# Quick preview-TimeBalance Framework



Dual teacher- Student framework for Semi-supervised Action Recognition

Teachers are pretrained with different types of *self-supervised video representations*: *temporally-invariant* and *temporally-distinctive* representations

Input video clip: processed by both teachers

Predictions combined: reweighting strategy

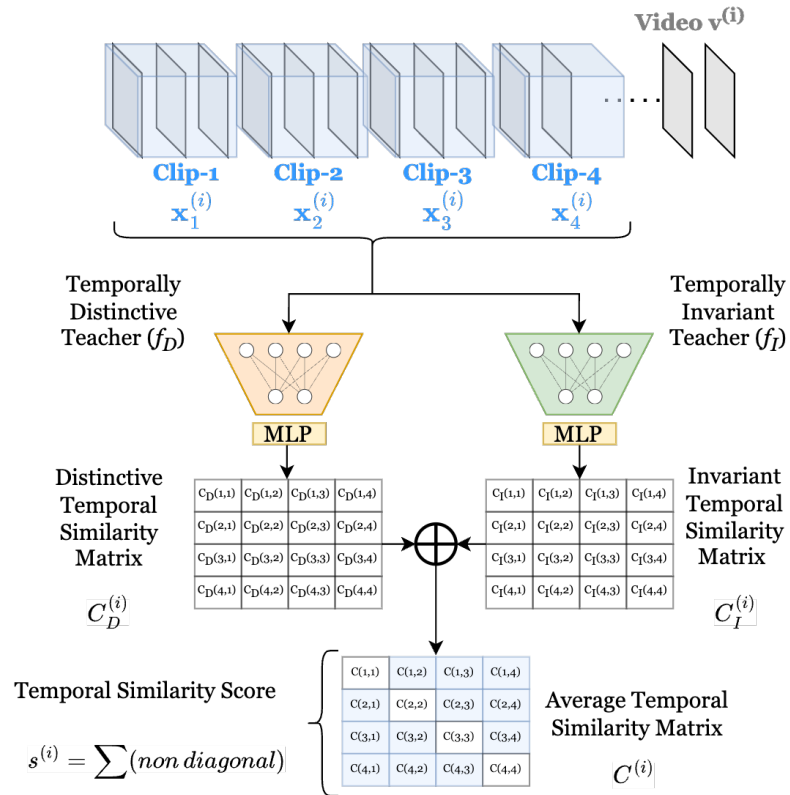Knowledge distillation: Provides unlabeled supervision to the student

Labeled supervision: Provided when labels are available

# Quick preview- Reweighting strategy



- Clips representation similarity ($\uparrow$) $\rightarrow$ $f_I$($\uparrow$)
- Clips representation similarity ($\downarrow$) $\rightarrow$ $f_D$ ($\uparrow$)

