

---

# Optical Flow Regularization of Implicit Neural Representations for Video Frame Interpolation

---

**Anonymous Author(s)**

Affiliation

Address

email

## Abstract

Recent works have shown the ability of Implicit Neural Representations (INR) to carry meaningful representations of signal derivatives. In this work, we leverage this property to perform video frame interpolation by explicitly constraining the derivatives of the INR to satisfy the optical flow constraint equation. We achieve state of the art video frame interpolation on limited motion ranges using only a target video and its optical flow, without learning the interpolation operator from additional training data. We further show that constraining the INR derivatives not only allows to interpolate intermediate frames but also improves the ability of narrow networks to fit observed frames, which suggests potential applications to INR optimization and video compression.

## 1 Introduction

Many core concepts across the fields of signal processing are defined in terms of continuous functions and their derivatives: surfaces are continuous manifolds in space, motion is a rate of change in space through time, etc. In contrast, the modern digital infrastructure is inherently discrete: digital sensors capture discrete observations of the world sampled in time and space; digital computers store and process discrete representations of signals. In order to model continuous notions on discrete signal representations, classical signal processing approaches have resorted to a variety of heuristics and assumptions, often taking the form of constant first or second derivatives of the signal between consecutive observations. The lack of generality of any such handcrafted heuristics, combined with the ever improving quantitative results of Machine Learning (ML) approaches, have led to the near ubiquitous use of ML approaches in recent signal processing research. These approaches leverage large collections of data to infer statistical properties of signals instead of hand-crafted heuristics.

In computer vision, Video Frame Interpolation (VFI) is one task representative of such development. VFI models aim to interpolate intermediate frames between the consecutive frames of a video. To do so, most successful approaches rely on the optical flow as an approximation of the motion field to guide the interpolation of pixel intensities from the grid of two consecutive frames onto the pixel grid of intermediate frames. Classical approaches formulate assumptions such as constant speed or acceleration of the motion field between consecutive frames [CITE]. The value of each pixel in the inferred intermediate frame is computed by first shifting the pixel intensities of the observed frames following the optical flow directions, and then interpolating the shifted pixel intensities onto the intermediate frame's pixel grid. These approaches suffer from the following two limitations:

- Optical flow constraint used to infer the optical flow holds for limited situations.
- Linear interpolation of pixel intensities along the optical flow directions does not hold in practice.

35 These limitations share a common root cause: discretization. Indeed, both the optical flow constraint  
 36 and the constant motion field assumption only truly hold at the infinitesimal scale, for much smaller  
 37 time deltas than typical FPS used in practice.

38 ML approaches [CITE] have instead proposed to learn the frame interpolation operator from large  
 39 video collections, without explicitly formulating any assumption on the optical flow. While these  
 40 approaches have achieved great success in terms of benchmark performance, they are prone to  
 41 generalization errors when applied to unseen videos. Indeed differences between the training set  
 42 distribution (i.e. VFI benchmark videos) and the target video distribution hinders the performance of  
 43 ML approaches: differences in the range of motion, exposure time and frame-per-second have been  
 44 shown to limit the generalization of state-of-the-art models to video frame interpolation in the wild  
 45 [CITE].

46 In the mean time, research on implicit representations seek better discrete representations of con-  
 47 tinuous signals. In recent years Implicit Neural Representations (INR), i.e. representing signals  
 48 as Neural Networks (NN) have been shown to offer several competitive advantages over explicit  
 49 representations, with notable early successes for 3D shape representations [CITE]. Of particular  
 50 interest to us is the work of SIREN [CITE], in which it has been shown that representing signals  
 51 using Multi Layer Perceptrons (MLP) with sine activation functions carry meaningful representations  
 52 of the signal derivatives. Inspired by this work, we question wether such approach may be used  
 53 to guide the interpolation process of VFI by controlling the exact derivatives of the signal rather  
 54 than finite differences, thus avoiding the discretization pitfalls of traditional approaches. We do so  
 55 by constraining the derivatives of SIREN representations to satisfy the optical flow constraint, i.e.,  
 56 to be orthogonal to the video’s optical flow (which we compute using existing state-of-the-art OF  
 57 models). We find that this approach outperforms most existing machine learning-based approaches on  
 58 small motion range benchmarks, without relying on machine learning for the interpolation operator:  
 59 we simply regularize the implicit representation to satisfy the definition of the optical flow. In this  
 60 sense, our approaches is most similar to classical VIF approaches, except that instead of wrapping  
 61 the OF on discrete explicit frame representations, we apply the optical flow constraint on the exact  
 62 gradient of the the INR. Our method is thus not subject to any mismatch between training and test  
 63 data. Furthermore, our approach can sample any number of frame in-between the observed frames  
 64 due to the continuous nature of the representation. In addition to its application to VFI, we also show  
 65 that constraining the gradient of the model also improves the ability of narrow MLPs to fit the signal,  
 66 suggesting potential applications in INR optimization and video compression.

