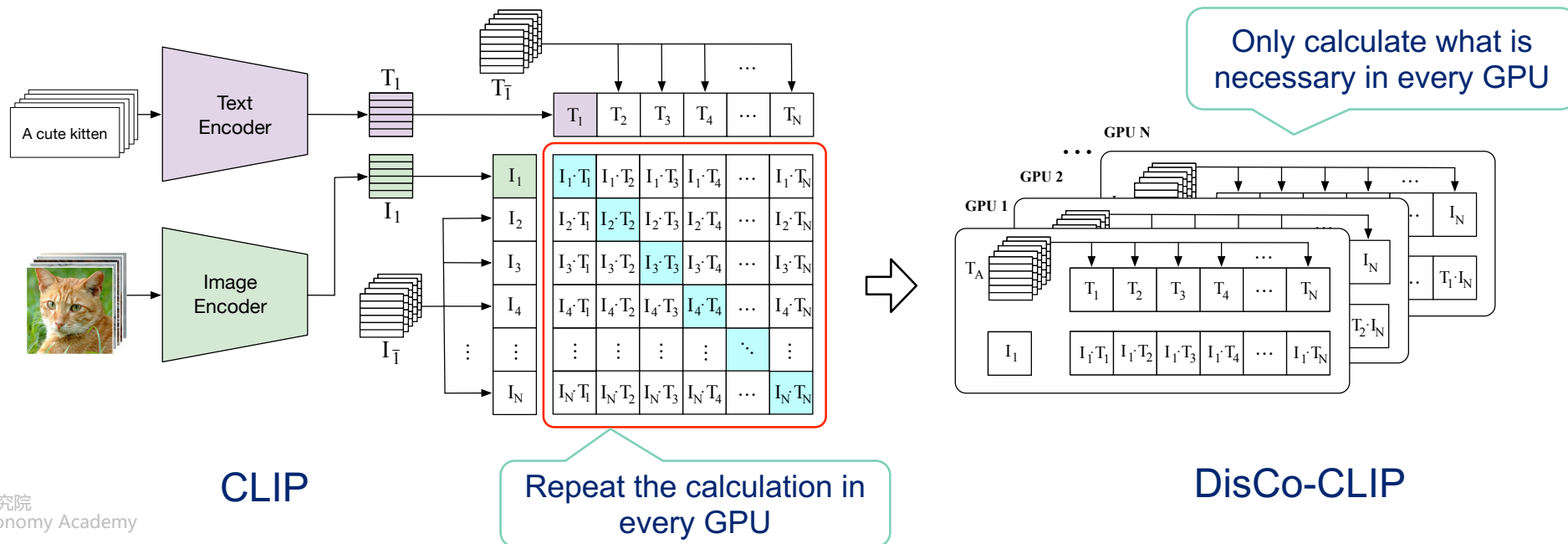


# DisCo-CLIP: A Distributed Contrastive Loss for Memory Efficient CLIP Training

Yihao Chen, Xianbiao Qi, Jianan Wang, Lei Zhang  
International Digital Economy Academy (IDEA)

Paper Tag: THU-PM-195

- We decompose the Contrastive Loss and reduce unnecessary redundant computations. We effectively reducing the computational complexity from  $O(N^2)$  to  $O(N^2/n)$ .
- Our work is mathematically equivalent to the original contrastive loss computation.
- We have observed that larger batch sizes are highly effective for contrastive learning.



Train CLIP more efficiently !

- Enable training with large batch size for better performance.
- Enable training with limited GPU resources.

- The original CLIP contrastive loss function:

$$\mathcal{L}_d = \mathcal{L}_1(I_A, T_A) + \mathcal{L}_2(T_A, I_A)$$

Image-to-Text  
Contrastive Loss

Text-to-Image  
Contrastive Loss

Where  $I_A$  and  $T_A$  denote **all** image and text features.

- Each GPU calculates the full  $I_A \times T_A$  similarity matrix for contrastive loss computation, causing huge computation waste.

➤ We split  $I_A$  and  $T_A$  as below :

$$I_A = [I_n, I_{\bar{n}}],$$

$$T_A = [T_n, T_{\bar{n}}]$$

Where  $I_n, T_n$  are image and text features on the  $n$ -th GPU, and  $I_{\bar{n}}$  and  $T_{\bar{n}}$  denote the collection of features on other GPUs.

➤ Then the  $\mathcal{L}_d$  can be decomposed and rewritten as below:

$$\mathcal{L}_d = \mathcal{L}_1(I_n, T_A) + \mathcal{L}_1(I_{\bar{n}}, T_A) + \mathcal{L}_2(T_n, I_A) + \mathcal{L}_2(T_{\bar{n}}, I_A)$$

➤ It is clear that  $\frac{\partial \mathcal{L}_1(I_{\bar{n}}, T_A)}{\partial I_n}$  and  $\frac{\partial \mathcal{L}_2(T_{\bar{n}}, I_A)}{\partial T_n}$  do not include gradients.

