on the features most critical for semantic understanding of 3D models.

More broadly, advancing multimodal embeddings of 3D models has implications that extend beyond retrieval. A richer semantic understanding of individual 3D models also lays the foundation for reasoning at the scene level. Since scenes are ultimately composed of multiple interacting models, the ability to capture fine-grained appearance and motion at the model level can extend to semantic scene understanding. This opens the possibility of retrieving not just single models but entire scenes based on high-level descriptions, as well as supporting context-aware composition where models are selected to fit naturally within larger environments. In addition, such embeddings may benefit 3D model generation, serving as a prior that constrains generative models toward more semantically coherent outputs. By bridging retrieval and generation, motion-aware multimodal representations could ultimately enable workflows where users specify high-level descriptions and systems populate scenes with appropriate, contextually consistent models.

## 7 CONCLUSION

This paper presents a method for representing 3D models that captures both visual appearance and motion. Because reliable 3D captions are scarce, we leverage pretrained image—and video—text encoders and train a multimodal fusion model. Compared to baselines, our approach achieves higher retrieval accuracy, scales more robustly, and is preferred by users over Sketchfab for semantic search, indicating practical value. Simple yet effective, it is suitable for large-scale retrieval and will continue to improve as encoders advance, without additional 3D data. This work takes a step toward making natural language a more effective interface for exploring and retrieving 3D content.

#### REFERENCES

- A. X. Chang, T. Funkhouser, L. Guibas, P. Hanrahan, Q. Huang, Z. Li,
  S. Savarese, M. Savva, S. Song, H. Su, J. Xiao, L. Yi, and F. Yu.
  Shapenet: An information-rich 3d model repository, 2015.
- [2] Chroma. Chroma the open-source embedding database, 2023. Accessed: 2025-03-30. 3
- [3] B. O. Community. Blender a 3D modelling and rendering package. Blender Foundation, Stichting Blender Foundation, Amsterdam, 2018. 3
- [4] M. Deitke, D. Schwenk, J. Salvador, L. Weihs, O. Michel, E. Vander-Bilt, L. Schmidt, K. Ehsani, A. Kembhavi, and A. Farhadi. Objaverse: A universe of annotated 3d objects, 2022. 1, 2, 3
- [5] Epic Games. Unreal engine. 1
- [6] M. Fisher and P. Hanrahan. Context-based search for 3d models. In ACM SIGGRAPH Asia 2010 Papers, SIGGRAPH ASIA '10. Association for Computing Machinery, New York, NY, USA, 2010. doi: 10. 1145/1866158.1866204
- [7] T. Funkhouser, P. Min, M. Kazhdan, J. Chen, A. Halderman, D. Dobkin, and D. Jacobs. A search engine for 3d models. 22(1):83–105, Jan. 2003. doi: 10.1145/588272.588279
