

End-to-End Trainable Video Super-Resolution Based on a New Mechanism for Implicit Motion Estimation and Compensation

Xiaohong Liu[†], Lingshi Kong[†], Yang Zhou[‡], Jiying Zhao[‡], Jun Chen[†]

[†] McMaster University, Ontario, Canada, [‡] University of Ottawa, Ontario, Canada

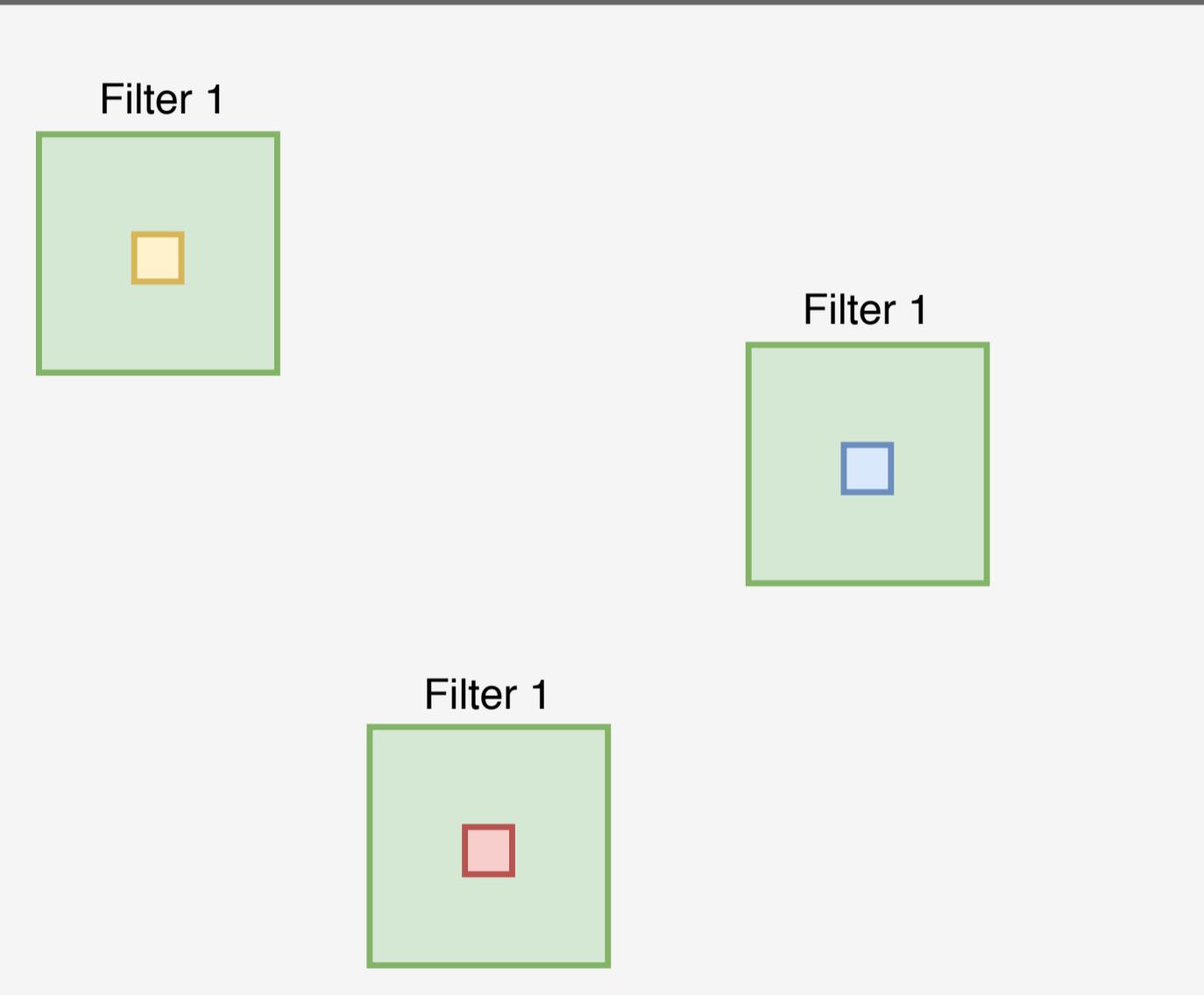


Introduction

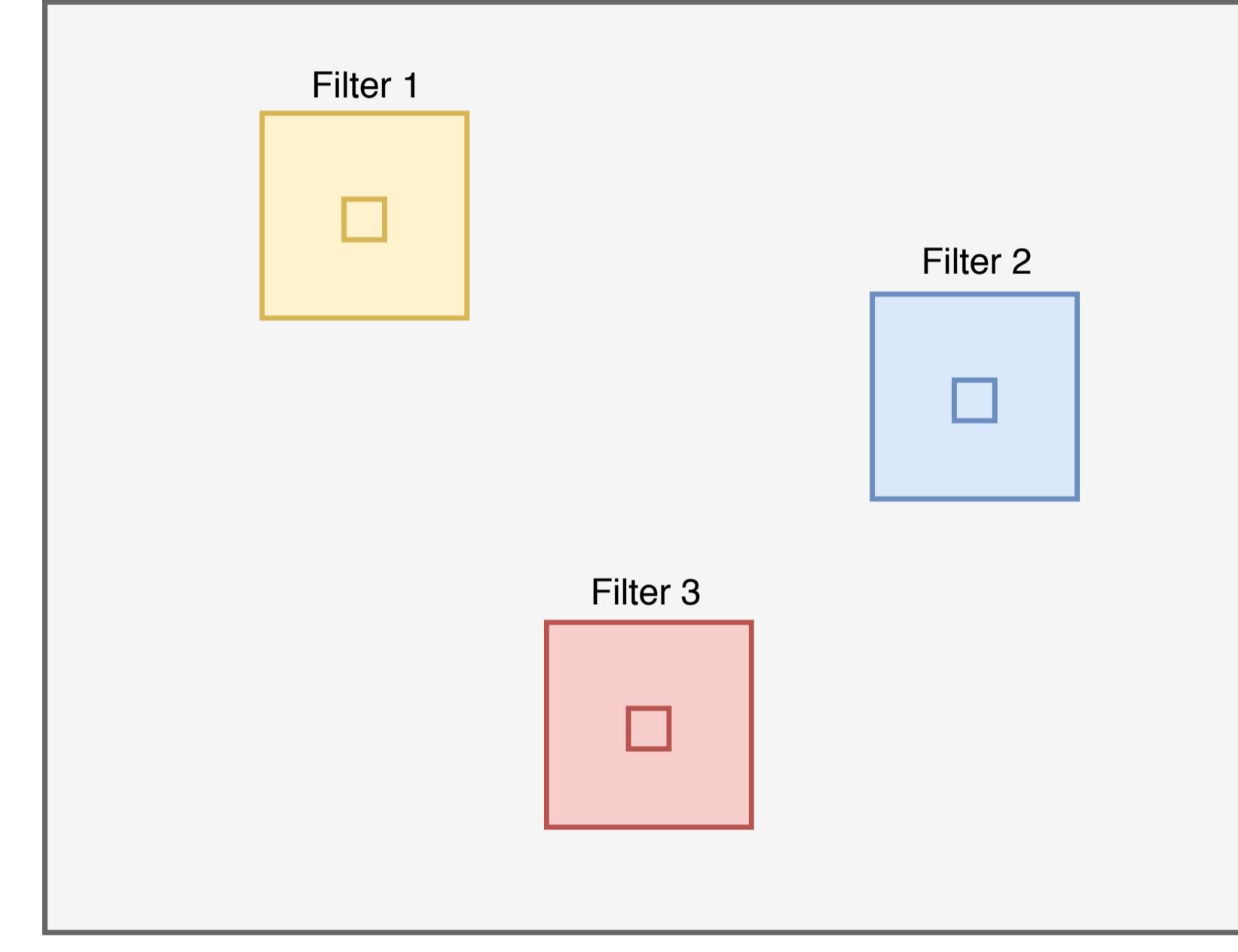
Video super-resolution aims at generating a high-resolution video from its low-resolution counterpart. With the rapid rise of deep learning, many recently proposed video super-resolution methods use convolutional neural networks in conjunction with explicit motion compensation to capitalize on statistical dependencies within and across low-resolution frames. Two common issues of such methods are noteworthy. Firstly, the quality of the final reconstructed HR video is often very sensitive to the accuracy of motion estimation. Secondly, the warp grid needed for motion compensation, which is specified by the two flow maps delineating pixel displacements in horizontal and vertical directions, tends to introduce additional errors and jeopardize the temporal consistency across video frames. To address these issues, we propose a novel dynamic local filter network to perform implicit motion estimation and compensation by employing, via locally connected layers, sample-specific and position-specific dynamic local filters that are tailored to the target pixels. We also propose a global refinement network based on ResBlock and autoencoder structures to exploit non-local correlations and enhance the spatial consistency of super-resolved frames. The experimental results demonstrate that the proposed method outperforms the state-of-the-art and validate its strength in terms of local transformation handling, temporal consistency as well as edge sharpness.

Proposed Method

In standard CNN layers, the filters (shown as green color) applied on three different pixel positions (shown as yellow, blue and red colors) have the same weights. In comparison, for LC layers, three different filters (shown as yellow, blue and red colors) are applied on the relevant pixel positions with unshared weights that are generated locally according to the imposed pixels.

