



---

# QUEST: Quadruple Multimodal Contrastive Learning with Constraints and Self-Penalization

---

Qi Song<sup>1\*</sup>, Tianxiang Gong<sup>2\*</sup>, Shiqi Gao<sup>2</sup>,  
Haoyi Zhou<sup>1,3†</sup>, Jianxin Li<sup>2,3</sup>

<sup>1</sup>School of Software, Beihang University

<sup>2</sup>School of Computer Science and Engineering, Beihang University

<sup>3</sup>Zhongguancun Laboratory, Beijing

{songqi23, gongtx, gaoshiqi, haoyi, lijx}@buaa.edu

---

\*Equal contribution

†Corresponding author.

## Background & Motivation

- Contrastive learning treats all negative samples equally, ignoring the potential semantic relationships between negative samples and the anchor.
- Contrastive learning often neglects significant portions of input information, leading to feature suppression and shortcut learning.



**Caption1:** A balding man wearing a red life jacket is sitting in a small boat.

**Caption2:** A man in a green shirt and red life jacket is sitting in a canoe drifting around the lake.

**Shared information**

A man, red life jacket, sitting, boat(canoe)

**Unique information**

**Caption1**  
balding (man), small (boat)

**Caption2**  
a green shirt, drifting, the lake

Figure1. Shared and unique Information in Multimodal Multi-view Scenario.

## Background & Motivation

- Current contrastive learning methods focus on maximizing mutual information between two views while ignoring unique information.

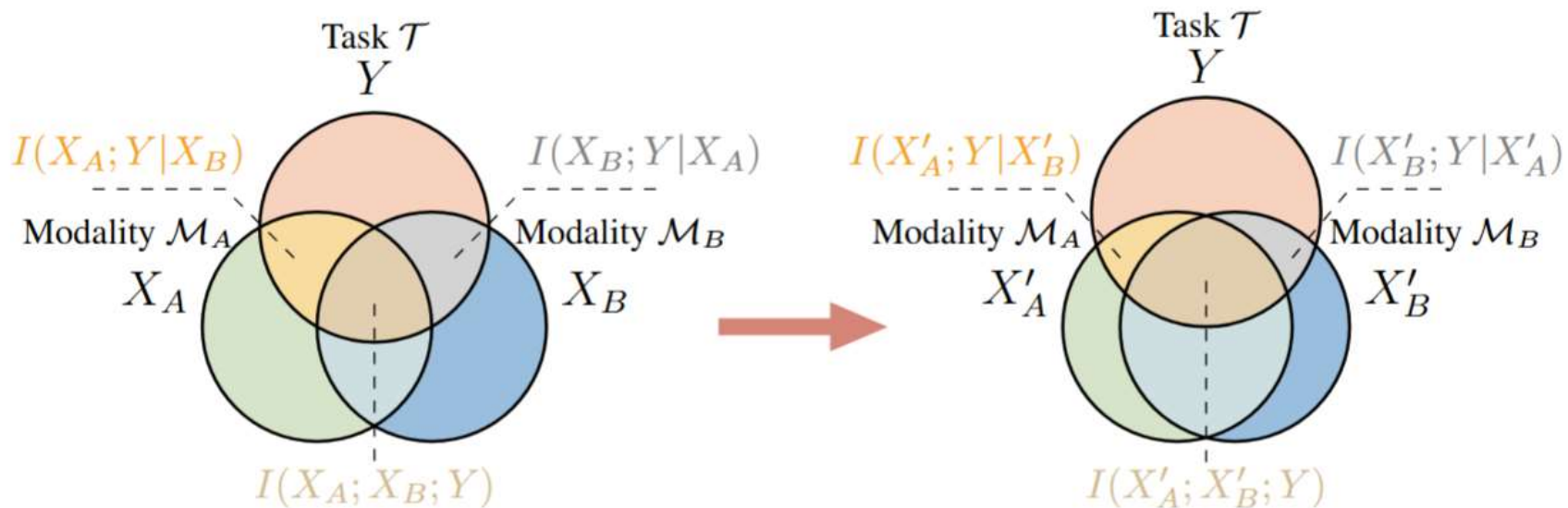


Figure 2: Feature suppression in multi-view contrastive learning. We define  $I(X_A; X_B; Y)$  as task-related shared information,  $I(X_A; Y|X_B)$  and  $I(X_B; Y|X_A)$  as task-related unique information related to task  $Y$  in modalities  $X_A$  and  $X_B$ , respectively. Contrastive losses, such as InfoNCE, tend to maximize the task-related shared information while suppressing the task-related unique information in each modality. Left: before training with InfoNCE. Right: after training with InfoNCE.

