

# Task Me Anything

Jieyu Zhang, Weikai Huang, Zixian Ma, Oscar Michel, Dong He,  
Tanmay Gupta, Wei-Chiu Ma, Ali Farhadi, Aniruddha Kembhavi,  
Ranjay Krishna



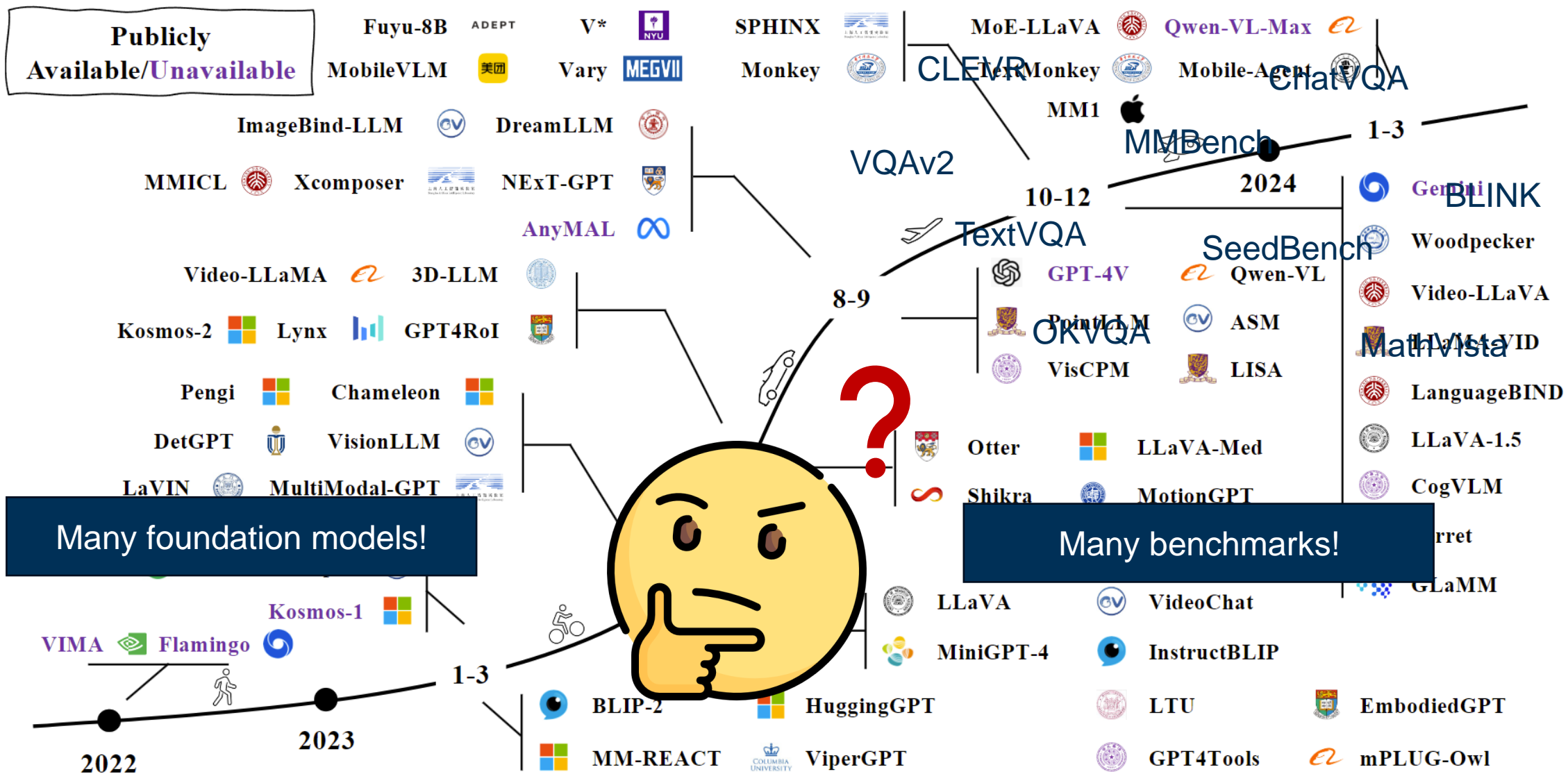
Presented by:

Arthur Nascimento, Manisha Natarajan, Milad Heydari

# Outline

- Problem Statement
- Background and Related Work
- Approach
- Experiments and Results
- Limitations , Societal Implications
- Summary of Strengths, Weaknesses, and relationship to other papers

# Problem Statement



# Introduction

Need for user-centric benchmark generation:

(Q1): Which model is best at recognizing different plants?

(Q2): Which types of attributes is model X (say GPT4o) bad at recognizing?

Task Me Anything uses procedural generation to generate user-centric benchmarks for evaluating multimodal language models (MLMs)

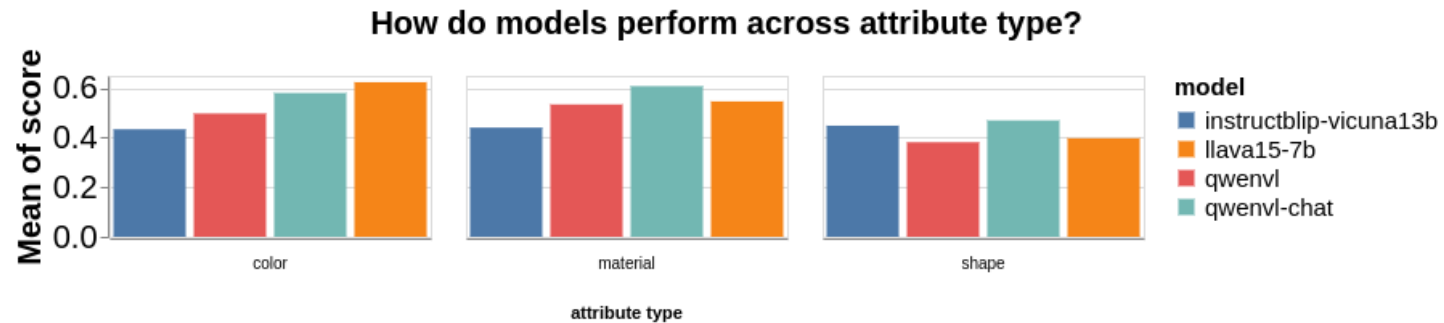
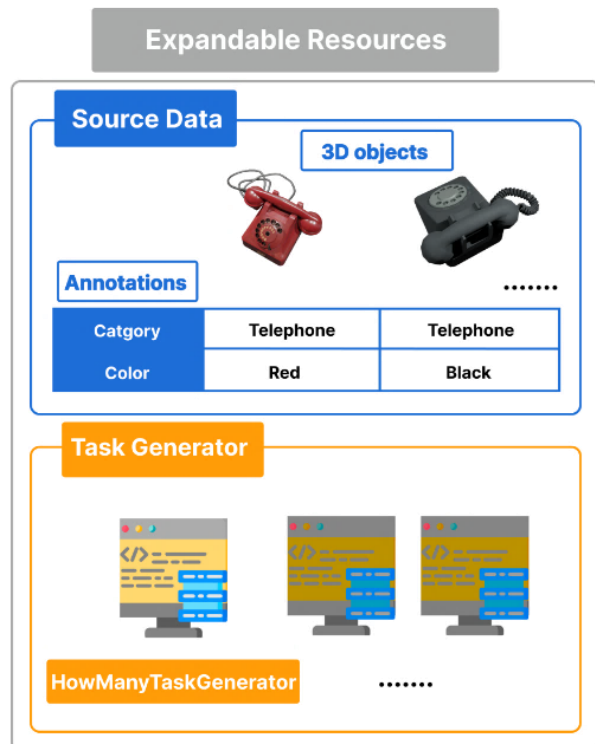
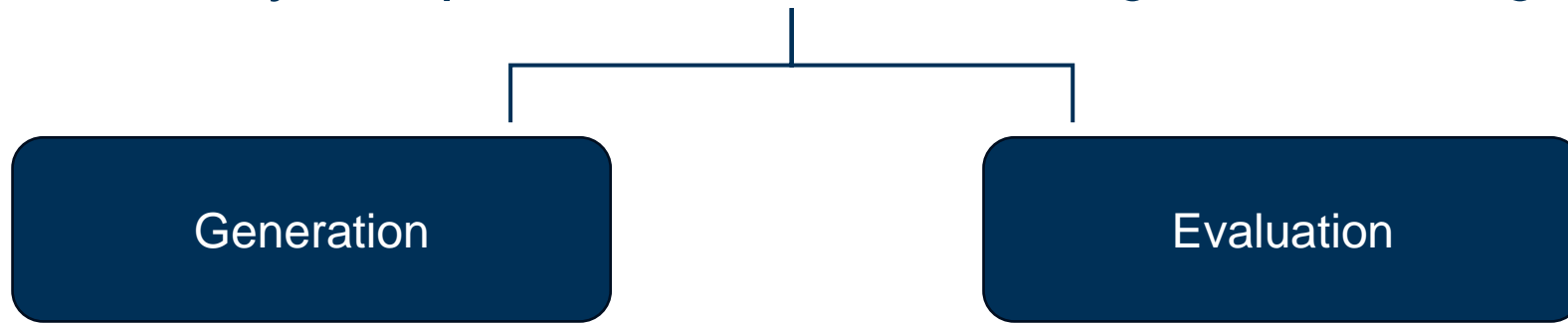
**(No AI MODELS INVOLVED !!!)**

## Contributions:

- ✓ User-centric benchmark generation
- ✓ Expandable Task space (~ 750M tasks)
- ✓ Supports fine-grained user query with budget approximation

# Introduction

Two key components for benchmark generation engines:



# Related Work





# Programmatic Task Generation

## Leveraging Scene Graphs:

GQA Dataset <sup>[1]</sup>:

Similar scene graph-based approach to VQA for more realistic scenes

Question Pattern:

What <type> is <Object>, <attribute> or <Attribute>?

