

# From Head to Tail: Towards Balanced Representation in Large Vision-Language Models through Adaptive Data Calibration

---

Poster: Jun 13th, 16:00-18:00 ExHall D Poster #387

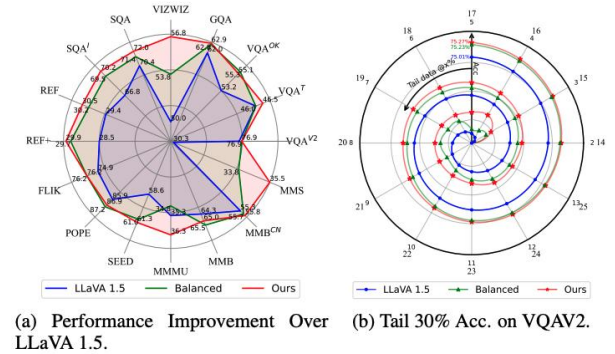
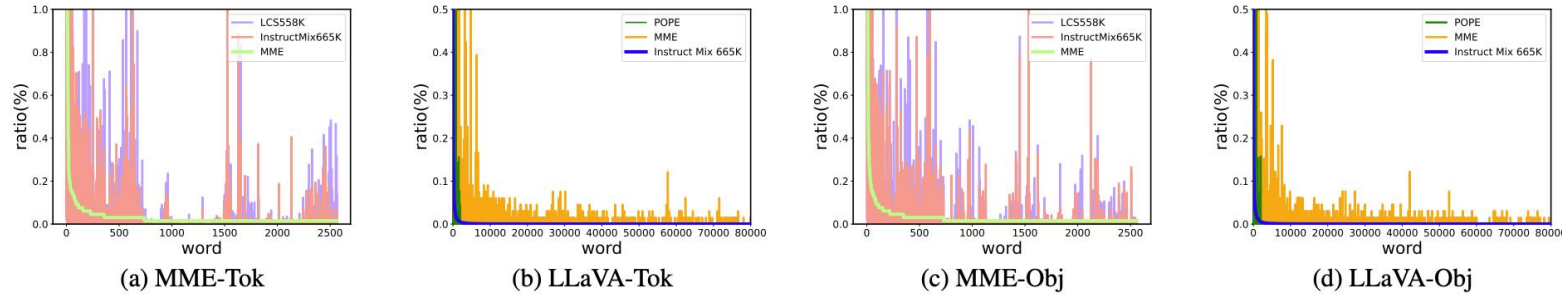
By: Mingyang Song



- **Motivation**
- **Pipeline**
- **Methods**
- **Experiments**

## 1. Current LVLM training data suffers from the “long-tail” problem (LT):

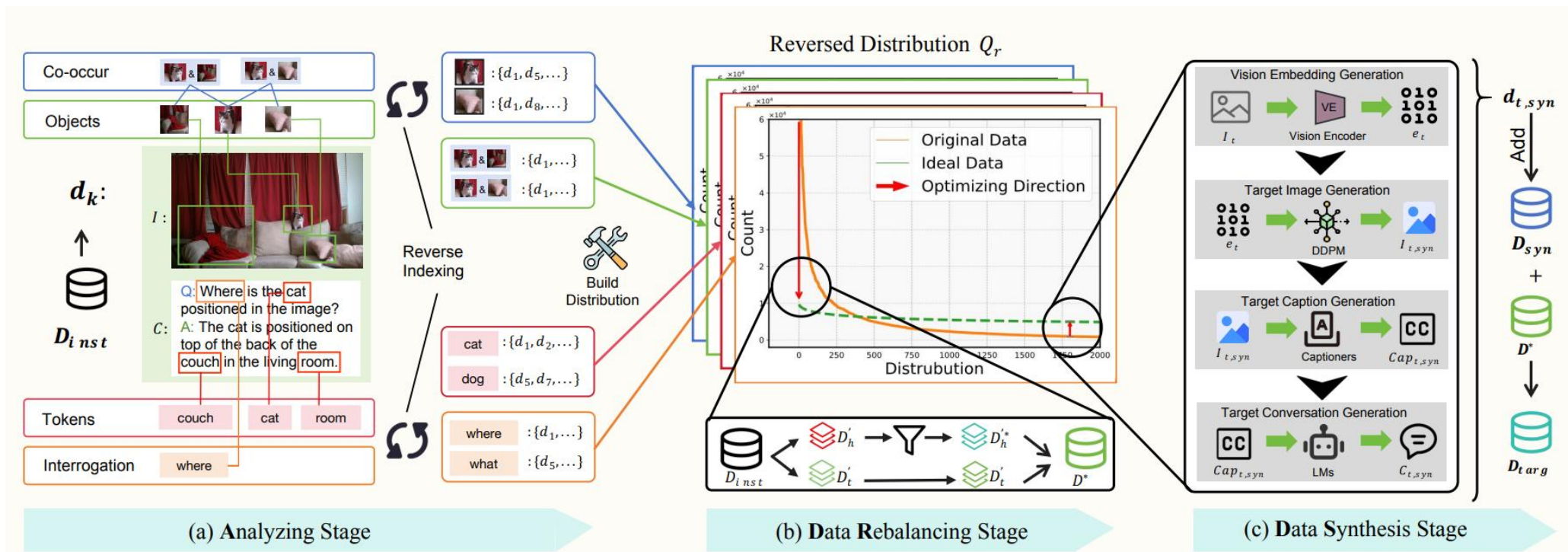
The training datasets present highly imbalanced distributions, with many tail concepts, and often **differ in distribution** from the test set.