## Architecture Overview

- Use extra unique decoder to capture unique information simultaneously with constraints and self-penalization.

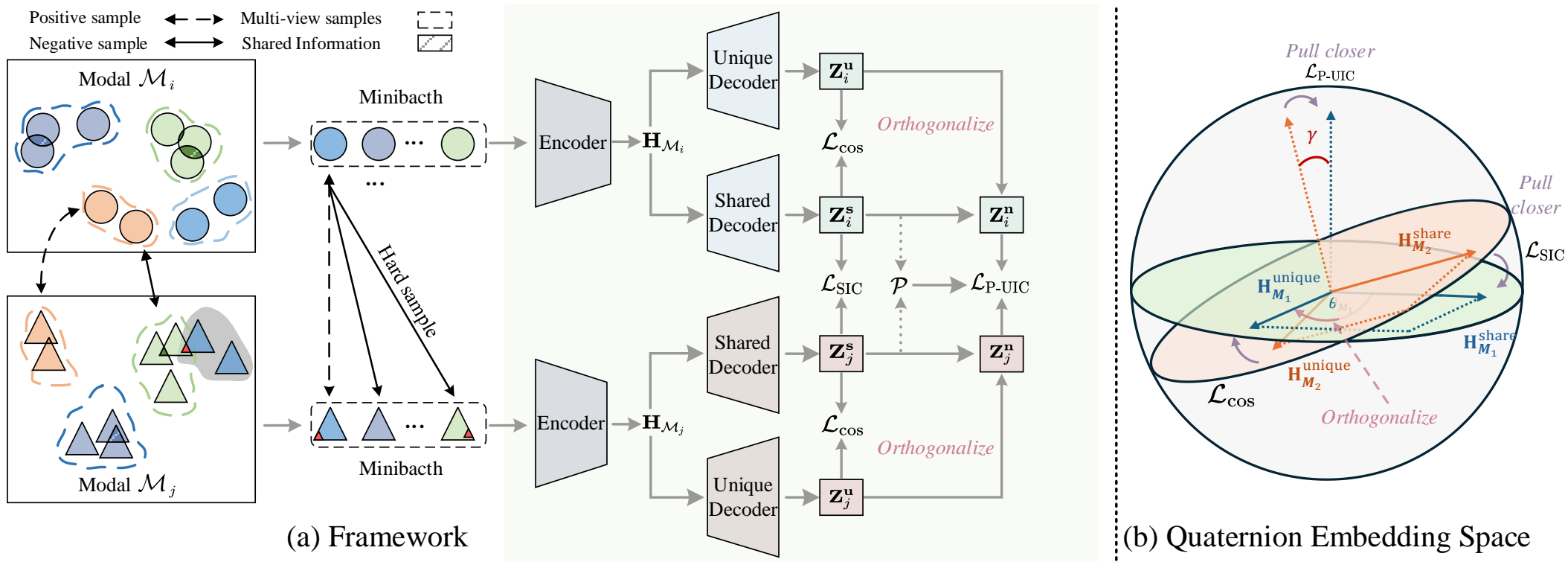


Figure 3. Left: An overview of our method's framework. Right: The quaternion embedding space, in which we orthogonalize the representation of high-dimensional data to extend the solution space.

## ■ Shared Decoder and Unique Decoder

- For each modality, input data  $\mathbf{X}_i$  is transformed by a modality-specific encoder  $\mathcal{F}_{\mathcal{M}_i}(\cdot)$  into a general representation  $\mathbf{H}_{\mathcal{M}_i}$ . Two decoders  $\mathcal{G}_{\mathcal{M}_i}^s(\cdot)$  and  $\mathcal{G}_{\mathcal{M}_i}^u(\cdot)$  then separate shared and unique information from  $\mathbf{H}_{\mathcal{M}_i}$ .

$$\begin{aligned}\mathbf{Z}_i^u &= \mathcal{G}_{\mathcal{M}_i}^u(\mathbf{H}_{\mathcal{M}_i}; \Phi_i) = \mathcal{G}_{\mathcal{M}_i}^u(\mathcal{F}_{\mathcal{M}_i}(\mathbf{X}_i; \Theta_i); \Phi_i), \\ \mathbf{Z}_i^s &= \mathcal{G}_{\mathcal{M}_i}^s(\mathbf{H}_{\mathcal{M}_i}; \Psi_i) = \mathcal{G}_{\mathcal{M}_i}^s(\mathcal{F}_{\mathcal{M}_i}(\mathbf{X}_i; \Theta_i); \Psi_i).\end{aligned}$$