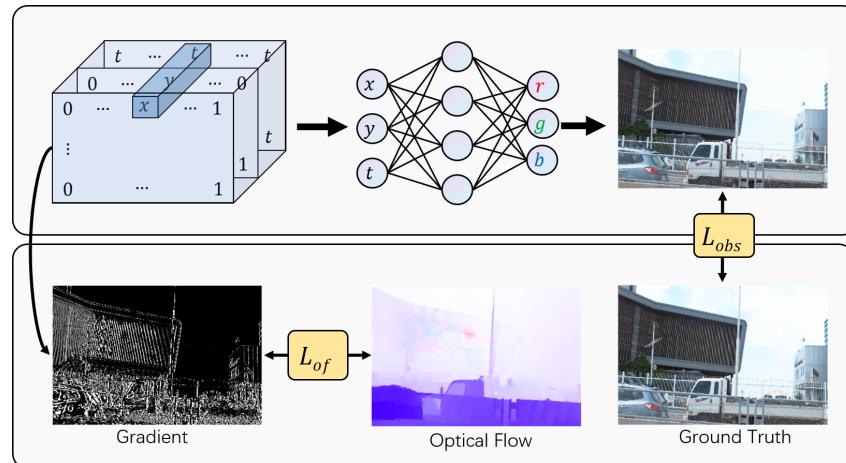


Figure 1: Illustration of our approach

67 To summarize, the contributions of this work are:

- 68 • We propose a regularization method for SIREN which achieve state-of-the-art video frame  
 69 interpolation on small motion ranges.
- 70 • In contrast to other state-of-the art approaches, our approach does not rely on training on a  
 71 large external training set. It only relies on the target video and its estimated optical flow.

- 72 • We show that our regularization approach not only helps generalizing to intermediate frame  
73 generalization but also helps narrow models fit the observed frames.

74 On the other hand, our approach (in its current form) presents important limitations:

- 75 • It relies on an input optical flow, which is computed using existing ML-based model and  
76 thus suffers the limitations of ML approaches.  
77 • Optimization of the INR is very time-consuming, which hinders our ability to work on full  
78 resolution videos for time constraints.  
79 • Our method currently only works on limited motion range. It does not match state-of-the art  
80 ML models on large motion ranges.

81 While we acknowledge the importance of the above limitations, we believe these to not be fundamental  
82 limitations of our approach but rather important future INR research directions. We discuss these  
83 limitations at length and present possible axis to tackle them in Section XXX. The remainder of this  
84 paper is organized as follows: We briefly present some related work in Section XXX, the detail of our  
85 method in Section XXX, and design several experiments to highlight the advantages of our approach  
86 in Section XXX.

## 87 2 Related Work

88 **Classical video interpolation.** The classical approach relies on optical flow to characterize the  
89 motion relationship between two frames. The quality of interpolation depends on the accuracy of the  
90 optical flow. Classical optical flow is usually sparse and can only work at one pixel or sub-pixel level  
91 of motion [1]. Sparse feature extraction [2] is used to enhance the optical flow correctly. In order  
92 to synthesize dense optical flow and work on large displacements, [3] uses sparse convolution and  
93 max-pooling. But instead of learning the model, the parameters are set manually.

94 **Deep learning video interpolation.** A number of deep learning models have been developed for  
95 video interpolation tasks. Almost all models can be categorized as: optical flow based, and kernel  
96 based.

97 *Optical Flow-Based.* Optical flow-based approaches are the most popular in video frame interpolation.  
98 The standard technique of video frame interpolation aims at explicitly estimating motion in the form  
99 of optical flow, warping two input frames to an intermediate frame, and synthesizing the occlusion  
100 region. The frames are constrained by the assumption of linear motion and constant luminance  
101 between them. However, video interpolation of video frames is heavily dependent on the accuracy of  
102 optical flow.

103 The Super-SloMo [4] proposed by Jiang et.al. is a non-negligible work in the task of optical flow-  
104 based video frame interpolation. Super-SloMo extends the U-Net architecture proposed by Liu et  
105 al [5]. The bilateral optical flow is calculated for the input two frames and approximates the key  
106 frame with the intermediate optical flow of the two frames. Then the frames of the input are warped  
107 according to the obtained intermediate optical flow.

108 RRIN [6] mentioned that the estimation of intermediate frames in Super-SloMo works poorly near the  
109 boundaries because the optical flow is not locally smooth in these regions. RRIN proposes to improve  
110 the accuracy of optical flow by residual learning. BMBC [7] adds two additional approximate vectors  
111 to Super-SloMo to make the bilateral motion estimation more accurate.

