



中国科学院自动化研究所

Institute of Automation, Chinese Academy of Sciences

# TaiSu: A 166M Large-scale High-Quality Dataset for Chinese Vision-Language Pre-training

Yulong Liu<sup>1,2</sup>, Guibo Zhu<sup>1,4,5</sup>\*, Bin Zhu<sup>3</sup>, Qi Song<sup>3</sup>, Guojing Ge<sup>1</sup>, Haoran Chen<sup>1,4</sup>, Guanhui Qiao<sup>1,4</sup>, Ru Peng<sup>2</sup>, Lingxiang Wu<sup>1</sup>, and Jinqiao Wang<sup>1,4,5</sup>

<sup>1</sup>Institute of Automation, Chinese Academy of Sciences (CASIA), Beijing, China

<sup>2</sup>Institute of Artificial Intelligence and Robotics, Xi'an Jiaotong University, Xi'an, China

<sup>3</sup>School of Artificial Intelligence, Beijing Normal University, Beijing, China

<sup>4</sup>School of Artificial Intelligence, University of Chinese Academy of Sciences, Beijing, China

<sup>5</sup>Wuhan AI Research

# Introduction to Vision-Language Pretraining (VLP)

- **Vision-Language pretraining**

Learning visual and linguistic representation from large-scale data pairs with weak correlation.

- **Downstream tasks**

- Image Captioning,
- Visual Question Answering,
- Cross-modal Retrieval
- Text to Image Generation
- .....

- **VLP model architectures**

- Dual-stream model
- Single-stream model

# VLP Datasets

- Most of the public large-scale multi-modal datasets are based on English corpus.
- Almost all the large-scale multi-modal datasets only provide one caption for each image.

Dataset	Language	Available	Size
CC3M[9]	English	Yes	3,000,000
M5Product[10]	English	Yes	5,000,000
JFT-300M[11]	English	No	300,000,000
WIT[12]	multilingual	Yes	11,500,000
CC12M[13]	English	Yes	12,000,000
YFCC100M[14]	English	Yes	99,200,000
LAION-400M[15]	English	Yes	400,000,000
JFT-3B[16]	English	No	3,000,000,000
IG-3.5B-17K[17]	English	No	3,500,000,000
LAION-5B[18]	multilingual	Yes	5,850,000,000
Product1M[19]	Chinese	Yes	1,000,000
WudaoMM[8]	Chinese	Yes	5,000,000
M6-Corpus[6]	Chinese	No	60,500,000
WuKong[7]	Chinese	Yes	101,483,885
<b>TaiSu(Ours)</b>	Chinese	Yes	<b>166,000,000</b>

**TaiSu is the largest public Chinese multi-modal dataset where one image may have more than one captions.**

# Main Contributions

- A pipeline for the construction of a large-scale multi-modal dataset.

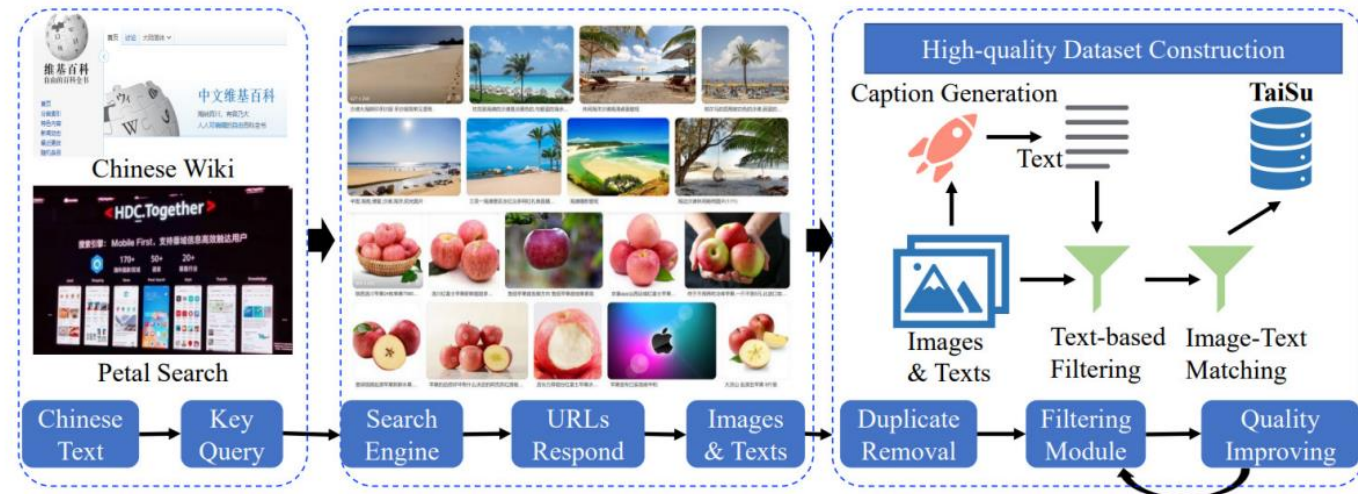
text-based filtering,  
image-text matching,  
text augmentation

