

---

# 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 neural implicit representations (NIR) 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 NIR 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 NIR 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 NIR 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 infer intermediate frames between 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 Neural Implicit Representations (NIR), 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 NIR. 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 NIR optimization and video compression.

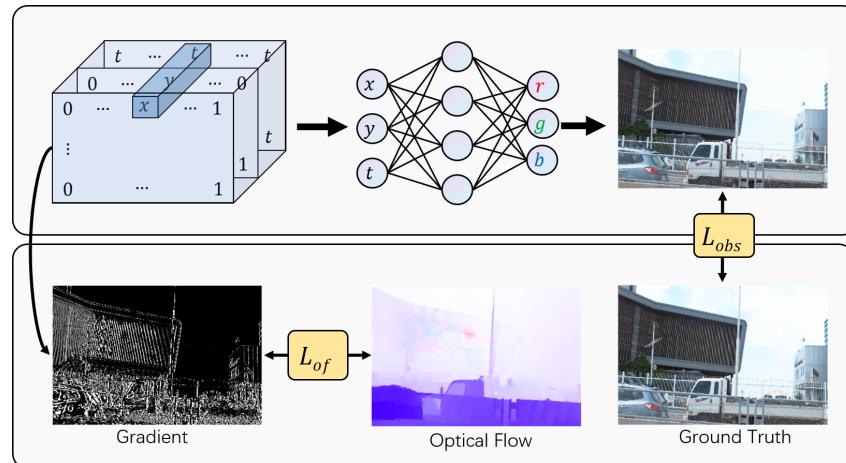


Figure 1: Illustration of our approach

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

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

97 The Super-SloMo ? proposed by Jiang et.al. is a non-negligible work in the task of optical flow-based  
98 video frame interpolation. Super-SloMo extends the U-Net architecture proposed by Liu et al ?. The  
99 bilateral optical flow is calculated for the input two frames and approximates the key frame with the  
100 intermediate optical flow of the two frames. Then the frames of the input are warped according to the  
101 obtained intermediate optical flow.

102 RRIN ? mentioned that the estimation of intermediate frames in Super-SloMo works poorly near  
103 the boundaries because the optical flow is not locally smooth in these regions. RRIN proposes to  
104 improve the accuracy of optical flow by residual learning. BMBC ? adds two additional approximate  
105 vectors to Super-SloMo to make the bilateral motion estimation more accurate.

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

114 Lee et al. proposed Adacof ?, which can use any pixel at any position for convolution operation,  
115 so that the convolution kernel is no longer limited to the local range. And many methods residing  
116 in optical flow are defined as a special case of Adacof. However, most kernel-based methods can  
117 only generate one intermediate frame, and if one wants to generate multiple intermediate frames, one  
118 needs to do it recursively. EDSC ? is the first kernel-based method proposed to generate multiple  
119 intermediate frames, but the results are not as good as the optical flow method.

## 120 Implicit Neural Network Representation. (INR)

121 INR use a neural network to represent an object approximately, which is essentially a way to  
 122 parameterize the signal. Since ?, ? was developed, INR has performed well in the areas of 3D vision  
 123 tasks, images, and video. The image and video tasks most relevant to this paper are around the  
 124 direction of image/video compression.

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

134 Most related to our approach is the concurrent work by XX et al. ?, which also uses INR for video  
 135 interpolation tasks. Their approach, CURE, uses machine learning. It requires visual features of the  
 136 video and does not fully map the pixel coordinates and frame positions of the video to RGB images.

### 137 3 Method

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

$$v : (x, y, t) \rightarrow (R, G, B) \\ v : \mathbb{R}^3 \rightarrow \mathbb{R}^3 \quad (1)$$

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

$$f_\theta : (x, y, t) \rightarrow (R, G, B) \\ s.t. \theta = \min_{\Theta} \iiint d(f_\theta(x, y, t), v(x, y, t)) dx dy dt \quad (2)$$

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

$$\mathcal{V} \in \mathbb{R}^{T \times H \times W \times 3} \\ s.t. \mathcal{V}_{xyt} = v(x, y, t) \quad (3)$$

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

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

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

154 The optical flow constraint equation states that for an infinitesimal lapse of time  $\delta t$ , the brightness of  
 155 physical points perceived by a camera at arbitrary coordinates  $(x, y, t)$  should remain constant. In

