

Improving Feature Stability during Upsampling – Spectral Artifacts and the Importance of Spatial Context

Shashank Agnihotri¹, Julia Grabinski^{1,2,3}, and Margret Keuper^{1,4}

¹ Data and Web Science Group, University of Mannheim

² Fraunhofer ITWM, Kaiserslautern

³ Institute for Machine Learning and Analytics (IMLA), Offenburg University

⁴ Max-Planck-Institute for Informatics, Saarland Informatics Campus

{shashank.agnihotri, julia.grabinski, keuper}@uni-mannheim.de

Abstract. Pixel-wise predictions are required in a wide variety of tasks such as image restoration, image segmentation, or disparity estimation. Common models involve several stages of data resampling, in which the resolution of feature maps is first reduced to aggregate information and then increased to generate a high-resolution output. Previous works have shown that resampling operations are subject to artifacts such as aliasing. During downsampling, aliases have been shown to compromise the prediction stability of image classifiers. During upsampling, they have been leveraged to detect generated content. Yet, the effect of aliases during upsampling has not yet been discussed w.r.t. the stability and robustness of pixel-wise predictions. While falling under the same term (*aliasing*), the challenges for correct upsampling in neural networks differ significantly from those during downsampling: when downsampling, some high frequencies can not be correctly represented and have to be removed to avoid aliases. However, when upsampling for pixel-wise predictions, we actually require the model to restore such high frequencies that can not be encoded in lower resolutions. The application of findings from signal processing is therefore a necessary but not a sufficient condition to achieve the desirable output. In contrast, we find that the availability of large spatial context during upsampling allows to provide stable, high-quality pixel-wise predictions, even when fully learning all filter weights.

1 Introduction

Most computer vision models addressing perceptual tasks such as image restoration [16, 89], semantic segmentation [7, 35, 70], optical flow estimation [20, 43, 83] and disparity estimation [6, 11, 50] in realistic scenarios are required to behave in a stable way, at least under mild corruptions. Interestingly, for the slightly simpler task of image classification, recent progress has shown that a model’s robustness does not only depend on its training but also on its architecture [29–32, 41, 42, 45, 58, 91, 95]. Specifically, *aliasing*, *i.e.* spectral artifacts that emerge from naïve image resampling, have shown to compromise prediction stability, in particular in the context of classical convolutional models [33, 37, 48, 52, 69, 80, 82] which predominantly use small filter kernels in combination with severe data aggregation during downsampling [30, 58]. Principled cures

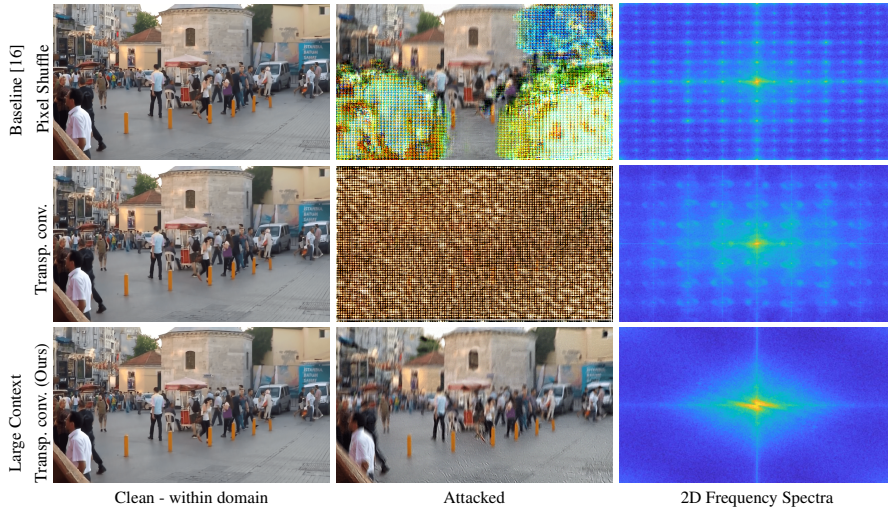


Fig. 1: Image restoration example using NAFNet [16] variants on GoPro [63]. Upsampling techniques like Pixel Shuffle [77] (first row) and transposed convolution [22] using small learnable filters (2×2 or 3×3) (second row) are used by most prior art. Both lead to spectral artifacts for which the model needs to compensate. The clean (in-domain) restored images look appealing - while adversaries (here 5-step PGD [49] attack) can leverage aliases such that artifacts become easily visible. When observed in the frequency domain, they manifest as repeating peaks all over the spectra. Based on sampling theoretic considerations, we propose Large Context Transposed Convolutions (7×7 or larger) (bottom row). They significantly increase the model’s stability during upsampling, observable in the restored image under attack and the frequency spectrum.

