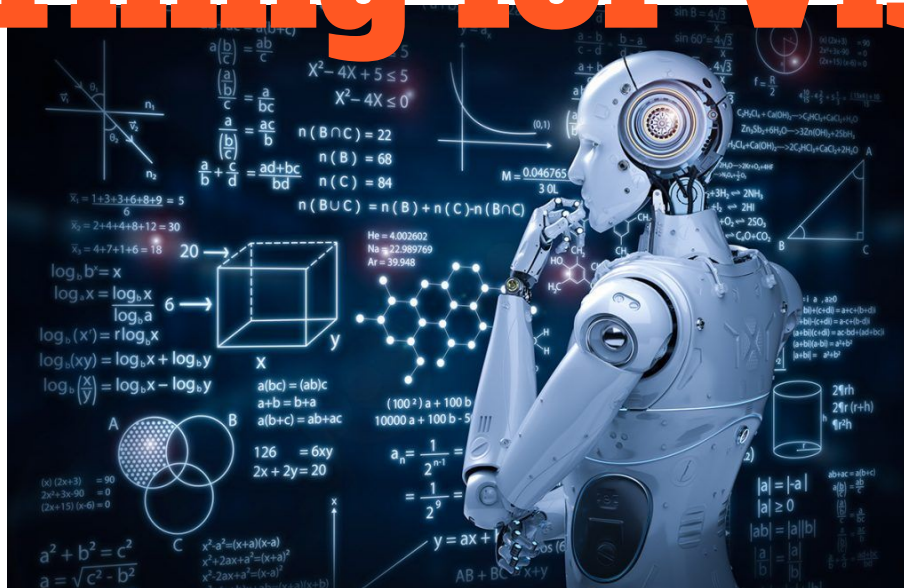


Self-Supervised Learning for Vision



Javier Selva Castelló



Universitat
de Barcelona



Motivation

- Supervision is **costly**!



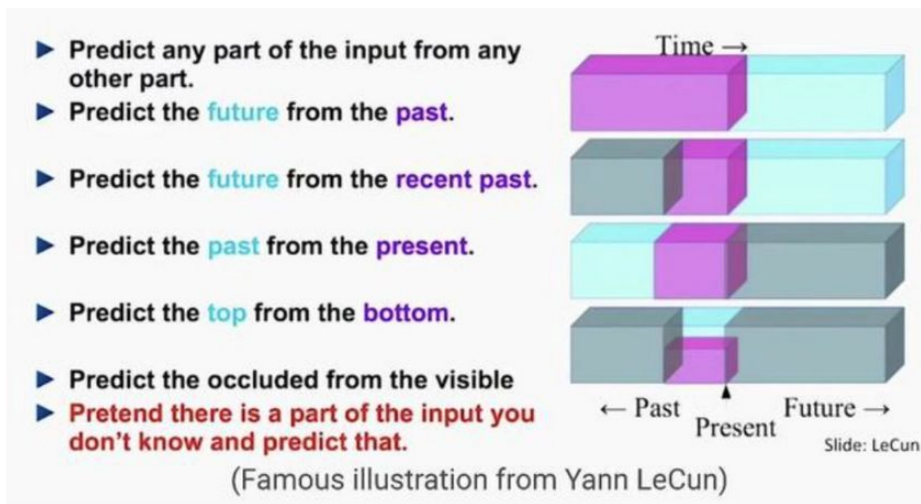
Motivation

- Supervision is **costly**!
- Billions of GB of **internet content**.



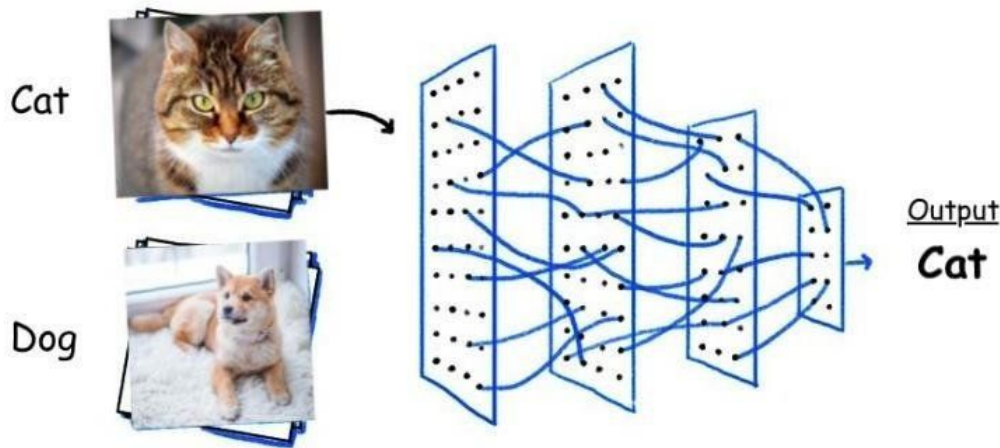
Motivation

- Supervision is **costly**!
- Billions of GB of **internet content**.
- Can we use the **data** itself as **supervision**?

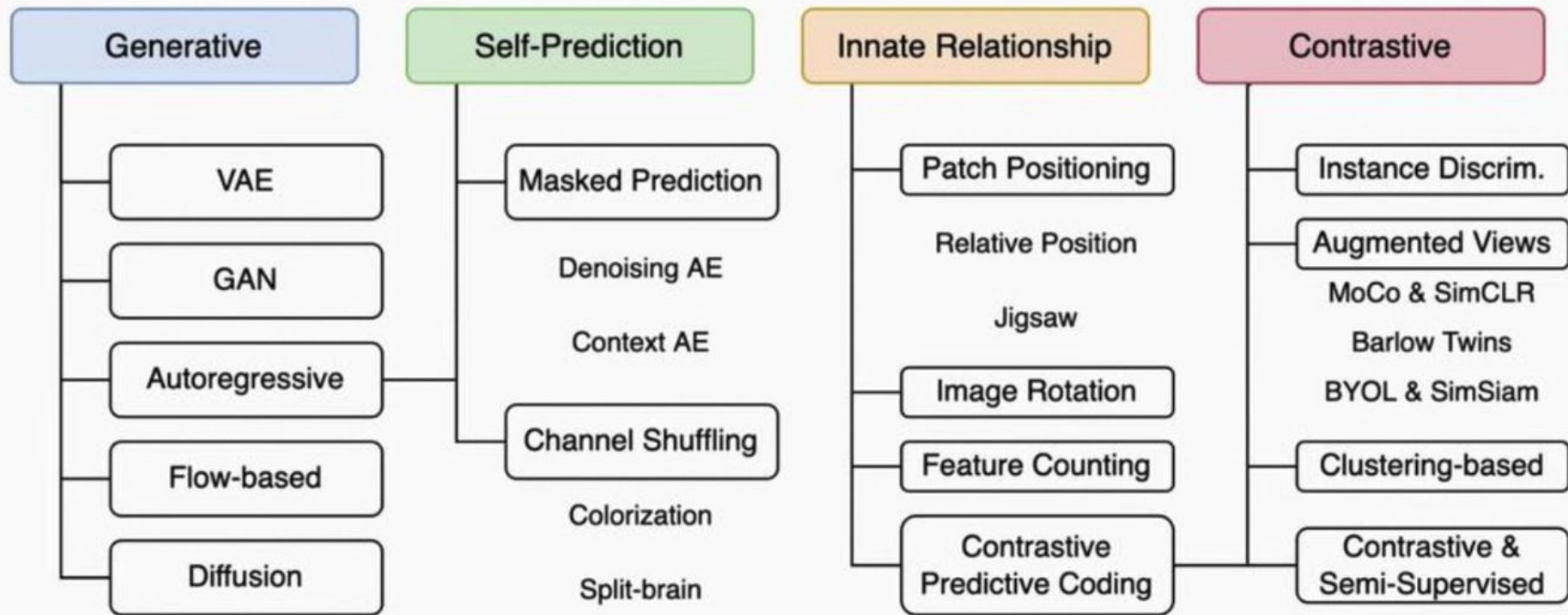


Motivation

- Supervision is **costly**!
- Billions of GB of **internet content**.
- Can we use the **data** itself **as supervision**?
- Fine-tune pre-trained model on **supervised downstream task**.



Many methods!!



**Some
background**

Basics: Contrastive Learning

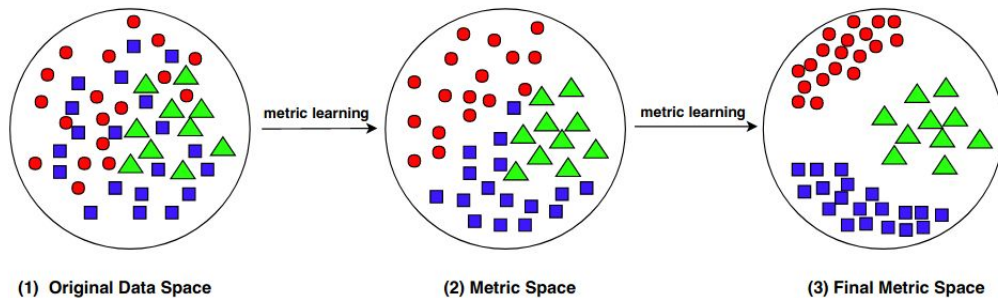
Similar samples should be **closer** in representation space.

