# Universal Vision-Language Dense Retrieval: Learning A Unified Representation Space for Multi-Modal Retrieval

Zhenghao Liu, Chenyan Xiong, Yuanhuiyi Lv, Zhiyuan Liu, and Ge Yu

ICLR

- Multi-modal retrieval focuses more on relevance modeling between queries and documents, single/cross modality matching, and modality routing.

## Cross-Modality Retrieval
## （Text/Image Matching）

**Query:** A woman wearing a net on her head cutting a cake.
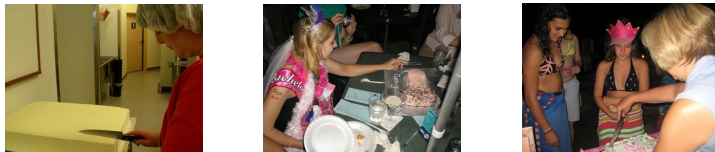**Retrieval Candidates:**

Image Retrieval ⇧

Text Retrieval ⇩

**Query:**

**Retrieval Candidates:**
**Text1:** A woman wearing a net on head cutting a cake.
**Text2:** A baker woman preparing bread dough on a tray with wax paper.

## Multi-Modal Retrieval
## （Relevance Modeling）

**Query:** What water-related object is sitting in front of the Torre del Reloj?
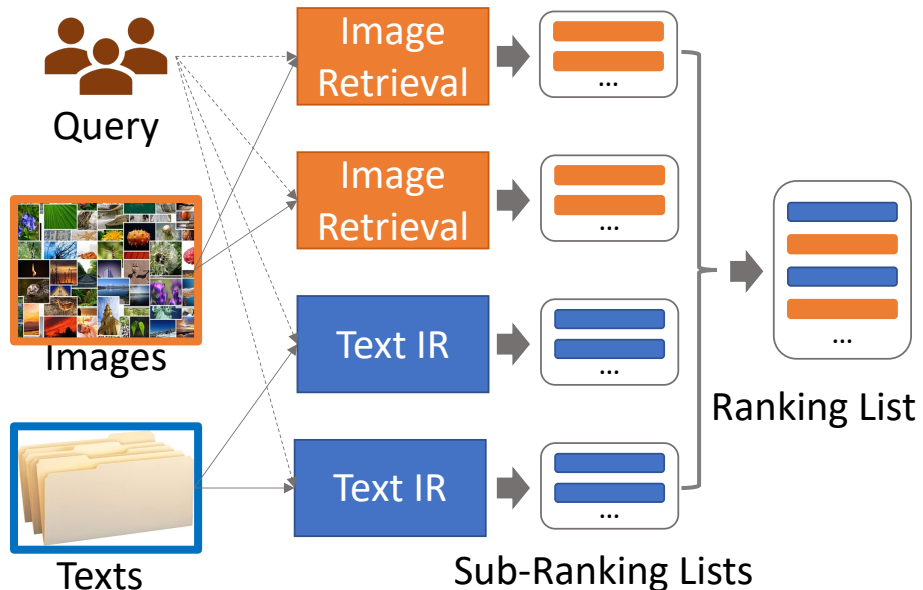**Retrieval Candidates:**

Image1    Image2    Image3

**Text1:** The Torre del Reloj Spanish is the main city gate of the historic center of Cartagena de Indias.
**Text2:** The Torre del Reloj is the clock tower, known as Arquillo Clock, and is one of the most emblematic buildings of Chiclana.
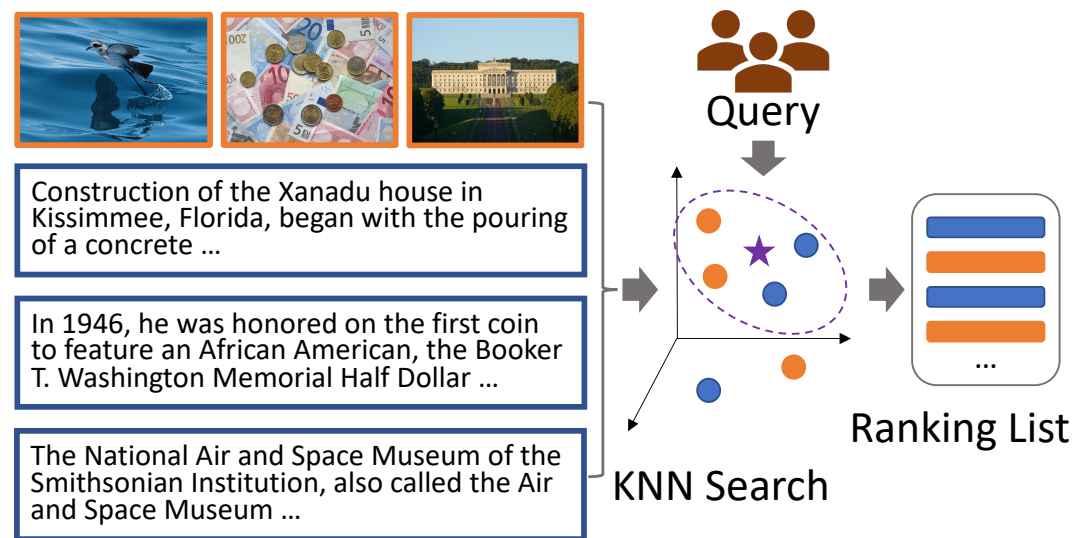**Text3:** Other landmarks in the city include the Torre del Reloj (Clock Tower).

- Divide-and-Conquer: Retrieving documents from multi-modalities and then fusion the retrieval results

- Universal Vision-Language Dense Retrieval: Leaning one universal embedding space for multi-modal retrieval, which unifies the relevance modeling and fusion stages
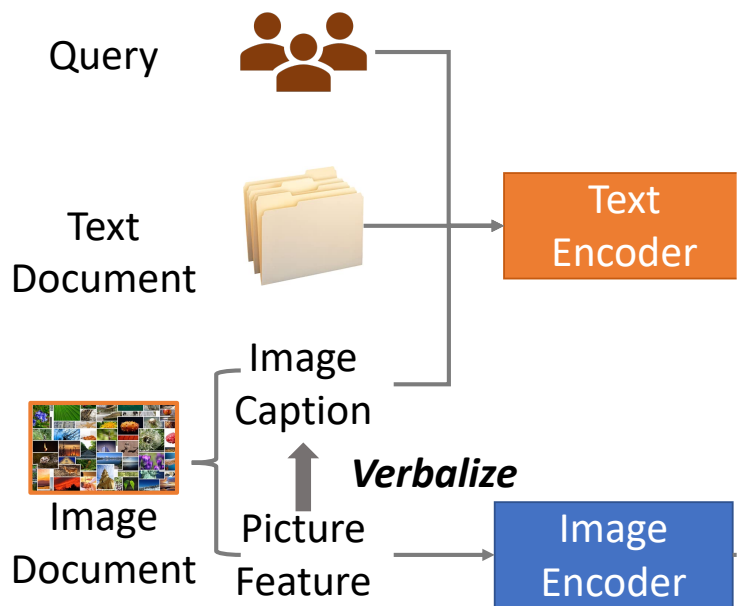


