



# Deep Learning for Vision & Language

Self-supervised Models for Computer Vision



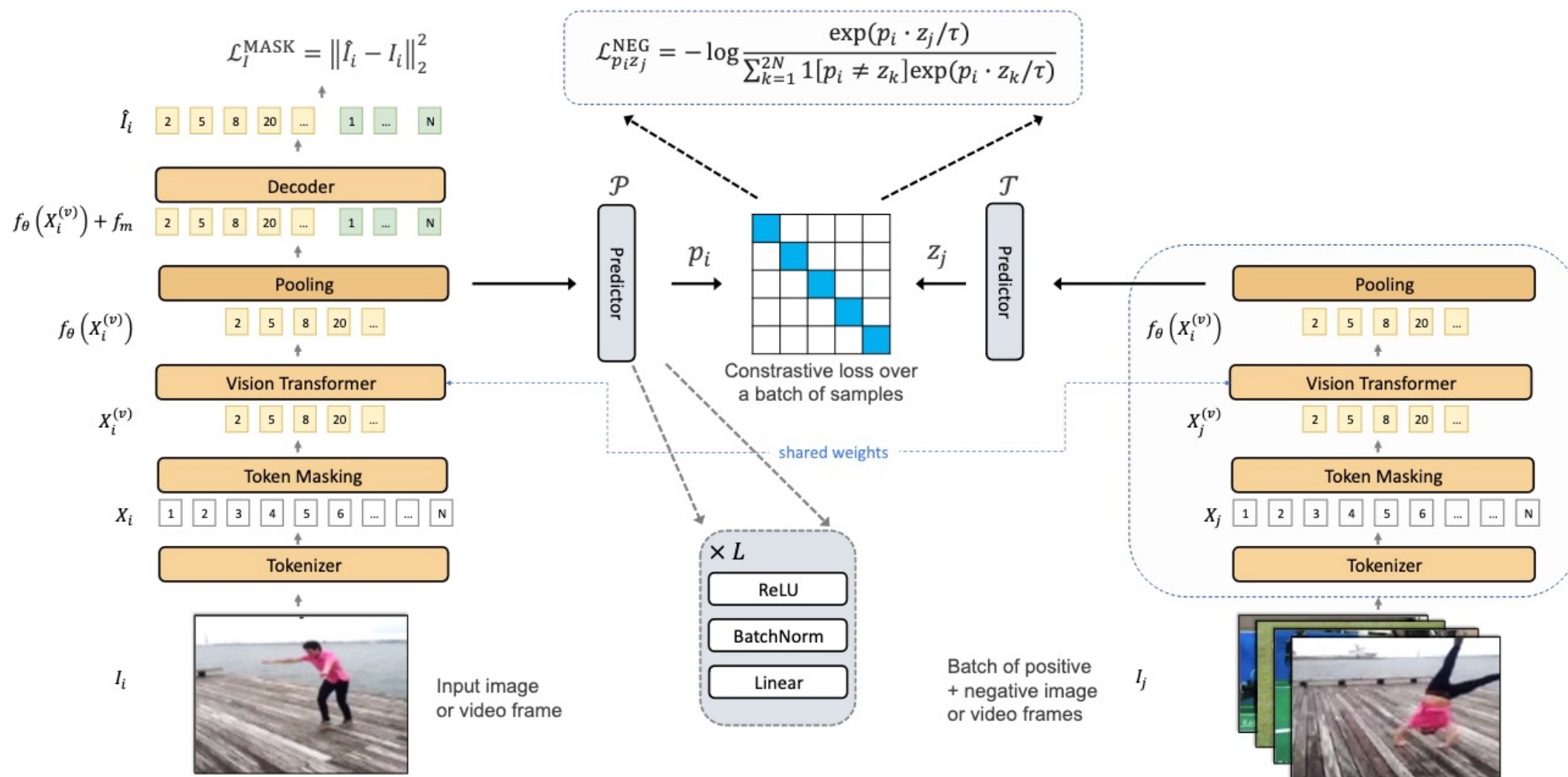
# **ViC-MAE: Self-Supervised Representation Learning from Images and Video with Contrastive Masked Autoencoders**

Jefferson Hernandez<sup>1</sup>, Ruben Villegas<sup>2</sup>, Vicente Ordonez<sup>1</sup>

<sup>1</sup>Rice University, <sup>2</sup>Google DeepMind

{jefebern, vicenteor}@rice.edu, rubville@google.com

# ViC-MAE



# Self-Supervision for Visual Model Learning

- Lots of data but no labels
- Labeling data is expensive

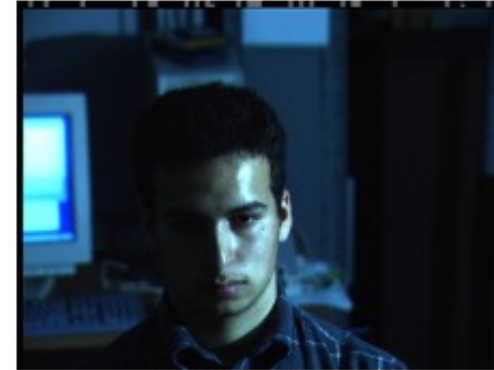
# Similarity Learning: Triplet Loss (Supervised)



$x_i^a$



$x_i^p$



$x_i^n$

$$\sum_i^N \left[ \|f(x_i^a) - f(x_i^p)\|_2^2 - \|f(x_i^a) - f(x_i^n)\|_2^2 + \alpha \right]_+$$

**FaceNet: A Unified Embedding for Face Recognition and Clustering**

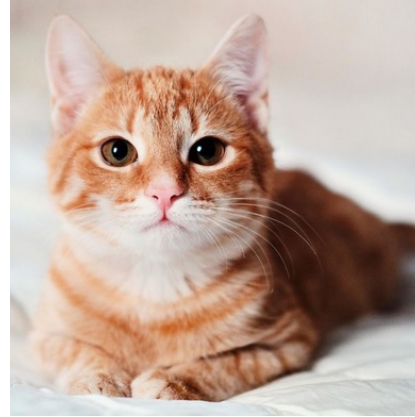
# Similarity Learning: Triplet Loss (Self Supervised)



