

# Video Object Segmentation-aware Video Frame Interpolation

## -Supplementary Material-

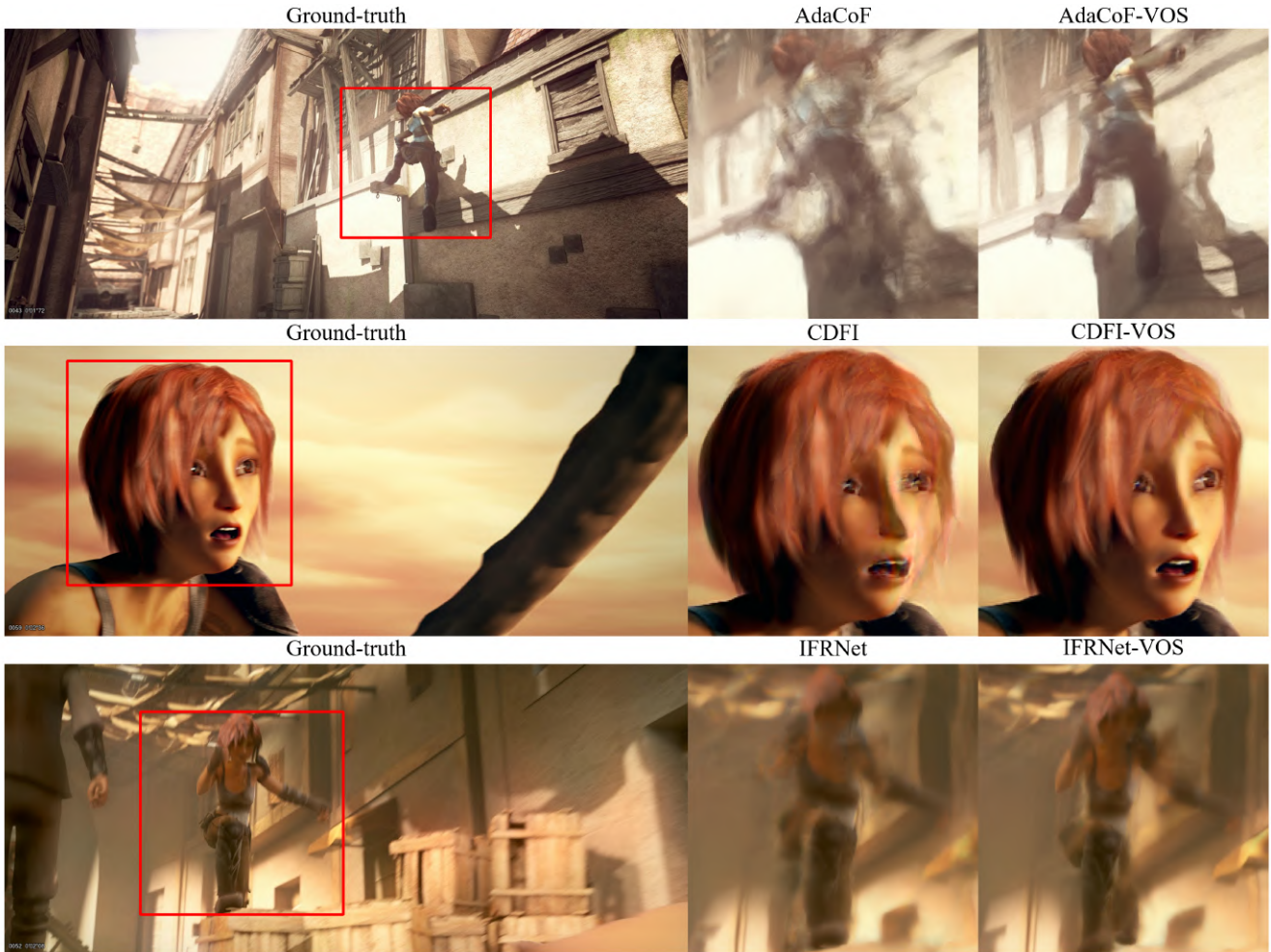


Figure S-1. Examples of  $4\times$  VFI results obtained by the baseline models, AdaCoF [5], CDFI [1], and IFRNet [4], trained with and without the proposed VOS-VFI.

