

Diagnosing and Preventing Instabilities in Recurrent Video Processing

Thomas Tanay, Aivar Sootla, Matteo Maggioni, Puneet K. Dokania, Philip Torr, Aleš Leonardis and Gregory Slabaugh

Abstract—Recurrent models are becoming a popular choice for video enhancement tasks such as video denoising. In this work, we focus on their stability as dynamical systems and show that they tend to fail catastrophically at inference time on long video sequences. To address this issue, we (1) introduce a diagnostic tool which produces adversarial input sequences optimized to trigger instabilities and that can be interpreted as visualizations of spatio-temporal receptive fields, and (2) propose two approaches to enforce the stability of a model: constraining the spectral norm or constraining the stable rank of its convolutional layers. We then introduce *Stable Rank Normalization of the Layers (SRNL)*, a new algorithm that enforces these constraints, and verify experimentally that it successfully results in stable recurrent video processing.



1 INTRODUCTION

LOW-LEVEL computer vision problems such as denoising, demosaicing or super-resolution can be formalised as inverse problems and approached with modern machine learning techniques: a degraded input is processed by a convolutional neural network (CNN) trained in a supervised way to produce a restored output. The input is typically a single frame [1], [2], [3], [4], [5], [6]—but significantly better results can be obtained by leveraging the temporal redundancy of sequential images [7], [8], [9], [10], [11], [12]. There are two main categories of video processing CNNs. *Feedforward models* operate in a sliding-window fashion and process multiple frames jointly to produce a current output. *Recurrent models* operate in a frame-by-frame fashion but have the ability to store information internally through feedback loops. Recurrent processing is appealing because it reuses information efficiently, potentially over a large number of frames. At the same time, Recurrent Neural Networks (RNNs) are dynamical systems that can exhibit complex and even chaotic behaviors [13]. In the context of sequence modelling for language or sound understanding for instance, RNNs are known to suffer from vanishing and exploding gradient issues at training time [14], and to be vulnerable to instabilities through positive feedback at inference time [15].

1.1 Motivation

In the context of video processing, recurrent CNNs are typically trained on short video sequences of up to 10 frames [10], [11], [12]. *At inference time*, they can be applied to sequences of arbitrary length—but it has been observed before that they sometimes suffer from instabilities on long video sequences [10]. To illustrate this phenomenon, we retrain the video denoiser of [10] on the Vimeo-90k dataset [16]

and we evaluate it on 4 sequences of 1600 frames (i.e. roughly one minute of video at 24 fps) downloaded from vimeo.com. We plot the performance of the model measured by the peak-signal-to-noise-ratio (PSNR) as a function of the frame number and we observe instabilities on all 4 sequences: the PSNR plunges suddenly and permanently at some unpredictable time in the sequence. Visually, these instabilities correspond to the formation of colorful artifacts at random locations, growing locally until the entire output frame is covered. Numerically, they correspond to diverging or saturating pixel values (see Figure 1).

Our working hypothesis is that recurrent connections create positive feedback loops prone to this type of divergent behaviour. As a proof of concept, we consider a backbone architecture made of five convolutional layers interleaved with ReLU non-linearities, and we augment it with various strategies for temporal processing. We consider single-frame and multi-frame inputs, and four types of temporal connections inspired from existing video processing works: feature-shifting [17], [18], feature-recurrence [10], [17], [19], frame-recurrence [11], [20] and recurrent latent space propagation (RLSP) [12] (see Figure 2a). We then initialize all the models randomly and feed them with random inputs. We see that feedforward architectures (single-frame, multi-frame, feature-shifting) produce stable outputs while recurrent architectures (feature-recurrence, frame-recurrence, RLSP) produce outputs that diverge (see Figure 2b). Note that feature-shifting is non-recurrent since, while information is temporally delayed, it is not fed back to the same processing block multiple times.

The instabilities described above are a serious concern for the deployment of recurrent video processing models in the real world. If catastrophic failures occur at unpredictable times in the sequence, affected models cannot be trusted, even on relatively short sequences. In fact, determining whether a given model is affected or not is already problematic: there is no guarantee *a priori* that a model that behaves in a stable manner on a set of video sequences will remain stable on all video sequences.

- T. Tanay, A. Sootla, M. Maggioni and A. Leonardis are with Huawei Technologies Ltd, Noah's Ark Lab. E-mail: thomas.tanay@huawei.com
- G. Slabaugh is with the Queen Mary University of London (work done while at Huawei). P. Torr and P. K. Dokania are with the Department of Engineering Science, University of Oxford.

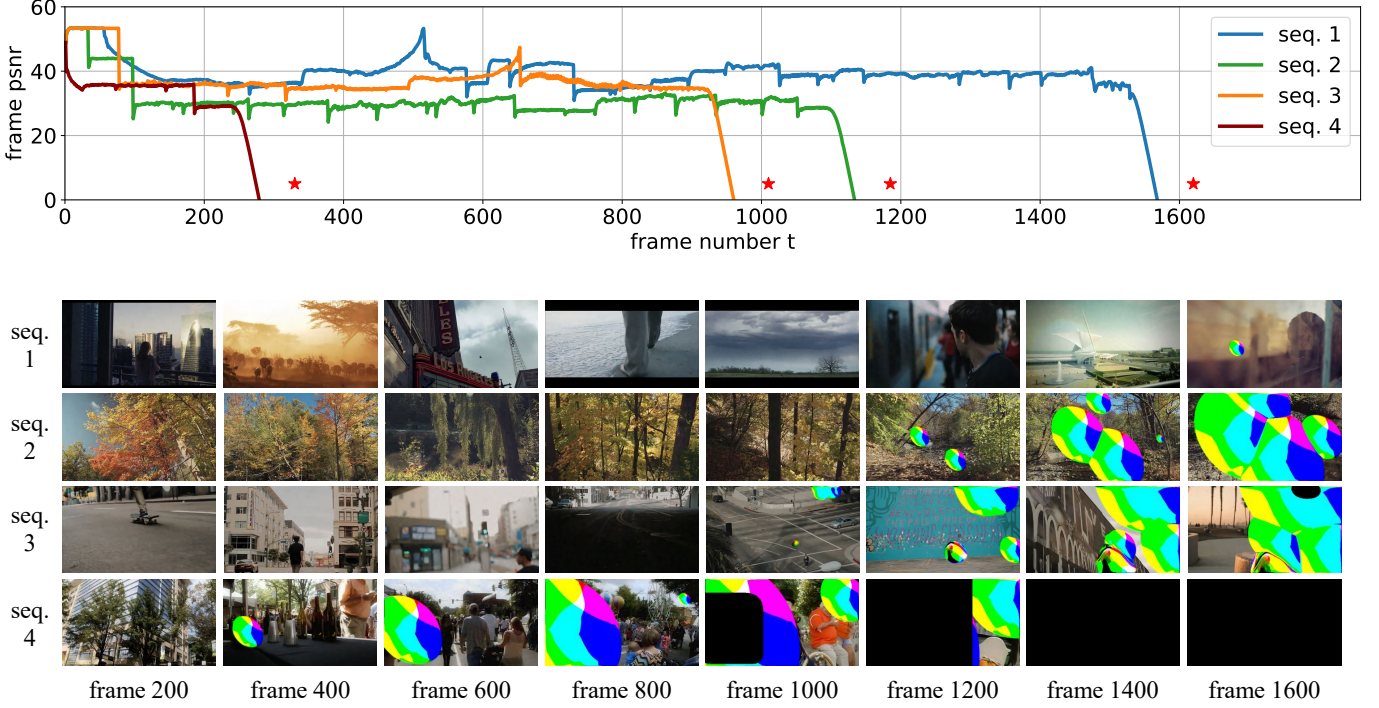
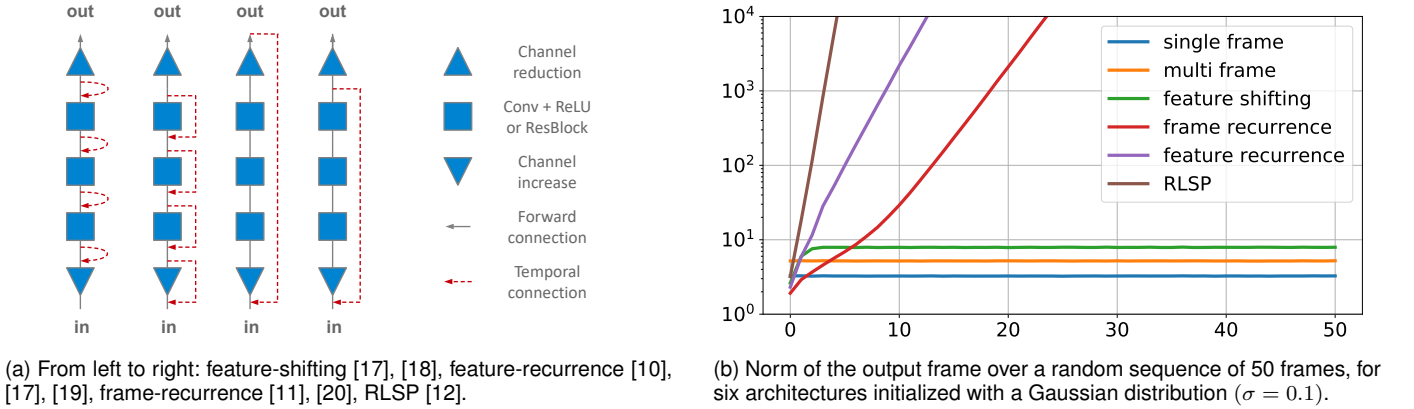


Fig. 1. The recurrent video denoiser from [10] is applied to four sequences of 1600 frames downloaded from vimeo.com. Above: The PSNR per frame is stable for a number of frames varying between 200 and 1500, before plunging below 0 on all 4 sequences (indicated by red stars). Below: The performance drop manifests itself in the form of strong colorful artifacts and black masks on the output image (see also Appendix A).



(a) From left to right: feature-shifting [17], [18], feature-recurrence [10], [17], [19], frame-recurrence [11], [20], RLSP [12].

(b) Norm of the output frame over a random sequence of 50 frames, for six architectures initialized with a Gaussian distribution ($\sigma = 0.1$).

Fig. 2. For untrained models over random inputs, feedforward architectures produce stable outputs (single-frame, multi-frame, feature-shifting) while recurrent architectures diverge (frame-recurrence, feature-recurrence, RLSP).

1.2 Contributions

The main contributions of this paper can be summarized as follows:

- We identify and characterize a serious vulnerability affecting recurrent networks for video processing: they can be unstable at inference time and fail catastrophically on long video sequences.
- To test stability, we introduce a fast and reliable diagnostic tool that produces adversarial input sequences optimized to trigger instabilities, and that can be interpreted as spatio-temporal receptive fields.
- We investigate two approaches to enforce the stability of recurrent video processing networks: con-

straining the spectral norm or constraining the stable rank of their convolutional layers.