## ■ Shared Information Constraint

- In multimodal and multi-view scenarios, maximizing the lower bound of mutual information (MI) between representations from different views encourages the shared decoder to learn task-related agreement.

$$\mathcal{L}_{\text{SIC}} = \sum_{i,j} \mathbf{1}_{\mathcal{M}_i \neq \mathcal{M}_j} \mathbb{E}_{\mathbf{Z}_i^s} \left[ -\log \frac{\exp(s(\mathbf{Z}_i^s, \mathbf{Z}_j^{s+})/\tau)}{\exp(s(\mathbf{Z}_i^s, \mathbf{Z}_j^{s+})/\tau) + \sum_{k=1}^m \mathbf{1}_{\hat{y}^-} \exp(s(\mathbf{Z}_i^s, \mathbf{Z}_{jk}^{s-})/\tau)} \right].$$



## ■ Unique Information Constraint

- In contrast to shared information, unique information is modality-specific and task-relevant, providing essential insights for downstream tasks.
- Firstly, we derive the representation space of normal vectors for shared and unique embedding spaces through cross-product calculations:

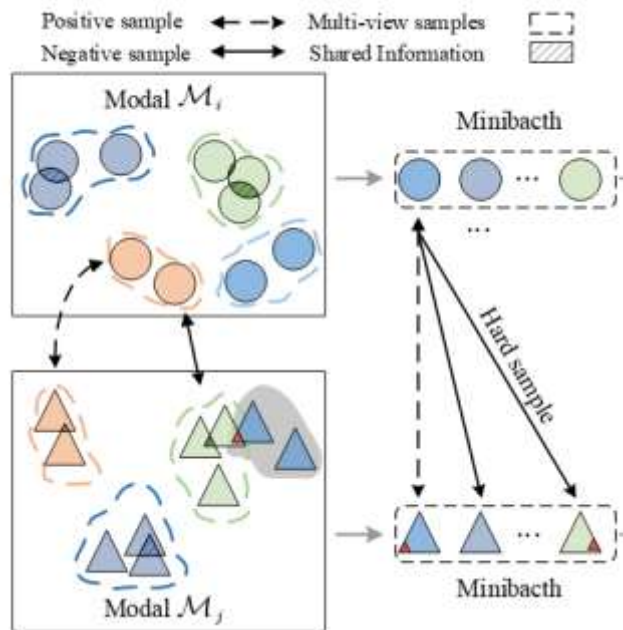
$$\mathbf{z}_i^n = \mathbf{z}_i^s \times \mathbf{z}_i^u$$

- In the newly projected space, our objectives aim to maximize the alignment of unique representation from different modalities within the plane spanned by the shared representation

$$\mathcal{L}_{\text{UIC}} = \sum_{i,j} \mathbb{1}_{\mathcal{M}_i \neq \mathcal{M}_j} \mathbb{E}_{\mathbf{z}_i^n} \left[ -\log \frac{\exp(s(\mathbf{z}_i^n, \mathbf{z}_j^{n+})/\tau)}{\exp(s(\mathbf{z}_i^n, \mathbf{z}_j^{n+})/\tau) + \sum_{k=1}^m \mathbb{1}_{\hat{y}^-} \exp(s(\mathbf{z}_i^n, \mathbf{z}_{jk}^{n-})/\tau)} \right] + \sum_i \sum_j \frac{\mathbf{z}_{ij}^s \cdot \mathbf{z}_{ij}^u}{\|\mathbf{z}_{ij}^s\| \|\mathbf{z}_{ij}^u\|}$$

## Self-Penalization Constraint

- Challenges in practical Contrastive learning implementations
  - Uniform treatment of all other samples within a batch B as negative examples
  - Misclassification of semantically similar samples as negative
- Use the intra-model shared information similar to penalization to term guide the optimization of unique information



Solution

### Self-Penalization Term

$$\mathcal{P} = \exp(\lambda[\mathbf{S} - \text{diag}(\mathbf{S}) + \mathbf{I}]), \quad \text{where } \mathbf{S} = \mathbf{Z}_i^s \mathbf{Z}_j^{sT}$$