### A. VFI with a scale factor of 4.

The proposed VOS-VFI can be applied to any VFI model. Consequently, arbitrary-rate VFI is feasible if VOS-VFI is applied to VFI models supporting such a capability. However, among the three representative baseline models chosen in our experiments, AdaCoF [5] and CDFI [1] do not support arbitrary-rate VFI. In addition, IFRNet [4] requires

dedicated weight parameters for multi-frame interpolation. Therefore, to evaluate the performance of VOS-VFI on another scale factor in addition to the standard scale factor 2, we performed  $2\times$  VFI twice to have  $4\times$  VFI results, which is a common strategy [3, 7]. Table S-1 shows the quantitative performance evaluation results on the DAVIS 2016 and 2017 datasets. The results showed a similar tendency as those obtained for the scale factor of 2. Specifically, the

| Model      | DAVIS 2016      |                            |               |             |             | DAVIS 2017      |                            |               |             |             |
|------------|-----------------|----------------------------|---------------|-------------|-------------|-----------------|----------------------------|---------------|-------------|-------------|
|            | PSNR $\uparrow$ | $\uparrow$ PSNR $\uparrow$ | $J&F\uparrow$ | $J\uparrow$ | $F\uparrow$ | PSNR $\uparrow$ | $\uparrow$ PSNR $\uparrow$ | $J&F\uparrow$ | $J\uparrow$ | $F\uparrow$ |
| AdaCoF [5] | <b>21.03</b>    | 26.72                      | 73.2          | 73.1        | 73.2        | <b>22.15</b>    | 26.63                      | 69.6          | 67.3        | 71.8        |
| AdaCoF-VOS | 20.96           | <b>26.81</b>               | <b>77.8</b>   | <b>77.6</b> | <b>78.0</b> | 22.13           | <b>26.71</b>               | <b>73.6</b>   | <b>71.0</b> | <b>76.2</b> |
| CDFI [1]   | 21.01           | 26.77                      | 78.6          | 77.8        | 79.4        | 21.98           | 26.69                      | 75.2          | 72.1        | 78.4        |
| CDFI-VOS   | <b>21.05</b>    | <b>26.88</b>               | <b>80.0</b>   | <b>79.2</b> | <b>80.8</b> | <b>22.03</b>    | <b>26.80</b>               | <b>76.1</b>   | <b>72.9</b> | <b>79.3</b> |
| IFRNet [4] | 21.32           | 26.83                      | 78.7          | 78.6        | 78.7        | 22.33           | 26.76                      | 75.2          | 72.5        | 78.0        |
| IFRNet-VOS | <b>21.36</b>    | <b>26.94</b>               | <b>79.3</b>   | <b>79.1</b> | <b>79.5</b> | <b>22.37</b>    | <b>26.82</b>               | <b>76.0</b>   | <b>73.1</b> | <b>78.9</b> |

Table S-1. Quantitative results of the three baseline models trained with and without VOS-VFI for  $4\times$  interpolation. The performance is evaluated in terms of the image quality PSNR and segmentation accuracy ( $J&F$ ,  $J$ , and  $F$ ) on the DAVIS 2016 and 2017 datasets.  $\uparrow$  represents PSNR scores on the foreground objects obtained by masking out the background using the ground-truth segmentation maps.

proposed VOS-VFI improved the segmentation accuracy of the baseline models by 4.6%, 1.4%, and 0.6% for AdaCoF, CDFI, and IFRNet for DAVIS 2016, respectively, and 4.0%, 0.9%, and 0.8% for AdaCoF, CDFI, and IFRNet for DAVIS 2017, respectively, in terms of  $J&F$ .

Fig. S-1 shows several results obtained for  $4\times$  VFI on the HD dataset. As can be seen, the proposed VOS-VFI contributed to the baseline models by improving the image quality of interpolated frames.

## B. More results

First, we experimented the method that uses the correspondence-wise loss (CoRR) [2] for achieving better visual quality, which can be adopted to any other VFI models as our approach. However, as shown in Table S-2, we could not obtain performance improvements using CoRR.