# Quick preview- Results

| Method | Venue | Backbone | Params (M) | Input | #F | UCF101 | | | | | HMDB51 | | | Kinetics400 | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | | | | 1% | 5% | 10% | 20% | 50% | 40% | 50% | 60% | 1% | 10% |
| PL | ICML'13 | 3D-ResNet18 | 13.5 | V | 16 | - | 17.6 | 24.7 | 37 | 47.5 | 27.3 | 32.4 | 33.5 | - | - |
| MT | NeuRIPS'17 | 3D-ResNet18 | 13.5 | V | 16 | - | 17.5 | 25.6 | 36.3 | 45.8 | 27.2 | 30.4 | 32.2 | - | - |
| S4L | ICCV'19 | 3D-ResNet18 | 13.5 | V | 16 | - | 22.7 | 29.1 | 37.7 | 47.9 | 29.8 | 31 | 35.6 | - | - |
| UPS | ICLR'21 | 3D-ResNet18 | 13.5 | V | 16 | - | - | - | 39.4 | 50.2 | - | - | - | - | - |
| SD | ICCV'19 | 3D-ResNet18 | 13.5 | V | 16 | - | 31.2 | 40.7 | 45.4 | 53.9 | 32.6 | 35.1 | 36.3 | - | - |
| MT+SD | WACV'21 | 3D-ResNet18 | 13.5 | V | 16 | - | 30.3 | 40.5 | 45.5 | 53 | 32.3 | 33.6 | 35.7 | - | - |
| VideoSemi | WACV'21 | 3D-ResNet18 | 13.5 | V | 16 | - | 32.4 | 42 | 48.7 | 54.3 | 32.7 | 36.2 | 37 | - | - |
| TG-FixMatch | CVPR'21 | 3D-ResNet18 | 13.5 | V | 8 | - | 44.8 | 62.4 | 76.1 | 79.3 | 46.5 | 48.4 | 49.7 | 9.8 | 43.8 |
| TCL | CVPR'21 | TSM-R18 | - | V | 8 | - | - | - | - | - | - | - | - | 11.6 | - |
| MvPL | ICCV'21 | 3D-ResNet18 | 13.5 | VFG | 8 | - | 41.2 | 55.5 | 64.7 | 65.6 | 30.5 | 33.9 | 35.8 | 5 | 36.9 |
| CMPL | CVPR'22 | 3D-ResNet18 | 13.5 | V | 8 | 23.8 | - | 67.6 | - | - | - | - | - | 16.5 | 53.7 |
| TACL | TSVT'22 | 3D-ResNet18 | 13.5 | V | 16 | - | 35.6 | 50.9 | 56.1 | 65.8 | 34.6 | 37.2 | 39.5 | - | - |
| TACL | TSVT'22 | 3D-ResNet18 | 13.5 | V | 16 | - | 43.7 | 55.6 | 59.2 | 67.2 | 38.7 | 40.2 | 41.7 | - | - |
| 3DRotNet | Arxiv'19 | 3D-ResNet18 | 13.5 | V | 16 | 15 | 31.5 | 40.4 | 47.1 | - | - | - | - | - | - |
| MemDPC | ECCV'20 | 3D-ResNet18 | 13.5 | V | 16 | - | - | 44.2 | 50.9 | 62.3 | - | - | - | - | - |
| MotionFit | ICCV'21 | 3D-ResNet18 | 13.5 | VF | 16 | - | - | - | 57.7 | 59 | - | - | - | - | - |
| TCLR | CVIU'22 | 3D-ResNet18 | 13.5 | V | 16 | 26.9 | - | 66.1 | 73.4 | 76.7 | - | - | - | - | - |
| TimeBalance | CVPR'23 | 3D-ResNet18 | 13.5 | V | 8 | **29.1** | **47.9** | **69.8** | **79.1** | **83.3** | **49.8** | **51.4** | **53.1** | **17.1** | **54.9** |
| ActorCM | Arxiv'21 | R(2+1)D-34 | 33.3 | V | 8 | - | 27 | 40.2 | 51.7 | 59.9 | 32.9 | 38.2 | 38.9 | - | - |
| ActorCM | Arxiv'21 | R(2+1)D-34 | 33.3 | V | 8 | - | 45.1 | 53 | 57.4 | 64.7 | 35.7 | 39.5 | 40.8 | - | - |
| FixMatch | NeuRIPS'20 | SlowFast-R50 | 60 | V | 8 | 16.1 | - | 55.1 | - | - | - | - | - | 10.1 | 49.4 |
| MvPL | ICCV'21 | 3D-ResNet50 | 31.8 | VFG | 8 | 22.8 | - | 80.5 | - | - | - | - | - | 17 | 58.2 |
| CMPL | CVPR'22 | 3D-ResNet50 | 31.8 | V | 8 | 25.1 | - | 79.1 | - | - | - | - | - | 17.6 | 58.4 |
| TimeBalance | CVPR'23 | 3D-ResNet50 | 31.8 | V | 8 | **30.1** | **53.5** | **81.1** | **83.3** | **85** | **52.6** | **53.9** | **54.5** | **19.6** | **61.2** |

# Overview of the details

- Motivation

- Temporal-Invariance vs Distinctiveness in video representations

- TimeBalance framework

- Temporal Similarity Based Reweighting

- Experiments

- Conclusion

# Semi-supervised Action Recognition

- Semi-supervised learning deals with enhancing performance on smaller labeled set leveraging large amount of unlabeled data.