- [8] R. Girdhar, A. El-Nouby, Z. Liu, M. Singh, K. V. Alwala, A. Joulin, and I. Misra. Imagebind: One embedding space to bind them all, 2023.
- [9] A. Hamdi, S. Giancola, and B. Ghanem. Mvtn: Multi-view transformation network for 3d shape recognition, 2021.
- [10] D. Hegde, J. M. J. Valanarasu, and V. M. Patel. Clip goes 3d: Lever-aging prompt tuning for language grounded 3d recognition, 2023. 2,
- [11] T. Huang, B. Dong, Y. Yang, X. Huang, R. W. H. Lau, W. Ouyang, and W. Zuo. Clip2point: Transfer clip to point cloud classification with image-depth pre-training, 2023. 2
- [12] L. Höllein, A. Cao, A. Owens, J. Johnson, and M. Nießner. Text2room: Extracting textured 3d meshes from 2d text-to-image models, 2023. 2
- [13] C. Jia, Y. Yang, Y. Xia, Y.-T. Chen, Z. Parekh, H. Pham, Q. V. Le, Y. Sung, Z. Li, and T. Duerig. Scaling up visual and vision-language representation learning with noisy text supervision, 2021. 2
- [14] B. Li, Y. Lu, A. Ghumman, B. Strylowski, M. Gutierrez, S. Sadiq, S. Forster, N. Feola, and T. Bugerin. 3d sketch-based 3d model retrieval. In *Proceedings of the 5th ACM on International Conference on Multimedia Retrieval*, ICMR '15, p. 555–558. Association for Computing Machinery, New York, NY, USA, 2015. doi: 10.1145/2671188 .2749349 2
- [15] H. Li, Y. Zhou, T. Tang, J. Song, Y. Zeng, M. Kampffmeyer, H. Xu, and X. Liang. Unigs: Unified language-image-3d pretraining with gaussian splatting, 2025. 2
- [16] D. Liu, X. Huang, Y. Hou, Z. Wang, Z. Yin, Y. Gong, P. Gao, and W. Ouyang. Uni3d-llm: Unifying point cloud perception, generation and editing with large language models, 2024. 2
- [17] M. Liu, R. Shi, K. Kuang, Y. Zhu, X. Li, S. Han, H. Cai, F. Porikli, and H. Su. Openshape: Scaling up 3d shape representation towards open-world understanding, 2023. 2, 4
- [18] W. Liu, Y. Du, T. Hermans, S. Chernova, and C. Paxton. Structdiffusion: Language-guided creation of physically-valid structures using unseen objects, 2023. 2
- [19] T. Luo, C. Rockwell, H. Lee, and J. Johnson. Scalable 3d captioning with pretrained models, 2023. 2, 4
- [20] P. Min, J. Chen, and T. Funkhouser. A 2d sketch interface for a 3d model search engine. In ACM SIGGRAPH 2002 Conference Abstracts and Applications, SIGGRAPH '02, p. 138. Association for Computing Machinery, New York, NY, USA, 2002. doi: 10.1145/1242073. 1242151 2
- [21] K. Mo, S. Zhu, A. X. Chang, L. Yi, S. Tripathi, L. J. Guibas, and H. Su. Partnet: A large-scale benchmark for fine-grained and hierarchical part-level 3d object understanding, 2018. 2
- [22] R. Ohbuchi, T. Minamitani, and T. Takei. Shape-similarity search of 3d models by using enhanced shape functions. In *Proceedings of Theory and Practice of Computer Graphics*, 2003., pp. 97–104, 2003. doi: 10.1109/TPCG.2003.1206936 2