COLLAPSE TO CONSTANT REPRESENTATION!

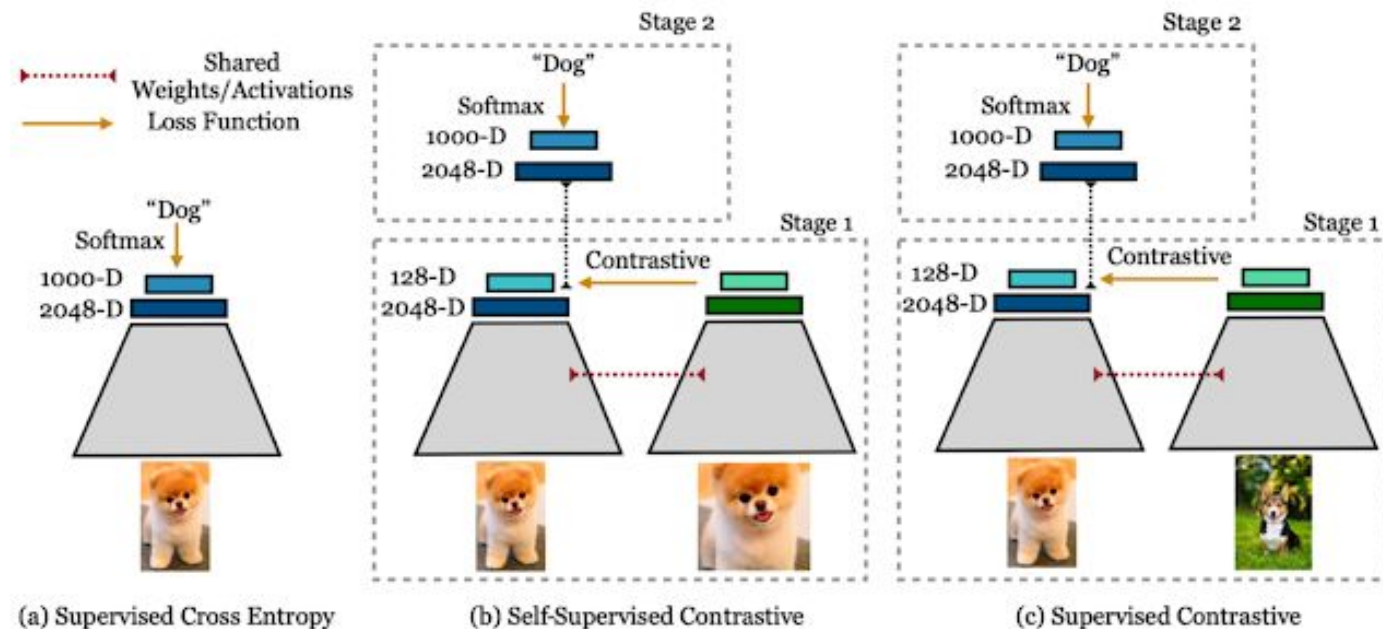
Different samples (negatives) should be **far away** in representation space.

Problems:

- **Size** of negative set
- **Quality** of negatives



Basics: Contrastive Learning

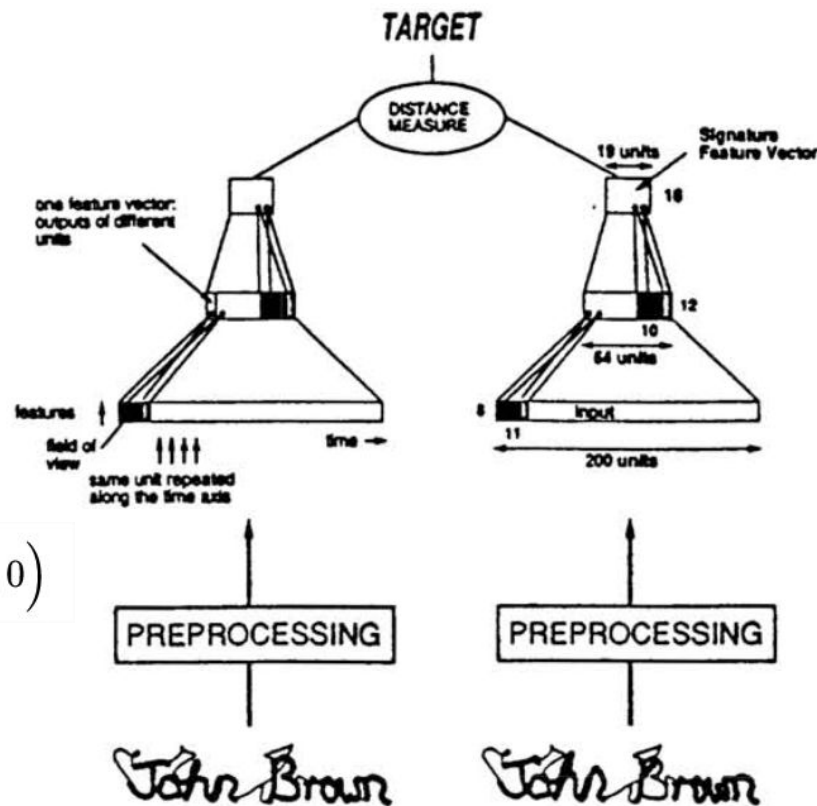


Basics: Siamese Networks

- Two sister networks (same weights)
- Different inputs
- Trained contrastively.

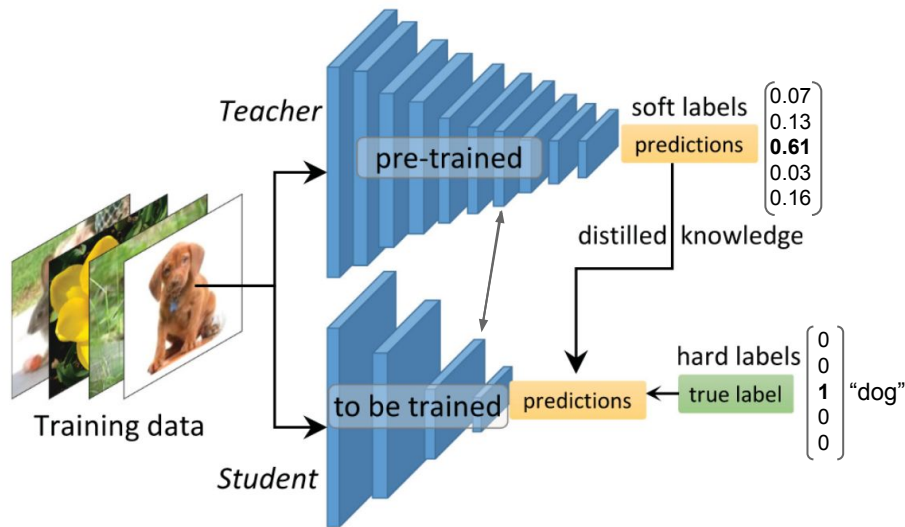
Triplet loss:

$$\mathcal{L}(A, P, N) = \max\left(\|f(A) - f(P)\|^2 - \|f(A) - f(N)\|^2 + \alpha, 0\right)$$



Basics: Knowledge Distillation

- Mimic behaviour of a larger network by **transferring learned representations** to a smaller one.
- Train large network.
- Train smaller network on output logits from the first network.

