

# Intermediate Layers Matter in Momentum Contrastive Self-Supervised Learning



Aakash Kaku



Sahana Upadhya

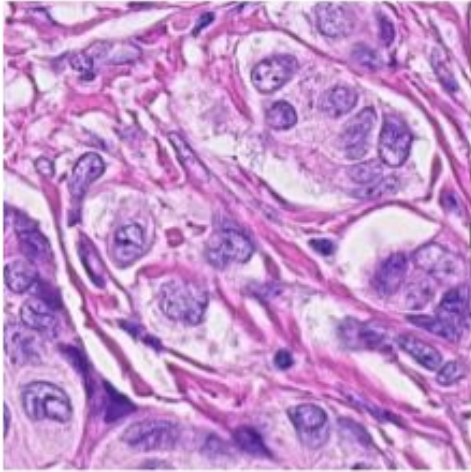


Narges Razavian

# Self-supervised learning can be useful to learn from unlabeled medical data

- SSL methods helps to learn useful representations from large unlabeled data.
  - Recent SOTA methods use contrastive based loss function:
    - MoCo [He et al., 2020]
    - BYOL [Grill et al., 2020]
    - Barlow Twins [Zbontar et al.,2021]
    - Dino [Caron et al., 2021]
- 
- In medical domain, with advent of technology, large amounts of data are collected.
  - Labeled data are limited as labeling is expensive and time consuming.
  - Ideal case for using SSL methods.

# MoCo has been widely used in the medical domain



[Ciga et al., 2020, Dehaene et al., 2020]



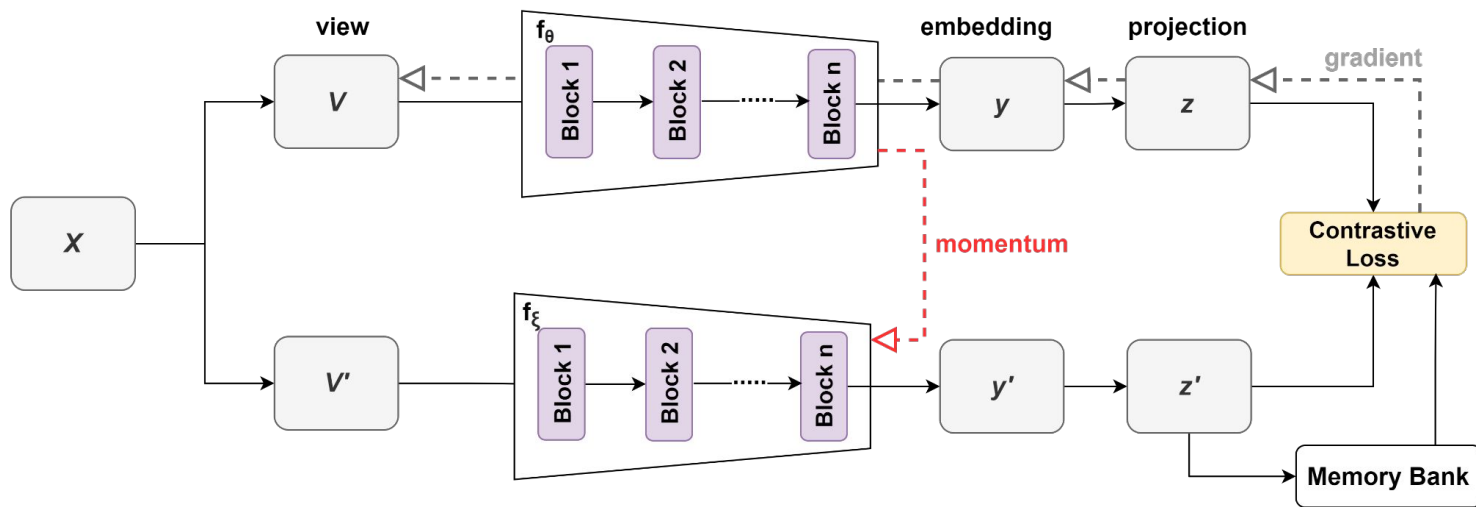
[Sowrirajan et al., 2021]