➤ We can decompose the gradients  $\frac{\partial \mathcal{L}_d}{\partial I_n}$  and  $\frac{\partial \mathcal{L}_d}{\partial T_n}$  as below :

$$\frac{\partial \mathcal{L}_d}{\partial I_n} = \frac{\partial \mathcal{L}_1(I_n, T_A)}{\partial I_n} + \frac{\partial \mathcal{L}_2(T_n, I_A)}{\partial I_n} + \frac{\partial \mathcal{L}_2(T_{\bar{n}}, I_A)}{\partial I_n}$$

$$\frac{\partial \mathcal{L}_d}{\partial T_n} = \frac{\partial \mathcal{L}_1(I_n, T_A)}{\partial T_n} + \frac{\partial \mathcal{L}_2(T_n, I_A)}{\partial T_n} + \frac{\partial \mathcal{L}_1(I_{\bar{n}}, T_A)}{\partial T_n}$$

**intra-GPU gradient**

**inter-GPU gradient**

- For **intra-GPU gradient** on the  $n$ -th GPU, we only need to calculate  $\mathcal{L}_1(I_n, T_A)$  and  $\mathcal{L}_2(T_n, I_A)$ .
- For **inter-GPU gradient** on the  $n$ -th GPU, It would have already been calculated as **intra-GPU gradient** on other GPUs.
- In order to obtain  $\frac{\partial \mathcal{L}_d}{\partial I_n}$  and  $\frac{\partial \mathcal{L}_d}{\partial T_n}$ , we use the *All\_Reduce* operation to perform gradient communication and to collect all gradients with respect to  $I_n$  and  $T_n$ .

## Algorithm 1: Pseudo-code for DisCo-CLIP

```
# image_encoder: ResNet or ViT
# text_encoder: text Transformer
# img[b,H,W,C]: minibatch of images
# text[b,L]: minibatch of texts
# t: temperature parameter
# D: feature dimension
# N: nums of GPU
# rank: global rank

# extract feature representations
i_e = image_encoder(img) # [b,D]
t_e = text_encoder(text) # [b,D]

# gather features from all gpus
I_E = all_gather(i_e) # [N,b,D]
T_E = all_gather(t_e) # [N,b,D]

# scaled dot product similarities
logits_i = dot(I_E[rank], T_E)*t # [b,N*b]
logits_t = dot(T_E[rank], I_E)*t # [b,N*b]

# image-to-text and text-to-image losses
labels = arange(b) + b * rank
loss_i = cross_entropy_loss(logits_i, labels)
loss_t = cross_entropy_loss(logits_t, labels)

# loss backward
loss = (loss_i + loss_t) / 2
loss.backward()

# reduce gradients and losses from all gpus
I_E.grad = all_reduce(I_E.grad, op=AVG)
T_E.grad = all_reduce(T_E.grad, op=AVG)
Loss = all_reduce(loss, op=AVG)

# backbone backward
i_e.backward(I_E.grad[rank])
t_e.backward(T_E.grad[rank])
```

$\mathcal{L}_1(I_n, T_A)$

$\mathcal{L}_2(T_n, I_A)$

Calculate the **intra-GPU gradient** and the **inter-GPU gradient** of other GPUs.

Collect **intra-GPU gradient** and **inter-GPU gradient** from all GPUs



Methods	GPUs	BS	Memory
CLIP	64× A100 40GB	32,768	27.4GB
CLIP	64× A100 40GB	65,536	OOM
DisCo-CLIP	64× A100 40GB	32,768	16.5GB
DisCo-CLIP	64× A100 40GB	<b>196,608</b>	38.1GB
DisCo-CLIP	<b>8</b> × A100 40GB	32,768	36.9GB

Table 3. Memory consumption of CLIP and DisCo-CLIP under different settings. “OOM” means out of memory.

Model	Data Sets	Epochs	Steps	Batch Size	INet [7]	INet-v2 [32]	INet-R [15]	INet-S [40]
CLIP [27]	CLIP WIT 400M	32	≈ 400 K	32,768	63.3	56.0	69.4	42.3
OpenCLIP [16, 34]	LAION-400M	32	≈ 400 K	32,768	62.9	55.1	73.4	49.4
DisCo-CLIP	LAION-400M*	32	≈ 400 K	32,768	63.2	55.2	73.4	50.6
DisCo-CLIP	LAION-400M*	32	≈ <b>200</b> K	<b>65,536</b>	<b>64.3</b>	<b>56.2</b>	<b>73.8</b>	<b>51.7</b>

Table 7. Comparison of DisCo-CLIP with vanilla CLIP and re-implemented CLIP by LAION group. We report top-1 zero-shot classification accuracy (%) on several data sets. All models are based on ViT-B/32. Our LAION-400M\* is a 400M subset of LAION-2B. Training resource: 64 A100 40GB GPUs.

- Disco-CLIP is mathematically equivalent to the original contrastive loss computation.
- DisCo-CLIP can enable a much larger batch size for contrastive learning compared to CLIP, using the same GPU resource.

# THANK YOU

中国深圳市福田区市花路5号  
长富金茂大厦1号楼3901&57层, 邮编518045

Room 3901 and 57/F , Building 1 ,  
Chang Fu Jin Mao Tower , 5 Shihua Road ,  
Futian District , Shenzhen , China , 518045

86-755-61610106 (3901)

86-755-83219017 (57楼)



[www.idea.edu.cn](http://www.idea.edu.cn)