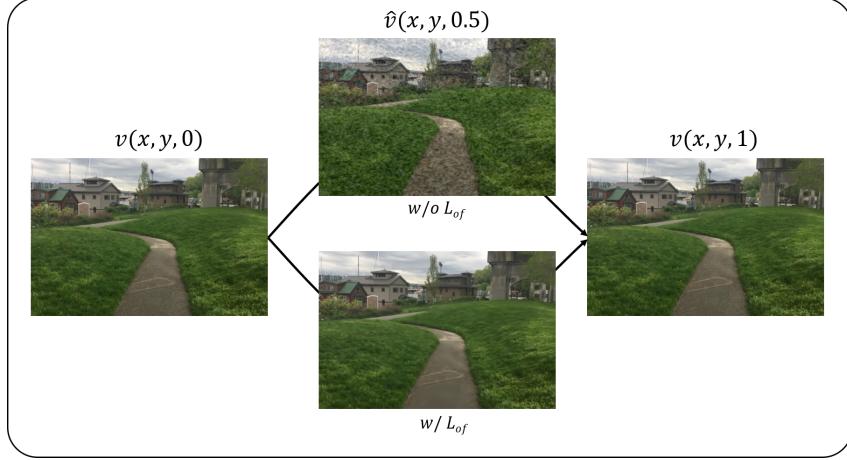


Figure 2: Illustration of NIR frame interpolation with and without optical flow regularization. Without regularization (middle top), intermediate frames show unnatural high-frequency variations. Regularizing the NIR to satisfy the optical flow constraint equation result in nicely interpolated frames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

156 other words, given the displacement  $(\delta x, \delta y)$  of a physical point in the image coordinate system, the  
157 image brightness  $v$  should remain constant:

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

158 Expressing movement as a ratio of displacement in time and abbreviating coordinates as  $x = (x, y, t)$ ,  
159 we can write the optical flow  $F$  and the above constraint as:

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

160 We leverage this optical flow constraint equation to regularize the NIR. Denoting the derivatives of  
161 the video signal as:

$$D(f, \theta, x, y, t) = \left( \frac{\delta f_\theta(x, y, t)}{\delta x}, \frac{\delta f_\theta(x, y, t)}{\delta y}, \frac{\delta f_\theta(x, y, t)}{\delta t} \right) \quad (7)$$

162 And the optical flow as:

163 we can now define the optical flow regularization loss

$$\mathcal{L}_{of} = \frac{1}{HWT} \sum_{x \in W} \sum_{y \in H} \sum_{t \in T} |D(f, \theta, x, y, t) \cdot F(x, y, t)| \quad (8)$$

164 This loss constrains the derivatives of the signal to be orthogonal to the optical flow and can be  
165 intuitively understood as keeping constant brightness along the optical flow trajectories. The total  
166 loss we use to optimize the NIR is a weighted sum of these two terms:

$$\mathcal{L} = \lambda \mathcal{L}_{obs} + (1 - \lambda) \mathcal{L}_{of} \quad (9)$$

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

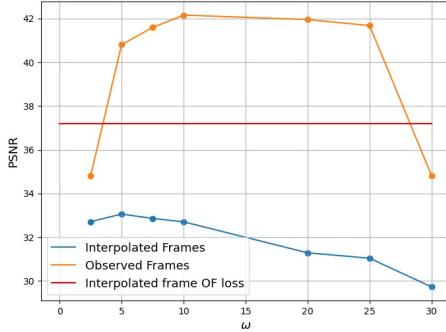


Figure 3: NEED TO SHOW TRAINING CURVE AND THRESHOLD ATTAINED BY OF

## 174 4 Experiments

175 Following previous works, we use the Adobe[CITE], X4K[CITE] and ND Scene[CITE] dataset  
 176 as benchmarks to compare to the state-of-the-art. We run all additional experiments on the  
 177 720p240fps1.mov video of Adobe dataset illustrated in Figure 2. Due to the time-consuming  
 178 operation of optimizing SIREN representations, we optimize and evaluate all models on a  $240 \times 360$   
 179 resolution. For each video in the Adobe dataset, we selected only the first 40 frames for our experi-  
 180 ments. For the Adobe dataset, we used the same dataset splitting approach as Super-SloMo, where  
 181 eight videos were used as the testing set.

