

Boosting Medical Vision-Language Pretraining via Momentum Self-Distillation under Limited Computing Resources

Phuc Pham*

Nhu Pham*

Ngoc Quoc Ly

Faculty of Information Technology, University of Science, Ho Chi Minh City, Vietnam
 Vietnam National University, Ho Chi Minh City, Vietnam
 {20120351, 20120153}@student.hcmus.edu.vn, lqngoc@fit.hcmus.edu.vn

Abstract

In medical healthcare, obtaining detailed annotations is challenging, highlighting the need for robust Vision-Language Models (VLMs). Pretrained VLMs enable fine-tuning on small datasets or zero-shot inference, achieving performance comparable to task-specific models. Contrastive learning (CL) is a key paradigm for training VLMs but inherently requires large batch sizes for effective learning, making it computationally demanding and often limited to well-resourced institutions. Moreover, with limited data in healthcare, it is important to prioritize knowledge extraction from both data and models during training to improve performance. Therefore, we focus on leveraging the momentum method combined with distillation to simultaneously address computational efficiency and knowledge exploitation. Our contributions can be summarized as follows: (1) leveraging momentum self-distillation to enhance multimodal learning, and (2) integrating momentum mechanisms with gradient accumulation to enlarge the effective batch size without increasing resource consumption. Our method attains competitive performance with state-of-the-art (SOTA) approaches in zero-shot classification, while providing a substantial boost in the few-shot adaption, achieving over 90% AUC-ROC and improving retrieval tasks by 2–3%. Importantly, our method achieves high training efficiency with a single GPU while maintaining reasonable training time. Our approach aims to advance efficient multimodal learning by reducing resource requirements while improving performance over SOTA methods. The implementation of our method is available at <https://github.com/phphuc612/MSD>.

1. Introduction

Medical imaging modalities are essential for diagnosing and managing a wide range of serious diseases. Integrating deep learning models into healthcare can significantly enhance early disease detection [19], enabling timely treatment and reducing health risks. However, fully supervised models require substantial annotated data, which is often time-consuming and costly to obtain. In contrast, raw radiograph-report data, such as the MIMIC-CXR [12] dataset with over 200,000 pairs, is abundant. This constraint has spurred interest in self-supervised and weakly supervised learning approaches, which reduce reliance on expensive annotations.

Contrastive learning has emerged as a key self-supervised learning strategy that accelerates the development of robust feature representations. It operates by learning representations that are invariant among augmented views of a sample while pushing apart representations of different samples. This is achieved by optimizing the noise contrastive estimation (NCE) loss. Early contrastive learning models focused on improving representation learning in the visual domain. For instance, SimCLR [4] improved performance through complex augmentation strategies and non-linear projection layers, while other models like SimSiam [5], BYOL [8], and MoCo [9] employed different parameter update strategies to enhance training efficiency.

Extending contrastive learning to multi-modality, vision-language models (VLMs) have recently gained traction in medical AI by leveraging paired image-text data to enhance feature representations. Models such as BioViL [2], MedCLIP [24], ConVIRT [28] and Gloria [10] have demonstrated strong performance in medical image-text alignment tasks. These methods align the latent spaces of two modalities by treating one modality as another view of the same data sample. However, these approaches often require large-scale datasets and extensive computational re-

*Contributed Equally.

Contact email: phphuc612@gmail.com

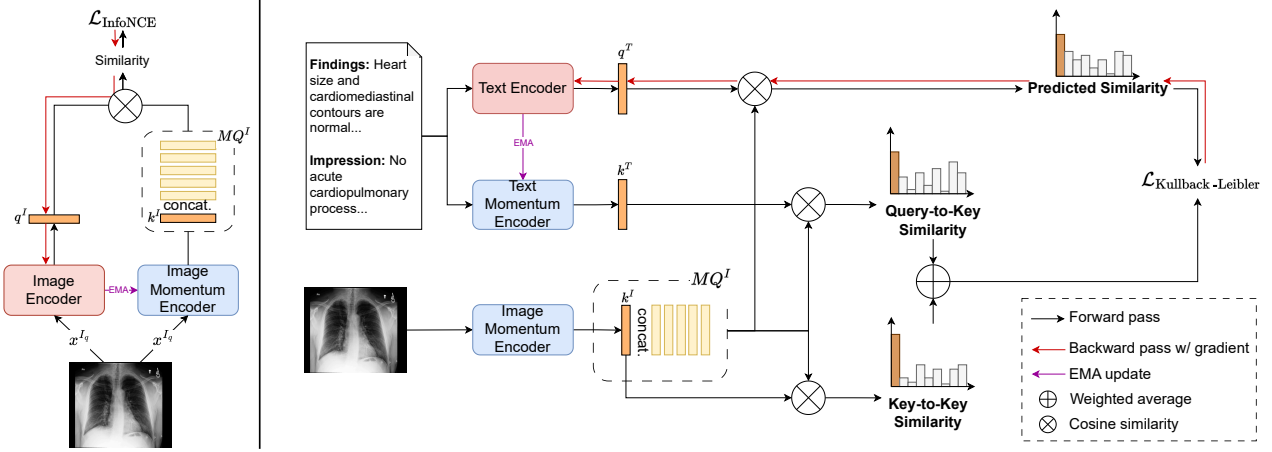


Figure 1. **Our overall framework.** **Left:** Uni-modal learning on images. **Right:** Multi-modal learning on text-to-image. For simplicity, we illustrate our method using a single sample. The same process applies to uni-modal learning on text and multi-modal learning on image-to-text by substituting the corresponding modules.

sources, which limits their practicality in medical AI settings. Additionally, some approaches, such as MedKLIP [25] and MAVL [18], integrate domain-specific knowledge to improve retrieval and classification tasks.