$x_i^a$



$x_i^p$

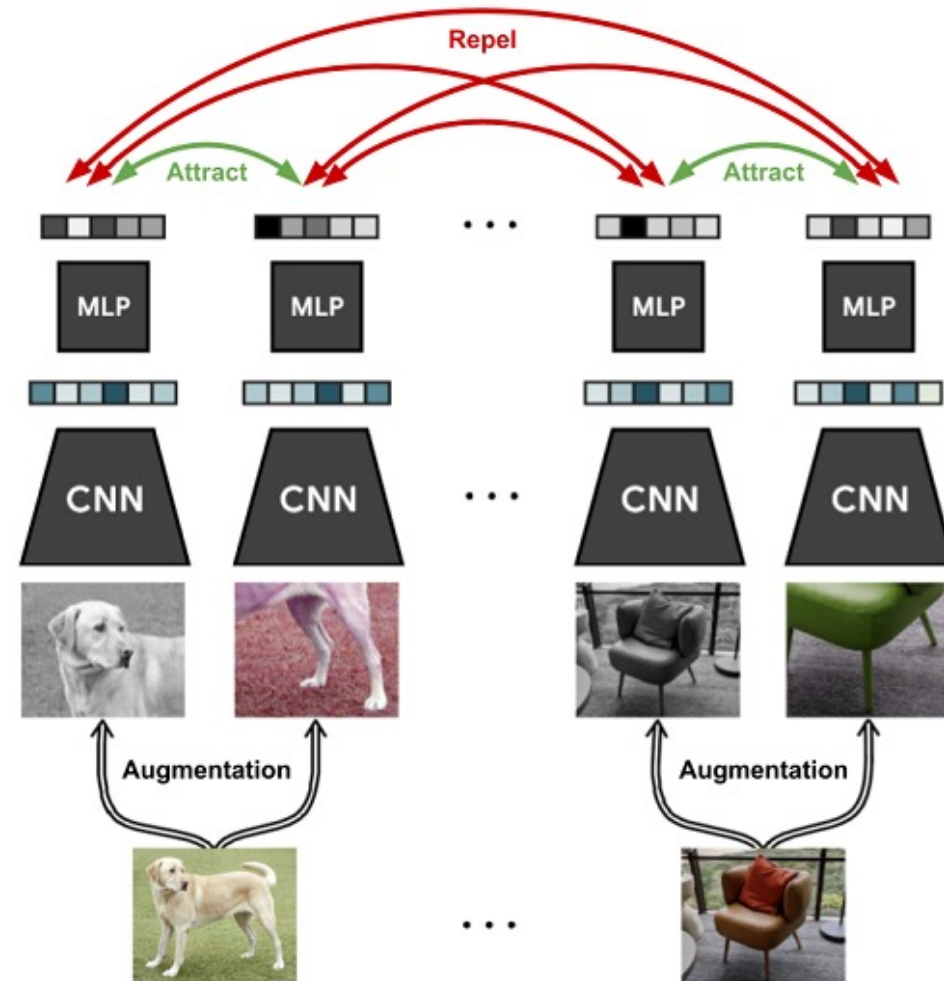


$x_i^n$

$$\sum_i^N \left[ \|f(x_i^a) - f(x_i^p)\|_2^2 - \|f(x_i^a) - f(x_i^n)\|_2^2 + \alpha \right]_+$$

## FaceNet: A Unified Embedding for Face Recognition and Clustering

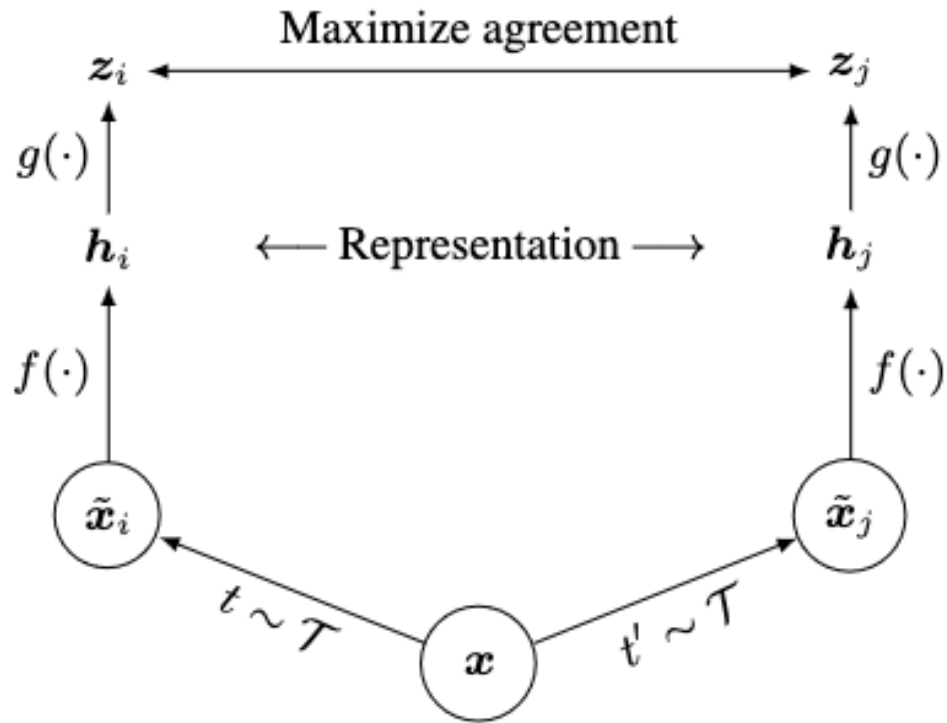
# SimCLR: Contrastive Learning



**A Simple Framework for Contrastive Learning of Visual Representations**

<https://arxiv.org/abs/2002.05709>

# Contrastive Learning



$$h_i = f(\tilde{x}_i) = \text{ResNet}(\tilde{x}_i)$$

$$z_i = g(h_i) = W^{(2)}\sigma(W^{(1)}h_i)$$

$$\ell_{i,j} = -\log \frac{\exp(\text{sim}(z_i, z_j)/\tau)}{\sum_{k=1}^{2N} \mathbb{1}_{[k \neq i]} \exp(\text{sim}(z_i, z_k)/\tau)}$$