### Penalized UIC

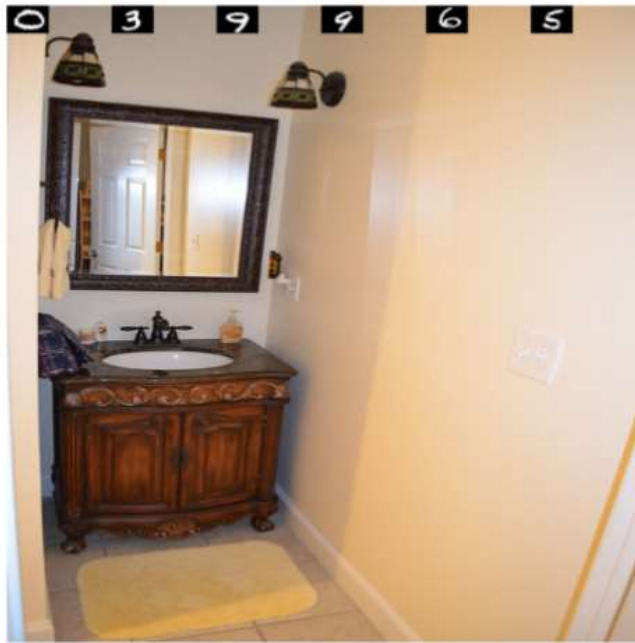
$$\mathcal{L}_{\text{P-UIC}} = \sum_{i,j} \mathbb{1}_{\mathcal{M}_i \neq \mathcal{M}_j} \mathbb{E}_{\mathbf{z}_i^n} \left[ -\log \frac{\exp(\mathcal{P}^+ \cdot s(\mathbf{z}_i^n, \mathbf{z}_j^{n+})/\tau)}{\exp(\mathcal{P}^+ \cdot s(\mathbf{z}_i^n, \mathbf{z}_j^{n+})/\tau) + \sum_{k=1}^m \mathbb{1}_{\hat{y}^-} \exp(\mathcal{P} \cdot s(\mathbf{z}_i^n, \mathbf{z}_{jk}^{n-})/\tau)} \right] + \sum_i \sum_j \frac{\mathbf{z}_{ij}^s \cdot \mathbf{z}_{ij}^u}{\|\mathbf{z}_{ij}^s\| \|\mathbf{z}_{ij}^u\|}$$

Figure 3(a). Potential correlations among data.

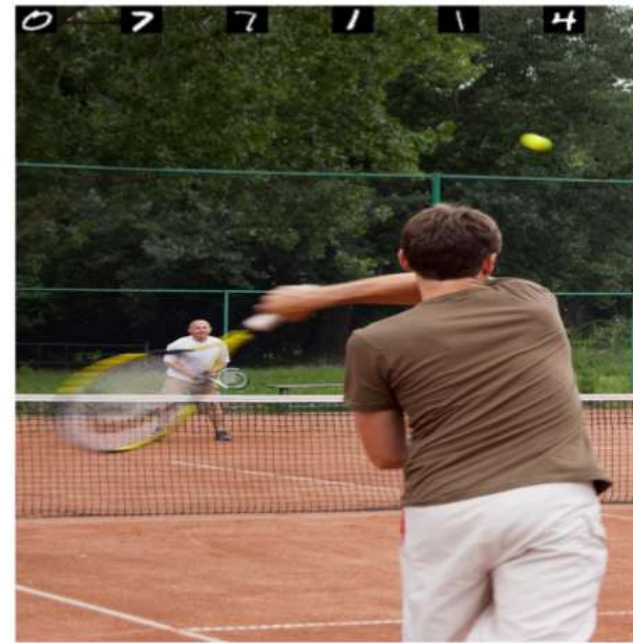


## ■ Experiment Setting

- Datasets: Flickr30k, MS-COCO, Flickr30k-shortcuts, MS-COCO-shortcuts
- Baseline: Vanilla InfoNCE, Latent Target Decoding (LTD), Implicit Feature Modification (IFM)



(a) Caption: "A bathroom sink with wood finish cabinets. 0 3 9 9 6 5."



(b) Caption: "A guy in a brown shirt has just hit a tennis ball. 0 7 7 1 1 4."

Figure 4<sup>[1]</sup>. Two random samples from the MS-COCO dataset including shortcuts added on both the image and caption.

[1] Demonstrating and Reducing Shortcuts in Vision-Language Representation Learning

## ■ Results on public benchmark

Table1:Result on Flickr30k and MS-COCO with varied method. sc denotes shortcut, we evaluate CLIP and VSE++ w/wo shortcut on i2t and i2i task. QUEST outperforms InfoNCE and achieve superior performance compare with other baselines in most cases. †denote use of ltd.