Next, we applied VOS-VFI to a more recent transformer-based VFI baseline [8]. VOS-VFI also introduced non-marginal improvements, especially for VOS metrics, on this latest baseline. Our source code can be found in our project page<sup>1</sup>. We expect that the proposed VOS-VFI training framework can be applied to upcoming VFI models to boost performance without increasing their number of parameters and inference time.

Lastly, Table S-3 provides the NIQE/PI/NIMA scores obtained from four datasets for clear performance comparisons. As can be seen, the proposed method improved these perceptual metrics on all datasets. Table S-4 provides the PSNR scores on the four datasets, where the foreground PSNRs were only measured for the DAVIS datasets using the ground-truth segmentation maps. Although the proposed VOS-VFI showed some improvements in foreground object synthesis, we could not obtain consistent performance improvements in terms of the PSNR. Note that the PSNR is not correlated with the human perceptual quality of interpolated frames [6], and VOS-VFI showed effectiveness

| Model            | DAVIS 2016   |                 |               |             |             |
|------------------|--------------|-----------------|---------------|-------------|-------------|
|                  | PSNR         | $\uparrow$ PSNR | $J&F\uparrow$ | $J\uparrow$ | $F\uparrow$ |
| AdaCoF           | <b>25.11</b> | 25.62           | 85.9          | 84.9        | 86.8        |
| AdaCoF-VOS       | 25.03        | <b>25.72</b>    | <b>87.0</b>   | <b>85.9</b> | <b>88.2</b> |
| AdaCoF-CoRR [2]  | 24.87        | 25.66           | 85.1          | 84.2        | 86.0        |
| EMA [8]          | <b>27.24</b> | 26.03           | 88.0          | 86.7        | 89.3        |
| EMA-VOS          | 27.19        | <b>26.06</b>    | <b>88.8</b>   | <b>87.6</b> | <b>90.1</b> |
| EMA [8]-CoRR [2] | 27.10        | 25.80           | 86.9          | 85.9        | 88.0        |

Table S-2. Evaluation using the additional methods [2, 8].  $\uparrow$  represents a foreground PSNR.

in the perceptual quality metrics, VOS performance (Table 1), video tracking performance (Table 3), object pose estimation performance (Table 5), and user studies (Section 4.2.1). Check our project page for more results.

<sup>1</sup><https://github.com/junsang7777/VOS-VFI>

| Model      | DAVIS 2016        |                 |                 | DAVIS 2017        |                 |                 | Vimeo90K (val)    |                 |                 | UCF101 (val)      |                 |                 |
|------------|-------------------|-----------------|-----------------|-------------------|-----------------|-----------------|-------------------|-----------------|-----------------|-------------------|-----------------|-----------------|
|            | NIQE $\downarrow$ | PI $\downarrow$ | NIMA $\uparrow$ | NIQE $\downarrow$ | PI $\downarrow$ | NIMA $\uparrow$ | NIQE $\downarrow$ | PI $\downarrow$ | NIMA $\uparrow$ | NIQE $\downarrow$ | PI $\downarrow$ | NIMA $\uparrow$ |
| AdaCoF [5] | 3.443             | 3.338           | 4.565           | 3.545             | 3.402           | 4.502           | 5.180             | 4.104           | 4.765           | 7.272             | 5.695           | 4.018           |
| AdaCoF-VOS | <b>3.431</b>      | <b>3.329</b>    | <b>4.586</b>    | <b>3.534</b>      | <b>3.400</b>    | <b>4.505</b>    | <b>5.153</b>      | <b>4.094</b>    | <b>4.771</b>    | <b>7.237</b>      | <b>5.678</b>    | <b>4.026</b>    |
| CDFI [1]   | 3.081             | 2.845           | 4.568           | 3.267             | 3.015           | 4.485           | 4.933             | 3.832           | 4.873           | 6.878             | 5.421           | 3.987           |
| CDFI-VOS   | <b>3.067</b>      | <b>2.658</b>    | <b>4.671</b>    | <b>3.254</b>      | <b>3.002</b>    | <b>4.492</b>    | <b>4.910</b>      | <b>3.822</b>    | <b>4.879</b>    | <b>6.875</b>      | <b>5.408</b>    | <b>3.991</b>    |
| IFRNet [4] | 3.534             | 3.304           | 4.407           | 3.668             | 3.494           | 4.351           | 5.062             | 3.969           | 4.820           | 7.191             | 5.665           | <b>4.023</b>    |
| IFRNet-VOS | <b>3.519</b>      | <b>3.294</b>    | <b>4.416</b>    | <b>3.651</b>      | <b>3.483</b>    | <b>4.361</b>    | <b>5.021</b>      | <b>3.935</b>    | <b>4.824</b>    | <b>7.115</b>      | <b>5.617</b>    | 4.020           |

