ACM multimedia
Chengdu, China OCT 20-24 2021

# DSANet: Dynamic Segment Aggregation Network for Video-Level Representation Learning

Wenhao Wu[1*†], Yuxiang Zhao[1,2*], Yanwu Xu[3], Xiao Tan[1], Dongliang He[1],
Zhikang Zou[1], Jin Ye[1], Yingying Li[1], Mingde Yao[1], Zichao Dong[1], Yifeng Shi[1]

[1] Baidu Inc.  [2] Shenzhen Institute of Advanced Technology, Chinese Academy of Sciences  [3] University of Pittsburgh
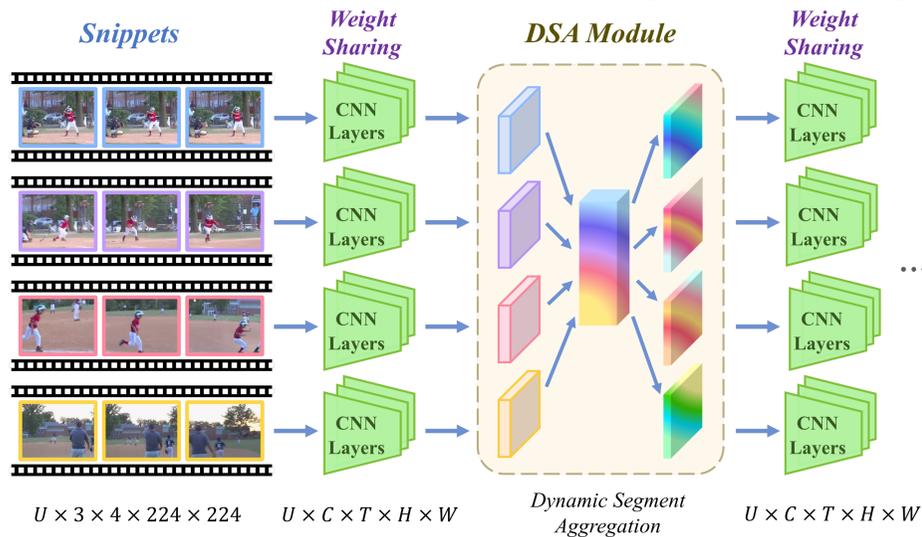
## Motivation

Previous pipeline for video-level representation :
1. (Training) Intra-clip modeling (*e.g.*, C3D/TSM/SlowFast, *etc*)
2. (Inference) Average the predictions of multiple clips

*Can we learn **video-level** representations **directly**?*

**Intra**-clip modeling (3D: T*H*W) → **Inter**-clip modeling (4D: U*T*H*W)
➤ 4D Convolution: effective but expensive
➤ TSN: temporal modeling unexplored but simple
➤ We focus on **efficient** and **effective** video-level representation learning



$U \times 3 \times 4 \times 224 \times 224$   $U \times C \times T \times H \times W$   *Dynamic Segment Aggregation*   $U \times C \times T \times H \times W$
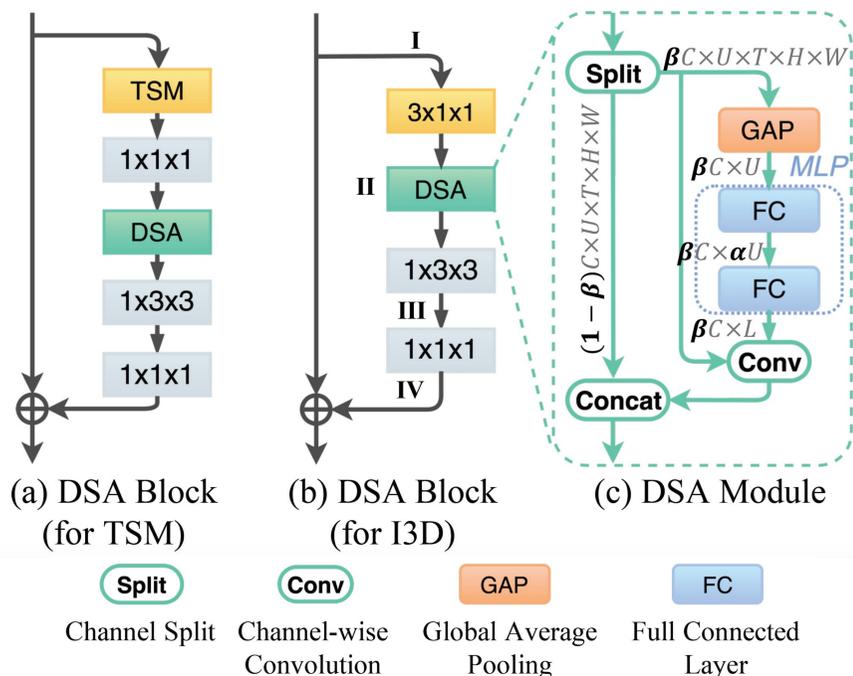
## Contribution