Method	sc	Flickr30k						RSUM	MS-COCO						RSUM	
		i2t			t2i				i2t			t2i				
		R@1	R@5	R@10	R@1	R@5	R@10		R@1	R@5	R@10	R@1	R@5	R@10		
<b>CLIP</b>																
$\mathcal{L}_{\text{InfoNCE}}$	✗	86.9 $\pm$ 0.1	97.4 $\pm$ 0.1	99.0 $\pm$ 0.0	72.4 $\pm$ 0.1	92.1 $\pm$ 0.0	95.8 $\pm$ 0.0	543.5 $\pm$ 1.1	63.8 $\pm$ 0.3	86.1 $\pm$ 0.2	92.3 $\pm$ 0.0	46.3 $\pm$ 0.3	74.8 $\pm$ 0.1	84.1 $\pm$ 0.2	447.5 $\pm$ 0.5	
$\mathcal{L}_{\text{InfoNCE+LTD}}$	✗	86.5 $\pm$ 0.6	97.1 $\pm$ 0.0	98.5 $\pm$ 0.0	72.4 $\pm$ 0.0	<b>92.3</b> $\pm$ 0.0	<b>95.9</b> $\pm$ 0.0	542.8 $\pm$ 0.8	63.8 $\pm$ 0.0	86.1 $\pm$ 0.0	92.3 $\pm$ 0.0	46.3 $\pm$ 0.0	74.7 $\pm$ 0.0	84.1 $\pm$ 0.0	447.4 $\pm$ 0.0	
$\mathcal{L}_{\text{InfoNCE+IFM}}$	✗	87.4 $\pm$ 0.1	97.4 $\pm$ 0.2	99.1 $\pm$ 0.0	73.2 $\pm$ 0.0	92.2 $\pm$ 0.0	95.6 $\pm$ 0.0	544.9 $\pm$ 0.2	63.0 $\pm$ 0.1	86.6 $\pm$ 0.1	92.6 $\pm$ 0.2	47.2 $\pm$ 0.0	75.6 $\pm$ 0.0	84.5 $\pm$ 0.0	449.5 $\pm$ 1.7	
$\mathcal{L}_{\text{QUEST(Ours)}}$	✗	<b>89.3</b> $\pm$ 0.3	<b>97.8</b> $\pm$ 0.2	<b>99.2</b> $\pm$ 0.3	<b>73.9</b> $\pm$ 0.5	91.5 $\pm$ 0.3	95.0 $\pm$ 0.3	<b>546.7</b> $\pm$ 1.9	<b>65.4</b> $\pm$ 0.5	<b>87.7</b> $\pm$ 0.2	<b>93.6</b> $\pm$ 0.4	<b>48.5</b> $\pm$ 0.2	<b>75.7</b> $\pm$ 0.5	<b>84.7</b> $\pm$ 0.6	<b>455.6</b> $\pm$ 2.4	
$\mathcal{L}_{\text{InfoNCE}}$	✓	57.2 $\pm$ 8.3	84.0 $\pm$ 4.8	91.0 $\pm$ 1.9	44.9 $\pm$ 4.5	74.9 $\pm$ 6.0	84.2 $\pm$ 2.5	436.2 $\pm$ 145.0	13.6 $\pm$ 0.9	31.5 $\pm$ 2.4	42.2 $\pm$ 3.7	7.3 $\pm$ 0.6	22.1 $\pm$ 1.0	32.7 $\pm$ 1.7	149.4 $\pm$ 32.7	
$\mathcal{L}_{\text{InfoNCE+LTD}}$	✓	64.0 $\pm$ 1.3	87.8 $\pm$ 0.9	93.2 $\pm$ 0.8	50.7 $\pm$ 0.6	79.8 $\pm$ 0.7	88.1 $\pm$ 0.5	463.6 $\pm$ 17.3	18.9 $\pm$ 0.1	41.8 $\pm$ 0.1	54.1 $\pm$ 0.1	16.5 $\pm$ 0.0	39.4 $\pm$ 0.0	52.6 $\pm$ 0.1	223.4 $\pm$ 0.2	
$\mathcal{L}_{\text{InfoNCE+IFM}}$	✓	73.8 $\pm$ 0.8	91.5 $\pm$ 0.5	95.6 $\pm$ 0.0	58.9 $\pm$ 0.1	84.4 $\pm$ 0.1	91.1 $\pm$ 0.2	495.2 $\pm$ 5.7	23.4 $\pm$ 1.5	46.5 $\pm$ 2.7	58.2 $\pm$ 2.5	17.1 $\pm$ 0.3	38.9 $\pm$ 0.9	51.3 $\pm$ 1.0	235.5 $\pm$ 43.8	
$\mathcal{L}_{\text{QUEST(Ours)}}$	✓	<b>84.2</b> $\pm$ 0.3	<b>96.0</b> $\pm$ 0.1	<b>97.7</b> $\pm$ 0.2	<b>67.6</b> $\pm$ 0.5	<b>88.9</b> $\pm$ 0.2	<b>93.4</b> $\pm$ 0.1	<b>527.8</b> $\pm$ 1.4	<b>50.8</b> $\pm$ 0.3	<b>75.4</b> $\pm$ 0.4	<b>84.1</b> $\pm$ 0.4	<b>37.9</b> $\pm$ 0.3	<b>65.1</b> $\pm$ 0.3	<b>76.1</b> $\pm$ 0.4	<b>389.4</b> $\pm$ 2.1	