While contrastive learning has proven effective in medical VLMs, further enhancements are needed to address computational constraints and the issue of false negatives. One promising technique to refine learned representations is self-distillation, which allows a model to transfer knowledge to itself for improved performance. Despite its demonstrated success in vision models such as BEiT [1] and DINO [3], self-distillation remains underutilized in contrastive learning, particularly in multimodal medical imaging settings.

Inspired by MoCo [9], we extend momentum contrastive learning to the multimodal domain by introducing dual momentum queues and maintaining separate momentum encoders for images and text, following the approach of Yuan et al. [27], referred to as multi-modal MoCo (MM-MoCo). Based on this architecture, we introduce two key innovations:

- **Momentum Self-Distillation:** We demonstrate that applying self-distillation with a momentum mechanism allows the model to *achieve strong performance even with small batch sizes*, thus alleviating the reliance on large-batch training.
- **Resource-Free Batch Enlargement:** We propose a novel method that exploits the non-gradient nature of momentum to *simulate large batch sizes without requiring additional computational resources*, leading to improved learning efficiency.

Building on these contributions, we propose a novel framework that offers a computationally efficient solution for medical vision-language representation learning. Our em-

pirical results demonstrate improved performance in medical image-text alignment tasks, validating the effectiveness of our approach.

2. Methodology

2.1. Problem Setting

We begin by defining the problem setting in our work. Given a dataset of size N containing image-text pairs, denoted as $D = \{(x_i^I, x_i^T)\}_{i=0}^{N-1}$, where x_i^I and x_i^T represent a medical image and its corresponding text report, respectively. This pairing is a natural characteristic of medical datasets thanks to the routine workflow of radiologists generating textual descriptions of images [12, 28].

Our goal is training image and text encoders so that their latent spaces align. We verify it by transferring the learned text and image embeddings to classification and retrieval tasks, following previous works [24–26].

2.2. Uni-modal Contrastive Learning

Traditionally, end-to-end contrastive learning requires two gradient update streams for both the query and key encoders [4, 20]. In contrast, MoCo [9] updates only the query branch through backpropagation, while the key branch is updated via an exponential moving average (EMA), thus, reducing computational costs. Denoting the parameters of the key branch as θ_k and those of the query branch as θ_q , the momentum update or EMA is: $\theta_k \rightarrow m\theta_k + (1-m)\theta_q$. The coefficient m typically set to a high value (e.g., $m = 0.995$) as suggested by MoCo’s experimental results to ensure gradual updates and minimize discrepancies across different versions of the momentum encoder. Thanks to this momentum mechanism, MoCo [9] enables matching a gradient-encoded query q with a large queue of momentum-encoded