112 *Kernel-Based.* To avoid explicit motion estimation and warping stages, the kernel-based approach  
113 performs a convolution operation on the input frames and the output of the convolution is used as  
114 the result of interpolating the frames. Niklaus et al. [8] proposed a fully convolutional deep neural  
115 network using a spatially adaptive convolutional kernel to perform the prediction of intermediate  
116 frames for two frames with consecutive inputs. Niklaus et al. [9] improved their method by using a  
117 separable convolution with spatially adaptive one-dimensional convolutional kernel pairs estimated  
118 for each pixel, in reducing the parameters of the model. The results of kernel-based methods for  
119 frame interpolation can be limited by the size of the kernel.

120 Lee et al. proposed Adacof [10], which can use any pixel at any position for convolution operation,  
121 so that the convolution kernel is no longer limited to the local range. And many methods residing

122 in optical flow are defined as a special case of Adacof. However, most kernel-based methods can  
 123 only generate one intermediate frame, and if one wants to generate multiple intermediate frames, one  
 124 needs to do it recursively. EDSC [11] is the first kernel-based method proposed to generate multiple  
 125 intermediate frames, but the results are not as good as the optical flow method.

### 126 Implicit Neural Representation. (INR)

127 INR use a neural network to represent an object approximately, which is essentially a way to  
 128 parameterize the signal. Since [12], [13] was developed, INR has performed well in the areas of 3D  
 129 vision tasks, images, and video. The image and video tasks most relevant to this paper are around the  
 130 direction of image/video compression.

131 COIN [14] first proposed the use of INR to compress images, mapping pixel coordinates to RGB  
 132 values. COIN++ [15] cooperated with the meta-learning approach for image compression work based  
 133 on COIN. In the field of video compression, NeRV [16] proposed by Chen et al. successfully encodes  
 134 the video into a neural network, i.e., the content of the video is saved using a neural network. Only the  
 135 frame index of the model needs to be provided, and the corresponding RGB picture is output. In other  
 136 words, this makes it possible to output infinite frames of video using a neural network. Although  
 137 NeRV briefly attempts the task of performing video frame interpolation, this is not NeRV's main  
 138 work. The NRFF [17] proposed by Rho et al., which uses optical flow and residuals information for  
 139 video compression, does not directly fit all frames.

140 Most related to our approach is the concurrent work by Shangguan et al. [18], which also uses  
 141 INR for video interpolation tasks. Their approach, CURE, uses machine learning. It requires visual  
 142 features of the video and does not fully map the pixel coordinates and frame positions of the video to  
 143 RGB images.

## 144 3 Method

145 We consider a ground-truth video as a continuous signal  $v$  mapping continuous spatial  $(x, y)$  and  
 146 temporal  $(t)$  coordinates to RGB values:

$$v : (x, y, t) \rightarrow (R, G, B) \quad (1)$$

$$v : \mathbb{R}^3 \rightarrow \mathbb{R}^3$$

147 Our goal is to find a continuous function  $f_\theta$ , parameterized by a finite parameter set  $\theta \in \Theta$ , with  
 148 minimum distance  $d$  to the ground-truth signal:

$$f_\theta : (x, y, t) \rightarrow (R, G, B) \quad (2)$$

$$\text{s.t. } \theta = \min_{\Theta} \iiint d(f_\theta(x, y, t), v(x, y, t)) dx dy dt$$

149 where the distance function  $d$  may either be the Peak Signal to Noise Ratio (PSNR) or the Structural  
 150 Similarity Index Measure (SSIM). To do so, we only have access to regularly sampled observation of  
 151 the signal  $v$  (i.e. the explicit representation of the video), which we denote as:

$$\mathcal{V} \in \mathbb{R}^{T \times H \times W \times 3} \quad (3)$$

$$\text{s.t. } \mathcal{V}_{xyt} = v(x, y, t)$$

152 where  $T$  represents the number of frames in the video, and  $H \times W$  the spatial resolution. We use  
 153 SIREN as parameterized function class  $f_\theta$ . The most straightforward way to approximate Equation 2  
 154 is to optimize the model parameters so as to fit the video frames, using the following loss function we  
 155 refer to as the observation loss:

$$\mathcal{L}_{obs} = \frac{1}{HWT} \sum_{x=1}^W \sum_{y=1}^H \sum_{t=1}^T \|f_\theta(x, y, t) - \mathcal{V}_{xyt}\|^2 \quad (4)$$