usually refer to basic concepts from signal processing such as anti-aliasing by blurring before downsampling [29, 91]. While this discussion on classifier (*i.e.* encoder) networks is insightful, it does not provide a recipe to counteract aliases emerging during *upsampling* for pixel-wise prediction tasks such as image restoration. Specifically, naïve upsampling introduces artifacts in the feature representation, such as grid artifacts [4, 65] or ringing artifacts [62]. As shown in Fig. 1, these artifacts, an inherent property of inadequate upsampling (refer Sec. 3) are not always visible to the human eye, are accentuated under adversarial attack such that they can also be seen with a human eye. We leverage this effect in our analysis. When observed in the frequency domain, these artifacts are apparent as multiple peaks, *i.e.* *aliases* of the original data.

While for downsampling, signal processing laws basically prescribe which part of the information can be retained at lower resolutions without aliases [76], “correct”, alias-free upsampling can not restore the original high-resolution information. Thus, learning to upsample feature maps such that the feature stability is not harmed is of paramount importance. In this paper, we therefore first provide a synopsis of different aliases that emerge from different upsampling techniques. Based on this work, we propose a simple, transposed convolution-based upsampling block. We study our proposed operation in the context of various models, from image restoration [16, 89] over semantic segmentation [70] to disparity estimation [50].

Our main contributions can be summarized as follows:

- Motivated by sampling theory [76], we study upsampling in models for diverse pixel-wise prediction tasks. We find that the availability of large kernels in transposed convolutions helps the feature stability and significantly improves over standard, small kernel transposed convolutions as well as pixel shuffle [77].
- While large kernels are required to allow for reduced aliasing and to provide the necessary spatial context for increasing the resolution, additional small kernels can add details and remain useful.
- We provide empirical evidence for our findings on diverse architectures (including vision transformer-based architectures) and downstream tasks such as image restoration, semantic segmentation, and depth estimation.
- We show empirically that our proposed upsampling operation complements other feature stability-increasing approaches like adversarial training.

2 Related Work

In the following, we discuss recent challenges for neural networks regarding artifacts introduced by spatial sampling methods [4, 62, 65]. Further, we review related work on the most recent use of large kernels in CNNs. Finally, we provide an overview of adversarial attacks to gauge the quality of representations learned by a network.

Spectral Artifacts. Several prior works have studied the effect of downsampling operations on model robustness, *e.g.* [2, 29, 30, 42, 46, 91, 95]. Inspired by [30], [29] propose an aliasing-free downsampling in the frequency domain which translates to an infinitely large blurring filter before downsampling in the spatial domain. Thus, for image classification, using large filter kernels has been shown to remove artifacts from downsampled representations and it leads to favorable robustness in all these cases [30, 42, 46]. However, all these works focus on improving the properties of encoder networks.

Models that use transposed convolutions in their decoders⁵ are widely used for tasks like image generation [27, 68] or segmentation [7, 55, 64, 70]. However, in simple transposed convolutions, the convolution kernels overlap based on the chosen stride and kernel size. If the stride is smaller than the kernel size, this will cause overlaps in the operation, leading to uneven contributions to different pixels in the upsampled feature map and thus to grid-like artifacts [4, 65]. Further, image resampling can lead to aliases that become visible as ringing artifacts [76]. In the context of deepFake detection, image generation, and deblurring, several works analyzed [14, 18, 21, 23, 38, 44, 47] and improved upsampling techniques [26, 46, 78, 87] to reduce visual artifacts.

Some architectures like PSPNet [93], PSANet [94], or PSMNet [15] simply use bilinear interpolation operations for upsampling the feature representations. While this reduces grid artifacts as bilinear interpolation smoothens out the feature maps, it also has major drawbacks as they sample incorrectly. These new artifacts are sometimes visible as overly smooth predictions, in particular, apparent in the PSPNet segmentation masks. The segmentation masks over-smooth around edges and often miss out on thin

⁵ For more details on Transposed Convolutions refer to [22].

details (predictions showing these are included in the Appendix B.4). This observation already shows why image encoding and decoding have to be considered separately when it comes to sampling artifacts. While during encoding, artifacts can be reduced by blurring, the main purpose of decoder networks is *reducing* blur in many applications, to create fine-granular, pixel-wise accurate outputs, which our approach facilitates.