- We extend a recently proposed weight normalization scheme called Stable Rank Normalization (SRN) [21] that simultaneously constrains the spectral norm and the stable rank of any linear mapping, to convolutional layers. We call it Stable Rank Normalization of the Layer (SRNL)—as opposed to stable rank normalization applied to the convolutional kernel.

1.3 Related Work

A number of approaches have been proposed in the literature for extending the connectivity of a CNN in the time

domain. In [22], the authors identify three classes of feedforward models: *early fusion* where the frames over a fixed time window are concatenated and processed simultaneously, *late fusion* where the frames are processed independently and their latent space representations are concatenated at a later stage, and *slow fusion* where intermediary features are concatenated at multiple levels such that higher layers get access to progressively more global information in both spatial and temporal dimensions. Variants of slow fusion were introduced a number of times under different names: *conditional convolutions* in [17], *3D convolutions* in [23], *progressive fusion* in [24] and *feature-shifting* in [18] all fuse features from different time steps at multiple levels of the network. For video restoration tasks, most standard models implement a form of early fusion [7], [8], [9], [25], [26] but late fusion [27] and two-level fusion [28] have also been used occasionally.

In contrast with the feedforward fusion approaches above, recurrent models contain feedback loops where the features are extracted at a layer l_1 at time $t - 1$ and fed back at a lower layer $l_2 < l_1$ at time t . To the best of our knowledge, a recurrent CNN was first applied to a video restoration task in [17]. The architecture proposed for video super-resolution used a heavy set of temporal connections, with forward ($x_{t-1} \rightarrow x_t \rightarrow x_{t+1}$) and backward ($x_{t+1} \rightarrow x_t \rightarrow x_{t-1}$) subnetworks, each using both *conditional* and *recurrent* connections corresponding to *feature-shifting* and *feature-recurrence* respectively according to our taxonomy from Figure 2a. Feature recurrence was used again for video denoising in [19] on a deep but non-convolutional RNN, and in [10] on the multi-frame branch of a hybrid architecture constituted of a single-frame denoiser and a multi-frame denoiser. Frame-recurrence, where the previous output frame is fed back as an additional input at the next time step, was introduced for video super-resolution in [11]. This type of recurrence was studied further in [20] where a connection was made with the concept of Kalman filtering [29]. Recently and still in super-resolution, recurrent latent space propagation (RLSP) was introduced [12]. RLSP can be interpreted as maximizing the depth and width of the recurrent connection, compared to feature-recurrence and frame-recurrence. Iterative approaches [30], [31], [32], [33] are conceptually similar to recurrent ones, but the feedback loop is part of a refinement mechanism that occurs for a fixed number of iterations, chosen as a hyperparameter by the user independently of the temporal length of the video sequence. Long Short-Term Memory (LSTM) [34] is a type of gated RNN that is compatible with convolutions and has been used for video description and captioning [35], [36], [37] but rarely for video restoration [38]. In this paper, we focus on standard recurrent architectures in their three variants: feature-recurrence, frame-recurrence and RLSP.

Training Recurrent Neural Networks (RNNs) is notoriously difficult due to the *vanishing and exploding gradients* problem: RNNs are trained by unrolling through time, which is effectively equivalent to training a very deep network [14], [39]. Relatedly, RNNs are *vulnerable to instabilities at inference time on long sequences*. This phenomenon was studied in the context of 1-layer fully connected networks in [40], and in the context of multi-layer and LSTM networks in [15], where it was shown that the RNN is stable if its Lip-

schitz constant is less than 1. In [15], it was proposed to enforce this stability constraint by projecting onto the spectral norm ball of the recurrence matrix (i.e. by clipping its singular values to 1) and a number of recent works have sought to avoid vanishing and exploding gradients by enforcing orthogonality (i.e. setting all the singular values to 1) [41], [42], [43], [44], [45], [46]. In the context of convolutional neural networks however, enforcing the Lipschitz constraint is challenging. In [47], it was proposed to clip singular values of the convolutional (but non-recurrent) layer, which was flattened into a matrix using the doubly block circulant matrix representation. However, the optimization method does not have formal convergence guarantees and requires computing all singular values of the flattened kernel. In [48], it was proposed to normalize the kernel of the convolutional (but non-recurrent) layer during training, the 4-D kernel being first reshaped into a 2-D matrix by flattening its first three dimensions. Normalization is performed by an elegant iterative scheme employing the power iteration estimating the maximal singular value of the flattened kernel. However, as discussed in [49] and [50], this approach is not suitable for Lipschitz regularization due to the invalid flattening operation used, and as a result is not suitable for stability enforcement using [15] either. To solve this issue, Gouk et al. [50] suggested replacing the 2-D matrix products in the power iteration with convolution and transpose convolution operations using the 4-D kernel tensor directly. This method was applied with success in [51] to train invertible ResNets.

Recently, Sanyal et al. (2020) [21] proposed Stable Rank Normalization (SRN), a provably optimal weight normalization scheme which minimizes the stable rank of a linear operator while constraining the spectral norm. They showed that SRN, while improving the classification accuracy, also improves generalization of neural networks and reduces memorization. However, SRN operates on a 2-D reshaping of the convolutional kernel, instead of operating on the convolutional layer as a whole.

2 STABILITY IN RECURRENT VIDEO PROCESSING

In this section, we mathematically define the notion of stability. We then introduce the Spatio-Temporal Receptive Field (STRF) diagnostic tool and the two stability constraints, before presenting our Stable Rank Normalization of the Layers algorithm.

2.1 Definitions

Partially reusing notations from [15], we define a recurrent video processing model as a non-linear dynamical system given by a Lipschitz continuous *recurrence map* $\phi_w : \mathbb{R}^n \times \mathbb{R}^d \rightarrow \mathbb{R}^n$ and an *output map* $\psi_w : \mathbb{R}^n \rightarrow \mathbb{R}^d$ parameterized by $w \in \mathbb{R}^m$. The hidden state $h_t \in \mathbb{R}^n$ and the output image $y_t \in \mathbb{R}^d$ evolve in discrete time steps according to the update rule¹

$$\begin{cases} h_t &= \phi_w(h_{t-1}, x_t) \\ y_t &= \psi_w(h_t) \end{cases} \quad (1)$$

1. The case where $y_t = h_t$ corresponds to the frame-recurrent architecture of [11].

where the vector $x_t \in [0, 1]^d$ is an arbitrary input image provided to the system at time t .

In Section 1.1, we showed examples of models that produced diverging outputs and called them “unstable”. In the following, we propose to use the notion *Bounded-Input Bounded-Output (BIBO) stability*.

Definition 1. A recurrent video processing model is *BIBO stable* if, for any admissible input $\{x_t\}_{t=0}^\infty$ for which there exist a constant C_1 such that $\sup_{t \geq 0} \|x_t\| \leq C_1$, there exists a constant C_2 such that $\sup_{t \geq 0} \|y_t\| \leq C_2$.

This definition is well suited for models using ReLU activation functions, and the diagnostic tool we introduce in the next section relies on it. However, it fails to capture a stricter notion of stability for models with bounded activation functions, which are BIBO stable by construction². Therefore, we will use the stricter notion of *Lipschitz stability* for stability enforcement, as in [15].

Definition 2. A recurrent video processing model is *Lipschitz stable* if its recurrence map ϕ_w is *contractive* in h , i.e. if there exists a constant $L < 1$ such that, for any states $h, h' \in \mathbb{R}^n$ and input $x \in \mathbb{R}^d$,

$$\|\phi_w(h, x) - \phi_w(h', x)\| \leq L\|h - h'\|. \quad (2)$$

The constant L is called the Lipschitz constant of ϕ_w . We show easily that Lipschitz stability implies BIBO stability, but the reciprocal is not always true.

2.2 Diagnosis

Consider a *trained* recurrent video processing model (ϕ_w, ψ_w) . Proving that it is *BIBO unstable* can be easy: one simply needs to find an *unstable sequence*, i.e. an input sequence $\{x_t\}_{t=0}^\infty$ such that the corresponding output sequence $\{y_t\}_{t=0}^\infty$ diverges. Proving that the model is *BIBO stable*, however, is hard: one needs to perform an exhaustive search over the (infinite) set of valid inputs, and check that none of them are unstable. Alternatively, one could try to show that the model is *Lipschitz stable* instead. Unfortunately, computing the Lipschitz constant for a neural network is, in general, NP-hard [49].

In practice, one realistic approach is to run the model on a very long video sequence and observe possible instabilities—effectively performing a *random search* for unstable sequences over the set of valid inputs—but it is not clear which video sequence should be used and how long the search should last before one can reliably conclude that the recurrent model is, indeed, stable. As Figure 1 shows, this is not a trivial question: instabilities do not occur after the same number of frames on all video sequences and it can easily take more than a thousand frames before an instability occurs. At the same time, guaranteeing stability is a fundamental safety concern in any real-world deployment.

Here, we propose to approach the problem in a different way and to *search for unstable sequence by gradient descent*. We fix a sequence length $2\tau + 1$ and an image size d , and also consider the finite input sequence $X = (x_{-\tau}, \dots, x_\tau)$ with

the corresponding finite output sequence $Y = (y_{-\tau}, \dots, y_\tau)$ such that, for all $t \in [-\tau, \tau]$:

$$\begin{aligned} h_{-\tau-1} &= 0 && \text{The initial hidden state is null} \\ h_t &= \phi_w(x_t, h_{t-1}) && \phi_w, \text{ is unrolled over the sequence} \\ y_t &= \psi_w(h_t) && \psi_w \text{ maps to the output image} \end{aligned}$$

And we search for unstable sequences by optimizing:

$$\max_{0 \leq X \leq 1} \mathcal{L}(Y) \quad (3)$$

where \mathcal{L} is a chosen objective function on the output sequence. We consider in particular the following two cases.

- When $\mathcal{L}(Y) = \|y_0\|$, the optimization problem (3) consists in finding an input sequence X such that the corresponding output sequence Y diverges maximally at time $t = 0$. Since this optimization process only affects the *frames* in X that have an effect on y_0 , the resulting input sequence X can be interpreted as a visualization of the *Temporal Receptive Field* for the output frame at time $t = 0$.
- When $\mathcal{L}(Y) = |p|$ where p is the pixel in the centre of y_0 , optimizing the problem (3) only affects the *pixels* in X that have an effect on the value of p , and hence, the resulting input sequence X can be interpreted as a visualization of the *Spatio-Temporal Receptive Field* for the output pixel p .