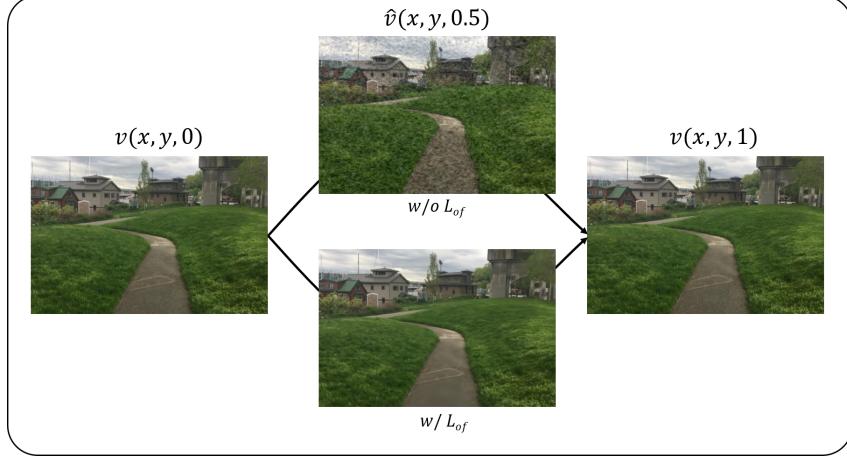


Figure 2: Illustration of INR frame interpolation with and without optical flow regularization. Without regularization (middle top), intermediate frames show unnatural high-frequency variations. Regularizing the INR to satisfy the optical flow constraint equation result in nicely interpolated frames (middle bottom).

156 However, we found that optimizing the INR to only minimize this observation loss leads to overfitting  
 157 the observation with high temporal frequencies: the intra-frame signal, which we aim to correctly  
 158 recover, shows important deviations from the observed frames, as illustrated in Figure XXX. This  
 159 observation has lead us to consider fitting not only the signal itself, but to also constrain its derivatives.  
 160 In particular, we regularize the model so as to respect the optical flow constraint.

161 The optical flow constraint equation states that for an infinitesimal lapse of time  $\delta t$ , the brightness of  
 162 physical points perceived by a camera at arbitrary coordinates  $(x, y, t)$  should remain constant. In  
 163 other words, given the displacement  $(\delta x, \delta y)$  of a physical point in the image coordinate system, the  
 164 image brightness  $v$  should remain constant:

$$v(x, y, t) = v(x + \delta x, y + \delta y, t + \delta t) \quad (5)$$

165 We introduce the vector notation  $\mathbf{x} = (x, y, t)$  for readability. Expressing movement as a ratio of  
 166 displacement in time and abbreviating, we can write the optical flow  $F$  and the above constraint as:

$$\begin{aligned} F(\mathbf{x}) &= \left( \frac{\delta x}{\delta t}, \frac{\delta y}{\delta t}, 1 \right) \\ v(\mathbf{x}) &= v(\mathbf{x} + F(\mathbf{x})) \end{aligned} \quad (6)$$

167 First order Taylor expansion allows us to rewrite expand equation XXX into the following

$$\begin{aligned} v(\mathbf{x}) &= v(\mathbf{x}) + \frac{\delta v}{\delta \mathbf{x}} \cdot F(\mathbf{x}) \\ \frac{\delta v}{\delta \mathbf{x}} \cdot F(\mathbf{x}) &= 0 \end{aligned} \quad (7)$$

168 which holds exactly in the limit of infinitesimal  $\delta t$ . We leverage this optical flow constraint equation  
 169 to regularize the INR. Denoting the derivatives of the SIREN as:

$$D(f, \theta, x, y, t) = \left( \frac{\delta f_\theta(x, y, t)}{\delta x}, \frac{\delta f_\theta(x, y, t)}{\delta y}, \frac{\delta f_\theta(x, y, t)}{\delta t} \right) \quad (8)$$

170 we can now define the optical flow regularization loss

$$\mathcal{L}_{of} = \frac{1}{HWT} \sum_{x \in W} \sum_{y \in H} \sum_{t \in T} |D(f, \theta, x, y, t) \cdot F(x, y, t)| \quad (9)$$

171 This loss constrains the derivatives of the signal to be orthogonal to the optical flow and can be  
172 understood as keeping constant brightness along the optical flow trajectories. The total loss we use to  
173 optimize the INR is a weighted sum of these two terms:

$$\mathcal{L} = \lambda \mathcal{L}_{obs} + (1 - \lambda) \mathcal{L}_{of} \quad (10)$$

174 where  $\lambda$  is a hyper-parameter taking values between 0 and 1 whose impact we investigate in the  
175 following section. The exactitude of the optical flow constraint at the infinitesimal scale plays in our  
176 favor: As we regularize the true derivative of the signal representation, we do not assume constant  
177 derivatives of the signal on any interval. We believe this is the main factor behind our positive results.  
178 On the other hand, the optical flow we used was estimated from discrete consecutive frames, and  
179 thus does not represent the true infinitesimal motion field but an estimation of finite differences. We  
180 discuss this limitation in Section XXX.