Large Kernels. For image classification, [54] showed that using large kernels like 7×7 in the CNN convolution layer can outperform self-attention based vision transformers [53, 84]. In [17, 33, 35, 51, 66], the receptive field of the convolution operations was further expanded by using larger kernels, up to 31×31 and 51×51 . These larger receptive fields provide more context to the **encoder**, leading to better performance on classification, segmentation, or object detection tasks. [17, 51] use a small kernel in parallel to capture the local context along with the global context. In contrast to these works, which are limited to exploring increased context only during encoding, we investigate if larger kernels can benefit upsampling when considering pixel-wise prediction tasks such as image restoration or segmentation.

Adversarial Attacks. The purpose of adversarial attacks is to reveal neural networks' weaknesses [3, 30, 74, 81] by perturbing pixel values in the input image [12, 28, 49]. These perturbations should lead to a false prediction even though the changes are hardly visible [28, 61, 81]. Especially attacks that have access to the network's architecture and weights, so-called white-box attacks, are a common approach to analyzing weaknesses within the networks' structure [12, 28]. They employ the gradient of the network to optimize the perturbation, which is bounded within an ϵ -ball of the original image, *i.e.* ϵ defines the strength of the attack. Most adversarial attacks are proposed to attack classification networks like the one-step Fast Gradient Sign Method (FGSM) [28] or the multi-step Projected Gradient Descent (PGD) attack [49]. However, they can be adapted to other tasks as *e.g.* in [59, 67, 88]. Furthermore, there are dedicated methods like SegPGD [34] for attacks on semantic segmentation models or PCFA [74] and [71, 73] for optical flow models and CosPGD [3] and others [72] for other pixel-wise prediction tasks. We evaluate the stability of upsampled features using adversarial attacks such as PGD and CosPGD for image restoration and FGSM and SegPGD for segmentation.

3 Spectral Upsampling Artifacts and How They Can Be Reduced

Following, we first theoretically review artifacts that are caused during upsampling from a signal processing aspect. We start by describing the spectral artifacts [76] induced by the bed of nails interpolation, similar to the discussion in [23], and then extend the theoretical analysis to further upsampling schemes. Second, we derive from this analysis two hypotheses for the prediction stability of encoder-decoder networks, depending on their architecture. These hypotheses will motivate the remainder of the manuscript.

Consider, w.l.o.g., a one-dimensional signal I and its discrete Fourier Transform $\mathcal{F}(I)$ with k being the index of discrete frequencies

$$\mathcal{F}(I)_k = \sum_{j=0}^{N-1} e^{-2\pi i \cdot \frac{jk}{N}} \cdot I_j, \quad \text{for } k = 0, \dots, N-1.$$

During decoding, we need to upsample the spatial resolution of I to get I^{up} . For example for an upsampling factor of 2 (often used in DNNs [1, 16, 19, 82, 89]) we have for $\bar{k} = 0, \dots, 2N - 1$

$$\mathcal{F}(I)_{\bar{k}}^{\text{up}} = \sum_{j=0}^{2N-1} e^{-2\pi i \cdot \frac{j\bar{k}}{2N}} \cdot I_j^{\text{up}} = \sum_{j=0}^{N-1} e^{-2\pi i \cdot \frac{2j\bar{k}}{2N}} I_j + \sum_{j=0}^{N-1} e^{-2\pi i \cdot \frac{(2j+1)\bar{k}}{2N}} \bar{I}_j, \quad (1)$$

where $\bar{I}_j = 0$ in **bed of nails interpolation**. Therefore, the second term in (1) can be dropped and the first term resembles the original $\mathcal{F}(I)$. Equivalently, we can rewrite Eq. (1), for $\bar{I}_j = 0$, using a Dirac impulse comb as

$$(1) = \sum_{j=0}^{2N-1} e^{-2\pi i \cdot \frac{j\bar{k}}{2N}} \cdot \sum_{t=-\infty}^{\infty} I_j^{\text{up}} \cdot \delta(j - 2t). \quad (2)$$

