
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 (VFI) by explicitly constraining the derivatives of the INR to satisfy the optical flow constraint equation. We achieve state of the art VFI 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 better interpolate intermediate frames but also improves the ability of narrow networks to fit the observed frames, which suggests potential applications to video compression and INR optimization.

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, modern digital hardware is inherently discrete: digital sensors capture discrete observations of the world regularly sampled in time and space; computers store and process discrete representations of signals. In order to model continuous notions on discrete signal representations, classical methods have used different simplifying 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) methods, have led to the near ubiquitous use of ML in recent signal processing research. These methods 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 methods rely on the optical flow to guide the interpolation of pixel intensities from the pixel grid of observed frames onto the pixel grid of intermediate frames. Classical methods formulate assumptions such as constant movement or acceleration fields between consecutive frames [1] [2] [7]. The value of each pixel in the inferred intermediate frame is computed by first shifting the pixel intensities of observed frames along the optical flow directions before interpolating the shifted intensities onto the intermediate frame's pixel grid. Such approaches suffer from the following two main limitations:

- The optical flow is prone to errors due to occlusions, external illumination variations, etc.
- Assumptions of constant motion field or its derivatives do not often hold true in practice.

These limitations share a common root cause: discretization. Indeed, both the constant brightness assumption, from which is derived the optical flow, and assumptions of constant motion field used by

36 the interpolation process, only truly hold at the infinitesimal scale, for time deltas typically much
 37 smaller than those of practically used Frames Per Second (FPS).
 38 ML approaches [8] [12] [18] [17] [11] have instead proposed to learn the frame interpolation operator
 39 from large video collections, without formulating explicit assumption on the signal. While these
 40 approaches have achieved great success in terms of benchmark performance, they are prone to
 41 generalization errors caused by domain shifts. Indeed differences between the training set distribution
 42 (i.e. VFI benchmark videos) and the target video distribution may hinder the performance of ML
 43 models, e.g.; differences stemming from the range of motion, exposure time, FPS and blur [29].
 44 In the mean time, research on implicit representations seeks better discrete representations of con-
 45 tinuous signals. In recent years Implicit Neural Representations (INR), i.e. representing signals
 46 as Neural Networks (NN) have been shown to offer several competitive advantages over explicit
 47 representations, with notable early successes for 3D shape representations [16]. Of particular interest
 48 to us is the work of SIREN [24], in which the authors have shown that representing images using
 49 Multi Layer Perceptrons (MLP) with sine activation functions allowed for meaningful representations
 50 of the signal derivatives. Inspired by this work, we question whether SIREN may be used to guide
 51 the interpolation process of VFI by controlling the exact derivatives of the signal instead of the finite
 52 differences between consecutive discrete frames, thus avoiding the pitfalls of traditional methods
 53 due to discretization. We do so by constraining the derivatives of SIREN representations to satisfy
 54 the optical flow constraint equation, i.e., to be orthogonal to the video’s optical flow (which we
 55 compute using existing state-of-the-art OF models). We find that this method outperforms most
 56 existing machine learning-based approaches on small motion range benchmarks, without relying on
 57 machine learning for the interpolation operator. In this sense, our method is most similar to classical
 58 VIF approaches, except that instead of wrapping the OF on discrete explicit frame representations, we
 59 apply the optical flow constraint on the exact gradient of the INR. Our method is thus not subject to
 60 any mismatch between training and test data. Furthermore, our approach can sample any number of
 61 frame in-between the observed frames due to the continuous nature of the representation. In addition
 62 to its application to VFI, we also show that constraining the gradient of the model also improves the
 63 ability of narrow MLPs to fit the signal, suggesting potential applications in INR optimization and
 64 video compression.

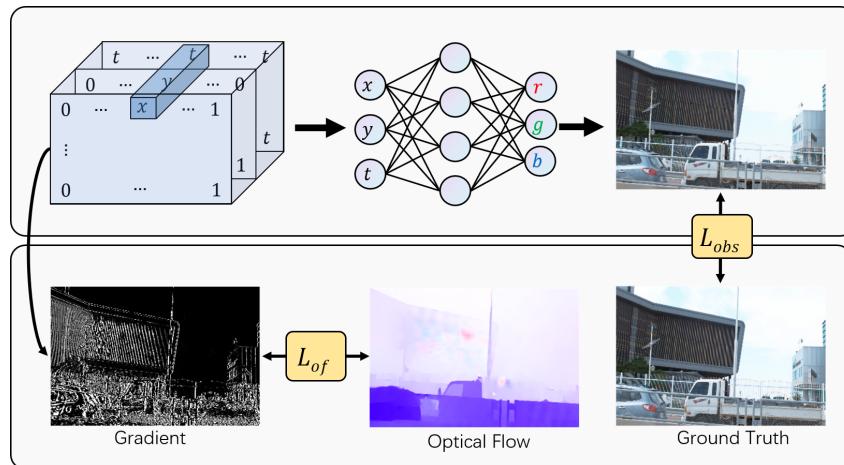


Figure 1: Illustration of our approach. We optimize SIREN to minimize the weighted sum of two losses: The observation loss measures the fit to the video frames, and the OF loss measures the orthogonality between the SIREN derivatives and the video’s optical flow.