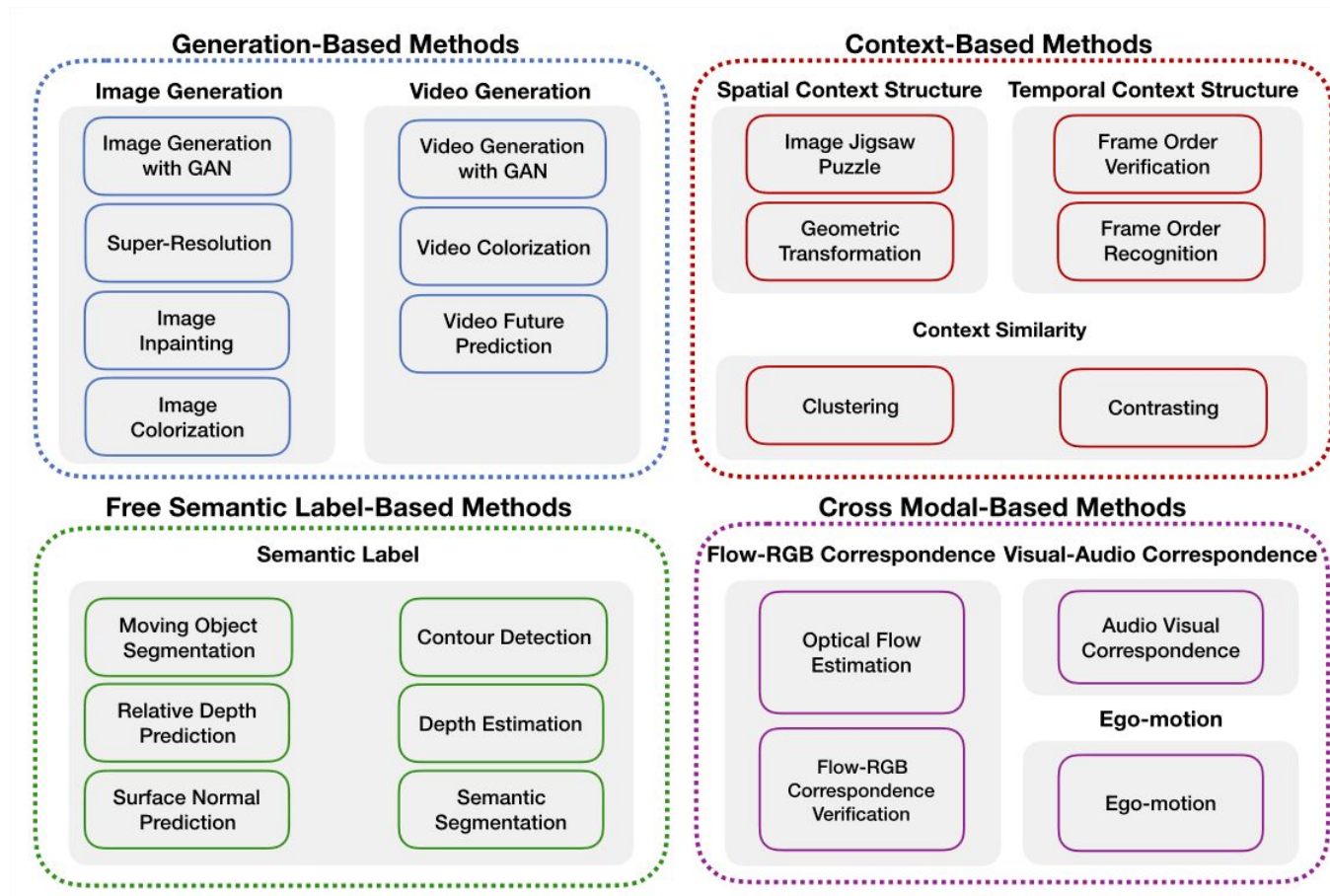


Classic **Self-supervised Learning for Video**

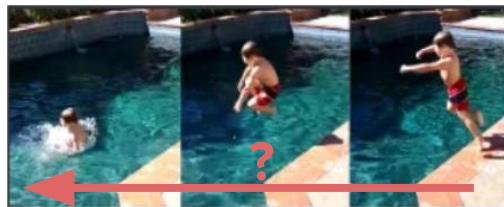
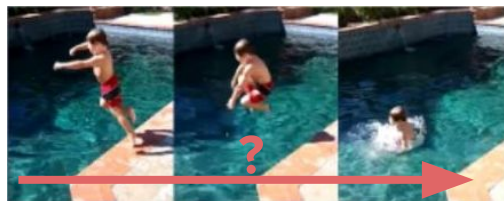


Disclaimer! Multi-modality not considered

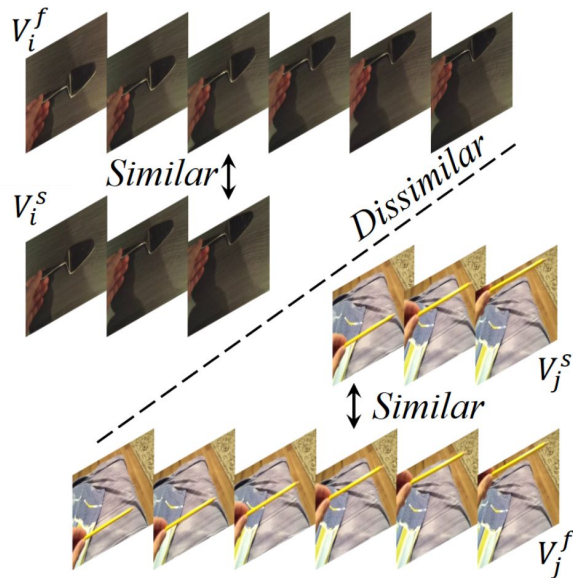
Classic SSL for Video



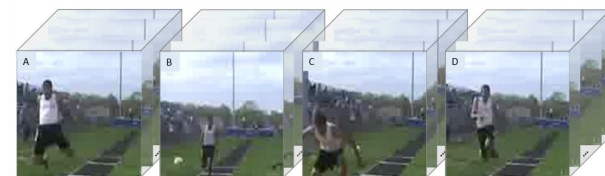
Temporal Consistency



Arrow of time

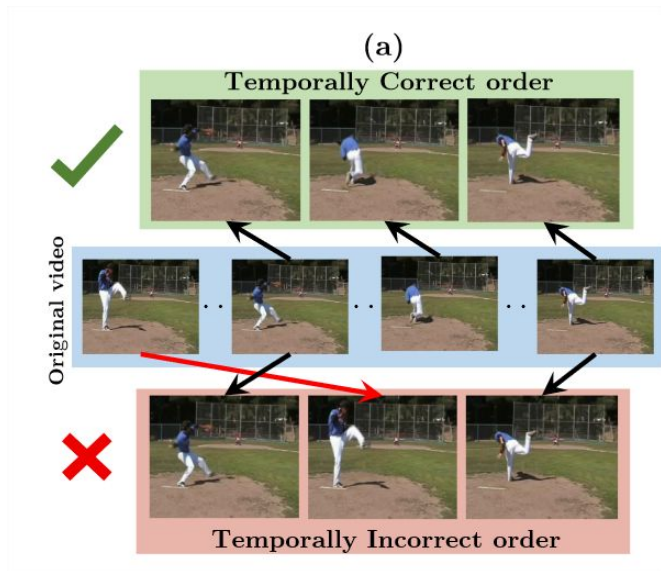


Tempo consistency

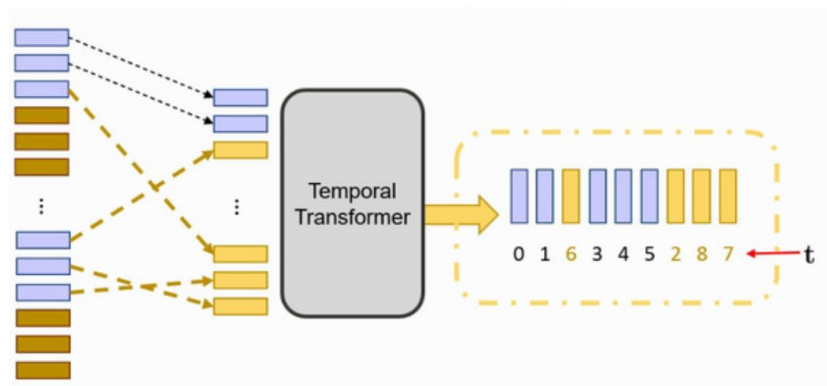


Cubic puzzles

Jigsaw Puzzle

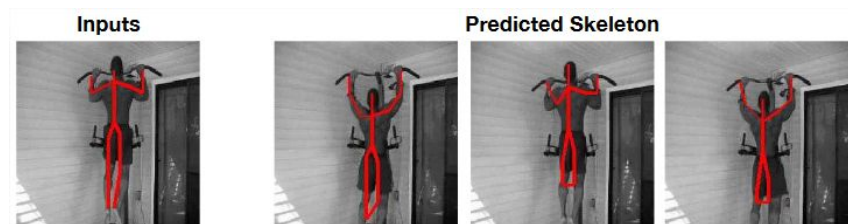
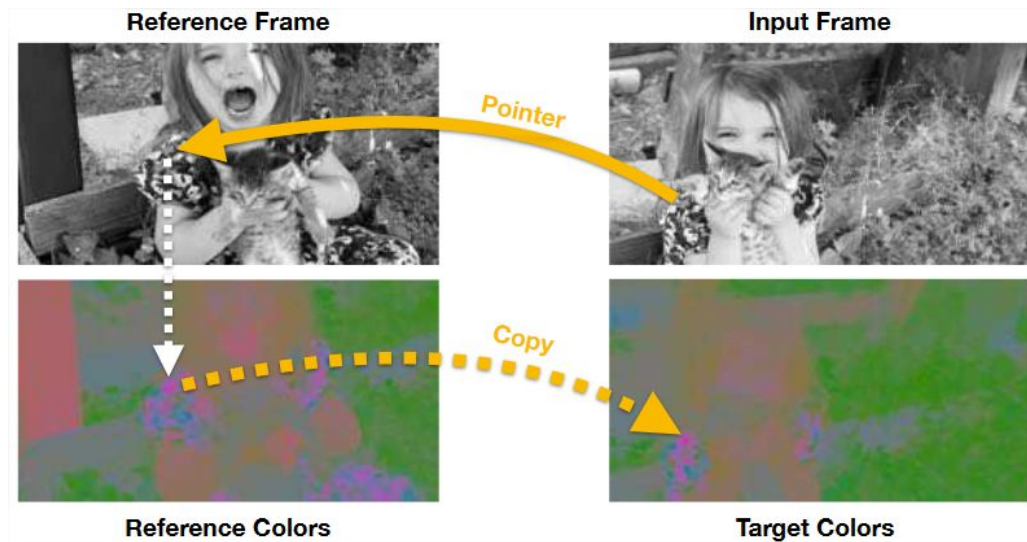