[Azizi et al., 2021]

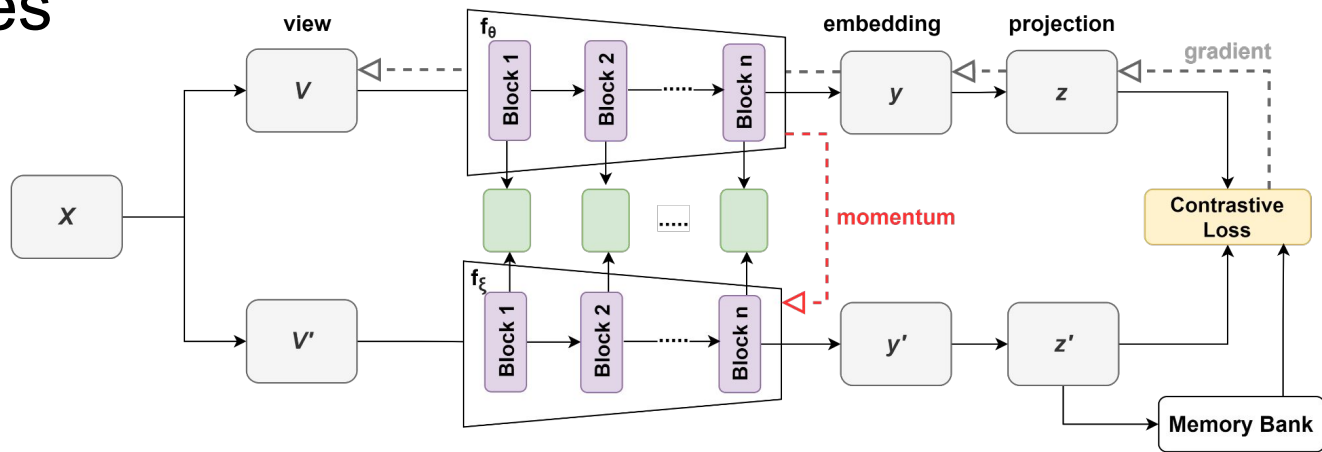
- Medical images are typically large (e.g. 1024 x 1024). MoCo allows to train to with small batches
- Used for histopathology classification task, chest x-ray interpretation, dermatology classification

# Review of standard MoCo



$$\mathcal{L}_{\text{InfoNCE}}(x) = -\log \frac{\exp(z \cdot z' / \tau)}{\exp(z \cdot z' / \tau) + \sum_{i=0}^K \exp(z \cdot z'_i \text{memory bank}) / \tau}$$

# Proposed MoCo: Enforcing similarity between intermediate features



 MSE/Barlow Twins Loss

$$\mathcal{L} = \mathcal{L}_{\text{InfoNCE}} + \mathcal{L}_{\text{MSE or BT}}$$

$$\mathcal{L}_{\text{BT}} = \underbrace{\sum_i (1 - C_{ii})^2}_{\text{invariance term}} + \lambda \underbrace{\sum_i \sum_{j \neq i} C_{ij}^2}_{\text{redundancy reduction term}}$$

$$C_{ij} = \frac{\sum_b z_{b,i}^A z_{b,j}^B}{\sqrt{\sum_b (z_{b,i}^A)^2} \sqrt{\sum_b (z_{b,j}^B)^2}}$$

Zbontar et al. [2021]

# Evaluation of models - Linear probing

- Linear classifier probing is the most standard way to evaluate SSL methodologies
  - A linear classifier is trained on the features generated by the SSL method. Rest of the network is frozen.
  - The test accuracy of the linear classifier is used as a proxy for the performance of SSL methodology.

