

FilterRAG: Zero-Shot Informed Retrieval-Augmented Generation to Mitigate Hallucinations in VQA



NOBIN SARWAR
CS PhD student



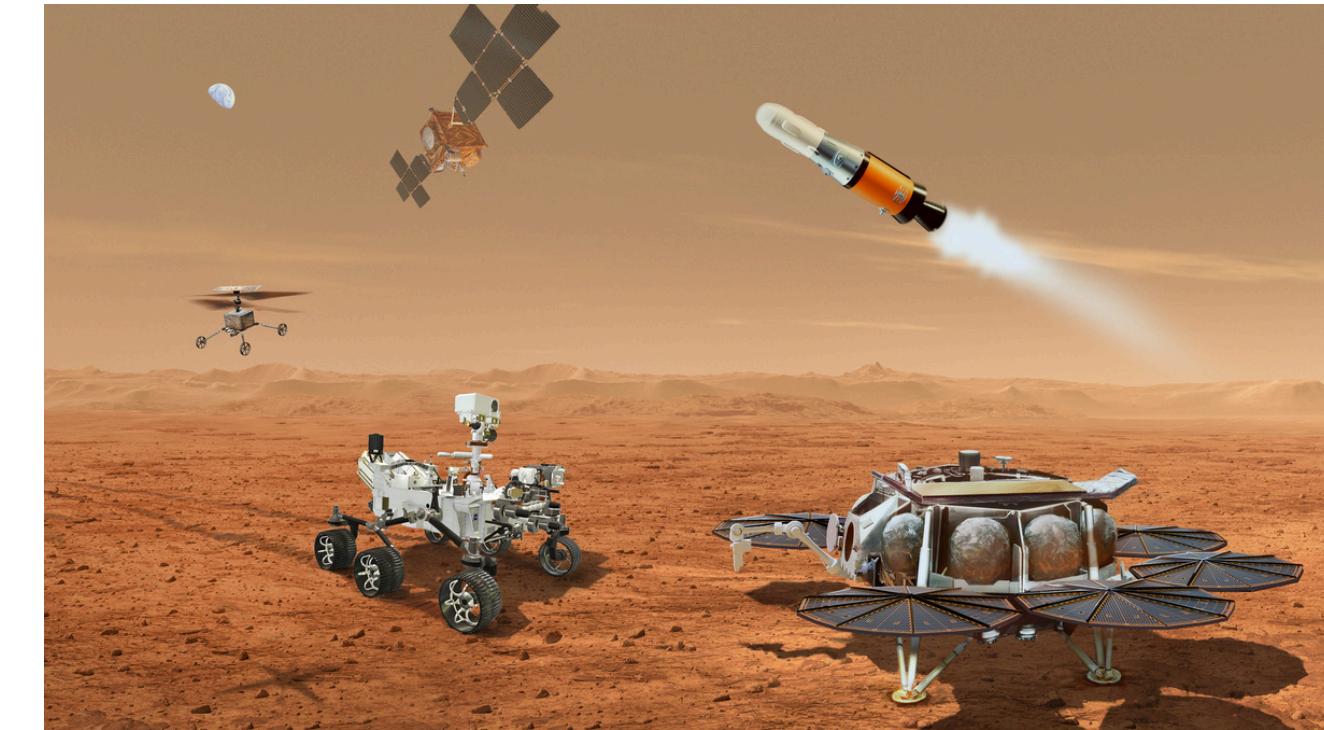
Existing problems

External knowledge



Example: What do they call running around the bases on a single hit?
- home run

Out-of-distribution, Hallucination



Example: What's the capital of Mars?
- The capital of Mars is **Muskland**.
(lack of grounding in reality)

External knowledge, Out-of-distribution, Hallucination

Roadmap

Part 1:

- **Preliminaries: Zero-Shot Learning, VLMs, VQA in VLMs, RAG, OOD, Hallucination**
- **Research Questions**
- **Datasets (VQA V2 and Ok-VQA)**

Part 2:

- Method: Architecture, Loss Function
- Experimental Results (Qualitative, Quantitative, Visual results, and Ablation Study)
- Final Discussion (Contributions and Limitations)

Visual question answering (VQA) in VLMs

- **Inputs:** Given an image (**I**) and a question (**Q**),
- **Goal:** Predict an answer (**A**) to the question (**Q**).

