
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 as
 48 Neural Networks (NN) have offered valuable alternative as representations for a variety of signals
 49 [CITE]. Of particular interest to us is the work of SIREN [CITE], in which it has been shown that
 50 representing signals using Multi Layer Perceptrons (MLP) with sine activation functions allowed
 51 for meaningful representations of the signal derivatives. Inspired by this work, we question whether
 52 such approach may be used to guide the interpolation process of VFI by applying the OF on the
 53 exact representation derivatives, thus avoiding the discretization pitfalls of traditional approaches.
 54 We do so by regularizing the derivatives of SIREN representations of videos to satisfy the optical
 55 flow constraint, i.e., to be orthogonal to their optical flow (which we compute using existing state-of-
 56 the-art OF models). We find that this approach outperforms most existing machine learning-based
 57 approaches on small motion range benchmarks, without relying on machine learning to learn the
 58 interpolation: we simply regularize the implicit representation using the definition of the optical flow
 59 and the optical flow constraint equation. In this sense, our approach is most similar to classical VIF
 60 approaches, except that instead of wrapping the OF on discrete explicit frame representations, we
 61 apply the optical flow constraint on the exact gradient of the the NIR. Our method is thus not subject
 62 to any mismatch between training and test data. Furthermore, our approach can sample any number of
 63 frame in-between the observed frames due to the continuous nature of the representation. In addition
 64 to its application to VFI, we also show that constraining the gradient of the model also improves the
 65 ability of narrow MLPs to fit the signal, suggesting potential applications in NIR optimization and
 66 video compression.

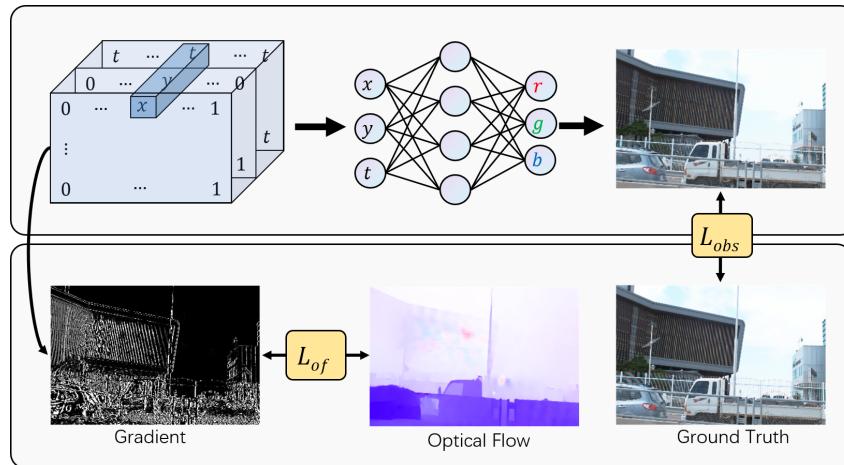


Figure 1: Illustration of our approach

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 videos as a continuous signal v mapping continuous spatial (x, y) and
 139 temporal (t) coordinates to RGB values:

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

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

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

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

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

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:

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

$$\text{s.t. } \mathcal{V}_{xyt} = v(x, y, t) \forall (x, y, t) \in \mathbb{N}^3$$

145 Following previous work on NIR (cite), we use the SIREN model (MLP with sine activation functions),
 146 The most straightforward way to solve Equation XXX is to optimize over the model parameters to fit
 147 the observations, using the following loss function

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

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

153 The optical flow represents the movement of brightness patterns in videos. For a given coordinate
 154 (x, y, t) in a video signal v , the optical is defined as the motion of this coordinate's brightness,
 155 The optical flow constraint equation states that for an infinitesimal lapse of time δt , the brightness



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 of a physical point perceived by the camera should remain constant. In other words, given the
 157 displacement $(\delta x, \delta y)$ of a physical point in the image coordinate system, the image brightness v
 158 should remain constant. This relationship is exact in the infinitesimal limit, as δt tends to zero, we
 159 have:

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

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

$$D(f, \theta, x, y, t) = \left(\frac{\delta f_\theta(x, y, t)}{\delta x}, \frac{\delta f_\theta(x, y, t)}{\delta y}, \frac{\delta f_\theta(x, y, t)}{\delta t} \right) \quad (6)$$

162 And the optical flow as:

$$F(x, y, t) = (u(x, y, t), v(x, y, t), 1) \quad (7)$$

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 λ 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 plays in our favor: As we regularize
 169 the true derivative of the signal representation, we do not assume constant derivatives of the signal on
 170 any interval. On the other hand, the optical flow we used was computed from discrete consecutive
 171 frames, and thus does not represent the true infinitesimal motion range. We discuss this limitation in
 172 Section XXX.