What color is the apple on the white table, red or green?

select: table → filter: white → relate(subject,on): apple → query: color



**Pattern:** What|Which <type> [do you think] <is> <dobject>, <attr> or <decoy>?

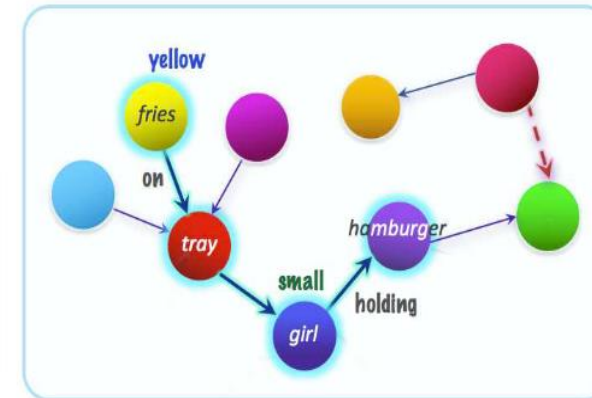
**Program:** Select: <dobject> → Choose <type>: <attr>|<decoy>

**Reference:** The food on the red object left of the small girl that is holding a hamburger

**Decoy:** brown

What color is the food on the red object left of the small girl that is holding a hamburger, yellow or brown?

Select: hamburger → Relate: girl, holding → Filter size: small → Relate: object, left → Filter color: red → Relate: food, on → Choose color: yellow | brown



### Graph Normalization

- Ontology construction
- Edge Pruning
- Object Augmentation
- Global Properties

### Question Generation

- Pattern Collection
- Compositional References
- Decoy Selection
- Probabilistic Generation

### Sampling and Balancing

- Distribution Balancing
- Type-Based Sampling
- Deduplication

### Entailment Relations

- Functional Programs
- Entailment Relations
- Recursive Reachability

### New Metrics

- Consistency
- Validity & Plausibility
- Distribution
- Grounding

[1] Hudson, D.A. and Manning, C.D., 2019. Gqa: A new dataset for real-world visual reasoning and compositional question answering. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition* (pp. 6700-6709).

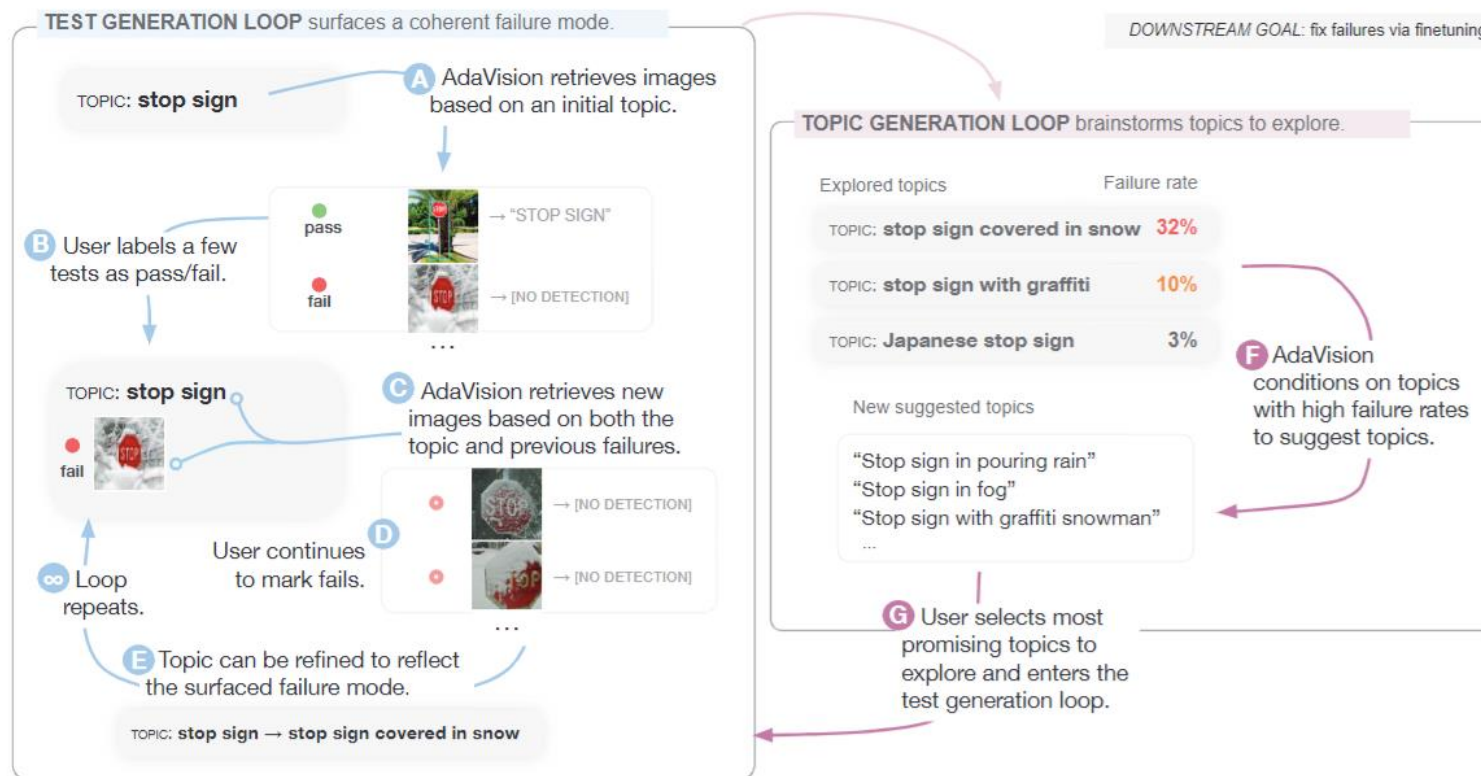


# Adaptive Evaluation



## Adaptive Testing and Debugging:

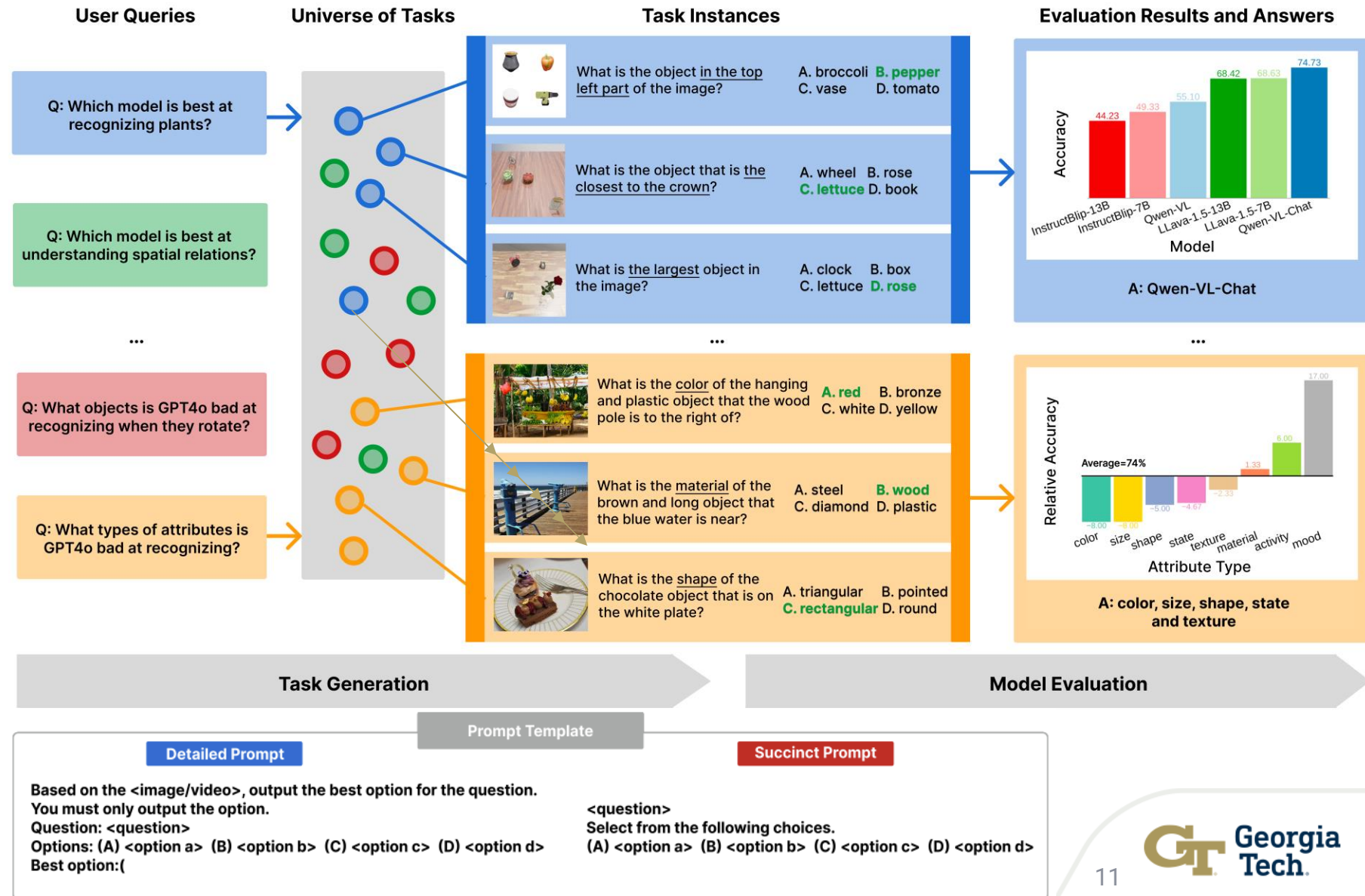
- Dynamically update test data (e.g., Dynabench, LatestEval)
- Adaptively identify task groups where the model underperforms