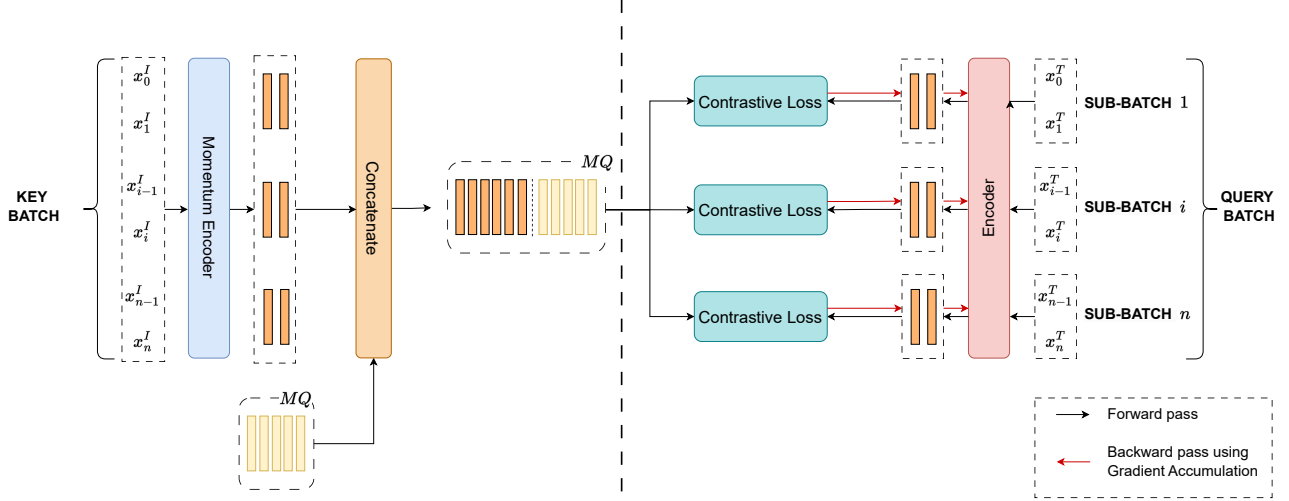


Figure 2. Our technique to increase batch size without additional resources. We split the primary batch into smaller sub-batches. The preparation of embedding vectors is divided into two separate steps: first, calculating and concatenating the momentum keys from sub-batches, and second, calculating the query vectors and optimizing the contrastive loss with prepared keys. In the second step, Gradient Accumulation is employed to achieve the effects of a large batch size.

keys $\{k_0, k_1, \dots\}$.

Without loss of generality, we focus on uni-modal contrastive learning for images, as illustrated in Figure 1. The same process applies to text by replacing the corresponding modules. Let x^{I_q} and x^{I_k} be two transformed views of an input image, encoded into a query vector q^I and a key vector k^I , respectively. Following prior works [4, 6, 7, 9], we define the InfoNCE loss as:

$$\mathcal{L}_{I2I} = -\log \frac{\exp(s(q^I, k^I)/\tau)}{\sum_{k_i^I \in MQ^I} \exp(s(q^I, k_i^I)/\tau)} \quad (1)$$

where τ is a learnable temperature parameter for the softmax function, s is cosine similarity, and MQ^I is the momentum queue storing both the previous key vectors and the current key vector. Similarly, the loss for text is derived in the same manner. The overall loss for uni-modal contrastive learning is then defined as the average of the two modal losses:

$$\mathcal{L}_{\text{uni}} = \frac{\mathcal{L}_{I2I} + \mathcal{L}_{T2T}}{2}. \quad (2)$$

2.3. Multi-modal Contrastive Learning

The multi-modal version of MoCo or MM-MoCo, introduced by Yuan et al. [27], utilizes one-hot or exact label of text-image pair for contrastive learning. However, we empirically observe that this approach performs inadequately under low batch size configurations (e.g., batch size = 16), as demonstrated in our ablation study. This suboptimal performance can be explained by the inherent ambiguity in

textual descriptions and the potential for multiple images within the dataset to match a given textual query. For example, a medical description such as "heart size and cardio-mediastinal contours are normal" could apply to numerous images across the dataset, thereby challenging the effectiveness of exact labels.

To address this limitation, we propose a momentum self-distillation strategy that replaces the exact labels with soft targets derived from similarity distributions. Specifically, we compute two distributions: (1) the *key-to-key similarity* distribution (p_{k2k}), representing similarities between the image key corresponding to the query text and other image keys in the momentum queue, and (2) the *query-to-key similarity* distribution (p_{q2k}), computed between the momentum-encoded query text vector and momentum-encoded image key vectors. This approach effectively establishes direct multi-modal correlations, significantly enhancing performance at low batch sizes and continuing to scale well with larger batch sizes. Crucially, our method reduces reliance on large batch sizes, thus alleviating GPU resource constraints during training.

Both similarity measures can be effectively combined by computing the Kullback-Leibler (KL) divergence between the predicted text-to-image similarity distribution p_{t2i} and the momentum-distilled similarities:

$$\mathcal{L}_{T2I} = \alpha \text{KL}(p_{q2k} \parallel p_{t2i}) + \beta \text{KL}(p_{k2k} \parallel p_{t2i}), \quad (3)$$

where we empirically set $\alpha = 0.3$ and $\beta = 0.7$. The corresponding image-to-text loss \mathcal{L}_{I2T} is defined analogously. The overall multi-modal contrastive loss is the average of

these two:

$$\mathcal{L}_{\text{multi}} = \frac{\mathcal{L}_{T2I} + \mathcal{L}_{I2T}}{2}. \quad (4)$$

Finally, the full pretraining objective is formulated as a weighted sum of uni-modal and multi-modal losses. We observed that the optimization inherently prioritized uni-modal optimization over multi-modal learning, leading to a significant decrease in total loss primarily driven by the uni-modal component. To balance this, we increased the weight of the multi-modal loss, setting $\omega_{\text{uni}} = 1$ and $\omega_{\text{multi}} = 10$:

$$\mathcal{L} = \frac{\omega_{\text{uni}} \cdot \mathcal{L}_{\text{uni}} + \omega_{\text{multi}} \cdot \mathcal{L}_{\text{multi}}}{\omega_{\text{uni}} + \omega_{\text{multi}}}. \quad (5)$$

2.4. Increasing batch size without additional resources

Another key contribution of our work is leveraging the gradient-free nature of the momentum key branch to efficiently increase the effective batch size without additional computational resources.