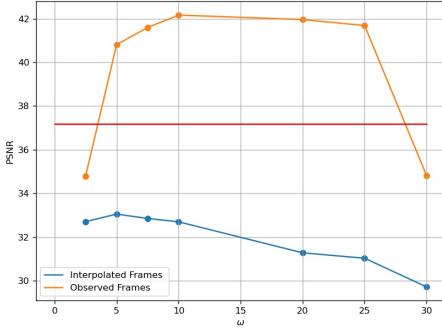


Figure 3: NEED TO SHOW TRAINING CURVE AND THRESHOLD ATTAINED BY OF

173 4 Experiments

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

180 Unless specified otherwise, all experiments are run with a SIREN of depth 6 and width 720. We use
 181 an omega of 25 and a lambda of 0.12. We optimize the models using the Adam optimizer using a
 182 cosine learning rate with maximum learning rate of 3.6e-5 during 15k epochs.

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

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

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

195 4.1 Optical Flow constraint and High Frequencies

196 Figure 2 illustrates the fact that applying the optical flow constraint smoothes the high-frequency
 197 variations of vanilla SIREN representation. We start by questioning whether the OF constraint does
 198 more than simply constraining the high frequency variations of the representation. To do so, we
 199 compare the results of vanilla SIREN representations geared towards different frequency and compare
 200 the best obtained results to OF-constrained representations. We constrain the SIREN frequency by
 201 varying their ω parameter, and report our comparison in Figure XXX.

202 While constraining the high frequency with low ω does improve the ability to interpolate intermediate
 203 frames, vanilla SIREN models remain well under the OF-constrained representations, confirming
 204 that the OF constraint provides more simply restricting the high temporal frequencies.

Table 1: XXX

	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	35.14/0.9397
Ours	36.52/0.9770	35.06/0.9441

Table 2: XXX

	ND Scene [9]
V-NF	23.30/0.7260
NSFF [10]	28.03/0.9250
CURE [11]	36.91/0.9843
Ours	29.22/0.9215

205 4.2 State of the art models

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

208 4.3 Ablation study

209 Next, we highlight the

210 4.4 Qualitative Analysis

211 4.5 Video fitting

212 5 Limitations

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

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

221 **Reliance on trained optical flow model.** SIREN models allow us to apply the optical flow on the
222 exact derivatives of the signal, thus bypassing the heuristics of classical approach without relying on
223 machine learning. The optical flow we use, however, is given by a trained ML model, which raises
224 two problems: it is subject to generalization error, and the flow is computed on discrete samples and
225 then subject to undesirable changes in illumination and occlusion. Future work will aim to bypass
226 our reliance on ML-based OF using proxy constraints on the exact derivatives.