- [23] Y. Pang, W. Wang, F. E. H. Tay, W. Liu, Y. Tian, and L. Yuan. Masked autoencoders for point cloud self-supervised learning, 2022. 2
- [24] D. Paschalidou, A. Kar, M. Shugrina, K. Kreis, A. Geiger, and S. Fidler. Atiss: Autoregressive transformers for indoor scene synthesis, 2021. 2
- [25] C. R. Qi, H. Su, K. Mo, and L. J. Guibas. Pointnet: Deep learning on point sets for 3d classification and segmentation, 2017. 2
- [26] C. R. Qi, L. Yi, H. Su, and L. J. Guibas. Pointnet++: Deep hierarchical feature learning on point sets in a metric space, 2017. 2
- [27] W. Qian, C. Gao, A. Sathya, R. Suzuki, and K. Nakagaki. Shape-it: Exploring text-to-shape-display for generative shape-changing behaviors with llms. In *Proceedings of the 37th Annual ACM Symposium on User Interface Software and Technology*, UIST '24. Association for Computing Machinery, New York, NY, USA, 2024. doi: 10.1145/3654777.3676348
- [28] A. Radford, J. W. Kim, C. Hallacy, A. Ramesh, G. Goh, S. Agarwal, G. Sastry, A. Askell, P. Mishkin, J. Clark, G. Krueger, and I. Sutskever. Learning transferable visual models from natural language supervision. 2021. 2
- [29] M. Savva, A. X. Chang, and M. Agrawala. Scenesuggest: Contextdriven 3d scene design, 2017. 2
- [30] Sketchfab. Sketchfab The best 3D viewer on the web. https://sketchfab.com/, 2025. Accessed: 9 September 2025. 1, 4
- [31] H. Su, S. Maji, E. Kalogerakis, and E. Learned-Miller. Multi-view convolutional neural networks for 3d shape recognition, 2015. 2
- [32] Unity Technologies. Unity, 2023. Game development platform. 1
- [33] Y. Wang, K. Li, Y. Li, Y. He, B. Huang, Z. Zhao, H. Zhang, J. Xu, Y. Liu, Z. Wang, S. Xing, G. Chen, J. Pan, J. Yu, Y. Wang, L. Wang, and Y. Qiao. Internvideo: General video foundation models via generative and discriminative learning, 2022. 2
- [34] Y. Wang, Y. Sun, Z. Liu, S. E. Sarma, M. M. Bronstein, and J. M. Solomon. Dynamic graph cnn for learning on point clouds, 2019.
- [35] Z. Wu, S. Song, A. Khosla, F. Yu, L. Zhang, X. Tang, and J. Xiao. 3d shapenets: A deep representation for volumetric shapes. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 1912–1920, 2015. 1, 2, 3
- [36] S. Xie, J. Gu, D. Guo, C. R. Qi, L. J. Guibas, and O. Litany. Point-contrast: Unsupervised pre-training for 3d point cloud understanding, 2020. 2
- [37] H. Xu, G. Ghosh, P.-Y. Huang, D. Okhonko, A. Aghajanyan, F. Metze, L. Zettlemoyer, and C. Feichtenhofer. Videoclip: Contrastive pretraining for zero-shot video-text understanding, 2021. 2
- [38] L. Xue, M. Gao, C. Xing, R. Martín-Martín, J. Wu, C. Xiong, R. Xu, J. C. Niebles, and S. Savarese. Ulip: Learning a unified representation of language, images, and point clouds for 3d understanding, 2023. 2, 4
- [39] Y. Yang, F.-Y. Sun, L. Weihs, E. VanderBilt, A. Herrasti, W. Han, J. Wu, N. Haber, R. Krishna, L. Liu, C. Callison-Burch, M. Yatskar, A. Kembhavi, and C. Clark. Holodeck: Language guided generation of 3d embodied ai environments, 2024. 2
- [40] X. Yu, L. Tang, Y. Rao, T. Huang, J. Zhou, and J. Lu. Point-bert: Pre-training 3d point cloud transformers with masked point modeling, 2022. 2
- [41] L. Zhang, J. Pan, J. Gettig, S. Oney, and A. Guo. Vrcopilot: Authoring 3d layouts with generative ai models in vr. In *Proceedings of the 37th Annual ACM Symposium on User Interface Software and Technology*, UIST '24. Association for Computing Machinery, New York, NY, USA, 2024. doi: 10.1145/3654777.3676451
- [42] R. Zhang, Z. Guo, W. Zhang, K. Li, X. Miao, B. Cui, Y. Qiao, P. Gao, and H. Li. Pointclip: Point cloud understanding by clip, 2021. 2
- [43] J. Zhou, J. Wang, B. Ma, Y.-S. Liu, T. Huang, and X. Wang. Uni3d: Exploring unified 3d representation at scale, 2023. 2
- [44] Q. Zhou and A. Jacobson. Thingi10k: A dataset of 10,000 3d-printing models, 2016. 2
- [45] B. Zhu, B. Lin, M. Ning, Y. Yan, J. Cui, H. Wang, Y. Pang, W. Jiang, J. Zhang, Z. Li, W. Zhang, Z. Li, W. Liu, and L. Yuan. Languagebind: Extending video-language pretraining to n-modality by language-based semantic alignment, 2024. 2, 3