65 To summarize the contributions of this work, we show that:

- 66
 - SIREN representations of videos can be constrained so as to satisfy the OF constraint in
 67 their exact derivatives.
 - Such representations reach state of the art VFI on limited motion ranges, without learning a
 69 residual flow nor interpolation operator.

- 70 • The OF constraint not only allows SIREN to generate intermediate frames, but also improve
71 the ability of narrow SIREN to fit observed frames.

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

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

79 Given these limitations, the aim of this paper is not to provide a standalone production ready VFI
80 system. Instead, we aim to present actionable insights on a simple method that can be either built
81 upon or integrated to existing models. The remainder of this paper is organized as follows: We
82 briefly present some related work in Section 2, the detail of our method in Section 3, and design
83 several experiments to highlight the merits of our approach in Section 4. Finally, we discuss current
84 limitations and present potential ways to address them in Section 5, before concluding in Section 6.

85 2 Related Work

86 **Implicit Neural Representations** have met early success in shape representation and 3D rendering
87 [19] [15] [16]. Since then, a number of works have attempted to apply INR to different signals
88 including audio [24] [10], images [5] [6], videos [3] [22] [21], medical imaging and climate data [6].
89 In [24] the authors have shown that MLP with sine activations could fit representations of images
90 with meaningful representations of their gradient, and that such models could be optimized to satisfy
91 constraints on their gradients. Combined, these two findings have motivated our idea to apply the
92 optical flow constraint to the gradient of SIREN representations of videos. A series of recent works
93 have applied INR to video compression [28] [3], with some works [3] even reporting higher PSNR
94 than practical codecs on high compression rates. Although closely related to video compression,
95 we differ from these works as we focus on VFI. Most related to ours is the concurrent work by
96 Shangguan et al. [22], which also uses INR for VFI. Their approach, CURE, differs from ours in
97 scope: they propose to learn a prior on the INR, while we only focus on leveraging INR to guide the
98 interpolation process using a given optical flow.

99 **Video Frame Interpolation** research has largely relied on optical flow to guide the video frame
100 interpolation process [1] [2] [7]. Most works have assumed uniform optical flow between consecutive
101 frames so as to linearly interpolate intermediate frames along the optical flow directions. One
102 exception is the work of [26], in which the authors propose to take into account acceleration to
103 perform the interpolation, leading to quadratic interpolation. Our work only constrains the first
104 derivatives of the signal. We differ from classical works in that we apply the OF to the exact
105 representation derivatives, so that we do not need to assume constancy of signal derivatives on
106 any time interval. Recent OF-based VFI leverages deep learning for optical flow estimation and
107 interpolation. Super-SloMo [8] is an important study of such methods. The authors use a deep
108 learning model to predict the forward and backward flows of intermediate frames, and warp the two
109 surrounding frames to obtain the intermediate frames. RRIN [12] uses residual learning to optimize
110 the performance of [8] at the motion estimation bound. AMBE [18], a current state-of-the-art VFI
111 method, proposes an asymmetric motion estimation method based on [17], which enhances the quality
112 of interpolated frames by loosening the linear motion constraint. Kernel-based approaches such as
113 AdaCof [11] avoid explicit separation of motion estimation and wrapping stages and instead directly
114 interpolate intermediate frames from consecutive observed ones.