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

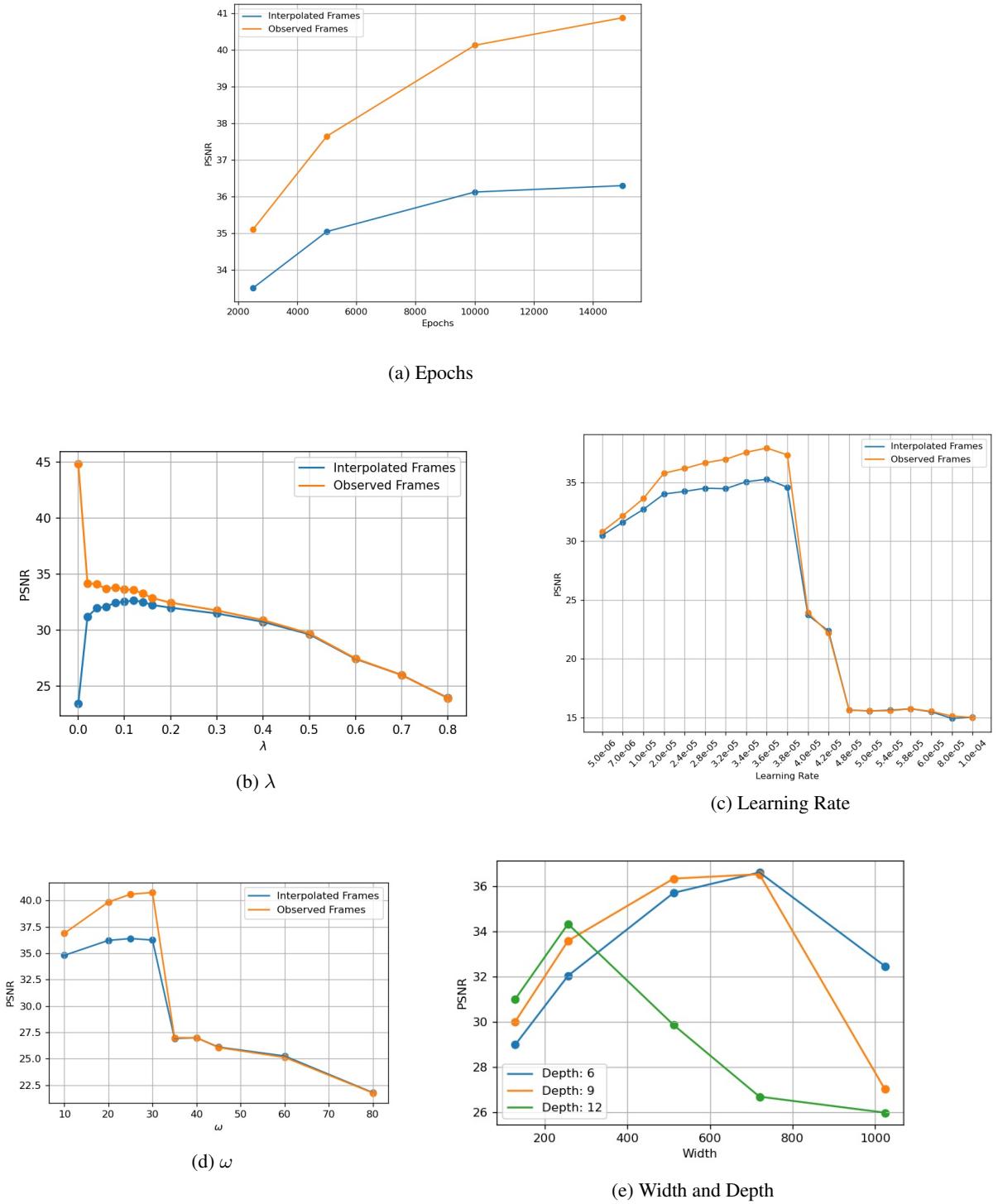


Figure 4: Ablation Experiments

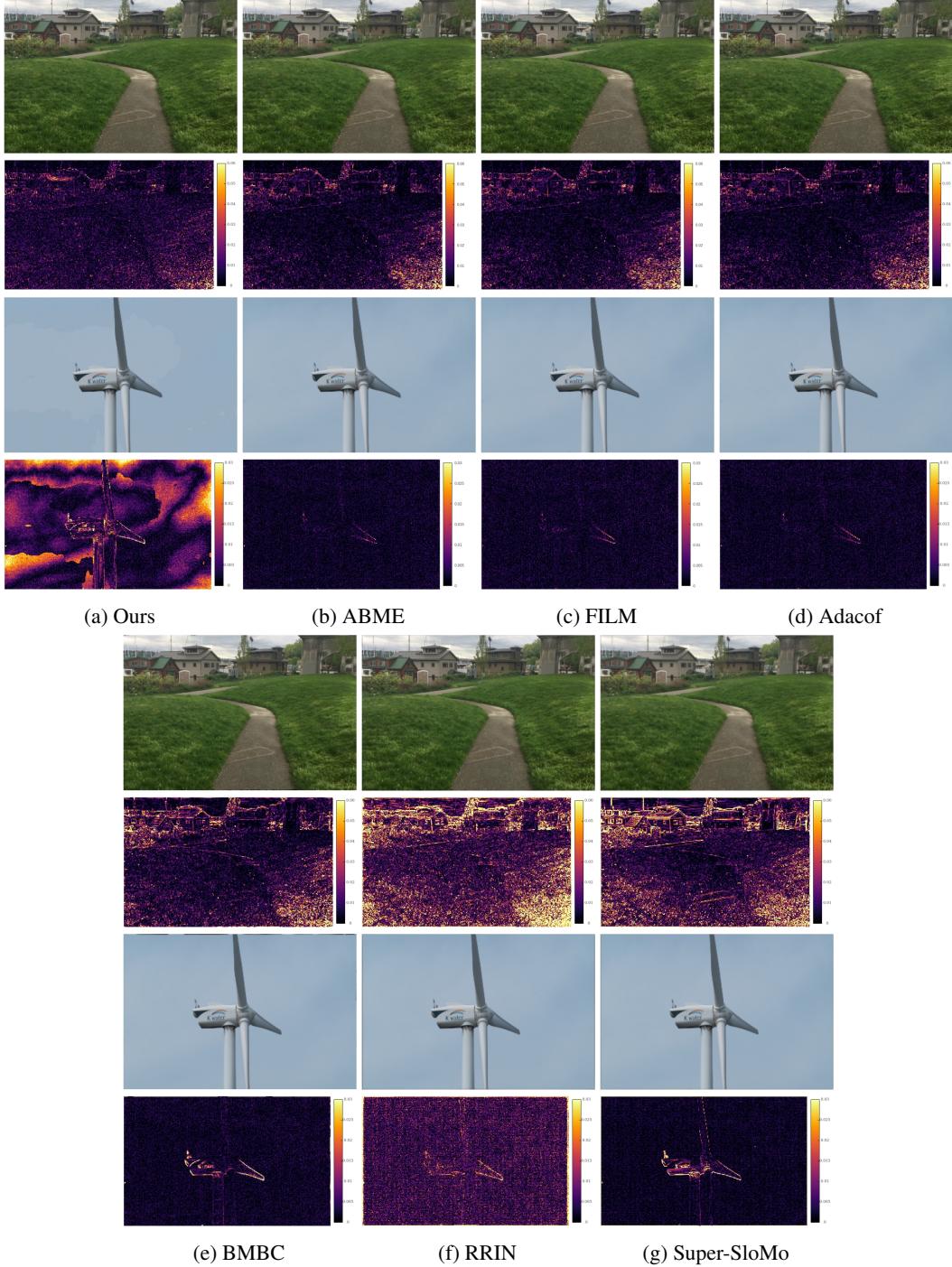


Figure 5: Small Motion Video Qualitative Analysis

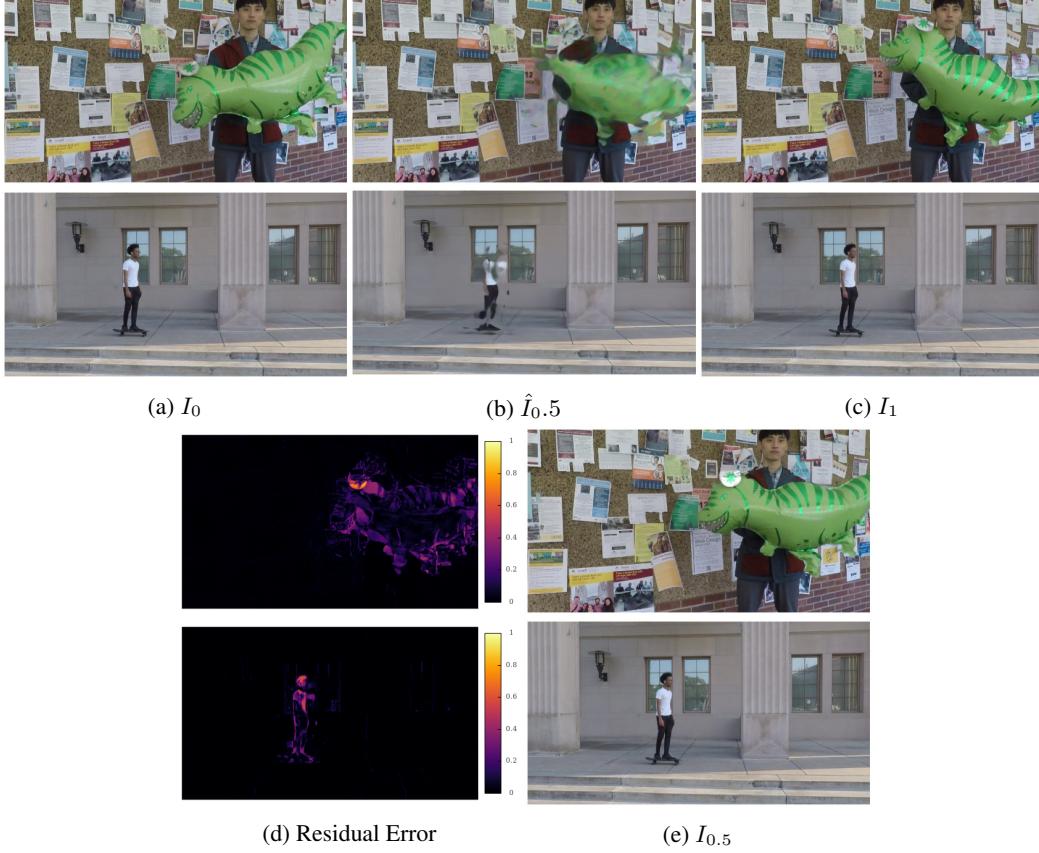


Figure 6: Large Motion Video Qualitative Analysis

232 6 Conclusion

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

236 References