## 181 4 Experiments

182 Following previous works, we use the Adobe[CITE], X4K[CITE] and ND Scene[CITE] dataset  
183 as benchmarks to compare to the state-of-the-art. We run all additional experiments on the  
184 720p240fps1.mov video of Adobe dataset illustrated in Figure 2. Due to the time-consuming  
185 operation of optimizing SIREN representations, we optimize and evaluate all models on a  $240 \times 360$   
186 resolution. For each video in the Adobe dataset, we selected only the first 40 frames for our experi-  
187 ments. For the Adobe dataset, we used the same dataset splitting approach as Super-SloMo, where  
188 eight videos were used as the testing set.

189 We show the effect of different hyperparameters on the model in Figure 4. We use the greedy  
190 algorithm to find the best combination of hyperparameters. We will end up using a SIREN model  
191 with 6 depth and 720 width. We use an omega of 25 and a lambda of 0.12. We optimize the models  
192 using the Adam optimizer using a cosine learning rate with maximum learning rate of 3.6e-5 during  
193 15k epochs.

194 We start by showing the impact of controlling the fit of high frequencies without the optical flow  
195 loss in section ???. We show that limiting the frequency fitted improves interpolation at the cost of  
196 degrading the fit of observed frames, but it does not allow to reach the same accuracy as optical flow  
197 regularization, showing that OF regularization does more than only regularizing the fitted frequency  
198 range. In Section 4.2, we quantitatively compare our results to state of the art models on standard  
199 datasets. We show that our approach achieves state-of-the-art results on low-range motion datasets,  
200 but underperforms existing methods for high-range motion videos. We present an ablation study in  
201 Section 4.3, providing insights and appropriate settings for the different model hyper-parameters, and  
202 a qualitative analysis of our results in Section 4.4.

203 Finally, we report an additional counter-intuitive result in Section 4.5. We show that our proposed  
204 optical flow regularization loss can help the lightweight SIREN model to fit the video better. This  
205 indicates that our method is potentially helpful for video compression.

### 206 4.1 Optical flow constraint and high frequencies

207 Figure 2 illustrates the fact that applying the optical flow constraint smoothes out the high-frequency  
208 variations from the intermediate frames of vanilla SIREN representations. We start by questioning  
209 whether the OF constraint does more than simply removing the high frequency variations of the  
210 representation. To do so, we compare the results of vanilla SIREN representations geared towards  
211 different frequencies and compare their interpolated frames PSNR to those of OF-constrained  
212 representations. We constrain the SIREN frequency by varying their  $\omega$  parameter, and report our  
213 comparison in Figure 3.

214 Constraining the high frequency with  $\omega$  down to 5 degrades the fitting of observed frames, but  
215 improves the quality of interpolated frames. This suggests that  $\omega$  behaves similarly to a regularization  
216 parameter controlling a regime of overfitting versus underfitting of the observed frames high  
217 frequencies. Figure 3 also shows that vanilla SIREN models remain well under the OF-constrained  
218 SIREN, confirming that the OF constraint provides more than high frequency regularization.

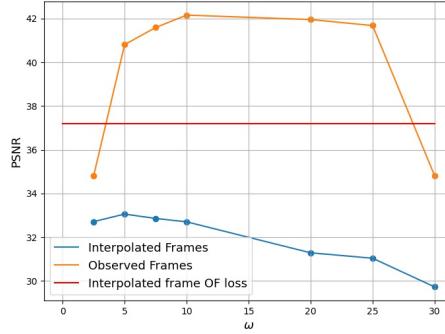


Figure 3: Evolution of the PSNR of observed and interpolated frames with  $\omega$  without OF loss. Limiting the high frequency fit alone does not reach the same interpolation accuracy as the OF loss.

Table 1: Quantitative comparison to state-of-the-art VFI on Standard benchmarks. Results are formatted as PSNR / SSMI.

| (a) Limited motion range |                             |                      |
|--------------------------|-----------------------------|----------------------|
|                          | Adobe-240FPS                | X4K                  |
| Super-SloMo [1]          | 27.77 / 0.886               | 27.38 / 0.852        |
| RRIN [XXX]               | 32.37 / 0.962               | 30.70 / 0.927        |
| BMBC [XXX]               | 27.83 / 0.917               | 27.42 / 0.858        |
| AdaCof [XXX]             | 35.50 / 0.968               | 34.61 / 0.921        |
| ABME [XXX]               | 35.28 / 0.966               | 34.30 / 0.919        |
| FILM [XXX]               | 35.97 / 0.971               | <b>35.14</b> / 0.939 |
| Ours                     | <b>36.52</b> / <b>0.977</b> | 35.06 / <b>0.944</b> |

| (b) Large motion range |                             |
|------------------------|-----------------------------|
|                        | ND Scene                    |
| V-NF [X]               | 23.30 / 0.726               |
| NSFF [X]               | 28.03 / 0.925               |
| CURE [X]               | <b>36.91</b> / <b>0.984</b> |
| Ours                   | 29.22 / 0.921               |

## 219 4.2 State of the art models

220 Table XXX quantitatively compares the results of our method to state-of the art VFI models on  
221 different datasets. Despite its simplicity, and without any training data, our approach outperforms  
222 existing models on limited motion ranges (Table XXX). However, as illustrated in Figure XXX,  
223 it falls short of state-of-the-art methods on the more complex ND Scene benchmark due to larger  
224 motion ranges. We provide further comparison in the qualitative analysis of Section XXX and section  
225 XXX discusses possible ways forward to bridging the gap performance on large motion datasets.