**Divide-and-Conquer Retrieval**

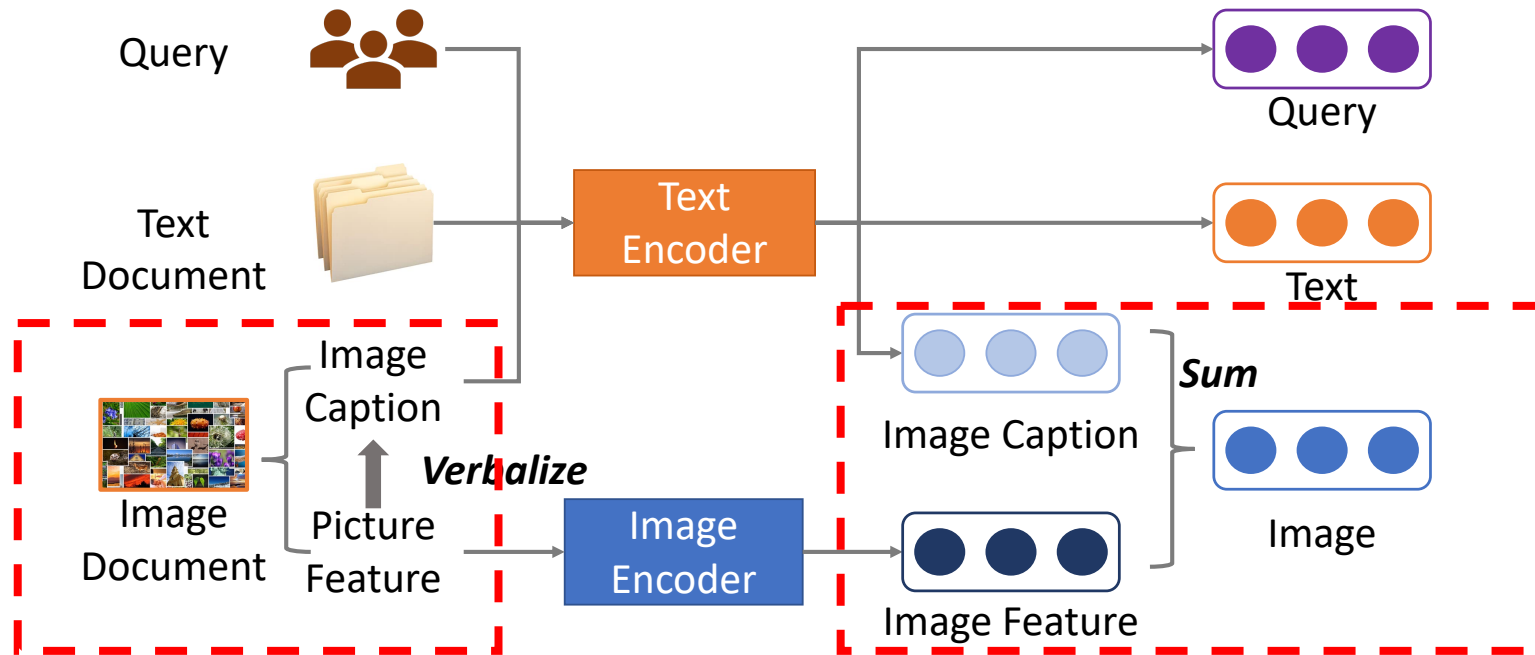**Universal Vision-Language Dense Retrieval**

- We use one Text Encoder to encode queries, text documents and image captions

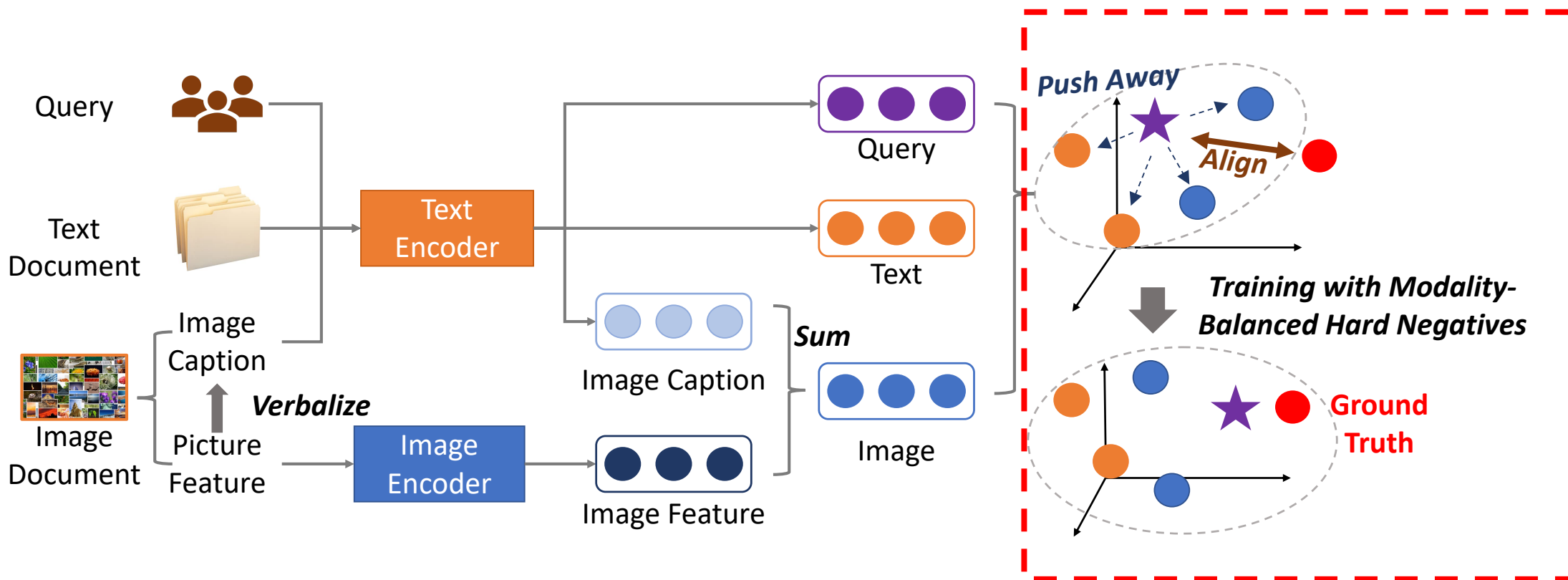- The Image Encoder is employed to encode image features as low-dimensional embeddings

Query

Text Document

Text Encoder

Image Caption

*Verbalize*

Image Document

Picture Feature

Image Encoder

- We enhance the representations of image documents by summing the representations of image captions and image features