2. The model is more **prone to making mistakes on tail Concepts**: We analyze the distribution of entities in failed cases by computing their positions, and compare them with those in correct cases. Failed cases tend to appear later in the distribution.

Methods	MME						POPE					
	Tok-C	Tok-W	Obj-C	Obj-W	Co-C	Co-W	Tok-C	Tok-W	Obj-C	Obj-W	Co-C	Co-W
Max	9738	10377	2708	3222	247315	257107	2242	2772	1085	1100	130043	141722
		+639		+514		+9792		+30		+15		+11679
Min	1	1	60	131	12732	20741	1	1	17	21	926	1033
		+0		+71		+8009		+0		+4		+107
Mean	1035	1068	842	1035	71123	79104	313	340	319	336	27457	30989
		+33		+193		+7981		+27		+17		+3532

3. LT in LVLMs poses **unique, under-explored challenges** due to cross-modal complexity, interactions, and distinct co-occurrence patterns.



An overview of our Adaptive Data Refinement Framework (ADR)

- Analyzing Stage:** we first extract **tokens, objects, co-occurrences, and interrogations** from the training instances, then construct corresponding distribution using a reverse-indexed mapping.
- Data Rebalancing stage:** we analyze the optimizing direction and adaptively rebalance the redundant data based on the entity distribution identified in the Analyzing stage.
- Data Synthesis stage,** we utilize DDPM and the latent representations of scarce image instances to synthesize the underrepresented data.

## 1. Analyzing Stage

### 1) Entity Distribution Construction

Specifically, we conduct the whole analysis procedure by constructing the frequency distribution of entities  $Q_e$  from these **four perspectives** among the whole training set.

Token:  $e_t = \{n | n \in \text{Noun} \wedge n \subseteq C \text{ for } (I, C) \text{ in } D\}$

Object:  $e_o = \{o | o \in I \text{ for } (I, C) \text{ in } D\}$

Co-occurrence:  $e_c = \{(o_1, o_2) | o_1 \in I \wedge o_2 \in I \text{ for } (I, C) \text{ in } D\}$

Interrogation:  $e_w = \{q | q \in Q \wedge q \in C \text{ for } (I, C) \text{ in } D\}$

Data	Level	thres	% E	% DI
LLaVA [27]	Tok	120	98.7	10.0
	Obj	304	98.0	10.0
	Co	24	92.7	25.0
	Int	4895	99.6	10.0
Avg.	-	-	<b>97.25</b>	<b>13.75</b>

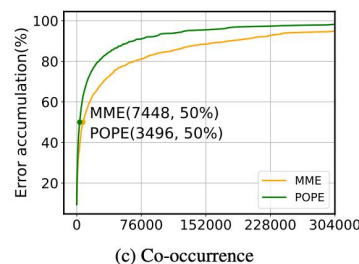
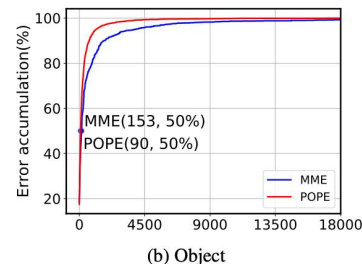
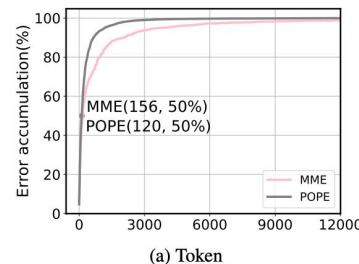
### 2) Reverse Indexing

Subsequently, we use the number of data instances corresponding to each entity as frequency to build the reversed distribution  $Q_r$ :  $Q_r = \{e_1 : N_{e_1}, e_2 : N_{e_2}, \dots, e_n : N_{e_n}\}$  where  $e_i$  means entity item and  $N_{e_i}$  means the number of corresponding data instances of  $e_i$ .