➤ Instead of snippet-level temporal modeling, we propose to exploit an effective and efficient video-level framework for learning video representation. To tackle this, the proposed **DSA module** provides a novel mechanism to adaptively aggregate snippet-level features.

➤ The DSA module works in a **plug-and-play** way and can be easily integrated into existing snippet-level methods. Without any bells and whistles, the DSA module brings consistent improvements when combined with both 2D CNN-based and 3D CNN-based networks (*e.g.*, TSM, I3D, *etc*).

➤ Extensive experiments on four public benchmark datasets demonstrate that the proposed DSA obtain an evident improvement over previous long-range modeling methods with a minimal computational cost.

## Method

### DSA: Dynamic Segment Aggregation



(a) DSA Block (for TSM)   (b) DSA Block (for I3D)   (c) DSA Module

Split — Channel Split   Conv — Channel-wise Convolution   GAP — Global Average Pooling   FC — Full Connected Layer

## Ablation Studies

### Ablation studies on Mini-Kinetics-200.

(a) Study on the effectiveness of DSA module. $T$ denotes the number of frames sampled from each video snippet, $U$ denotes the number of snippets. Backbone: I3D R18.

| Model | $T_{train} \times U_{train}$ | $T_{infer} \times U_{infer} \times$ #crop | Top-1 | Top-5 | Params |
|---|---|---|---|---|---|
| I3D R18 | $4 \times 1$ | $4 \times 10 \times 3$ | 72.2 | 91.2 | 32.3M |
| I3D R18 | $16 \times 1$ | $16 \times 10 \times 3$ | 73.4 | 91.1 | 32.3M |
| TSN+I3D R18 | $4 \times 4$ | $4 \times 10 \times 3$ | 73.0 | 91.3 | 32.3M |
| V4D+I3D R18 | $4 \times 4$ | $4 \times 10 \times 3$ | 75.6 | 92.7 | 33.1M |
| DSA+I3D R18 | $4 \times 4$ | $4 \times 8 \times 3$ | **77.3** | **93.9** | 32.3M |

(b) Study on different position to insert DSA module. Setting: I3D R50, $\alpha=2$, $\beta=1$, stage: res$_5$.

| Position | Top-1 | Top-5 |
|---|---|---|
| I | 81.4 | **95.4** |
| II | **81.5** | 95.2 |
| III | 80.8 | 95.2 |
| IV | 81.4 | 95.1 |

(c) Parameter choices of $\alpha$. Setting: I3D R50, Position II, $\beta=1$, inserted stage: res$_5$.

| Setting | Top-1 | Top-5 |
|---|---|---|
| $\alpha=1$ | 81.0 | 95.1 |
| $\alpha=2$ | **81.5** | **95.2** |
| $\alpha=4$ | 81.2 | 95.0 |
| $\alpha=8$ | 81.3 | 95.0 |

(d) The DSA blocks in different stage of I3D R50. Setting: Position II, $\alpha=2$, $\beta=1$.

| Stage | Top-1 | Top-5 |
|---|---|---|
| res{2} | 81.4 | 94.7 |
| res{3} | 81.3 | 95.1 |
| res{4} | 81.3 | **95.3** |
| res{5} | **81.5** | 95.2 |

(e) Parameter choices of $\beta$. Setting: I3D R50, Position II, $\alpha=2$, inserted stage: res$_5$.

