

# Revisiting Adaptive Convolutions for Video Frame Interpolation

## — Supplementary Material —

Simon Niklaus  
Adobe Research

Long Mai  
Adobe Research

Oliver Wang  
Adobe Research

### 1. Middlebury Benchmark

A screenshot of the interpolation error category from the Middlebury benchmark for optical flow [1] is shown Figure 1, where our SepConv++ currently ranks second among all published methods. Please note that DCM is an unpublished paper that, as per the title, distills a cheating model.

### 2. Structural Similarity

We focus on PSNR in our main paper since SSIM [4] is subject to unexpected and unintuitive results [3]. However, we provide relevant results with SSIM instead of PSNR in Tables 1–3. These results support our claims and are generally aligned with PSNR in terms of relative improvements.

### References

- [1] Simon Baker, Daniel Scharstein, J. P. Lewis, Stefan Roth, Michael J. Black, and Richard Szeliski. A Database and Evaluation Methodology for Optical Flow. *International Journal of Computer Vision*, 92(1):1–31, 2011. 1, 2
- [2] Ziwei Liu, Raymond A. Yeh, Xiaoou Tang, Yiming Liu, and Aseem Agarwala. Video Frame Synthesis Using Deep Voxel Flow. In *IEEE International Conference on Computer Vision*, 2017. 2
- [3] Jim Nilsson and Tomas Akenine-Möller. Understanding SSIM. *arXiv/2006.13846*, 2020. 1
- [4] Zhou Wang, Alan C. Bovik, Hamid R. Sheikh, and Eero P. Simoncelli. Image Quality Assessment: From Error Visibility to Structural Similarity. *IEEE Transactions on Image Processing*, 13(4):600–612, 2004. 1
- [5] Tianfan Xue, Baian Chen, Jiajun Wu, Donglai Wei, and William T. Freeman. Video Enhancement With Task-Oriented Flow. *International Journal of Computer Vision*, 127(8):1106–1125, 2019. 2

