# CLIPLoss and Norm-Based Data Selection Methods for Multimodal Contrastive Learning
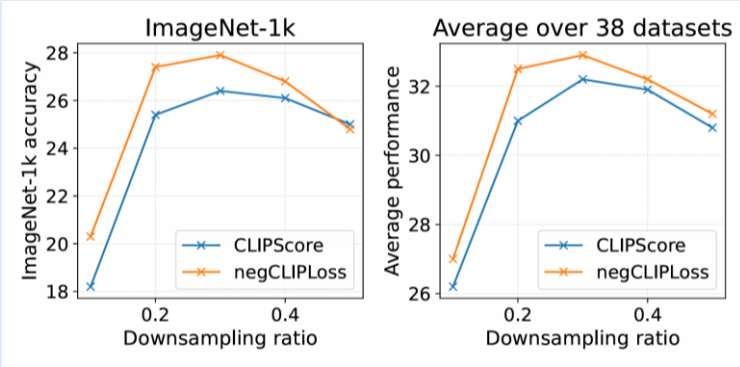
Yiping Wang[1,*], Yifang Chen[1,*], Wendan Yan[1,3], Alex Fang[1], Wenjing Zhou[2], Kevin Jamieson[1], Simon Shaolei Du[1]

[1] University of Washington, [2] University of Michigan, [3] Microsoft Corporation, [*] Equal Contribution

## Background

- Data selection is critical for large-scale multimodal pretraining, particularly for noisy web-curated datasets.

- Current best data selection approaches for CLIP pretraining leverage external non-CLIP models [2,3] (like BLIP and OCR models, etc) or the external large-scale high-quality pretraining datasets (like HQITP-350M[1]), while the potential of CLIP embedding is under-explored.



ImageNet-1k / Average over 38 datasets (CLIPScore vs negCLIPLoss plots against Downsampling ratio)

## Key Takeaway

- **CLIP Loss is better than CLIPScore in image-text data selection.** Simply replacing CLIPScore with negative CLIP Loss (negCLIPLoss) can consistently produce better quality measurement in data selection for CLIP pretraining.

- **New SOTA on DataComp-medium.** Our methods can achieve a new SOTA on the DataComp benchmark by combining them with the current top techniques.

- **High Efficiency and Universality.** Our methods don't need external non-CLIP models or the external large-scale pretraining datasets as in previous works. They are very efficient and also universal to different CLIP teacher models.

Code: https://github.com/ypwang61/negCLIPLoss_NormSim
DataComp samples and checkpoint:
https://huggingface.co/ypwang61/negCLIPLoss_NormSim/tree/main/medium_scale
DataComp Benchmark: https://www.datacomp.ai/dcclip/leaderboard.html

## negCLIPLoss ourperforms CLIPScore

- The difference between CLIPScore (similarity of one image-text pair) and the negative of CLIP Loss (negCLIPLoss) is only a normalization term.

- This term calculates the similarity among **cross**-image-text pairs, which penalizes the image/text features **shared** by different data (Like monotonous patterns/colors in images, or "Photo"/"Image" in texts), and rewards more distinctive data.

- Compared with CLIPScore, negCLIPLoss consistently provides better estimates of data quality.

### Quality Metric for data $i$

$$\text{CLIPScore} := f_i^T g_i \qquad \text{negCLIPLoss} := f_i^T g_i - \mathfrak{N} \propto \text{negative CLIP loss for data } i$$

where the normalization term $\mathfrak{N} := \frac{\tau}{2}\left[\log\left(\sum_{j\in B} e^{f_i^T g_j/\tau}\right) + \log\left(\sum_{j\in B} e^{f_j^T g_i/\tau}\right)\right]$

$f/g$ : Image/text embedding     $\tau$: temperature     $B$: batch



CLIPScore can underestimate the quality

"San Juan Islands Friday Harbor" — CLIPScore: Top 78%, negCLIPLoss: Top 34%, $\mathfrak{N}$: Top 100%
"American football" — CLIPScore: Top 56%, negCLIPLoss: Top 37%, $\mathfrak{N}$: Top 99%
"CIMG5175 Woolly sheep at Dunk's Green" — CLIPScore: Top 83%, negCLIPLoss: Top 42%, $\mathfrak{N}$: Top 100%