115 3 Method

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

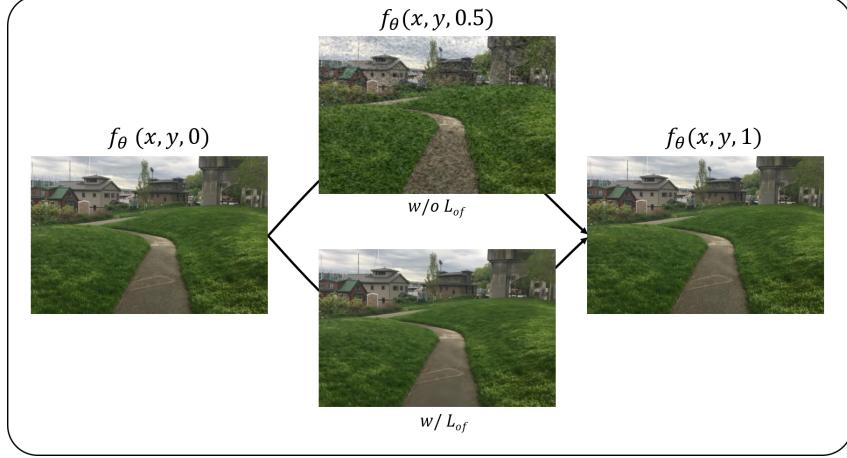


Figure 2: Illustration of INR frame interpolation with and without optical flow regularization. Without regularization (middle top), intermediate frames show unnatural high-frequency variations. Regularizing the INR to satisfy the optical flow constraint equation result in nicely interpolated frames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

$$\begin{aligned} v : (x, y, t) &\rightarrow (R, G, B) \\ v : \mathbb{R}^3 &\rightarrow \mathbb{R}^3 \end{aligned} \tag{1}$$

118 Our goal is to find a continuous function f_θ , parameterized by $\theta \in \Theta$, with minimum distance d to
119 the ground-truth signal:

$$\begin{aligned} f_\theta : (x, y, t) &\rightarrow (R, G, B) \\ \text{s.t. } \theta = \min_{\Theta} \int \int \int d(f_\theta(x, y, t), v(x, y, t)) dx dy dt \end{aligned} \tag{2}$$

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

$$\begin{aligned} \mathcal{V} &\in \mathbb{R}^{T \times H \times W \times 3} \\ \text{s.t. } \mathcal{V}_{xyt} &= v(x, y, t) \end{aligned} \tag{3}$$

123 where T represents the number of frames in the video, and $H \times W$ the spatial resolution. We use
124 SIREN as parameterized function f_θ . The most straightforward way to approximate Equation 2 is to
125 optimize the model parameters so as to fit the video frames, using the following loss function (we
126 refer to as the observation loss):

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

127 However, we found that optimizing the INR to only minimize this observation loss leads to overfitting
128 the observation with high temporal frequencies: the intra-frame signal, which we aim to correctly
129 recover, shows important deviations from the observed frames, as illustrated in Figure 2. This
130 observation has lead us to consider fitting not only the signal itself, but to also constrain its derivatives.
131 In particular, we regularize the model so as to respect the optical flow constraint. The optical flow
132 constraint equation states that for an infinitesimal lapse of time δt , the brightness of physical points
133 perceived by a camera at arbitrary coordinates (x, y, t) should remain constant. In other words, given
134 the displacement $(\delta x, \delta y)$ of a physical point in the image coordinate system, the image brightness v

135 should remain constant:

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

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

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

138 First order Taylor expansion of Equation 6 gives the following

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

139 which holds exactly in the limit of infinitesimal δt . We constrain the SIREN derivatives to obey the
140 constraint of Equation 7. Denoting the derivatives of the SIREN as:

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

141 we can now define the optical flow regularization loss

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

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

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

145 where λ is a hyper-parameter taking values between 0 and 1, whose impact we investigate in the
146 following section. The exactness of the optical flow constraint at the infinitesimal scale plays in our
147 favor: As we regularize the true derivative of the signal representation, we do not assume constant
148 motion on any time interval. We believe this to be the main factor behind our positive results. On the
149 other hand, the optical flow we use was estimated from discrete consecutive frames, and thus does
150 not represent the true infinitesimal motion field but an estimation of finite differences. We discuss
151 potential alternatives in Section 5.

152 4 Experiments

153 Following previous works, we use the Adobe[25], X4K[23] and ND Scene[27] datasets as benchmark
154 to compare our method to state-of-the-art models. We use every two frames of each video as
155 observations, and evaluate the ability of SIREN to interpolate on every other (intermediate) frame.
156 For the Adobe dataset, we evaluate our method on the eight videos test split proposed in previous
157 works [8]. We run all additional experiments on the 720p240fps1.mov video of the Adobe dataset
158 (illustrated in Figure 2). Due to the time-consuming operation of optimizing SIREN representations,
159 we optimize and evaluate all models on a 240×360 pixel resolution, and we restrict the Adobe
160 dataset videos to their first 40 frames. Unless specified otherwise, we use the following default
161 parameters: SIREN model with depth 9, width 512 and an ω of 30. We optimize the models with
162 the Adam optimizer using a cosine learning rate with maximum learning rate of 10^{-5} during 5000
163 epochs. We use $\lambda = 0.12$ for the loss function. We compute the optical flow of videos in original
164 resolution using the GMA [9] OF model.