- The largest public multi-modal dataset for Chinese VLP.

166M images

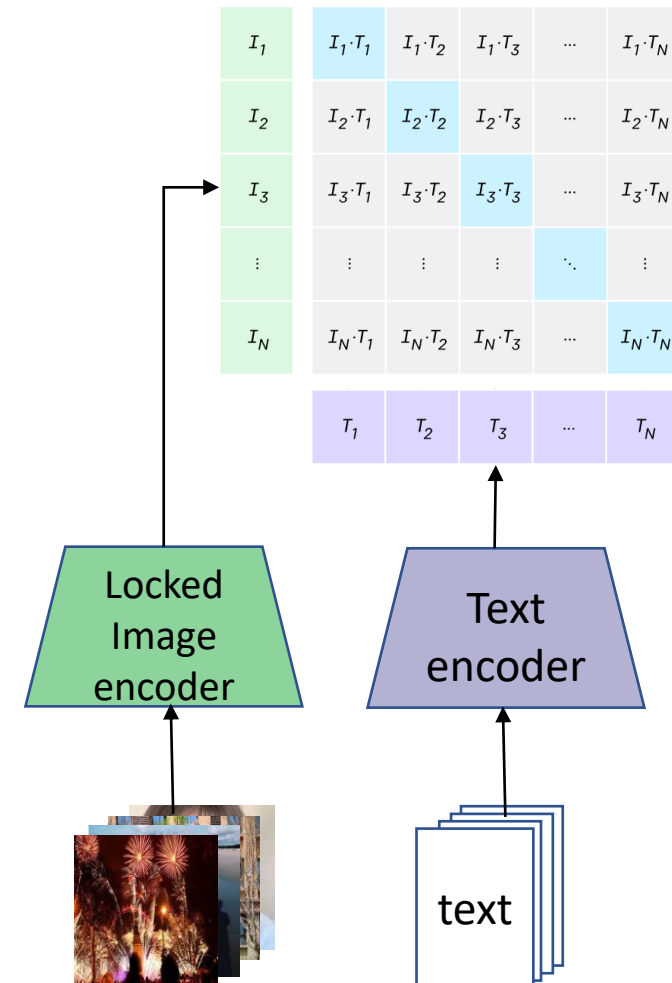
219M Chinese captions

- Our experiments show that the pretrained models can benefit from the filtering process and the combination of the web data and the generated captions.



# Summary of the pretraining methods

- **Dual-stream architecture**  
image encoder + text encoder.
- **Contrastive loss with global similarity**  
Follow the settings of CLIP
- **LiT-tuning**  
The pretrained image encoder → **frozen**,  
The text encoder → **need to be trained**.
- **Three kinds of data**  
Raw web data  
Filtered web data  
Filtered web data + generated data