- The image features are verbalized to enhance image documents in the raw text space
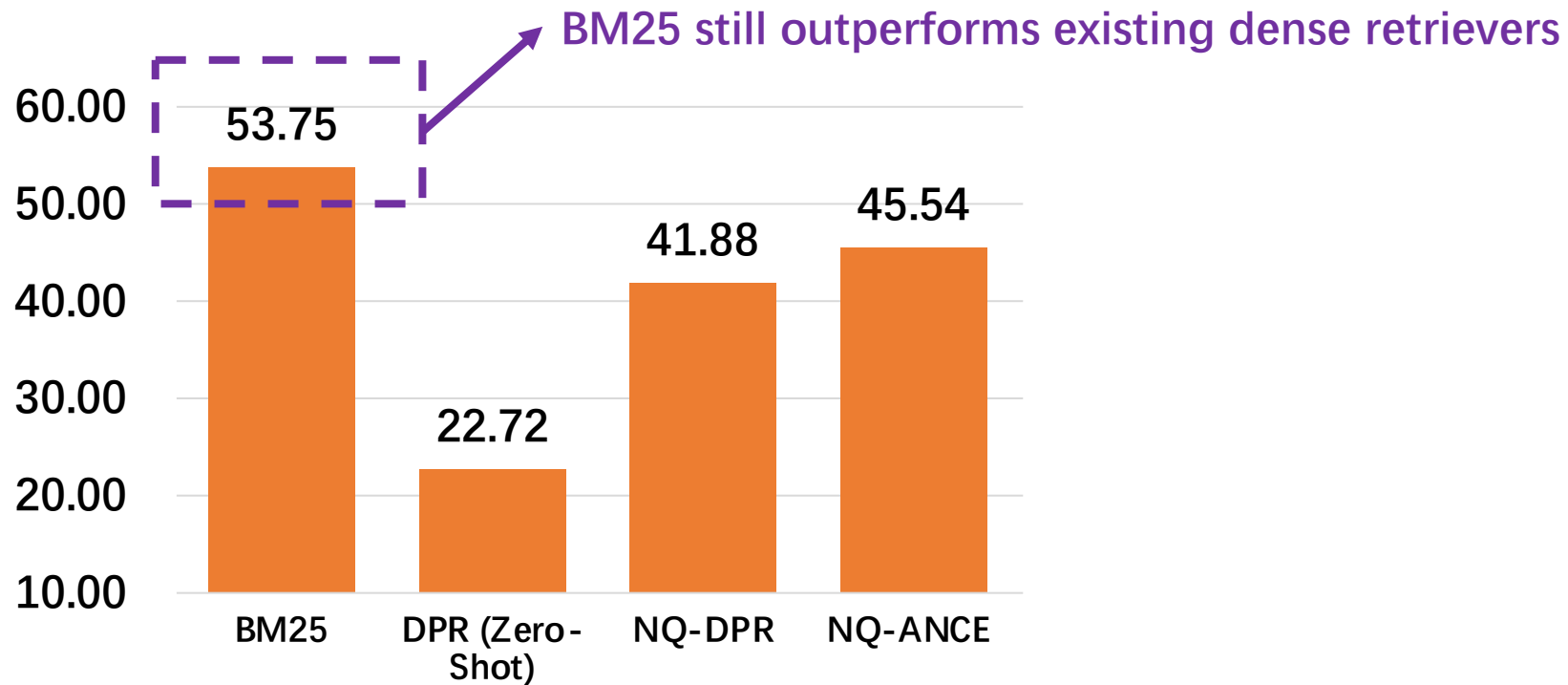
- The queries, text documents and image documents are mapped in one universal embedding space

- UniVL-DR designs a modality-balanced hard negative training strategies to train retrieval models

- We first conduct multi-modal retrieval using text retrieval models (***Multi-modal->Single modal***)

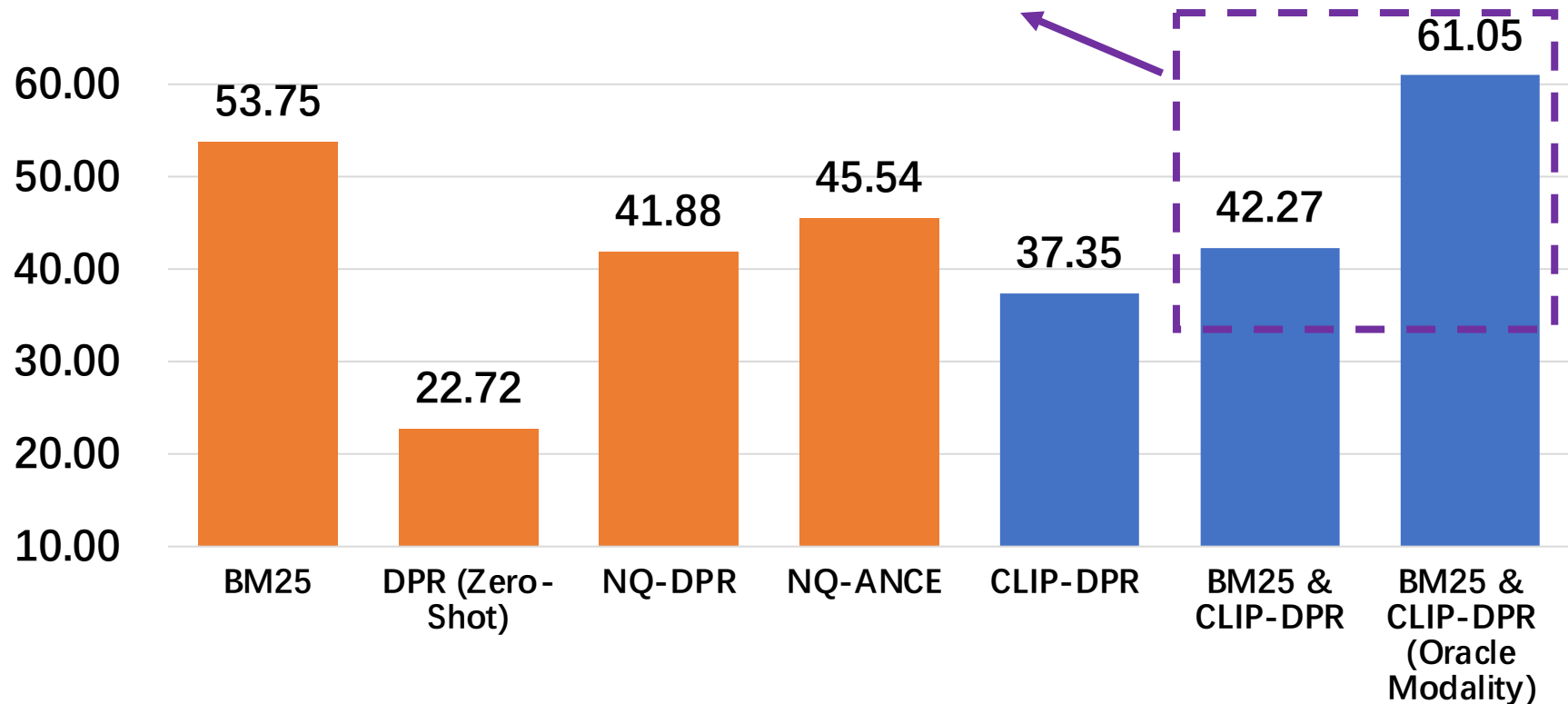  - Encoding image captions as the representations of image documents