| Setting | Top-1 | Top-5 |
|---|---|---|
| $\beta=1$ | 81.5 | 95.2 |
| $\beta=1/2$ | **81.7** | **95.4** |
| $\beta=1/4$ | 81.6 | 95.0 |
| $\beta=1/8$ | 81.5 | 95.0 |

(f) The number of DSA block inserted into I3D R50. Setting: Position II, $\alpha=2$, $\beta=1/8$.

| Stages | Blocks | Top-1 | Top-5 |
|---|---|---|---|
| res{5} | 1 | 81.5 | 95.2 |
| res{4,5} | 4 | 81.5 | 95.3 |
| res{3,4} | 5 | **81.8** | **95.4** |
| res{2,3} | 3 | 81.4 | 95.1 |

(g) Different short-term temporal structure for DSA module.

| Model | Top-1 | Top-5 |
|---|---|---|
| TSM R50 | 77.4 | 93.4 |
| DSA+TSM R50 | 80.4 | 95.0 |
| I3D R50 | 78.0 | 93.9 |
| DSA+I3D R50 | **81.8** | **95.4** |

(h) Study on the effectiveness of DSA module with different backbones (I3D R18, I3D R50). SENet+I3D uses SE module to replace the DSA module in DSANet.

| Arch. | I3D | SENet+I3D | DSA+I3D |
|---|---|---|---|
| ResNet18 | 72.2 | 73.8 | **77.3** |
| ResNet50 | 78.0 | 78.5 | **81.8** |

(i) Training FLOPs. Comparison with V4D, the extra computation cost brought by the DSA module is close to zero.

| Model | Input size | FLOPs |
|---|---|---|
| TSN+I3D R50 | $4 \times 4 \times 224^2 \times 3$ | 83.8G |
| V4D+I3D R50 | $4 \times 4 \times 224^2 \times 3$ | 143.0G |
| DSA+I3D R50 | $4 \times 4 \times 224^2 \times 3$ | 83.8G |

## Comparison with SOTAs

### Results on **Mini-Kinetics-200** dataset

| Method | Backbone | $T_{train} \times U_{train}$ | $T_{infer} \times U_{infer} \times$ #crop | Top-1 | Top-5 |
|---|---|---|---|---|---|
| S3D [37] | S3D Inception | $64 \times 1$ | N/A | 78.9% | - |
| I3D [38] | 3D ResNet50 | $32 \times 1$ | $32 \times 10 \times 3$ | 75.5% | 92.2% |
| I3D [38] | 3D ResNet101 | $32 \times 1$ | $32 \times 10 \times 3$ | 77.4% | 93.2% |
| I3D+NL [38] | 3D ResNet50 | $32 \times 1$ | $32 \times 10 \times 3$ | 77.5% | 94.0% |
| I3D+CGNL [38] | 3D ResNet50 | $32 \times 1$ | $32 \times 10 \times 3$ | 78.8% | 94.4% |
| I3D+NL [38] | 3D ResNet101 | $32 \times 1$ | $32 \times 10 \times 3$ | 79.2% | 93.2% |
| I3D+CGNL [38] | 3D ResNet101 | $32 \times 1$ | $32 \times 10 \times 3$ | 79.9% | 93.4% |
| V4D+I3D [39] | 3D ResNet18 | $4 \times 4$ | $4 \times 10 \times 3$ | 75.6% | 92.7% |
| V4D+I3D [39] | 3D ResNet50 | $4 \times 4$ | $4 \times 10 \times 3$ | 80.7% | 95.3% |
| DSA+I3D (Ours) | 3D ResNet18 | $4 \times 4$ | $4 \times 8 \times 3$ | **77.3%** | **93.9%** |
| DSA+I3D (Ours) | 3D ResNet50 | $4 \times 4$ | $4 \times 8 \times 3$ | **81.8%** | **95.4%** |