Positive pair

Negative pairs

**A Simple Framework for Contrastive Learning of Visual Representations**

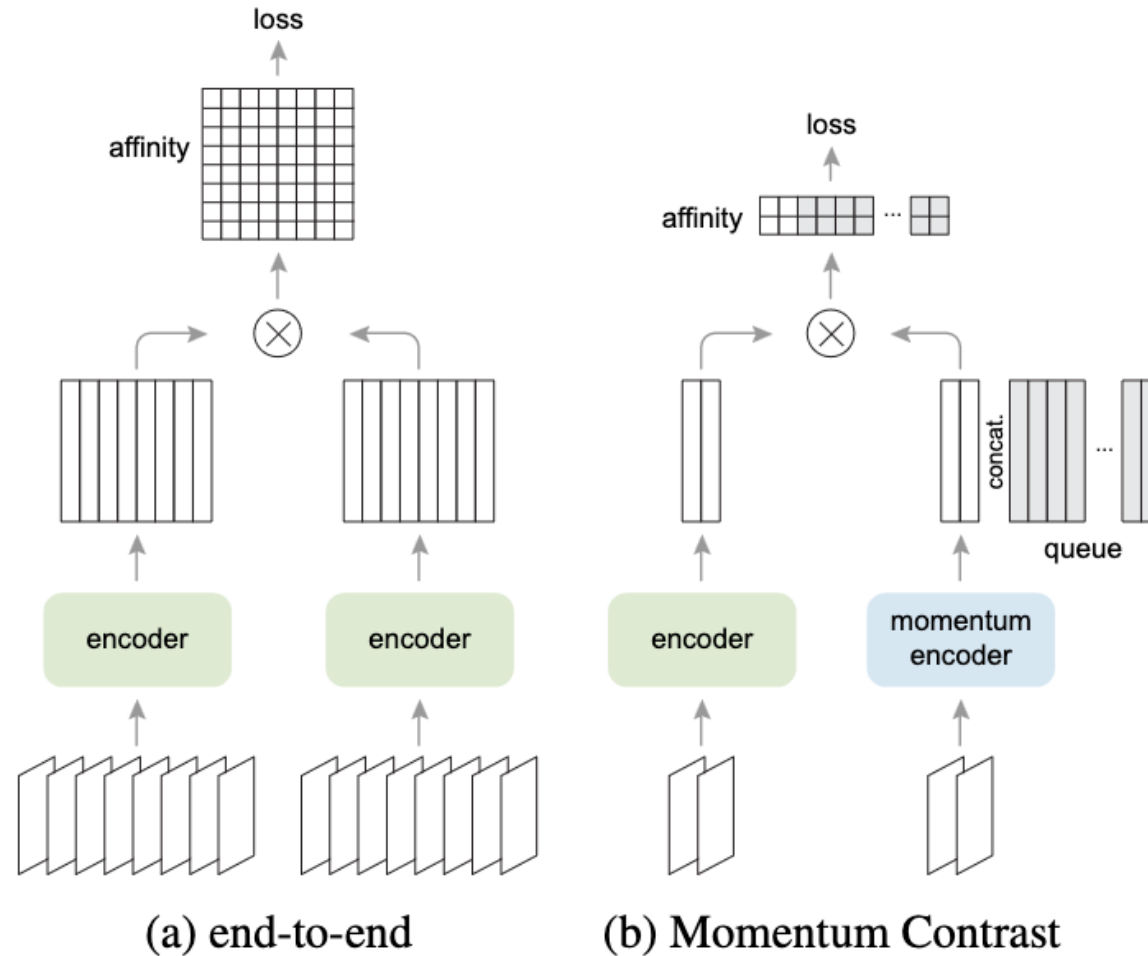
<https://arxiv.org/abs/2002.05709>



# Issue with Contrastive Learning

- Large number of negative examples in the denominator are needed
- In practice this is approximated through batches – all other elements in the batch are negatives.
- Random negative examples are easy
- Large number of negative examples means larger batch sizes which occupy more memory
- GPU memory is expensive  
(What is the max memory a GPU has these days vs your PC?)

