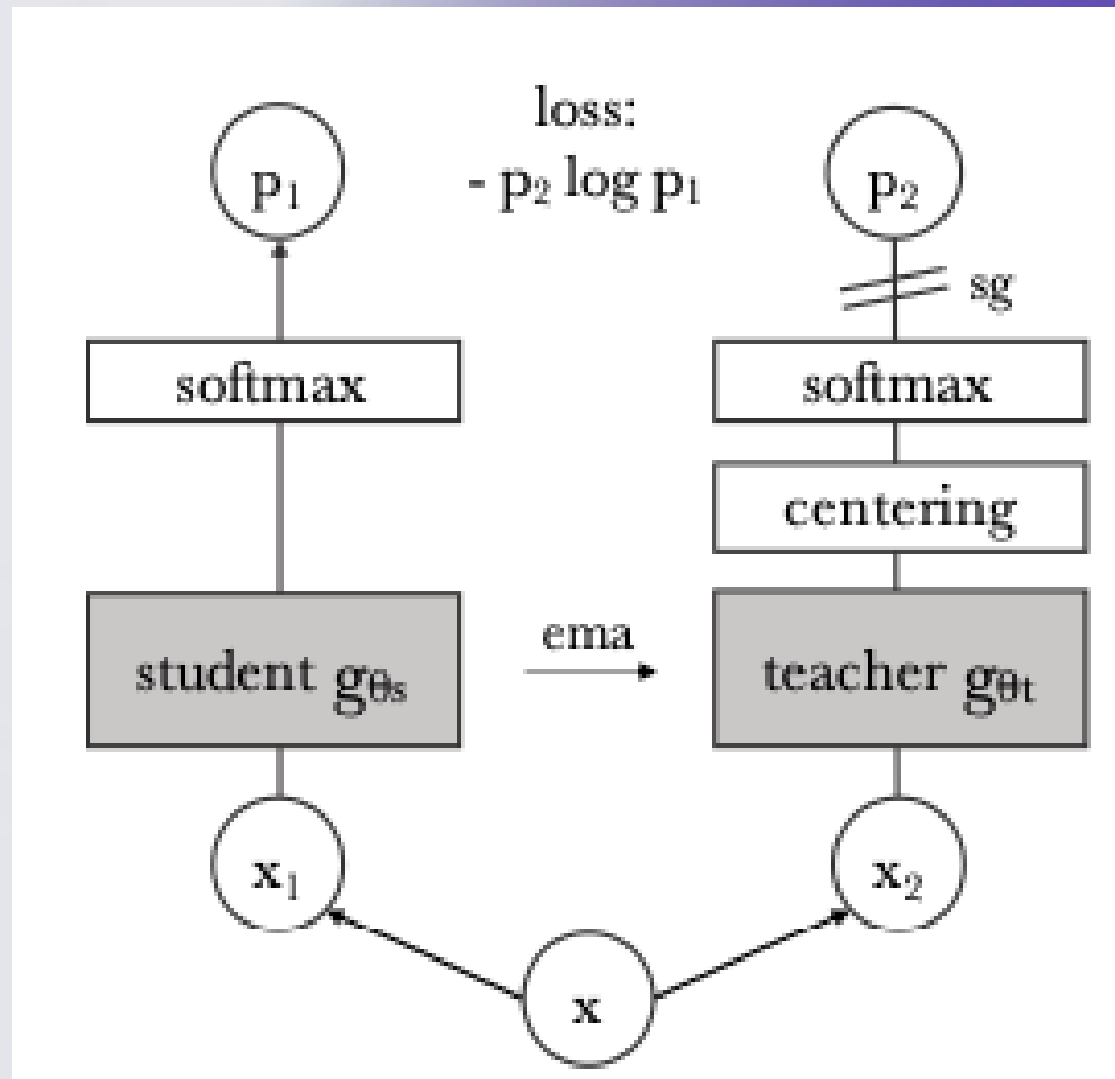


# DINO

CLASS 9/3/24



# Idea

Image-Text grounding vs Self-Supervised

# Algorithm

---

## Algorithm 1 DINO PyTorch pseudocode w/o multi-crop.

---

```
# gs, gt: student and teacher networks
# C: center (K)
# tps, tpt: student and teacher temperatures
# l, m: network and center momentum rates
gt.params = gs.params
for x in loader: # load a minibatch x with n samples
    x1, x2 = augment(x), augment(x) # random views

    s1, s2 = gs(x1), gs(x2) # student output n-by-K
    t1, t2 = gt(x1), gt(x2) # teacher output n-by-K

    loss = H(t1, s2)/2 + H(t2, s1)/2
    loss.backward() # back-propagate

    # student, teacher and center updates
    update(gs) # SGD
    gt.params = l*gt.params + (1-l)*gs.params
    C = m*C + (1-m)*cat([t1, t2]).mean(dim=0)

def H(t, s):
    t = t.detach() # stop gradient
    s = softmax(s / tps, dim=1)
    t = softmax((t - C) / tpt, dim=1) # center + sharpen
    return - (t * log(s)).sum(dim=1).mean()
```

---



# Architecture

- Both student and teacher share the same architecture, just different params
- Both teacher and student outputs a  $K$  dimensional vector, which is normalized with a temperature softmax
- EMA updates to prevent model collapse
  - Centering+Sharpening to prevent feature collapse

Algorithm 1 DINO PyTorch pseudocode w/o multi-crop.

```
# gs, gt: student and teacher networks
# C: center (K)
# tps, tpt: student and teacher temperatures
# l, m: network and center momentum rates
gt.params = gs.params
for x in loader: # load a minibatch x with n samples