### Results on **Kinetics-400** dataset.

| Method | Backbone | $T_{infer} \times U_{infer} \times$ #crop | GFLOPs | Top-1 | Top-5 |
|---|---|---|---|---|---|
| TSM [17] | ResNet-50 | $8 \times 10 \times 3$ | $33 \times 30 = 990$ | 74.1% | 91.2% |
| TEINet [19] | ResNet-50 | $8 \times 10 \times 3$ | $33 \times 30 = 990$ | 74.9% | 91.8% |
| TEA [16] | ResNet-50 | $8 \times 10 \times 3$ | $35 \times 30 = 1050$ | 75.0% | 91.8% |
| TANet [20] | ResNet-50 | $8 \times 10 \times 3$ | $43 \times 30 = 1290$ | 76.1% | 92.3% |
| MVFNet [34] | ResNet-50 | $8 \times 10 \times 3$ | $33 \times 30 = 990$ | 76.0% | 92.4% |
| NL+I3D [31] | 3D ResNet-50 | $32 \times 10 \times 3$ | $70.5 \times 30 = 2115$ | 74.9% | 91.6% |
| NL+I3D [31] | 3D ResNet-50 | $128 \times 10 \times 3$ | $282 \times 30 = 8460$ | 76.5% | 92.6% |
| X3D-L [7] | - | $16 \times 10 \times 3$ | $24.8 \times 30 = 744$ | 77.5% | 92.9% |
| Slowfast [8] | 3D R50+3D R50 | $(4+32) \times 10 \times 3$ | $36.1 \times 30 = 1083$ | 75.6% | 92.1% |
| Slowfast [8] | 3D R50+3D R50 | $(8+32) \times 3 \times 10$ | $65.7 \times 30 = 1971$ | 77.0% | 92.6% |
| Slowfast [8] | 3D R101+3D R101 | $(8+32) \times 3 \times 10$ | $106 \times 30 = 3180$ | 77.9% | 93.2% |
| Slowonly [8] | 3D ResNet-50 | $8 \times 10 \times 3$ | $41.9 \times 30 = 1257$ | 74.9% | 91.5% |
| V4D+I3D [39] | 3D ResNet-50 | $8 \times 10 \times 3$ | $286.1 \times 2.5 \times 3 = 2146^*$ | 77.4% | 93.1% |
| DSA+I3D (Ours) | 3D ResNet-50 | $4 \times 8 \times 3$ | $83.8 \times 2 \times 3 = 503$ | 77.7% | 93.1% |
| DSA+I3D (Ours) | 3D ResNet-50 | $8 \times 8 \times 3$ | $167.7 \times 2 \times 3 = 1006$ | 78.2% | 93.2% |
| DSA+I3D (Ours) | 3D ResNet-50 | $(4+8) \times 8 \times 3$ | $251.5 \times 2 \times 3 = 1509$ | 79.0% | 93.7% |

Results on **ActivityNet 1.3** dataset.

| Model | Backbone | mAP |
|---|---|---|
| TSN [30] | BN-Inception | 79.7% |
| TSN [30] | Inception V3 | 83.3% |
| TSN-Top3 [30] | Inception V3 | 84.5% |
| V4D+I3D [39] | 3D ResNet50 | 88.9% |
| DSA+I3D (Ours) | 3D ResNet50 | **90.5%** |

Results on **Something-Something V1** dataset.

| Method | Backbone | Top-1 |
|---|---|---|
| MultiScale TRN [40] | BN-Inception | 34.4% |
| ECO [41] | BN-Inception+3D ResNet 18 | 46.4% |
| S3D-G [37] | S3D Inception | 45.8% |
| Nonlocal+GCN [32] | 3D ResNet50 | 46.1% |
| TSM [17] | ResNet50 | 47.2% |
| I3D (our impl.) | 3D ResNet50 | 48.7% |
| V4D+I3D [39] | 3D ResNet50 | 50.4% |
| DSA+I3D (Ours) | 3D ResNet50 | **51.8%** |

## Accuracy-FLOPs Curve



Accuracy-computation trade-off on Kinetics-400 for different methods in the **inference phase**

## Contact

■ Wenhao Wu
Baidu Inc.
wuwenhao17@mails.ucas.edu.cn

■ Yuxiang Zhao
SIAT, CAS
zhaoyuxiang19@mails.ucas.edu.cn

## Codes

https://github.com/whwu95/DSANet