## 226 4.3 Ablation study

227 We describe in this section the method of searching for the best hyperparameter combination using  
228 the greedy algorithm. Figure 4 illustrates the results of the hyperparameter searching. We will start  
229 with a SIREN model of depth 9 and width 512. The baseline experiment is set up with a learning rate  
230 of 1e-5, 5000 epochs, and 30  $\omega$ . The order of our hyperparameter search is the lambda controlling  
231 the loss balance  $\lambda$ , the learning rate, epochs, omega, and the depth and width of SIREN. For each  
232 new hyperparameter searching, the best combination of previous hyperparameters is used.

## 233 4.4 Qualitative analysis

234 Do qualitative analysis

## 235 4.5 Video fitting

236 Fit video.

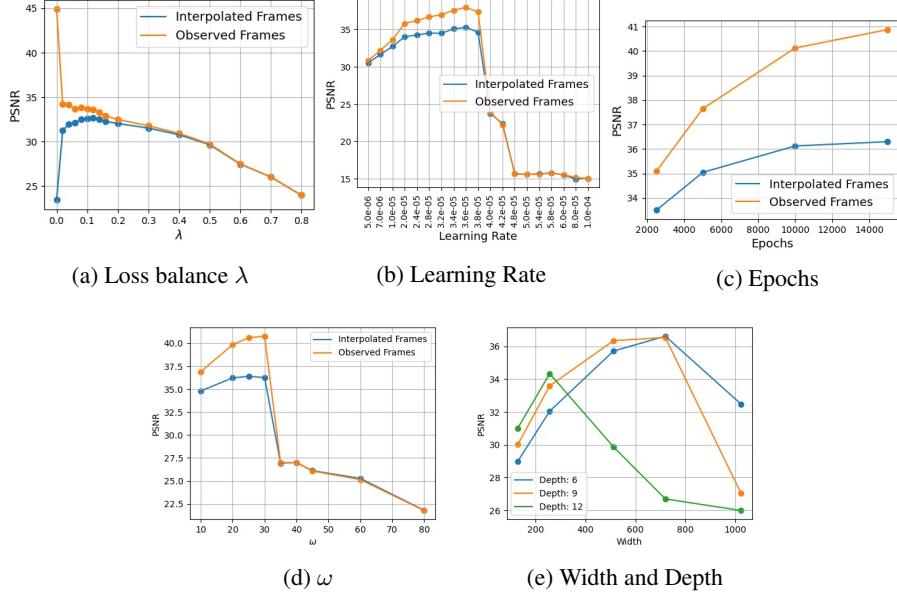


Figure 4: Hyperparameter Searching and Ablation Experiments.

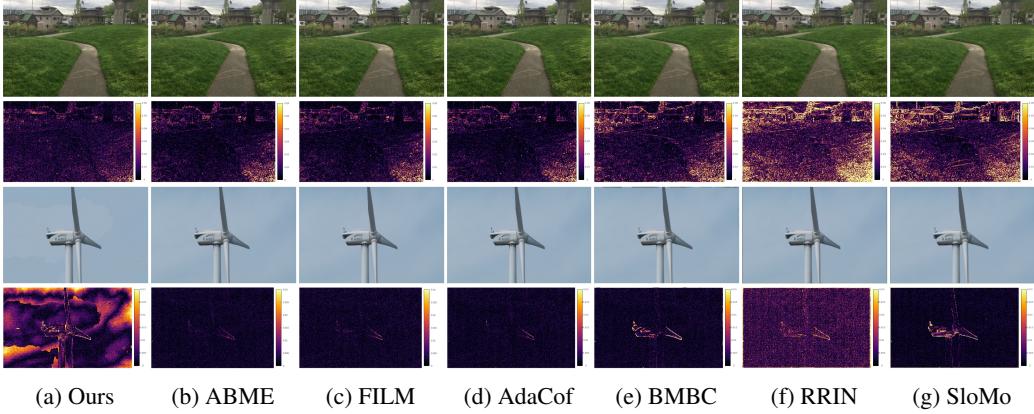


Figure 5: Small Motion Video Qualitative Analysis. The interpolated frame results and the residual heat map of the two videos are shown. Our proposed method can fit high frequency details well (e.g., grass), but fits low frequency information poorly (e.g., sky). The BMBC and RRIN lose the pixels at the edges.