182 We show the effect of different hyperparameters on the model in Figure 4. We use the greedy  
 183 algorithm to find the best combination of hyperparameters. We will end up using a SIREN model  
 184 with 6 depth and 720 width. Unless specified otherwise, all experiments are run with a SIREN of  
 185 depth 6 and width 720. We use an omega of 25 and a lambda of 0.12. We optimize the models using  
 186 the Adam optimizer using a cosine learning rate with maximum learning rate of 3.6e-5 during 15k  
 187 epochs.

188 We start by showing the impact of controlling the fit to high frequency without the optical flow loss in  
 189 section 4.1. We show that while limiting the frequency fitted does improve generalization, it does not  
 190 allow to reach the same accuracy as optical flow regularization, showing that OF regularization does  
 191 more than just limiting the fitted frequencies.

192 In Section 4.2, we compare our results to state of the art quantitatively on standard benchmarks.  
 193 We show that our approach achieves state-of-the-art results on low-range motion datasets, but  
 194 underperforms existing methods on the high-range motion dataset. We present an ablation in  
 195 Section 4.3, providing insight and appropriate settings on the different model hyperparameters and a  
 196 qualitative analysis of our results in Section 4.4.

197 Finally, we report a surprising additional result in Section 4.5. We show that our proposed optical flow  
 198 regularization loss can help the lightweight SIREN model to fit the video better. This indicates that  
 199 our method is potentially helpful for video compression.

### 200 4.1 Optical Flow constraint and High Frequencies

201 Figure illustrates the fact that applying the optical flow constraint smoothes out the high-frequency  
 202 variations from the intermediate frames of vanilla SIREN representations. We start by questioning  
 203 whether the OF constraint does more than simply removing the high frequency variations of the  
 204 representation. To do so, we compare the results of vanilla SIREN representations geared towards  
 205 different frequency and compare the best obtained results to OF-constrained representations. We  
 206 constrain the SIREN frequency by varying their  $\omega$  parameter, and report our comparison in Figure  
 207 XXX.

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

	Adobe-240FPS [XXX]	X4K [XXX]
Super-SloMo [XXX]	27.77 / 0.8866	27.38 / 0.8527
RRIN [XXX]	32.37 / 0.9624	30.70 / 0.9270
BMBC [XXX]	27.83 / 0.9172	27.42 / 0.8585
AdaCof [XXX]	35.50 / 0.9684	34.61 / 0.9218
ABME [XXX]	35.28 / 0.9669	34.30 / 0.9195
FILM [XXX]	35.97 / 0.9710	<b>35.14 / 0.9397</b>
Ours	<b>36.52 / 0.9770</b>	35.06 / <b>0.9441</b>

Table 2: Quantitative comparison to state-of-the-art VFI on large motion range benchmarks. Results are formatted as PSNR / SSMI.

	ND Scene [XXX]
V-NF	23.30 / 0.7260
NSFF [10]	28.03 / 0.9250
CURE [11]	<b>36.91 / 0.9843</b>
Ours	29.22 / 0.9215

208 While constraining the high frequency with low  $\omega$  does improve the ability to interpolate intermediate  
 209 frames, vanilla SIREN models remain well under the OF-constrained representations, confirming  
 210 than the OF constraint provides more simply restricting the high temporal frequencies.

## 211 4.2 State of the art models

212 Table XXX quantitatively compare the results of our model to state-of the art VFI models on different  
 213 datasets. We show that

## 214 4.3 Ablation study

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

221 Next, we highlight the

## 222 4.4 Qualitative Analysis

## 223 4.5 Video fitting

## 224 5 Limitations

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

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

233 **Reliance on trained optical flow model.** SIREN models allow us to apply the optical flow on the  
 234 exact derivatives of the signal, thus bypassing the heuristics of classical approach without relying on  
 235 machine learning. The optical flow we use, however, is given by a trained ML model, which raises