Surprisingly, among four perspectives, an average of **97.25%** entries account for only **13.75%** data instances on average,

### 3) Discovery

- Tail data accounts for **more failed cases**.
- Distribution **varies** between train and test data.





## 2. Data Rebalancing Stage

We construct entity-wise sampling probabilities to downsample redundant data and select a core dataset based on entity frequency and coverage.

## 3. Data Synthesis Stage

### 1) Language Data Synthesis

We rewrite head-biased data by replacing dominant concepts with tail synonyms via LLM-guided paraphrasing.

### 2) Visual Data Synthesis

- given a tail instance  $d_t = (I_t, C_t)$  our objective is to generate an image similar to  $I$  and produce corresponding instruction data
- We leverage ControlNet to generate images that retain the original content while exhibiting a similar visual style  $I_{t,syn} = G(I_t, p_{def})$ .
- We use an off-the-shelf vision captioner to generate captions for the synthetic image, which are then expanded into full conversations using LLMs for visual instruction tuning.

#### Algorithm 1 Pseudo Code for Data Resampling

```

1  # D: raw training set;
2  # C: target perspectives list
3  # tau: the threshold for entities;
4  # D_bal: the rebalanced data, a.k.a. D*;
5  # n_p, alpha: hyperparameters
6  D_bal=[]
7  for pers in C:      # build prob dict
8      entity_dist =
9          entity_distribution_construction
10         (D,pers)
11     prob_dict[pers] = {ent:tau[pers]/
12                       entry_dist[ent] for ent in
13                       entry_dist.keys()}
14 for instance in D: # data rebalancing
15     pass_cnt = 0
16     for pers in C:
17         for entity in instance['entity']
18             [pers]:
19                 if random.random() <
20                     prob_dict[pers][entity]:
21                         pass_cnt += 1
22                         break
23     if pass_cnt > n_p and random.random
24         () < alpha:
25             D_bal.append(instance)

```

## ► Quantitative Analysis of Different Methods on Popular Benchmarks

Method	IT*	VQA <sup>OK</sup>	SEED <sup>2</sup>	QB <sup>2</sup>	MMS	MME <sup>P</sup>	SQA <sup>I</sup>	MMMU	VQA <sup>T</sup>	GQA	MMB	VQA <sup>v2</sup>
LLaVA 1.5	665.0K	53.2	48.7	47.3	33.5	1510.7	69.3	35.3	46.0	61.9	64.3	76.6
<b>+DR</b>	581.0K	55.3	57.2	46.8	33.8	1470.6	69.5	34.8	46.0	62.8	<b>65.5</b>	76.9
<b>+DR +DS</b>	665.0K	<b>57.4</b>	<b>57.4</b>	<b>49.6</b>	<b>35.5</b>	<b>1512.8</b>	<b>70.4</b>	<b>36.7</b>	<b>47.2</b>	<b>62.9</b>	65.0	<b>76.9</b>
ShareGPT4V	1246.0K	54.0	59.6	44.2	34.7	1560.4	68.9	35.1	50.2	63.3	68.0	78.6
<b>+DR</b>	1168.0K	56.7	59.6	44.9	35.0	1542.3	68.6	35.7	<b>50.9</b>	<b>63.9</b>	67.9	78.7
<b>+DR +DS</b>	1246.0K	<b>57.9</b>	<b>59.9</b>	<b>45.7</b>	<b>35.5</b>	<b>1564.9</b>	<b>69.4</b>	<b>36.1</b>	50.9	63.7	<b>68.8</b>	<b>78.7</b>

Comparison with models trained with different methods on different benchmarks.

## ► Quantitative Analysis among Different Data Rebalancing Methods

Method	IT	GQA	SEED	SEED <sup>v2</sup>	POPE	MMB
Baseline	665K	62.0	61.0	57.2	87.2	65.5
EL2N	581K	62.5	53.6	47.4	87.2	65.2
Perplexity	581K	62.3	53.4	47.4	86.8	63.7
CLIP Score	581K	62.5	53.0	47.0	87.0	64.5
COINCIDE	133K	59.8	-	-	86.1	63.1
Ours-DR	581K	<u>62.8</u>	<u>61.0</u>	<u>57.2</u>	<u>87.2</u>	<b>65.5</b>
Ours	665K	<b>62.9</b>	<b>61.3</b>	<b>57.4</b>	<b>87.4</b>	<u>65.0</u>

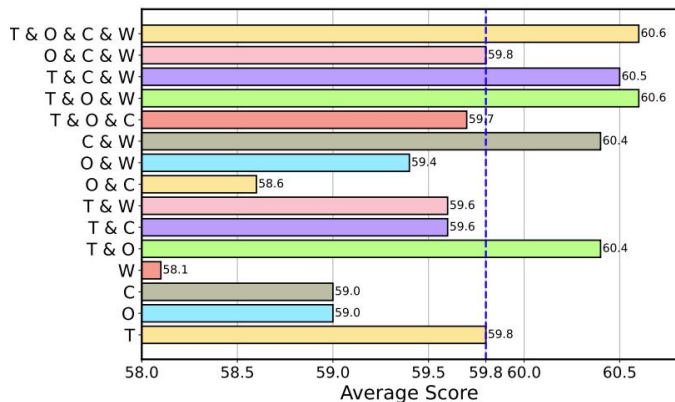
Popular Data Rebalancing Methods.