Our approach is applicable to both uni-modal and multi-modal contrastive learning frameworks, as both share a common structure consisting of a gradient-based query branch and a gradient-free key branch. Specifically, we utilize the momentum encoder’s gradient-free property to compute a large batch of key vectors simultaneously. For the query branch, which requires gradient calculations, we handle large batch sizes by segmenting them into smaller sub-batches that fit within GPU memory constraints.

Once all momentum keys are prepared, we apply gradient accumulation across these smaller sub-batches. Gradients from each sub-batch are accumulated sequentially, and the optimizer updates the model parameters only after processing all sub-batches, as illustrated in Fig. 2. This approach ensures each query interacts simultaneously with the complete set of keys, effectively simulating the behavior of a large batch size.

By structuring computations in this manner, our method significantly scales the batch size without incurring additional computational overhead. Coupled with the momentum self-distillation technique introduced earlier, this strategy achieves performance comparable to current SOTA methods, making it highly practical for resource-constrained environments.

3. Experiments

Following CXR-CLIP [26] as our baseline, we also compare and analyze the performance of our VLM with popular contrastive-learning-based VLMs on three main tasks: zero-shot classification, few-shot classification and retrieval on multiple datasets.

3.1. Implementation details for evaluations

Augmentation A radiographic report’s “Findings” and a summarized “Impression” are treated as augmented views of each other. To enhance data diversity, we use back-translation in six languages using Helsinki-NLP models [23] and rearrange sentences for more context flexibility. For images, we apply transformations—Zoom In, Zoom Out, and Equalize, Cropping, Translation, Rotation—alongside with different image views in the report as augmentations.

Configurations Each modality has query and key versions, with one randomly using the original or augmented sample and the other using the augmented sample. For fair comparison, we follow CXR-CLIP [26] using Swin-Tiny [14] for images and pretrained BERT [2] for text. Momentum queue size is 4096. We achieve the batch size of 512 thanks to our proposed technique. The model is trained for 50 epochs, with AdamW [16] as the optimizer and a cosine-annealing scheduler [15] with learning rate 1e-6. All experiments are conducted on a single NVIDIA RTX 4090 GPU.

Evaluation details We evaluate our method on three primary tasks for foundation models: image-to-text retrieval, zero-shot classification, and few-shot classification.

For **image-to-text retrieval**, performance is measured using the Recall@K (R@K) score, which assesses the model’s ability to capture semantic relationships between images and text by retrieving the exact report within the top- K candidates for a given image. Results are presented in Tab. 4.

For **zero-shot classification**, we examine the capability of the models to identify anomalies across different datasets without any fine-tuning. The evaluation metric is the Area Under the Curve (AUC), as shown in Tab. 2. For CheXpert5x200, we take average performance of different sampling runs as prior work [13] and report accuracy and F1-score.

For **few-shot classification**, each method is fine-tuned using 10% of the training data and evaluated on the test set with AUC as the primary metric. This task assesses the adaptability of the models to new tasks when only limited annotated data is available.

3.2. Datasets

Data Split	Pre-training MIMIC-CXR	Evaluation		
		VinDR	RSNA	SIIM
Train	222,628	14550	18,678	8,422
Valid	1,808	450	4,003	1,808
Test	3,264	3000	4,003	1,807

Table 1. The number of studies for each dataset and split.

Method	Published	Pretraining Dataset(s)	VinDR-CXR		RSNA		SIIM	
			ZS	FS	ZS	FS	ZS	FS
GLoRIA [10]	ICCV 2021	C*	78.0	73.0	80.6	88.2	84.0	91.5
BioViL [2]	ECCV 2022	M	-	-	84.1	86.0	70.3	79.5
ConVIRT [28]	MLHC 2022	M	-	-	79.2	85.4	64.3	80.4
MedCLIP [24]	EMNLP 2022	M, C	82.4	84.9	81.9	88.9	89.0	90.4
MedKLIP [25]	ICCV 2023	M	-	-	86.6	87.1	89.8	89.9
MAVL [18]	CVPR 2024	M	-	-	86.9	87.9	92.0	93.0
CXR-CLIP [26]	MICCAI 2023	M	78.3	84.9	81.3	88.0	85.5	86.9
CXR-CLIP [26]	MICCAI 2023	M, C	82.7	86.1	84.5	88.1	87.9	89.6
Ours	-	M	83.8	91.3	87.1	89.3	82.2	93.4

Table 2. Comparison with SOTA on zero-shot (ZS) and few-shot (FS) classification. M, C, and C* mean MIMIC-CXR, CheXpert with multi-class labels only, and CheXpert with textual medical reports using for pretraining.

The datasets used for evaluations are MIMIC-CXR [12], VinDr-CXR[17], RSNA [21] and SIIM [22]. Our pretraining dataset is MIMIC-CXR. The statistics of datasets is show in Tab. 1.

MIMIC-CXR [12] is a large dataset includes chest x-ray studies. Each study contains one or more image and free-form text report pairs. We use the official training split for pretraining the and test split for image-to-text retrieval. **VinDr-CXR** [17] is an image-bounding box dataset with 22 local labels and 6 global labels. We do not use the label "Other disease", "Other lesion" and labels with less than 10 samples for evaluation, following the work of You et al..