This is expressed as:

$$P(\hat{A}) = \arg \max_{A \in \mathcal{A}} P(A | I, Q)$$

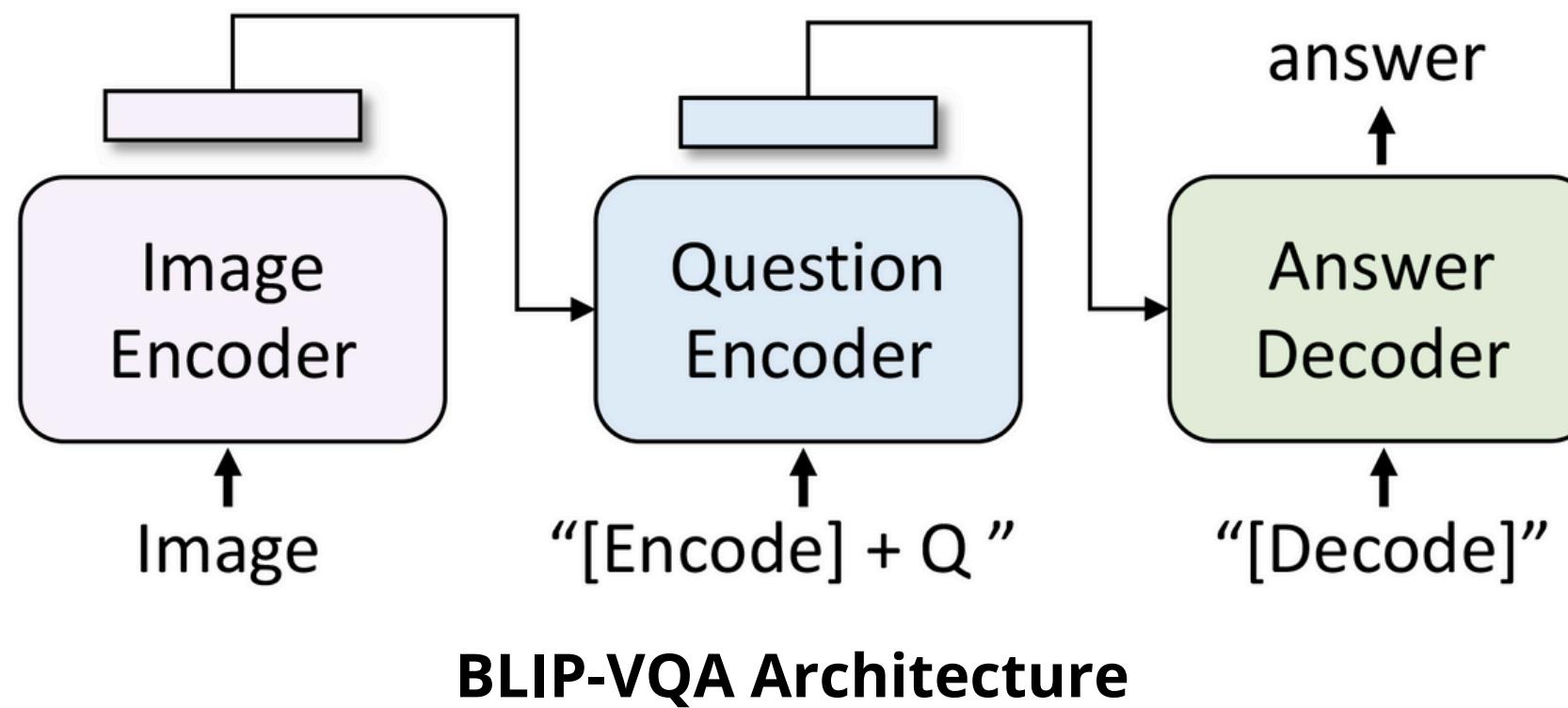
In VLMs, **answer (A)** - an open-ended sequence
(e.g., free text)

$$P(\hat{A}) = \prod_{t=1}^T P(a_t | a_{1:t-1}, I, Q)$$



Is this at a salt water beach or a lake?
- Salt water beach, Salt water, Lake, Beach

Vision language models (VLMs)