Average interpolation error	avg. rank	Mequon (Hidden texture)				Schefflera (Hidden texture)				Urban (Synthetic)				Teddy (Stereo)				Backyard (High-speed camera)				Basketball (High-speed camera)				Dumprtruck (High-speed camera)				Evergreen (High-speed camera)																			
		im0		GT		im1		im0		GT		im1		im0		GT		im1		im0		GT		im1		im0		GT		im1																			
		all	disc	all	untxt	all	disc	all	untxt	all	disc	all	untxt	all	disc	all	untxt	all	disc	all	disc	all	untxt	all	disc	all	disc	all	untxt	all	disc	all	untxt																
SoftSplat [169]	3.4	<u>2.06</u>	<u>3.06</u>	<u>1.14</u>	<u>6</u>	<u>2.80</u>	<u>3.91</u>	<u>1.24</u>	<u>1</u>	<u>1.99</u>	<u>2.73</u>	<u>1.21</u>	<u>4</u>	<u>3.84</u>	<u>4.64</u>	<u>2.69</u>	<u>1</u>	<u>8.10</u>	<u>10.0</u>	<u>13</u>	<u>2.96</u>	<u>1</u>	<u>4.10</u>	<u>7.53</u>	<u>1.98</u>	<u>3</u>	<u>5.49</u>	<u>4</u>	<u>12.14</u>	<u>1.39</u>	<u>4</u>	<u>5.40</u>	<u>2</u>	<u>8.33</u>	<u>2</u>	<u>1.50</u>	<u>4</u>												
DCM [185]	8.5	<u>2.16</u>	<u>3.30</u>	<u>2.16</u>	<u>8</u>	<u>2.82</u>	<u>3.83</u>	<u>1.27</u>	<u>9</u>	<u>1.89</u>	<u>2.58</u>	<u>1.10</u>	<u>1</u>	<u>4.09</u>	<u>4.98</u>	<u>2.81</u>	<u>3</u>	<u>8.59</u>	<u>10.7</u>	<u>16</u>	<u>3.09</u>	<u>13</u>	<u>4.87</u>	<u>17</u>	<u>9.18</u>	<u>17</u>	<u>2.00</u>	<u>6</u>	<u>6.17</u>	<u>14</u>	<u>13.5</u>	<u>14</u>	<u>1.44</u>	<u>12</u>	<u>6.14</u>	<u>9</u>	<u>9.45</u>	<u>9</u>	<u>1.65</u>	<u>22</u>									
SepConv++ [186]	10.7	<u>2.39</u>	<u>18</u>	<u>4.17</u>	<u>20</u>	<u>1.20</u>	<u>21</u>	<u>2.98</u>	<u>5</u>	<u>4.21</u>	<u>6</u>	<u>1.28</u>	<u>16</u>	<u>3.34</u>	<u>20</u>	<u>3.23</u>	<u>5</u>	<u>2.20</u>	<u>83</u>	<u>4.49</u>	<u>10</u>	<u>5.81</u>	<u>13</u>	<u>2.87</u>	<u>5</u>	<u>7.64</u>	<u>4</u>	<u>9.42</u>	<u>2</u>	<u>2.97</u>	<u>2</u>	<u>3.77</u>	<u>2</u>	<u>6.80</u>	<u>2</u>	<u>1.96</u>	<u>1</u>	<u>5.26</u>	<u>3</u>	<u>11.63</u>	<u>1.36</u>	<u>1</u>	<u>5.71</u>	<u>6</u>	<u>8.86</u>	<u>1.45</u>	<u>1</u>		
FGME [158]	10.8	<u>2.08</u>	<u>2</u>	<u>3.34</u>	<u>4</u>	<u>0.98</u>	<u>1</u>	<u>3.32</u>	<u>18</u>	<u>4.43</u>	<u>10</u>	<u>1.63</u>	<u>109</u>	<u>2.46</u>	<u>4</u>	<u>3.28</u>	<u>6</u>	<u>1.41</u>	<u>14</u>	<u>4.08</u>	<u>2</u>	<u>4.85</u>	<u>2</u>	<u>3.05</u>	<u>15</u>	<u>7.36</u>	<u>1</u>	<u>9.06</u>	<u>1</u>	<u>3.03</u>	<u>6</u>	<u>4.17</u>	<u>6</u>	<u>7.62</u>	<u>6</u>	<u>2.06</u>	<u>17</u>	<u>4.95</u>	<u>2</u>	<u>10.72</u>	<u>1.44</u>	<u>12</u>	<u>5.45</u>	<u>3</u>	<u>8.41</u>	<u>3</u>	<u>1.57</u>	<u>13</u>	
BMBC [172]	11.9	<u>2.30</u>	<u>11</u>	<u>3.40</u>	<u>6</u>	<u>1.20</u>	<u>21</u>	<u>3.07</u>	<u>6</u>	<u>4.25</u>	<u>7</u>	<u>1.41</u>	<u>96</u>	<u>3.17</u>	<u>16</u>	<u>4.19</u>	<u>26</u>	<u>1.66</u>	<u>35</u>	<u>4.24</u>	<u>5</u>	<u>5.28</u>	<u>5</u>	<u>3.14</u>	<u>19</u>	<u>7.79</u>	<u>5</u>	<u>9.62</u>	<u>5</u>	<u>3.14</u>	<u>19</u>	<u>4.08</u>	<u>3</u>	<u>7.47</u>	<u>3</u>	<u>2.02</u>	<u>8</u>	<u>5.63</u>	<u>6</u>	<u>12.4</u>	<u>1.40</u>	<u>6</u>	<u>5.54</u>	<u>4</u>	<u>8.58</u>	<u>4</u>	<u>1.61</u>	<u>19</u>	