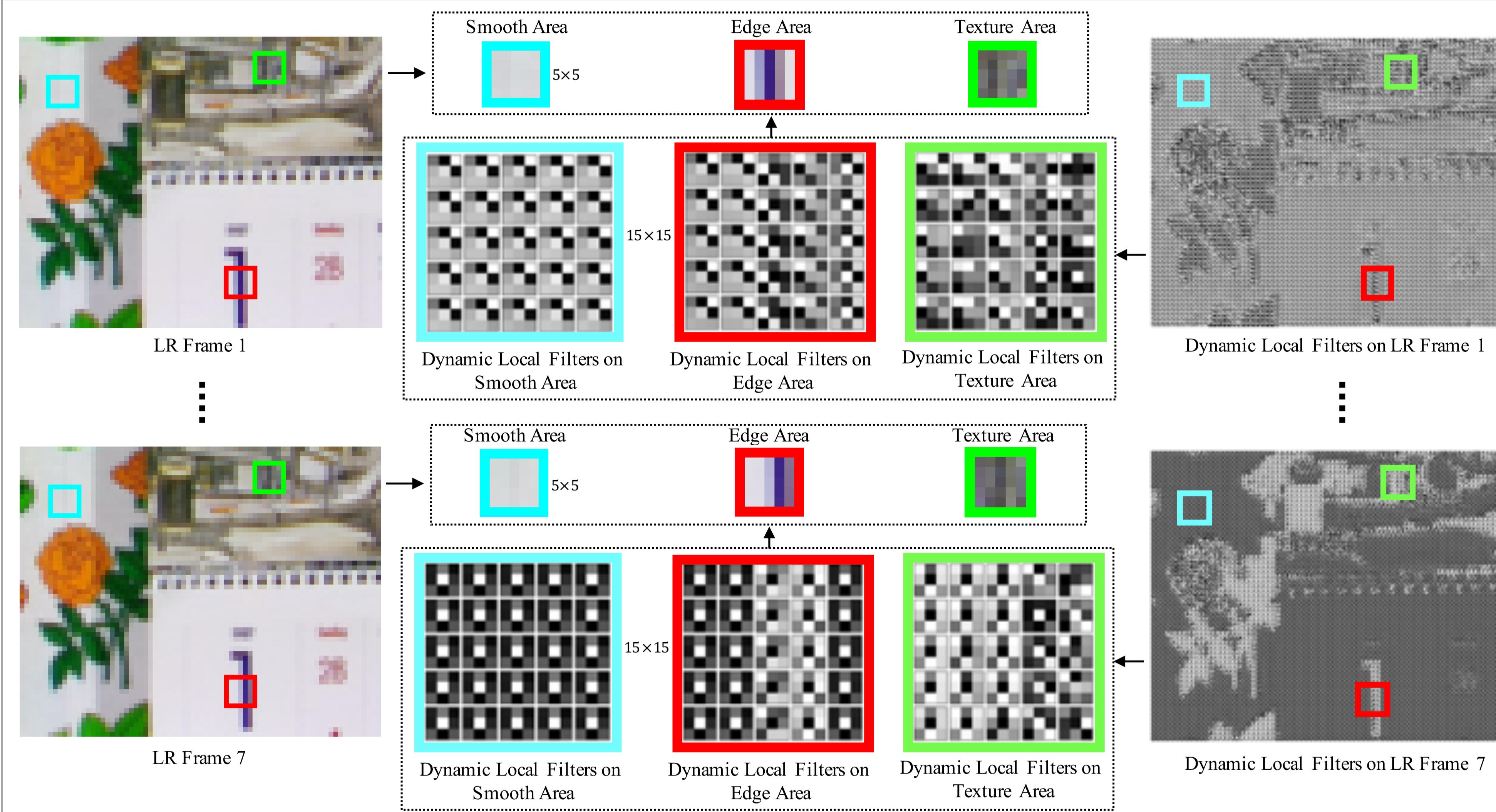


(a) Standard CNN layers



(b) Locally connected layers

Visualization of Dynamical Local Filters



The formulation of Dynamic Local Filter Network:

$$\hat{Y}_l^{(i,j)} = \sum_{m=i-d}^{i+d} \sum_{n=j-d}^{j+d} \sum_{k=t-T}^{t+T} \Theta_{i,j,l}^{(m-i+d+1, n-j+d+1, k-t+T+1)} \cdot Y_k^{(m,n)}$$