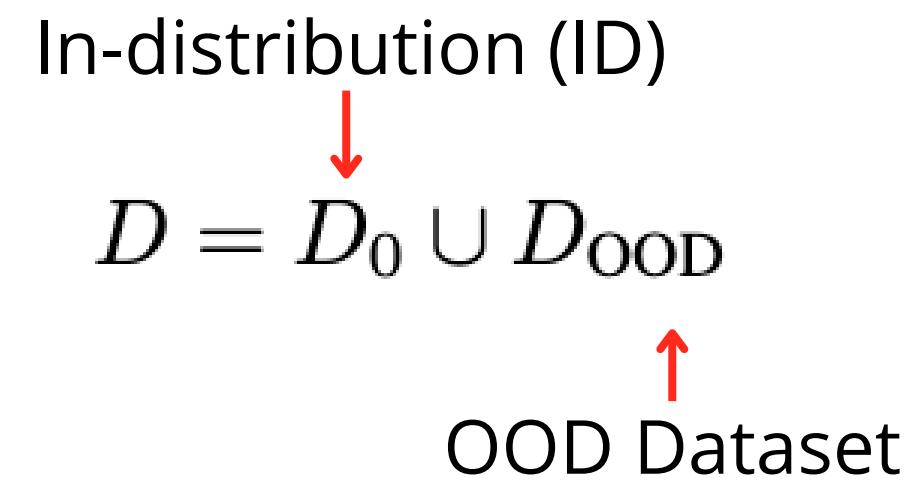


BLIP Pre-training Dataset

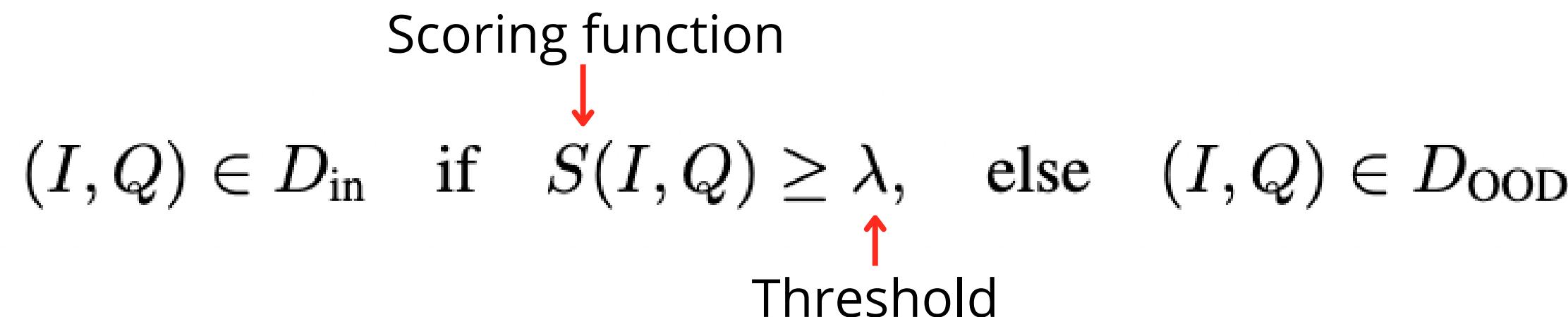
	COCO	VG	SBU	CC3M	CC12M	LAION
# image	113K	100K	860K	3M	10M	115M
# text	567K	769K	860K	3M	10M	115M

BLIP Fine-tuning: VQA V2 (83k/41k/81k images for training/validation/test)