**BM25 still outperforms existing dense retrievers**

- Multi-modal retrieval can be implemented by divide-and-conquer models (***Multi-modal->Single/Cross modality retrieval & Fusion***)

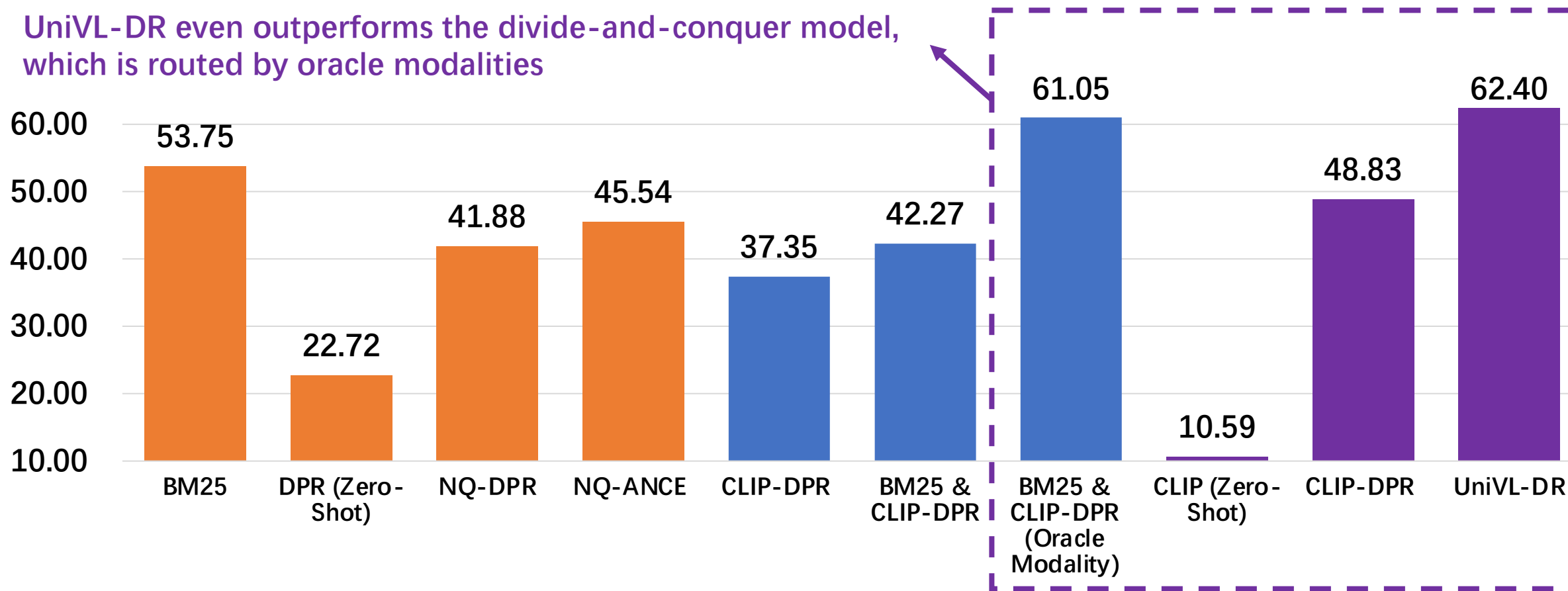**Modality routing still challenges the multi-modal retrieval tasks**

- Learn one universal embedding space for queries, text documents and image documents (***Multi-modal -> Universal Dense Retrieval***)

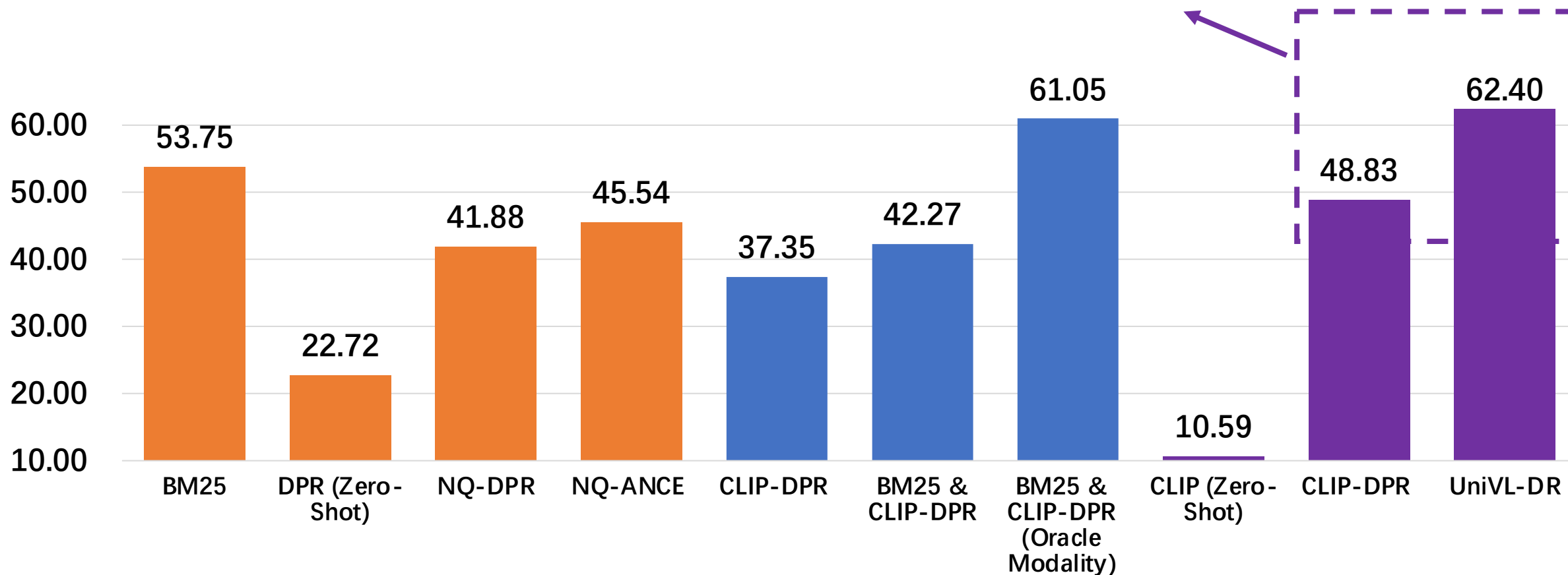UniVL-DR even outperforms the divide-and-conquer model, which is routed by oracle modalities



| Model | Value |
|---|---|
| BM25 | 53.75 |
| DPR (Zero-Shot) | 22.72 |
| NQ-DPR | 41.88 |
| NQ-ANCE | 45.54 |
| CLIP-DPR | 37.35 |
| BM25 & CLIP-DPR | 42.27 |
| BM25 & CLIP-DPR (Oracle Modality) | 61.05 |
| CLIP (Zero-Shot) | 10.59 |
| CLIP-DPR | 48.83 |
| UniVL-DR | 62.40 |