CLIPScore can overestimate the quality

"Caruba Step-down verloopring 67-46" — CLIPScore: Top 27%, negCLIPLoss: Top 39%, $\mathfrak{N}$: Top 17%
"17596 Green Willow Place - Photo 25" — CLIPScore: Top 16%, negCLIPLoss: Top 30%, $\mathfrak{N}$: Top 14%
"Listing Image 7" — CLIPScore: Top 36%, negCLIPLoss: Top 52%, $\mathfrak{N}$: Top 7%
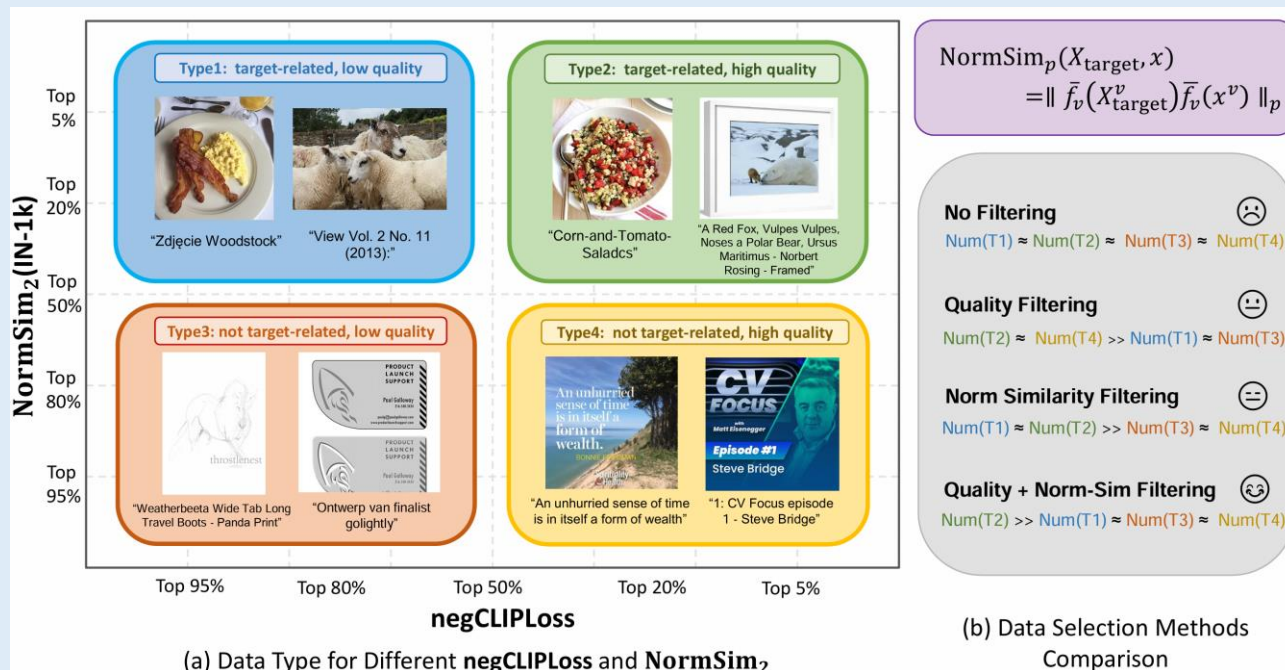
## NormSim: Estimate source-target similarity

- Quality and source-target similarity are orthogonal metrics. OCR-style image-text pairs can also have high-quality scores.

- Using $p$-norm similarity on the image embeddings to evaluate the source-target relevance. We use the validation data from the 24 downstream tasks as target. When $p = 2$, it's theoretically provable that $\text{NormSim}_2$ is optimal under a linear setting. And in experiments, $p = \infty$ has the best performance.

- We also have a variant $\text{NormSim}_2$-D (Dynamic) which is target-data-free.



$$\text{NormSim}_p(X_{\text{target}}, x) = \| \tilde{f}_v(X_{\text{target}}^v)\bar{f}_v(x^v) \|_p$$

No Filtering :( Num(T1) ≈ Num(T2) ≈ Num(T3) ≈ Num(T4)
Quality Filtering :| Num(T2) ≈ Num(T4) >> Num(T1) ≈ Num(T3)
Norm Similarity Filtering :| Num(T1) ≈ Num(T2) >> Num(T3) ≈ Num(T4)
Quality + Norm-Sim Filtering :) Num(T2) >> Num(T1) ≈ Num(T3) ≈ Num(T4)

(a) Data Type for Different **negCLIPLoss** and **NormSim₂**
(b) Data Selection Methods Comparison

## Experiment Results

**Setup:** Following **DataComp-medium** pipeline. It contains 128M low-quality data for filtering (successfully download 110M). After obtaining subsets with some data filtering strategies, it will train a CLIP-B/32 model with a fixed budget.

### E1: methods Just utilizing CLIP embedding

| Strategy | IN-1k | VTAB | Avg. |
|---|---|---|---|
| CLIPScore | 26.4 | 32.6 | 32.2 |
| negCLIPLoss (Ours) | 27.9 | 33.2 | 32.9 |
| negCLIPLoss ∩ NormSim₂-D (Ours) | 29.8 | 34.8 | 34.1 |
| negCLIPLoss ∩ NormSim∞ (Ours) | 31.7 | 36.0 | 35.0 |

### E2: Comparing All methods

| Strategy | IN-1k | VTAB | Avg. |
|---|---|---|---|
| DFN [1] | 36.0 | 36.2 | 35.4 |
| Devil [2] | 31.0 | 35.9 | 34.5 |
| DFN ∪ HYPE [3] | 36.4 | 38.5 | 36.8 |
| DFN ∪ Ours | 36.4 | 38.6 | 37.6 |
| DFN ∪ HYPE ∪ Ours | 37.3 | 38.5 | 37.7 |

### E3: Universality : Our methods applied for OpenAI's CLIP-B/32, CLIP-L/14, and the public version of DFN models.

More details in our papers. In Leadboard we are applying our methods on the complete data pool (i.e., 128M data pool)

[1] Fang, A., Jose, A. M., Jain, A., Schmidt, L., Toshev, A., & Shankar, V. (2023). Data filtering networks. *arXiv preprint arXiv:2309.17425.*
[2] Yu, H., Tian, Y., Kumar, S., Yang, L., & Wang, H. (2023). The devil is in the details: A deep dive into the rabbit hole of data filtering. *arXiv preprint arXiv:2309.15954.*
[3] Kim, W., Chun, S., Kim, T., Han, D., & Yun, S. (2024). HYPE: Hyperbolic Entailment Filtering for Underspecified Images and Texts. arXiv preprint arXiv:2404.17507.