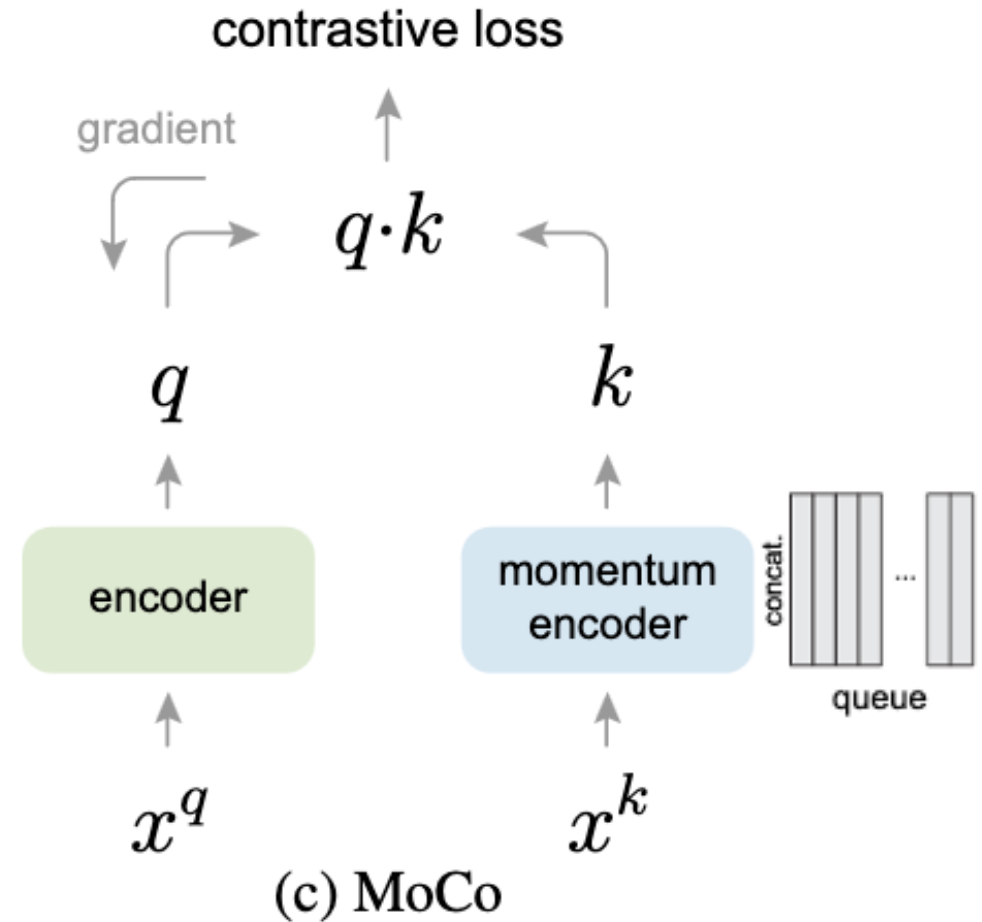
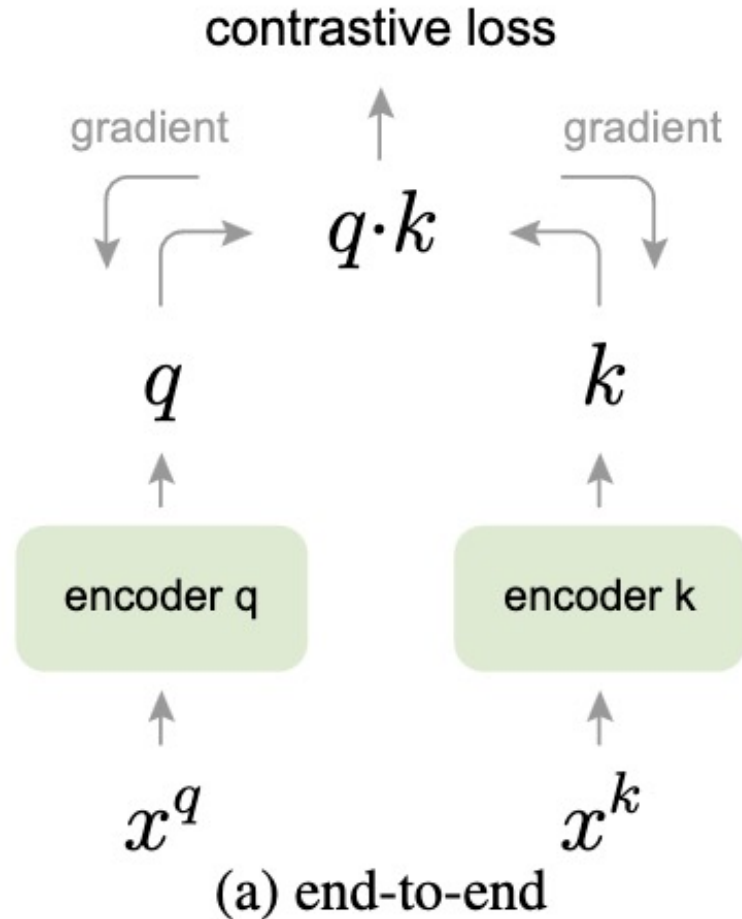
# Momentum Contrastive Learning (MoCo)



## Improved Baselines with Momentum Contrastive Learning

<https://arxiv.org/abs/2003.04297>    <https://arxiv.org/abs/1911.05722>

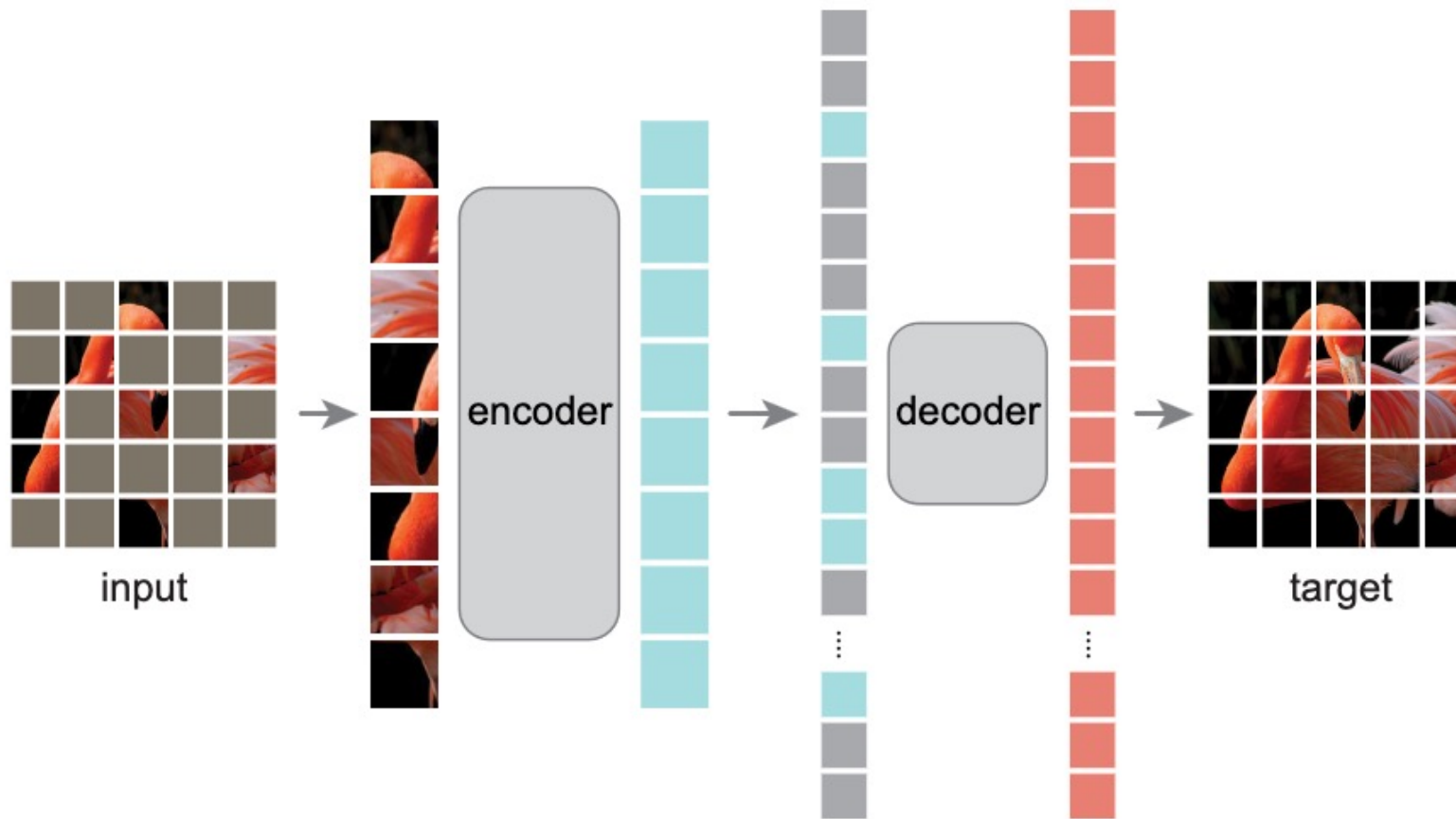
# Momentum Contrastive Learning (MoCo)