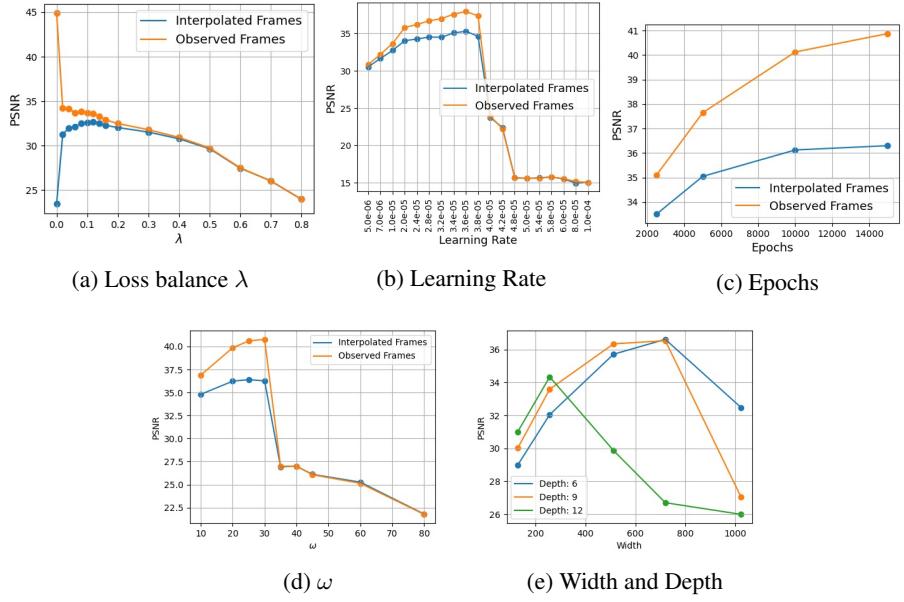


Figure 4: Hyperparameter Searching and Ablation Experiments.

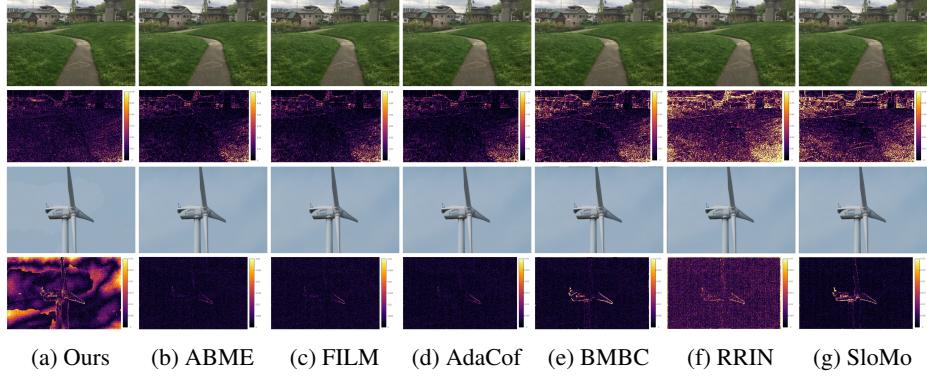


Figure 5: Small Motion Video Qualitative Analysis. The interpolated frame results and the residual heat map of the two videos are shown. Our proposed method can fit high frequency details well (e.g., grass), but fits low frequency information poorly (e.g., sky). The BMBC and RRIN lose the pixels at the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