    x1, x2 = augment(x), augment(x) # random views

    s1, s2 = gs(x1), gs(x2) # student output n-by-K
    t1, t2 = gt(x1), gt(x2) # teacher output n-by-K

    loss = H(t1, s2)/2 + H(t2, s1)/2
    loss.backward() # back-propagate

    # student, teacher and center updates
    update(gs) # SGD
    gt.params = l*gt.params + (1-l)*gs.params
    C = m*C + (1-m)*cat([t1, t2]).mean(dim=0)

def H(t, s):
    t = t.detach() # stop gradient
    s = softmax(s / tps, dim=1)
    t = softmax((t - C) / tpt, dim=1) # center + sharpen
    return - (t * log(s)).sum(dim=1).mean()
```

```
t = softmax((t - C) / tpt, dim=1) # center + sharpen
```

$$g_t(x) \leftarrow g_t(x) + c$$

$$c \leftarrow mc + (1 - m) \frac{1}{B} \sum_{i=1}^B g_{\theta_t}(x_i)$$

# Loss

$$P_s(x)^{(i)} = \frac{\exp(g_{\theta_s}(x)^{(i)} / \tau_s)}{\sum_{k=1}^K \exp(g_{\theta_s}(x)^{(k)} / \tau_s)}$$

$$\min_{\theta_s} H(P_t(x), P_s(x))$$

$$H(a, b) = -a \log b$$

$$\min_{\theta_s} \sum_{x \in \{x_1^g, x_2^g\}} \sum_{\substack{x' \in V \\ x' \neq x}} H(P_t(x), P_s(x'))$$

- 2 global views of 224x224 covering  $\geq 50\%$  of  $x$
- Several local views of 96x96  $< 50\%$  of  $x$

## Exponential Moving Average (EMA)

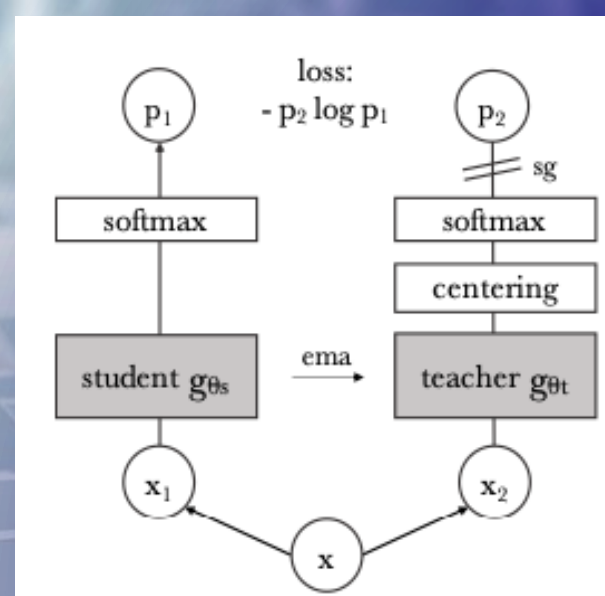
$$\theta_t \leftarrow \lambda \theta_t + (1 - \lambda) \theta_s$$

$$\eta_t = \eta_T + \frac{\eta_0 - \eta_T}{2} (1 + \cos(\pi t / T))$$

$\lambda$  follows a cosine scheduling from 0.996 to 1

$$h: g = h \circ f$$

$h$  is a 3 layer MLP - hidden dimension 2048, L2 normalization, weight normalized fully connected layer