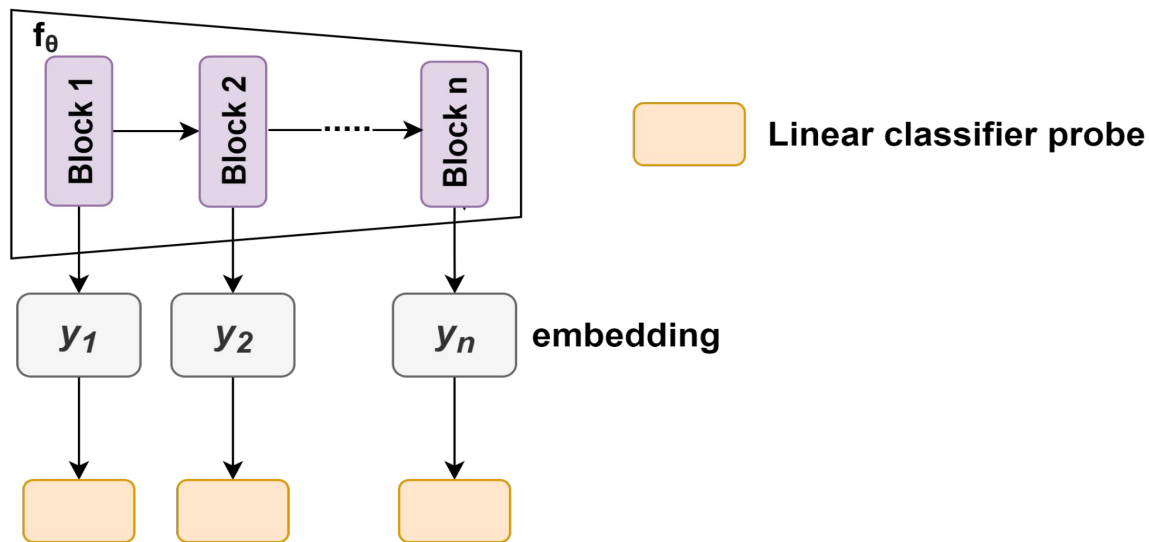
## **Limitations:**

- For most of the medical imaging applications, model is typically fully fine tuned. The network is not frozen.
- It is shown that for such applications, unlike natural images, linear probing yield models that have very low performance and hence become meaningless as a proxy for evaluating SSL methods.

# Beyond linear probing - Feature reuse

- Taking inspiration from transfer learning, we propose feature reuse as a metric to evaluate SSL methods that are fully fine-tuned [Neyshabur et al., 2021].
- Feature reuse is measured by measuring similarity of features before and after fully fine tuning the model.
- Intuitively, if two ssl methods achieve similar performance on the downstream task after fully fine-tuning the model, the model with higher feature similarity has learned useful features during the SSL phase.
- Similarity is measured using Centered kernel alignment [Kornblith et al., 2019]
  - Invariant to orthogonal transformations.
  - Invariant to invertible linear transformation
  - Invariant to isotronic scaling.

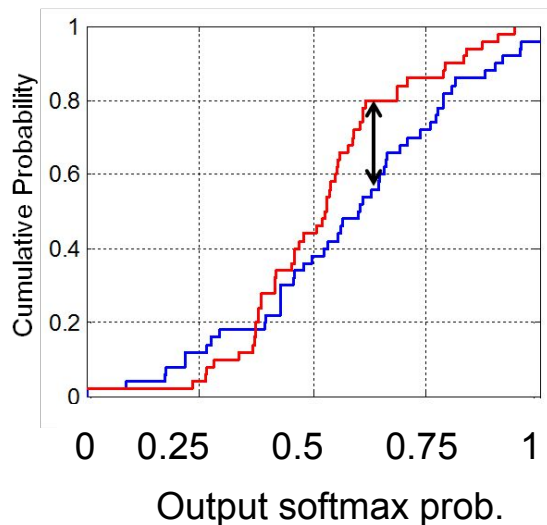
# Beyond linear probing - Layer-wise probing



- Probing the intermediate layers helps to understand how informative the intermediate features are for the downstream task.
- We expect the proposed approach to learn more informative intermediate features as compared to the standard MoCo



# Beyond linear probing - KS distance



- KS distance is the greatest separation between the two cumulative distribution functions (CDFs)
- We compare the distribution of output softmax prob. of a SSL model fine-tuned on fraction of labeled data (like 1% or 6%) and a SSL model fine-tuned on the entire labeled data.
- Lower KS distance signifies both the models are similar in performance (classification and calibration)