- Semi supervised learning is more crucial for video data due to **higher cost** of video annotations

- Compared to images, videos provide additional **temporal dimension** to exploit to leverage unlabeled data

# Prior work

Prior work have been relied on:

- Two-stream networks: RGB + Optical Flow/ Temporal Gradient

- Two-stream skip rate: Fast + slow playback

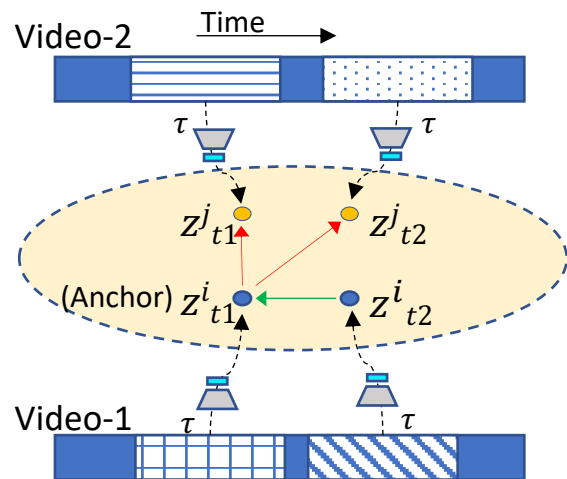Prior work have input specific inductive-bias

- Requires precomputation of datasets. e.g., computing flow takes about a week for Kinetics400

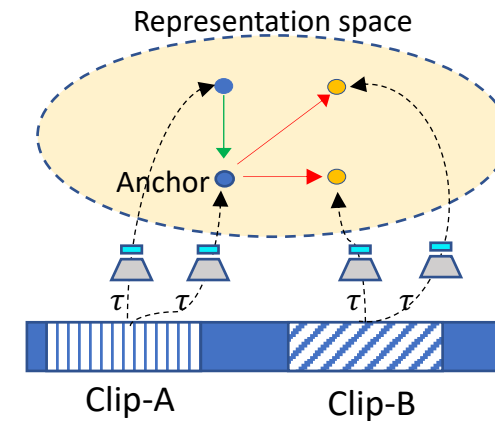Prior work have not systematically leveraged SSL video representations

UCF CENTER FOR RESEARCH IN COMPUTER VISION

# Self-supervised Representations for Videos



**Temporally Invariant ($f_I$)**

**Temporally Distinctive ($f_D$)**

Legend:
- MLP Projection Head
- 3D-CNN backbone
- Attract
- Repel
- $\tau$   Random Augmentation

Temporally Invariant diagram:
- Video-2, Time →
- $z^j_{t1}$, $z^j_{t2}$
- (Anchor) $z^i_{t1}$, $z^i_{t2}$
- Video-1

