

MVFNet: Multi-View Fusion Network for Efficient Video Recognition

Wenhao Wu¹, Dongliang He¹, Tianwei Lin¹, Fu Li¹, Chuang Gan², Errui Ding¹

¹Department of Computer Vision Technology (VIS), Baidu Inc. ²MIT-IBM Watson AI Lab



AAAI 2021



Task

Video Recognition: classify the short clip or untrimmed video into pre-defined class.





Task

Video Recognition: classify the short clip or untrimmed video into pre-defined class.



- More than simply recognizing objects
- Complex person-person interaction & people-object interactions
- Videos bring motions



Key Observations

- Efficient spatial-temporal modeling is the key to action recognition
- Classical C2D :
 temporal modeling unexplored but simple
- 3D CNN, e.g., SlowFast or SlowOnly: effective but expensive
- TSM enables C2D to model temporal relationship at nearly zero cost
 - *fixed* channel-wise 3x1x1 conv
 - kernel of [0,0,1] for forward shift and [1,0,0] for backward shift





Is TSM our ultimate choice?

NO! We CAN have better choice:

- From the **regular viewpoint** of HW-T: TSM can be improved to have arbitrary learnable *"shift"* kernels
- Why not model **relationships from other viewpoints** of WT-H and TH-W?
- With careful designing, better **effectiveness**efficiency trade-off is possible





Key Innovation





Why MVFNet will work

MVFNet is a **generalized** architecture of existing frameworks

- $\alpha = 0$, MVFNet specializes to be C2D
- $\alpha = 1, \beta_H = \beta_W = 0$, MVFNet is a channel-wise 3x1x1 Conv version of SlowOnly/C3D
- $\alpha = 1/4$, $\beta_H = \beta_W = 0$, and half of channel-wise 3x1x1 conv kernels are [0,0,1] and the rest kernels are [1,0,0], then MVFNet becomes TSM





Design choice of α : *MVFBlock is inserted into res_4 and res_5*

| Setting | Sth-sth v1 | | | | Kinetics-400 | | | | |
|----------------|------------|-------|-------|--------|--------------|-------|-------|--------|--|
| Setting | #F | Top-1 | Top-5 | FLOPs | #F | Top-1 | Top-5 | FLOPs | |
| $\alpha = 0$ | 8 | 17.12 | 43.46 | 32.88G | 4 | 71.87 | 90.02 | 16.44G | |
| $\alpha = 1/8$ | 8 | 49.74 | 78.09 | 32.90G | 4 | 74.21 | 91.34 | 16.45G | |
| $\alpha = 1/4$ | 8 | 49.24 | 77.91 | 32.92G | 4 | 74.18 | 91.46 | 16.46G | |
| $\alpha = 1/2$ | 8 | 50.48 | 79.14 | 32.96G | 4 | 74.21 | 91.42 | 16.48G | |
| $\alpha = 1$ | 8 | 49.73 | 77.94 | 33.04G | 4 | 73.75 | 91.40 | 16.52G | |

(a) Parameter choices of α . Backbone: R-50.



Design choice of how many and where *MVFBlocks are inserted:* $\alpha = 1/2$ and 1/8 for Sth-v1 and K400 , respectively

| Stores | Placks | Sth-sth v1, $\alpha = 1/2$ | | | Kinetics-400, <i>α</i> =1/8 | | | | |
|--------------|--------|----------------------------|-------|-------|-----------------------------|----|-------|-------|--------------|
| Stages | DIOCKS | #F | Top-1 | Top-5 | FLOPs | #F | Top-1 | Top-5 | FLOPs |
| None | 0 | 8 | 17.12 | 43.46 | 32.88G | 4 | 71.87 | 90.02 | 16.44G |
| $res{5}$ | 3 | 8 | 46.02 | 75.60 | 32.90G | 4 | 73.46 | 91.09 | 16.44G |
| res{4,5} | 9 | 8 | 50.48 | 79.14 | 32.96G | 4 | 74.21 | 91.34 | 16.45G |
| res{3,4,5} | 13 | 8 | 49.72 | 78.82 | 33.04G | 4 | 74.08 | 91.51 | 16.46G |
| res{2,3,4,5} | 16 | 8 | 49.95 | 77.96 | 33.12G | 4 | 74.22 | 91.56 | 16.47G |

(b) The number of MVF Blocks inserted into R-50.