Out-of-distribution (OOD) detection



OOD detection in VQA setting

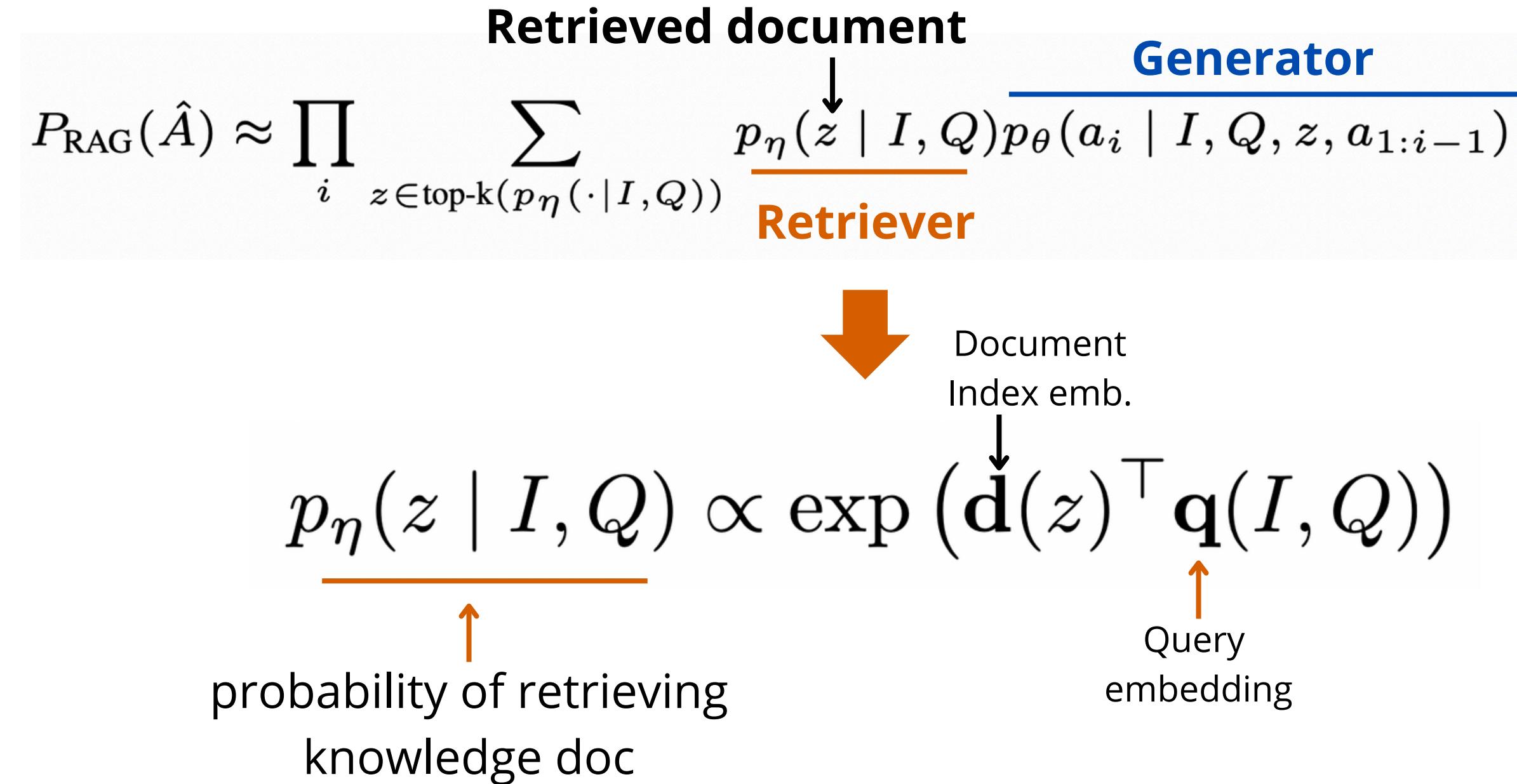


Zero-shot learning

- **Inputs:** Given an image (**I**) and a question (**Q**)
- **Goal:** To enable a model to generalize to unseen tasks or domains.

Example: A **BLIP-VQA** model, $f(I, Q)$, is trained on the **VQA V2 dataset** but will be evaluated on the **OK-VQA dataset** without task-specific fine-tuning on OK-VQA.

Retrieval-augmented generation (RAG)



Hallucination detection

Grounding Score: $g_{\text{mean}}(\hat{A}) = \frac{1}{n} \sum_{i=1}^n \frac{\mathbf{v}_{\text{pred}} \cdot \mathbf{v}_{\text{gt}}^i}{\|\mathbf{v}_{\text{pred}}\| \|\mathbf{v}_{\text{gt}}^i\|}$

Predicted answer
embedding
Ground Truth
embedding

Cosine Similarity

Hallucination if $g_{\text{mean}}(\hat{A}) < \tau$

↑
Threshold

Research questions

RQ1: How can **zero-shot learning** improve **retrieval** and **VQA** accuracy to address **hallucination** in multimodal **RAG** systems?

RQ2: How does zero-shot learning contribute to better **OOD performance** in VQA models?

Dataset: VQA V2



Is this person trying to hit a ball?
What is the person hitting the ball with?



What is the animal in the water?
How many people are present?

VQA V2

- Images: **MS-COCO**
- **1.1M** questions
- **11.1M** ground truth answers



Prof. Devi Parikh

Dataset: OK-VQA



What city is this?

Answer: Washington dc



What was the first movie was the character in this image first featured?

Answer: Star wars

Outside Knowledge VQA (OK-VQA)

- Images: **MS-COCO**
- **14,055** open-ended Qs
- **5 ground truth** ans per Qs

 ALLEN INSTITUTE
for ARTIFICIAL INTELLIGENCE

Carnegie
Mellon
University 

Prepare OK-VQA dataset in OOD setting

- **One: Vehicles and Transportation**
- **Two: Brands, Companies and Products**
- Three: Objects, Material and Clothing
 - For in domain setting
 - **For OOD setting**
- **Four: Sports and Recreation**
- Five: Cooking and Food
- Six: Geography, History, Language and Culture
- Seven: People and Everyday Life
- Eight: Plants and Animals
- **Nine: Science and Technology**
- **Ten: Weather and Climate**
- Other: Other

Roadmap

Part 1:

- Preliminaries: Zero-Shot Learning, VLMs, VQA in VLMs, RAG, OOD, Hallucination
- Research Questions
- Datasets (VQA V2 and Ok-VQA)

Part 2:

- **Method: Architecture, Loss Function**
- **Experimental Results (Qualitative, Quantitative, Visual results, and Ablation Study)**
- **Final Discussion (Contributions and Limitations)**



Overall RAG Pipeline



 **BLIP-VQA**

Step 1



Retrieval

Step 2

Wiki: Search-Based Retrieval
Summarization-Based Extraction
DBpedia: SPARQL Query

Retrieved knowledge
combines I-Q pair

Augmentation

Step 3

Evaluation

Step 5

Generation

Step 4

GPT-Neo 1.3B
(decoder-only)



Loss function

Binary cross-entropy loss

Predicted probability

$$\mathcal{L} = -\frac{1}{n} \sum_{i=1}^n [y_i \cdot \log(p_i) + (1 - y_i) \cdot \log(1 - p_i)]$$