<b>VSE++</b>																
$\mathcal{L}_{\text{InfoNCE}}$	✗	52.6 $\pm$ 1.1	79.8 $\pm$ 0.1	87.8 $\pm$ 0.1	39.5 $\pm$ 0.3	69.8 $\pm$ 0.0	79.4 $\pm$ 0.1	409.0 $\pm$ 4.0	42.2 $\pm$ 0.1	72.7 $\pm$ 0.1	83.2 $\pm$ 0.1	30.9 $\pm$ 0.0	61.2 $\pm$ 0.1	73.5 $\pm$ 0.1	363.8 $\pm$ 2.3	
$\mathcal{L}_{\text{InfoNCE+LTD}}$	✗	54.1 $\pm$ 0.1	81.1 $\pm$ 0.8	88.6 $\pm$ 0.1	42.5 $\pm$ 0.0	71.9 $\pm$ 0.1	81.3 $\pm$ 0.0	419.6 $\pm$ 0.1	43.6 $\pm$ 0.1	73.5 $\pm$ 0.0	83.7 $\pm$ 0.0	32.4 $\pm$ 0.1	62.5 $\pm$ 0.0	74.7 $\pm$ 0.0	370.5 $\pm$ 0.1	
$\mathcal{L}_{\text{InfoNCE+IFM}}$	✗	52.4 $\pm$ 0.2	76.9 $\pm$ 0.1	85.3 $\pm$ 0.0	39.1 $\pm$ 0.0	68.8 $\pm$ 0.1	78.2 $\pm$ 0.1	400.7 $\pm$ 0.0	40.2 $\pm$ 0.0	70.8 $\pm$ 0.1	81.6 $\pm$ 0.1	30.8 $\pm$ 0.0	61.5 $\pm$ 0.0	74.3 $\pm$ 0.0	359.3 $\pm$ 1.1	
$\mathcal{L}_{\text{QUEST(Ours)}}$	✗	<b>54.7</b> $\pm$ 0.2	<b>81.3</b> $\pm$ 0.4	<b>88.8</b> $\pm$ 0.3	<b>42.9</b> $\pm$ 0.1	<b>72.3</b> $\pm$ 0.4	<b>81.6</b> $\pm$ 1.1	<b>421.6</b> $\pm$ 2.5	<b>45.3</b> $\pm$ 0.1	<b>75.5</b> $\pm$ 0.5	<b>85.4</b> $\pm$ 0.4	<b>34.1</b> $\pm$ 0.1	<b>64.5</b> $\pm$ 0.2	<b>76.3</b> $\pm$ 0.2	<b>381.1</b> $\pm$ 1.5	
$\mathcal{L}_{\text{InfoNCE}}$	✓	0.1 $\pm$ 0.0	0.4 $\pm$ 0.0	0.8 $\pm$ 0.0	0.1 $\pm$ 0.0	0.4 $\pm$ 0.0	1.0 $\pm$ 0.0	2.9 $\pm$ 0.0	0.0 $\pm$ 0.0	0.1 $\pm$ 0.0	0.2 $\pm$ 0.0	0.0 $\pm$ 0.0	0.1 $\pm$ 0.0	0.2 $\pm$ 0.0	0.6 $\pm$ 0.0	
$\mathcal{L}_{\text{InfoNCE+LTD}}$	✓	24.7 $\pm$ 0.5	<b>51.8</b> $\pm$ 0.7	<b>65.6</b> $\pm$ 1.4	<b>20.7</b> $\pm$ 1.0	<b>49.2</b> $\pm$ 0.6	<b>62.6</b> $\pm$ 1.2	<b>274.6</b> $\pm$ 4.6	3.9 $\pm$ 0.0	13.7 $\pm$ 0.6	21.6 $\pm$ 0.9	3.1 $\pm$ 0.2	11.0 $\pm$ 1.6	18.1 $\pm$ 3.0	71.4 $\pm$ 3.6	
$\mathcal{L}_{\text{InfoNCE+IFM}}$	✓	0.0 $\pm$ 0.0	0.6 $\pm$ 0.1	0.9 $\pm$ 0.2	0.1 $\pm$ 0.0	0.5 $\pm$ 0.0	1.0 $\pm$ 0.0	3.2 $\pm$ 0.8	0.0 $\pm$ 0.0	0.1 $\pm$ 0.0	0.2 $\pm$ 0.0	0.0 $\pm$ 0.0	0.1 $\pm$ 0.0	0.2 $\pm$ 0.0	0.7 $\pm$ 0.0	
$\mathcal{L}_{\text{QUEST(Ours)}}$ †	✓	<b>24.9</b> $\pm$ 0.4	48.4 $\pm$ 0.3	61.1 $\pm$ 0.5	17.5 $\pm$ 0.3	43.4 $\pm$ 0.6	56.5 $\pm$ 0.8	251.8 $\pm$ 2.9	<b>10.5</b> $\pm$ 0.6	<b>27.9</b> $\pm$ 0.3	<b>40.6</b> $\pm$ 0.9	<b>9.4</b> $\pm$ 0.5	<b>29.0</b> $\pm$ 1.4	<b>42.6</b> $\pm$ 2.1	<b>160.0</b> $\pm$ 5.8	