# Preventing Collapse



# Other Details

- Patch sizes: 8x8, 16x16
- $h$  is attached to the "[CLS]" token
- Backbone (ResNet, ViT) pretrained on ImageNet

# Evaluation

- Linear evaluation, Fine-tuning, K-NN (20-NN)
- K-NN needs no hyperparam tuning

Table 2: **Linear and  $k$ -NN classification on ImageNet.** We report top-1 accuracy for linear and  $k$ -NN evaluations on the validation set of ImageNet for different self-supervised methods. We focus on ResNet-50 and ViT-small architectures, but also report the best results obtained across architectures. \* are run by us. We run the  $k$ -NN evaluation for models with official released weights. The throughput (im/s) is calculated on a NVIDIA V100 GPU with 128 samples per forward. Parameters (M) are of the feature extractor.

Method	Arch.	Param.	im/s	Linear	$k$ -NN
Supervised	RN50	23	1237	79.3	79.3
SCLR [11]	RN50	23	1237	69.1	60.7
MoCov2 [13]	RN50	23	1237	71.1	61.9
InfoMin [54]	RN50	23	1237	73.0	65.3
BarlowT [66]	RN50	23	1237	73.2	66.0
OBoW [21]	RN50	23	1237	73.8	61.9
BYOL [23]	RN50	23	1237	74.4	64.8
DCv2 [9]	RN50	23	1237	75.2	67.1
SwAV [9]	RN50	23	1237	<b>75.3</b>	65.7
DINO	RN50	23	1237	<b>75.3</b>	<b>67.5</b>
Supervised	ViT-S	21	1007	79.8	79.8
BYOL* [23]	ViT-S	21	1007	71.4	66.6
MoCov2* [13]	ViT-S	21	1007	72.7	64.4
SwAV* [9]	ViT-S	21	1007	73.5	66.3
DINO	ViT-S	21	1007	<b>77.0</b>	<b>74.5</b>
<i>Comparison across architectures</i>					
SCLR [11]	RN50w4	375	117	76.8	69.3
SwAV [9]	RN50w2	93	384	77.3	67.3
BYOL [23]	RN50w2	93	384	77.4	—
DINO	ViT-B/16	85	312	78.2	76.1
SwAV [9]	RN50w5	586	76	78.5	67.1
BYOL [23]	RN50w4	375	117	78.6	—
BYOL [23]	RN200w2	250	123	79.6	73.9
DINO	ViT-S/8	21	180	79.7	<b>78.3</b>
SCLRv2 [12]	RN152w3+SK	794	46	79.8	73.1
DINO	ViT-B/8	85	63	<b>80.1</b>	77.4



Table 3: **Image retrieval.** We compare the performance in retrieval of off-the-shelf features pretrained with supervision or with DINO on ImageNet and Google Landmarks v2 (GLDv2) dataset. We report mAP on revisited Oxford and Paris. Pretraining with DINO on a landmark dataset performs particularly well. For reference, we also report the best retrieval method with off-the-shelf features [46].

Pretrain	Arch.	Pretrain	ROx		RPar	
			M	H	M	H
Sup. [46]	RN101+R-MAC	ImNet	49.8	18.5	74.0	<b>52.1</b>
Sup.	ViT-S/16	ImNet	33.5	8.9	63.0	37.2
DINO	ResNet-50	ImNet	35.4	11.1	55.9	27.5
DINO	ViT-S/16	ImNet	41.8	13.7	63.1	34.4
DINO	ViT-S/16	GLDv2	<b>51.5</b>	<b>24.3</b>	<b>75.3</b>	51.6

Table 5: **DAVIS 2017 Video object segmentation.** We evaluate the quality of frozen features on video instance tracking. We report mean region similarity  $\mathcal{J}_m$  and mean contour-based accuracy  $\mathcal{F}_m$ . We compare with existing self-supervised methods and a supervised ViT-S/8 trained on ImageNet. Image resolution is 480p.