If we now apply the pointwise multiplication with the Dirac impulse comb as convolution in the Fourier domain (assuming periodicity) [25], it is

$$\begin{aligned} \mathcal{F}(I)_{\bar{k}}^{\text{up}} &= \frac{1}{2} \sum_{t=-\infty}^{\infty} \left(\sum_{j=-\infty}^{\infty} e^{-2\pi i \cdot \frac{j\bar{k}}{2N}} I_j^{\text{up}} \right) \left(\bar{k} - \frac{t}{2} \right) \\ &\stackrel{(1)}{=} \frac{1}{2} \sum_{t=-\infty}^{\infty} \left(\sum_{j=-\infty}^{\infty} e^{-2\pi i \cdot \frac{j\bar{k}}{2N}} \cdot I_j \right) \left(\bar{k} - \frac{t}{2} \right) = \frac{1}{2} \sum_{t=-\infty}^{\infty} \mathcal{F}(I)_{\bar{k}} \left(\bar{k} - \frac{t}{2} \right). \end{aligned} \quad (3)$$

We can see that such upsampling creates high-frequency replica of the signal at $\frac{t}{2}$ for t in $-\infty, \dots, \infty$ in $\mathcal{F}(I)^{\text{up}}$ and spatial frequencies apparent beyond array positions $\frac{N}{2}$ will be impacted by spectral artifacts if no appropriate countermeasures are taken.

A standard countermeasure is interpolation of the inserted values with $\bar{I}_j = \frac{I_{j-1} + I_j}{2}$ for **linear interpolation** in Eq. (1). Linear interpolation (and in consequence bi-linear interpolation in 2D signals) corresponds to a convolution with a triangular impulse with width 2, which can be represented as the convolution of two rectangle functions with width 1. Accordingly, the Fourier response for frequency ℓ , \mathcal{F}_{ℓ} of the triangular impulse is a squared sinc function ($\text{sinc}^2(\ell)$) with $\text{sinc}(\ell) = \frac{\sin(\pi\ell)}{\pi\ell}$. Since the output signal after interpolation is still discrete, i.e. sampled with sampling rate $\frac{1}{2}$, a replica of the interpolation function, the sinc^2 function, will appear with rate 2 in the resulting spectrum (see also Fig. 2). The resulting interpolated signal is not optimal for several reasons. Most importantly, the spectrum of the interpolation function is not flat although the estimated values appear overly smooth (see Fig. 3.). This is arguably suboptimal for, for example, image restoration or segmentation tasks, where fine structural details are supposed to emerge in the upsampled data.

Note that, in Eq. (1), **pixel shuffle** [77] will set \bar{I}_j to completely unrelated values of a different feature map channel, leading to a highly non-smooth signal with frequencies at the band limit. The resulting issues in the spectrum are similar to the ones caused by the bed of nails interpolation. These spectral artifacts can be visually observed in Fig. 3.

Therefore, in **transposed convolutions**, the interpolation function is not fixed to a predefined smoothing kernel but learned so that the resulting signal can represent fine

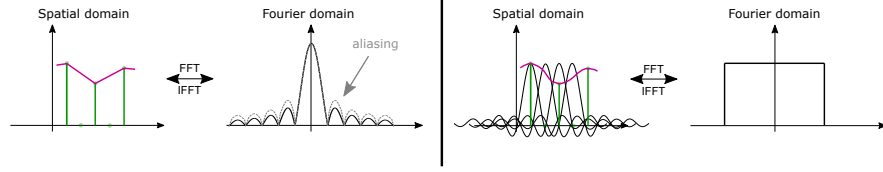


Fig. 2: (Left) Linear interpolation (pink) of the samples (green) causes aliases. (Right) Optimal signal reconstruction (pink) is achieved by sinc interpolation. In practice our spatial context is limited and the interpolation function is discrete. Yet, increasing the kernel size enables the approximation of larger sinc-like structures.

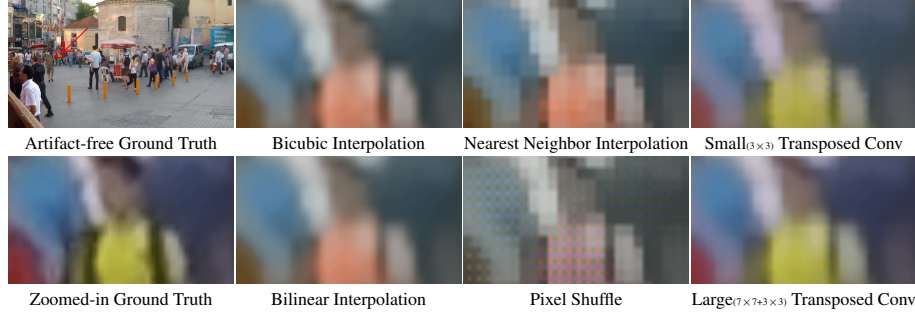


Fig. 3: An image from GoPro [63] downsampled with 3×3 MaxPooling and then upsampled using various upsampling techniques. The resulting artifacts are compared on **zoomed-in red box regions** for better visibility. Bilinear interpolation causes over-smoothing. Bicubic interpolation causes overestimation along image boundaries while Pixel Shuffle and Nearest Neighbor cause strong grid artifacts along with discoloration. Small kernel transposed convolutions cause grid artifacts, however, on increasing kernel size we start getting better upsampling.