Table S-3. Evaluation in terms of the three representative perceptual quality metrics.

| Model      | DAVIS 2016      |                  | DAVIS 2017      |                  | Vimeo90K (val)  | UCF101 (val)    |
|------------|-----------------|------------------|-----------------|------------------|-----------------|-----------------|
|            | PSNR $\uparrow$ | †PSNR $\uparrow$ | PSNR $\uparrow$ | †PSNR $\uparrow$ | PSNR $\uparrow$ | PSNR $\uparrow$ |
| AdaCoF [5] | <b>25.11</b>    | 25.62            | <b>26.23</b>    | 26.13            | <b>34.34</b>    | 35.16           |
| AdaCoF-VOS | 25.03           | <b>25.72</b>     | 26.21           | <b>26.22</b>     | 34.26           | 35.16           |
| CDFI [1]   | 25.68           | 25.74            | 26.71           | 26.24            | 35.17           | 35.21           |
| CDFI-VOS   | <b>25.75</b>    | <b>25.80</b>     | <b>26.79</b>    | <b>26.30</b>     | <b>35.28</b>    | <b>35.25</b>    |
| IFRNet [4] | 26.70           | 25.91            | 27.57           | 26.44            | 35.73           | 35.26           |
| IFRNet-VOS | <b>26.74</b>    | <b>25.98</b>     | <b>27.60</b>    | <b>26.50</b>     | <b>35.80</b>    | <b>35.28</b>    |

Table S-4. Evaluation in terms of the PSNR. † represents a foreground PSNR.

## References

- [1] Tianyu Ding, Luming Liang, Zihui Zhu, and Ilya Zharkov. CDFI: Compression-driven network design for frame interpolation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8001–8011, 2021.
- [2] Daniel Geng, Max Hamilton, and Andrew Owens. Comparing correspondences: Video prediction with correspondence-wise losses. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3365–3376, 2022.
- [3] Zhewei Huang, Tianyuan Zhang, Wen Heng, Boxin Shi, and Shuchang Zhou. Real-time intermediate flow estimation for video frame interpolation. In *Proceedings of the European Conference on Computer Vision*, pages 624–642, 2022.
- [4] Lingtong Kong, Boyuan Jiang, Donghao Luo, Wenqing Chu, Xiaoming Huang, Ying Tai, Chengjie Wang, and Jie Yang. IFRNet: Intermediate feature refine network for efficient frame interpolation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1969–1978, 2022.
- [5] Hyeongmin Lee, Taeh Kim, Tae-young Chung, Daehyun Pak, Yuseok Ban, and Sangyoun Lee. AdaCoF: Adaptive collaboration of flows for video frame interpolation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5316–5325, 2020.
- [6] Hui Men, Hanhe Lin, Vlad Hosu, Daniel Maurer, Andres Bruhn, and Dietmar Saupe. Technical report on visual quality assessment for frame interpolation. *arXiv preprint arXiv:1901.05362*, 2019.
- [7] Junheum Park, Keunsoo Ko, Chul Lee, and Chang-Su Kim. BMBC: Bilateral motion estimation with bilateral cost volume for video interpolation. In *Proceedings of the European Conference on Computer Vision*, pages 109–125, 2020.
- [8] Guozhen Zhang, Yuhan Zhu, Haonan Wang, Youxin Chen, Gangshan Wu, and Limin Wang. Extracting motion and appearance via inter-frame attention for efficient video frame interpolation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5682–5692, 2023.