Method	Data	Arch.	$(\mathcal{J}\&\mathcal{F})_m$	$\mathcal{J}_m$	$\mathcal{F}_m$
<i>Supervised</i>					
ImageNet	INet	ViT-S/8	66.0	63.9	68.1
STM [39]	I/D/Y	RN50	81.8	79.2	84.3
<i>Self-supervised</i>					
CT [58]	VLOG	RN50	48.7	46.4	50.0
MAST [33]	YT-VOS	RN18	65.5	63.3	67.6
STC [30]	Kinetics	RN18	67.6	64.8	70.2
DINO	INet	ViT-S/16	61.8	60.2	63.4
DINO	INet	ViT-B/16	62.3	60.7	63.9
DINO	INet	ViT-S/8	<b>69.9</b>	<b>66.6</b>	<b>73.1</b>
DINO	INet	ViT-B/8	<b>71.4</b>	<b>67.9</b>	<b>74.9</b>

Table 4: **Copy detection.** We report the mAP performance in copy detection on Copydays “strong” subset [18]. For reference, we also report the performance of the multigrain model [4], trained specifically for particular object retrieval.

Method	Arch.	Dim.	Resolution	mAP
Multigrain [4]	ResNet-50	2048	224 <sup>2</sup>	75.1
Multigrain [4]	ResNet-50	2048	largest side 800	82.5
Supervised [56]	ViT-B/16	1536	224 <sup>2</sup>	76.4
DINO	ViT-B/16	1536	224 <sup>2</sup>	81.7
DINO	ViT-B/8	1536	320 <sup>2</sup>	<b>85.5</b>

Table 6: **Transfer learning by finetuning pretrained models on different datasets.** We report top-1 accuracy. Self-supervised pretraining with DINO transfers better than supervised pretraining.

	Cifar <sub>10</sub>	Cifar <sub>100</sub>	INat <sub>18</sub>	INat <sub>19</sub>	Flwrs	Cars	INet
<i>ViT-S/16</i>							
Sup. [56]	<b>99.0</b>	89.5	70.7	76.6	98.2	92.1	79.9
DINO	<b>99.0</b>	<b>90.5</b>	<b>72.0</b>	<b>78.2</b>	<b>98.5</b>	<b>93.0</b>	<b>81.5</b>
<i>ViT-B/16</i>							
Sup. [56]	99.0	90.8	<b>73.2</b>	77.7	98.4	92.1	81.8
DINO	<b>99.1</b>	<b>91.7</b>	72.6	<b>78.6</b>	<b>98.8</b>	<b>93.0</b>	<b>82.8</b>

# More Results

# Ablations

Table 7: **Important component for self-supervised ViT pre-training.** Models are trained for 300 epochs with ViT-S/16. We study the different components that matter for the  $k$ -NN and linear (“Lin.”) evaluations. For the different variants, we highlight the differences from the default DINO setting. The best combination is the momentum encoder with the multicrop augmentation and the cross-entropy loss. We also report results with BYOL [23], MoCo-v2 [13] and SwAV [9].

	Method	Mom.	SK	MC	Loss	Pred.	$k$ -NN	Lin.
1	DINO	✓	✗	✓	CE	✗	72.8	76.1
2		✗	✗	✓	CE	✗	0.1	0.1
3		✓	✓	✓	CE	✗	72.2	76.0
4		✓	✗	✗	CE	✗	67.9	72.5
5		✓	✗	✓	MSE	✗	52.6	62.4
6		✓	✗	✓	CE	✓	71.8	75.6
7	BYOL	✓	✗	✗	MSE	✓	66.6	71.4
8	MoCov2	✓	✗	✗	INCE	✗	62.0	71.6
9	SwAV	✗	✓	✓	CE	✗	64.7	71.8

SK: Sinkhorn-Knopp, MC: Multi-Crop, Pred.: Predictor  
CE: Cross-Entropy, MSE: Mean Square Error, INCE: InfoNCE