Binary



Multi-class

Colorization

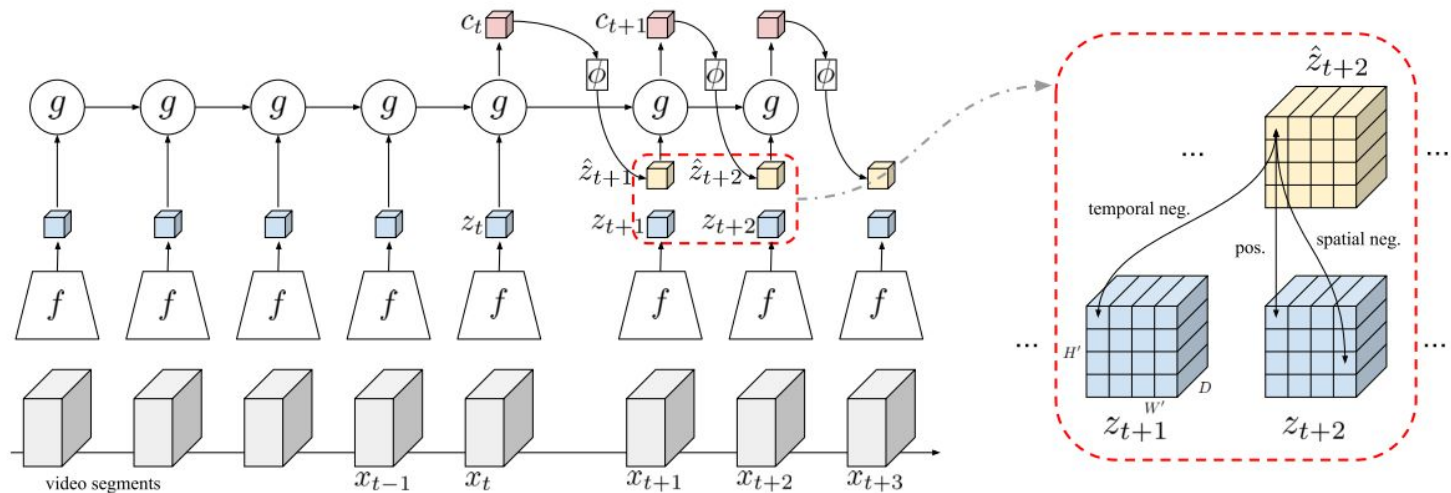


Video Prediction

Given a video sequence, generate the next frames.



Dense Predictive Coding



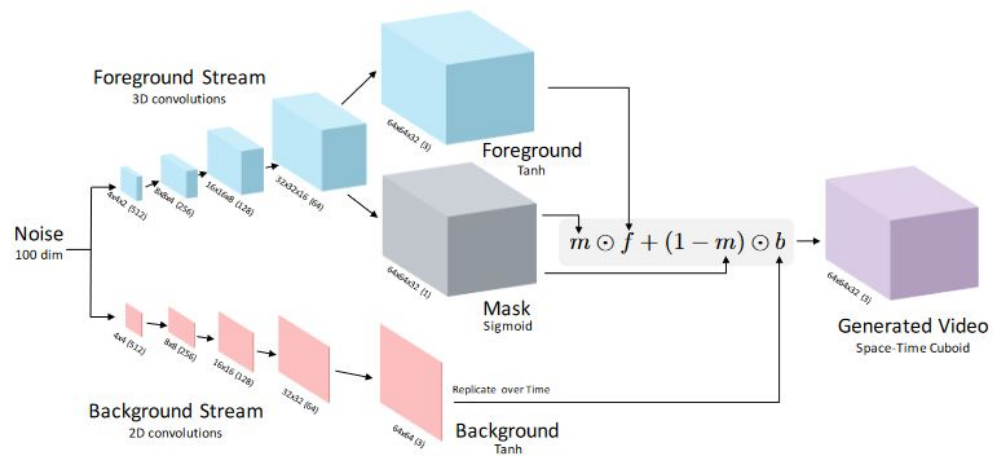
Predictive Coding: anticipate the future clip representations given a context of multiple clips.

Sample negatives from:

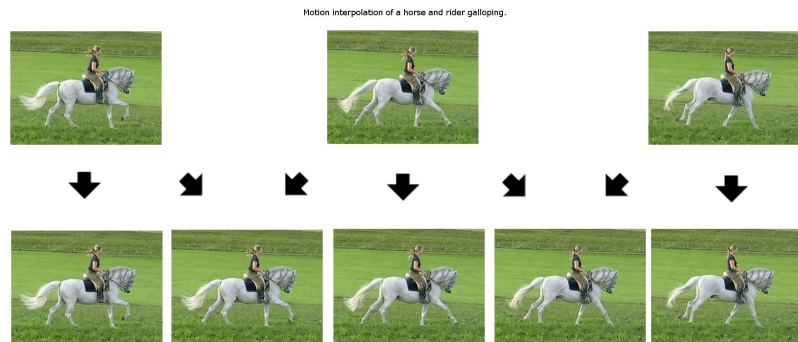
- other videos
- same video on different time
- same time on different positions.

$$\mathcal{L} = - \sum_{i,k} \left[\log \frac{\exp(\hat{z}_{i,k}^\top \cdot z_{i,k})}{\sum_{j,m} \exp(\hat{z}_{i,k}^\top \cdot z_{j,m})} \right]$$

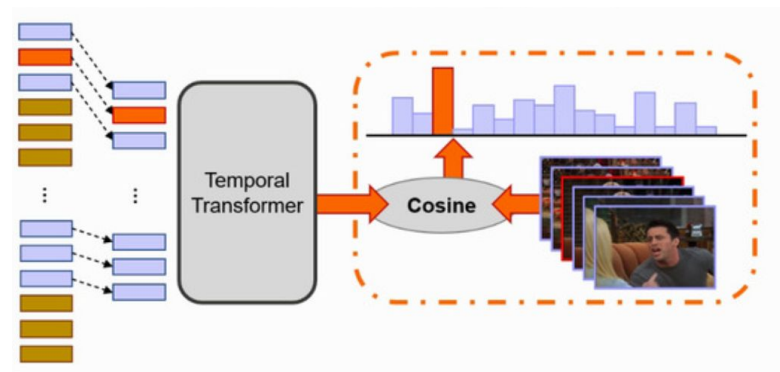
Generative Video Models



Interpolation



Generative



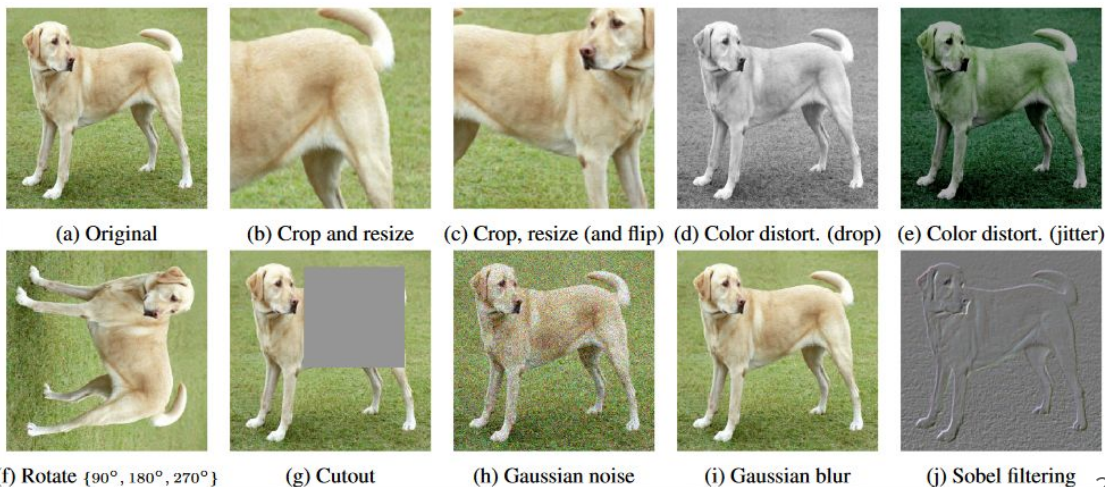
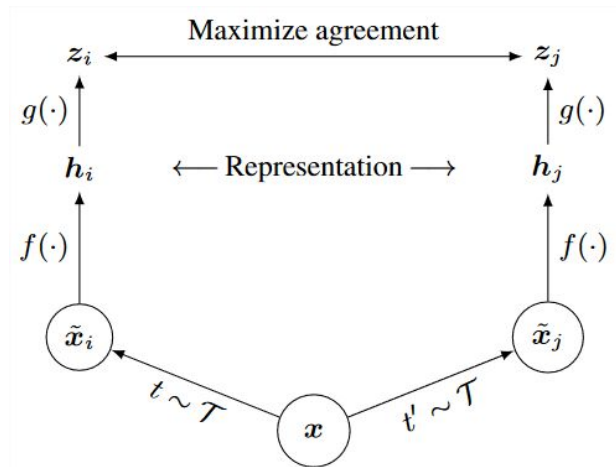
Contrastive

Current Trends on Visual SSL

Contrastive: SimCLR

- Siamese setting.
- Data augmentation.
- Contrastive loss.
- Negatives from batch.