details after the initial bed of nails interpolation and potentially learn to add fine details. One issue is that the learned convolution kernels may overlap based on the chosen stride and kernel size. If the stride is smaller than the kernel size, this will cause overlaps in the operation, leading to uneven contributions to different pixels in the upsampled feature map and thus to grid-like artifacts [4, 65]. Besides this rather technical aspect, transposed convolutions, if sufficiently large (thus also containing more context), could in principle learn to approximate correct upsampling functions. This can be understood when again looking at the Fourier representation. When interpolating, we want to increase the signal array size so that all the original information is preserved and the model can easily learn additional details. Such upsampling to preserve the information from the original low-resolution data is most easily achieved by transforming the signal to the Fourier domain, then padding the missing high-frequency parts with zeros and transforming the resulting array back to the spatial domain [79]. In the Fourier domain, this padding operation can be understood as a point-wise multiplication of the desired full spectrum with a rectangle function with width N (denoted rect_N). Conversely, this operation corresponds to a convolution with $\mathcal{F}^{-1}(\text{rect}_N) = \frac{1}{N}\text{sinc}(xN)$ in the spatial domain. While the *sinc* function drops off as x increases, it never drops to zero. When applied for interpolation, its crests and the troughs cancel out the aliasing to a large

extent as shown in Fig. 2. Thus, in order to allow the approximation of the optimal interpolation function, the kernel size in transposed convolutions has to be chosen as large as possible. This is, however, at odds with the “learnability” of suitable filter weights. Note that for pixel-wise predictions, models not only need to correctly interpolate, but they also need to “fill in” the missing details, which requires global as well as local context. Therefore, we expect a trade-off on the kernel size of transposed convolutions, where larger kernels improve the stability of the upsampled features and thus can reduce artifacts while the absolute prediction quality can suffer from very large learnable kernels. Sufficiently but not overly large kernels provide sufficient spatial context and are appropriate to allow for the model to learn when to blur and when to preserve/sharpen upsampled features. We illustrate this in Fig. 12 in Appendix C.4.

From this theoretical analysis of common upsampling methods, we derive the following hypotheses that we deem relevant for encoder-decoder architectures:

Hypothesis 1 (H1): *Large Context Transposed Convolutions (LCTC)* *i.e. Large kernels in transposed convolution operations provide more context and reduce spectral artifacts and can therefore be leveraged by the network to facilitate better and more robust pixel-wise predictions.*

Hypothesis 2 (H2, Null Hypothesis): *To leverage prediction context and reduce spectral artifacts, it is crucial to increase the size of the transposed convolution kernels (upsample using large filters). Increasing the size of normal (i.e. non-upsampling) decoder convolutions does not have this effect.*

In the following, we show the proposed, simple, and principled architecture changes that allow for studying the above hypotheses and improving robustness by improving feature stability.

4 Upsampling using Large Context Transposed Convolutions

Driven by the observations on upsampling artifacts, we investigate the advantage of larger kernel sizes during upsampling, for applications such as semantic segmentation or disparity estimation. Therefore, we keep the models’ encoder part fixed and exclusively change operations in the architecture of the decoder part of the model. There, we have two design choices: **Upsampling** – The kernel size for the transposed convolution operations that learn upsampling, and **Decoder Block** – The kernel size in the convolution operations of blocks that learn to decode the features. Probing options for **Upsampling** works towards proving H1 while a combination of both options proves H2, *i.e.* shows that a pure increase in the decoder parameters does not have the desired effect. This is considered in our ablation study in Sec. 5.2.

Figure 4 summarizes the studied options for an abstract encoder-decoder architecture like [70]. The model decoder is depicted in the green box. Operations that we consider to be executed along the red upwards arrows (**Upsampling Operators**) are detailed in the top right part of the figure (operations a) to c)). Operations that we consider to be executed along the blue sideways arrows (**Decoder Building Blocks**) are depicted in the bottom right (operations d) to f)).