165 In Section 4.1, we start by highlighting a trade-off akin to underfitting vs overfitting of the signal
 166 high frequencies in vanilla SIREN representations. We show that OF-regularized SIREN outperform
 167 the best performing vanilla SIREN, showing that the impact of our proposed OF regularization goes
 168 beyond high frequency regularization. In Section 4.2, we quantitatively compare our method to state
 169 of the art models on standard datasets. We show that our method achieves state-of-the-art results on
 170 videos with limited motion ranges, but underperforms recent methods for videos with large motion
 171 ranges. We present an ablation study in Section 4.3, providing insights and appropriate settings for
 172 the main hyper-parameters, and a qualitative analysis of our results in Section 4.4. Finally, Section
 173 4.5 presents a surprising and counter-intuitive result: we show that our OF loss helps SIREN converge
 174 to higher PSNR on the observed frames, opening new potential perspectives for INR optimization
 175 and video compressions.

176 **4.1 Optical flow constraint and signal frequencies**

177 Figure 2 illustrates the fact that the OF constraint smooths out high-frequency noise in the interpolated
 178 frames of vanilla SIREN representations. Healthy skepticism leads us to question whether the
 179 impact of the OF constraint is limited to dampening high frequency components of vanilla SIREN
 180 representations. To do so, we analyze the representations of vanilla SIREN geared towards different
 181 frequency ranges, and compare them to OF-constrained SIREN representations. We constrain the
 182 vanilla SIREN frequencies by varying their ω parameter, and report our comparison in Figure 3, with
 183 low ω values corresponding to lower frequency ranges.

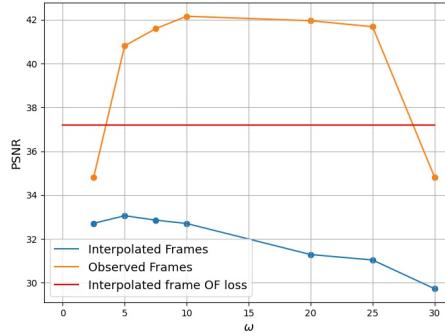


Figure 3: Evolution of the PSNR of observed and interpolated frames with ω without OF loss. Limiting the high frequency fit alone does not reach the same interpolation accuracy as the OF loss.

184 Constraining the frequency range of vanilla SIREN with ω down to 5 degrades the fit to observed
 185 frames but improves interpolation. This suggests that ω behaves similarly to a regularization parameter
 186 by controlling a regime of overfitting to the observed frames high frequencies (for high ω values),
 187 versus underfitting (for low ω values). Figure 3 further shows that OF-constrained SIREN achieve far
 188 higher interpolation PSNR than the best performing vanilla SIREN, confirming that the OF constraint
 189 impact goes beyond dampening of the high frequency noise. Note that we did not vary the ω of the OF
 190 constrained SIREN in this figure in order to better illustrate our point, the red line represents results
 191 for the best performing ω . The impact of the ω parameter on OF-constrained SIREN is illustrated
 192 separately in Figure 4d.

193 **4.2 State of the art models**

194 Table 1 quantitatively compares the results of our method to state-of the art VFI models on different
 195 datasets. Despite its simplicity, and without any training data, our method outperforms most existing
 196 models on limited motion ranges (Table 1). However, as illustrated in Figure 6, it falls short of
 197 state-of-the-art methods on the more complex ND Scene benchmark due to larger motion ranges. We
 198 provide further comparison in the qualitative analysis of Section 4.4 and Section 5 discusses possible
 199 ways forward to bridging the gap performance on large motion datasets.

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

	(a) Limited motion range	(b) Large motion range	
	Adobe-240FPS	X4K	ND Scene
Super-SloMo [8]	27.77 / 0.886	27.38 / 0.852	V-NF [16] 23.30 / 0.726
RRIN [12]	32.37 / 0.962	30.70 / 0.927	NSFF [13] 28.03 / 0.925
BMBC [17]	27.83 / 0.917	27.42 / 0.858	CURE [22] 36.91 / 0.984
AdaCof [11]	35.50 / 0.968	34.61 / 0.921	Ours 29.22 / 0.921
ABME [18]	35.28 / 0.966	34.30 / 0.919	
FILM [20]	35.97 / 0.971	35.14 / 0.939	
Ours	36.52 / 0.977	35.06 / 0.944	