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

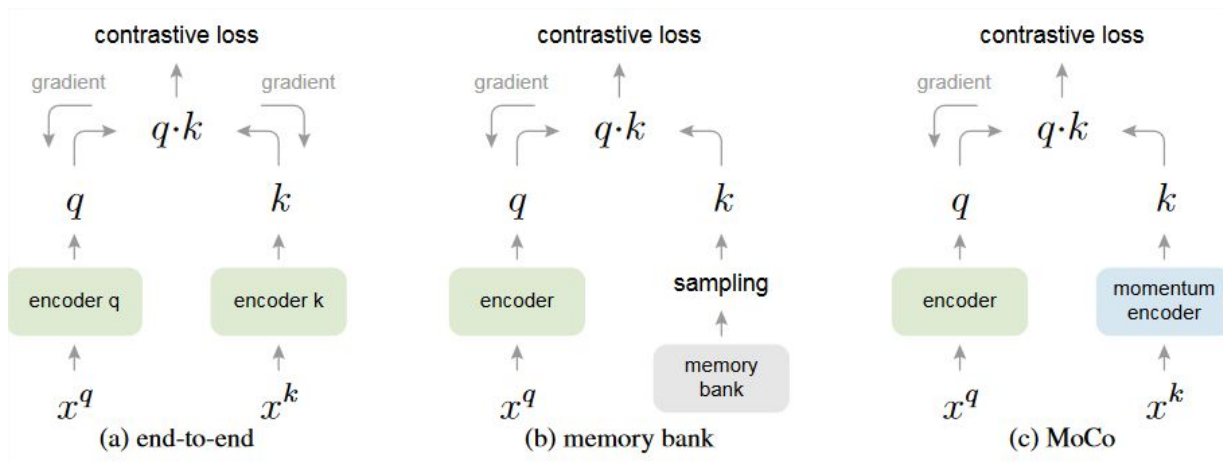


Contrastive: MoCo

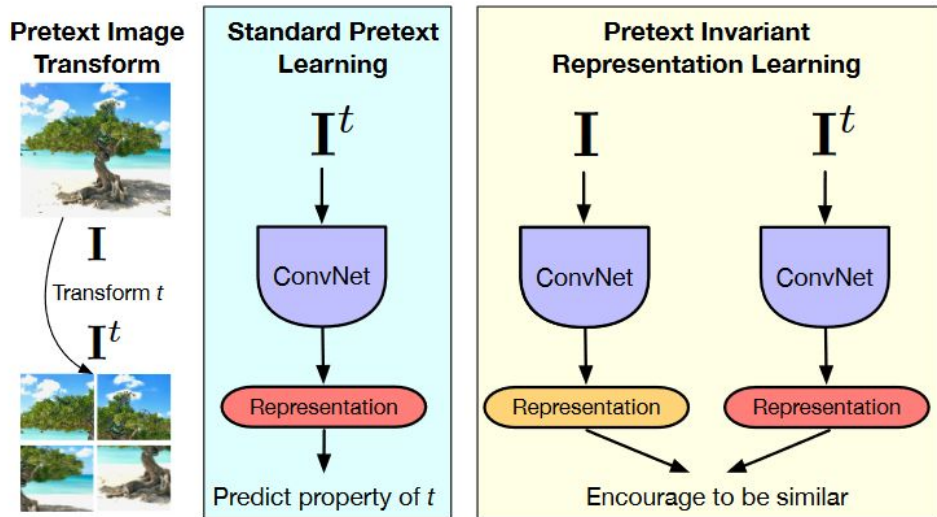
- Similar to SimCLR
- Uses queue to store past batches
- Draw negatives from queue

- Unstable training!
- Use momentum encoder

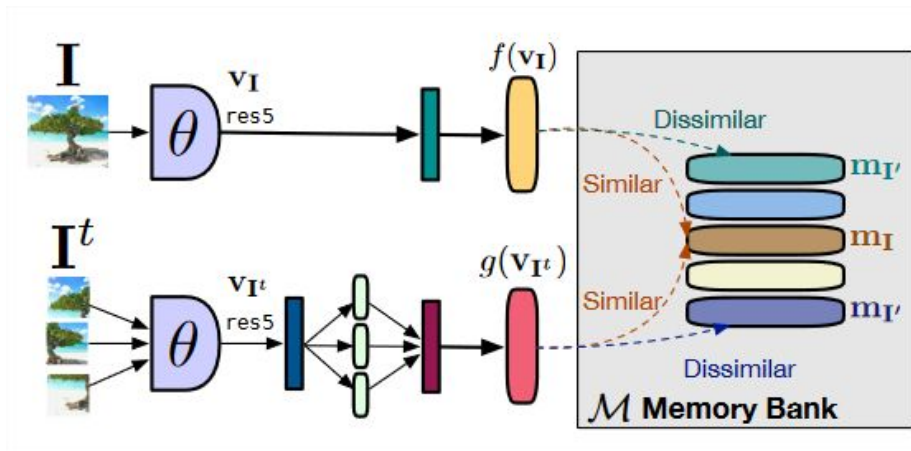
$$\theta_k \leftarrow m\theta_k + (1 - m)\theta_q$$



Contrastive: PIRL



Pretext tasks → **covariant** representations

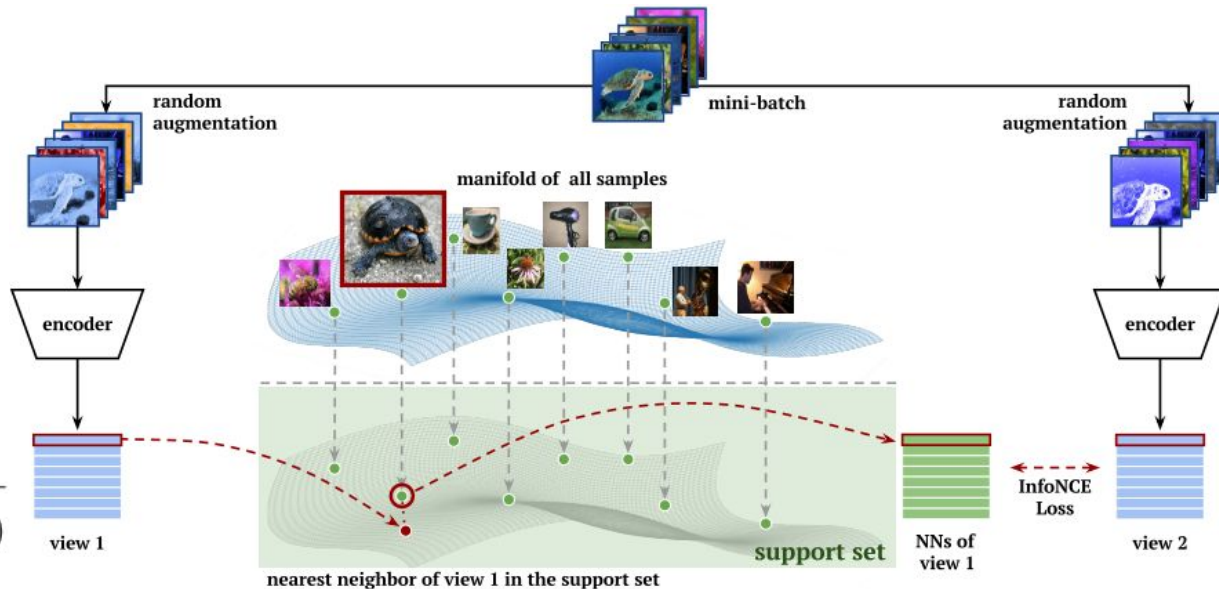


Use **siamese** and **contrastive** to make them **invariant**.