Temporally Distinctive diagram:
- Representation space
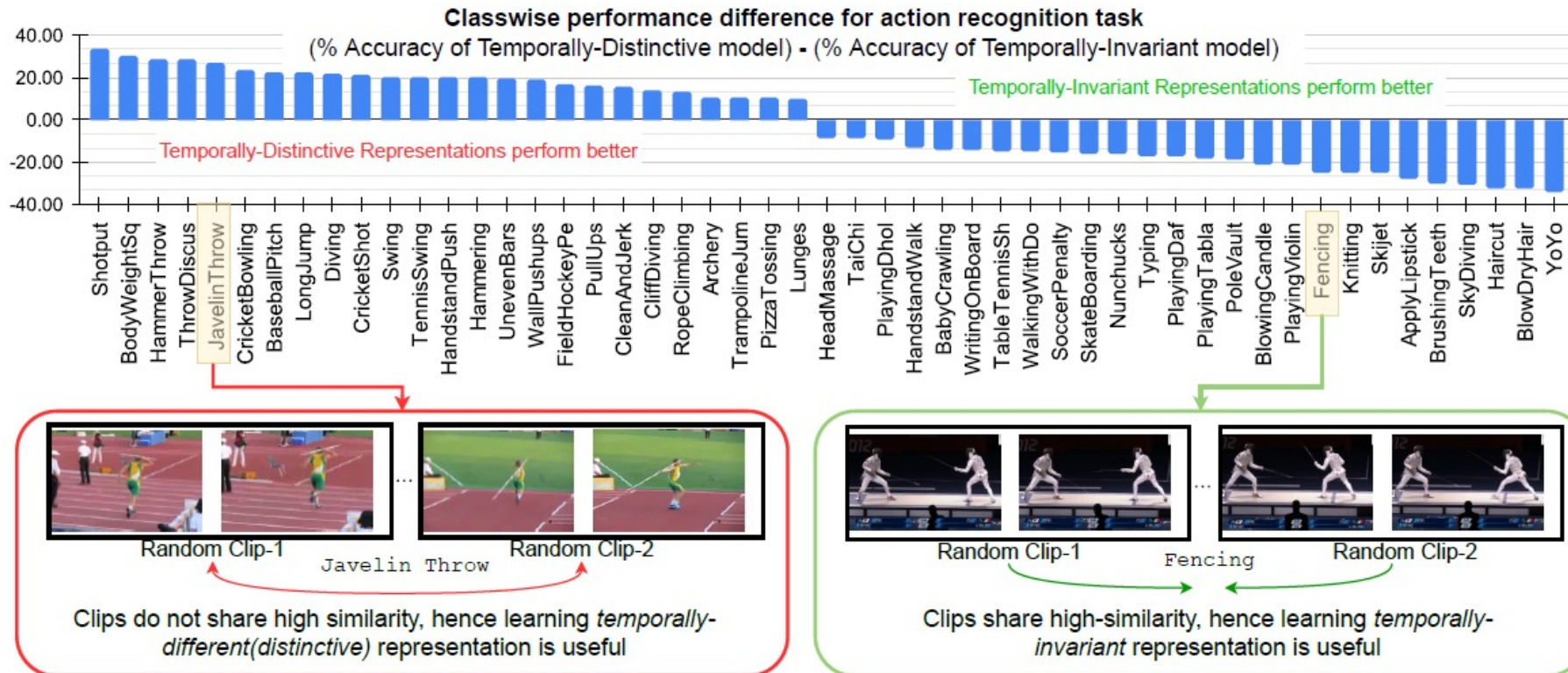- Anchor
- Clip-A, Clip-B

Contrastive loss/ SSL loss encourages **two clip** of video to have **same** representation