## ► The Impact of ADR on Tail concepts

Methods	IT	ScienceQA					
		@5	@10	@15	@20	H@80	Overall
LLaVA 1.5	665.0K	67.9	70.0	67.9	68.5	74.6	69.3
<b>+DR</b>	581.0K	69.2	69.7	67.8	68.5	76.2	69.5
<b>+DR +DS</b>	665.0K	<b>70.1</b>	<b>70.5</b>	<b>68.3</b>	<b>69.0</b>	<b>78.6</b>	<b>70.2</b>

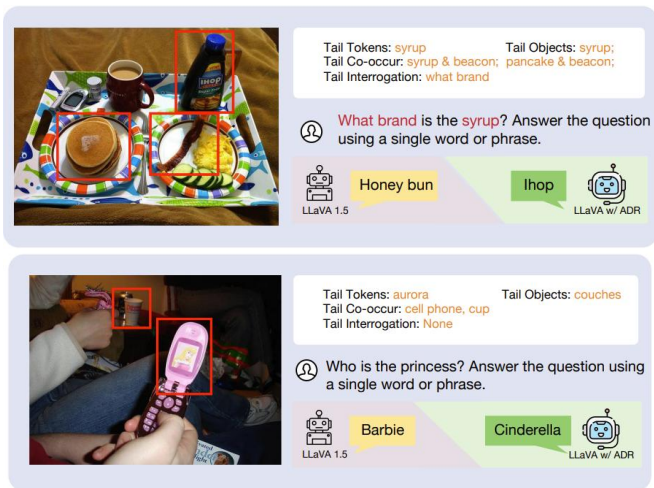
Tail concept prediction accuracy on ScienceQA-IMG

## ► Ablation Study

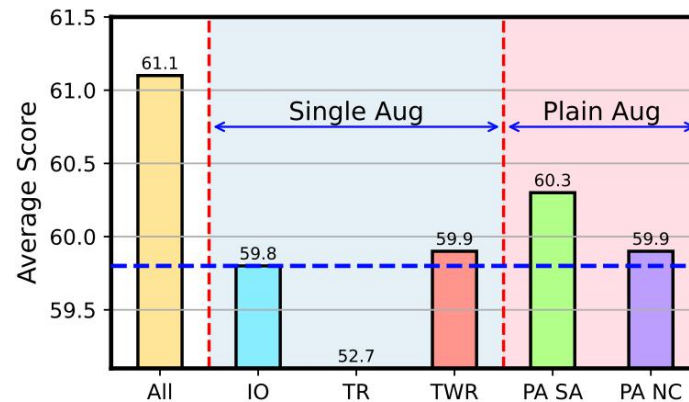


Ablation study on data rebalancing combinations.

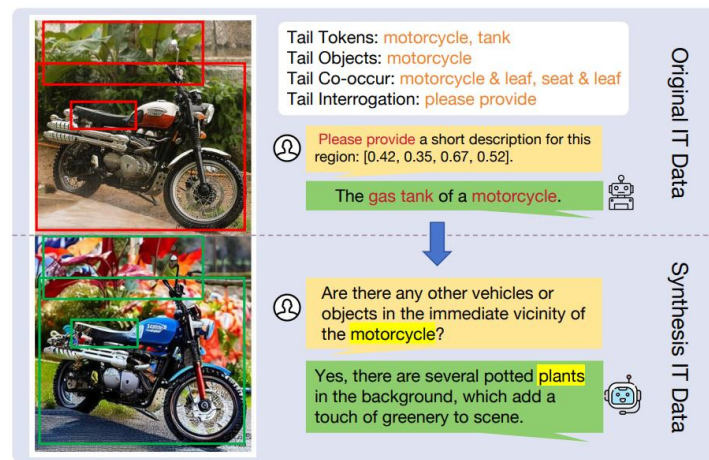
## ► Visualization



Qualitative comparison between LLaVA 1.5 and LLaVA 1.5 w/ ADR



Ablation study on data synthesis methods.



Comparison between the original instruction-tuning (IT) data and our synthesized IT data.



# Thanks!

---

Poster: Jun 13th, 16:00-18:00 ExHall D Poster #387