- Learn one universal embedding space for queries, text documents and image documents (***Multi-modal ->***

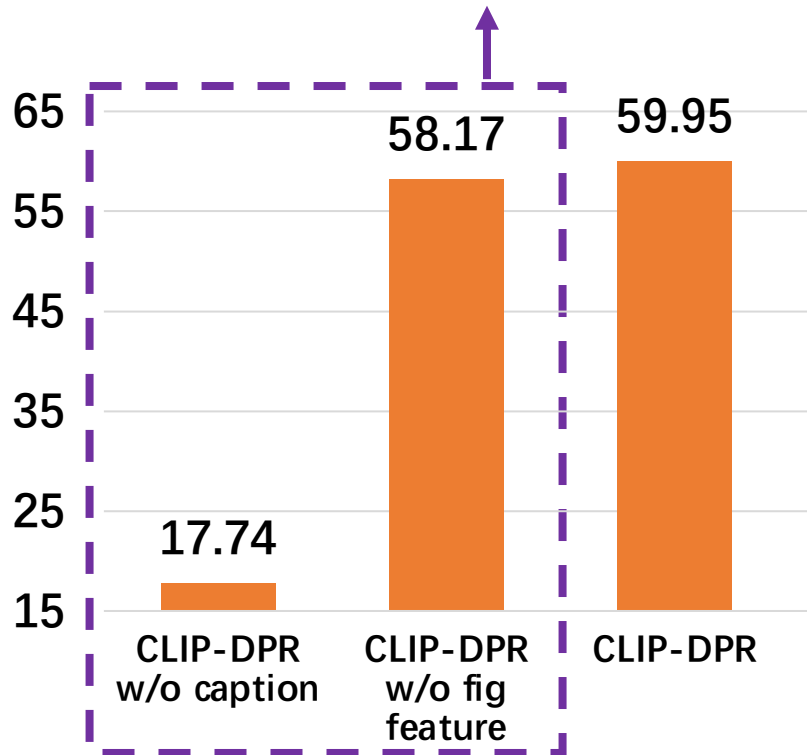  ***Universal Dense Retrieval***)

**UniVL-DR outperforms CLIP-DPR**



| | |
|---|---|
| BM25 | 53.75 |
| DPR (Zero-Shot) | 22.72 |
| NQ-DPR | 41.88 |
| NQ-ANCE | 45.54 |
| CLIP-DPR | 37.35 |
| BM25 & CLIP-DPR | 42.27 |
| BM25 & CLIP-DPR (Oracle Modality) | 61.05 |
| CLIP (Zero-Shot) | 10.59 |
| CLIP-DPR | 48.83 |
| UniVL-DR | 62.40 |

- We conduct experiments on the image retrieval task to show how to represent image documents

The image captions play a critical role in modeling relevance modeling between queries and image documents



**Query:** Does a Minnetonka Rhododendron flower have petals in a cup shape?
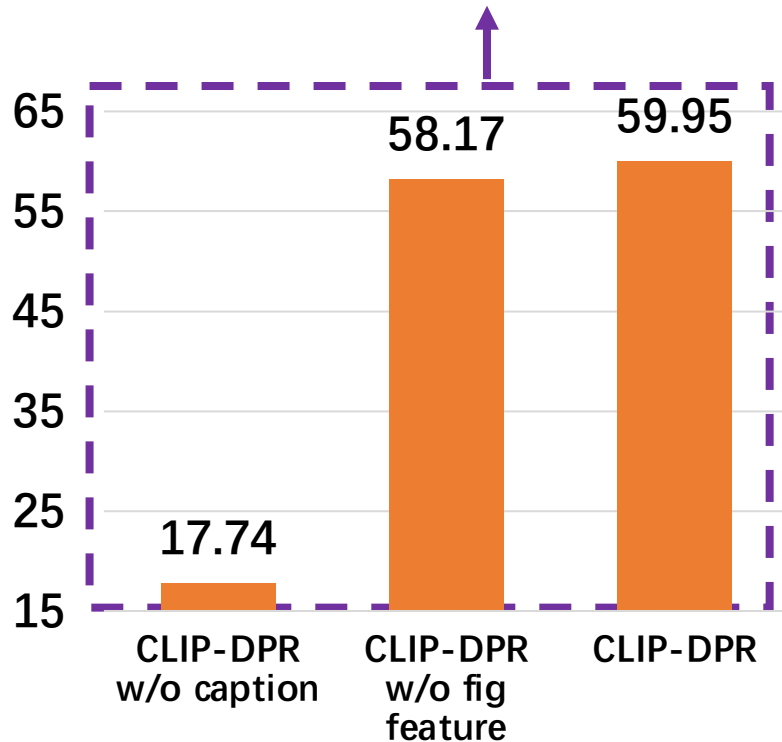
**Figure Features:**



**Image Captions:**
Minnetonka Rhododendron flower along Tranquility Court …

- We conduct experiments on the image retrieval task to show how to represent image documents