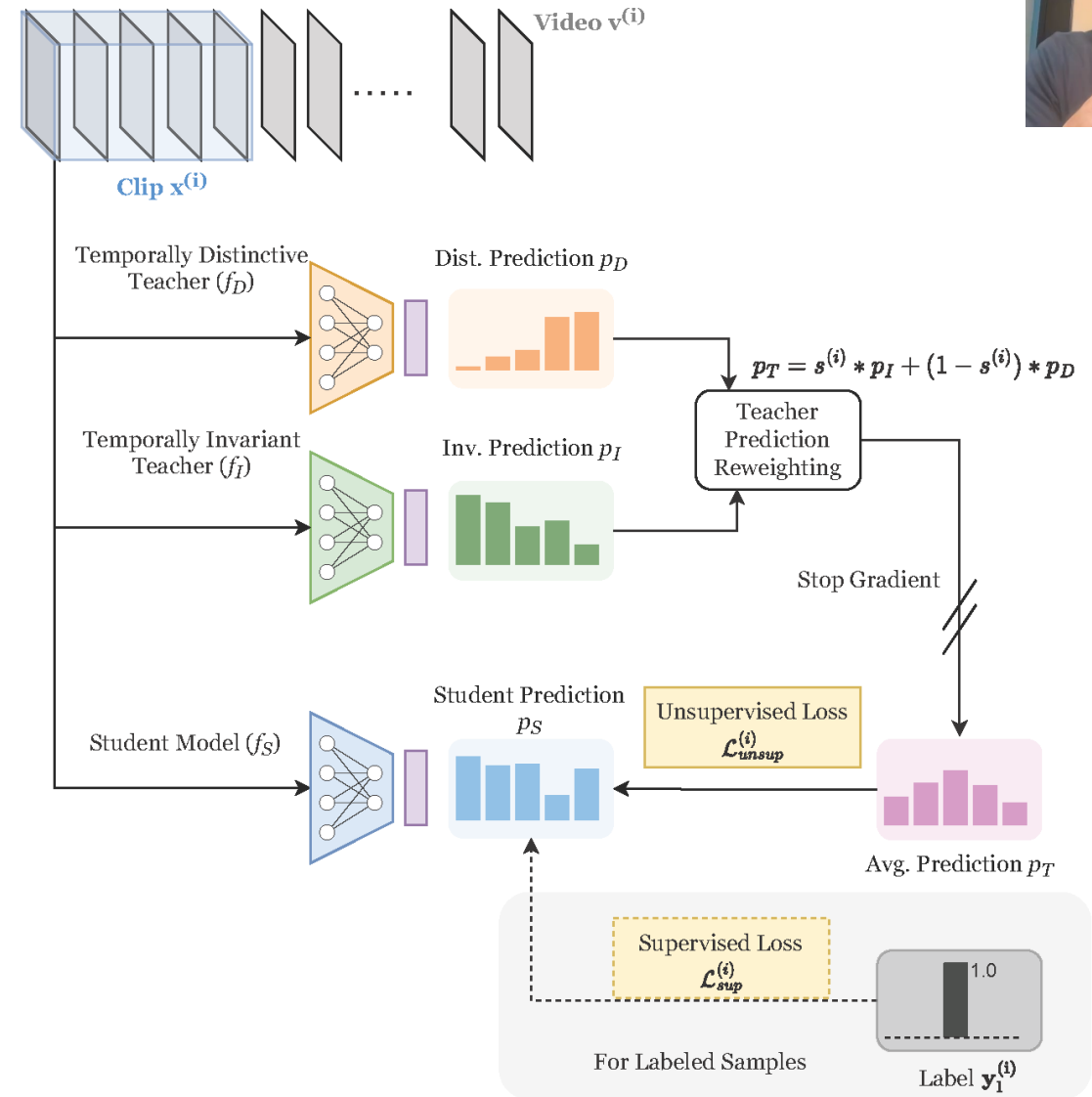Contrastive loss/ SSL loss encourages **two clip** of video to have **different** representation

UCF CENTER FOR RESEARCH IN COMPUTER VISION

# Complementary behavior of $f_I$ and $f_D$



**Classwise performance difference for action recognition task**
(% Accuracy of Temporally-Distinctive model) - (% Accuracy of Temporally-Invariant model)

Temporally-Invariant Representations perform better

Temporally-Distinctive Representations perform better

Javelin Throw

Clips do not share high similarity, hence learning *temporally-different(distinctive)* representation is useful

Fencing

Clips share high-similarity, hence learning *temporally-invariant* representation is useful

# TimeBalance Framework

- Dual teacher- Student framework

- Both teachers are pretrained in self-supervised way and finetuned with the limited labeled data

- Prediction of both teachers are combined with proposed temporal-similarity based reweighting scheme



$$p_T = s^{(i)} * p_I + (1 - s^{(i)}) * p_D$$

UCF CENTER FOR RESEARCH IN COMPUTER VISION

# Temporal Similarity Based Reweighting

- Consecutive clips are sampled from a video

- Clips are passed through the $f_I$ and $f_D$ models to get their SSL representation

- Temporal self-similarity matrix is computed using cosine distance

- Obtain a temporal self-similarity score

- Clips representation similarity ($\uparrow$) $\rightarrow$ $f_I$($\uparrow$)
- Clips representation similarity ($\downarrow$) $\rightarrow$ $f_D$ ($\uparrow$)



$$\mathbf{p}_T^{(i)} = s^{(i)} \cdot \mathbf{p}_I^{(i)} + (1 - s^{(i)}) \cdot \mathbf{p}_D^{(i)}$$