**Improved Baselines with Momentum Contrastive Learning**

<https://arxiv.org/abs/2003.04297> <https://arxiv.org/abs/1911.05722>

# Alternative: Masked AutoEncoders (MAE)

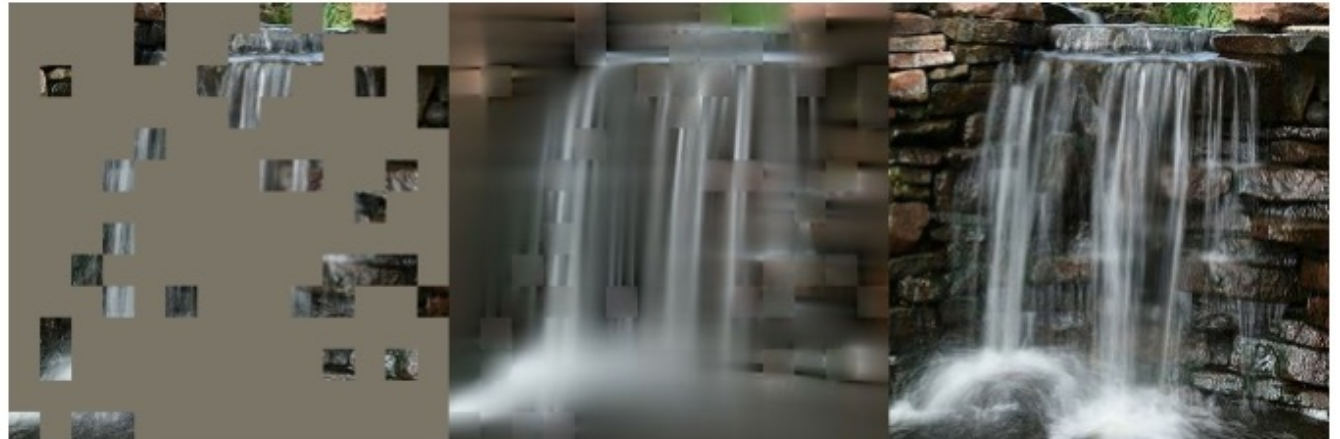
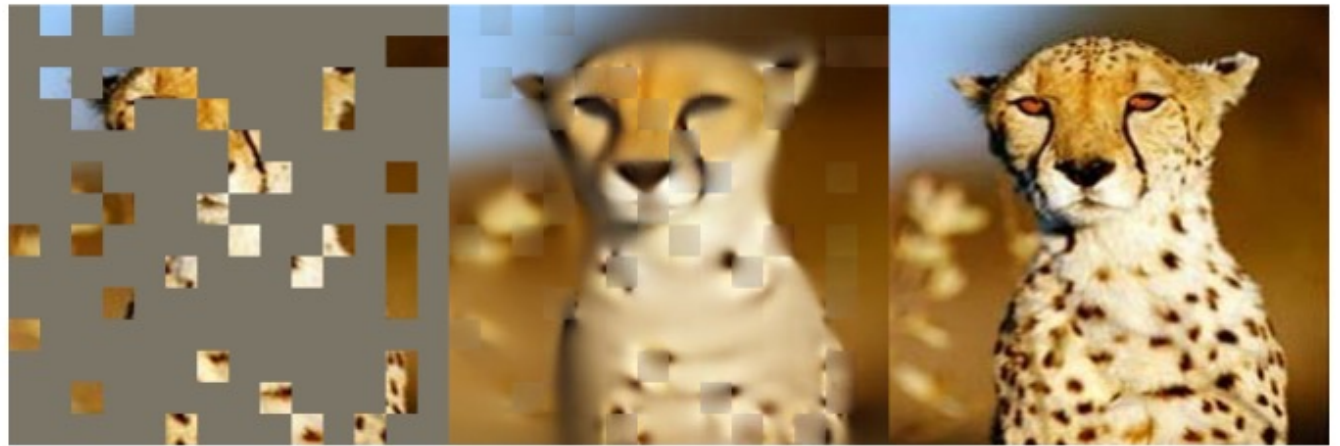