200 4.3 Ablation study

201 Figure 4 summarizes the impact of the main parameters of our method. In (a) we observe a trade-off
 202 between the observed and interpolated frames quality in the low λ ranges. The quality of interpolated
 203 frames peaks at $\lambda = 0.12$, beyond which point the interpolated frames quality is limited by the
 204 quality of the fit to the observed frames, in a similar way to the classical overfitting/underfitting
 205 trade-off. However, it should be noted that this trade-off differs widely depending on the SIREN’s
 206 width. Indeed, as we will show in Section 4.5, the OF constraints actually improves the fit to
 207 observed frames for narrow models. In (b) and (c) we observe that both higher learning rates and
 208 longer fitting times improve both observed and interpolated frames. The learning rate is limited in
 209 amplitude by instabilities of the optimization procedure, while the fitting time is limited by practical
 210 time constraints. Large ω (d) also improve the accuracy up to 30, after which instabilities in the
 211 optimization see the accuracy drop abruptly. Width and depth (e) show interesting co-dependencies:
 212 Increasing width improves interpolation up to a peak after which it degrades. The peak width gets
 213 smaller with increasing depth.

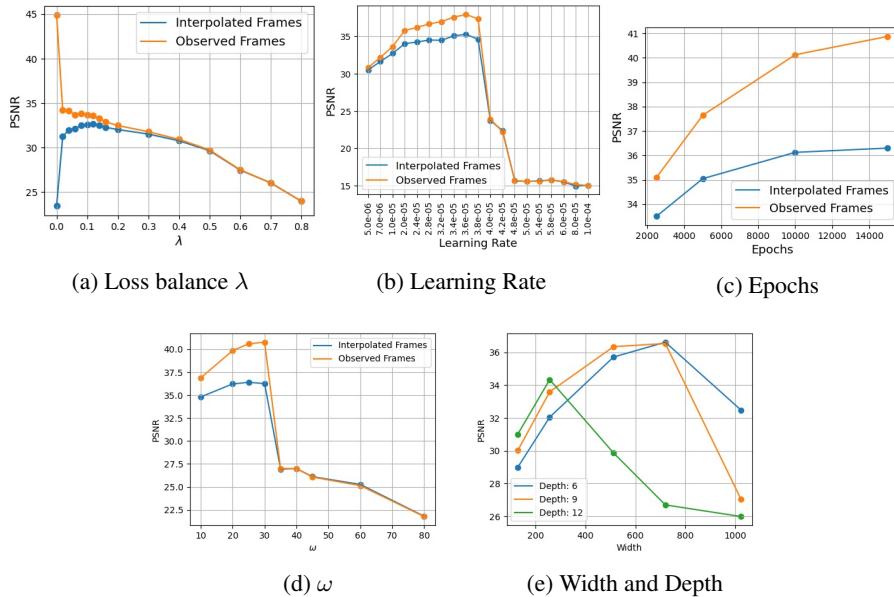


Figure 4: Impact of our method’s main parameters. Plots from (a) to (d) show both the observed and interpolated frames PSNR while plot (e) only shows the interpolated frames PSNR.

214 Based on these experiments, our final results, as reported in Table 1 were computed with a SIREN
 215 model with depth 6, width 720 and $\omega = 25$. We used $\lambda = 0.12$ for the loss, and optimized using
 216 Adam with a maximum learning rate of 3.6e-5 during 15k epochs.

217 **4.4 Qualitative analysis**

218 Figures 5 and 6 provide a qualitative illustration to the results presented in Section 4.2. The upper
 219 frame in Figure 5 shows that our method tends to outperform other methods on videos with limited
 220 motion range. In particular it seems to better catch high spatial frequency regions (grass, sharp edges
 221 of the building). In contrast, large motion as illustrated in Figure 6 shows ghosting effects that the OF
 222 regularization is not able to address.

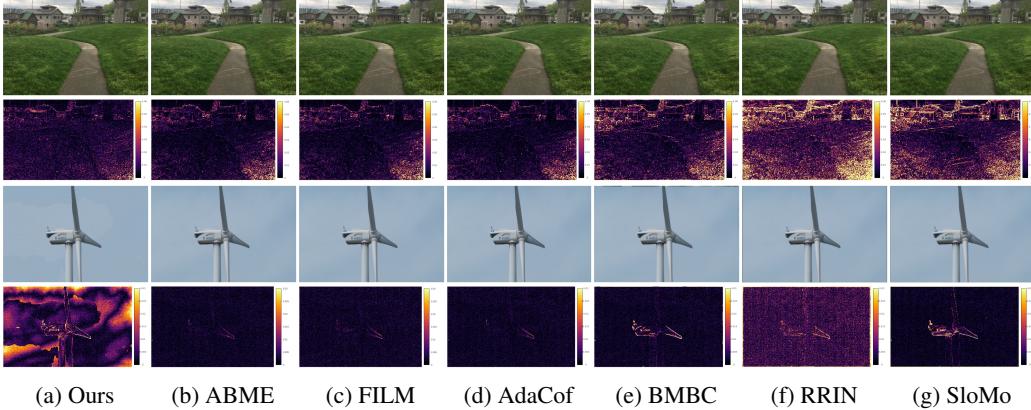


Figure 5: Small Motion Video Qualitative Analysis. The interpolated frame results are shown above their residual heat map. The upper frames show a successfully interpolated frames, the lower one shows a rare failure case.