Design choice of fusing multiple viewpoints: $\alpha = 1/2$ and 1/8 for Sth-v1 and K400, respectively; MVFblocks in res_4, res_5

| Views | S | th v1 | K400 | | |
|-----------|----|-------|------|-------|--|
| VIEWS | #F | Top-1 | #F | Top-1 | |
| Т | 8 | 49.13 | 4 | 73.72 | |
| T-H | 8 | 49.22 | 4 | 74.01 | |
| T-W | 8 | 49.31 | 4 | 73.88 | |
| T-H-W | 8 | 50.48 | 4 | 74.21 | |
| T-H-W (S) | 8 | 47.21 | 4 | 73.81 | |

(c) Study on the different viewsof MVF module. Backbone: R-50. S denotes weight sharing.

Fusing multi-view information is beneficial

Channel-wise 3x1x1 temporal / horizontal / vertical convolution must have independent kernels



Impact of MVFBlocks when different backbones are used: $\alpha = 1/2$ and 1/8 for Sth-v1 and K400, respectively; MVFblocks in res_4, res_5

| | #F | Top-1 | FLOPs |
|-------|----|-------|--------------|
| | 4 | 74.21 | 16.45G |
| R-50 | 8 | 75.99 | 32.90G |
| | 16 | 77.04 | 65.81G |
| | 4 | 75.98 | 31.36G |
| R-101 | 8 | 77.46 | 62.72G |
| | 16 | 78.42 | 125.45G |

| | Model | Top-1 | FLOPs |
|----------|-------|-------|--------------|
| ML VO | C2D | 64.4 | 1.25G |
| WID- V 2 | MVF | 67.5 | 1.25G |
| D 50 | C2D | 71.9 | 16.44G |
| K-30 | MVF | 74.2 | 16.48G |
| | | | |

(e) Advanced backbones for MVFNet on Kinetics-400. (f) **Different backbones for MVFNet on Kinetics-400**. Mb-V2 denotes MobileNet-V2.



Comparison with Similar Variants

 $\alpha = 1/2$ and 1/8 for Sth-v1 and K400, respectively; MVFblocks in res_4, res_5

| Method | Sth v1 | K400 | FI ODe | Params | |
|---------------|--------|-------|--------|--------|--|
| Methou | Top-1 | Top-1 | TLOFS | | |
| C2D | 17.1 | 71.4 | 32.9G | 24.3M | |
| TSM | 47.2 | 74.1 | 32.9G | 24.3M | |
| SlowOnly | | 74.9 | 41.9G | 32.4M | |
| CoST* | - | - | 45.8G | 24.3M | |
| MVFNet | 50.5 | 76.0 | 32.9G | 24.3M | |

(d) Study on the effectiveness of **MVFNet**. Backbone: R-50, 8f input. * indicates our implementation.