**Masked Autoencoders Are Scalable Vision Learners**

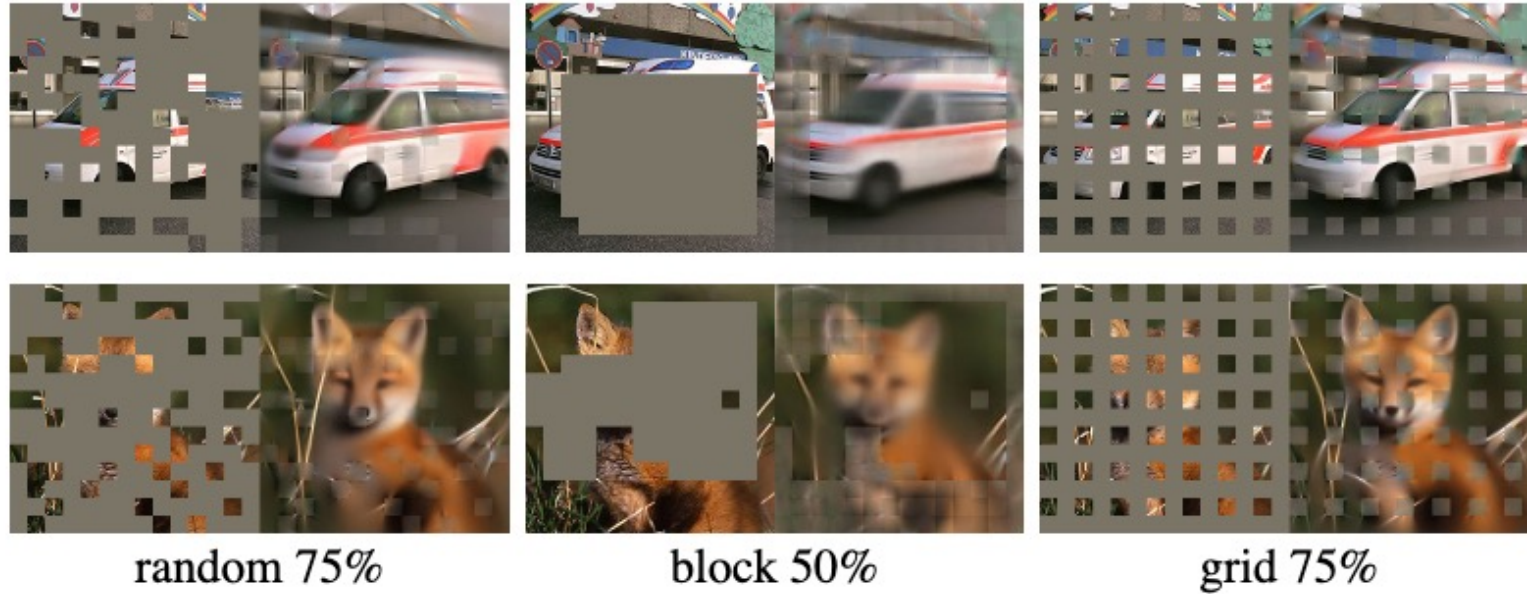
<https://arxiv.org/abs/2111.06377>

# Examples

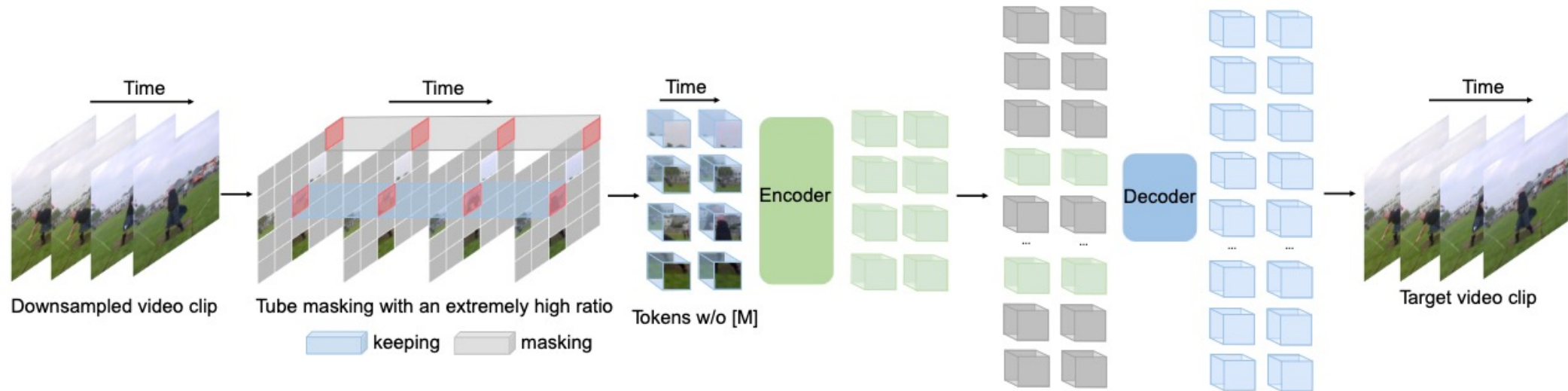
- After training, the model learns to fill-in-the-blanks for images.
- Similar to text masked image modeling
- The model can be finetuned for any other task.