This type of optimization on the output of a model with regards to its input is inspired from the research on *adversarial examples* in image classification [52], [53], [54], [55], [56] and an unstable sequence X can be viewed as an *adversarial sequence*. It has also been used for *feature visualization* before [52], [57], [58], [59], [60], although never in the context of recurrent networks to the best of our knowledge. In practice, we initialise X randomly, we choose a sequence length $\tau = 40$ and an image size $d = 64 \times 64$. We then solve the optimization problem using the Adam optimizer for 1000 iterations. We iterate multiple times to avoid vanishing gradient issues but we found experimentally that the number of iterations does not need to be very high.

Computing the Spatio-Temporal Receptive Field (STRF) of a given model can then be used as a diagnostic tool for its stability, and we observe two possible behaviours.

- **The STRF is not temporally bounded.** Input frames in the distant past have an effect on p and output frames in the distant future diverge (see Figure 3a). The input sequence X constitutes an unstable sequence and we can conclude with certainty that the model is unstable.
- **The STRF is temporally bounded.** Input frames in the distant past have no effect on p and output frames in the distant future remain unaffected (see Figure 3b). No unstable sequence has been found and we can conclude with reasonable confidence that the model is stable.

Running the model on 2h30m of video (random search) and computing the STRF (gradient descent search) cannot constitute formal proofs of stability, but we show in the experimental section that they give consistent answers on the stability of various models. STRFs also help visualize

². Simply applying a sigmoid function to the output of an unstable model technically makes it BIBO stable, yet in practice, the model still suffers from instabilities and its output simply saturates.

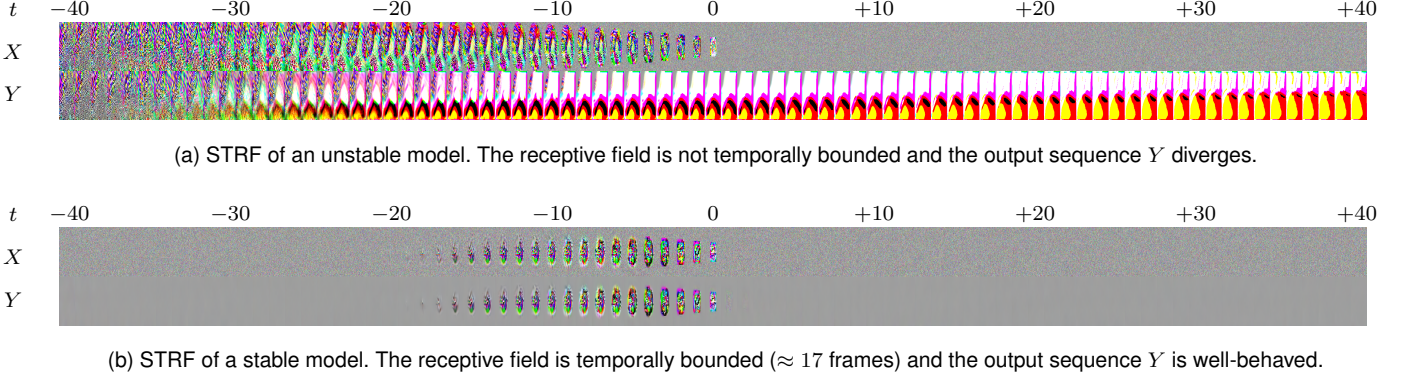


Fig. 3. Spatio-Temporal Receptive Field (STRF) as a diagnostic tool. The input sequence X is optimized to trigger instabilities in the output sequence Y (here we use $\mathcal{L}(Y) = |p|$). The sequences have been horizontally compressed to fit the page width. In the rest of the paper, we plot the STRFs every 5 frames for convenience, but the optimization is always performed on sequences of 81 frames as $\tau = 40$.

the temporal window of influence of a model, or how long information can stay in memory, and therefore illustrates the relationship between stability and memory in RNNs.

2.3 Prevention

Now, consider an *untrained* recurrent video processing model (ϕ_w, ψ_w) . In order to prevent instabilities from occurring at inference time, we want to enforce a stability constraint into the model at training time. As discussed in Section 2.1, this can be achieved by ensuring that ϕ_w is contractive with respect to the recurrent variable.

Suppose that ϕ_w is made of l convolutional layers separated by ReLU non-linearities. Each convolution can be represented by its 4D kernel tensors \mathbf{K} , or by a corresponding 2D matrix \mathbf{W} obtained from \mathbf{K} as a block matrix of doubly-block-circulant matrices [47]. Then for a layer \mathbf{W} with singular values $\{\sigma_k\}$ assumed to be sorted, the spectral norm is $\|\mathbf{W}\| = \sigma_1$, the Frobenius norm is $\|\mathbf{W}\|_F = \sqrt{\sum_k \sigma_k^2}$ and the *stable rank*³ is defined as [21], [61]:

$$\text{srnk}(\mathbf{W}) = \frac{\|\mathbf{W}\|_F^2}{\|\mathbf{W}\|^2} = \frac{\sum_k \sigma_k^2}{\sigma_1^2}. \quad (4)$$

It is a scale independent quantity and can be interpreted as an *area-under-the-curve* for the normalized singular value spectrum. Now, let L be the Lipschitz constant of ϕ_w . Since the Lipschitz constant of the ReLU non-linearity is 1, we know that L is upper-bounded by the product of the spectral norms of the linear layers [49], [52].

Proposition 1. For a recurrent model ϕ_w constituted of l linear layers with weight matrices $\mathbf{W}_1, \dots, \mathbf{W}_l \in \mathbb{R}^{n \times n}$ interspaced with ReLU non-linearities, the Lipschitz constant L of ϕ_w satisfies:

$$L \leq \prod_{i=1}^l \|\mathbf{W}_i\|. \quad (5)$$

Using this upper-bound, one can guarantee that ϕ_w is contractive (i.e. $L < 1$) with the following approach.

3. The stable rank is a soft, numerical approximation of the rank operator. It is *stable* under small perturbations of the matrix—the name has nothing to do *a priori* with the notion of stability studied here.

Approach 1. Hard Lipschitz Constraint

For all $i \in [1, l]$, we enforce $\|\mathbf{W}_i\| < 1$.

This approach has the advantage of providing a theoretical guarantee of stability. However, it is overly restrictive because the upper-bound (5) tends to significantly overestimate the Lipschitz constant L [21]. To illustrate why this is the case, suppose that ϕ_w contains two layers \mathbf{W}_1 and \mathbf{W}_2 . Then the only situation in which we have $L = \|\mathbf{W}_1\| \|\mathbf{W}_2\|$ is when the first right singular vector of \mathbf{W}_1 is aligned with the first left singular vector of \mathbf{W}_2 . In other situations, L depends on the rest of the singular value spectra of \mathbf{W}_1 and \mathbf{W}_2 and hence, on their stable ranks. These considerations lead us to a second approach to enforce $L < 1$.

Approach 2. Soft Lipschitz Constraint

For all $i \in [1, l]$, we fix $\|\mathbf{W}_i\| = \alpha$ and minimize $\text{srnk}(\mathbf{W}_i)$.

This approach does not offer any theoretical guarantee of stability for $\alpha > 1$. However, we verify empirically in Section 3 that it is also successful at promoting stability.

2.4 Stable rank normalization of the layers

A couple of tools have been proposed before to enforce the constraints of Approaches 1 and 2. *Spectral normalization* (SN), introduced by Miyato et al. [48] and popularized in GAN training [62], [63], allows one to fix the spectral norm of convolutional layers to a desired value α . *Stable rank normalization* (SRN), introduced by Sanyal et al. [21], builds on top of the previous work and allows one to also control the stable rank with a parameter $\beta \in [0, 1]$ (see Algorithm 1). However, as observed before in [49], [50], there is an issue with SN and by extension with SRN: they operate on a 2D reshaping of the kernel tensor \mathbf{K} , instead of operating on the matrix of the convolutional layer \mathbf{W} . At the same time, operating on \mathbf{W} directly is impossible: the matrix is too large to be expressed explicitly⁴. Below, we introduce a version of SRN that operates on \mathbf{W} indirectly, using \mathbf{K} . To distinguish between the two versions, we refer to our algorithm as *Stable Rank Normalization of the Layer* or SRNL (see Algorithm 2).

4. For a kernel size k , a number of input and output channels m and an image size n , the dimension of \mathbf{K} is $[k, k, m, m]$ (typically around 10^3 parameters) while the dimension of \mathbf{W} is $[nm, nm]$ (typically around 10^{13} sparse parameters).

Algorithm 1: SRN- α - β (Sanyal et al. (2020) [21])

Input: Number of iterations N , learning rate η ,
 number of channels m , image size n ,
 initial $\mathbf{K} \in \mathbb{R}^{k \times k \times m \times m}$, initial $\mathbf{u} \in \mathbb{R}^m$.

Parameters: Spectral norm α , stable rank β .

begin

```

1   for  $i = 1, \dots, N$  do
2        $\tilde{\mathbf{K}} = \text{Reshape}(\mathbf{K}, [kkm, m])^T$ 
3       Power iteration:
4        $\mathbf{v} = \tilde{\mathbf{K}}^T \mathbf{u} / \|\tilde{\mathbf{K}}^T \mathbf{u}\|_2$ 
5        $\mathbf{u} = \tilde{\mathbf{K}} \mathbf{v} / \|\tilde{\mathbf{K}} \mathbf{v}\|_2$ 
6       Spectral normalization:
7        $\tilde{\mathbf{K}} = \tilde{\mathbf{K}} / (\mathbf{u}^T (\tilde{\mathbf{K}} \mathbf{v}) + \varepsilon)$ 
8       Stable rank ( $\beta < 1$ ):
9        $\mathbf{S}_1 = \mathbf{u} \mathbf{v}^T$ 
10       $\mathbf{S}_2 = \tilde{\mathbf{K}} - \mathbf{S}_1$ 
11       $\gamma = \sqrt{\beta m - 1} / \|\mathbf{S}_2\|_F$ 
12      if  $\gamma \leq 1$  then
13           $\tilde{\mathbf{K}} = \mathbf{S}_1 + \gamma \mathbf{S}_2$ 
14       $\tilde{\mathbf{K}} = \text{Reshape}(\tilde{\mathbf{K}}^T, [k, k, m, m])$ 
15      Training step:
16       $\mathbf{K} = \mathbf{K} - \eta \nabla_{\mathbf{K}} L(\alpha \tilde{\mathbf{K}})$ 

```

Algorithm 2: SRNL- α - β (Our)

Input: Number of iterations N , learning rate η ,
 number of channels m , image size n ,
 initial $\mathbf{K} \in \mathbb{R}^{k \times k \times m \times m}$, initial $\mathbf{u} \in \mathbb{R}^{n \times n \times m}$.