**The figure features can help better understand the semantics of image documents**



**Query:** Does a Minnetonka Rhododendron flower have petals in a cup shape?
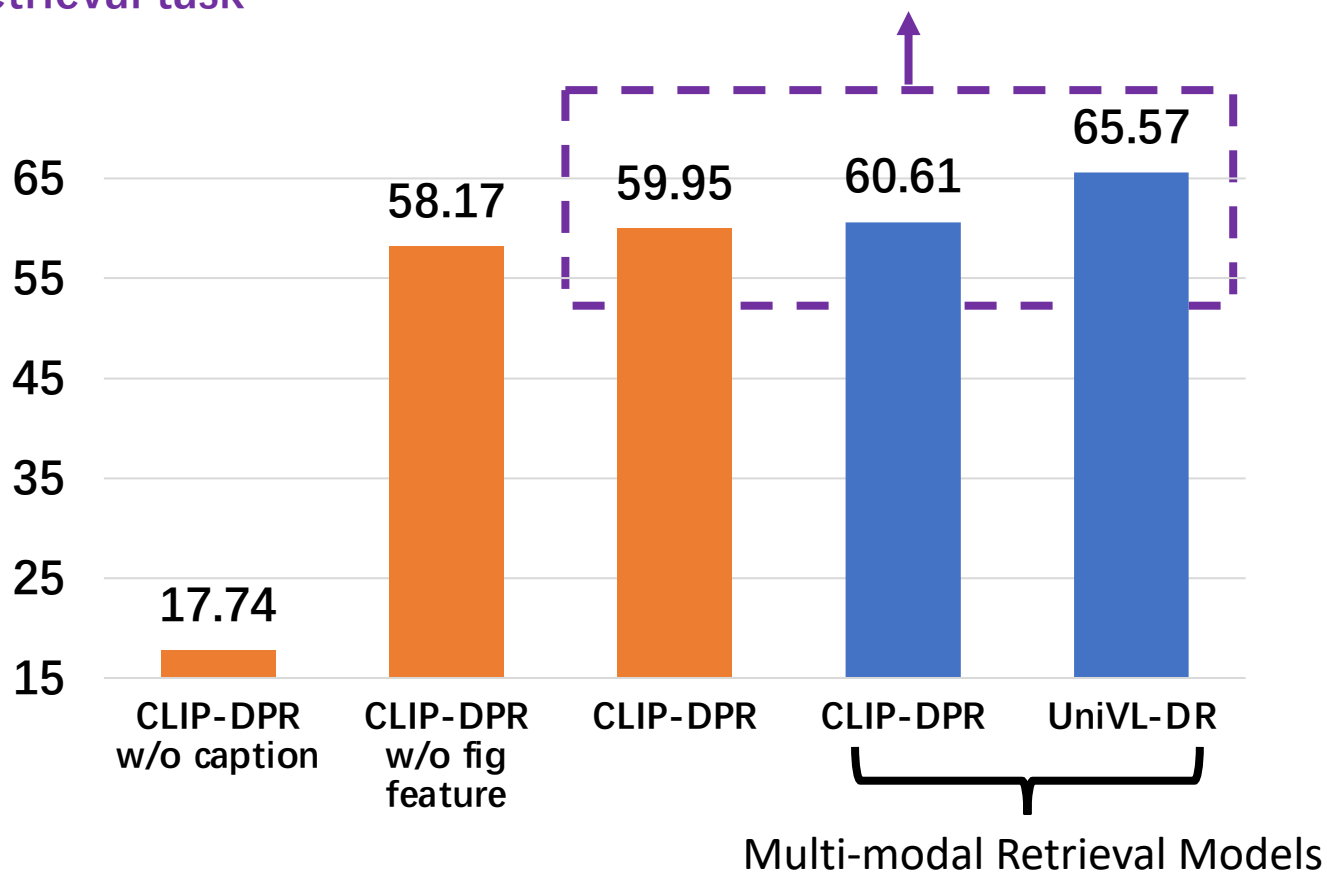
**Figure Features:**



**Image Captions:** Minnetonka Rhododendron flower along Tranquility Court …

- Then we further show the image retrieval performance of multi-modal retrieval models

**The text document retrieval tasks can also benefit the image retrieval task**



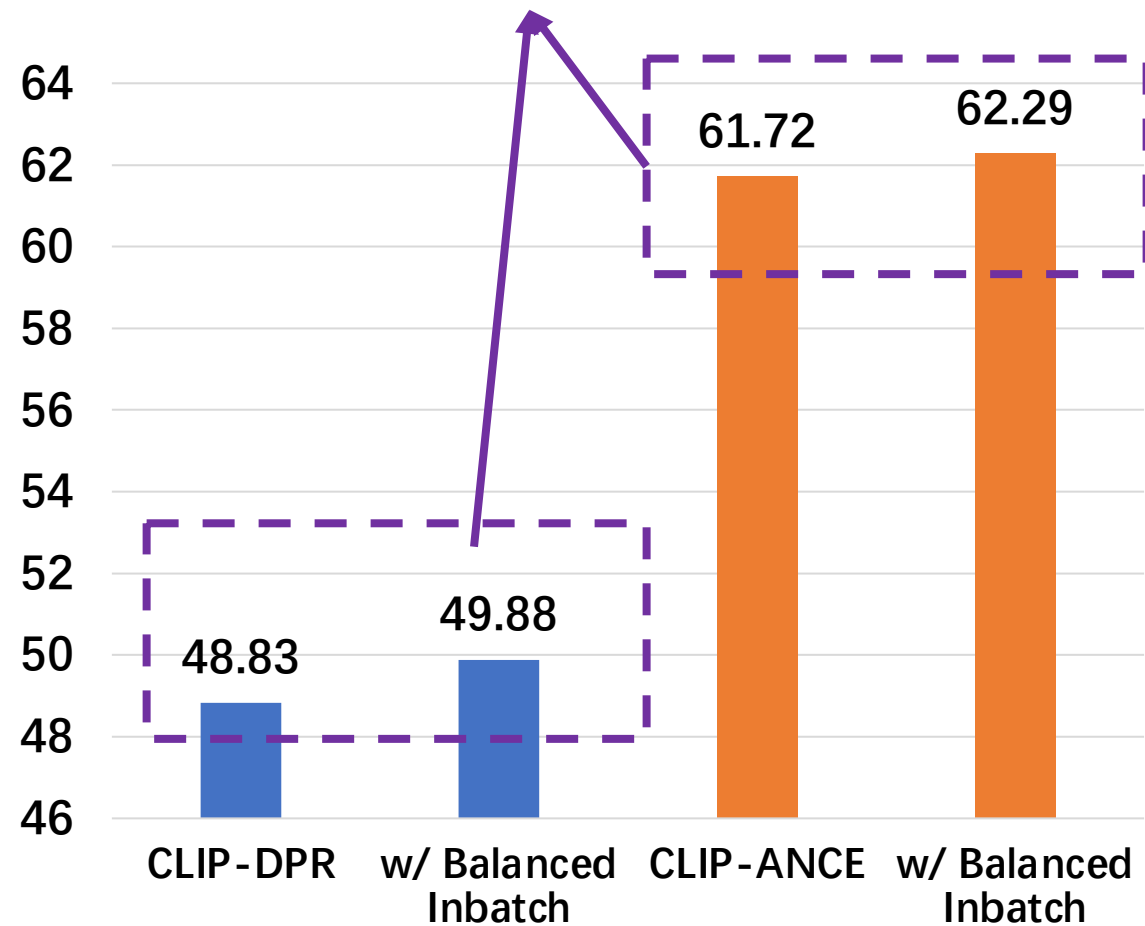**Query:** Does a Minnetonka Rhododendron flower have petals in a cup shape?