RSNA Pneumonia [21] is a binary image-label dataset, with Pneumonia and Normal as two classes. Train, valid and test set are split by 70%, 15%, 15%, following the work of Huang et al. [10]. This dataset serves as external dataset for comparison of performance between models on classification tasks.

SIIM Pneumothorax (2019) is also a binary label dataset, with Pneumothorax and Normal as classess. Similar to RSNA dataset, we split the dataset into train, valid and test set according to the work of Huang et al. [10].

CheXpert5x200 is a standardized benchmark derived from the CheXpert dataset [11]. It contains 200 training examples per class for five selected thoracic pathologies (Atelectasis, Cardiomegaly, Consolidation, Edema, and Pleural Effusion).

3.3. Comparison with state-of-the-arts (SOTA)

Zero-shot Classification Tab. 2 and Tab. 3 present results on zero-shot and few-shot, which is consistent with previous works [18, 26]. The results indicate that our method achieves competitive performance compared to prior contrastive-based vision-language models on zero-shot classification tasks across most of datasets. Notably, we obtain 83.8% AUC on VinDr-CXR, 87.1% on RSNA, and 82.2% on SIIM, highlighting the strong generalization of our momentum self-distillation framework even without

task-specific fine-tuning. These results validate the effectiveness of replacing exact labels with soft similarity distributions, which better capture the semantic ambiguity of radiology reports.

Method	Published	Accuracy	F1-score
GLoRIA [10]	ICCV 2021	0.50	0.48
CXR-CLIP [26]	MICCAI 2023	0.53	0.51
eCLIP [13]	ECCV 2024	0.57	0.57
Ours	-	0.59	0.58

Table 3. Comparison with SOTA for classification performance on CheXpert5x200 dataset.

Few-shot Classification In Tab. 2, our method also demonstrates substantial performance gains over the zero-shot baseline, confirming the adaptability of the learned representations with limited supervision. Specifically, we achieve 91.3% AUC on VinDr-CXR, 89.3% on RSNA, and 93.4% on SIIM, corresponding to improvements of +7.5%, +2.2%, and +11.2% respectively. These findings show the potentials of momentum self-distillation combined with resource-free batch enlargement in producing highly transferable representations that can be efficiently adapted to new datasets with minimal labeled data.

Method	Pretraining Dataset(s)	R@1	R@5	R@10
GLoRIA [10]	C*	7.2	20.6	30.3
MedCLIP [24]	M, C	1.1	1.4	5.5
CXR-CLIP [26]	M	21.6	48.9	60.2
CXR-CLIP [26]	M, C	19.6	44.2	57.1
Ours	M	23.0	49.2	61.6

Table 4. Comparison with SOTA for image-to-text retrieval on MIMIC-CXR. M, C mean MIMIC-CXR, CheXpert using for pre-training.

Method	Update	MSD	RFBE	BS	#GPUs	Requirements	Peak VRAM	Epoch Time
End-to-End [26]	Grad	-	-	16	1	RTX 4090 24GB	~22 GB	~26 mins
MM-MoCo [27]	Momen	-	-	16	1	RTX 4090 24GB	~9 GB	~30mins
Ours	Momen	✓	-	16	1	RTX 4090 24GB	~9 GB	~30mins
Ours	Momen	✓	-	16	1	RTX 2080Ti 11GB	~9 GB	~100mins
End-to-End [26]	Grad	-	-	128	4	A100 40GB	~32 GB	~26 mins
MM-MoCo [27]	Momen	-	✓	512	1	RTX 4090 24GB	~9 GB	~30mins
Ours	Momen	✓	✓	512	1	RTX 4090 24GB	~9 GB	~30mins
Ours	Momen	✓	✓	512	1	RTX 2080Ti 11GB	~9 GB	~105mins

Table 5. Training configurations for ablation study and reported computational efficiency for methods directly relevant to our work. MSD is Momentum Self-Distillation; RFBE is Resource-Free Batch-Enlargement; BS is training batch size.

Method	Training batch size	VinDr		RSNA		SIIM	
		ZS	FS	ZS	FS	ZS	FS
End-to-End [26]	16	77.3	79.2 (+1.9)	81.1	82.5 (+1.4)	67.4	68.0 (+0.6)
MM-MoCo [27]	16	83.3	83.4 (+0.1)	85.7	85.7 (0.0)	65.4	65.4 (0.0)
Ours	16	82.6	90.3 (+7.7)	87.2	87.3 (+0.1)	81.2	92.5 (+11.3)
End-to-End [26]	128	78.3	84.9 (+6.6)	81.3	88.0 (+6.7)	85.5	86.9 (+1.4)
MM-MoCo [27]	512	81.8	91.5 (+9.7)	86.4	89.2 (+2.8)	73.9	92.6 (+18.7)
Ours	512	83.8	91.3 (+7.5)	87.1	89.3 (+2.2)	82.2	93.4 (+11.2)

Table 6. Ablation study on the effectiveness of Momentum Self-distillation for small and large batch size on classification task. Each few-shot (FS) value is annotated with the difference from its zero-shot (ZS) value; dark green for positive, red for zero or negative difference.