223 The lower part of Figure 5 shows a rare failure case of our method on limited motion ranges: some
 224 artificial stain-like patterns appear in the sky background, suggesting additional care may be needed
 225 especially in low frequency regions. Despite this rare exception, the overall quality of interpolation
 226 on limited motion range videos performs on par with the best existing methods.

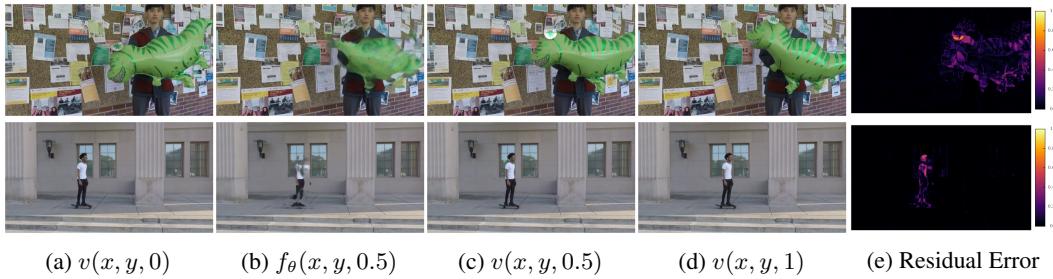


Figure 6: Large Motion Video Qualitative Analysis

227 **4.5 Video fitting**

228 Figure 7 shows an unexpected side-effect of the OF regularization observed for narrow networks. As
 229 \mathcal{L}_{obs} explicitly maximizes the PSNR of observed frames, we expected the addition of the \mathcal{L}_{of} term
 230 to negatively impact the PSNR of observed frames, especially for capacity-limited SIREN which
 231 should have to compromise between satisfying both loss terms. It turns out that, for width up to 50,
 232 optimizing the SIREN with the additional OF constraint actually improves the fit to observed frames.

233 Although a complete investigation of this phenomenon is out of the scope of this work, we highlight
 234 how this observation may prove interesting for future works: From a practical standpoint, improving
 235 the fit of low-capacity INR is the key challenge towards practical INR video compression. It remains
 236 to be seen whether this phenomenon can be replicated on more practical architectures (i.e. [3]).
 237 From a theoretical standpoint, increasing width has been shown to help optimization by alleviating
 238 second order effects [14] and guarantee convergence of gradient descent to global minima [4]. As
 239 the understanding of gradient descent dynamics in the high curvature low width setting is currently

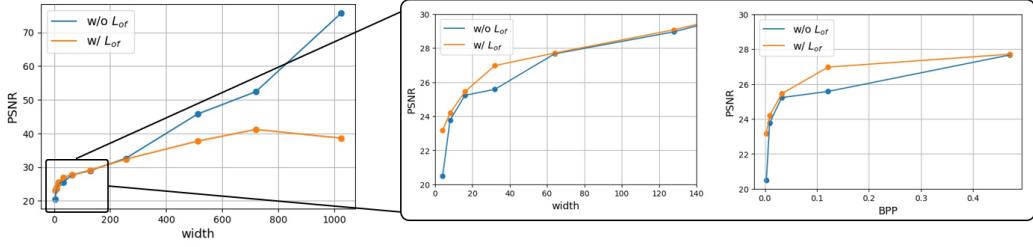


Figure 7: Evolution of the **observed** frames PSNR with depth, with and without OF regularization. Left: Trend from very narrow to very wide models. Right: Zoom on the low width regime with the x axis expressed either in number of neurons or corresponding Bits Per Pixel measure.

240 an elusive question, understanding how the OF constraint helps optimization may provide useful
241 insights into gradient descent dynamics in narrow models.

242 5 Current Limitations and Future Work

243 While our method does reach state of the art interpolation results on limited motion ranges, this work
244 is not meant to deliver a production ready VFI system, which would require the ability to interpolate
245 high resolution and large motion range videos. Instead, we aim to provide actionable insights for
246 future works on both VFI and INR to integrate and build upon. Towards that goal, we discuss below
247 what we see as the three main limitations of our method in its current form, and possible ways to
248 address these limitations.

249 **Slow optimization process.** Fitting 20 frames of a video at 240×360 resolution currently takes 15
250 hours on a $4 \times 2080\text{Ti}$ GPU using Pytorch. This computation time is an important drawback as it
251 limits our ability to process full resolution video, as well as to explore different hyper parameters and
252 variations of the method within realistic times. We expect advances in INR optimization to be very
253 beneficial to this line of research. Given recent successes of INR in signal compression [28][5] [6]
254 [19] [15] [3], we hopefully expect to see such development in the near future.