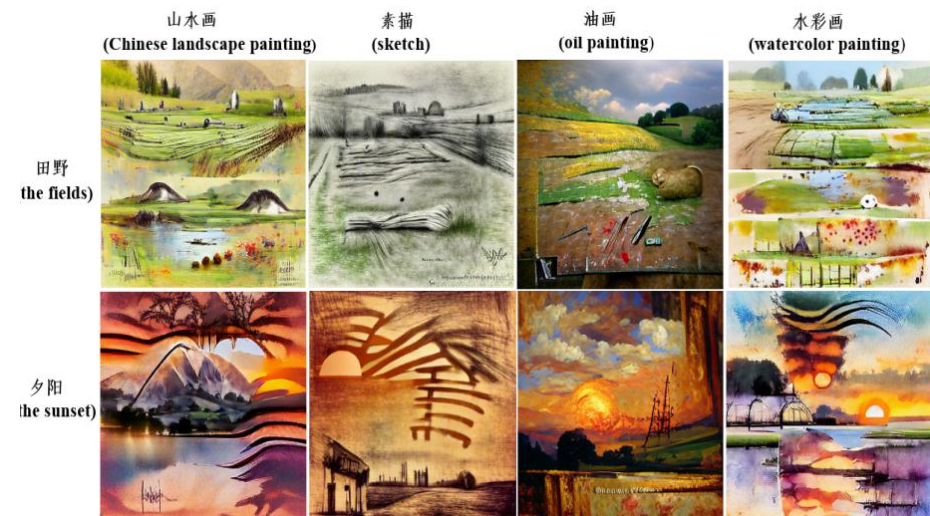
# Zero-shot Image-text Retrieval & Zero-shot Image Classification & Text-to-image Generation

## Zero-shot Image-text Retrieval

Dataset	Method	Image-to-Text			Text-to-Image			MR
		R@1	R@5	R@10	R@1	R@5	R@10	
Flickr8K-CN	BriVL[44]	13.4	31.2	40.7	8.0	20.7	29.5	23.9
	Wukong <sub>ViT-B</sub> [7]	55.4	82.3	90.0	43.2	71.3	81.3	70.6
	Wukong <sub>Swin-L</sub> [7]	47.2	78.8	87.6	36.6	64.8	76.2	65.2
	Ours <sub>RN101</sub>	55.1	82.6	<b>90.9</b>	44.9	74.2	<b>84.3</b>	72.0
	Ours <sub>ViT-B</sub>	<b>57.6</b>	<b>83.4</b>	90.6	<b>45.4</b>	<b>74.4</b>	84.1	<b>72.6</b>
Flickr30K-CN	BriVL[44]	17.7	42.3	54.3	10.3	27.5	37.9	31.7
	Wukong <sub>ViT-B</sub> [7]	<b>66.2</b>	88.7	94.3	45.7	73.8	82.2	75.1
	Wukong <sub>Swin-L</sub> [7]	58.7	86.7	92.7	40.9	68.0	78.4	70.9
	Ours <sub>RN101</sub>	65.3	88.6	94.1	<b>51.2</b>	<b>79.1</b>	<b>89.5</b>	77.6
	Ours <sub>ViT-B</sub>	65.6	<b>90.1</b>	<b>94.9</b>	49.9	78.9	87.0	<b>77.7</b>
COCO-CN	BriVL[44]	17.1	41.7	57.5	14.8	39.0	54.2	37.4
	Wukong <sub>ViT-B</sub> [7]	48.3	77.8	88.8	49.2	79.4	87.9	71.9
	Wukong <sub>Swin-L</sub> [7]	47.3	78.0	88.3	46.4	77.0	87.6	70.8
	Ours <sub>RN101</sub>	<b>54.1</b>	<b>82.8</b>	<b>91.8</b>	<b>54.3</b>	82.7	<b>92.4</b>	<b>76.4</b>
	Ours <sub>ViT-B</sub>	52.5	81.5	91.4	53.6	<b>83.7</b>	<b>92.4</b>	75.9
MUGE	BriVL[44]	-	-	-	12.7	30.9	41.8	28.5
	Wukong <sub>ViT-B</sub> [7]	-	-	-	33.4	59.3	69.7	54.1
	Wukong <sub>Swin-L</sub> [7]	-	-	-	<b>34.5</b>	<b>60.6</b>	<b>71.2</b>	<b>55.5</b>
	Ours <sub>RN101</sub>	-	-	-	27.5	53.9	64.8	48.7
	Ours <sub>ViT-B</sub>	-	-	-	29.7	57.0	67.4	51.4