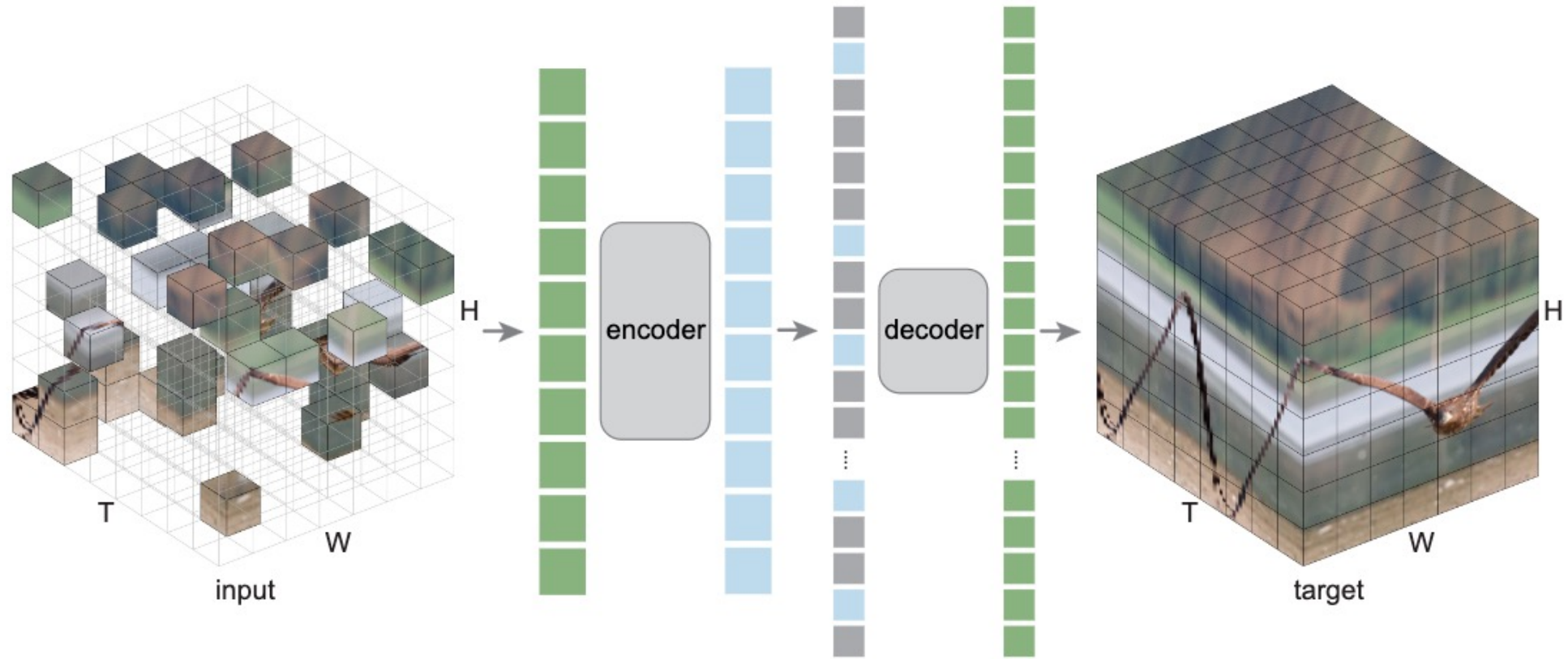
# Several Analysis in the Paper



# Video MAE

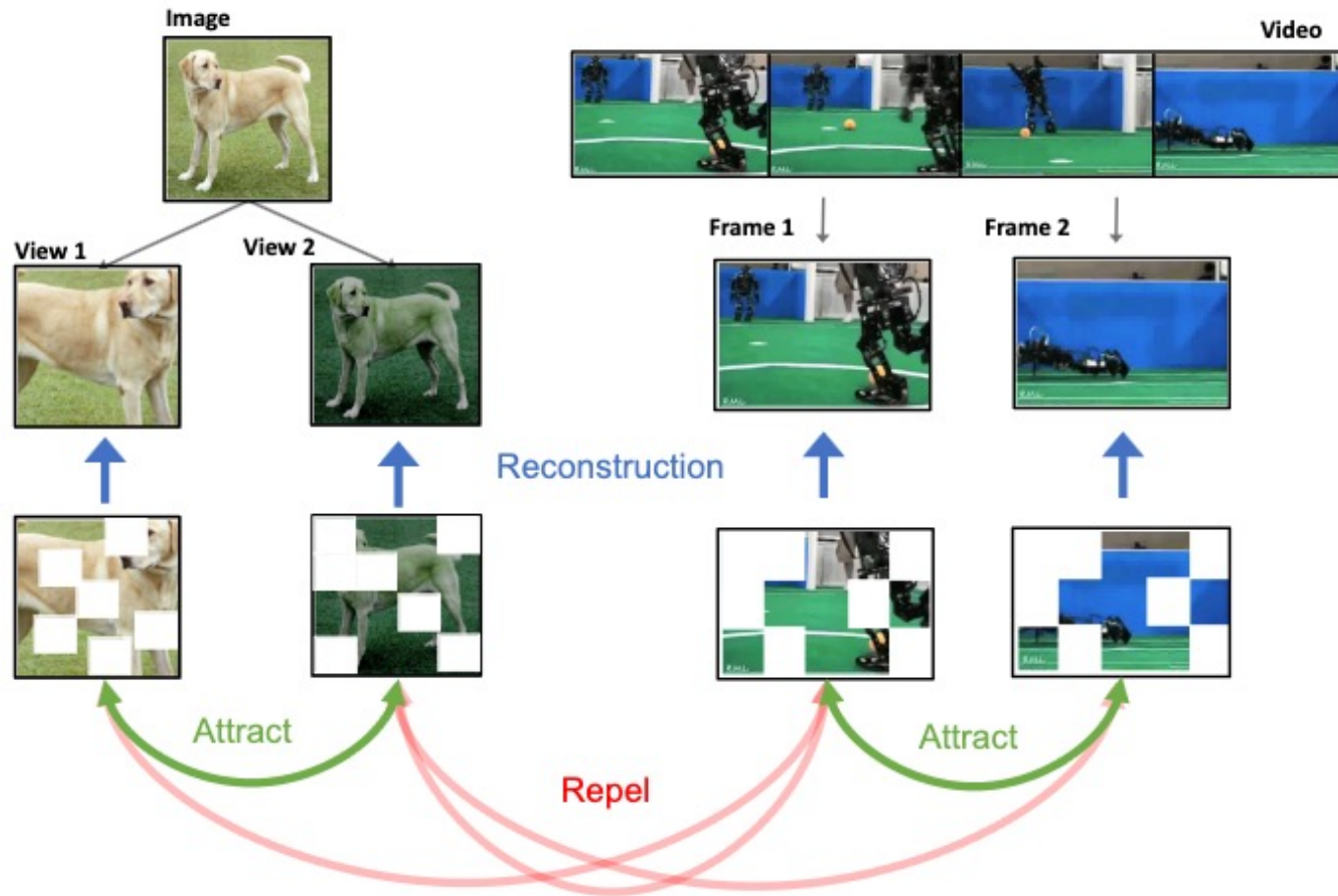


# ST MAE





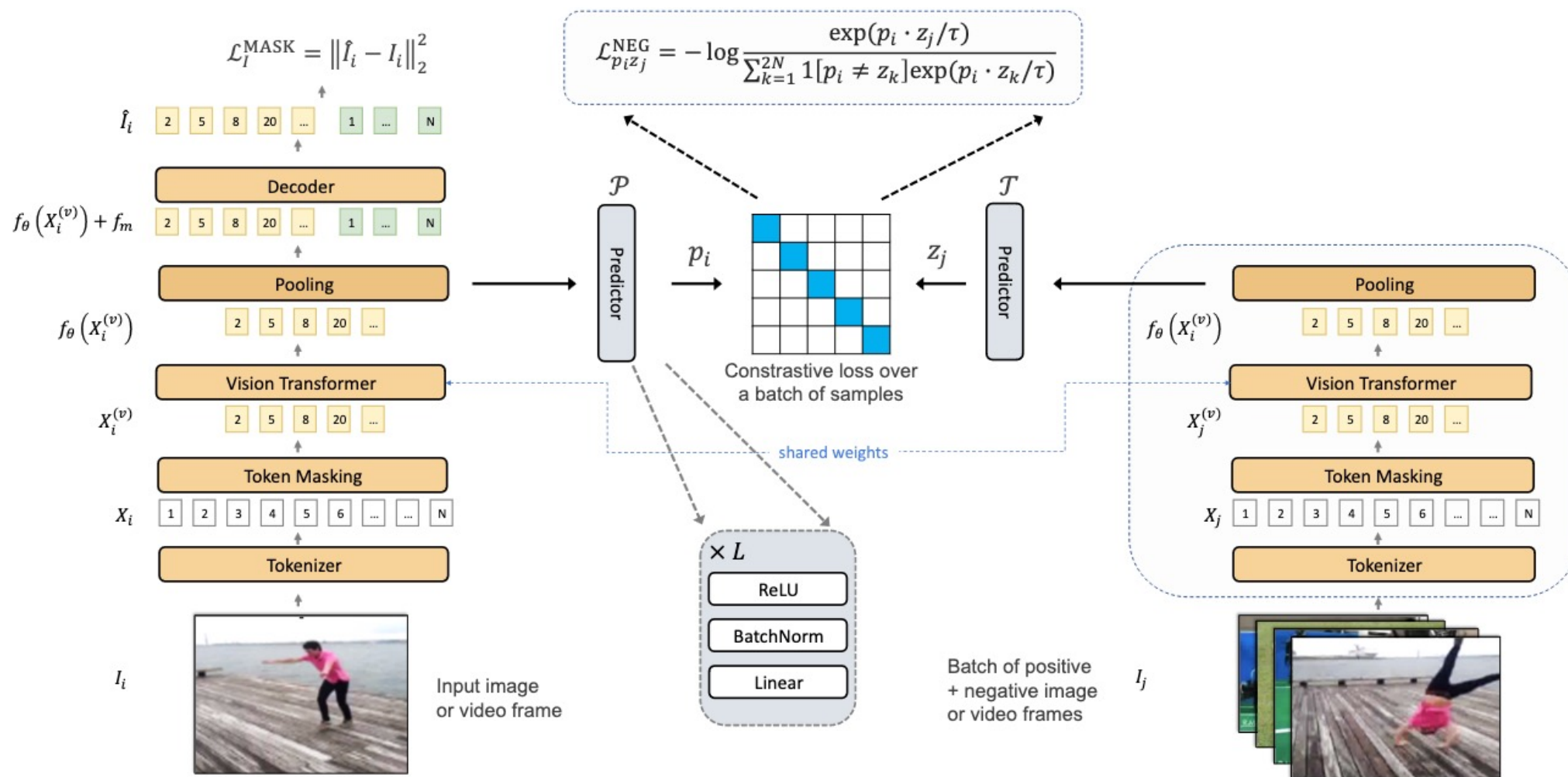
# ViC-MAE



ViC-MAE: Self-Supervised Representation Learning from Images and Video with Contrastive Masked Autoencoders

<https://arxiv.org/abs/2303.12001>

# ViC-MAE



# Results

	Method	Arch.	Pre-training Data	In-Domain		Out-of-Domain	
				IN1K	K400	Places-365	SSv2
Supervised	ViT [22] <i>ICML'20</i>	ViT-B	IN1K	82.3	68.5	57.0	61.8
	ViT [22] <i>ICML'20</i>	ViT-L	IN1K	82.6	78.6	58.9	66.2
	OMNIVORE [27] <i>CVPR'22</i>	ViT-B	IN1K + K400 + SUN RGB-D	84.0	83.3	<b>59.2</b>	68.3
	OMNIVORE [27] <i>CVPR'22</i>	ViT-L	IN1K + K400 + SUN RGB-D	<b>86.0</b>	84.1	–	–
	TubeViT [63] <i>CVPR'23</i>	ViT-B	K400 + IN1K	81.4	88.6	–	–
	TubeViT [63] <i>CVPR'23</i>	ViT-L	K400 + IN1K	–	<b>90.2</b>	–	<b>76.1</b>