## 237 5 Limitations

238 While we believe our results to be very encouraging, the proposed approach is not yet practical. Here,  
239 we discuss what we believe to be the three main limitations of, and possible solutions to, our approach

240 **Slow optimization process.** Fitting XXX frames of a video at XXX resolution currently takes XXX  
241 hours on a XXX GPU using Pytorch. This computation time is a huge draw back as it limits our  
242 ability to process full resolution video as well as to explore different hyper parameters and variants of  
243 the methods. We expect new methods speeding up the convergence of video INR to be very beneficial  
244 to this line of research. Given recent successes of INR approaches to high impact applications (i.e.,  
245 video compression [CITE]), We hopefully expect to see advances in INR optimization research.

246 **Reliance on trained optical flow model.** SIREN models allow us to apply the optical flow on the  
247 exact derivatives of the signal, bypassing the heuristics of classical approach without relying on  
248 machine learning. The optical flow we use, however, is given by a ML model trained on discrete  
249 representations, which raises two problems: it is subject to the same generalization errors as ML

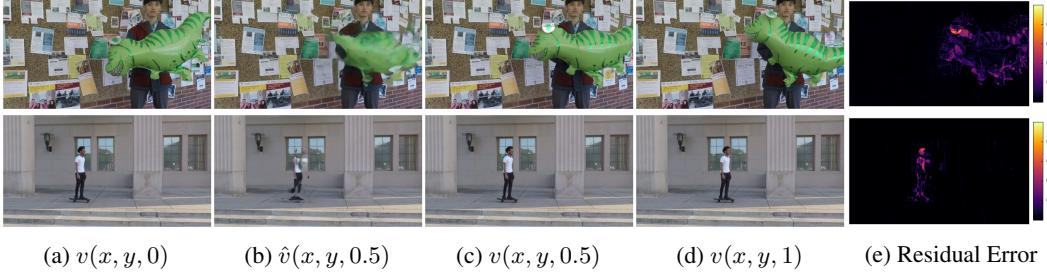


Figure 6: Large Motion Video Qualitative Analysis

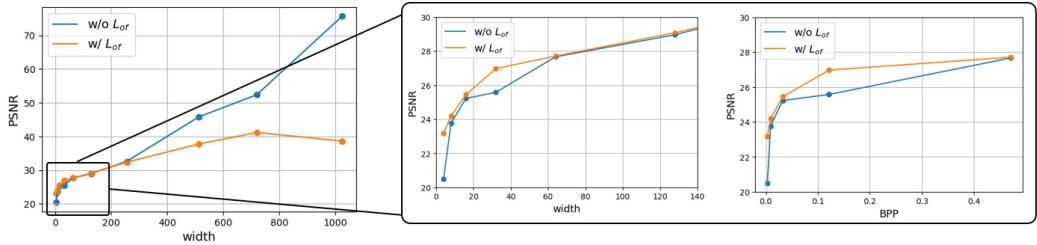


Figure 7: Our proposed optical flow regularization loss can help the lightweight SIREN model to fit the video better.

250 approaches, and is subject to finite difference errors, i.e. noisy illumination variations and occlusions.  
 251 Future work will aim to bypass our reliance on ML-based OF using alternative constraints on the  
 252 exact derivatives.

253 **Inability to interpolate high motion range videos.** In its current form, we only apply the optical  
 254 flow constraint on the observed frames of the video. This has proven sufficient to reach state-of-the  
 255 art on low motion ranges but is not sufficient for large motions. A promising axis of improvement  
 256 would be to apply further regularization on the interpolates frames (i.e for non-inter input times like  
 257  $t = 0.5$ ). Possible regularization may include interpolated optical flows, or texture constraints as has  
 258 been proposed in related works [COTE], which may prevent the ghosting effects illustrated in Figure  
 259 XXX.

## 260 6 Conclusion

261 In this paper, we have shown that regularizing INR using the optical flow constraint equation enabled  
 262 VFI without relying on ML to perform the interpolation step. We have shown that this approach  
 263 is sufficient to reach state-of-the-art interpolation on low motion range videos, without resorting to  
 264 learning-based interpolation. This method is not yet practical for high resolution large motion range  
 265 videos due to the three limitations highlighted in the previous section. Nevertheless, these limitations  
 266 can be tackled and we hope the insights presented in this paper can serve prove useful in the future,  
 267 either on their own or combined with related approaches.