# Approach

# Approach: Overview

- User brings benchmarking queries
- Program Identify relevant tasks
- Generates relevant VQA Task Instances
- Evaluate model on Scene Graphs & Generated Images
- Expandable processes



# Approach: Terminology

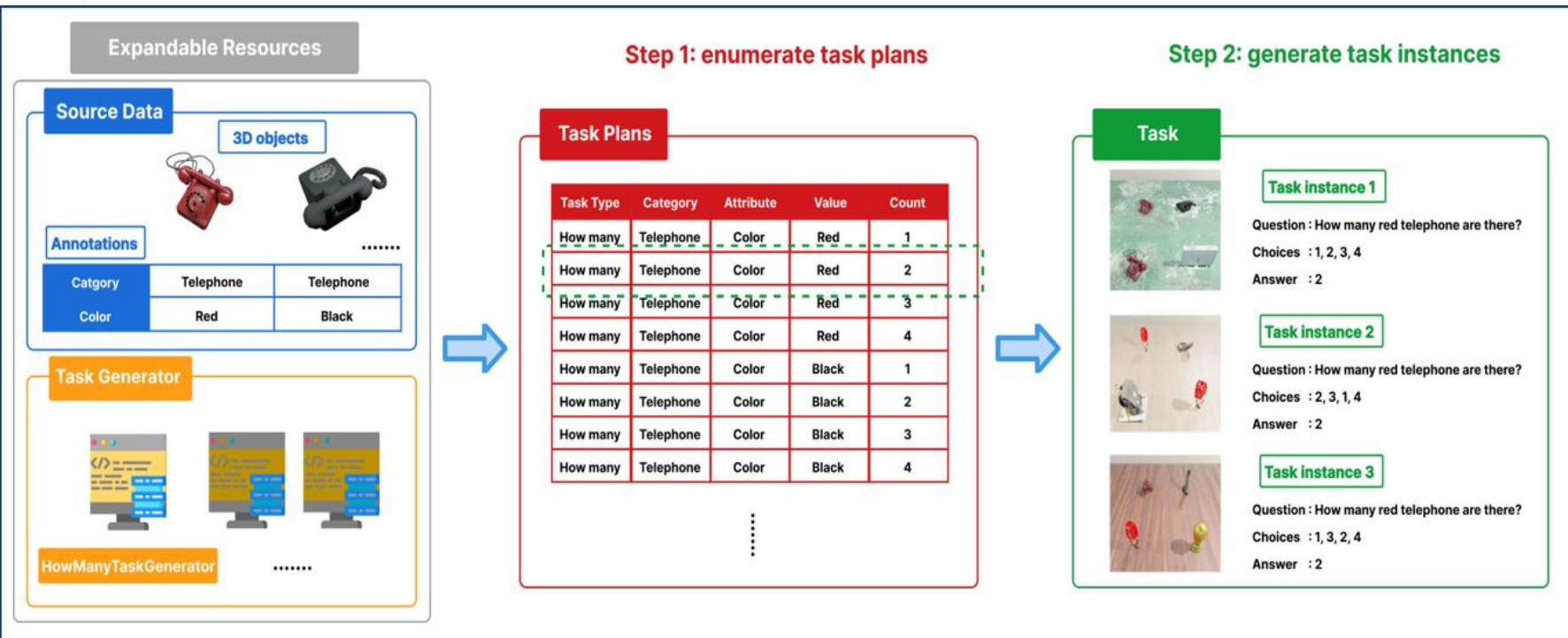
## Example :

- **Source Data:** Images, scenes, attributes
- **Query:** What model is best at counting red Telephones?



## Task generator (program):

1. Uses the benchmarking **query** (input) to enumerate all relevant **task plans** (table)
2. Finds relevant **tasks** (questions + mc)
3. Generates **task instances** (Images) for VQA



**Taxonomy:** Task Space / Assets

**Task Generator:** program

**Task Plans:** meta data

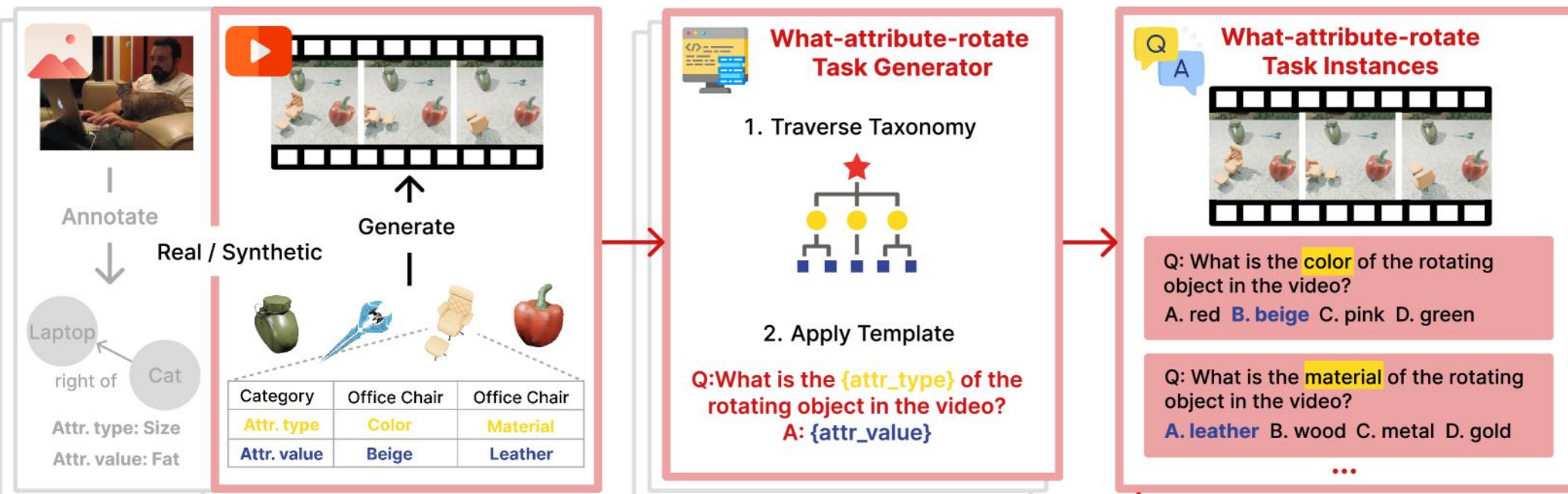
**Task :** The question + MC

**Task Instance:** VQA + Image



# Approach: Generating Tasks

- Images/videos programmatically rendered using Blender / image renderer  
OR
- Real images used with pre-annotated scene graph
- Questions programmatically generated from templates + taxonomy/attributes
- False but plausible answers generated using Adversarial Filtering (LLM, but no image input)




# Approach: Scene Types

## Skills

- counting
- spatial
- relation
- attribute
- 3D attribute
- object
- temp. attr
- action


