# (a) Mixed Modality Search

🔍

**Query q**

Mountain Fuji

**Text Document $d_1$**

Reviewer #2: This is a very interesting paper! I really like it!

**Image Document $d_2$**

**Text and Image Document $d_3$**

Mount Rainier is a large, active stratovolcano..

# (b) Embedding Method

$d_3^T$  Mount Rainier is a large, active stratovolcano..  → CLIP Text Encoder → $\overrightarrow{e_3^T}$

Fusion Weight

$\alpha$

$\overrightarrow{e_3}$

$d_3^I$  → CLIP Image Encoder → $\overrightarrow{e_3^I}$

$1 - \alpha$

# (c) Modality Gap



Modality Gap

Text Embeddings

Fused Embeddings

Image Embeddings

$\overrightarrow{e_q}$

$\overrightarrow{e_1}$

$\overrightarrow{e_3}$

$\overrightarrow{e_2}$

UMAP 1

UMAP 2

# (d) Cosine Similarity across Modalities

| q | | Cosine Similarity | Rank | Ground Truth |
|---|---|---|---|---|
| | $d_1$ | 0.53 | 1 | 3 |
| | $d_2$ | 0.26 | 3 | 1 |
| | $d_3$ | 0.35 | 2 | 2 |



Cosine Similarity

Text-Text    Image-Image    Image-Text

# (e) Performance on MixBench



NDCG@10

GFLOPS

- 🟢 CLIP-B/16
- 🟡 CLIP-L/14
- 🔵 OpenCLIP-B/16
- 🟣 OpenCLIP-L/14
- ⚪ SigLIP-400m
- 🟩 GR-CLIP-B/16
- 🟧 GR-CLIP-L/14
- 🟦 GR-OpenCLIP-B/16
- 🟪 GR-OpenCLIP-L/14
- ⬛ GR-SigLIP-400m
- 🔺 VLM2Vec(Qwen)
- 🔺 VLM2Vec(LLaVA-Next)

0.63  0.60  0.66  0.63  0.67  0.63  0.62

+0.09  +0.11  +0.15  +0.15  +0.26

0.54  0.55  0.48  0.45  0.41