Ground truth

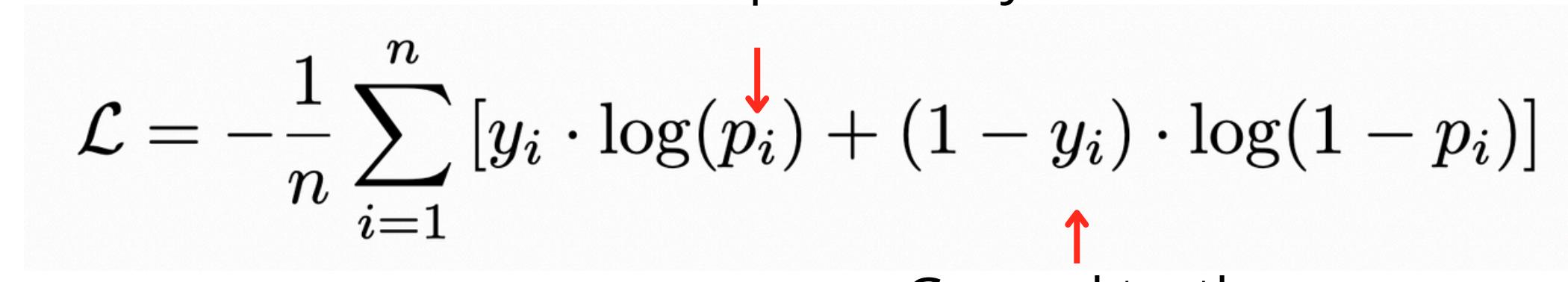
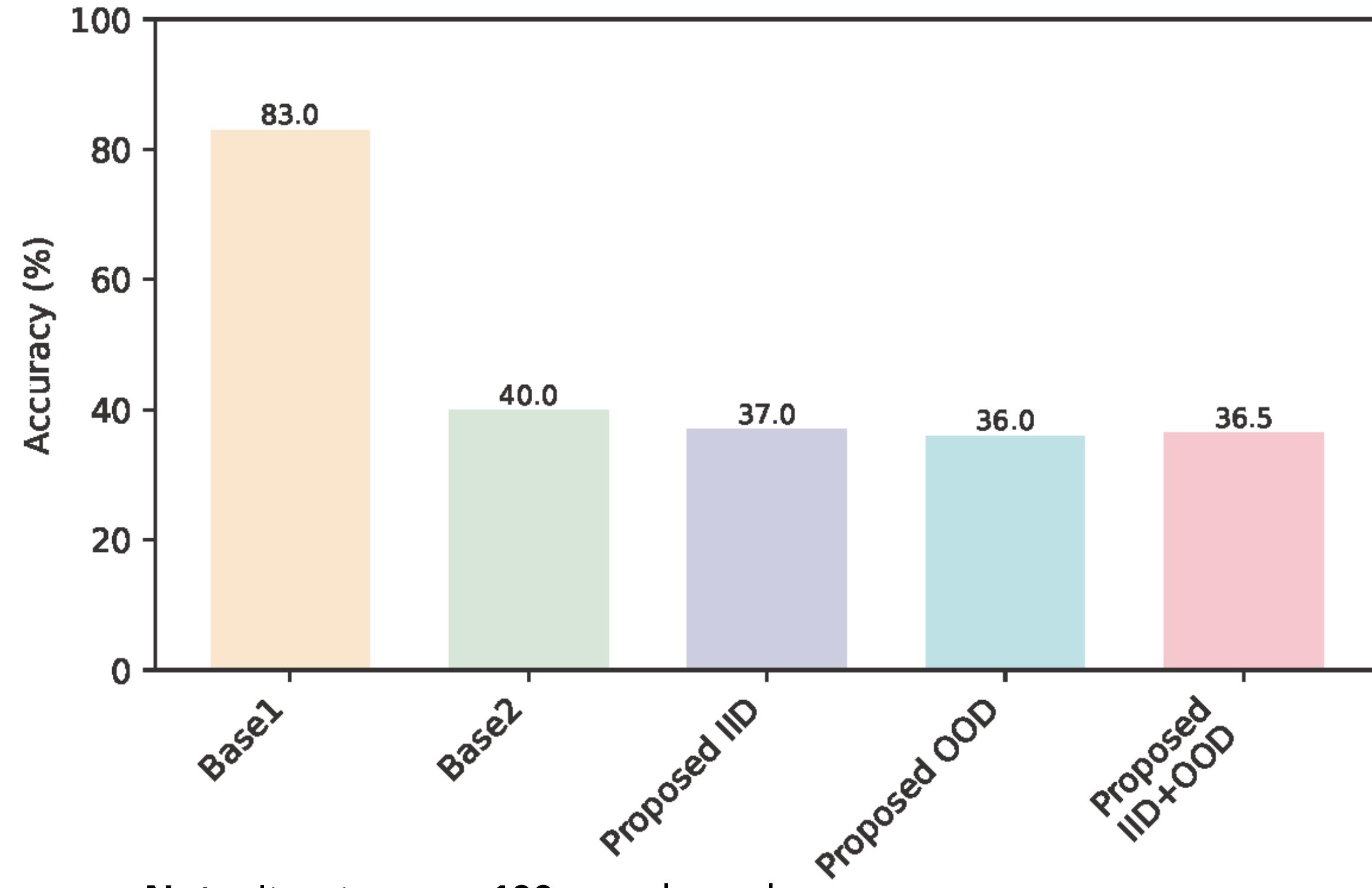