## Image QA

**Visual Input 1:**  
2D Sticker  
Image



Task Generator		Number of Tasks
how many	<span style="color: red;">■</span>	32,487
what	<span style="color: blue;">■</span>	58,689,137
where	<span style="color: orange;">■</span>	58,689,137
what attribute	<span style="color: green;">■</span>	47,541,884
where attribute	<span style="color: orange;">■</span> <span style="color: green;">■</span>	47,541,884
<b>Total</b>		<b>212,494,529</b>


**Visual Input 2:**  
3D Tabletop  
Scene



Task Generator		Number of Tasks
how many	<span style="color: red;">■</span>	32,487
what	<span style="color: blue;">■</span>	58,689,137
where	<span style="color: orange;">■</span>	58,689,137
what attribute	<span style="color: green;">■</span>	47,541,884
where attribute	<span style="color: orange;">■</span> <span style="color: green;">■</span>	47,541,884
<b>Total</b>		<b>431,525,153</b>

Task Generator		Number of Tasks
what size	<span style="color: lightblue;">■</span> <span style="color: blue;">■</span>	19,688
what attribute size	<span style="color: lightblue;">■</span> <span style="color: green;">■</span>	15,968
where size	<span style="color: lightblue;">■</span> <span style="color: orange;">■</span>	54,164,744
what distance	<span style="color: lightblue;">■</span> <span style="color: blue;">■</span>	58,657,144
what attribute distance	<span style="color: lightblue;">■</span> <span style="color: green;">■</span>	47,515,936
where distance	<span style="color: lightblue;">■</span> <span style="color: orange;">■</span>	58,657,144

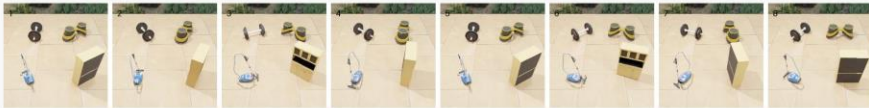
**Visual Input 3:**  
Scene Graph  
with Real Image



Task Generator		Number of Tasks
what object	<span style="color: blue;">■</span>	25,173
what attribute	<span style="color: green;">■</span>	20,546
what relation	<span style="color: teal;">■</span>	23,241
<b>Total</b>		<b>68,960</b>


## Video QA

**Visual Input 4:**  
3D Tabletop  
Scene



Task Generator		Number of Tasks
what rotate video	<span style="color: magenta;">■</span> <span style="color: blue;">■</span>	39,376
what attribute rotate video	<span style="color: magenta;">■</span> <span style="color: green;">■</span>	31,936
where rotate video	<span style="color: magenta;">■</span> <span style="color: orange;">■</span>	108,329,488
what move video	<span style="color: magenta;">■</span> <span style="color: blue;">■</span>	78,752
what attribute move video	<span style="color: magenta;">■</span> <span style="color: green;">■</span>	63,872
where move video	<span style="color: magenta;">■</span> <span style="color: orange;">■</span>	78,752
<b>Total</b>		<b>108,622,176</b>

**Visual Input 5:**  
Scene Graph  
with Real Video



Task Generator		Number of Tasks
what object video	<span style="color: blue;">■</span>	428,342
what relation video	<span style="color: teal;">■</span>	428,342
what action video	<span style="color: purple;">■</span>	335,386
<b>Total</b>		<b>1,192,070</b>



# Approach: Task Space

- **133k** Images, 10k videos, 2k 3D objects
- **365** Object Categories
- **655** Attributes (color, texture, size)
- **335** Relationships (spatial, modeled within scene graphs)
- **28** Task Generators (how many, what color)
- **5** Types of Visual Input (2D tabletop, 3D, video)
- **750M** possible image/video question-answering pairs
- Much larger than comparable datasets: GQA(22M), (CLEVR 100k)
- Can be specialized for specific tasks

# Approach: Query Types

## Top-K

Find Top-5 objects that GPT-4o is worst at recognizing when rotating

## Threshold

Identify colors that GPT-4o recognizes with less than 30% average accuracy

## Model Compare

Compare with LLaVA-Next to determine which objects GPT-4o performs better at recognizing

## Model Debug

Identify tasks for which InstructBLIP performs significantly worse than its average performance

- Querys are not open language.
- Must be written 'SQL' style language using the set of attributes, query types, relationships, etc.

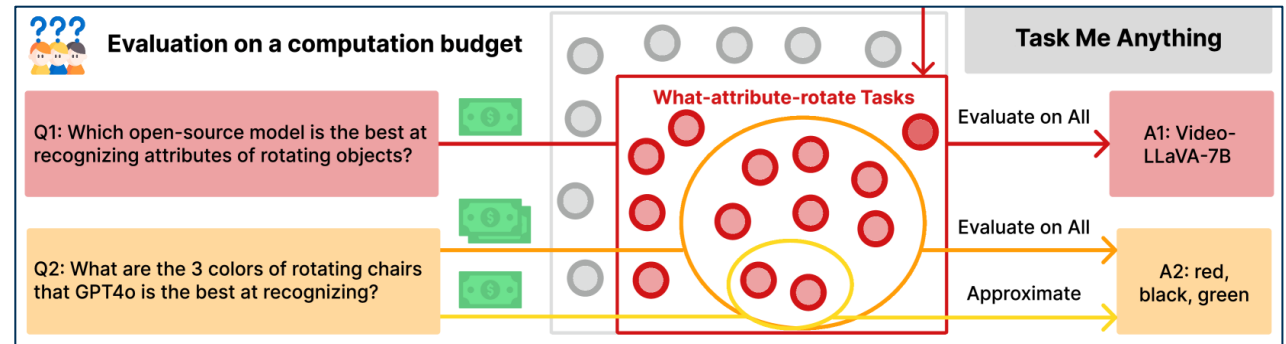
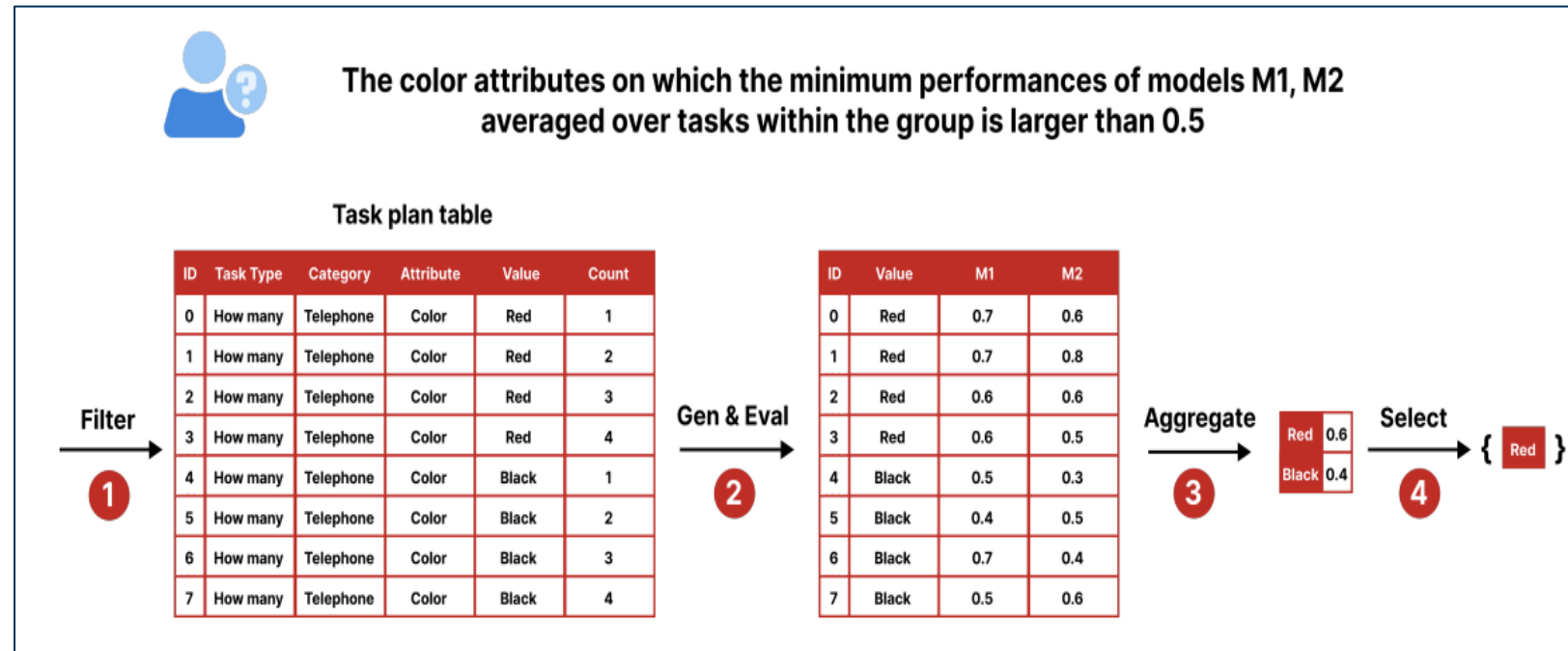
# Approach: Query Execution / Cost Optimization

- **Random Sampling:**

Randomly selects a subset of task instances to evaluate MLMs.