- 237 [1] Alexander, J.A. & Mozer, M.C. (1995) Template-based algorithms for connectionist rule extraction. In
 238 G. Tesauro, D.S. Touretzky and T.K. Leen (eds.), *Advances in Neural Information Processing Systems 7*, pp.
 239 609–616. Cambridge, MA: MIT Press.
- 240 [2] Bower, J.M. & Beeman, D. (1995) *The Book of GENESIS: Exploring Realistic Neural Models with the
 241 GEneral NEural SImulation System*. New York: TELOS/Springer-Verlag.
- 242 [3] Hasselmo, M.E., Schnell, E. & Barkai, E. (1995) Dynamics of learning and recall at excitatory recurrent
 243 synapses and cholinergic modulation in rat hippocampal region CA3. *Journal of Neuroscience* **15**(7):5249-5262.

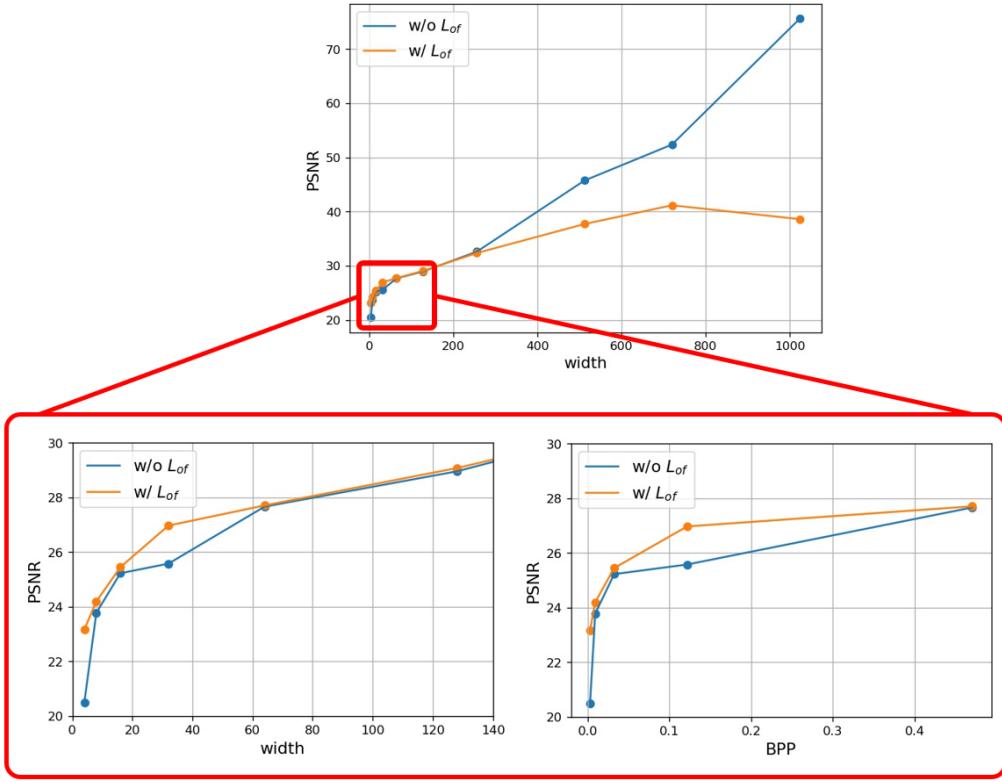


Figure 7: Our proposed optical flow regularization loss can help the lightweight SIREN model to fit the video better.