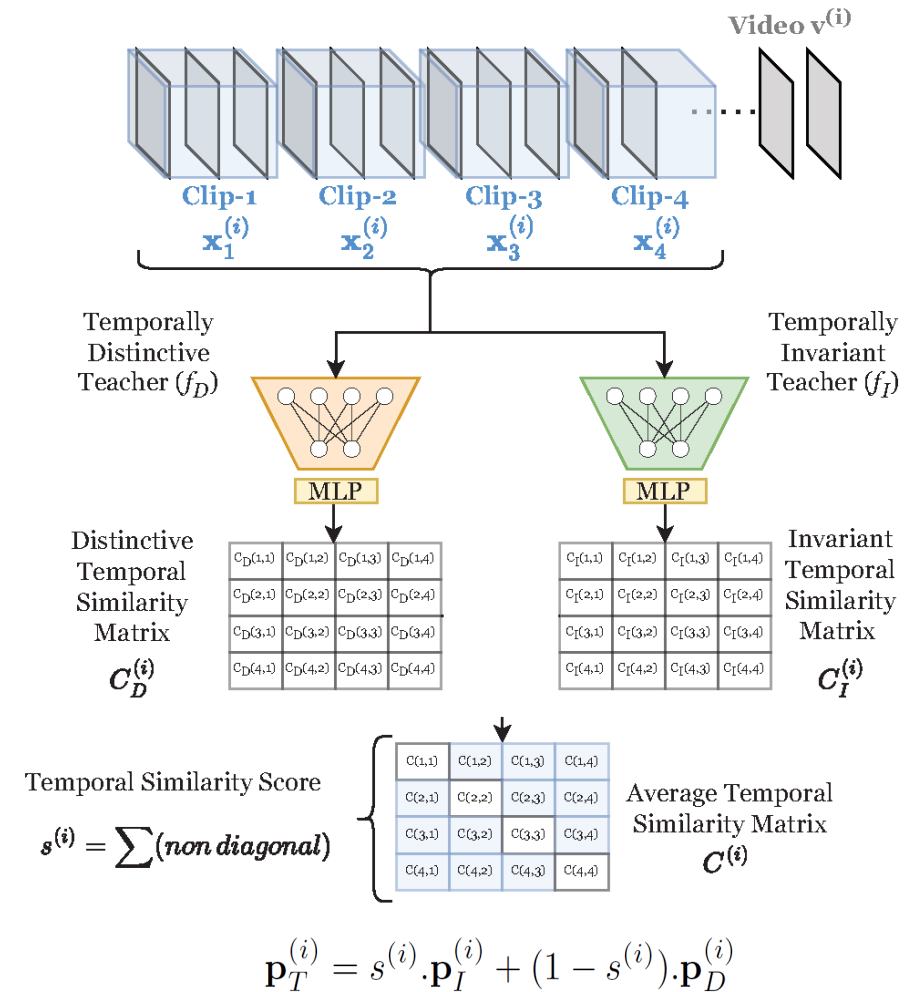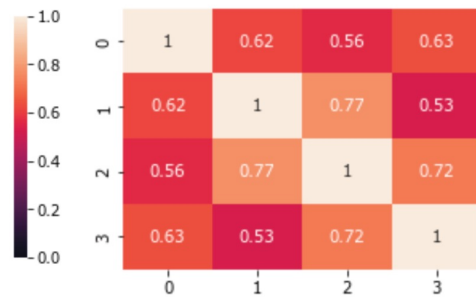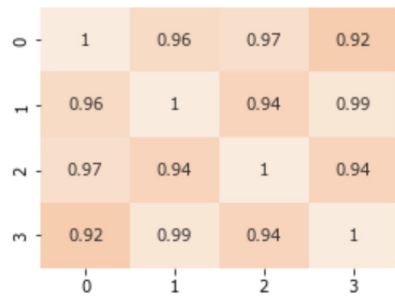UCF CENTER FOR RESEARCH IN COMPUTER VISION