- **Fitting:** Trains regressor to predict MLM performance based on past samples and task metadata.

- **Active Learning:** Iteratively refines the regressor by sampling most uncertain task instances for improved predictions.



# Approach: Output/Contributions

- **TASK-ME-ANYTHING:** Task Generator process/code itself. Expandable with new datasets/features and can be used to generate new, custom benchmarks
- **TASK-ME-ANYTHING-RANDOM:** 100 random tasks, 5700 ImageQA and 2700 VideoQA instances. Evaluated 18 MLMs with detailed and succinct prompts
- **TASK-ME-ANYTHING-DB:** Over 100K tasks generating 1M+ instances. 13 MLMs evaluated across 24.24 million evaluation pairs. Results aid in model performance prediction
- **TASK-ME-ANYTHING-UI:** Graphical interface with tabs for model performance, task embedding visualization, performance anomalies, and detailed query investigations

# Experiments and Results

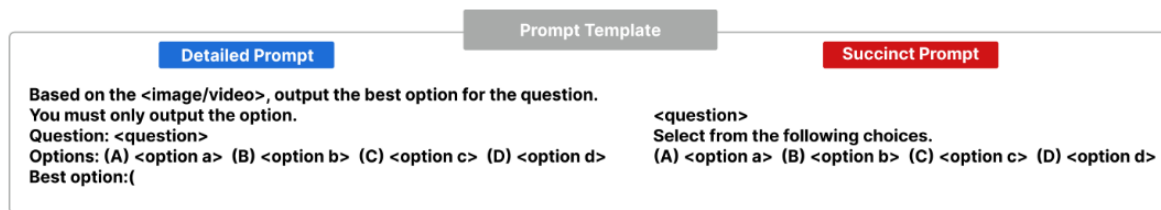
# Experiments and Results: Overview

Pipeline for automatic benchmark generation paper



Benchmarking of previous work!

- Evaluation of 18 different VLMs across 8400 tasks instances made public



100 random tasks  
3 instances each

100k random tasks  
15 instances each

- 1) Proposed benchmark: random subset of 2700 (IQA)/5700 (VQA) tasks instances (TMA-Random)
- 2) Open source VLMs (13): over 1M tasks instances (TMA-DB)
- 3) UI to explore TMA-DB with different queries (TMA-UI)



# Welcome to TaskMeAnything-UI!

Overview Embedding Query Surprisingness

## Visualize the overall task distribution and model performance

scenario

imageqa-2d-sticker  imageqa-3d-tabletop  imageqa-scene-graph  videoqa-3d-tabletop  videoqa-scene-graph

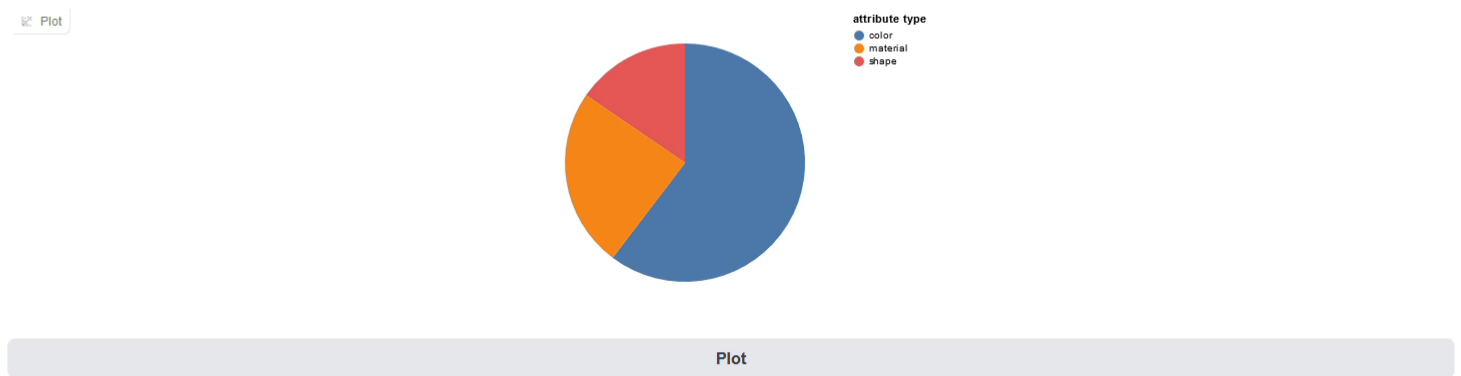
task space of the following task generator

2d-what-attribute

### Overall task metadata distribution

task metadata

attribute type



### Models' overall performance by task metadata

model

qwenvl-chat  qwenvl  llava15-7b  llava15-13b  instructblip-vicuna13b  instructblip-vicuna7b