Comparison with Other Solutions on Kinetics400

| | Method | Backbone | Frames × Crops × Clips | GFLOPs | Top-1 | Top-5 |
|-----|---|------------------|---|-------------------|-------|---------------|
| | I3D (Carreira et al. 2017) | Inception V1 | $64 \times N/A \times N/A$ | 108×N/A | 72.1% | 90.3% |
| | S3D-G (Xie et al. 2018) | Inception V1 | 64×3×10 | 71.4×30 | 74.7% | 93.4% |
| | TSN (Wang et al. 2016) | Inception V3 | $25 \times 10 \times 1$ | 80×10 | 72.5% | 90.2% |
| | $ECO-RGB_{En}$ (Zolfaghari et al. 2018) | BNIncep+Res3D-18 | $92 \times 1 \times 1$ | 267×1 | 70.0% | -% |
| | R(2+1)D (Tran et al. 2018) | ResNet-34 | $32 \times 1 \times 10$ | 152×10 | 74.3% | 91.4% |
| | X3D-M (Feichtenhofer 2020) | - | $16 \times 3 \times 10$ | 6.2×30 | 76.0% | 92.3% |
| | STM (Jiang et al. 2019) | ResNet-50 | 16×3×10 | 67×30 | 73.7% | 91.6% |
| | TSM (Lin, Gan, and Han 2019) | ResNet-50 | $8 \times 3 \times 10$ | 33×30 | 74.1% | 91.2% |
| | SlowOnly (Feichtenhofer et al. 2019) | ResNet-50 | $8 \times 3 \times 10$ | 41.9×30 | 74.9% | 91.5% |
| | TEINet (Liu et al. 2020) | ResNet-50 | 8×3×10 | 33×30 | 74.9% | 91.8% |
| | TEA (Li et al. 2020b) | ResNet-50 | $8 \times 3 \times 10$ | 33×30 | 75.0% | 91.8% |
| | Slowfast (Feichtenhofer et al. 2019) | R50+R50 | (4+32)×3×10 | 36.1×30 | 75.6% | 92.1% |
| | NL+I3D (Wang et al. 2018b) | ResNet-50 | $32 \times 3 \times 10$ | 70.5×30 | 74.9% | 91.6% |
| | NL+I3D (Wang et al. 2018b) | ResNet-50 | $128 \times 3 \times 10$ | 282×30 | 76.5% | 92.6% |
| 8 | MVFNet | ResNet-50 | 8×3×10 | 32.9×30 | 76.0% | 92.4 % |
| | MVFNet | ResNet-50 | $16 \times 3 \times 10$ | 65.8×30 | 77.0% | 92.8% |
| 00- | ip-CSN (Tran et al. 2019) | ResNet-101 | 32×3×10 | 82×30 | 76.7% | 92.3% |
| | SmallBig (Li et al. 2020a) | ResNet-101 | $32 \times 3 \times 4$ | 418×12 | 77.4% | 93.3% |
| | SlowOnly (Feichtenhofer et al. 2019) | ResNet-101 | $16 \times 3 \times 10$ | 185×30 | 77.2% | -% |
| | NL+I3D (Wang et al. 2018b) | ResNet-101 | $128 \times 3 \times 10$ | 359×30 | 77.7% | 93.3% |
| | Slowfast (Feichtenhofer et al. 2019) | R101+R101 | (8+32)×3×10 | 106×30 | 77.9% | 93.2% |
| | Slowfast (Feichtenhofer et al. 2019) | R101+R101 | (16+64)×3×10 | 213×30 | 78.9% | 93.5% |
| | TPN (Yang et al. 2020) | ResNet-101 | $32 \times 3 \times 10$ | 374×30 | 78.9% | 93.9% |
| | MVFNet | ResNet-101 | 8×3×10 | 62.7×30 | 77.4% | 92.9% |
| | MVFNet | ResNet-101 | $16 \times 3 \times 10$ | 125.4×30 | 78.4% | 93.4% |
| | $MVFNet_{En}$ | R101+R101 | $(16+8) \times 3 \times 10$ | 188.1×30 | 79.1% | 93.8% |