Clustering: NNCLR

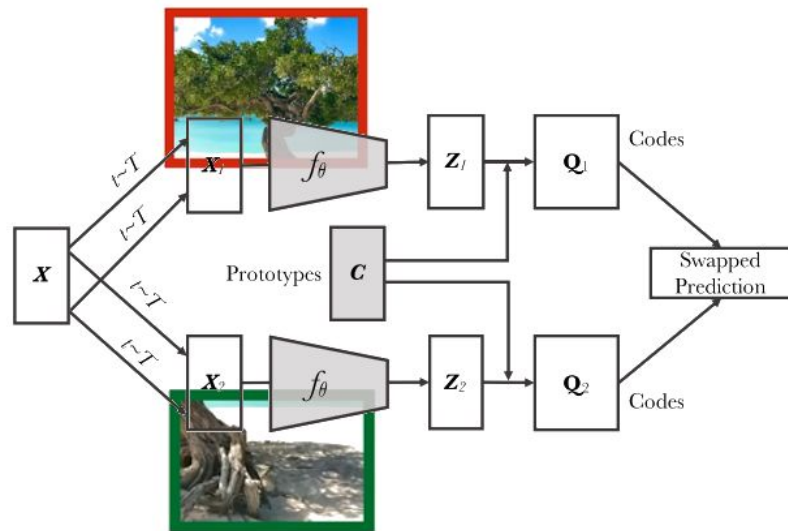
- **Contrastive** setting
- Negatives from mini-batch
- Positive is a **similar sample** from a **memory bank**

$$\mathcal{L}_i^{\text{NNCLR}} = -\log \frac{\exp(\text{NN}(z_i, Q) \cdot z_i^+ / \tau)}{\sum_{k=1}^n \exp(\text{NN}(z_i, Q) \cdot z_k^+ / \tau)}$$

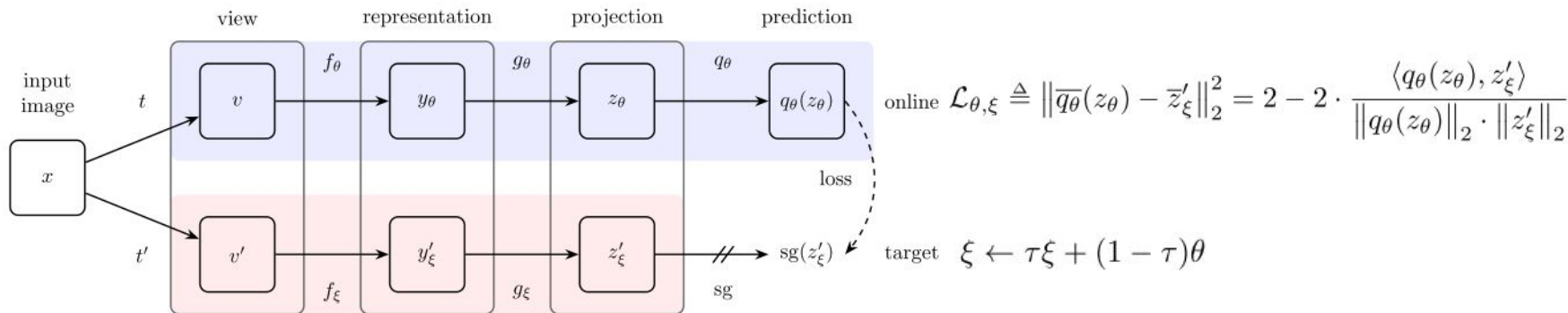


Clustering: SwAv

- NO MORE NEGATIVES!
- Learnable prototypes
- Collapse!
 - Cluster using **Sinkhorn-Knopp**
 - **Uniform sample distribution**
 - **Soft assignment**
- Swapped prediction.
- **Multi-crop**



Distillation: BYOL



- Knowledge **distillation** setting
- **Asymmetric** Siamese architecture
- Both networks **start from scratch**
- Teacher is a **moving average**
- Minimize distance between augmentations

WHY DOES THIS
NOT COLLAPSE??



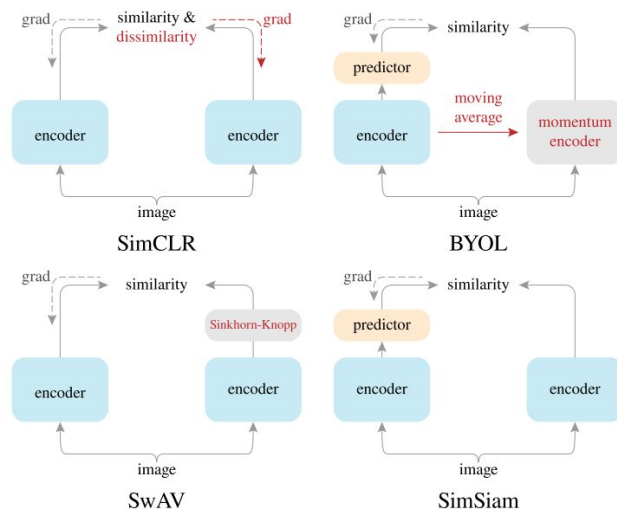
Distillation: SimSiam

Which components are **crucial for avoiding collapse**?

Stop-gradient

Momentum encoder and **predictor** improve accuracy.

- **Expectation Maximization** algorithm:
- E: the teacher generates various samples (n^t)
 - M: fitting the parameters of the online net to approximate those representations.



$$\mathcal{L}(\theta, \eta) = \mathbb{E}_{x, \mathcal{T}} \left[\left\| \mathcal{F}_{\theta}(\mathcal{T}(x)) - \eta_x \right\|_2^2 \right]$$

$$\theta^t \leftarrow \arg \min_{\theta} \mathcal{L}(\theta, \eta^{t-1})$$

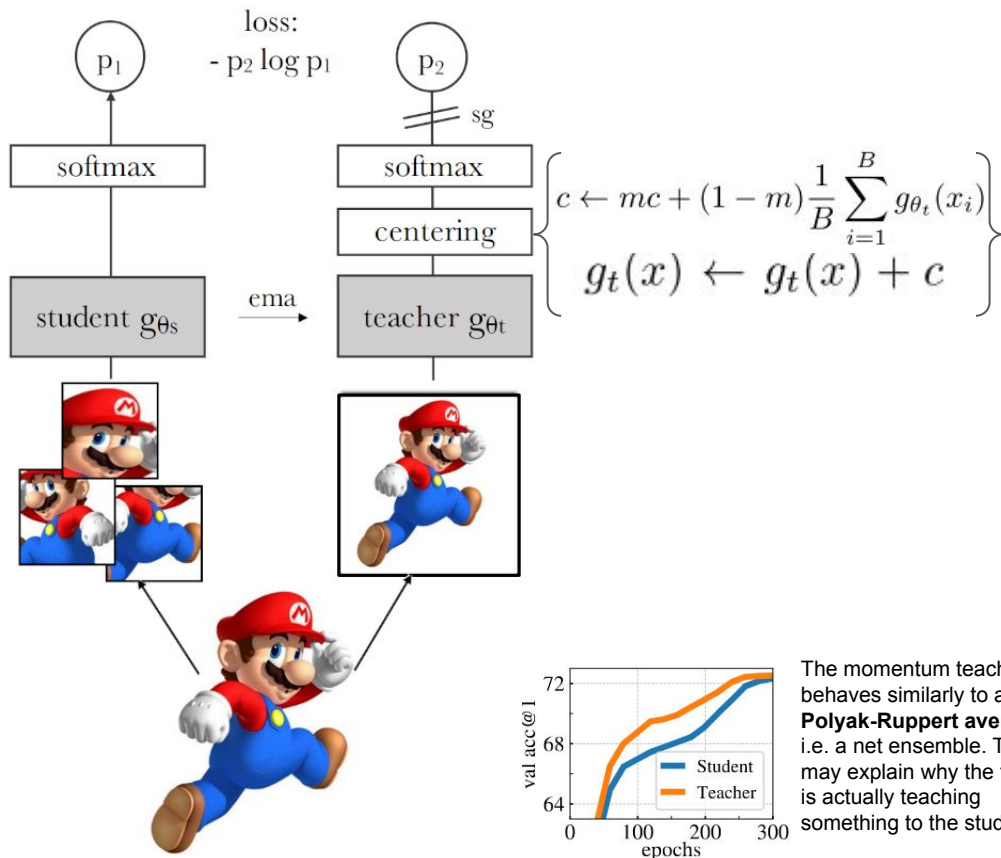
$$\eta^t \leftarrow \arg \min_{\eta} \mathcal{L}(\theta^t, \eta)$$

$$\eta_x^t \leftarrow \mathcal{F}_{\theta^t}(\mathcal{T}'(x)).$$

$$\theta^{t+1} \leftarrow \arg \min_{\theta} \mathbb{E}_{x, \mathcal{T}} \left[\left\| \mathcal{F}_{\theta}(\mathcal{T}(x)) - \mathcal{F}_{\theta^t}(\mathcal{T}'(x)) \right\|_2^2 \right]$$

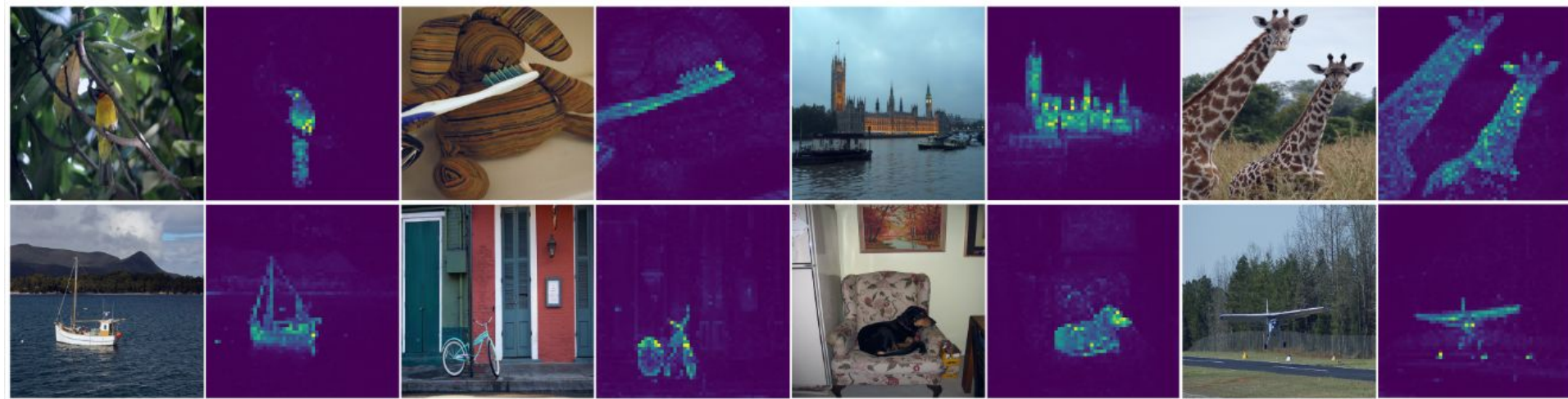
Distillation: DINO

- Instance classification (CE)
- Multi Crop
- Transformer Architecture
- Centering & Sharpening