Image-to-text retrieval We compare our work with contrastive models without momentum contrast and distillation. As shown in Tab. 4, our proposed method consistently outperforms existing SOTA vision-language models on the MIMIC-CXR retrieval benchmark. Specifically, our approach improves Recall@1 and Recall@10 by around 1.4%, clearly highlighting the effectiveness of the momentum self-distillation mechanism in capturing richer multimodal representations. Unlike MedCLIP, which decouples image-text pairs and thus achieves relatively lower retrieval performance, our method maintains a strong multimodal alignment throughout training, ensuring superior retrieval.

3.4. Ablation studies

3.4.1. Ablation configurations and computational resources comparison

To evaluate the efficiency of our proposed approach, we compare training cost and scalability across three representative settings: (1) an end-to-end baseline with gradient updates on both branches (CXR-CLIP [26]), (2) a multi-modal MoCo baseline without distillation (MM-MoCo [27]), and (3) our method that combines momentum self-distillation (MSD) with resource-free batch enlargement (RFBE). RFBE leverages the gradient-free property of the momentum encoder to precompute keys and applies gradient accumulation on the query branch, thereby simulating large batch sizes under strict memory budgets. Table 5 summarizes configurations and epoch times on different GPUs.

Additionally, our method achieves an effective batch size of 512 on a single RTX 4090 (24 GB) with epoch time comparable to MM-MoCo, while remaining feasible even on a lower-spec RTX 2080Ti (11 GB). Importantly, our method peaks at only **~9 GB VRAM** usage even under the largest batch configuration, while the end-to-end approach typically requires **at least 32 GB VRAM per GPU** (e.g., on 4×A100) to reach similar batch sizes. For fairness, we also apply RFBE to MM-MoCo, which isolates the batching effect from our distillation mechanism. Together with downstream results reported in Tables 6 and 8, these comparisons demonstrate that MSD consistently boosts performance under both small- and large-batch configurations, while RFBE ensures scalability without additional hardware, validating the practicality of our framework for resource-constrained environments.

3.4.2. The effectiveness of momentum self-distillation

We perform ablation studies on both classification and retrieval tasks to isolate the contribution of MSD under different batch size configurations.

Classification Our model consistently demonstrates substantial performance gains when transitioning from zero-shot to few-shot conditions. Notably, on the SIIM dataset, our method with a small batch size (16) achieves an impressive 11.3% improvement in AUC from zero-shot to few-shot learning. In contrast, other methods such as MM-MoCo [27] and the End-to-End [26] approach exhibit poor

Query-to-key ratio (α)	Key-to-key ratio (β)	VinDr		RSNA		SIIM	
		ZS	FS	ZS	FS	ZS	FS
0.0	1.0	77.9	86.1	85.7	87.6	77.6	92.4
0.3	0.7	83.8	91.3	87.1	89.3	82.2	93.4
0.5	0.5	79.3	87.4	86.5	88.1	78.9	93.2
0.7	0.3	81.2	89.2	86.7	88.1	80.1	92.3
1.0	0.0	Training Failed					

Table 7. Ablation study on query-to-key and key-to-key ratio for classification task.

classification performance at small batch sizes, with **negligible or no improvement even after fine-tuning in the few-shot setting**. With the larger batch size (512), our approach continues to excel, achieving the highest few-shot classification accuracy of 93.4% AUC on the SIIM dataset. These results underline the superior adaptability and effectiveness of our momentum self-distillation and batch enlargement techniques in enhancing model performance, even with limited labeled data.

Image-to-text retrieval The results once again emphasize the critical role of momentum self-distillation in enhancing model performance, particularly under challenging scenarios with limited computational resources. Specifically, at a small batch size of 16, our proposed method achieves a significant performance improvement, attaining a Recall@1 of 22.7%. This result represents a notable advancement compared to traditional methods such as MM-MoCo [27], which only achieves 3.5%, and the End-to-End [26] approach, which reaches 10.9%. Moreover, the effectiveness of our momentum self-distillation technique becomes even more pronounced at larger batch sizes (512), where our method attains the highest retrieval performance with Recall@1 of 23.0%, clearly demonstrating its consistent effectiveness across varying batch sizes.

Method	Training batch size	R@1	R@5	R@10
End-to-End [26]	16	10.9	27.6	37.2
MM-MoCo [27]	16	3.5	11.6	17.3
Ours	16	22.7	48.4	59.6
End-to-End [26]	128	21.6	48.9	60.2
MM-MoCo [27]	512	22.1	49.2	59.5
Ours	512	23.0	49.2	61.6

Table 8. Ablation study on the effectiveness of Momentum Self-distillation for small and large batch size in the retrieval task on the MIMIC-CXR dataset.