Self-Supervised	MAE [35] <i>CVPR'22</i>	ViT-B	IN1K	83.4	–	57.9	59.6
	MAE [35] <i>CVPR'22</i>	ViT-L	IN1K	85.5	82.3	59.4	57.7
	ST-MAE [26] <i>NeurIPS'22</i>	ViT-B	K400	81.3	81.3	57.4	69.3
	ST-MAE [26] <i>NeurIPS'22</i>	ViT-L	K400	81.7	84.8	58.1	73.2
	VideoMAE [68] <i>NeurIPS'22</i>	ViT-B	K400	81.1	80.0	–	69.6
	VideoMAE [68] <i>NeurIPS'22</i>	ViT-L	K400	–	85.2	–	74.3
	OmniMAE [29] <i>CVPR'23</i>	ViT-B	K400 + IN1K	82.8	80.8	58.5	69.0
	OmniMAE [29] <i>CVPR'23</i>	ViT-L	K400 + IN1K	84.7	84.0	59.4	73.4
	ViC-MAE	ViT-L	K400	85.0	85.1	59.5	73.7
	ViC-MAE	ViT-L	MiT	85.3	84.9	59.7	73.8
	ViC-MAE	ViT-B	K400 + IN1K	83.0	80.8	58.6	69.5
	ViC-MAE	ViT-L	K400 + IN1K	86.0	86.8	60.0	75.0
	ViC-MAE	ViT-B	K400 + K600 + K700 + MiT + IN1K	83.8	80.9	59.1	69.8
	ViC-MAE	ViT-L	K400 + K600 + K700 + MiT + IN1K	<b>87.1</b>	<b>87.8</b>	<b>60.7</b>	<b>75.9</b>

# Questions