## Zero-shot Image classification

Dataset	Metric	BriVL[44]	RN101			ViT-B		
			RD	TS <sub>Web</sub>	TS <sub>all</sub>	RD	TS <sub>Web</sub>	TS <sub>all</sub>
CIFAR10	R@1	72.3	79.8	81.0	74.1	87.9	<b>89.2</b>	85.5
CIFAR100	R@1	35.9	43.3	43.0	43.0	55.5	<b>56.7</b>	55.4
Caltech101	R@1	72.0	71.4	71.8	73.5	72.6	73.2	<b>74.2</b>
Caltech256	R@1	58.0	68.5	69.9	70.2	69.1	71.0	<b>72.0</b>
DTD	R@1	18.8	24.9	26.2	<b>29.8</b>	26.5	29.3	28.7
Flowers	R@1	18.4	24.3	23.3	21.3	23.8	<b>30.6</b>	22.4
EuroSAT	R@1	25.5	29.9	23.6	39.7	36.0	29.5	<b>52.6</b>
ImageNet	R@1	24.3	33.3	34.2	33.4	33.7	<b>35.4</b>	34.0
AVG		40.65	46.93	46.63	48.13	50.64	51.86	<b>53.10</b>



# Ablation study

Dataset	Model	Image-to-Text			Text-to-Image			MR
		R@1	R@5	R@10	R@1	R@5	R@10	
Flickr8k	RD-RN101	41.0	69.0	80.0	28.3	54.4	66.8	56.6
	TS <sub>Web</sub> -RN101	44.2	69.2	80.2	29.4	56.0	68.3	57.9
	TS <sub>all</sub> -RN101	55.1	82.6	<b>90.9</b>	44.9	74.2	<b>84.3</b>	72.0
	RD-ViT	35.5	63.3	76.8	26.9	52.4	64.7	53.8
	TS <sub>Web</sub> -ViT	40.8	70.1	80.6	31.3	57.4	69.0	58.2
	TS <sub>all</sub> -ViT	<b>57.6</b>	<b>83.4</b>	90.6	<b>45.4</b>	<b>74.4</b>	84.1	<b>72.6</b>
Flickr30k	RD-RN101	47.0	75.5	85.3	31.9	61.5	72.9	62.4
	TS <sub>Web</sub> -RN101	48.5	77.4	86.2	35.2	63.5	73.8	64.1
	TS <sub>all</sub> -RN101	65.3	88.6	94.1	<b>51.2</b>	<b>79.1</b>	<b>89.5</b>	77.6
	RD-ViT	42.3	71.6	82.8	29.7	57.1	68.6	58.7
	TS <sub>Web</sub> -ViT	49.8	77.3	87.2	34.7	63.0	73.0	64.2
	TS <sub>all</sub> -ViT	<b>65.6</b>	<b>90.1</b>	<b>94.9</b>	49.9	78.9	87.0	<b>77.7</b>
MUGE	RD-RN101	-	-	-	26.7	52.4	63.8	47.6
	TS <sub>Web</sub> -RN101	-	-	-	<b>30.1</b>	56.4	<b>67.4</b>	51.3
	TS <sub>all</sub> -RN101	-	-	-	27.5	53.9	64.8	48.7
	RD-ViT	-	-	-	27.6	53.6	65.1	48.8
	TS <sub>Web</sub> -ViT	-	-	-	29.7	<b>57.0</b>	<b>67.4</b>	<b>51.4</b>
	TS <sub>all</sub> -ViT	-	-	-	28.5	54.2	65.2	49.3
COCO-CN	RD-RN101	41.8	71.6	84.2	39.8	69.9	81.5	64.8
	TS <sub>Web</sub> -RN101	43.4	75.1	85.7	41.8	70.6	84.8	66.9
	TS <sub>all</sub> -RN101	<b>54.1</b>	<b>82.8</b>	<b>91.8</b>	<b>54.3</b>	82.7	<b>92.4</b>	<b>76.4</b>
	RD-ViT	40.8	70.3	83.3	37.2	68.6	82.2	63.7
	TS <sub>Web</sub> -ViT	44.4	73.3	84.8	39.8	72.3	83.9	66.4
	TS <sub>all</sub> -ViT	52.5	81.5	91.4	53.6	<b>83.7</b>	<b>92.4</b>	75.9