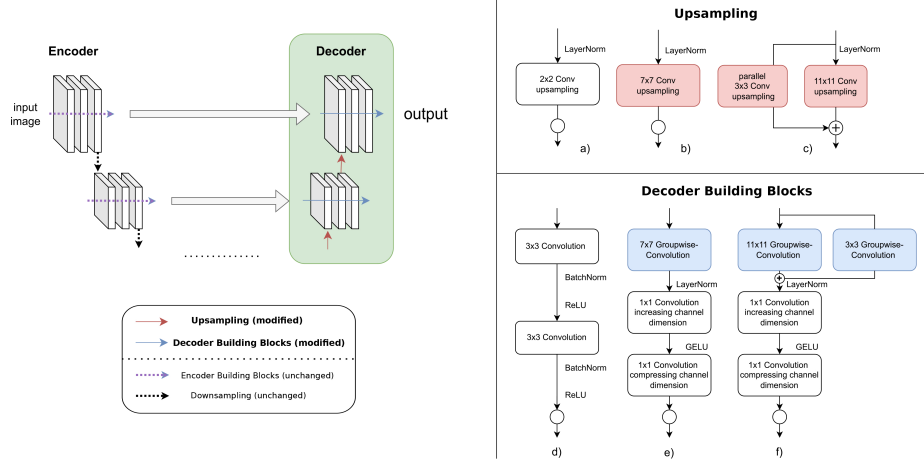


Fig. 4: Abstract representation of an encoder-decoder architecture. While for different tasks, the implementation of the model encoder varies (including transformer-based encoders), our study focuses on the *model decoder* (in green). The backbone for the decoder is commonly a ResNet-like structure for feature extraction [7, 70], additionally we also used a ConvNeXt-like [54] structure. We investigate variants of different upsampling operations (the operations along the **red arrows** in the decoder) for fixed decoder blocks. We consider, as a probe for H1, the baseline transposed deconvolution (a) in the top right), and for LTC an increased convolution kernel size (b) in the top right), and an increased convolution kernel with a second path using a small convolution kernel (c) in the top right). To test whether the plain increase in parameters is responsible for improved results (zero hypotheses, H2), we also ablate on the increase of convolution kernel size in the decoder block (operations along the **blue arrows** in the green block), as shown on the bottom right. We consider the common ResNet-like decoder building block structure (in d)) and two ConvNeXt-like structured backbones for the decoder building block in e) and f), where f) has an additional small convolution applied in parallel, analog to c).

Model Details. Here, we provide details on the studied models. All implementation details are given in the Appendix A.

Transposed Convolution Kernels for Upsampling. The upsampling operation is typically performed with small kernels (2×2 or 3×3) in the transposed convolution operations [8, 13, 70]. We aim to increase the spatial context during upsampling and to reduce grid artifacts. Thus we use **Large Context Transposed Convolutions (LTC)**. We either use 7×7 transposed convolutions or 11×11 transposed convolutions with a parallel 3×3 transposed convolution. Adding a parallel 3×3 kernel is motivated by [17], as large convolution kernels tend to lose local context, and thus adding a parallel small kernel helps to overcome this potential drawback (see Appendix B.3).

Decoder Building Blocks. To verify that the measurable effects are due to the improved upsampling and not due to merely increasing the decoder capacity, we ablate on decoder convolution blocks similar to convolution blocks used in the ConvNeXt [54] basic block for encoding. While the standard ConvNeXt block uses a 7×7 depth-wise convolution, we consider 7×7 and 11×11 group-wise convolutions, followed by layers present in a ConvNeXt basic block to analyze the importance of the receptive field within the block.

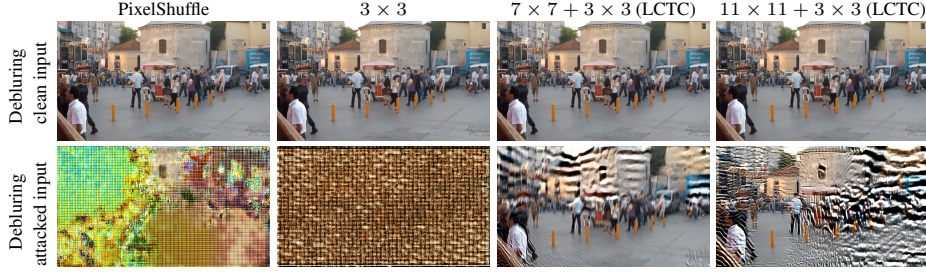


Fig. 5: NAFNet, as proposed, uses Pixel Shuffle for upsampling. We modify only the upsampling operations to transposed convolution with kernel size (3×3) and LCTC (Ours) for comparisons. We observe, for example, under a 10-step PGD attack with $\epsilon \approx \frac{8}{255}$ our proposed H1 gains validity. More examples for [16, 89] using different attacks and budgets are in Appendix C.3.