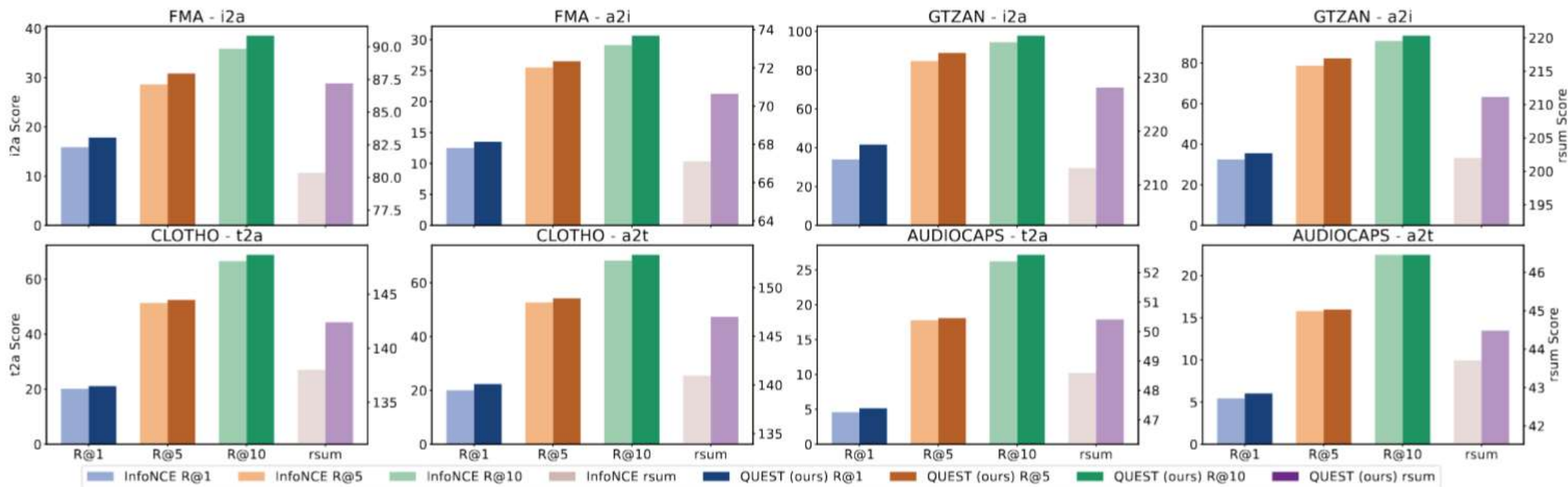
## ■ Ablation Study

Table2: ablation study on image caption retrieval task with different training objectives. D1 and D2 denote decoders in the architecture. Decoder with all  $\times$  beneath are omitted, while those with  $\checkmark$  indicate optimization with corresponding objective functions. Bold and underlined numbers indicate the best and second-best results, respectively.