Image	Question, Prediction, GT	Accuracy (%)		Grounding Score (%)	
		Base	FilterRAG	Base	FilterRAG
	<p>What type of plane is that?</p> <p>Predicted Answer: commercial</p> <p>Ground Truth (GT) Answers: commercial, passenger, quanta, md 80</p>	40.0	36.5	71.70	70.37

Accuracy: baseline vs. FilterRAG

- **Base1:** BLIP VQA (model) + VQA V2 (Dataset)
- **Base2:** BLIP VQA (model) + Ok-VQA (Dataset)
- **Proposed:** BLIP VQA (model) + RAG + Ok-VQA (Dataset) + OOD

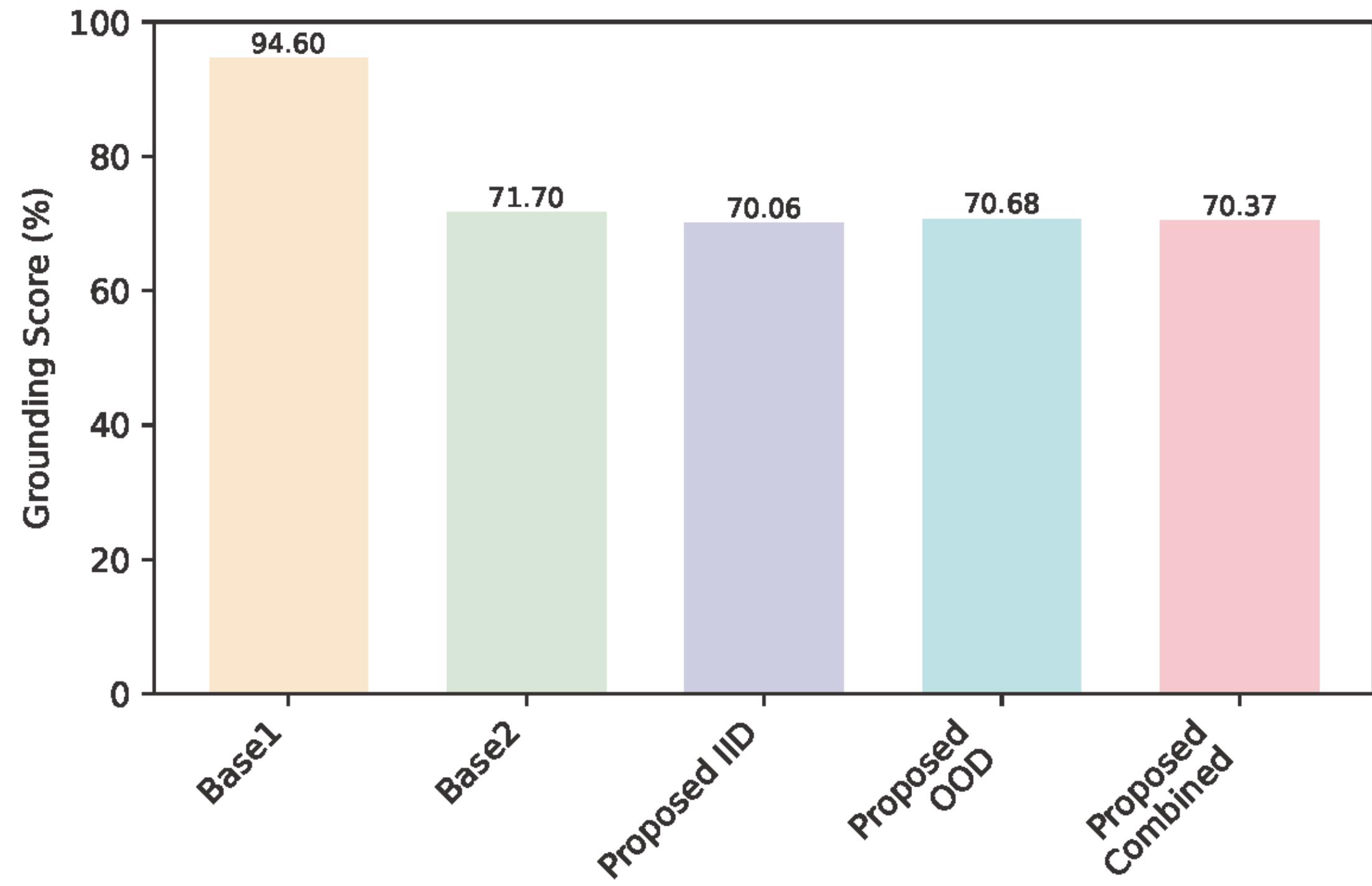


Accuracy: SOTA vs FilterRAG (OK-VQA)

Method	Knowledge Resources	Acc
BAN (Marino et al., 2019)	-	25.1
MUTAN (Marino et al., 2019)	-	26.41
KRISP (Marino et al., 2021)	Wikipedia+ConceptNet	38.35
MAVEx (Wu et al., 2022)	Wikipedia+ConceptNet+Google Images	39.4
KAT (Gui et al. 2022)	Wikidata+Frozen GPT-3 (175B)	54.41
FilterRAG (Proposed)	Wikidata + DBpedia + GPT-Neo 1.3B	36.5

Grounding score: baseline vs. FilterRAG

- **Base1:** BLIP VQA (model) + VQA V2 (Dataset)
- **Base2:** BLIP VQA (model) + Ok-VQA (Dataset)
- **Proposed:** BLIP VQA (model) + RAG + Ok-VQA (Dataset) + OOD



Note: Iterates over 100 samples only

Prediction visualization (ID case)



A center affixed unit like this one in a kitchen is called a what?