255 **Reliance on trained optical flow model.** SIREN models allow us to apply the optical flow on the
256 exact derivatives of the signal, bypassing the heuristics of classical methods without relying on
257 machine learning. The optical flow we use, however, is given by a ML model trained on discrete
258 representations, which raises two problems: it is subject to generalization errors, and is subject to
259 finite difference errors such as occlusions. Bypassing this reliance on ML-based OF using alternative
260 constraints on the exact derivatives of the representation is another interesting way forward.

261 **Inability to interpolate large motion range videos.** In its current form, we only apply the optical
262 flow constraint on the observed frames of the video. This has proven sufficient to reach state-of-the
263 art on limited motion ranges, but is not sufficient for large motions. A promising axis of improvement
264 would be to apply additional constraints to the interpolated frames (e.g. for intra-frame time indices
265 $t = 0.5$). Possible regularization methods may include constraints on intra-frame texture, as proposed
266 in recent works [20], or interpolated optical flows, which may prevent the ghosting effects illustrated
267 in Figure 6.

268 6 Conclusion

269 In this paper, we have shown that SIREN representations of videos can be constrained to satisfy
270 the OF constraint equation in their exact derivatives. We have seen that OF-constrained SIREN
271 reach state of the art VFI on limited motion ranges, without relying on ML based residual flow
272 and interpolation. We have also shown that the OF constraint not only allows SIREN to generate
273 intermediate frames, but can also improve the ability of narrow SIREN to fit observed frames. We
274 have discussed the limitations of our approach in its current form and outlined potentially impactful
275 way forwards for future research.

276 **References**

- 277 [1] S. Baker, D. Scharstein, J. Lewis, S. Roth, M. J. Black, and R. Szeliski. A database and
278 evaluation methodology for optical flow. *International journal of computer vision*, 92(1):1–31,
279 2011.
- 280 [2] J. L. Barron, D. J. Fleet, and S. S. Beauchemin. Performance of optical flow techniques.
281 *International journal of computer vision*, 12(1):43–77, 1994.
- 282 [3] H. Chen, B. He, H. Wang, Y. Ren, S. N. Lim, and A. Shrivastava. Nerv: Neural representations
283 for videos. *Advances in Neural Information Processing Systems*, 34, 2021.
- 284 [4] S. S. Du, X. Zhai, B. Poczos, and A. Singh. Gradient descent provably optimizes over-
285 parameterized neural networks. *arXiv preprint arXiv:1810.02054*, 2018.
- 286 [5] E. Dupont, A. Goliński, M. Alizadeh, Y. W. Teh, and A. Doucet. Coin: Compression with
287 implicit neural representations. *arXiv preprint arXiv:2103.03123*, 2021.
- 288 [6] E. Dupont, H. Loya, M. Alizadeh, A. Goliński, Y. W. Teh, and A. Doucet. Coin++: Data
289 agnostic neural compression. *arXiv preprint arXiv:2201.12904*, 2022.
- 290 [7] E. Herbst, S. Seitz, and S. Baker. Occlusion reasoning for temporal interpolation using optical
291 flow. *Department of Computer Science and Engineering, University of Washington, Tech. Rep.*
292 *UW-CSE-09-08-01*, 2009.
- 293 [8] H. Jiang, D. Sun, V. Jampani, M.-H. Yang, E. Learned-Miller, and J. Kautz. Super slomo: High
294 quality estimation of multiple intermediate frames for video interpolation. In *Proceedings of the*
295 *IEEE conference on computer vision and pattern recognition*, pages 9000–9008, 2018.
- 296 [9] S. Jiang, D. Campbell, Y. Lu, H. Li, and R. Hartley. Learning to estimate hidden motions
297 with global motion aggregation. In *Proceedings of the IEEE/CVF International Conference on*
298 *Computer Vision*, pages 9772–9781, 2021.
- 299 [10] J. Kim, Y. Lee, S. Hong, and J. Ok. Learning continuous representation of audio for arbitrary
300 scale super resolution. In *ICASSP 2022-2022 IEEE International Conference on Acoustics,*
301 *Speech and Signal Processing (ICASSP)*, pages 3703–3707. IEEE, 2022.
- 302 [11] H. Lee, T. Kim, T.-y. Chung, D. Pak, Y. Ban, and S. Lee. Adacof: Adaptive collaboration of
303 flows for video frame interpolation. In *Proceedings of the IEEE/CVF Conference on Computer*
304 *Vision and Pattern Recognition*, pages 5316–5325, 2020.
- 305 [12] H. Li, Y. Yuan, and Q. Wang. Video frame interpolation via residue refinement. In *ICASSP 2020-*
306 *2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*,
307 pages 2613–2617. IEEE, 2020.
- 308 [13] Z. Li, S. Niklaus, N. Snavely, and O. Wang. Neural scene flow fields for space-time view
309 synthesis of dynamic scenes. In *Proceedings of the IEEE/CVF Conference on Computer Vision*
310 *and Pattern Recognition*, pages 6498–6508, 2021.
- 311 [14] C. Liu, L. Zhu, and M. Belkin. On the linearity of large non-linear models: when and why
312 the tangent kernel is constant. *Advances in Neural Information Processing Systems*, 33:15954–
313 15964, 2020.
- 314 [15] L. Mescheder, M. Oechsle, M. Niemeyer, S. Nowozin, and A. Geiger. Occupancy networks:
315 Learning 3d reconstruction in function space. In *Proceedings of the IEEE/CVF Conference on*
316 *Computer Vision and Pattern Recognition*, pages 4460–4470, 2019.
- 317 [16] B. Mildenhall, P. P. Srinivasan, M. Tancik, J. T. Barron, R. Ramamoorthi, and R. Ng. Nerf:
318 Representing scenes as neural radiance fields for view synthesis. In *European conference on*
319 *computer vision*, pages 405–421. Springer, 2020.
- 320 [17] J. Park, K. Ko, C. Lee, and C.-S. Kim. Bmbc: Bilateral motion estimation with bilateral cost
321 volume for video interpolation. In *European Conference on Computer Vision*, pages 109–125.
322 Springer, 2020.