EAFI [171]	12.0	<u>2.22</u>	<u>5</u>	<u>3.69</u>	<u>9</u>	<u>1.15</u>	<u>7</u>	<u>3.16</u>	<u>9</u>	<u>4.44</u>	<u>11</u>	<u>1.50</u>	<u>63</u>	<u>2.12</u>	<u>3</u>	<u>2.96</u>	<u>3</u>	<u>1.15</u>	<u>2</u>	<u>4.28</u>	<u>7</u>	<u>5.30</u>	<u>6</u>	<u>2.92</u>	<u>6</u>	<u>8.71</u>	<u>19</u>	<u>10.8</u>	<u>17</u>	<u>3.03</u>	<u>6</u>	<u>4.97</u>	<u>19</u>	<u>9.45</u>	<u>19</u>	<u>1.99</u>	<u>4</u>	<u>5.94</u>	<u>12</u>	<u>13.3</u>	<u>13</u>	<u>1.38</u>	<u>3</u>	<u>5.84</u>	<u>7</u>	<u>9.03</u>	<u>7</u>	<u>1.55</u>	<u>11</u>
STAR-Net [164]	14.0	<u>2.18</u>	<u>4</u>	<u>3.37</u>	<u>5</u>	<u>1.21</u>	<u>39</u>	<u>3.46</u>	<u>27</u>	<u>4.88</u>	<u>26</u>	<u>1.47</u>	<u>75</u>	<u>3.04</u>	<u>14</u>	<u>3.53</u>	<u>11</u>	<u>1.58</u>	<u>28</u>	<u>4.41</u>	<u>9</u>	<u>5.44</u>	<u>9</u>	<u>2.76</u>	<u>2</u>	<u>7.51</u>	<u>2</u>	<u>9.27</u>	<u>2</u>	<u>2.98</u>	<u>4</u>	<u>4.65</u>	<u>3</u>	<u>8.72</u>	<u>9</u>	<u>1.99</u>	<u>4</u>	<u>6.21</u>	<u>15</u>	<u>13.4</u>	<u>14</u>	<u>1.41</u>	<u>7</u>	<u>6.17</u>	<u>10</u>	<u>9.45</u>	<u>9</u>	<u>1.49</u>	<u>3</u>
EDSC [174]	15.2	<u>2.32</u>	<u>15</u>	<u>3.90</u>	<u>14</u>	<u>1.16</u>	<u>8</u>	<u>3.10</u>	<u>7</u>	<u>4.38</u>	<u>9</u>	<u>1.51</u>	<u>85</u>	<u>2.98</u>	<u>12</u>	<u>3.54</u>	<u>12</u>	<u>1.36</u>	<u>12</u>	<u>4.49</u>	<u>10</u>	<u>5.74</u>	<u>11</u>	<u>3.16</u>	<u>28</u>	<u>8.05</u>	<u>12</u>	<u>9.96</u>	<u>12</u>	<u>3.08</u>	<u>11</u>	<u>4.89</u>	<u>18</u>	<u>9.28</u>	<u>18</u>	<u>2.02</u>	<u>8</u>	<u>5.55</u>	<u>5</u>	<u>12.3</u>	<u>5</u>	<u>1.41</u>	<u>7</u>	<u>6.42</u>	<u>17</u>	<u>9.99</u>	<u>18</u>	<u>1.55</u>	<u>11</u>
AdaCoF [165]	19.7	<u>2.41</u>	<u>20</u>	<u>4.10</u>	<u>19</u>	<u>1.26</u>	<u>129</u>	<u>3.10</u>	<u>7</u>	<u>4.32</u>	<u>8</u>	<u>1.43</u>	<u>62</u>	<u>3.48</u>	<u>24</u>	<u>3.31</u>	<u>7</u>	<u>1.78</u>	<u>51</u>	<u>4.84</u>	<u>18</u>	<u>5.94</u>	<u>19</u>	<u>2.93</u>	<u>7</u>	<u>8.68</u>	<u>18</u>	<u>10.8</u>	<u>17</u>	<u>3.14</u>	<u>19</u>	<u>4.13</u>	<u>5</u>	<u>7.59</u>	<u>5</u>	<u>1.97</u>	<u>2</u>	<u>5.77</u>	<u>10</u>	<u>12.9</u>	<u>11</u>	<u>1.37</u>	<u>2</u>	<u>5.60</u>	<u>5</u>	<u>8.67</u>	<u>5</u>	<u>1.48</u>	<u>2</u>
DSepConv [162]	23.3	<u>2.47</u>	<u>21</u>	<u>4.39</u>	<u>26</u>	<u>1.21</u>	<u>39</u>	<u>3.32</u>	<u>18</u>	<u>4.60</u>	<u>18</u>	<u>1.72</u>	<u>128</u>	<u>3.28</u>	<u>17</u>	<u>3.66</u>	<u>13</u>	<u>1.50</u>	<u>21</u>	<u>5.11</u>	<u>25</u>	<u>6.36</u>	<u>23</u>	<u>3.23</u>	<u>62</u>	<u>7.85</u>	<u>6</u>	<u>9.69</u>	<u>6</u>	<u>3.11</u>	<u>16</u>	<u>4.68</u>	<u>11</u>	<u>8.78</u>	<u>11</u>	<u>2.04</u>	<u>15</u>	<u>5.65</u>	<u>7</u>	<u>12.57</u>	<u>1.44</u>	<u>12</u>	<u>6.54</u>	<u>21</u>	<u>10.2</u>	<u>21</u>	<u>1.58</u>	<u>16</u>	