Table 1: Comparison of performances of different **upsampling** methods in *SotA* Image Restoration Networks on the GoPro dataset. The architectures use Pixel Shuffle for Upsampling, we propose replacing the Pixel Shuffle with Large Context Transposed Convolutions (LCTC). We report additional results using adversarial training in Tab. 15. Note, that some trade-off between clean performance and robustness is expected [85, 90].

Network	Upsampling Method	Test Accuracy		CosPGD ($\epsilon \approx \frac{8}{255}$) attack iterations						PGD ($\epsilon \approx \frac{8}{255}$) attack iterations					
		PSNR	SSIM	PSNR	SSIM	PSNR	SSIM	PSNR	SSIM	PSNR	SSIM	PSNR	SSIM	PSNR	SSIM
Restormer	Pixel Shuffle	31.99	0.9635	11.36	0.3236	9.05	0.2242	7.59	0.1548	11.41	0.3256	9.04	0.2234	7.58	0.1543
	Transposed Conv 3×3	9.68	0.095	8.24	0.0452	8.53	0.0628	8.44	0.0631	7.66	0.0464	7.72	0.0577	8.64	0.0527
	LCTC: $7 \times 7 + 3 \times 3$ (Ours)	29.51	0.9337	13.69	0.4186	11.53	0.3136	10.16	0.2484	13.69	0.4183	11.54	0.3137	10.16	0.2483
	LCTC: $11 \times 11 + 3 \times 3$ (Ours)	29.44	0.9324	14.65	0.4251	12.83	0.3438	11.48	0.29	14.65	0.4253	12.84	0.3445	11.48	0.2893
NAFNet	Pixel Shuffle	32.87	0.9606	8.67	0.2264	6.68	0.1127	5.81	0.0617	10.27	0.3179	8.66	0.2282	5.95	0.0714
	Transposed Conv 3×3	31.02	0.9422	6.15	0.0332	5.95	0.0258	5.87	0.0233	6.15	0.0332	5.95	0.0258	5.87	0.0234
	LCTC: $7 \times 7 + 3 \times 3$ (Ours)	31.12	0.9430	14.54	0.4827	11.05	0.3220	9.06	0.2213	14.53	0.4823	11.03	0.3201	9.08	0.2224
	LCTC: $11 \times 11 + 3 \times 3$ (Ours)	30.77	0.9392	14.34	0.4492	11.41	0.3244	9.54	0.2411	14.34	0.45	11.4	0.3236	9.55	0.2398

Figure 4 (bottom right e) and f)) shows the structure of a ConvNeXt-style building block used in our work. First, a group-wise convolution is performed, followed by a LayerNorm [5] and two 1×1 convolutions which, similar to [54], creates an inverted bottleneck by first increasing the channel dimension and after a GELU [40] activation compressing the channel dimension again. We consider the ResNet-style building block (Figure 4, d)), with 3×3 convolution, yet without skip connection, as our baseline when studying this architectural design choice.

5 Experiments

In the following, we evaluate the effect of the considered upsampling operators in several applications. We start by evaluating the effect on the upsampled feature stability of recent *state-of-the-art* (SotA) image restoration models [16, 89], then provide results on semantic segmentation using more generic convolutional architectures that allow us to provide compulsory ablations. Last, we show that our results also extend to disparity estimation [50]. We provide details on the used adversarial attacks, datasets, reported metrics, and other experimental details in Appendix A.

In all cases, we observe that Large Context Transposed Convolutions (LCTC) improve the results of the respective pixel-wise prediction task in terms of stability under attack, showing that H1 holds. Further, our extensive ablation on image segmentation shows that increasing the convolution kernel in the decoder building blocks does not have this beneficial effect, providing experimental evidence for our hypothesis H2 and confirming the impact of spectral artifacts on pixel-wise predictions.

5.1 Image Restoration

For image restoration, we consider the Vision Transformer-based Restormer [89] and NAFNet [16]. Both originally use the Pixel Shuffle [77] for upsampling. Here, we compare the reconstructions from these proposed architectures to their variants using the proposed operators with large transposed convolution filters. We use the same metrics as [16, 89], Peak Signal-to-Noise Ratio (PSNR), and structural similarity index measure (SSIM) [86]. We perform our experiments on the GoPro [63] image deblurring dataset, following the experimental setup in [1].