Comparison with Other Solutions on Sth-Sth-v1/v2

| Method | Backbone | Frames×Crons×Clins | FLOPs | Pre-train | V1 Val | V2 Val |
|-------------------------------------|-------------------|-------------------------|--------------------------|--------------|-----------|------------------|
| Wiethou | Duckbone | Traines ~ erops ~ enps | TLOIS | TTe-train | Top-1 (%) | Top-1 (%) |
| I3D (Wang et al. 2018) | 3D ResNet50 | | $153G \times 3 \times 2$ | ImageNet | 41.6 | - |
| NL I3D (Wang et al. 2018) | 3D ResNet50 | $32 \times 3 \times 2$ | 168G×3×2 | + | 44.4 | - |
| NL I3D+GCN (Wang et al. 2018) | 3D ResNet50+GCN | | $303G \times 3 \times 2$ | K400 | 46.1 | - |
| ECO (Zolfaghari et al. 2018) | DNIncon 12D Dec19 | 8×1×1 | $32G \times 1 \times 1$ | V 400 | 39.6 | - |
| ECO_{En} (Zolfaghari et al. 2018) | BNIICep+5D Kes18 | 92×1×1 | 267G×1×1 | K 400 | 46.4 | - |
| S3D-G (Xie et al. 2018) | Inception | 64×1×1 | $71G \times 1 \times 1$ | K400 | 48.2 | - |
| TSN (Wang et al. 2016) | ResNet50 | $8 \times 3 \times 2$ | $33G \times 3 \times 2$ | ImageNet | 20.5 | 30.4 |
| TEM (Lin et al. 2010) | DecNet50 | 8×3×2 | $33G \times 3 \times 2$ | ImagaNat | 47.2 | 61.2 |
| 15M (Lin et al. 2019) | Resinet50 | $16 \times 3 \times 2$ | $65G \times 3 \times 2$ | Imagenet | 48.4 | 63.1 |
| STM (liong at al. 2010) | DecNet50 | 8×3×10 | 33G×3×10 | ImagaNat | 49.2 | 62.3 |
| STM (Jiang et al. 2019) | Resiletou | $16 \times 3 \times 10$ | 67G×3×10 | Imagenet | 50.7 | 64.3 |
| TEINet (Lin et al. 2020) | DecNet50 | 8×3×10 | $33G \times 3 \times 10$ | ImagaNat | 48.8 | 64.0 |
| TEINet (Liu et al. 2020) | Resiletou | $16 \times 3 \times 10$ | 66G×3×10 | ImageNet | 51.0 | 64.7 |
| TEA (Li at al. 2020b) | PosNot50 | 8×3×10 | 35G×3×10 | ImagaNat | 51.7 | - |
| TEA (LI et al. 20200) | Resiletou | $16 \times 3 \times 10$ | 70G×3×10 | Imagenet | 52.3 | - |
| | | 8×1×1 | 33G×1×1 | | 48.8 | 60.8 |
| | | $8 \times 3 \times 2$ | $33G \times 3 \times 2$ | | 50.5 | 63.5 |
| MVFNet | ResNet50 | $16 \times 1 \times 1$ | 66G×1×1 | ImageNet | 51.0 | 62.9 |
| | | $16 \times 3 \times 2$ | 66G×3×2 | - | 52.6 | 65.2 |
| | | (16+8)×3×2 | 99G×3×2 | | 54.0 | 66.3 |



Transfer Learning on UCF101/HMDB51

Mean class accuracy of RGB modality is reported, RGB models are pretrained on Kinetics400

| Method | Backbone | UCF-101 | HMDB-51 |
|------------|------------------|----------------|---------|
| ECO_{En} | BNIncep+Res3D-18 | 94.8% | 72.4% |
| ARTNet | ResNet-18 | 94.3% | 70.9% |
| I3D | Inception V1 | 95.6% | 74.8% |
| R(2+1)D | Inception V1 | 96.8% | 74.5% |
| S3D-G | Inception V1 | 96.8% | 75.9% |
| TSN | BNInception | 91.1% | - |
| StNet | ResNet-50 | 93.5% | - |
| TSM | ResNet-50 | 95.9% | 73.5% |
| STM | ResNet-50 | 96.2% | 72.2% |
| TEINet | ResNet-50 | 96.7% | 72.1% |
| MVFNet | ResNet-50 | 96.6% | 75.7% |



Conclusion

- Upgrading fixed shift kernels of TSM to be learnable is more flexible
- Relationship modeling from multiple viewpoints is a strong boost
- MVFNet consistently outperforms existing solutions on Kinetics400, Something-Something-v1/v2
- Codes & models will be available

https://github.com/whwu95/MVFNet



Thank you!

MVFNet: Multi-View Fusion Network for Efficient Video Recognition

Contact: Wenhao Wu <u>wuwenhao17@mails.ucas.edu.cn</u> <u>https://github.com/whwu95/MVFNet</u>