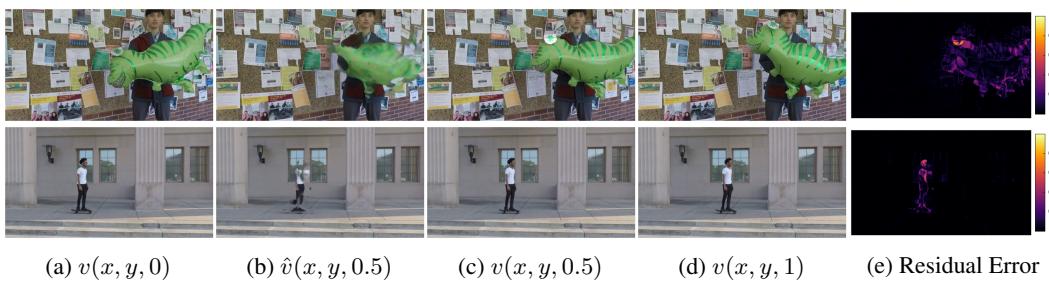


Figure 6: Large Motion Video Qualitative Analysis

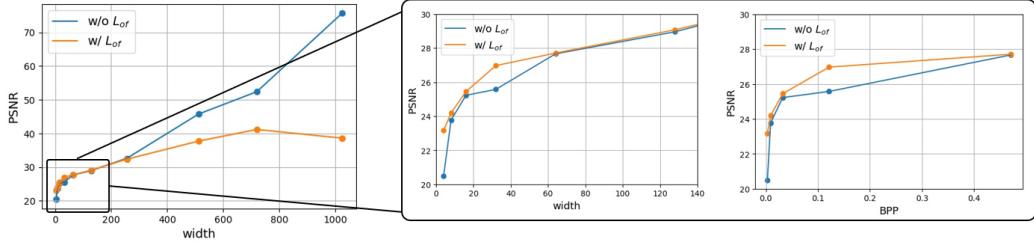


Figure 7: Our proposed optical flow regularization loss can help the lightweight SIREN model to fit the video better.

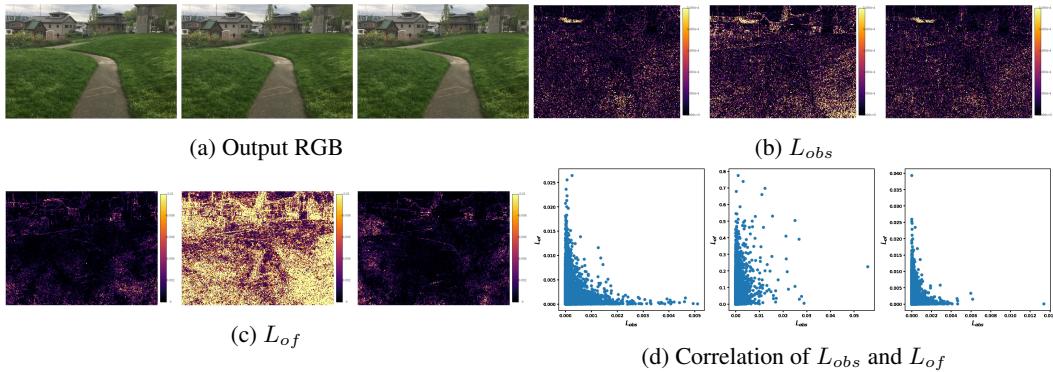


Figure 8: Figures from left to right are: Output Observed Frame  $I_0$ , Output Interpolated Frame  $I_{0.5}$ , Output Observed Frame  $I_1$

236 two problems: it is subject to generalization error, and the flow is computed on discrete samples and  
 237 then subject to undesirable changes in illumination and occlusion. Future work will aim to bypass  
 238 our reliance on ML-based OF using proxy constraints on the exact derivatives.

239 **Inability to interpolate high motion range videos.** In its current form, our approach only regu-  
 240 larizes observed frames of the video. This has proven sufficient to reach state-of-the art on low  
 241 motion ranges but is not sufficient for large motions. For larger motions several improvements can be  
 242 considered, most notably by regularizing intermediate frames. Texture conservation in intermediate  
 243 frames, interpolated optical flows.

## 244 6 Conclusion

245 In this paper, we have shown that regularizing NIR using the optical flow constraint equation enabled  
 246 VFI without relying on ML to perform the interpolation step. We show that this approach is sufficient  
 247 to reach state-of-the-art interpolation on low motion ranges

## 248 References

- 249 [1] Alexander, J.A. & Mozer, M.C. (1995) Template-based algorithms for connectionist rule extraction. In  
 250 G. Tesauro, D.S. Touretzky and T.K. Leen (eds.), *Advances in Neural Information Processing Systems 7*, pp.  
 251 609–616. Cambridge, MA: MIT Press.
- 252 [2] Bower, J.M. & Beeman, D. (1995) *The Book of GENESIS: Exploring Realistic Neural Models with the  
 253 GEneral NEural SImulation System*. New York: TELOS/Springer–Verlag.
- 254 [3] Hasselmo, M.E., Schnell, E. & Barkai, E. (1995) Dynamics of learning and recall at excitatory recurrent  
 255 synapses and cholinergic modulation in rat hippocampal region CA3. *Journal of Neuroscience* **15**(7):5249–5262.