Parameters: Spectral norm α , stable rank β .

begin

```

1   for  $i = 1, \dots, N$  do
2        $\tilde{\mathbf{K}} = \mathbf{K}$ 
3       Power iteration:
4        $\mathbf{v} = \tilde{\mathbf{K}}^T * \mathbf{u} / \|\tilde{\mathbf{K}}^T * \mathbf{u}\|_2$ 
5        $\mathbf{u} = \tilde{\mathbf{K}} * \mathbf{v} / \|\tilde{\mathbf{K}} * \mathbf{v}\|_2$ 
6       Spectral normalization:
7        $\tilde{\mathbf{K}} = \tilde{\mathbf{K}} / (\mathbf{u}^T (\tilde{\mathbf{K}} * \mathbf{v}) + \varepsilon)$ 
8       Stable rank ( $\beta < 1$ ):
9        $\mathbf{S}_1 = \nabla_{\tilde{\mathbf{K}}} (\mathbf{u}^T (\tilde{\mathbf{K}} * \mathbf{v}))$ 
10       $\mathbf{S}_2 = \tilde{\mathbf{K}} - \mathbf{S}_1$ 
11       $\gamma = \sqrt{\beta m - 1} / n^2 / \|\mathbf{S}_2\|_F$ 
12      if  $\gamma \leq 1$  then
13           $\tilde{\mathbf{K}} = \mathbf{S}_1 + \gamma \mathbf{S}_2$ 
14      Training step:
15       $\mathbf{K} = \mathbf{K} - \eta \nabla_{\mathbf{K}} L(\alpha \tilde{\mathbf{K}})$ 

```

The two algorithms are structurally identical—they consist in a power iteration to compute the spectral norm (steps 2,3), a normalization (step 4) and a re-weighting of a rank one matrix \mathbf{S}_1 and a residual matrix \mathbf{S}_2 (steps 5, 6, 7, 8)—but they present a number of key differences. In SRNL, the random vector \mathbf{u} has two more dimensions and is the size of a full input feature map $([1, n, n, m])$. The kernel is not flattened (steps 1, 9). The power iteration is performed using a convolution $(\tilde{\mathbf{K}} * \cdot)$ and a transposed convolution $(\tilde{\mathbf{K}}^T * \cdot)$ as suggested in [50], based on the observations that:

$$\begin{aligned} \mathbf{v} = \tilde{\mathbf{W}}^T \mathbf{u} &\Leftrightarrow \mathbf{v} = \tilde{\mathbf{K}}^T * \mathbf{u} \quad (\text{step 2}) \quad \text{and} \\ \mathbf{u} = \tilde{\mathbf{W}} \mathbf{v} &\Leftrightarrow \mathbf{u} = \tilde{\mathbf{K}} * \mathbf{v} \quad (\text{step 3}). \end{aligned} \quad (6)$$

The spectral normalization is also performed using a convolution (step 4). The rank one matrix $\mathbf{S}_1 = \mathbf{u} \mathbf{v}^T$ is expressed as a 4D kernel tensor through the gradient of $\mathbf{u}^T (\tilde{\mathbf{K}} * \mathbf{v})$ with respect to $\tilde{\mathbf{K}}$ (step 5), based on the observation that:

$$\mathbf{u} \mathbf{v}^T = \nabla_{\tilde{\mathbf{W}}} (\text{trace}(\tilde{\mathbf{W}} \mathbf{v} \mathbf{u}^T)) = \nabla_{\tilde{\mathbf{W}}} (\mathbf{u}^T \tilde{\mathbf{W}} \mathbf{v}). \quad (7)$$

Finally, writing $\|\tilde{\mathbf{W}}\|_F$ explicitly yields $\|\tilde{\mathbf{W}}\|_F = n \|\tilde{\mathbf{K}}\|_F$ and therefore (step 7):

$$\gamma = \frac{\sqrt{\beta n m - 1}}{n \|\mathbf{S}_2\|_F} = \frac{\sqrt{\beta m - 1/n^2}}{\|\mathbf{S}_2\|_F}. \quad (8)$$

When $\beta = 1$, SRN and SRNL are equivalent to performing spectral normalization on \mathbf{K} and \mathbf{W} respectively. When $\beta < 1$, they also have an effect on the stable rank of their respective matrices. We found experimentally that SRN multiplies the training time by a factor of ≈ 1.8 and SRNL multiplies the training time by a factor of ≈ 2.2 . At inference time, the weights are fixed and normalized convolutions have the same complexity as standard convolutions.

3 EXPERIMENTS

In this section, we illustrate our diagnostic tool on a number of video denoising models trained in a standard way without constraints. We then show that our stable rank normalization algorithm successfully enforces stability via the two approaches discussed.

3.1 Unconstrained Models

To reflect the variety of architectures used for video processing tasks, we consider two backbone networks and three types of recurrence. The two backbone networks consist in a DnCNN-like [1] stack of 10 convolutions and ReLU non-linearities (VDnCNN), and a ResNet-like [64] stack of 5 residual blocks containing two convolutions each separated by a ReLU (VResNet). The three types of recurrences are the ones considered in Section 1.1, namely, feature-recurrence [10], [17], [19], frame-recurrence [11], [20] and recurrent latent space propagation (RLSP) [12]. More architectural details are provided in Figure 4.

Here, we chose to focus on video denoising and keep the evaluation of other video processing tasks for future work. We train our models using the Vimeo-90k septuplet dataset [16], consisting of about 90k 7-frame RGB sequences with a resolution of 448×256 downloaded from vimeo.com. We generate clean-noisy training pairs by applying Gaussian noise with standard deviation $\sigma = 30$. The recurrent networks are trained using backpropagation through time on sequences of 7 frames—making use of the full length of the Vimeo-90k sequences—on image crops of 64×64 pixels. We train using the Adam optimizer with a batch size of 32 for 600k steps. For comparison, we also consider the traditional non-recurrent methods BM3D [65] and VNLB [66] and the recurrent deep burst denoising network of [10], which we

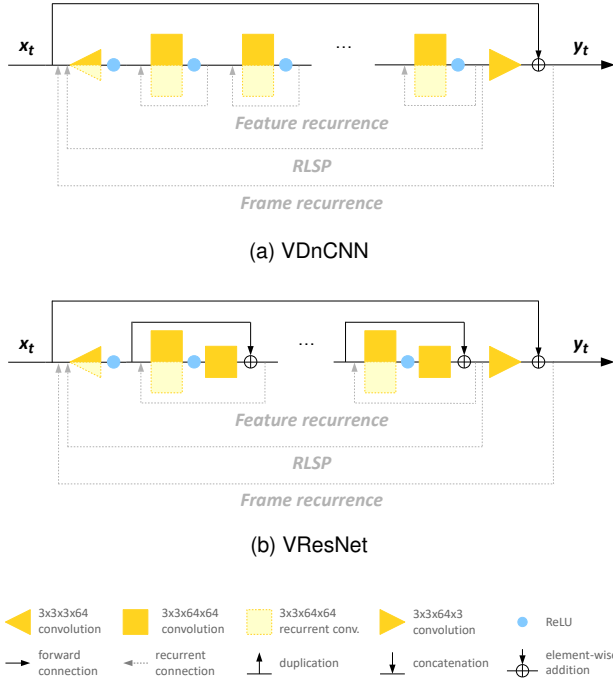


Fig. 4. The two architectures and three types of recurrence considered.

refer to as DBDNet and retrain in the same conditions as the other recurrent models. In Table 1, we show the numbers of parameters and processing speeds⁵ (fps) of each method, as well as their denoising performances as measured by the PSNR on the first frame (PSNR_1), last frame (PSNR_7) and averaged over all the frames ($\text{PSNR}_{\text{mean}}$) on the first 1024 validation sequences of the Vimeo-90k septuplet dataset. Since BM3D processes sequences frame by frame, VNLB processes sequences all at once, and recurrent networks reuse information from one frame to the next, the performance of BM3D is flat over each sequence, it peaks on the median frame for VNLB and it peaks on the last frame for recurrent networks. VNLB outperforms other methods on average, but recurrent models have a higher PSNR on the seventh frame, and they are orders of magnitude faster to run. The performance on the last frame (PSNR_7) is of special interest as it can be regarded as an approximation of the plateaued performance of an online video denoiser applied to a very long sequence. For that reason, we focus on this metric in the rest of the paper. Among the different recurrent models, the VResNet backbone systematically outperforms the VDnCNN one, possibly partly because VDnCNN is slower to converge. Frame-recurrent architectures are the lightest and fastest, but feature-recurrence and RLSP yield better performance. Interestingly, VResNet-RLSP has about 60% less parameters than DBDNet [10], is significantly faster and performs better (+0.33dB on PSNR_7).

To evaluate the stability of each recurrent model, we apply them to one long video sequence lasting approximately 2h30m consisting of several clips downloaded from vimeo.com and concatenated together (2×10^5 frames). Each

TABLE 1

Size, processing speed and performance of the different video denoising methods considered, measured on the first frame (PSNR_1), last frame (PSNR_7), and averaged over all the frames ($\text{PSNR}_{\text{mean}}$) on the Vimeo-90k septuplet dataset.

	# param.	fps	PSNR_1	PSNR_7	$\text{PSNR}_{\text{mean}}$
BM3D [65]	n/a	2	33.86	33.83	33.85
VNLB [66]	n/a	0.02	35.24	35.17	35.78
DBDNet [10]	965k	30	34.16	35.47	35.16
VDnCNN-frame	375k	70	33.94	34.84	34.68
VDnCNN-feat	741k	40	34.05	35.02	34.79
VDnCNN-RLSP	410k	60	33.95	34.98	34.77
VResNet-frame	375k	70	34.23	35.47	35.21
VResNet-feat	557k	50	34.35	35.74	35.41
VResNet-RLSP	410k	60	34.25	35.80	35.42

time the PSNR drops below 0, we consider that an instability occurs and we *reset* the recurrent features to 0. We call *instability onset* the number of frames leading to an instability. In Table 2, we report the 1st and 9th decile of the instability onsets for each model (∞ means no instability observed). According to this test, VDnCNN-frame, VResNet-frame and VResNet-RLSP are stable while VDnCNN-feat, VDnCNN-RLSP and VResNet-feat are unstable. For VDnCNN-feat in particular, the instability onset is above 5709 frames in 10% of the cases, highlighting the necessity to run this test on very long video sequences. We show examples of output sequences for the 3 unstable models in Appendix A. In an attempt at explaining why certain models are stable and others are not, we compute the singular value spectra of all their convolutional layers using the code provided by Sedghi et al. in [47]—this gives us access to their spectral norms in particular, or maximum singular value (leftmost value in each spectrum). Unfortunately, the spectra all have comparable profiles and are therefore uninformative, with spectral norms typically around 6: too high to conclude anything about the contractiveness of the recurrent maps using the upper-bound (5). Finally, we compute the spatio-temporal receptive fields of each model using the approach described in Section 2.2. In accordance with the previous test, VDnCNN-feat, VDnCNN-RLSP and VResNet-feat exhibit the characteristic behaviour of unstable models, with long-range temporal dependencies accumulating in the input sequences X , resulting in unstable output sequences Y diverging at frame +40. The spatio-temporal receptive fields of VDnCNN-frame, VResNet-frame and VResNet-RLSP on the other hand, are well-behaved: the information flow is limited to a finite temporal window, input frames in the distant past have no influence on the current frame and future output frames do not diverge. Frame-recurrent models appear to have a short temporal receptive field (≈ 10 frames) compared to other models, possibly making this type of recurrence more stable in practice.