Figure 1: Screenshot of our IE-ranking in the Middlebury benchmark (taken on the 25th of September, currently private).

	training dataset	Middlebury Baker <i>et al.</i> [1]		Vimeo-90k Xue <i>et al.</i> [5]		UCF101 - DVF Liu <i>et al.</i> [2]		Xiph - 1K (4K resized to 1K)		Xiph - 2K (4K resized to 2K)	
		SSIM	relative	SSIM	relative	SSIM	relative	SSIM	relative	SSIM	relative
		↑	improvement	↑	improvement	↑	improvement	↑	improvement	↑	improvement
original SepConv	proprietary	0.959	—	0.956	—	0.947	—	0.959	—	0.929	—
reimplementation	Vimeo-90k	0.958	—	0.953	—	0.949	—	0.959	—	0.925	—
+ delayed padding	— ” —	0.959	+ 0.001	0.956	+ 0.003	0.950	+ 0.001	0.960	+ 0.001	0.926	+ 0.001
+ input normalization	— ” —	0.962	+ 0.003	0.957	+ 0.001	0.950	+ 0.000	0.955	+ 0.005	0.925	+ 0.001
+ improved network	— ” —	0.963	+ 0.001	0.957	+ 0.000	0.950	+ 0.000	0.962	+ 0.007	0.930	+ 0.005
+ normalized kernels	— ” —	0.967	+ 0.004	0.960	+ 0.003	0.950	+ 0.000	0.963	+ 0.001	0.932	+ 0.002
+ contextual training	— ” —	0.968	+ 0.001	0.961	+ 0.001	0.950	+ 0.000	<u>0.964</u>	+ 0.001	0.933	+ 0.001
+ self-ensembling	— ” —	<u>0.969</u>	+ 0.001	<u>0.962</u>	+ 0.001	<u>0.951</u>	+ 0.001	<u>0.964</u>	+ 0.000	<u>0.934</u>	+ 0.001

Table 1: Ablation experiments to quantitatively analyze the effects of our proposed techniques. In short, they each positively affect the interpolation quality across different dataset as long as the inter-frame motion does not exceed the kernel size.

	training dataset	Middlebury Baker <i>et al.</i> [1]		Vimeo-90k Xue <i>et al.</i> [5]		UCF101 - DVF Liu <i>et al.</i> [2]		Xiph - 1K (4K resized to 1K)		Xiph - 2K (4K resized to 2K)	
		SSIM	relative	SSIM	relative	SSIM	relative	SSIM	relative	SSIM	relative
		↑	improvement	↑	improvement	↑	improvement	↑	improvement	↑	improvement
SepConv - $\mathcal{L}_1$	proprietary	0.959	—	0.956	—	0.947	—	0.959	—	0.929	—
Ours - $\mathcal{L}_{C_{Tx}}$	Vimeo-90k	0.968	+ 0.009	0.961	+ 0.005	0.950	+ 0.003	<u>0.964</u>	+ 0.005	0.933	+ 0.004
Ours - $\mathcal{L}_{C_{Tx}} - 8\times$	— ” —	<u>0.969</u>	+ 0.001	<u>0.962</u>	+ 0.001	<u>0.951</u>	+ 0.001	<u>0.964</u>	+ 0.000	<u>0.934</u>	+ 0.001