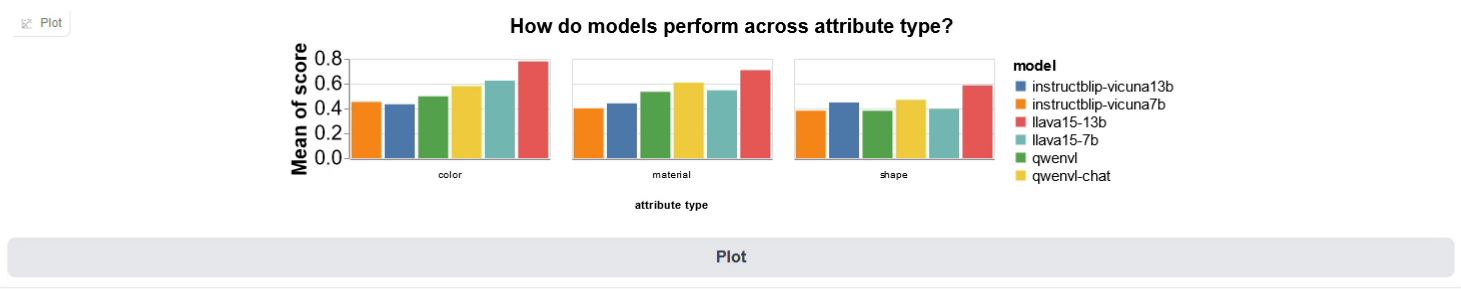
aggregate models' accuracy by

mean  median  min  max

task metadata

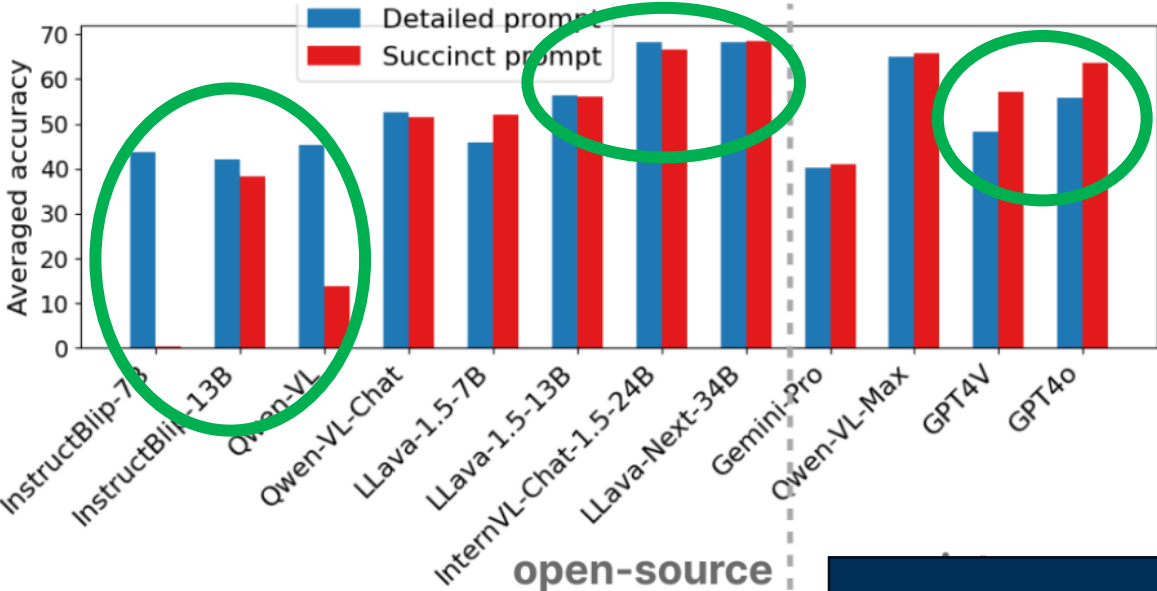
attribute type

Optional: second task metadata

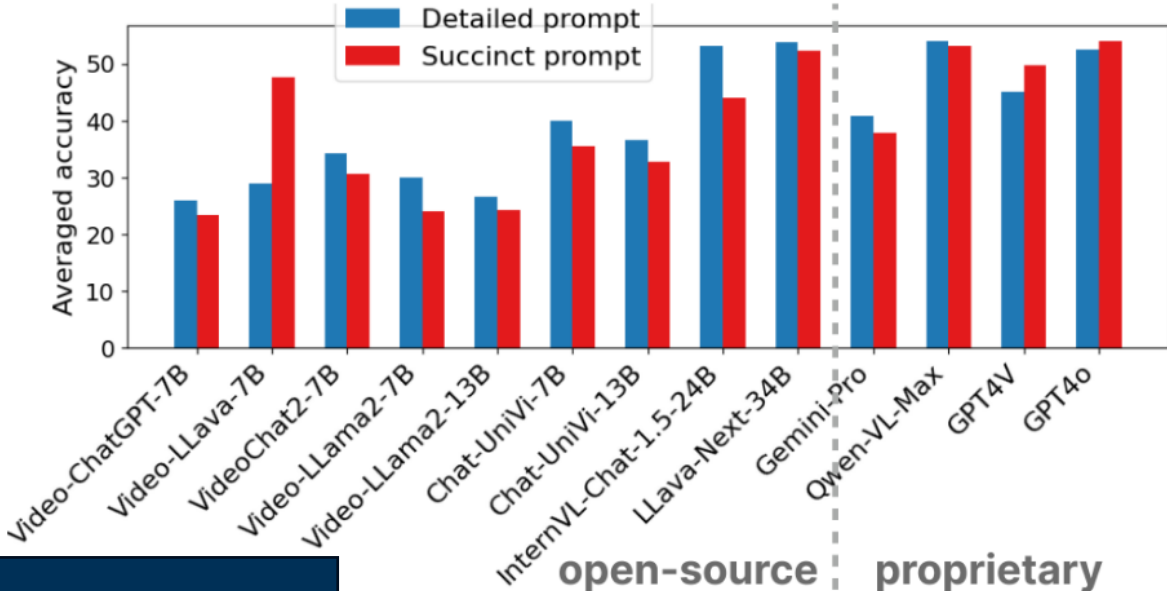


# Q1: How do models perform over a random subset of all possible questions?

ImageQA



VideoQA



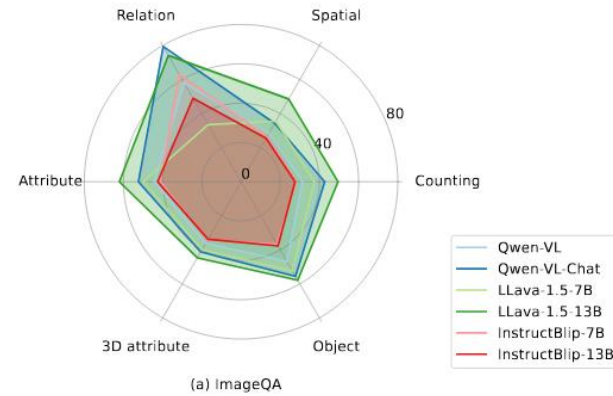
Each model might present different prompt sensitivity

Detailed Prompt	Succinct Prompt
Based on the <image/video>, output the best option for the question. You must only output the option. Question: <question> Options: (A) <option a> (B) <option b> (C) <option c> (D) <option d> Best option: {	<question> Select from the following choices. (A) <option a> (B) <option b> (C) <option c> (D) <option d>

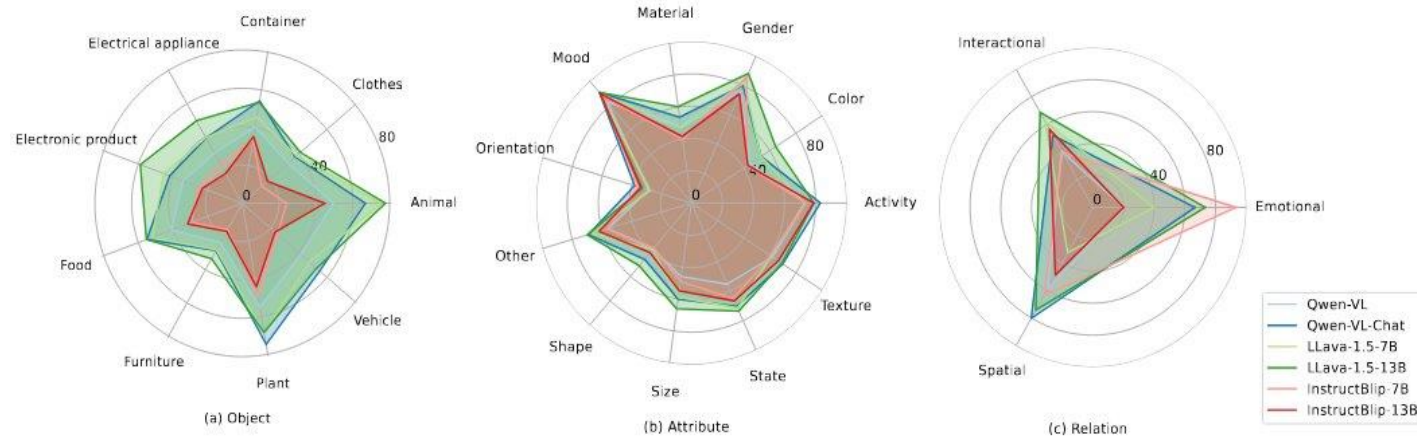
# Q2: What is the best MLM for each specific skill? (IQA)

Different models have different expertise

## High-level skills



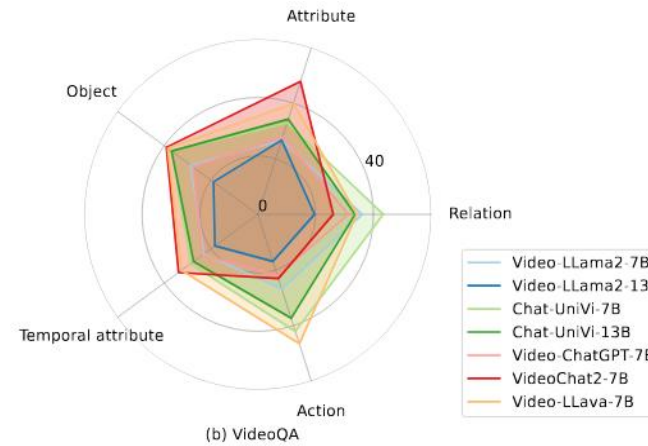
## Fine-grained



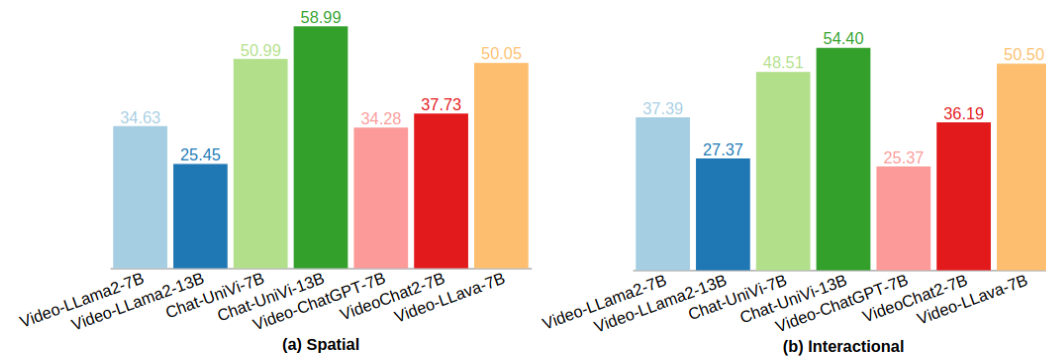
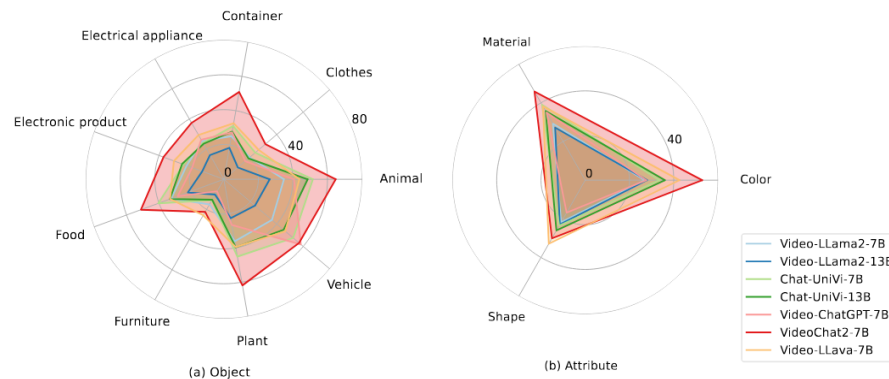
# Q3: What is the best MLM for each specific skill? (VQA)