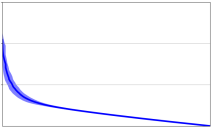
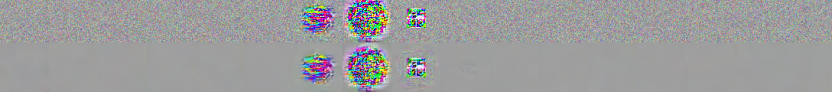
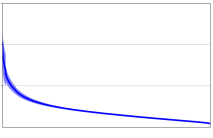
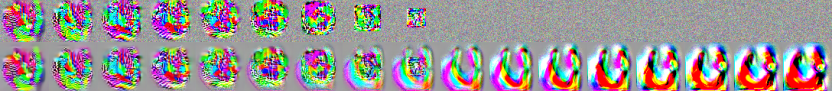
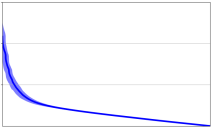
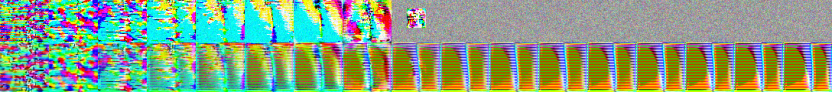
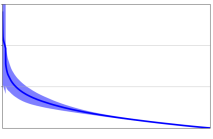
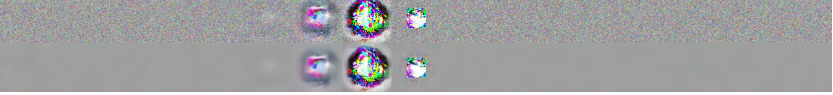
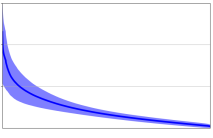
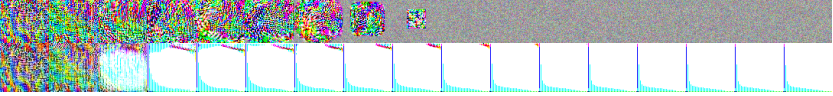
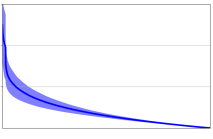
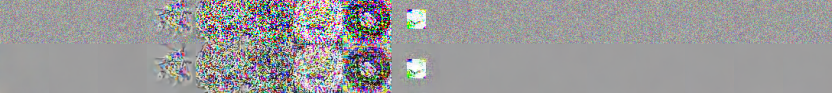
3.2 Constrained models

We saw in the previous section that various models with different backbone architectures and types of recurrence trained in standard conditions are unstable on long video sequences at inference time. We have also discussed in the introduction how the instabilities observed constitute

5. We use the authors' implementations of BM3D (Matlab) and VNLB (C++). All our recurrent networks are implemented in TensorFlow. The processing speeds are indicative of an order of magnitude only.

TABLE 2

Instabilities in 6 models with 2 backbone architectures and 3 types of recurrences. For each model, we show the performance on the 7th frame of the Vimeo-90k validation dataset (PSNR_7), the 1st and 9th deciles of the instability onsets on a sequence of about 2h30min (∞ means no instabilities observed). We also show the singular value spectrum averaged over the convolutions of the model, computed using the code from [47], and the temporal receptive field computed using our method.

model	PSNR_7 1 st dec. 9 th dec.	Average Singular Value Spectrum	Spatio-Temporal Receptive Field									
VDnCNN -frame	34.84 ∞ ∞		t	-40	-30	-20	-10	0	+10	+20	+30	+40
			X									
VDnCNN -feat	35.02 157 5709		t	-40	-30	-20	-10	0	+10	+20	+30	+40
			X									
VDnCNN -RLSP	34.98 74 271		t	-40	-30	-20	-10	0	+10	+20	+30	+40
			X									
VResNet -frame	35.47 ∞ ∞		t	-40	-30	-20	-10	0	+10	+20	+30	+40
			X									
VResNet -feat	35.74 29 75		t	-40	-30	-20	-10	0	+10	+20	+30	+40
			X									
VResNet -RLSP	35.80 ∞ ∞		t	-40	-30	-20	-10	0	+10	+20	+30	+40
			X									

catastrophic failures that are a serious concern for real-world deployment. Now, we show that inference-time stability can be enforced during training, with the help of our *stable rank normalization of the layers* algorithm (SRNL). We focus on the VResNet-feat architecture, as it appeared to be the most vulnerable to instabilities with 80% of the onsets happening between 29 and 75 frames only.

First, let us consider the model trained with SRN-1.0-1.0 in the first line of Table 3. According to the average singular value spectrum computed using the code from [47], its convolutional layers have spectral norms that are significantly larger than 1 at around 2.5, and which vary significantly (± 0.2). This observation confirms that normalizing a 2D reshaping of the convolutional kernel \mathbf{K} is a poor approximation of normalizing the convolutional layer \mathbf{W} : SRN fails to set the spectral norm of \mathbf{W} to the desired value α and this motivates the introduction of SRNL.

Now, let us consider the models trained with SRNL- α -1.0 for $\alpha \in \{2.0, 1.5, 1.0, 0.5\}$ in lines 2 to 5 of Table 3. As expected, the spectral norms of all the convolutional layers

are now precisely set to their respective values of α . Our test on the long video sequence and our spatio-temporal receptive field diagnostic then show that the models trained with $\alpha > 1$ are unstable while the models trained with $\alpha < 1$ are stable. This observation confirms that our Hard Lipschitz Constraint is effective at enforcing stability. Interestingly, reducing $\alpha < 1$ shortens the temporal length of the receptive field, a side effect of the recurrence map becoming more contractive. However, reducing α also hurts performance, as measured by the PSNR_7 (-0.4dB from $\alpha = 2.0$ to $\alpha = 1.0$), and this motivates the introduction of the Soft Lipschitz Constraint.

Finally, let us consider the models trained with SRNL-2.0- β for $\beta \in \{0.4, 0.2, 0.1, 0.05\}$ in lines 6 to 9 of Table 3. As expected, varying β has no effect on the spectral norm of the convolutional layers, but it has an effect on their stable rank, or the area-under-the-curve of their singular value spectra. Again, our test on the long video sequence and our spatio-temporal receptive field diagnostic show that there is a value of β for which the stability of the model

changes: models trained with $\beta > 0.1$ are unstable while the models trained with $\beta < 0.1$ are stable. This observation confirms that our Soft Lipschitz Constraint is also effective at promoting stability. Interestingly, reducing β also shortens the temporal length of the receptive field but the effect is softer than with α , suggesting that controlling the stable rank of the linear layers of a model has a softer effect on its Lipschitz constant than controlling their spectral norms. Importantly, the cost of stability in terms of performance is now lower, and we obtain a stable model that performs better than with the Hard Lipschitz Constraint approach (+0.18dB between $\alpha = 1.0, \beta = 1.0$ and $\alpha = 2.0, \beta = 0.1$).

3.3 Discussion

In the previous sections, we showed that inference-time stability can be enforced during training by constraining the Lipschitz constant of the model to be lower than 1. In this section, we discuss possible alternative strategies.

Given a pre-trained model known to be unstable, one could consider ways to operate it in a stable manner. One approach would for instance consist in running the model burst by burst—either without overlap (e.g. running frames 1 to 10, then 11 to 20, then 21 to 30, etc.), with partial overlap (e.g. frames 1-10, 5-15, 10-20, etc.) or with full overlap (e.g. frames 1-10, 2-11, 3-12, etc.)—while resetting the recurrent features to zero between each burst. This strategy would prevent instabilities from building up, but it presents a number of issues. Without overlap, the performance fluctuates, the model constantly having to go through a new burn-in period, starting from its performance on a single frame. With overlap, the approach becomes redundant and computationally inefficient, therefore losing one of the main advantages of using a recurrent model in the first place. Moreover, running the model burst by burst introduces discontinuities which are likely to be visually perceptible in the form of flickering artefacts. Another approach would consist in *dampening the recurrent features* by a factor $\lambda < 1$, allowing a smooth transition between a stable, single frame regime ($\lambda = 0$), and an unstable fully recurrent regime ($\lambda = 1$). This approach is illustrated on a sequence of 700 frames in Figure 5, and we study it further in Appendix B, where we confirm that the price of stability in terms of performance is much higher than with our Hard and Soft Lipschitz Constraints (see also the summary Table 4).

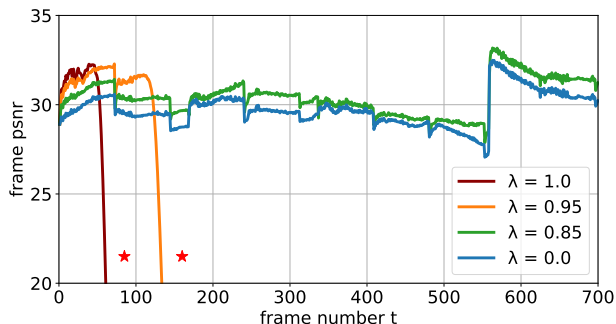


Fig. 5. Feature dampening of VResNet-feat for four dampening factors λ on a video sequence of 700 frames. Decreasing λ improves stability but has a strong negative impact on performance.

TABLE 4

Summary Table. We compare the performances of VResNet-feat in different scenarios. The soft Lipschitz constraint offers the best performance for a stable model.

VResNet-feat with ...	PSNR ₇	Stable
No Constraint	35.74	✗
Feature Dampening ($\lambda = 0.55$)	34.62	✓
Hard Lipschitz Constraint (SRNL-1.0-1.0)	35.31	✓
Soft Lipschitz Constraint (SRNL-2.0-0.1)	35.59	✓

The instabilities studied in this paper could also be interpreted as a domain adaptation problem: models trained on short sequences fail to generalize to sequences of several hundred frames. To test this hypothesis however, we face computational and data constraints. It is unrealistic to train large recurrent video processing models on sequences of more than 10 to 20 frames—the training process involves backpropagation through time, which has large memory requirements—and even if it was possible, collecting the required data would quickly become impractical. To work around these issues, we perform experiments on a small VDnCNN model where the number of internal convolutions has been reduced to only one, allowing us to unroll the model up to 56 times through time during training, and we generate long sequences with synthetic motion from single frames. We show in Appendix C that, not only does the model trained on sequences of 56 frames still suffer from instabilities at inference time, but it also suffers from instabilities at training time due to exploding gradients.

