# Taskverse: A Benchmark Generation Engine for Multi-modal Language Model

**Jieyu Zhang[1], Weikai Huang[1],* Zixian Ma[1],* Oscar Michel[2], Dong He[1],**
**Tanmay Gupta[2], Wei-Chiu Ma[2], Ali Farhadi[1,2], Aniruddha Kembhavi[2], Ranjay Krishna[1,2]**
[1]University of Washington, [2]Allen Institute for Artificial Intelligence
https://www.task-me-anything.org

## Abstract

Benchmarks for large multimodal language models (MLMs) now serve to simultaneously assess the general capabilities of models, rather than evaluating a specific capability. As a result, when a developer seeks to identify which models to use for their application, they are often overwhelmed by the number of benchmarks and remain uncertain about which benchmark results are most reflective of their specific use case. This paper introduces TASKVERSE, a benchmark generation engine that produces benchmarks tailored to different user needs. TASKVERSE maintains an extendable taxonomy of visual assets, including 113K images, 10K videos, 2K 3D object assets, over 365 object categories, 655 attributes, and 335 relationships, and can programmatically generate over 750 million Image/VideoQA questions. Additionally, it algorithmically addresses user queries regarding MLM performance efficiently within a relative computational budget by employing interactive learning query approximation algorithms. With TASKVERSE, we can answer specific and fine-grained user queries like: "Which model is the best VideoQA model for recognizing material within 1000 inference times?"
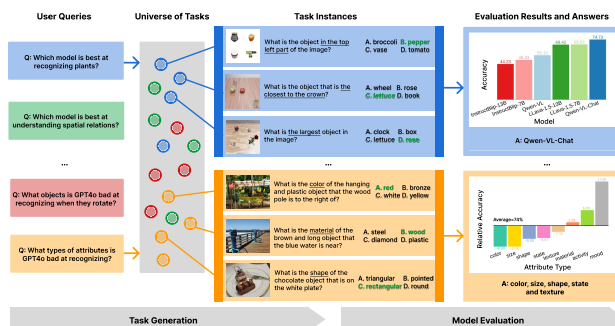
## 1 Introduction



Figure 1: We present examples of user queries, corresponding task instances generated by TASK-VERSE as well as the evaluation results on them that answer the queries.

When a developer wants to identify which models to use for their application, they remain uncertain about which benchmark results are most aligned with their specific use case.

---

* The authors contribute equally to this work.

Although MLMs are released with evaluations on benchmarks like MMBench, MMMU, BLINK, and SeedBench [29, 39, 25, 24, 9], their performance across these holistic benchmarks do not pinpoint which fine-grained capabilities are lacking.

There is a need for a principled benchmark generation process that answers task-specific user queries: "(**Q1**) Which model is the best at recognizing the shape of objects?" or "(**Q2**) what are the model's weaknesses that we can further improve on?".

To address this, we present TASKVERSE, a benchmark generation engine that curates a custom benchmark given a user query (Figure 1).

**Programmatic and scalable task generation.** The task generation process of TASKVERSE includes 2 key components: taxonomy and task generators. Task generators are Python codes that take in taxonomy and generate ImageQA/VideoQA tasks.

First, TASKVERSE maintains an extendable taxonomy with corresponding visual assets (*e.g.* images with scene graphs [23], 3D object assets [7], videos with spatio-temporal annotations [18], rendering softwares [5], *etc.*). It is implemented as an extendable library where new concepts and their corresponding assets and annotations can be easily added.

Second, TASKVERSE contains programmatic task generators that sub-select from the taxonomy to curate a large number of ImageQA/VideoQA questions. TASKVERSE contains 28 different task generators focused on different capabilities of ImageQA/VideoQA models. With our current taxonomy, TASKVERSE can generate over 750 million questions with real or synthetic images/videos.

**Efficient and fine-grained evaluation.** TASKVERSE support supports 4 types of user query: Top-K query, Threshold query, Model comparison query, Model debugging query. For evaluation efficiency, we also implement 3 different ML algorithms to approximate the results of user queries via predicting the model performance across numerous questions. For example, With *Active* algorithm, we can achieve a 70% hit rate in Top-K query with under 10% of model inference times compared with evaluation on all questions.

**What can we do with TASKVERSE.** With TASKVERSE, we answer several critical user queries about MLMs, like what are the best VideoQA models in specific capability (e.g. recognizing, reasoning)? We also released useful resources including TASKVERSE-RANDOM benchmark, TASKVERSE-2024 benchmark, TASKVERSE-DB, TASKVERSE-UI.