In VQA, bigger difference in expertise per model

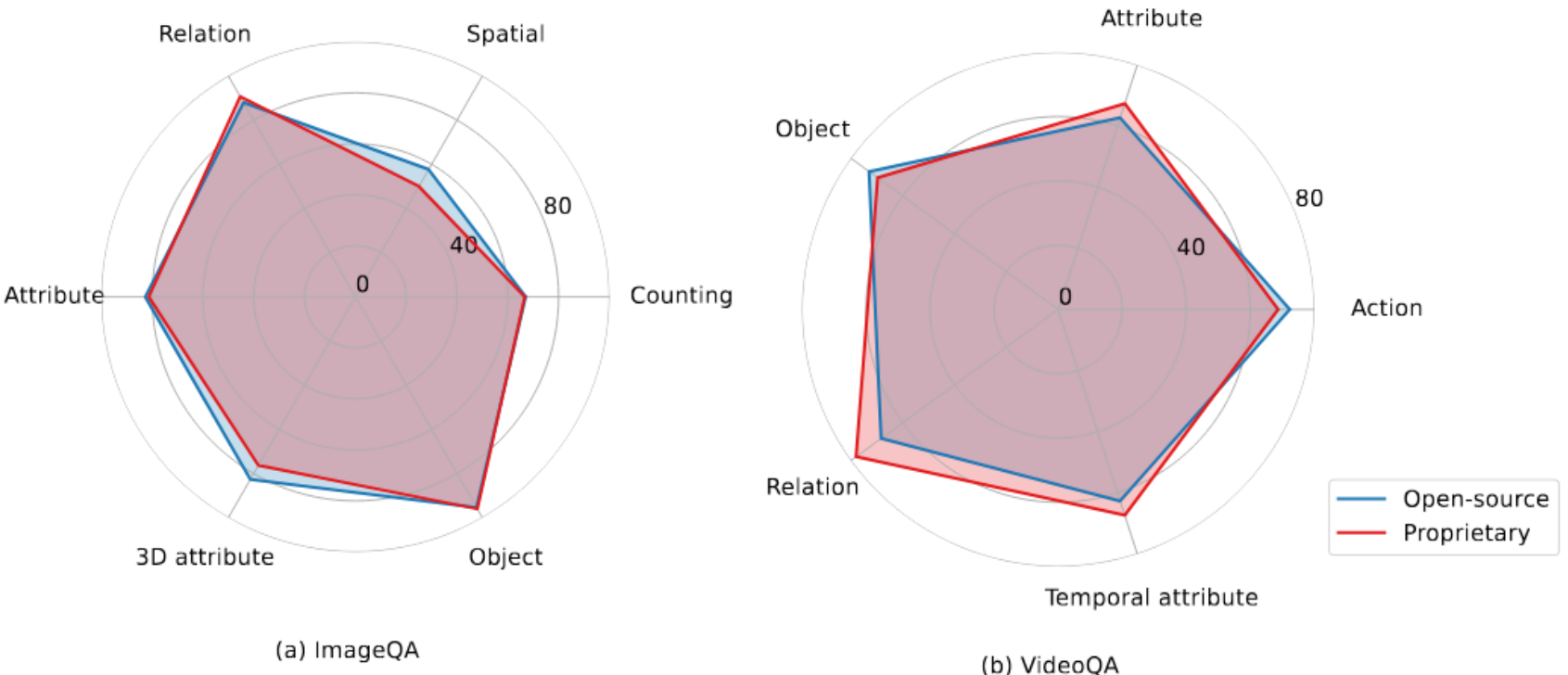
## High-level skills



## Fine-grained



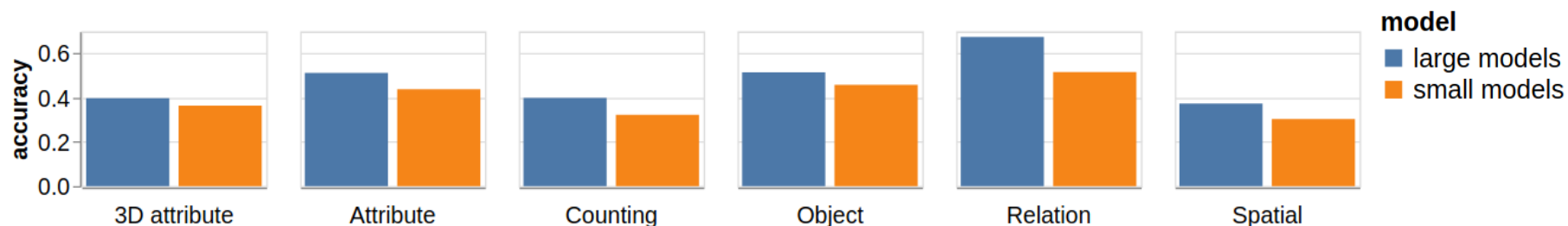
# Q4: How does the best open-source model compare against the best proprietary model across skills?



Open and closed source models are overall competitive

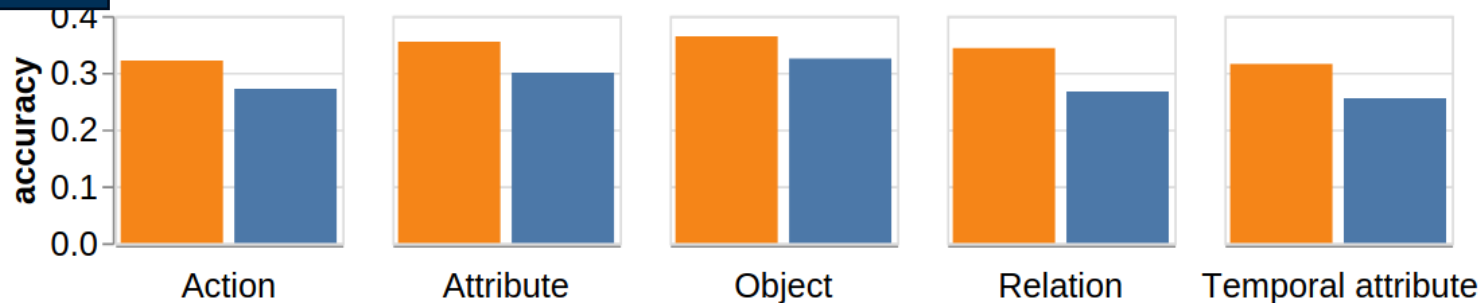
# Q5: How do small models compare against large models?