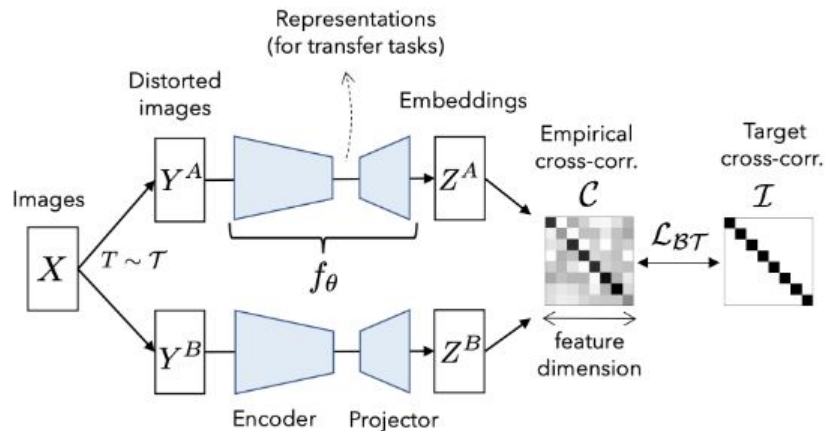
The momentum teacher behaves similarly to a **Polyak-Ruppert averaging**, i.e. a net ensemble. This may explain why the teacher is actually teaching something to the student.

Distillation: DINO



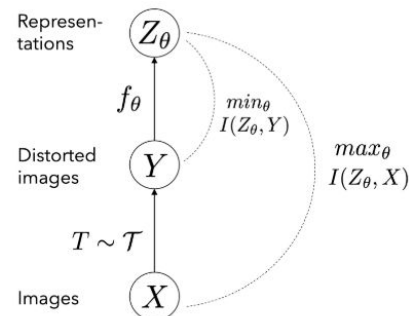
Redundancy Reduction: Barlow Twins

- Symmetric Siamese
- Correlation matrix
 - **Invariance** (same feature should be similar in both embeddings)
 - **Reduces redundancy** (no two features should be similar)



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

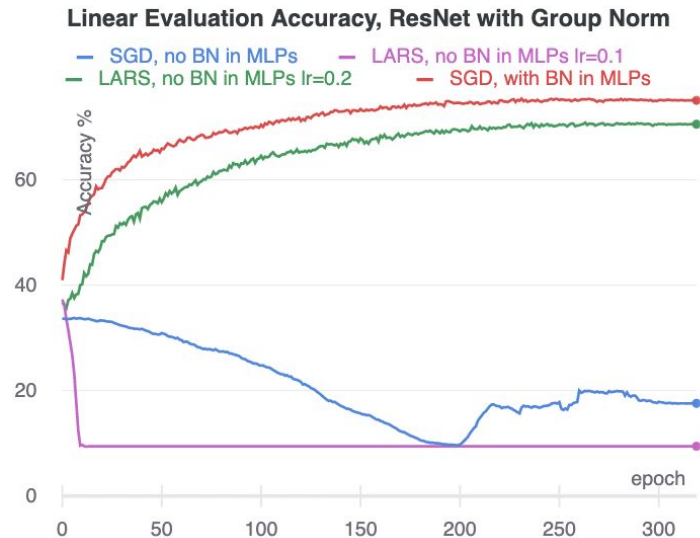
Information Bottleneck
Principle: Maximize input-output mutual information while being invariant to distortions



Information Bottleneck Principle

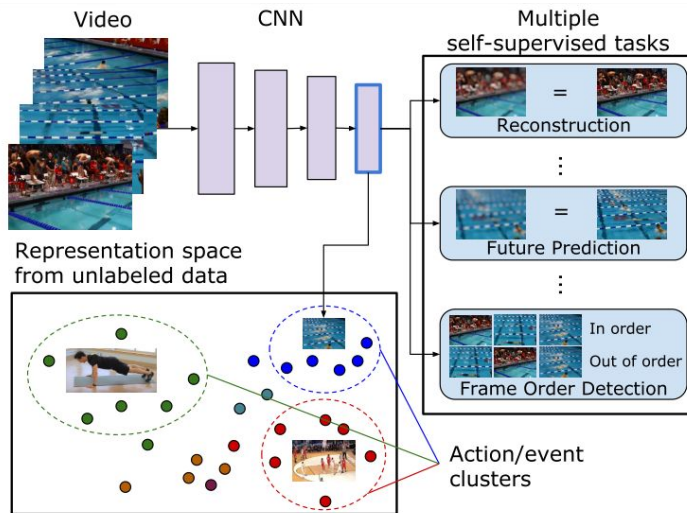
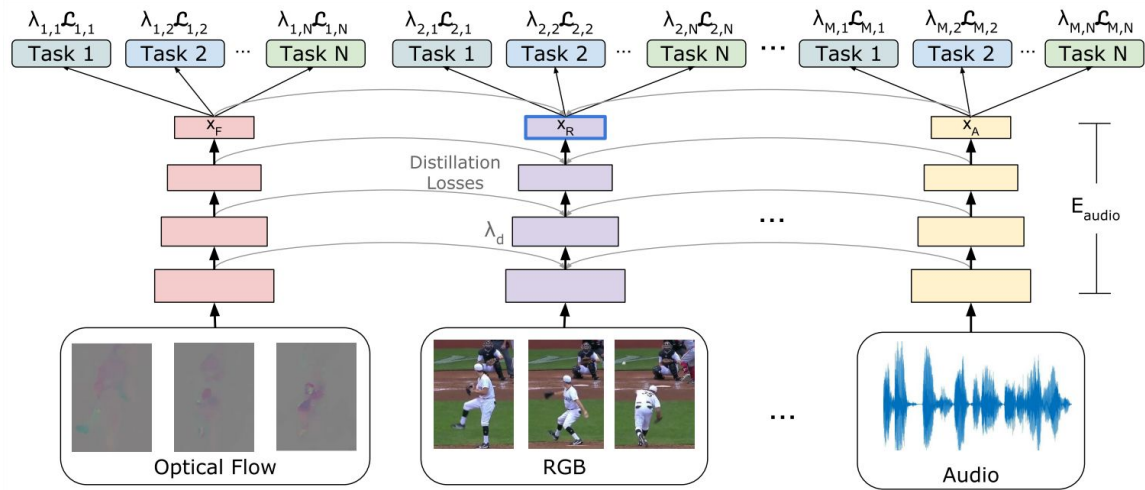
Intuitions

- Batch Normalization
- LARS optimizer
- Asymmetry (projector/predictor with large LR)
- Momentum Network
- Weight Decay
- Stop-gradient



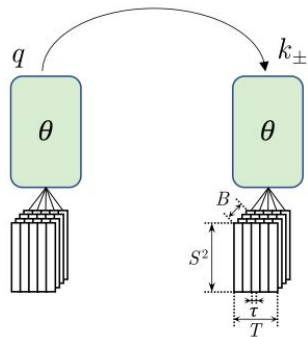
**What about
video?**

Evolved SSL Losses

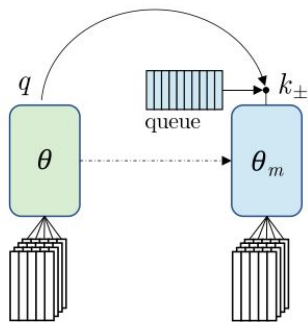


- **Multiple streams** (flow, grayscale, audio...)
- **Distilled** at different depths onto the RGB stream.
- **Multi-task** (reconstruction, prediction, temporal ordering, multi-modal contrastive/alignment).
- **Evolutionary algorithm** to weight the losses.
- **Clustered** output forced to follow **Zipf's law** through **KL divergence**