## 268 References

- 269 [1] Lei Chen, Hua Yang, Takeshi Takaki, and Idaku Ishii. Real-time optical flow estimation using  
 270 multiple frame-straddling intervals. *Journal of Robotics and Mechatronics*, 24(4):686–698,  
 271 2012.
- 272 [2] Ce Liu, Jenny Yuen, Antonio Torralba, Josef Sivic, and William T Freeman. Sift flow: Dense  
 273 correspondence across different scenes. In *European conference on computer vision*, pages  
 274 28–42. Springer, 2008.

- 275 [3] Philippe Weinzaepfel, Jerome Revaud, Zaid Harchaoui, and Cordelia Schmid. Deepflow: Large  
 276 displacement optical flow with deep matching. In *Proceedings of the IEEE international*  
 277 *conference on computer vision*, pages 1385–1392, 2013.
- 278 [4] Huaizu Jiang, Deqing Sun, Varun Jampani, Ming-Hsuan Yang, Erik Learned-Miller, and  
 279 Jan Kautz. Super slomo: High quality estimation of multiple intermediate frames for video  
 280 interpolation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*,  
 281 pages 9000–9008, 2018.
- 282 [5] Ziwei Liu, Raymond A Yeh, Xiaou Tang, Yiming Liu, and Aseem Agarwala. Video frame  
 283 synthesis using deep voxel flow. In *Proceedings of the IEEE International Conference on*  
 284 *Computer Vision*, pages 4463–4471, 2017.
- 285 [6] Haopeng Li, Yuan Yuan, and Qi Wang. Video frame interpolation via residue refinement. In  
 286 *ICASSP 2020-2020 IEEE International Conference on Acoustics, Speech and Signal Processing*  
 287 (*ICASSP*), pages 2613–2617. IEEE, 2020.
- 288 [7] Junheum Park, Keunsoo Ko, Chul Lee, and Chang-Su Kim. Bmhc: Bilateral motion estimation  
 289 with bilateral cost volume for video interpolation. In *European Conference on Computer Vision*,  
 290 pages 109–125. Springer, 2020.
- 291 [8] Simon Niklaus, Long Mai, and Feng Liu. Video frame interpolation via adaptive convolution.  
 292 In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages  
 293 670–679, 2017.
- 294 [9] Simon Niklaus, Long Mai, and Feng Liu. Video frame interpolation via adaptive separable  
 295 convolution. In *Proceedings of the IEEE International Conference on Computer Vision*, pages  
 296 261–270, 2017.
- 297 [10] Hyeongmin Lee, Taeoh Kim, Tae-young Chung, Daehyun Pak, Yuseok Ban, and Sangyoun Lee.  
 298 Adacof: Adaptive collaboration of flows for video frame interpolation. In *Proceedings of the*  
 299 *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5316–5325, 2020.
- 300 [11] Xianhang Cheng and Zhenzhong Chen. Multiple video frame interpolation via enhanced  
 301 deformable separable convolution. *IEEE Transactions on Pattern Analysis and Machine*  
 302 *Intelligence*, 2021.
- 303 [12] Ben Mildenhall, Pratul P Srinivasan, Matthew Tancik, Jonathan T Barron, Ravi Ramamoorthi,  
 304 and Ren Ng. Nerf: Representing scenes as neural radiance fields for view synthesis. In *European*  
 305 *conference on computer vision*, pages 405–421. Springer, 2020.
- 306 [13] Vincent Sitzmann, Julien Martel, Alexander Bergman, David Lindell, and Gordon Wetzstein. Im-  
 307 plicit neural representations with periodic activation functions. *Advances in Neural Information*  
 308 *Processing Systems*, 33:7462–7473, 2020.
- 309 [14] Emilien Dupont, Adam Goliński, Milad Alizadeh, Yee Whye Teh, and Arnaud Doucet. Coin:  
 310 Compression with implicit neural representations. *arXiv preprint arXiv:2103.03123*, 2021.
- 311 [15] Emilien Dupont, Hrushikesh Loya, Milad Alizadeh, Adam Goliński, Yee Whye Teh, and Arnaud  
 312 Doucet. Coin++: Data agnostic neural compression. *arXiv preprint arXiv:2201.12904*, 2022.
- 313 [16] Hao Chen, Bo He, Hanyu Wang, Yixuan Ren, Ser Nam Lim, and Abhinav Shrivastava. Nerv:  
 314 Neural representations for videos. *Advances in Neural Information Processing Systems*, 34,  
 315 2021.
- 316 [17] Daniel Rho, Junwoo Cho, Jong Hwan Ko, and Eunbyung Park. Neural residual flow fields for  
 317 efficient video representations. *arXiv preprint arXiv:2201.04329*, 2022.
- 318 [18] Wentao Shangguan, Yu Sun, Weijie Gan, and Ulugbek S Kamilov. Learning cross-video neural  
 319 representations for high-quality frame interpolation. *arXiv preprint arXiv:2203.00137*, 2022.