3.4.3. Analysis of the momentum self-distillation formula

The MSD loss in Eq. (3) is defined as a weighted sum of two KL divergences: between the predicted similarity dis-

Query-to-key ratio (α)	Key-to-key ratio (β)	R@1	R@5	R@10
0.0	1.0	23.7	49.0	60.3
0.3	0.7	23.0	49.2	61.6
0.5	0.5	21.7	47.7	59.9
0.7	0.3	19.3	46.3	58.8
1.0	0.0	Training Failed		

Table 9. Ablation study on query-to-key and key-to-key ratio for retrieval task.

tribution and (1) the *query-to-key* distribution p_{q2k} , and (2) the *key-to-key* distribution p_{k2k} . The coefficients α and β control the relative influence of these two signals.

Tables 7 and 9 highlight two important observations. First, we find that the key-to-key signal alone is sufficient for stable training, whereas the query-to-key signal alone causes model divergence. This can be explained by the fact that the key-to-key signal acts as a soft version of hard labels: text-to-text or image-to-image pairs naturally achieve the highest cosine similarity when and only when they correspond to the same instance. In contrast, the query-to-key signal does not inherently possess this property, making it unstable if used independently. However, once the model begins to learn meaningful alignments, p_{q2k} provides additional multimodal supervision that can be fed back as a self-distillation signal when combined with p_{k2k} at an appropriate ratio. These results suggest that p_{k2k} offers a stable global structure anchored in the smoothly updated momentum encoder, while p_{q2k} contributes finer multimodal alignment but requires stabilization from p_{k2k} . The asymmetric weighting ($\beta > \alpha$) therefore achieves the best balance, ensuring stable convergence and maximizing both classification and retrieval performance.

4. Conclusion

In this work, we introduced a novel, resource-efficient contrastive learning framework designed for medical vision-language representation learning. By combining momentum self-distillation with a resource-free batch-enlargement strategy, we demonstrated significant gains in retrieval and classification performance across various datasets. Our

method effectively mitigates the limitations posed by small batch sizes and computational constraints, enabling strong performance even with minimal GPU resources. Extensive experiments confirm that the proposed approach outperforms existing baselines, particularly under few-shot and low-resource settings.

Nevertheless, our approach is positioned as a foundational model rather than a directly deployable application. While it effectively captures rich multimodal representations, further work is needed to tailor these representations for specific downstream clinical tasks. Additionally, future research should explore more extensive fine-tuning on downstream targets, and the integration of generative objectives to enhance model interpretability and clinical utility.

Our work is aimed to lay a foundation for building scalable and efficient multimodal models in medical AI and opens new avenues for further exploration in low-resource and domain-adaptive learning.

Acknowledgements

This research is supported by research funding from Faculty of Information Technology, University of Science, Viet Nam National University – Ho Chi Minh City.