# Illustration of teacher reweighting



Temporal Similarity Matrix from Invariant teacher $C_I$

Temporal Similarity Matrix from Distinctive teacher $C_D$

Instance Similarity score $s = 0.795$

Repetitive action

High similarity score

More weight given to $f_I$

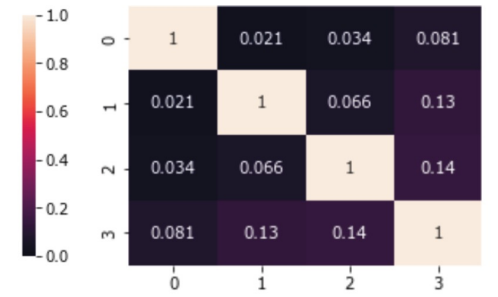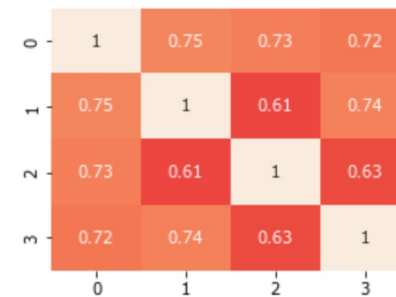Temporal Similarity Matrix from Invariant teacher $C_I$

Temporal Similarity Matrix from Distinctive teacher $C_D$

Instance Similarity score $s = 0.387$

Non-Repetitive action

Low similarity score

More weight given to $f_D$

UCF CENTER FOR RESEARCH IN COMPUTER VISION

# Results

| Method | Venue | Backbone | Params (M) | Input | #F | UCF101 | | | | | HMDB51 | | | Kinetics400 | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | | | | 1% | 5% | 10% | 20% | 50% | 40% | 50% | 60% | 1% | 10% |
| PL | ICML'13 | 3D-ResNet18 | 13.5 | V | 16 | - | 17.6 | 24.7 | 37 | 47.5 | 27.3 | 32.4 | 33.5 | - | - |
| MT | NeuRIPS'17 | 3D-ResNet18 | 13.5 | V | 16 | - | 17.5 | 25.6 | 36.3 | 45.8 | 27.2 | 30.4 | 32.2 | - | - |
| S4L | ICCV'19 | 3D-ResNet18 | 13.5 | V | 16 | - | 22.7 | 29.1 | 37.7 | 47.9 | 29.8 | 31 | 35.6 | - | - |
| UPS | ICLR'21 | 3D-ResNet18 | 13.5 | V | 16 | - | - | - | 39.4 | 50.2 | - | - | - | - | - |
| SD | ICCV'19 | 3D-ResNet18 | 13.5 | V | 16 | - | 31.2 | 40.7 | 45.4 | 53.9 | 32.6 | 35.1 | 36.3 | - | - |
| MT+SD | WACV'21 | 3D-ResNet18 | 13.5 | V | 16 | - | 30.3 | 40.5 | 45.5 | 53 | 32.3 | 33.6 | 35.7 | - | - |
| VideoSemi | WACV'21 | 3D-ResNet18 | 13.5 | V | 16 | - | 32.4 | 42 | 48.7 | 54.3 | 32.7 | 36.2 | 37 | - | - |
| TG-FixMatch | CVPR'21 | 3D-ResNet18 | 13.5 | V | 8 | - | 44.8 | 62.4 | 76.1 | 79.3 | 46.5 | 48.4 | 49.7 | 9.8 | 43.8 |
| TCL | CVPR'21 | TSM-R18 | - | V | 8 | - | - | - | - | - | - | - | - | 11.6 | - |
| MvPL | ICCV'21 | 3D-ResNet18 | 13.5 | VFG | 8 | - | 41.2 | 55.5 | 64.7 | 65.6 | 30.5 | 33.9 | 35.8 | 5 | 36.9 |
| CMPL | CVPR'22 | 3D-ResNet18 | 13.5 | V | 8 | 23.8 | - | 67.6 | - | - | - | - | - | 16.5 | 53.7 |
| TACL | TSVT'22 | 3D-ResNet18 | 13.5 | V | 16 | - | 35.6 | 50.9 | 56.1 | 65.8 | 34.6 | 37.2 | 39.5 | - | - |
| TACL | TSVT'22 | 3D-ResNet18 | 13.5 | V | 16 | - | 43.7 | 55.6 | 59.2 | 67.2 | 38.7 | 40.2 | 41.7 | - | - |
| 3DRotNet | Arxiv'19 | 3D-ResNet18 | 13.5 | V | 16 | 15 | 31.5 | 40.4 | 47.1 | - | - | - | - | - | - |
| MemDPC | ECCV'20 | 3D-ResNet18 | 13.5 | V | 16 | - | - | 44.2 | 50.9 | 62.3 | - | - | - | - | - |
| MotionFit | ICCV'21 | 3D-ResNet18 | 13.5 | VF | 16 | - | - | - | 57.7 | 59 | - | - | - | - | - |
| TCLR | CVIU'22 | 3D-ResNet18 | 13.5 | V | 16 | 26.9 | - | 66.1 | 73.4 | 76.7 | - | - | - | - | - |
| TimeBalance | CVPR'23 | 3D-ResNet18 | 13.5 | V | 8 | 29.1 | 47.9 | 69.8 | 79.1 | 83.3 | 49.8 | 51.4 | 53.1 | 17.1 | 54.9 |
| ActorCM | Arxiv'21 | R(2+1)D-34 | 33.3 | V | 8 | - | 27 | 40.2 | 51.7 | 59.9 | 32.9 | 38.2 | 38.9 | - | - |
| ActorCM | Arxiv'21 | R(2+1)D-34 | 33.3 | V | 8 | - | 45.1 | 53 | 57.4 | 64.7 | 35.7 | 39.5 | 40.8 | - | - |
| FixMatch | NeuRIPS'20 | SlowFast-R50 | 60 | V | 8 | 16.1 | - | 55.1 | - | - | - | - | - | 10.1 | 49.4 |
| MvPL | ICCV'21 | 3D-ResNet50 | 31.8 | VFG | 8 | 22.8 | - | 80.5 | - | - | - | - | - | 17 | 58.2 |
| CMPL | CVPR'22 | 3D-ResNet50 | 31.8 | V | 8 | 25.1 | - | 79.1 | - | - | - | - | - | 17.6 | 58.4 |
| TimeBalance | CVPR'23 | 3D-ResNet50 | 31.8 | V | 8 | 30.1 | 53.5 | 81.1 | 83.3 | 85 | 52.6 | 53.9 | 54.5 | 19.6 | 61.2 |