Predicted Answer: island

Ground Truth Answers: island



Is this at a salt water beach or a lake?

Predicted Answer: beach

Ground Truth Answers: salt water beach, salt water, lake, beach

Prediction visualization (OOD case)



What is the name of the board he is on?

Predicted Answer: surfboard

Ground Truth Answers: surf board, surfboard, surf

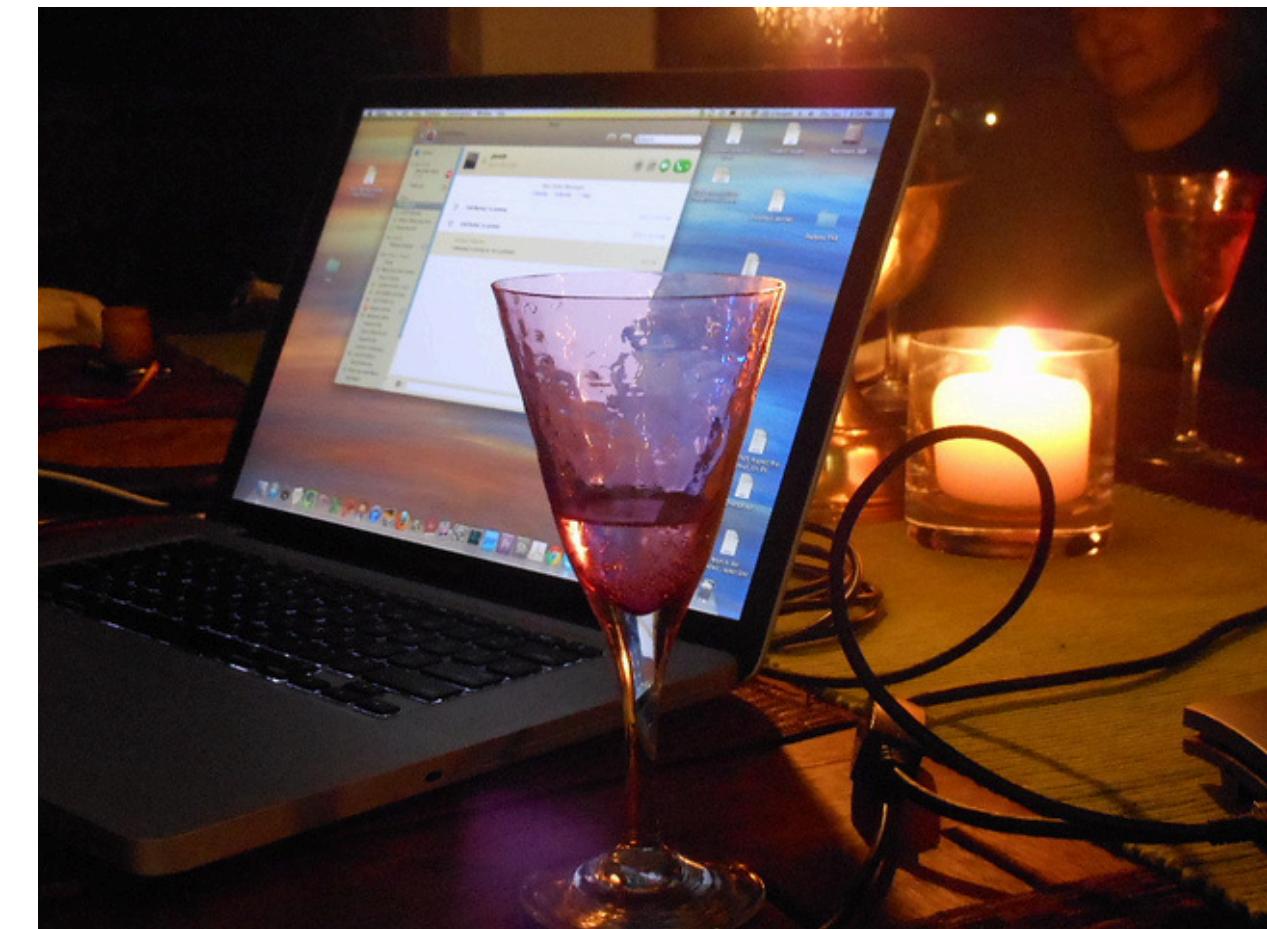


What type of plane is that?