- 323 [18] J. Park, C. Lee, and C.-S. Kim. Asymmetric bilateral motion estimation for video frame
324 interpolation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*,
325 pages 14539–14548, 2021.
- 326 [19] J. J. Park, P. Florence, J. Straub, R. Newcombe, and S. Lovegrove. Deepsdf: Learning continuous
327 signed distance functions for shape representation. In *Proceedings of the IEEE/CVF Conference*
328 *on Computer Vision and Pattern Recognition*, pages 165–174, 2019.
- 329 [20] F. Reda, J. Kontkanen, E. Tabellion, D. Sun, C. Pantofaru, and B. Curless. Film: Frame
330 interpolation for large motion. *arXiv preprint arXiv:2202.04901*, 2022.
- 331 [21] D. Rho, J. Cho, J. H. Ko, and E. Park. Neural residual flow fields for efficient video representa-
332 tions. *arXiv preprint arXiv:2201.04329*, 2022.
- 333 [22] W. Shangguan, Y. Sun, W. Gan, and U. S. Kamilov. Learning cross-video neural representations
334 for high-quality frame interpolation. *arXiv preprint arXiv:2203.00137*, 2022.
- 335 [23] H. Sim, J. Oh, and M. Kim. Xvfi: Extreme video frame interpolation. In *Proceedings of the*
336 *IEEE/CVF International Conference on Computer Vision*, pages 14489–14498, 2021.
- 337 [24] V. Sitzmann, J. Martel, A. Bergman, D. Lindell, and G. Wetzstein. Implicit neural representations
338 with periodic activation functions. *Advances in Neural Information Processing Systems*, 33:7462–
339 7473, 2020.
- 340 [25] S. Su, M. Delbracio, J. Wang, G. Sapiro, W. Heidrich, and O. Wang. Deep video deblurring for
341 hand-held cameras. In *Proceedings of the IEEE Conference on Computer Vision and Pattern*
342 *Recognition*, pages 1279–1288, 2017.
- 343 [26] X. Xu, L. Siyao, W. Sun, Q. Yin, and M.-H. Yang. Quadratic video interpolation. *Advances in*
344 *Neural Information Processing Systems*, 32, 2019.
- 345 [27] J. S. Yoon, K. Kim, O. Gallo, H. S. Park, and J. Kautz. Novel view synthesis of dynamic scenes
346 with globally coherent depths from a monocular camera. In *Proceedings of the IEEE/CVF*
347 *Conference on Computer Vision and Pattern Recognition*, pages 5336–5345, 2020.
- 348 [28] Y. Zhang, T. van Rozendaal, J. Brehmer, M. Nagel, and T. Cohen. Implicit neural video
349 compression. *arXiv preprint arXiv:2112.11312*, 2021.
- 350 [29] Y. Zhang, C. Wang, and D. Tao. Video frame interpolation without temporal priors. *Advances*
351 *in Neural Information Processing Systems*, 33:13308–13318, 2020.