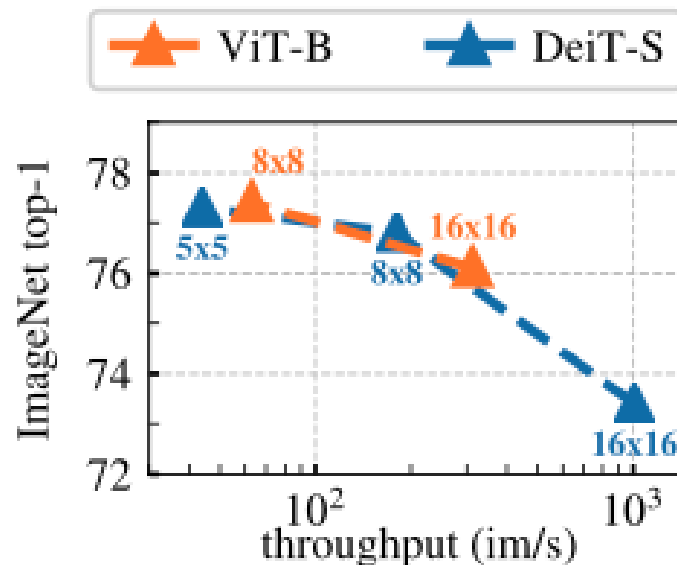


Figure 5: **Effect of Patch Size.**  $k$ -NN evaluation as a function of the throughputs for different input patch sizes with ViT-B and ViT-S. Models are trained for 300 epochs.

Attendance Time ...





Threshold to retain 60% of attention mass

# Qualitative Results



Figure 3: **Attention maps from multiple heads.** We consider the heads from the last layer of a ViT-S/8 trained with DINO and display the self-attention for [CLS] token query. Different heads, materialized by different colors, focus on different locations that represents different objects or parts (more examples in Appendix).

# DINOv2

- Curating a large, diverse dataset
- Stabilizing large scale training on this dataset
- Trained a 1B ViT from which smaller models are distilled





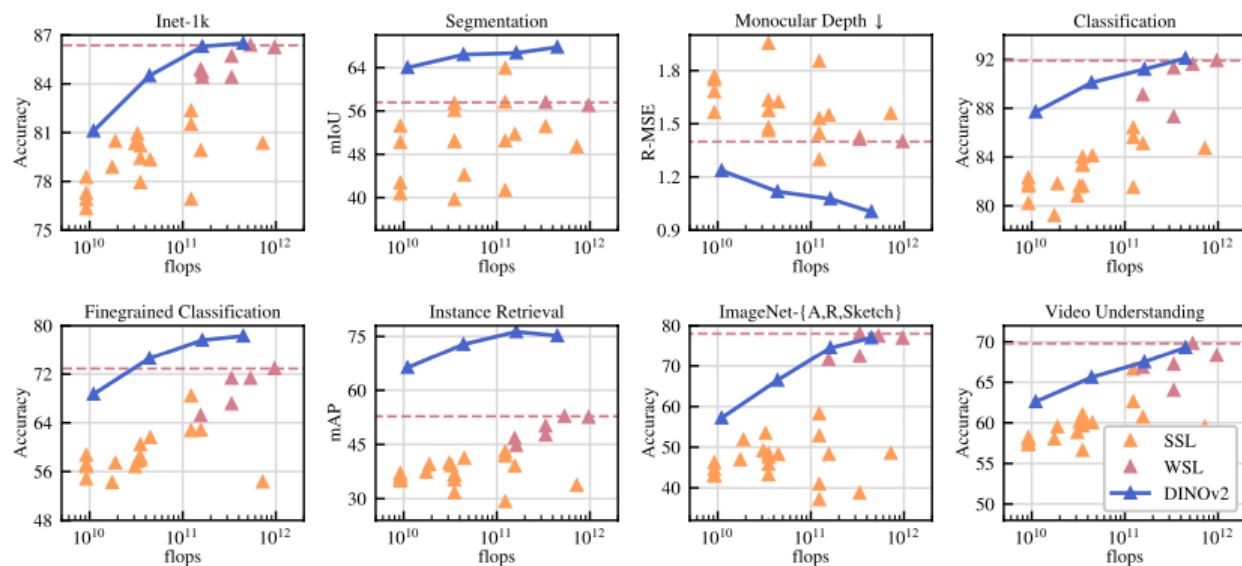


Figure 2: **Evolution of performance when scaling in parameters.** We show performance on eight types of vision tasks, as presented in Sec. 7 and average metrics with each type. Features are extracted from our self-supervised encoders, DINOv2 (dark blue), and we compare them with self-supervised methods (pale orange), as well as weakly-supervised methods (dark pink). We report the best-performing weakly-supervised model's performance as a dashed horizontal line. Our family of models drastically improves over the previous state of the art in self-supervised learning and reaches performance comparable with weakly-supervised features. See Sec. 7 for a detailed analysis.