Predicted Answer: commercial

Ground Truth Answers: commercial, passenger, quanta, md 80

Prediction visualization (OOD Case - failure)



What is this surf trick called?

Predicted Answer: **riding wave**

Ground Truth Answers: ride, tube ride, ollie, wave runner

Why is this plugged in?

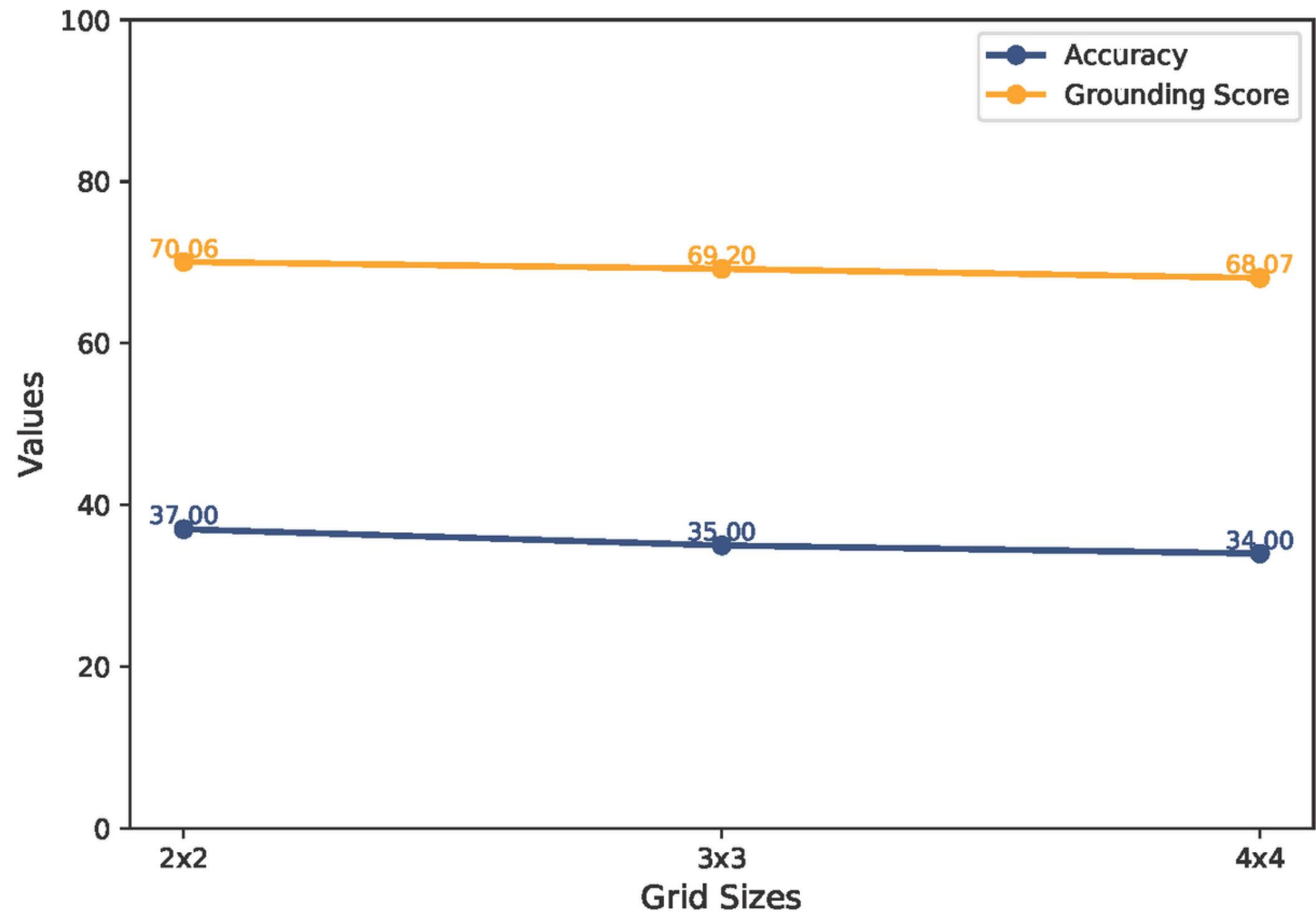
Predicted Answer: **plug**

Ground Truth Answers: charge, to have power and work, power, outlet

Ablation study



Grid Size: 2 x 2



Note: Iterates over **100** samples only

Final discussion

Contributions:

- Eliminates fine-tuning through zero-shot learning
- Uses external knowledge to address OOD cases beyond image-based reasoning
- Ensures reliable hallucination evaluation for VQA tasks

Limitations/Future works:

- **Optimize** generation modules (**LLM/VLM**) through **fine-tuning** for better outputs
- Explore OK-VQA **like datasets** for comprehensive OOD representation
- Use fine-tuning to create **synthetic** questions for **underrepresented OOD** cases

Thank You!