**Figure Features:**
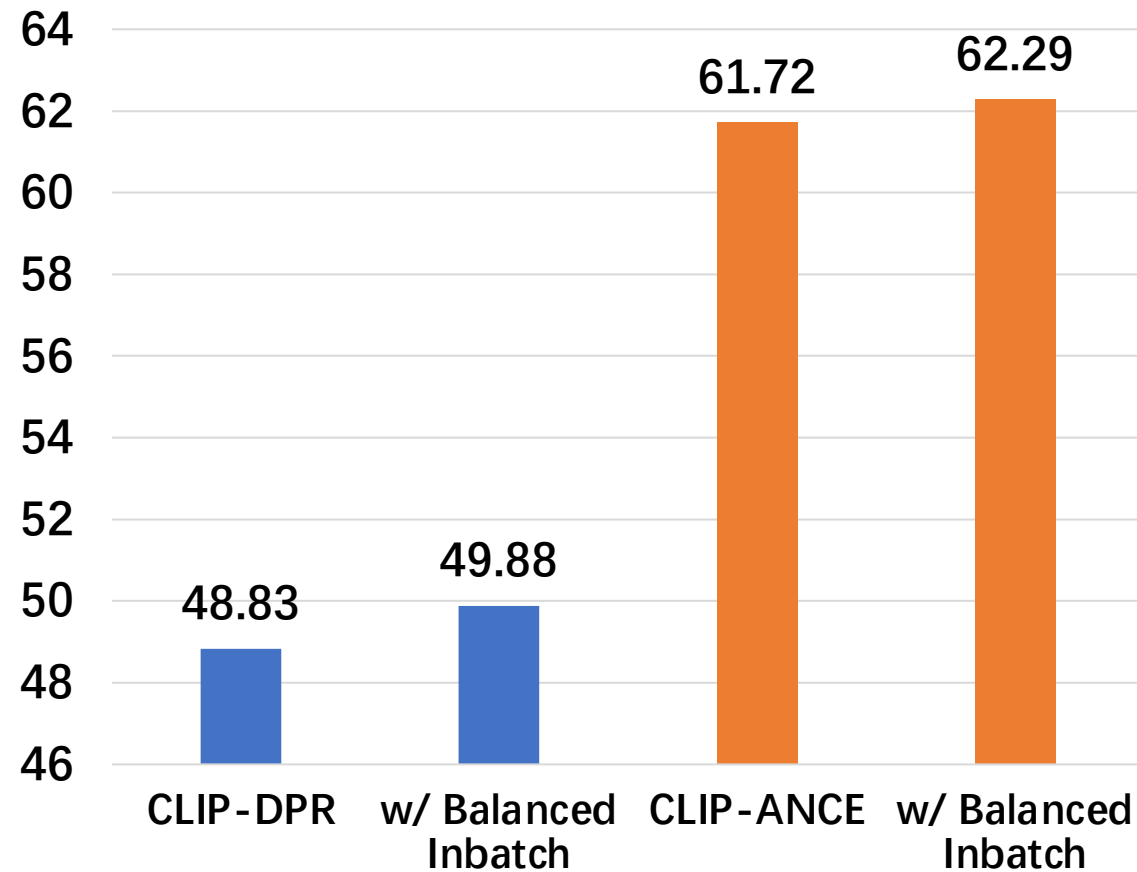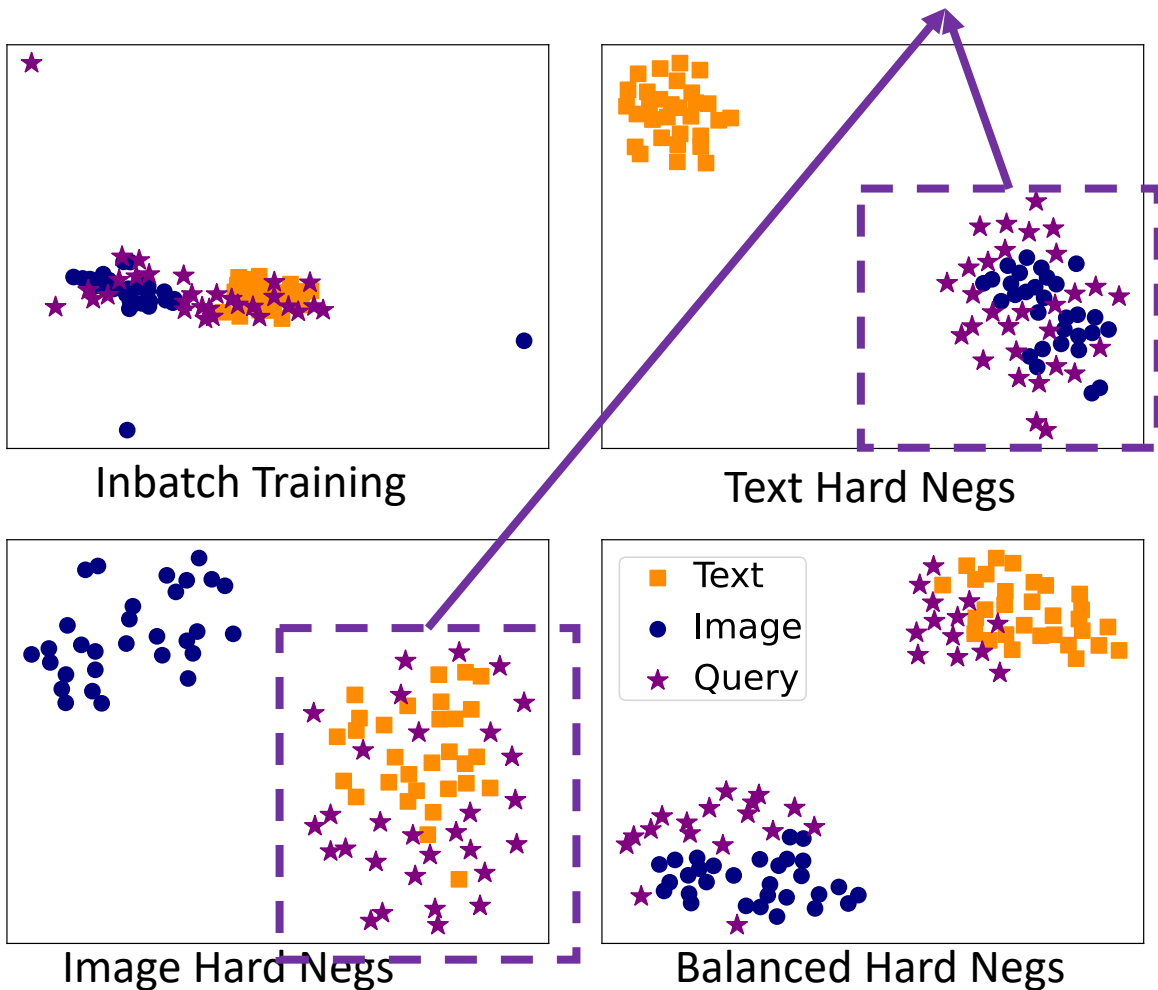


**Image Captions:**
Minnetonka Rhododendron flower along Tranquility Court …

Balancing the modality of negatives can improve the retrieval performance of inbatch and hard negative trained models