4 CONCLUSION

In this paper, we have identified and characterized a serious vulnerability affecting recurrent networks for video restoration tasks: they can be unstable and fail catastrophically on long video sequences. To avoid problems in practice, we recommend to follow some guidelines. (1) The stability of the model should always be tested, either by evaluating it on hours of video data, or preferably by actively looking for unstable sequences, using our spatio-temporal receptive field diagnostic tool. (2) In safety-critical applications, stability can be guaranteed by applying a Hard Lipschitz Constraint on the spectral norms of the convolutional layers (SRNL with $\alpha < 1$ and $\beta = 1$). (3) In non safety-critical applications, stability can be obtained with minimal performance loss by applying a Soft Lipschitz Constraint on the stable rank of the convolutional layers (SRNL with $\alpha > 1$ and $\beta < 1$).

REFERENCES

- [1] K. Zhang, W. Zuo, Y. Chen, D. Meng, and L. Zhang, “Beyond a Gaussian denoiser: Residual learning of deep CNN for image denoising,” *IEEE Transactions on Image Processing*, vol. 26, no. 7, pp. 3142–3155, 2017.
- [2] T. Brooks, B. Mildenhall, T. Xue, J. Chen, D. Sharlet, and J. T. Barron, “Unprocessing images for learned raw denoising,” in *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019.
- [3] M. Gharbi, G. Chaurasia, S. Paris, and F. Durand, “Deep joint demosaicking and denoising,” *ACM Transactions on Graphics (TOG)*, vol. 35, no. 6, pp. 1–12, 2016.

- [4] F. Kokkinos and S. Lefkimmiatis, "Iterative joint image demosaicking and denoising using a residual denoising network," *IEEE Transactions on Image Processing*, vol. 28, no. 8, 2019.
- [5] B. Lim, S. Son, H. Kim, S. Nah, and K. Mu Lee, "Enhanced deep residual networks for single image super-resolution," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops (CVPR workshops)*, 2017.
- [6] T. Dai, J. Cai, Y. Zhang, S.-T. Xia, and L. Zhang, "Second-order attention network for single image super-resolution," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019.
- [7] Z. Liu, L. Yuan, X. Tang, M. Uyttendaele, and J. Sun, "Fast burst images denoising," *ACM Transactions on Graphics (TOG)*, vol. 33, no. 6, 2014.
- [8] Y. Jo, S. Wug Oh, J. Kang, and S. Joo Kim, "Deep video super-resolution network using dynamic upsampling filters without explicit motion compensation," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018.
- [9] X. Wang, K. C. Chan, K. Yu, C. Dong, and C. Change Loy, "Edvr: Video restoration with enhanced deformable convolutional networks," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops (CVPR Workshops)*, 2019.
- [10] C. Godard, K. Matzen, and M. Uyttendaele, "Deep burst denoising," in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2018.
- [11] M. S. M. Sajjadi, R. Vemulapalli, and M. Brown, "Frame-recurrent video super-resolution," in *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018.
- [12] D. Fuoli, S. Gu, and R. Timofte, "Efficient video super-resolution through recurrent latent space propagation," in *Proceedings of the IEEE International Conference on Computer Vision Workshops (IC-CVW)*, 2019.
- [13] T. Laurent and J. von Brecht, "A recurrent neural network without chaos," in *International Conference on Learning Representations (ICLR)*, 2017.
- [14] R. Pascanu, T. Mikolov, and Y. Bengio, "On the difficulty of training recurrent neural networks," in *International Conference on Machine Learning (ICML)*, 2013.
- [15] J. Miller and M. Hardt, "Stable recurrent models," in *International Conference on Learning Representations (ICLR)*, 2019.
- [16] T. Xue, B. Chen, J. Wu, D. Wei, and W. T. Freeman, "Video enhancement with task-oriented flow," *International Journal of Computer Vision (IJCV)*, vol. 127, no. 8, pp. 1106–1125, 2019.
- [17] Y. Huang, W. Wang, and L. Wang, "Bidirectional recurrent convolutional networks for multi-frame super-resolution," in *Advances in Neural Information Processing Systems (NeurIPS)*, 2015.
- [18] J. Lin, C. Gan, and S. Han, "Tsm: Temporal shift module for efficient video understanding," in *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, 2019.
- [19] X. Chen, L. Song, and X. Yang, "Deep rnns for video denoising," in *Applications of Digital Image Processing XXXIX*, vol. 9971. International Society for Optics and Photonics, 2016, p. 99711T.
- [20] P. Arias and J.-M. Morel, "Kalman filtering of patches for frame-recursive video denoising," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops (CVPRW)*, 2019.
- [21] A. Sanyal, P. H. Torr, and P. K. Dokania, "Stable rank normalization for improved generalization in neural networks and gans," in *International Conference on Learning Representations (ICLR)*, 2020.
- [22] A. Karpathy, G. Toderici, S. Shetty, T. Leung, R. Sukthankar, and L. Fei-Fei, "Large-scale video classification with convolutional neural networks," in *Proceedings of the IEEE conference on Computer Vision and Pattern Recognition (CVPR)*, 2014.
- [23] D. Tran, L. Bourdev, R. Fergus, L. Torresani, and M. Paluri, "Learning spatiotemporal features with 3d convolutional networks," in *Proceedings of the IEEE international conference on computer vision (ICCV)*, 2015.
- [24] P. Yi, Z. Wang, K. Jiang, J. Jiang, and J. Ma, "Progressive fusion video super-resolution network via exploiting non-local spatiotemporal correlations," in *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, 2019.
- [25] J. Caballero, C. Ledig, A. Aitken, A. Acosta, J. Totz, Z. Wang, and W. Shi, "Real-time video super-resolution with spatio-temporal networks and motion compensation," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017.
- [26] B. Mildenhall, J. T. Barron, J. Chen, D. Sharlet, R. Ng, and R. Carroll, "Burst denoising with kernel prediction networks," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018.
- [27] M. Tassano, J. Delon, and T. Veit, "Dvdnet: A fast network for deep video denoising," in *IEEE International Conference on Image Processing (ICIP)*, 2019.
- [28] —, "Fastdvdnet: Towards real-time deep video denoising without flow estimation," in *Proceedings of the IEEE conference on computer vision and pattern recognition (CVPR)*, 2020.
- [29] R. E. Kalman, "A new approach to linear filtering and prediction problems," *Journal of Basic Engineering*, vol. 82, no. 1, pp. 35–45, 1960.
- [30] M. Irani and S. Peleg, "Improving resolution by image registration," *CVGIP: Graphical models and image processing*, vol. 53, no. 3, pp. 231–239, 1991.
- [31] A. R. Zamir, T.-L. Wu, L. Sun, W. B. Shen, B. E. Shi, J. Malik, and S. Savarese, "Feedback networks," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017.
- [32] M. Haris, G. Shakhnarovich, and N. Ukita, "Deep back-projection networks for super-resolution," in *Proceedings of the IEEE conference on computer vision and pattern recognition (CVPR)*, 2018.
- [33] —, "Recurrent back-projection network for video super-resolution," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019.
- [34] S. Hochreiter and J. Schmidhuber, "Long short-term memory," *Neural computation*, vol. 9, no. 8, pp. 1735–1780, 1997.
- [35] J. Donahue, L. Anne Hendricks, S. Guadarrama, M. Rohrbach, S. Venugopalan, K. Saenko, and T. Darrell, "Long-term recurrent convolutional networks for visual recognition and description," in *Proceedings of the IEEE conference on computer vision and pattern recognition (CVPR)*, 2015.
- [36] H. Yu, J. Wang, Z. Huang, Y. Yang, and W. Xu, "Video paragraph captioning using hierarchical recurrent neural networks," in *Proceedings of the IEEE conference on computer vision and pattern recognition (CVPR)*, 2016.
- [37] J. Johnson, A. Karpathy, and L. Fei-Fei, "Densecap: Fully convolutional localization networks for dense captioning," in *Proceedings of the IEEE conference on computer vision and pattern recognition (CVPR)*, 2016.
- [38] X. Tao, H. Gao, R. Liao, J. Wang, and J. Jia, "Detail-revealing deep video super-resolution," in *Proceedings of the IEEE International Conference on Computer Vision (CVPR)*, 2017.
- [39] Y. Bengio, P. Simard, and P. Frasconi, "Learning long-term dependencies with gradient descent is difficult," *IEEE transactions on neural networks*, vol. 5, no. 2, pp. 157–166, 1994.
- [40] L. Jin, P. N. Nikiforuk, and M. M. Gupta, "Absolute stability conditions for discrete-time recurrent neural networks," *IEEE Transactions on Neural Networks*, vol. 5, no. 6, pp. 954–964, 1994.
- [41] M. Arjovsky, A. Shah, and Y. Bengio, "Unitary evolution recurrent neural networks," in *International Conference on Machine Learning (ICML)*, 2016.
- [42] S. Wisdom, T. Powers, J. Hershey, J. Le Roux, and L. Atlas, "Full-capacity unitary recurrent neural networks," in *Advances in neural information processing systems (NeurIPS)*, 2016.
- [43] Z. Mhammedi, A. Hellicar, A. Rahman, and J. Bailey, "Efficient orthogonal parametrisation of recurrent neural networks using householder reflections," in *International Conference on Machine Learning (ICML)*, 2017.
- [44] E. Vorontsov, C. Trabelsi, S. Kadoury, and C. Pal, "On orthogonality and learning recurrent networks with long term dependencies," in *International Conference on Machine Learning (ICML)*, 2017.
- [45] C. Jose, M. Cissé, and F. Fleuret, "Kronecker recurrent units," in *International Conference on Machine Learning (ICML)*, 2018.
- [46] J. Zhang, Q. Lei, and I. S. Dhillon, "Stabilizing gradients for deep neural networks via efficient svd parameterization," in *International Conference on Machine Learning (ICML)*, 2018.
- [47] H. Sedghi, V. Gupta, and P. M. Long, "The singular values of convolutional layers," in *International Conference on Learning Representations (ICLR)*, 2019.
- [48] T. Miyato, T. Kataoka, M. Koyama, and Y. Yoshida, "Spectral normalization for generative adversarial networks," in *International Conference on Learning Representations (ICLR)*, 2018.
- [49] K. Scaman and A. Virmaux, "Lipschitz regularity of deep neural networks: analysis and efficient estimation," in *Advances in Neural Information Processing Systems (NeurIPS)*, 2018.
- [50] H. Gouk, E. Frank, B. Pfahringer, and M. Cree, "Regularisation of neural networks by enforcing lipschitz continuity," in *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018.