The produced l th feature map

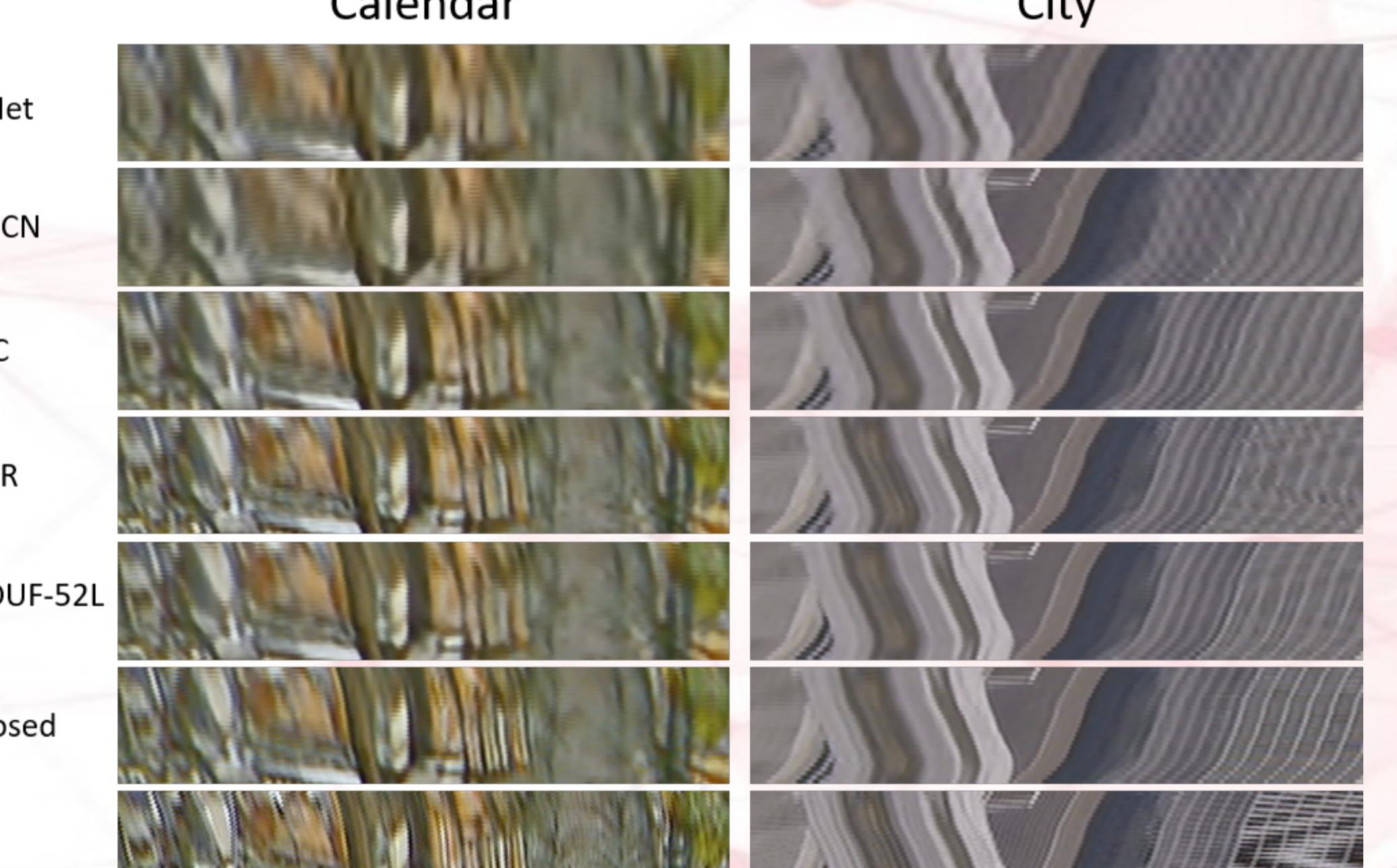
Dynamic local filters The k th input LR frame

Experimental Results

Qualitative comparisons on Vid4 dataset:



Analysis of Temporal Consistency:



Quantitative comparisons and ablation study:

Vid4	Metric	Bicubic	Bayesian [24]	VSRNet [18]	VESPCN [1]	$B_{1,2,3} + T$ [25]	SPMC [33]	FRVSR [31]	DUF-16L [17]	DUF-52L [17]	Proposed
x3	PSNR	25.28	25.82	26.79	27.25	-	27.49	-	28.90	-	29.51
	SSIM	0.7329	0.8323	0.8098	0.8447	-	0.8400	-	0.8898	-	0.8964

Table 1: Quantitative comparisons on the Vid4 dataset with $r = 3, 4$.

Vid4	Metric	SPMC [33]	DUF-16L [17]	DUF-52L [17]	Proposed
x3	PSNR	32.10	-	-	33.91
	SSIM	0.9000	-	-	0.9358

Table 2: Quantitative comparisons on the SPMCS dataset with $r = 3, 4$.

Vid4	Metric	w/o LC Layers	w/o GRN	w/U-Net [30]	Our full model
x3	PSNR	27.27	28.13	29.20	29.51
	SSIM	0.8471	0.8752	0.8896	0.8964

Table 3: Quantitative comparisons on the Vid4 dataset for different variants of the proposed method with $r = 3, 4$.