## 2 TASKVERSE

### 2.1 TASKVERSE paradigm

Consider a user who wants to know "Which open-sourced MLM is best at recognizing objects even if the object is rotating?". TASKVERSE provides an interface for the user to pose such questions and provides them with an answer (Figure 2). It contains a taxonomy to symbolically represent visual content. A query identifies the relevant portion of the Taxonomy required to answer the query. It also contains task generators that create input-output pairs that test for a specific capability. The Taxonomy subset is used to select the appropriate task generator. We adopt the common input-output format used in existing benchmarks, *i.e.*, all the task instances in TASKVERSE contain an image/video, a question, and multiple options with one ground truth answer. MLMs will be evaluated on these generated task instances and the results will be returned back to the user. Finally, it also supports efficient approximation algorithms to estimate model performance without evaluating all the task space. Unlike most existing procedural data systems, we design TASKVERSE so that the generation space of tasks can be expanded by adding new source data and/or task generator code. More details about taxonomy, task generation, supported user queries, and results approximation algorithms are in Appendix B and C.

### 2.2 Stats of TASKVERSE

TASKVERSE can enable evaluations with up to 113K+ realistic images (from Visual Genome), 9K+ videos (from Charades), and infinite synthetic images and videos.
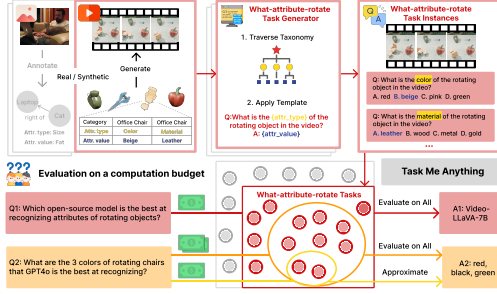
Figure 2: The top part illustrates the task generation process with an example video synthesized with 3D objects and their annotations. The bottom part depicts the model evaluation process, which selects the relevant tasks based on the user's query and their budget and performs either full evaluation or results approximation to answer the query.
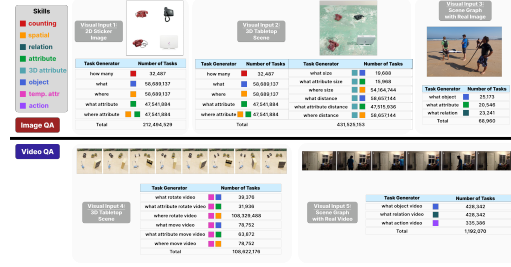


Figure 3: The statistics of generatable tasks of each task generator and example image/video in TASKVERSE. Each task generator uses high-level perceptual skills, and this collection of task generators can collectively generate over 750M VQA tasks.

We include 28 different task generators across 5 types of visual inputs: 2D sticker images (2D), 3D tabletop scene images/videos (3D), and real images/videos with manually-annotated scene graphs. More stats of TASKVERSE are in Appendix D.

## 3 Analysing Video-Language Models with TASKVERSE

**Example user query 1: what is the best video-language model for each specific skill?** From Figure 4, VIDEO-LLAVA-7B and CHAT-UNIVI-7B are relatively well-rounded, positioning in the top 3 models across all skills except for Attribute understanding. On the other hand, while VIDEOCHAT2-7B specializes in object, attribute, and temporal attribute understanding.

**Example user query 2: how do current video-language models perform in recognizing different types of concepts?** From Figure 5, we can see that all models perform poorly at recognizing concepts related to orientation and materials, but perform well with activities, mood, etc. This may indicate a data bias in current model training, highlighting the need for more balanced data curation.

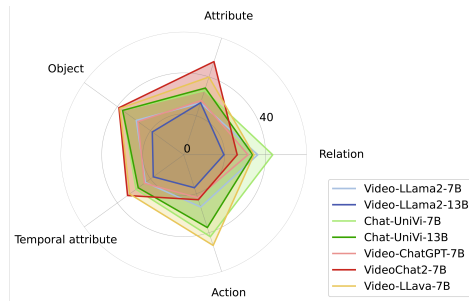More fine-grained user queries about MLMs are in Appendix H.



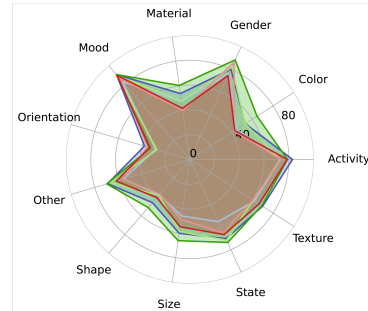Figure 4: **Video-language model, high-level skills**



Figure 5: **Video-language model, fine-grained object and attribute skills**

## 4 Conclusion

In this work, we introduce TASKVERSE, a task generation and evaluation system designed to address user queries with different evaluation objectives.