- [51] J. Behrmann, W. Grathwohl, R. T. Chen, D. Duvenaud, and J.-H. Jacobsen, "Invertible residual networks," in *International Conference on Machine Learning (ICML)*, 2019.
- [52] C. Szegedy, W. Zaremba, I. Sutskever, J. Bruna, D. Erhan, I. Goodfellow, and R. Fergus, "Intriguing properties of neural networks," in *International Conference on Learning Representations (ICLR)*, 2014.
- [53] I. J. Goodfellow, J. Shlens, and C. Szegedy, "Explaining and harnessing adversarial examples," in *International Conference on Learning Representations (ICLR)*, 2015.
- [54] A. Kurakin, I. Goodfellow, and S. Bengio, "Adversarial machine learning at scale," in *International Conference on Learning Representations (ICLR)*, 2017.
- [55] A. Ilyas, S. Santurkar, D. Tsipras, L. Engstrom, B. Tran, and A. Madry, "Adversarial examples are not bugs, they are features," in *Advances in Neural Information Processing Systems (NeurIPS)*, 2019.
- [56] B. Biggio and F. Roli, "Wild patterns: Ten years after the rise of adversarial machine learning," *Pattern Recognition*, vol. 84, pp. 317–331, 2018.
- [57] D. Erhan, Y. Bengio, A. Courville, and P. Vincent, "Visualizing higher-layer features of a deep network," *University of Montreal*, vol. 1341, no. 3, p. 1, 2009.
- [58] A. Mahendran and A. Vedaldi, "Understanding deep image representations by inverting them," in *Proceedings of the IEEE conference on computer vision and pattern recognition (CVPR)*, 2015.
- [59] A. Nguyen, J. Yosinski, and J. Clune, "Deep neural networks are easily fooled: High confidence predictions for unrecognizable images," in *Proceedings of the IEEE conference on computer vision and pattern recognition (CVPR)*, 2015.
- [60] C. Olah, A. Mordvintsev, and L. Schubert, "Feature visualization," *Distill*, 2017, <https://distill.pub/2017/feature-visualization>.
- [61] M. Rudelson and R. Vershynin, "Sampling from large matrices: An approach through geometric functional analysis," *Journal of the ACM (JACM)*, vol. 54, no. 4, pp. 21–es, 2007.
- [62] H. Zhang, I. Goodfellow, D. Metaxas, and A. Odena, "Self-attention generative adversarial networks," in *International Conference on Machine Learning (ICML)*, 2019.
- [63] A. Brock, J. Donahue, and K. Simonyan, "Large scale gan training for high fidelity natural image synthesis," in *International Conference on Learning Representations (ICLR)*, 2019.
- [64] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proceedings of the IEEE conference on computer vision and pattern recognition (CVPR)*, 2016.
- [65] K. Dabov, A. Foi, V. Katkovnik, and K. Egiazarian, "Image denoising by sparse 3-d transform-domain collaborative filtering," *IEEE Transactions on image processing*, vol. 16, no. 8, pp. 2080–2095, 2007.
- [66] P. Arias and J.-M. Morel, "Towards a bayesian video denoising method," in *International Conference on Advanced Concepts for Intelligent Vision Systems*. Springer, 2015, pp. 107–117.
- [67] G. Boracchi and A. Foi, "Modeling the performance of image restoration from motion blur," *IEEE Transactions on Image Processing*, vol. 21, no. 8, pp. 3502–3517, 2012.



Thomas Tanay is a research scientist at Huawei Technologies Ltd, Noah's Ark Lab (UK). He received his PhD in Machine Learning after following a MRes in Modelling Biological Complexity at University College London (UK). He previously received a MSc in Cognitive Computing from Goldsmiths, University of London (UK) and a Mechanical Engineering degree from Supméca Paris (France).



Aivar Sootla is a senior research scientist at Huawei R&D UK. He has received his MSc in applied mathematics from Lomonosov Moscow State University (Russia) and his PhD in control engineering from Lund University (Sweden). He has held research positions at the Department of Bioengineering, Imperial College London (UK); at the Montefiore Institute, University of Liège (Belgium); and at the Department of Engineering Science, University of Oxford (UK).



Matteo Maggioni received B.Sc. and M.Sc. degrees in computer science from Politecnico di Milano (Italy), a Ph.D. degree in image processing from Tampere University of Technology (Finland), and a Post-Doc with the EEE Department in Imperial College London (UK). He is currently a research scientist with Huawei R&D (London, UK). His research interests include nonlocal transform-domain filtering, adaptive signal-restoration techniques, and deep learning methods for computer vision and image restoration.



Puneet K. Dokania is a senior researcher at the University of Oxford and a principal research scientist at Five AI (UK). Prior to this, he received his Master of Science in Informatics from Grenoble INP (ENSIMAG), and a Ph.D. in machine learning and applied mathematics from École Centrale Paris and INRIA, France. Puneet's current research revolves around developing reliable and efficient algorithms with natural intelligence using deep learning.



Philip Torr did his PhD (DPhil) at University of Oxford. He left Oxford to work for six years as a research scientist for Microsoft Research, first in Redmond USA in the Vision Technology Group, then in Cambridge UK founding the vision side of the Machine learning and perception group. He then became a Professor in Computer Vision and Machine Learning at Oxford Brookes University. In 2013 he returned to the University of Oxford as full professor. His group has worked mostly in the area of computer vision and machine learning, and has won several awards including the Marr prize, best paper CVPR, best paper ECCV. He is a Fellow of the Royal Academy of Engineering. He has founded companies Alstetic and Oxsight and is Chief Scientific Advisor to FiveAI.



Aleš Leonardis is a Senior Research Scientist and Computer Vision Team Leader at Huawei Technologies Noah's Ark Lab (London, UK). He is also Chair of Robotics at the School of Computer Science, University of Birmingham and Professor of Computer and Information Science at the University of Ljubljana. Previously, he held research positions at GRASP Laboratory at the University of Pennsylvania and at Vienna University of Technology. He is a Fellow of the IAPR.



Gregory Slabaugh is Professor of Computer Vision and AI and Director of the Digital Environment Research Institute (DERI) at Queen Mary University of London. Previous appointments include Huawei, City University of London, Medicsight, and Siemens. He received the PhD degree in Electrical Engineering from Georgia Institute of Technology.

APPENDIX A: UNSTABLE SEQUENCES

We reported in Section 3.1 that VDnCNN-feat, VDnCNN-RLSP and VResNet-feat are unstable on long video sequences at inference time. To illustrate their behaviour in more details, we apply them in Figure 6 to the 4 sequences of 1600 frames already used in Figure 1. We can make a number of observations:

- 1) Each model seems to be characterized by a distinct “instability pattern” that appears at random locations and grows locally until the entire frame is covered.
- 2) VDnCNN-RLSP and VResNet-feat are unstable on all four sequences after around 120 frames only. VDnCNN-feat is unstable on sequences 1 and 2 after ten times more frames, and it is stable on sequences 3 and 4. In this context, it would be easy—but dangerous—to mistake VDnCNN-feat for a stable model.
- 3) Motion is not necessary to trigger instabilities. The beginning of sequence 2 consists in the static title “fall’s arrival” and we see that this is enough to trigger an instability on VResNet-feat.

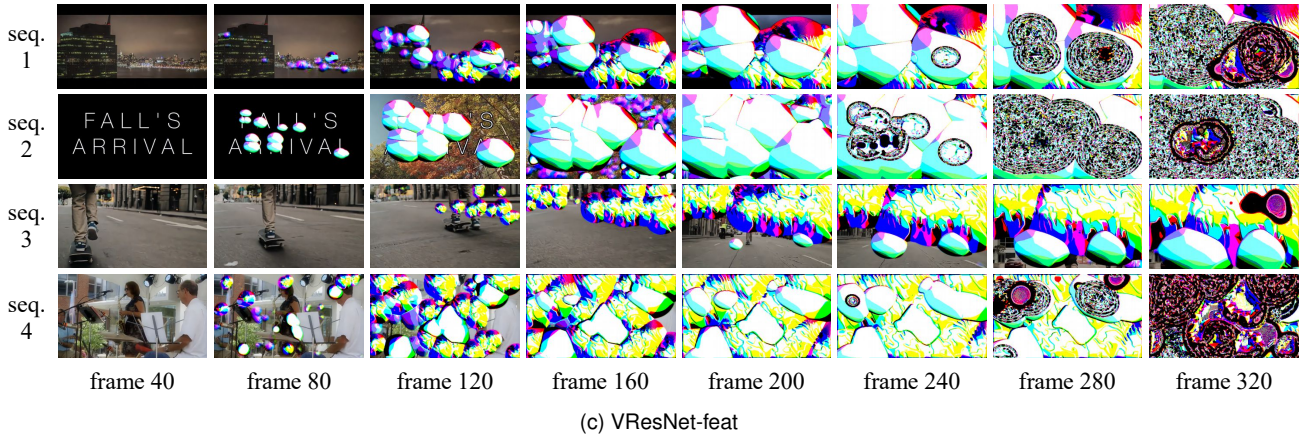
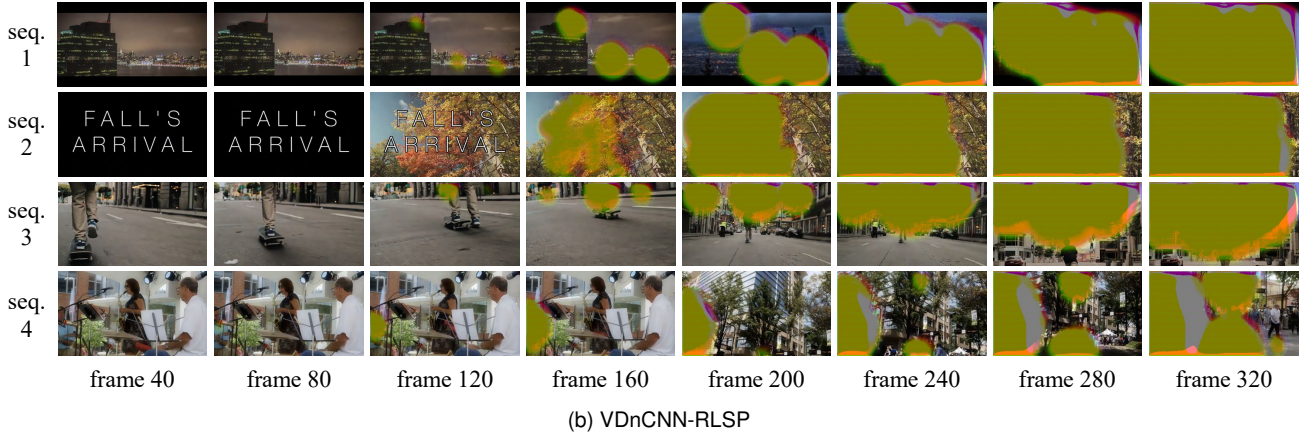
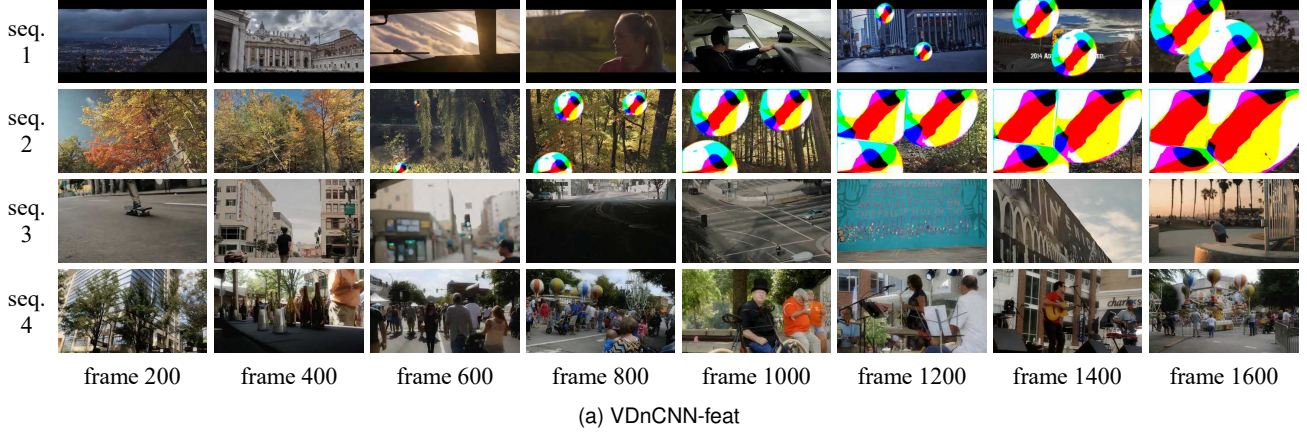