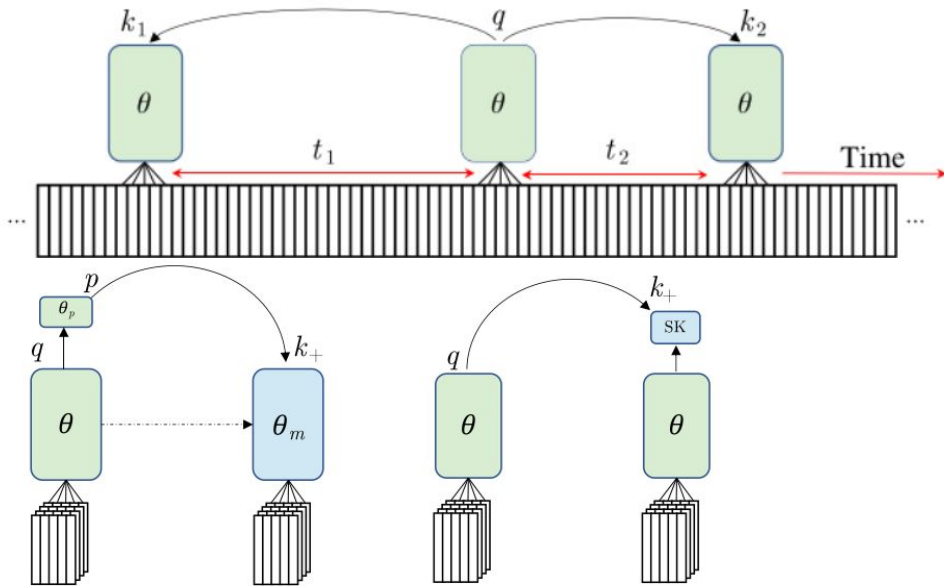
Large-scale study



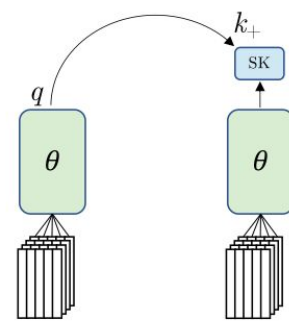
(a) SimCLR



(b) MoCo



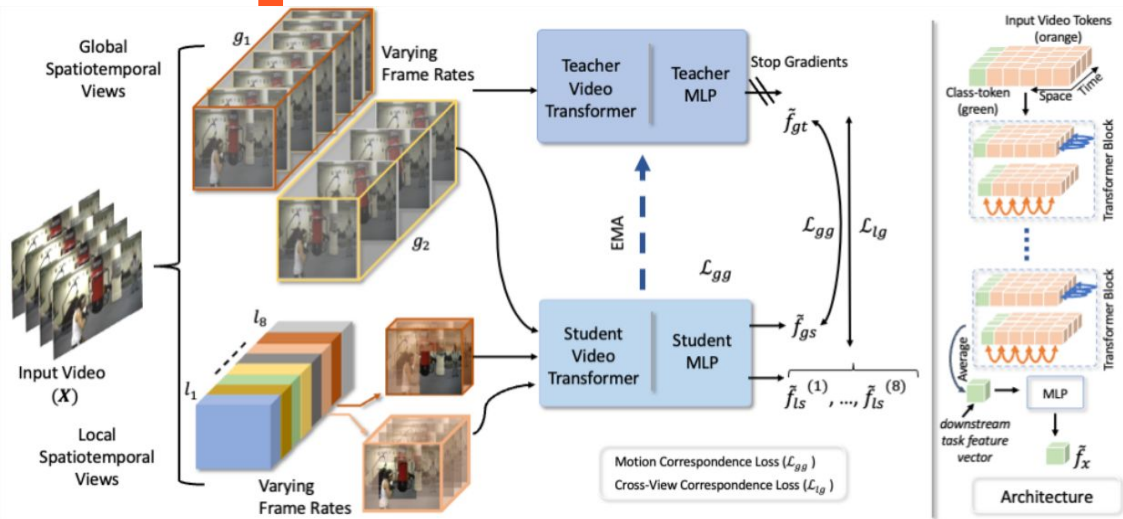
(c) BYOL



(d) SwAV

method	pre-train	linear protocol	finetuning accuracy			
		K400	UCF101	AVA (mAP)	Charades (mAP)	SSv2
supervised	scratch	74.7	68.8	11.7	7.4	48.8
supervised	K400-240K	-	94.8	22.2	34.7	52.8
SimCLR	K400-240K	62.0 (-12.7)	87.9 (-6.9)	17.6 (-4.6)	11.4 (-23.3)	52.0 (-0.8)
SwAV		62.7 (-11.5)	89.4 (-5.4)	18.2 (-4.0)	10.7 (-24.0)	51.7 (-1.1)
BYOL		68.3 (-6.4)	93.8 (-1.0)	23.4 (+1.2)	21.0 (-13.7)	55.8 (+3.0)
MoCo		67.3 (-7.4)	92.8 (-2.0)	20.3 (-1.9)	33.5 (-1.2)	54.4 (+1.8)

Self-supervised Video Transformer

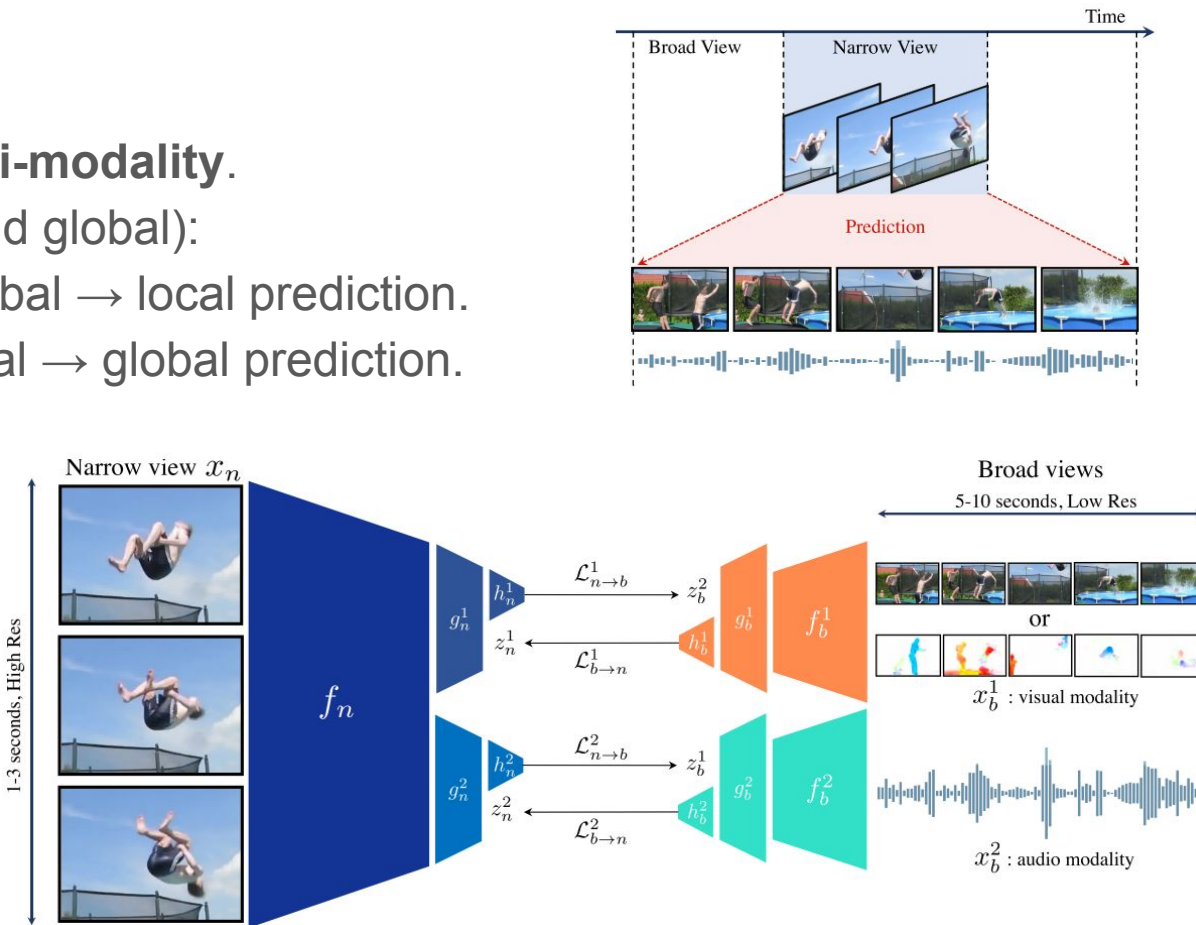


Extension to **DINO** for video. Use of multiple **global views** (varying frame-rate) and multiple **local views**. Matching at global-global (motion) and local-global (cross-view).

$l \rightarrow g$	$g \rightarrow g$	$l \rightarrow l$	$g \rightarrow l$	UCF-101	HMDB-51
✓	✗	✗	✗	84.11	50.72
✗	✓	✗	✗	81.95	49.04
✓	✓	✗	✗	84.64	52.17
✓	✓	✓	✗	83.11	51.23
✓	✓	✗	✓	84.71	51.88
✓	✓	✓	✓	83.69	51.71

Brave

- Focus on leveraging **multi-modality**.
- **Temporal crops** (local and global):
 - Teacher performs global \rightarrow local prediction.
 - Student performs local \rightarrow global prediction.
- Pairwise predictors.
- Synchronization matters!



Thanks!

QUESTIONS?