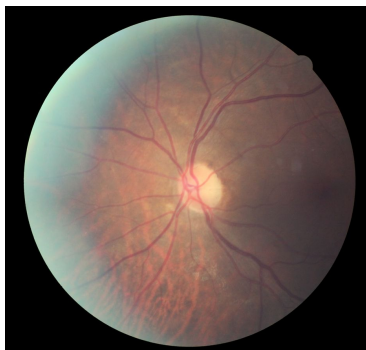
# Datasets

NIH Chest X Ray



# = 112,120; 14 classes  
Multi-class classification;  
Metric = mac - AUC

Diabetic Retinopathy



# = 88,702; Binary  
classification;  
Metric = mac - AUC

Breast Cancer Histopathology



# = 277,524; Binary classification;  
Metric = F1-measure

- We divide the data into training (70%), validation (20%) and testing (10%) set
- For training during SSL phase, we use the entire training set.
- We fine tune the model using different fraction of labeled training data i.e. 1%, 6% and 100%

# Results: Linear probing results

Dataset / Method	MoCo	MoCo + MSE	MoCo + Barlow Twins	Supervised
NIH Chest X-ray (AUC (95% CI))	74.4 (73.9-75.0)	<b>74.8</b> <b>(74.2-75.4)</b>	73.5 (72.9-74.0)	79.8 (79.2-80.3)
Diabetic Retinopathy (AUC (95% CI))	74.6 (74.5-74.7)	<b>84.8</b> <b>(84.6-85.0)</b>	79.7 (79.6-79.7)	94.1 (94.1-94.2)
Breast Cancer Histopathology (F1-score (95% CI))	80.7 (80.4-81.1)	<b>82.5</b> <b>(82.2 - 82.9)</b>	82.3 (82.0-82.7)	82.7 (82.4-83.1)

# Results: Fully-fine tuned results

Label fraction	Supervised	MoCo	MoCo + MSE	MoCo + Barlow Twins
NIH Chest X-ray (AUC (95% CI))				
100%	79.8 (79.2-80.3)	<b>82.4 (81.7-83.0)</b>	81.5 (80.9-82.1)	80.0 (79.5-80.7)
6%	65.2 (64.6-65.8)	69.8 (69.3-70.4)	<b>70.5 (69.9-71.0)</b>	70.0 (69.2-70.6)
1%	57.8 (57.2-58.4)	59.2 (58.6-59.9)	61.4 (60.7-62.0)	<b>62.9 (62.3-63.5)</b>
Diabetic Retinopathy (AUC (95% CI))				
100%	94.1 (94.1-94.2)	94.6 (94.3-94.6)	<b>96.6 (96.6-96.7)</b>	95.7 (95.7-95.8)
6%	69.1 (69.0-69.2)	92.4 (92.2-92.6)	<b>95.1 (94.8-95.2)</b>	94.0 (94.0-94.3)
1%	65.5 (65.4-65.6)	88.1 (88.1-88.4)	<b>93.6 (93.2-93.6)</b>	92.5 (92.2-92.7)
Breast Cancer Histopathology (F1-score (95% CI))				
100%	82.7 (82.4-83.1)	82.9 (82.6-83.3)	85.7 (85.4-86.0)	<b>86.4 (86.1-86.7)</b>
6%	82.7 (82.4-83.1)	82.8 (82.4-83.2)	<b>84.6 (84.2-84.9)</b>	84.5 (84.2-84.8)
1%	80.6 (80.3-81.0)	82.8 (82.5-83.2)	<b>85.1 (84.7-85.4)</b>	84.4 (84.1-84.7)