# Quick Results



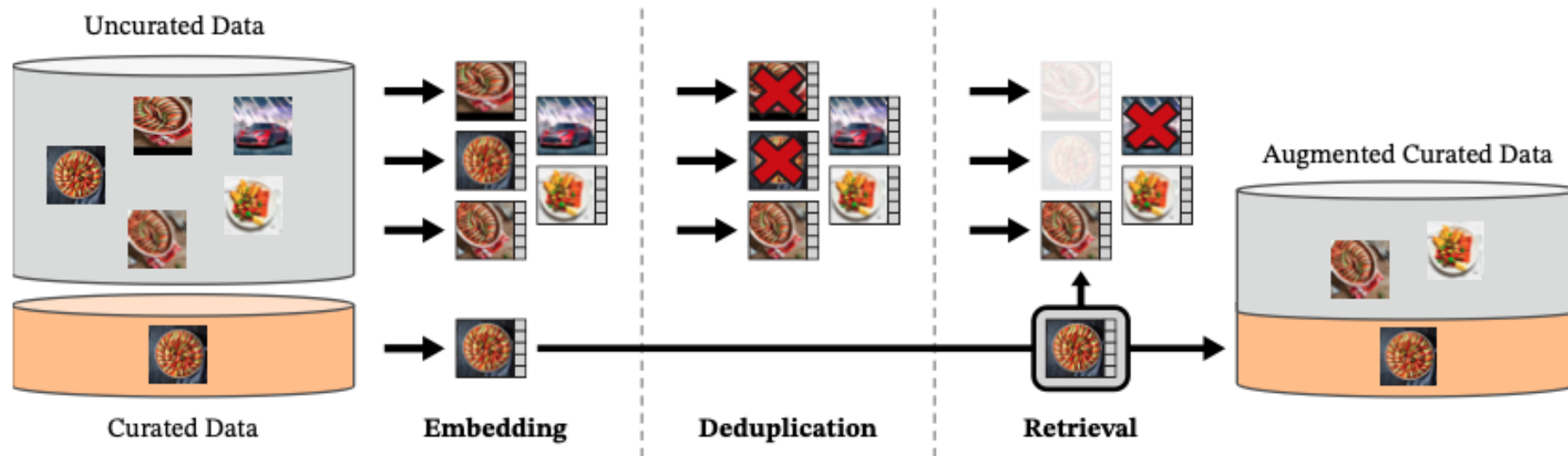


Figure 3: **Overview of our data processing pipeline.** Images from curated and uncurated data sources are first mapped to embeddings. Uncurated images are then deduplicated before being matched to curated images. The resulting combination augments the initial dataset through a self-supervised retrieval system.

Curated data is from ImageNet22k/1k, Google Landmarks, fine-grained datasets -> 142M

# Data Curation Pipeline

# Main Components

Arch	Method	INet-1k	Segm.	Depth↓	Classif.
ViT-g/14	Scratch	86.5	73.4	1.00	92.1
ViT-L/14	Scratch	84.5	72.2	1.10	90.2
ViT-L/14	Distill	<b>86.3</b>	<b>73.3</b>	<b>1.08</b>	<b>91.2</b>

Arch	Method	Finegr.	Retriev.	ARSketch	Video
ViT-g/14	Scratch	78.3	75.2	77.0	69.3
ViT-L/14	Scratch	75.8	71.3	69.5	67.3
ViT-L/14	Distill	<b>77.6</b>	<b>76.3</b>	<b>74.5</b>	<b>67.5</b>

- Image level objective - Dino
- Pixel level objective - random masking of patches
- Untying heads between teacher and student objective
- Sinkhorn-Knopp batch normalization in lieu of centering
- Koleo regularizer
- 518x518 at end of pretraining
- FlashAttention - HRAM vs SRAM (fastest)
- Same forward pass global and local crops
- Stochastic depth - randomly drop residual blocks with skip
- FSDP - Fully Sharded Data Parallel
- Model distillation from ViT-1B