## ImageQA



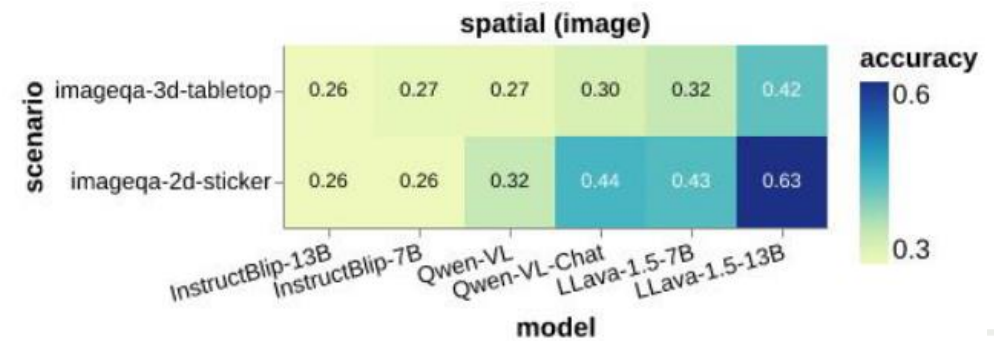
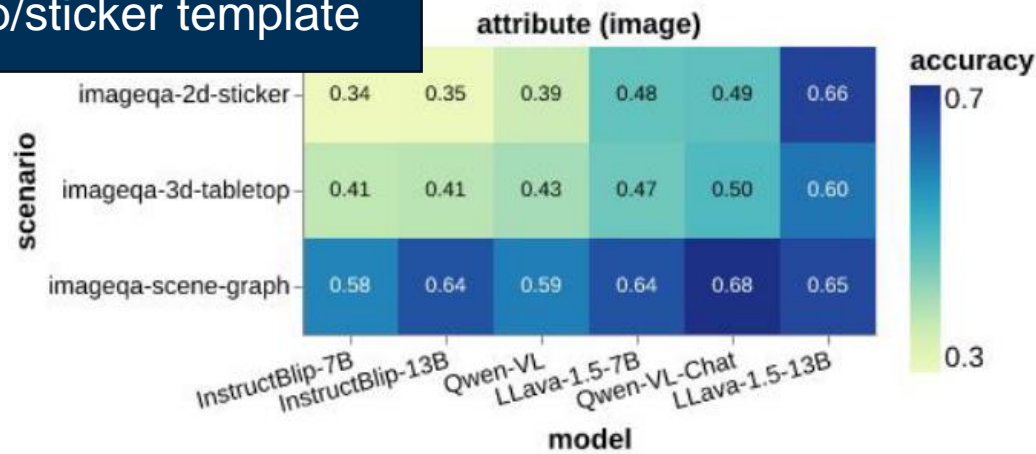
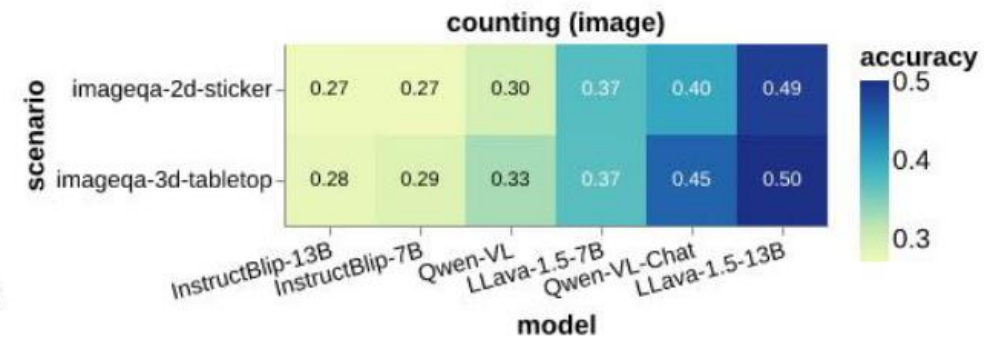
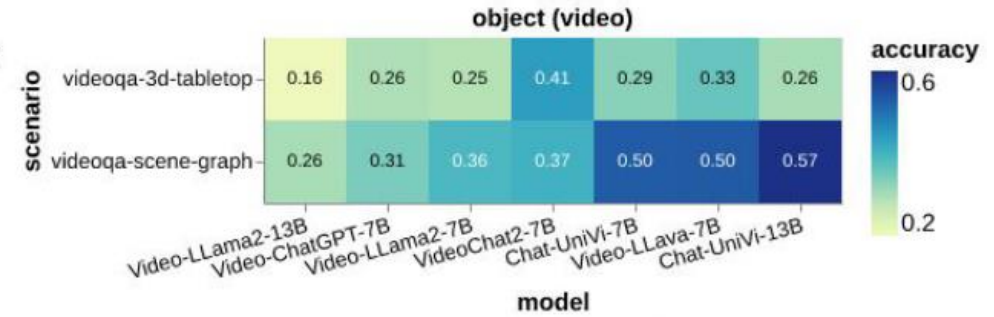
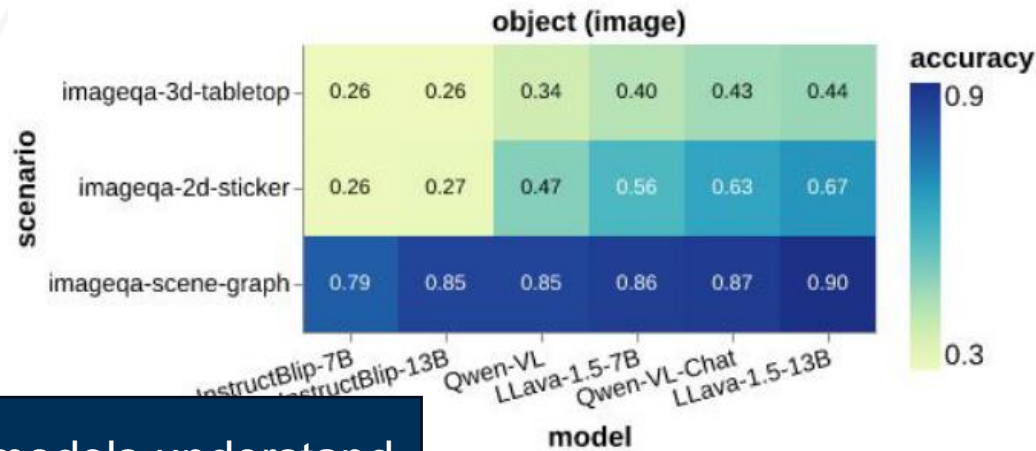
ImageQA: the larger the better (overall)  
VideoQA: the other way around (!!!)

## VideoQA





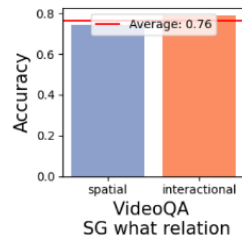
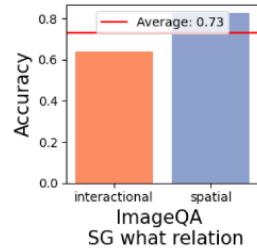
# Q6: Are models' strengths and weaknesses consistent across visual inputs?



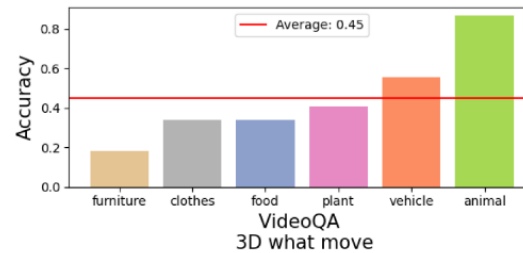
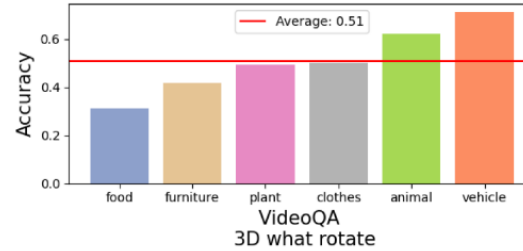
Overall, models understand scene-graph better than 3D tabletop/sticker template

# Q7: What is today's popular proprietary model (GPT4o) bad at?

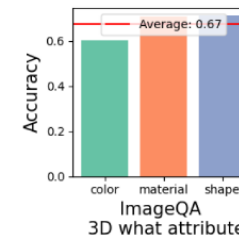
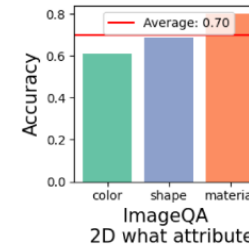
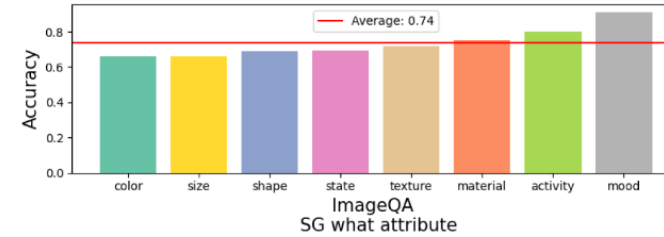
(Q1) What relations are GPT4o bad at understanding?



(Q2) What objects are GPT4o bad at recognizing when rotating/moving?



(Q3) What attribute of objects are GPT4o bad at recognizing?



Spatial relations are harder to maintain as objects/camera angle changes

Models "expect" certain temporal behaviors for different objects

Size of color dictionary is significantly larger than mood dictionary

Question	Task	Relations/attributes	$\Delta$ Perf. (%)
what relations are GPT4o bad at understanding?	VideoQA 3D what rotate	duct, hamper, tool, furniture, air conditioner	-21.67
what objects are GPT4o bad at recognizing when rotating/moving?	VideoQA 3D what move	g, pushing, pushed by, carrying, above	-19.33
what attributes of objects are GPT4o bad at recognizing?	ImageQA 2D what attribute	purple, brown, red, gray, beige	-51.05
	ImageQA 3D what attribute	stone, rubber, textile, leather, plastic	-16.66
	ImageQA SG what attribute	crooked, power, lower, steep, glowing	-5.33
			-10.67
			-45.45

# Discussion

# Limitations and Societal Implications

- **Generated tasks can be unrealistic and biased:** might not capture the nuances of real-world scenarios
- **Designing task space is challenging:** identifying the relevant attributes for each task type might require domain knowledge
- **Adding new task generators requires technical expertise**
- **Inaccuracies in Query results:** Efficient query results approximation within certain budgets might sometimes yield inaccurate results, especially when the budget limits are constrained

# Limitations and Societal Implications

- **Misuse for malicious benchmarks:** create adversarial examples to trick or expose vulnerabilities in AI systems.
- **Reinforcing biases and discrimination:** if task generators are not carefully designed, they could perpetuate biases in source data
- **Overreliance on synthetic tasks:** might create a false sense of progress and hinder the development of AI models that effectively address real-world challenges.
- **Data contamination:** Models might learn to exploit patterns in synthetic data and fail to generalize to real-world scenarios.
- **Access and fairness:** Requires technical expertise to create new task generators

# Summary of Strengths and Weaknesses

- Strengths:

- Provide a systematic approach for evaluating different MLMs for user-specific task requirements.
- Enable users to provide fine-grained queries (e.g., top-k, threshold) for task generation.
- Provide a database with different open-source and proprietary MLMs evaluated on several benchmarks.
- Enable users to evaluate different MLMs on fixed computational budget using approaches such as fitting and active learning.

- Weaknesses:

- Synthetically generated data might not capture the nuances of real-world scenarios.
- MLMs might learn to exploit specific patterns in synthetic data (especially when trained at scale) and may not generalize well for practical applications.
- Inaccuracies in estimating model performance for different queries under constrained budget.