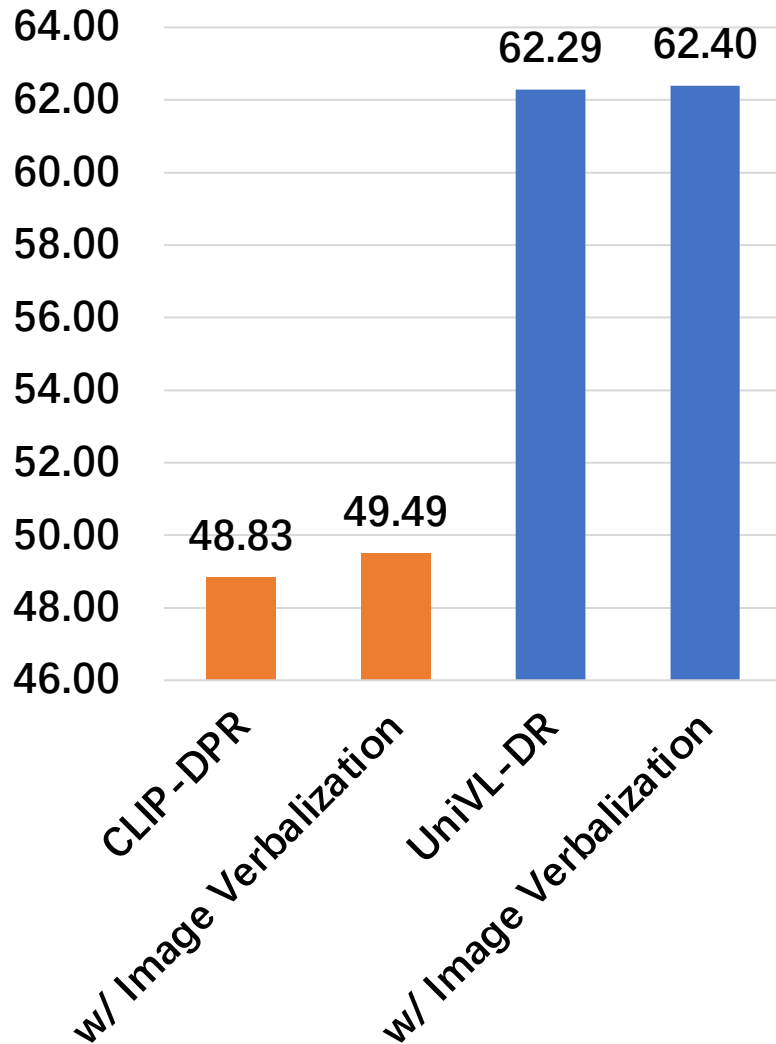**Balancing the modality of hard negatives can alleviate the modality preference**



Inbatch Training

Text Hard Negs

Image Hard Negs

Balanced Hard Negs

Legend:
- Text
- Image
- Query

Bar chart values:
- CLIP-DPR: 48.83
- w/ Balanced Inbatch: 49.88
- CLIP-ANCE: 61.72
- w/ Balanced Inbatch: 62.29

| Figures | Text |
|---|---|
| | **Query:** Does a Minnetonka Rhododendron flower have petals in a cup shape? |
| | Manual Caption: 2020-05-08 15 17 05 Minnetonka Rhododendron flower along Tranquility Court in the Franklin Farm section of Oak Hill, Fairfax County, Virginia Minnetonka Rhododendron flower along Tranquility Court in the Franklin Farm section of Oak Hill, Fairfax County, Virginia |
| | Verbalized Caption: a purple flower with green leaves and purple flowers |
| | Verbalized Query: what shape are the petals of the minnetonka rhododendron flower? |
| | **Query:** Are the heads of Iranian women covered in traditional clothing? |
| | Manual Caption: Iranian family, gathered together wearing traditional clothes - Nishapur - Nowruz2014 Iranian family, gathered together wearing traditional clothes |
| | Verbalized Caption: a group of people in costumes standing in a park. |
| | Verbalized Query: how many people are wearing hats in the group of iranian family members? |
| | **Query:** At the 1928 Amsterdam Olympics, what is the maximum number of buttons that you can get on the Egyptian men's uniform? |
| | Manual Caption: Egyptische atleten bij OS Amsterdam 1928 - Egyptian Olympic athletes, Amsterdam 1928 (6941436605) http://www.spaarnestadphoto.nl/component/option,com memorix ... |
| | Verbalized Caption: a group of men in suits and hats standing in a field |
| | Verbalized Query: did all the men in the egyptian olympic athletes wear the same type of caps? |

**These verbalized queries usually contain matched entities and bridge the modality gap between images and texts**
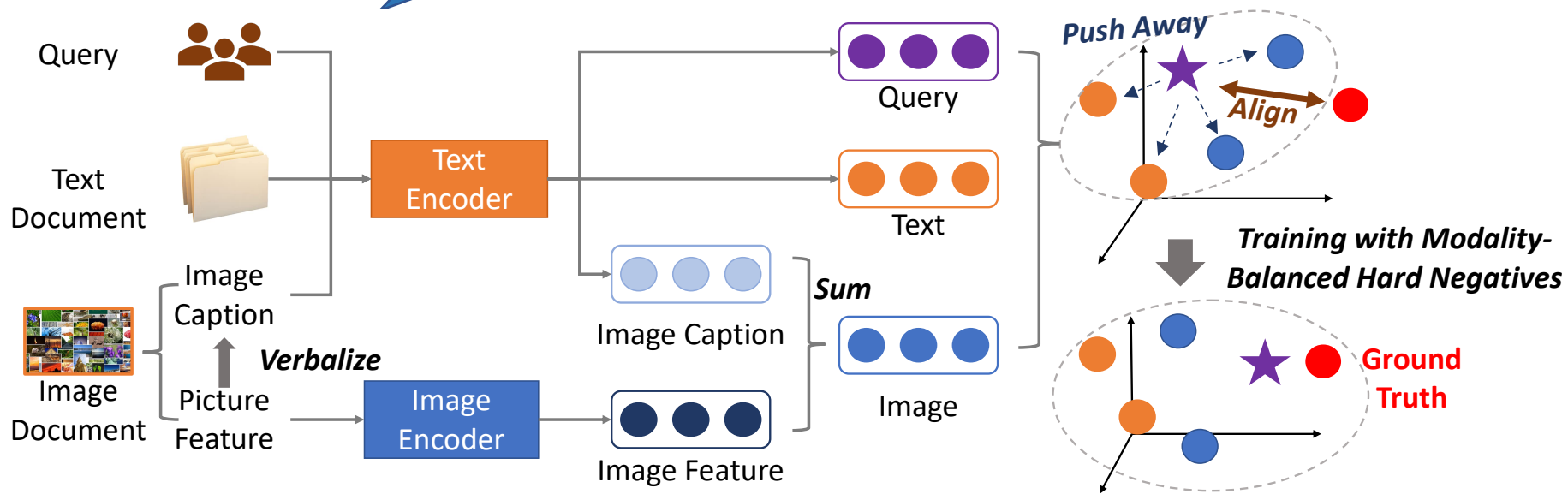
All codes and data are available at **https://github.com/OpenMatch/UniVL-DR**.

# Thank you!

NEU CSE