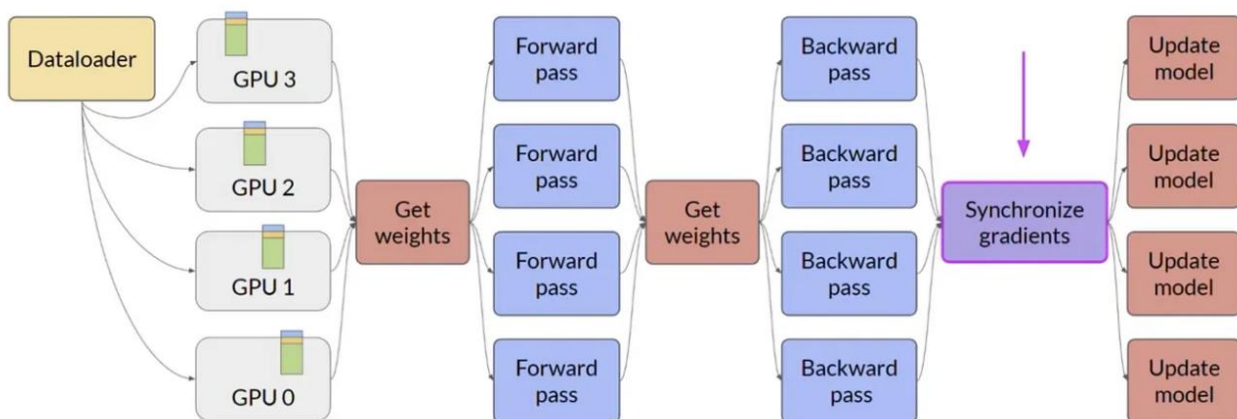
# DinoV2 Efficiency Playbook

- FlashAttention!!
- Sequence packing
  - Crops produce different tokens length (patches are different)
  - These tokens therefore have to be forwarded separately
  - But, pack them into one long sequence, put special separator tokens ... prevent attention between different sequences
- Stochastic depth
  - Original stochastic depth is randomly dropping layers (identity matrix)
  - Skips the computation of the residuals
- FSDP (Fully Sharded Data Parallel)
- Model distillation
  - Smaller ViT distilled from ViT-1B

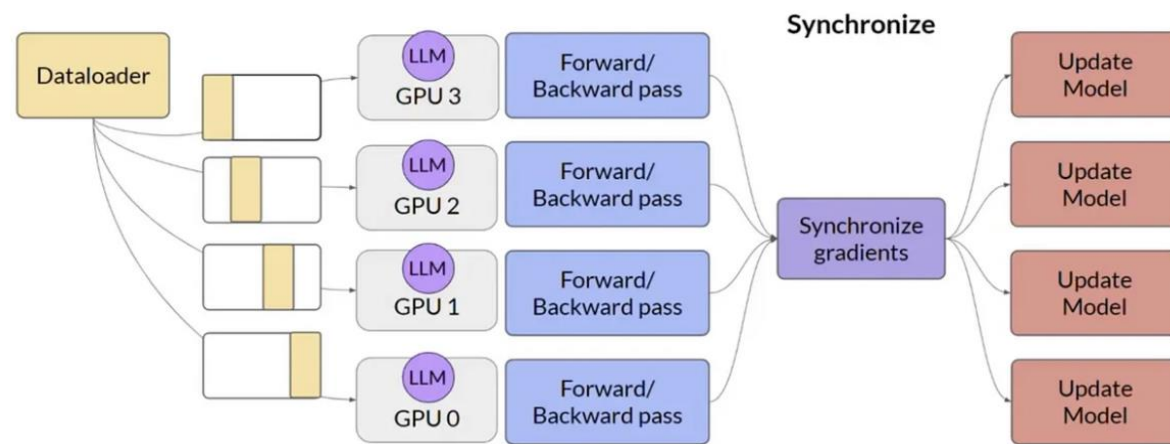


# FSDP vs DDP (Distributed Data Parallel)

## Fully Sharded Data Parallel (FSDP)



## Distributed Data Parallel (DDP)



# Results, Assignment Due 9/17

- Train DinoV2 on Human Action Recognition dataset
  - Same train/test split
  - Use ViT-S
  - 8 patches
- Evaluation
  - Test set classification (linear)
  - Speed: Turn off FlashAttention vs FlashAttention on (both training and testing speed)
- Report, code, video demo