Fig. 6. Images generated by the three unstable recurrent video denoisers studied in Section 3.1, when applied to four sequences of 1600 frames downloaded from vimeo.com. VDnCNN-RLSP and VResNet-feat are unstable on all four sequences after around 120 frames only. VDnCNN-feat is unstable on sequences 1 and 2 after 1200 frames and 800 frames respectively, and it is stable on sequences 3 and 4.

APPENDIX B: FEATURE DAMPENING

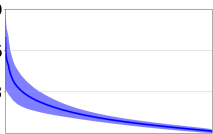

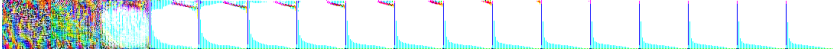
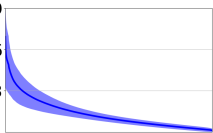

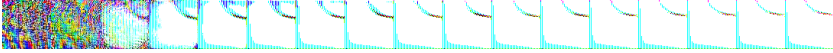
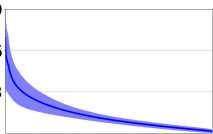

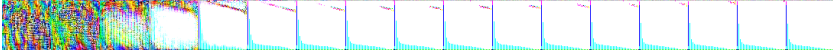
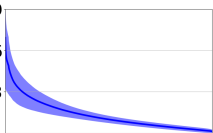


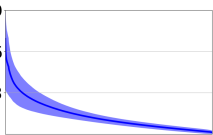

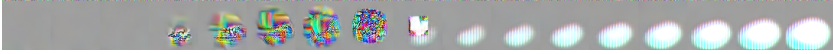
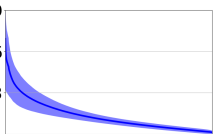

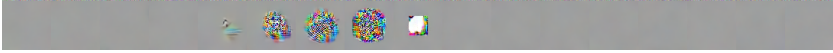
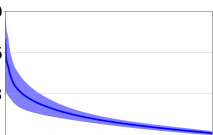

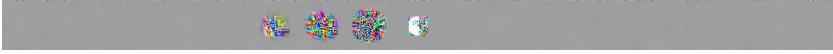
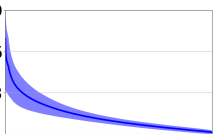


Given a trained model ϕ_w with Lipschitz constant L , one brute-force approach to enforce $L < 1$ is to reduce the magnitude of the recurrent weights $\mathbf{K} \leftarrow \lambda \mathbf{K}$ for some $\lambda < 1$. Interestingly, this is equivalent to reducing the magnitude of the recurrent features $h_{t-1} \leftarrow \lambda h_{t-1}$ in the convolutions:

$$(\lambda \mathbf{K}) * \mathbf{y}_{t-1} = \mathbf{K} * (\lambda \mathbf{y}_{t-1}).$$

For this reason, we refer to this approach as *feature dampening*. Table 5 shows the effect of feature dampening on the number of instabilities measured on a long sequence, and on the spatio-temporal receptive field. VResNet-feat is stable for $\lambda \leq 0.55$, and the recurrence is turned off completely for $\lambda = 0.0$.

TABLE 5

Feature Dampening by a factor λ with VResNet-feat. For each model, we show the performance on the 7th frame of the Vimeo-90k validation dataset (PSNR_7), the 1st and 9th deciles of the instability onsets on a sequence of about 2h30min (∞ means no instabilities observed). We also show the singular value spectrum averaged over the convolutions of the model, computed using the code from [47], and the temporal receptive field computed using our method.

model	PSNR ₇ 1 st dec. 9 th dec.	Average Singular Value Spectrum	Spatio-Temporal Receptive Field										
$\lambda = 1.0$	35.74 29 75		t	-40	-30	-20	-10	0	+10	+20	+30	+40	
			X										
			Y										
$\lambda = 0.95$	35.31 30 93		t	-40	-30	-20	-10	0	+10	+20	+30	+40	
			X										
			Y										
$\lambda = 0.85$	34.93 38 257		t	-40	-30	-20	-10	0	+10	+20	+30	+40	
			X										
			Y										
$\lambda = 0.75$	34.78 59 1024		t	-40	-30	-20	-10	0	+10	+20	+30	+40	
			X										
			Y										
$\lambda = 0.65$	34.69 181 38097		t	-40	-30	-20	-10	0	+10	+20	+30	+40	
			X										
			Y										
$\lambda = 0.55$	34.62 ∞ ∞		t	-40	-30	-20	-10	0	+10	+20	+30	+40	
			X										
			Y										
$\lambda = 0.45$	34.56 ∞ ∞		t	-40	-30	-20	-10	0	+10	+20	+30	+40	
			X										
			Y										
$\lambda = 0.0$	34.32 ∞ ∞		t	-40	-30	-20	-10	0	+10	+20	+30	+40	
			X										
			Y										

APPENDIX C: TRAINING ON LONG SEQUENCES

In their brief discussion of the instabilities affecting their model, Godard et al. [10] suggested that they were due to an inability of the recurrent models to generalize beyond their training length sequences. To test this hypothesis however, we face computational and data constraints. It is unrealistic to train large recurrent video denoising models on sequences of more than 10 to 20 frames—the training process involves backpropagation through time, which has large memory requirements—and even if it was possible, collecting the required data would quickly become impractical. To work around these issues, we perform experiments on a small VDnCNN model where the number of internal convolutions has been reduced to only 1. This allows us to unroll the model up to 56 times through time during training. We also generate long sequences with synthetic motion from single frames in Vimeo-90k using the technique described in [67]. We train our models on gray-scale patches of 32×32 pixels using Gaussian noise with standard deviation $\sigma = 20$, for 300k training steps. We show in Figure 7 the training curves of four models trained on sequences of 7, 14, 28 and 56 frames. Their profile is similar for the first three models but we observe sharp drops in the training curve of the model trained on sequences of 56 frames, likely due to the onset of instabilities during training resulting in gradient explosions. In Table 6, we report the performance and stability of each model. Even the model trained on sequences of 56 frames is vulnerable to instabilities.

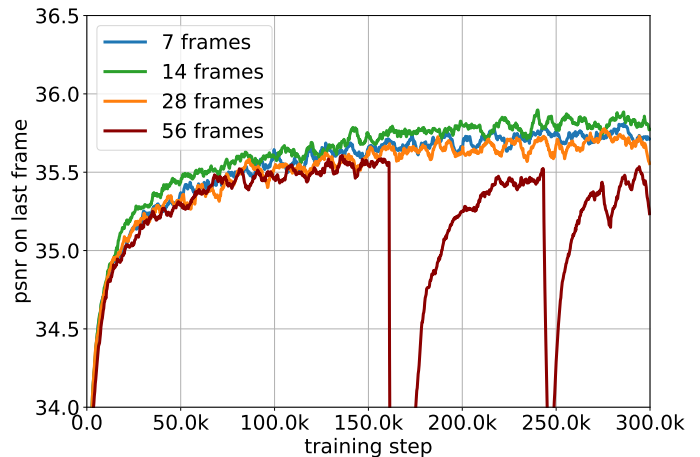
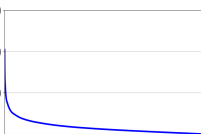

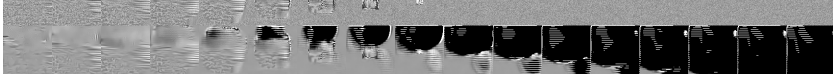
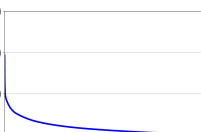


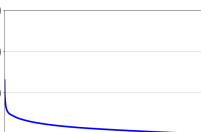

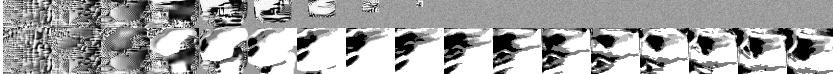
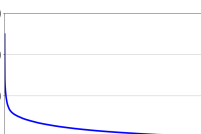



Fig. 7. Validation PSNR (on synthetic motion sequences) as a function of the training step for four models trained on sequences of varying lengths. The model trained on sequences of 56 frames suffers from training instabilities.

TABLE 6

Influence of the length of the training sequence. For each model, we show the performance on the 7th frame of the Vimeo-90k validation dataset (PSNR₇), the 1st and 9th deciles of the instability onsets on a sequence of about 2h30min (∞ means no instabilities observed). We also show the singular value spectrum averaged over the convolutions of the model, computed using the code from [47], and the temporal receptive field computed using our method.

model	PSNR ₇ 1 st dec. 9 th dec.	Average Singular Value Spectrum	Spatio-Temporal Receptive Field										
7 frames	34.72 57 75		t	-40	-30	-20	-10	0	+10	+20	+30	+40	
			X										
			Y										
14 frames	34.78 50 7121		t	-40	-30	-20	-10	0	+10	+20	+30	+40	
			X										
			Y										
28 frames	34.73 57 234		t	-40	-30	-20	-10	0	+10	+20	+30	+40	
			X										
			Y										
56 frames	34.58 261 12356		t	-40	-30	-20	-10	0	+10	+20	+30	+40	
			X										
			Y	