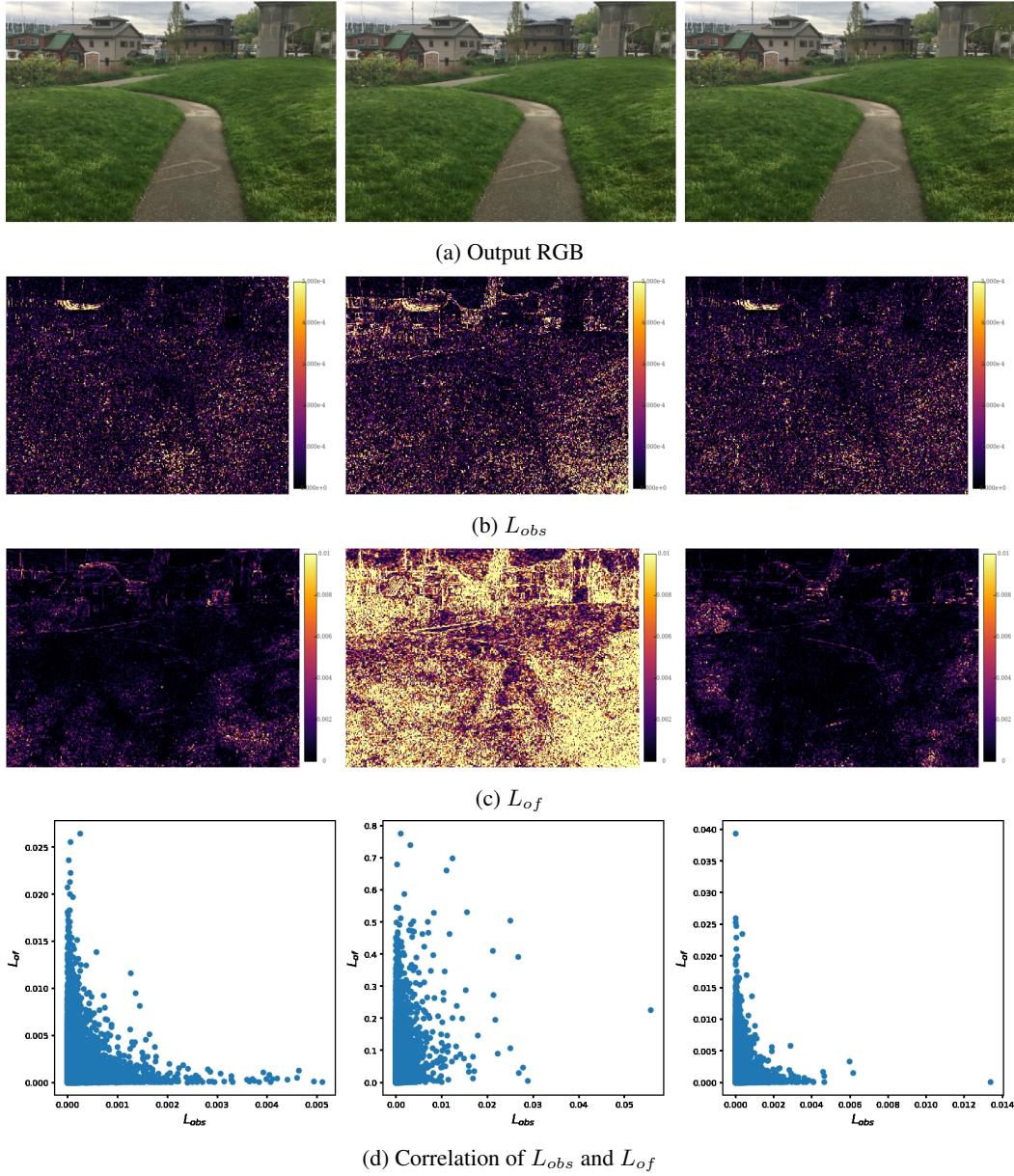


Figure 8: Figures from left to right are: Output Observed Frame I_0 , Output Interpolated Frame $I_{0.5}$, Output Observed Frame I_1