# Ablations

| | $f_S$ (rand. init.) | $f_I$ | $f_D$ | Teacher Reweighting | UCF101 % Labels 5% | 20% |
|---|---|---|---|---|---|---|
| (a) | ✓ | ✗ | ✗ | ✗ | 25.60 | 46.20 |
| (b) | ✓ | ✓ | ✗ | ✗ | 43.94 | 74.85 |
| (c) | ✓ | ✗ | ✓ | ✗ | 44.30 | 75.22 |
| (d) | ✓ | ✓ | ✓ | ✗ | 49.57 | 80.06 |
| (e) | ✓ | ✓ | ✓ | ✓ | 53.10 | 83.00 |

# Ablation: Different Teacher Combinations

| | Teacher-1 | Teacher-2 | UCF101 % labels 5% | 20% |
|---|---|---|---|---|
| (a) | Inv1 | Inv2 | 48.33 | 78.76 |
| (b) | Dist1 | Dist2 | 49.15 | 80.49 |
| (c) | Inv1 | Dist1 | 52.14 | 82.02 |
| (d) | Inv1 | Dist2 | 51.78 | 81.43 |

Different Teacher Combinations

# Different Initializations of student

# Conclusion

- TimeBalance, a teacher-student framework for semi-supervised action recognition

- We utilize the complementary strengths of *temporally-invariant* and *temporally-distinctive* representations to leverage unlabeled videos

- State-of-the-art for semi-supervised action recognition: UCF101, HMDB51, Kinetics400

UCF CENTER FOR RESEARCH IN COMPUTER VISION