Table 2: Quantitative comparison with SepConv. We list two separate results of our proposed approach, one without and one with self-ensembling. The self-ensembling is denoted by  $8\times$  as it represents a combination of eight independent estimates.

	venue	Middlebury Baker <i>et al.</i> [1]		Vimeo-90k Xue <i>et al.</i> [5]		UCF101 - DVF Liu <i>et al.</i> [2]		Xiph - 1K (4K resized to 1K)		Xiph - 2K (4K resized to 2K)	
		SSIM	absolute	SSIM	absolute	SSIM	absolute	SSIM	absolute	SSIM	absolute
		↑	rank	↑	rank	↑	rank	↑	rank	↑	rank
SepConv - $\mathcal{L}_1$	ICCV 2017	0.959	8 <sup>th</sup> of 10	0.956	9 <sup>th</sup> of 10	0.947	10 <sup>th</sup> of 10	0.959	10 <sup>th</sup> of 10	0.929	7 <sup>th</sup> of 10
CtxSyn - $\mathcal{L}_{Lap}$	CVPR 2018	0.964	7 <sup>th</sup> of 10	0.961	5 <sup>th</sup> of 10	0.949	9 <sup>th</sup> of 10	0.963	5 <sup>th</sup> of 10	0.936	3 <sup>rd</sup> of 10
DAIN	CVPR 2019	0.965	5 <sup>th</sup> of 10	<u>0.964</u>	2 <sup>nd</sup> of 10	0.950	3 <sup>rd</sup> of 10	<u>0.965</u>	2 <sup>nd</sup> of 10	<u>0.939</u>	2 <sup>nd</sup> of 10
CAIN	AAAI 2020	0.951	10 <sup>th</sup> of 10	0.959	8 <sup>th</sup> of 10	0.950	3 <sup>rd</sup> of 10	0.961	8 <sup>th</sup> of 10	0.936	3 <sup>rd</sup> of 10
EDSC - $\mathcal{L}_C$	arXiv 2020	0.967	4 <sup>th</sup> of 10	0.961	5 <sup>th</sup> of 10	0.950	3 <sup>rd</sup> of 10	0.963	5 <sup>th</sup> of 10	OOM	OOM
AdaCoF	CVPR 2020	0.959	8 <sup>th</sup> of 10	0.956	9 <sup>th</sup> of 10	0.950	3 <sup>rd</sup> of 10	0.960	9 <sup>th</sup> of 10	0.927	8 <sup>th</sup> of 10
SoftSplat - $\mathcal{L}_{Lap}$	CVPR 2020	<u>0.971</u>	1 <sup>st</sup> of 10	<u>0.970</u>	1 <sup>st</sup> of 10	<u>0.952</u>	1 <sup>st</sup> of 10	<u>0.969</u>	1 <sup>st</sup> of 10	<u>0.944</u>	1 <sup>st</sup> of 10
BMBC	ECCV 2020	0.965	5 <sup>th</sup> of 10	<u>0.964</u>	2 <sup>nd</sup> of 10	0.950	3 <sup>rd</sup> of 10	0.963	5 <sup>th</sup> of 10	OOM	OOM
Ours - $\mathcal{L}_{C_{Tx}}$	N/A	0.968	3 <sup>rd</sup> of 10	0.961	5 <sup>th</sup> of 10	0.950	3 <sup>rd</sup> of 10	0.964	3 <sup>rd</sup> of 10	0.933	6 <sup>th</sup> of 10
Ours - $\mathcal{L}_{C_{Tx}} - 8\times$	— ” —	<u>0.969</u>	2 <sup>nd</sup> of 10	0.962	4 <sup>th</sup> of 10	<u>0.951</u>	2 <sup>nd</sup> of 10	0.964	3 <sup>rd</sup> of 10	0.934	5 <sup>th</sup> of 10

Table 3: Quantitative comparison with recent approaches for video frame interpolation. In addition to highlighting the best result by underlining it, we emphasize the second-best result via a dotted underline. Note that some methods were unable to run on 2K footage due to exceeding the 16 gigabytes of memory available on our graphics card (denoted as “OOM”).