# Results: Feature reuse

1% labeled data					
Method	Block 1	Block 2	Block 3	Block 4	Performance
NIH Chest X-ray (Performance in AUC)					
MoCo	0.81	0.80	0.57	0.41	59.2
MoCo + MSE	0.97	0.83	0.65	<b>0.42</b>	61.4
MoCo + Barlow Twins	<b>0.99</b>	<b>0.98</b>	<b>0.76</b>	0.38	62.9
Diabetic Retinopathy (Performance in AUC)					
MoCo	0.87	0.80	<b>0.51</b>	0.19	88.1
MoCo + MSE	0.96	0.78	0.33	<b>0.26</b>	93.6
MoCo + Barlow Twins	<b>0.98</b>	<b>0.83</b>	0.58	0.24	92.5
Histopathology (Performance in F1-score)					
MoCo	0.50	0.55	<b>0.98</b>	0.16	82.8
MoCo + MSE	<b>0.77</b>	<b>0.82</b>	0.58	<b>0.42</b>	85.1
MoCo + Barlow Twins	<b>0.77</b>	0.74	0.54	0.36	84.4

# Results: Layer-wise probing

	Block 1	Block 2	Block 3	Block 4
	NIH Chest X-ray (AUC (95% CI))			
MoCo	<b>58.8 (58.4-59.3)</b>	59.5 (59.0-60.0)	65.3 (64.8-65.8)	74.4 (73.9-75.0)
MoCo + MSE	57.6 (56.9-58.3)	<b>59.9 (59.4-60.4)</b>	<b>69.2 (68.70-69.7)</b>	<b>74.8 (74.2-75.4)</b>
MoCo + Barlow Twins	56.6 (56.2-57.0)	56.5 (56.1-56.9)	64.2 (63.7-64.6)	73.5 (72.9-74.0)
	Diabetic Retinopathy (AUC (95% CI))			
MoCo	68.1 (68.0-68.1)	68.2 (68.2-68.3)	69.2 (69.2-69.5)	74.6 (74.5-74.7)
MoCo + MSE	<b>68.3 (68.2-68.3)</b>	<b>70.1 (70.0-70.1)</b>	<b>71.2 (71.1-71.3)</b>	<b>84.8 (84.6-85.0)</b>
MoCo + Barlow Twins	67.2 (67.2-67.3)	68.6 (68.5-68.7)	69.9 (69.4-69.9)	79.7 (79.6-79.7)
	Breast Cancer Histopathology (F1-score (95% CI))			
MoCo	80.9 (80.5-81.3)	81.1 (80.8-81.5)	81.1 (80.7-81.5)	80.7 (80.4-81.1)
MoCo + MSE	80.6 (80.2-81.0)	<b>81.3 (81.0-81.7)</b>	<b>82.7 (82.4-83.0)</b>	<b>82.5 (82.2-82.9)</b>
MoCo + Barlow Twins	<b>81.0 (90.7-81.4)</b>	81.1 (80.7-81.5)	82.2 (81.9-82.6)	82.3 (82.0-82.7)



# Results: KS distance

Label fraction	Supervised	MoCo	MoCo + MSE	MoCo + Barlow Twins
NIH Chest X-ray (Compared to MoCo - Fine tuned on 100% labeled data)				
6%	0.040	<b>0.028</b>	0.039	0.034
1%	0.260	0.244	<b>0.094</b>	0.104
Diabetic Retinopathy (Compared to MoCo + MSE - Fine tuned on 100% labeled data)				
6%	0.37	0.21	<b>0.12</b>	0.30
1%	0.60	0.31	<b>0.18</b>	0.22
Breast Cancer Histopathology (Compared to MoCo + Barlow twins - Fine tuned on 100% labeled data)				
6%	0.040	0.082	0.036	<b>0.026</b>
1%	0.155	0.043	0.082	<b>0.033</b>

# Future work

- SSL approaches are typically designed for natural images. In our work, we tried to built models exclusively for medical datasets.
- Having said that, the method proposed in our work is general and can be adapted to different model architectures, SSL methods and datasets.
- In future, we would like to investigate the effectiveness of our proposed method for other datasets, and SSL methods.



# Intermediate Layers Matter in Momentum Contrastive Self-Supervised Learning

Aakash Kaku, Sahana Upadhyaya, Narges Razavian

Github link:

[https://github.com/aakashrkaku/intermdiate\\_layer\\_matter\\_ssl](https://github.com/aakashrkaku/intermdiate_layer_matter_ssl)



NYU