References

- [1] Hangbo Bao, Li Dong, Songhao Piao, and Furu Wei. Beit: Bert pre-training of image transformers. *arXiv preprint arXiv:2106.08254*, 2021. 2
- [2] Benedikt Boecking, Naoto Usuyama, Shruthi Bannur, Daniel C Castro, Anton Schwaighofer, Stephanie Hyland, Maria Wetscherek, Tristan Naumann, Aditya Nori, Javier Alvarez-Valle, et al. Making the most of text semantics to improve biomedical vision–language processing. In *European conference on computer vision*, pages 1–21. Springer, 2022. 1, 4, 5
- [3] Mathilde Caron, Hugo Touvron, Ishan Misra, Hervé Jégou, Julien Mairal, Piotr Bojanowski, and Armand Joulin. Emerging properties in self-supervised vision transformers. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 9650–9660, 2021. 2
- [4] Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework for contrastive learning of visual representations. In *International conference on machine learning*, pages 1597–1607. PMLR, 2020. 1, 2, 3
- [5] Xinlei Chen and Kaiming He. Exploring simple siamese representation learning. 2021 ieee. In *CVF conference on computer vision and pattern recognition (CVPR)*, pages 15745–15753, 2020. 1
- [6] Xinlei Chen, Haoqi Fan, Ross Girshick, and Kaiming He. Improved baselines with momentum contrastive learning. *arXiv preprint arXiv:2003.04297*, 2020. 3
- [7] Tianyu Gao, Xingcheng Yao, and Danqi Chen. Simcse: Simple contrastive learning of sentence embeddings. *arXiv preprint arXiv:2104.08821*, 2021. 3
- [8] Jean-Bastien Grill, Florian Strub, Florent Altché, Corentin Tallec, Pierre Richemond, Elena Buchatskaya, Carl Doersch, Bernardo Avila Pires, Zhaohan Guo, Mohammad Gheshlaghi Azar, et al. Bootstrap your own latent-a new approach to self-supervised learning. *Advances in neural information processing systems*, 33:21271–21284, 2020. 1
- [9] Kaiming He, Haoqi Fan, Yuxin Wu, Saining Xie, and Ross Girshick. Momentum contrast for unsupervised visual representation learning. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 9729–9738, 2020. 1, 2, 3
- [10] Shih-Cheng Huang, Liyue Shen, Matthew P Lungren, and Serena Yeung. Gloria: A multimodal global-local representation learning framework for label-efficient medical image recognition. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 3942–3951, 2021. 1, 5
- [11] Jeremy Irvin, Pranav Rajpurkar, Michael Ko, Yifan Yu, Silvana Ciurea-Ilcus, Chris Chute, Henrik Marklund, Behzad Haghighi, Robyn Ball, Katie Shpanskaya, et al. Chexpert: A large chest radiograph dataset with uncertainty labels and expert comparison. In *Proceedings of the AAAI conference on artificial intelligence*, number 01, pages 590–597, 2019. 5
- [12] Alistair EW Johnson, Tom J Pollard, Seth J Berkowitz, Nathaniel R Greenbaum, Matthew P Lungren, Chih-ying Deng, Roger G Mark, and Steven Horng. Mimic-cxr, a de-identified publicly available database of chest radiographs with free-text reports. *Scientific data*, 6(1):317, 2019. 1, 2, 5
- [13] Yogesh Kumar and Pekka Marttinen. Improving medical multi-modal contrastive learning with expert annotations. In *European Conference on Computer Vision*, pages 468–486. Springer, 2024. 4, 5
- [14] Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining Guo. Swin transformer: Hierarchical vision transformer using shifted windows. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 10012–10022, 2021. 4
- [15] Ilya Loshchilov and Frank Hutter. Sgdr: Stochastic gradient descent with warm restarts. *arXiv preprint arXiv:1608.03983*, 2016. 4
- [16] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. *arXiv preprint arXiv:1711.05101*, 2017. 4
- [17] Ha Q Nguyen, Khanh Lam, Linh T Le, Hieu H Pham, Dat Q Tran, Dung B Nguyen, Dung D Le, Chi M Pham, Hang TT Tong, Diep H Dinh, et al. Vindr-cxr: An open dataset of chest x-rays with radiologist’s annotations. *Scientific Data*, 9(1):429, 2022. 5
- [18] Vu Minh Hieu Phan, Yutong Xie, Yuankai Qi, Lingqiao Liu, Liyang Liu, Bowen Zhang, Zhibin Liao, Qi Wu, Minh-Son To, and Johan W Verjans. Decomposing disease descriptions for enhanced pathology detection: A multi-aspect vision-language pre-training framework. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11492–11501, 2024. 2, 5
- [19] Chunli Qin, Demin Yao, Yonghong Shi, and Zhijian Song. Computer-aided detection in chest radiography based on ar-

- tificial intelligence: a survey. *Biomedical engineering online*, 17:1–23, 2018. 1
- [20] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pages 8748–8763. PMLR, 2021. 2
- [21] George Shih, Carol C Wu, Safwan S Halabi, Marc D Kohli, Luciano M Prevedello, Tessa S Cook, Arjun Sharma, Judith K Amorosa, Veronica Arteaga, Maya Galperin-Aizenberg, et al. Augmenting the national institutes of health chest radiograph dataset with expert annotations of possible pneumonia. *Radiology: Artificial Intelligence*, 1(1): e180041, 2019. 5
- [22] SIIM. Society for imaging informatics in medicine: Siim-acr pneumothorax segmentation. <https://www.kaggle.com/c/siim-acr-pneumothorax-segmentation>, 2019. 5
- [23] Jörg Tiedemann, Mikko Aulamo, Daria Bakshandaeva, Michele Boggia, Stig-Arne Grönroos, Tommi Nieminen, Alessandro Raganato, Yves Scherrer, Raul Vazquez, and Sami Virpioja. Democratizing neural machine translation with OPUS-MT. *Language Resources and Evaluation*, (58): 713–755, 2023. 4
- [24] Zifeng Wang, Zhenbang Wu, Dinesh Agarwal, and Jimeng Sun. Medclip: Contrastive learning from unpaired medical images and text. *arXiv preprint arXiv:2210.10163*, 2022. 1, 2, 5
- [25] Chaoyi Wu, Xiaoman Zhang, Ya Zhang, Yanfeng Wang, and Weidi Xie. Medklip: Medical knowledge enhanced language-image pre-training for x-ray diagnosis. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 21372–21383, 2023. 2, 5
- [26] Kihyun You, Jawook Gu, Jiyeon Ham, Beomhee Park, Jiho Kim, Eun K Hong, Woonhyuk Baek, and Byungseok Roh. Cxr-clip: Toward large scale chest x-ray language-image pre-training. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, pages 101–111. Springer, 2023. 2, 4, 5, 6, 7
- [27] Xin Yuan, Zhe Lin, Jason Kuen, Jianming Zhang, Yilin Wang, Michael Maire, Ajinkya Kale, and Baldo Faieta. Multimodal contrastive training for visual representation learning. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 6995–7004, 2021. 2, 3, 6, 7
- [28] Yuhao Zhang, Hang Jiang, Yasuhide Miura, Christopher D Manning, and Curtis P Langlotz. Contrastive learning of medical visual representations from paired images and text. In *Machine Learning for Healthcare Conference*, pages 2–25. PMLR, 2022. 1, 2, 5