Dataset	Metric	BriVL[44]	RN101			ViT-B		
			RD	TS <sub>Web</sub>	TS <sub>all</sub>	RD	TS <sub>Web</sub>	TS <sub>all</sub>
CIFAR10	R@1	72.3	79.8	81.0	74.1	87.9	<b>89.2</b>	85.5
CIFAR100	R@1	35.9	43.3	43.0	43.0	55.5	<b>56.7</b>	55.4
Caltech101	R@1	72.0	71.4	71.8	73.5	72.6	73.2	<b>74.2</b>
Caltech256	R@1	58.0	68.5	69.9	70.2	69.1	71.0	<b>72.0</b>
DTD	R@1	18.8	24.9	26.2	<b>29.8</b>	26.5	29.3	28.7
Flowers	R@1	18.4	24.3	23.3	21.3	23.8	<b>30.6</b>	22.4
EuroSAT	R@1	25.5	29.9	23.6	39.7	36.0	29.5	<b>52.6</b>
ImageNet	R@1	24.3	33.3	34.2	33.4	33.7	<b>35.4</b>	34.0
AVG		40.65	46.93	46.63	48.13	50.64	51.86	<b>53.10</b>

- The filtering based on image-text matching can improve the performances of the pretrained models both on the zero-shot image-text retrieval task and on the zero-shot image classification task.
- The addition of the generated captions can significantly improve the performance on the zero-shot image-text retrieval task.

# Comparison with CLIP + translator

Dataset	Method	Image-to-Text			Text-to-Image			MR
		R@1	R@5	R@10	R@1	R@5	R@10	
Flickr8K-CN	CLIP <sub>RN101</sub> [1]	50.2	77.4	86.9	32.2	59.0	70.0	62.6
	CLIP <sub>ViT-B/32</sub>	51.1	77.1	85.9	34.3	60.5	71.2	63.4
	Ours <sub>RN101</sub>	55.1	82.6	<b>90.9</b>	44.9	74.2	<b>84.3</b>	72.0
	Ours <sub>ViT-B</sub>	<b>57.6</b>	<b>83.4</b>	90.6	<b>45.4</b>	<b>74.4</b>	84.1	<b>72.6</b>
Flickr30K-CN	CLIP <sub>RN101</sub> [1]	57.2	83.8	91.1	34.2	58.8	68.5	65.6
	CLIP <sub>ViT-B/32</sub>	58.5	83.9	90.2	34.4	60.1	69.1	66.0
	Ours <sub>RN101</sub>	65.3	88.6	94.1	<b>51.2</b>	<b>79.1</b>	<b>89.5</b>	77.6
	Ours <sub>ViT-B</sub>	<b>65.6</b>	<b>90.1</b>	<b>94.9</b>	49.9	78.9	87.0	<b>77.7</b>
COCO-CN	CLIP[1]	37.8	66.2	77.6	33.2	62.6	75.8	58.9
	CLIP <sub>ViT-B/32</sub>	37.5	65.0	76.8	34.3	63.1	75.3	58.7
	Ours <sub>RN101</sub>	<b>54.1</b>	<b>82.8</b>	<b>91.8</b>	<b>54.3</b>	82.7	<b>92.4</b>	<b>76.4</b>
	Ours <sub>ViT-B</sub>	52.5	81.5	91.4	53.6	<b>83.7</b>	<b>92.4</b>	75.9
MUGE	CLIP <sub>RN101</sub> [1]	-	-	-	8.8	18.2	24.1	17.0
	CLIP <sub>ViT-B/32</sub>	-	-	-	8.4	18.3	23.5	16.7
	Ours <sub>RN101</sub>	-	-	-	27.5	53.9	64.8	48.7
	Ours <sub>ViT-B</sub>	-	-	-	<b>29.7</b>	<b>57.0</b>	<b>67.4</b>	<b>51.4</b>

CLIP's training data:  
400M English  
image-text pairs

On the Chinese image-text retrieval task, TaiSu's models obtain better performance than CLIP's models. This can reveal **the necessity of the construction of a large-scale Chinese multi-modal dataset**