Methods				Flickr30k						MS-COCO							
D1		D2		<i>i2t</i>			<i>t2i</i>			RSUM	<i>i2t</i>			<i>t2i</i>			RSUM
$\mathcal{L}_{SIC}$	$\mathcal{L}_{SIC}$	$\mathcal{L}_{UIC}$	$\mathcal{L}_{P-UIC}$	<i>R@1</i>	<i>R@5</i>	<i>R@10</i>	<i>R@1</i>	<i>R@5</i>	<i>R@10</i>		<i>R@1</i>	<i>R@5</i>	<i>R@10</i>	<i>R@1</i>	<i>R@5</i>	<i>R@10</i>	
				<b>CLIP</b>						<b>CLIP</b>							
$\checkmark$	$\times$	$\times$	$\times$	<u>86.9</u>	<u>97.4</u>	<u>98.8</u>	<u>72.4</u>	<b>92.1</b>	<b>95.8</b>	<u>543.4</u>	63.8	86.1	92.3	46.3	74.8	<u>84.1</u>	447.4
$\times$	$\times$	$\checkmark$	$\times$	80.8	94.4	96.5	66.9	88.4	92.9	519.9	55.1	81.6	89.4	43.3	72.5	82.9	424.8
$\times$	$\times$	$\times$	$\checkmark$	85.5	96.6	98.1	70.0	87.3	90.7	528.2	63.0	85.9	91.7	45.2	70.4	79.2	435.4
$\checkmark$	$\times$	$\checkmark$	$\times$	81.9	96.2	98.0	69.5	90.0	94.2	529.8	<u>64.9</u>	<u>86.6</u>	<u>92.8</u>	<u>47.3</u>	<u>75.1</u>	83.9	<u>450.6</u>
$\checkmark$	$\checkmark$	$\times$	$\times$	76.6	92.3	95.9	59.0	84.4	90.9	499.1	54.4	79.2	86.7	39.9	68.4	78.9	407.5
$\checkmark$	$\times$	$\times$	$\checkmark$	<b>89.3</b>	<b>97.8</b>	<b>99.2</b>	<b>73.9</b>	<u>91.5</u>	<u>95.0</u>	<b>546.7</b>	<b>65.4</b>	<b>87.7</b>	<b>93.6</b>	<b>48.5</b>	<b>75.7</b>	<b>84.7</b>	<b>455.6</b>
				<b>VSE++</b>						<b>VSE++</b>							
$\checkmark$	$\times$	$\times$	$\times$	52.6	<u>79.8</u>	<u>87.8</u>	39.5	<u>69.8</u>	<u>79.4</u>	<u>408.9</u>	42.2	72.7	83.2	30.9	61.2	73.5	363.7
$\times$	$\times$	$\checkmark$	$\times$	47.8	72.8	80.8	36.7	63.2	73.3	374.6	40.7	71.2	82.1	30.3	60.5	73.0	357.8
$\times$	$\times$	$\times$	$\checkmark$	49.0	74.1	81.3	36.4	64.1	73.1	378.0	40.9	71.4	82.4	30.8	60.6	73.2	359.3
$\checkmark$	$\times$	$\checkmark$	$\times$	<u>53.3</u>	<u>79.8</u>	87.6	<u>40.5</u>	68.1	78.0	407.3	<u>44.9</u>	<u>74.1</u>	<u>84.4</u>	<u>32.3</u>	<u>62.8</u>	<u>74.7</u>	<u>373.2</u>
$\checkmark$	$\times$	$\times$	$\checkmark$	<b>54.7</b>	<b>80.3</b>	<b>88.2</b>	<b>42.0</b>	<b>70.3</b>	<b>79.6</b>	<b>415.1</b>	<b>45.3</b>	<b>75.5</b>	<b>85.4</b>	<b>34.1</b>	<b>64.5</b>	<b>76.3</b>	<b>381.1</b>

## ■ Extend To More Modalities

Figure 5: Performance comparison of InfoNCE and QUEST methods with additional audio modality on image-to-audio (i2a) and audio-to-image (a2i) retrieval tasks across FMA, GTZAN, CLOTHO, and AUDIOCAPS datasets.





## ■ Visualization

Figure 6: Case Study: (a) Image-to-text retrieval, where the results of  $L_{QUEST}$  and  $L_{InfoNCE}$  are denoted by italics and underlines, respectively. (b) Text-to-image retrieval, where red and green borders indicate the top-5 retrievals using  $L_{QUEST}$ , while blue borders represent those using  $L_{InfoNCE}$ .



(a)

(b)

## ■ Conclusions

- We propose QUEST, a novel multimodal contrastive learning method that effectively captures and preserves more task-relevant unique information from individual modalities.
- Quaternions indirectly optimize unique and shared information simultaneously, expanding the unique representational space.
- Unique decoders extract unique and shared information through constraints and self-penalization, outperforming SOTA methods.





---

# QUEST: Quadruple Multimodal Contrastive Learning with Constraints and Self-Penalization

---

Qi Song<sup>1\*</sup>, Tianxiang Gong<sup>2\*</sup>, Shiqi Gao<sup>2</sup>,  
Haoyi Zhou<sup>1,3†</sup>, Jianxin Li<sup>2,3</sup>

<sup>1</sup>School of Software, Beihang University

<sup>2</sup>School of Computer Science and Engineering, Beihang University

<sup>3</sup>Zhongguancun Laboratory, Beijing

{songqi23, gongtx, gaoshiqi, haoyi, lijx}@buaa.edu

---

\*Equal contribution

†Corresponding author.