Results on Image Restoration. We first consider qualitative results on NAFNet [16] in Figure 5 and Restormer [89] in Fig. 10, Fig. 11 (in Appendix C.3), where we see that the proposed upsampling operators allow for visually good results in image deblurring on clean data (similar to pixel shuffle). Yet, in contrast to pixel shuffle and the baseline small transposed convolution filters, the proposed Large Context Transposed Convolutions (LCTC) significantly reduces artifacts that arise on attacked images (in this case, 10-step PGD with $\epsilon \approx \frac{8}{255}$), attacks with varying numbers of steps.

In Table 1, we report the average PSNR and SSIM values of the reconstructed images from the GoPro test set. These results confirm that at filter size 3×3 , the performance of the transposed convolution variant of both the considered networks is significantly worse than the originally proposed Pixel Shuffle variant, justifying the community’s extensive use of Pixel Shuffle. However, we observe on increasing context by increasing the kernel size to 7×7 that the performance of the transposed convolution variants significantly improves, especially making the networks more stable when facing adversarial attacks. This boost in performance is further accentuated by increasing the kernel size to 11×11 (both with parallel small kernels). These results provide evidence for Hypothesis 1.

Note that the slightly reduced performance on clean images, seen in Table 1, is expected to some degree: here, we only investigate sampling in the decoder, while pixel unshuffle is used in the encoder, potentially causing a mismatch. Further, previous works have shown that there exists a trade-off between adversarial robustness and clean performance [85, 90]. However, we do not observe this trade-off for matching encoder-decoder architectures, *e.g.* in semantic segmentation.

5.2 Semantic Segmentation

As baseline architecture for semantic segmentation, we consider a UNet-like architecture [70] with encoder backbone layers from ConvNeXt [54] (see Appendix B.2 on the choice of encoder). This generic architecture facilitates providing a thorough ablation

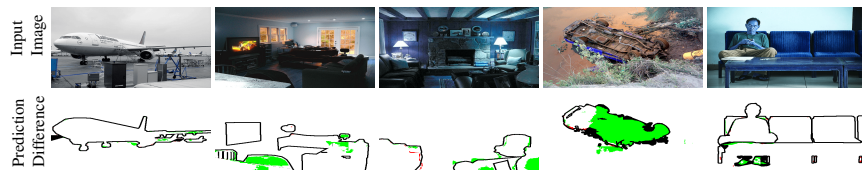


Fig. 6: A comparison of semantic segmentation mask predictions for the shown input images. The row labeled “Prediction Difference” shows the difference in predictions between the baseline model and the model with **Large Context Transposed Convolutions** ($11 \times 11 + 3 \times 3$ kernels). On white pixels, both models agree. Red pixels indicate that the baseline model predicts correctly but our modified model predicts incorrectly. Green pixels indicate that our modified model predicts correctly but the baseline does not. The ground truth segmentation boundaries are drawn in black. Our modification improves the segmentation result along object boundaries, which can be attributed to spectral artifact removal, but also in more extended regions, where the context plays a more crucial role.

Table 2: Semantic Segmentation performance on the PASCAL VOC2012 validation set for UNet with ConvNeXt encoder, and the baseline UNet decoder (see Figure 4) with differently sized **kernels in transposed convolution for feature map upscaling** while keeping rest of the architecture fixed. Additional results are provided in Tab. 7 and Tab. 8 in Appendix B.1.

Transposed Convolution Kernels	Clean Test Accuracy			FGSM attack epsilon						SegPGD ($\epsilon \approx \frac{8}{255}$) attack iterations					
	mIoU	mAcc	allAcc	mIoU	mAcc	allAcc	mIoU	mAcc	allAcc	mIoU	mAcc	allAcc	mIoU	mAcc	allAcc
2×2 (baseline)	78.34	86.89	95.15	53.54	70.96	86.08	47.02	65.41	82.78	23.06	46.51	45.30	5.54	18.79	23.72
LCTC: 7×7 (Ours)	78.92	88.06	95.23	56.02	74.13	86.45	49.24	68.89	82.87	26.53	53.05	61.16	7.17	23.05	27.52
LCTC: 11×11 + 3×3 (Ours)	79.33	87.81	95.41	58.04	74.93	87.80	51.25	69.31	84.64	27.49	53.08	64.13	7.08	23.30	26.82

on all considered blocks in the decoder network. Our experiments are conducted on the PASCAL VOC 2012 dataset [24]. We report the mean Intersection over Union (mIoU) of the predicted and ground truth segmentation mask, the mean accuracy over all pixels (mAcc), and the mean accuracy over all classes (allAcc).

Results on Semantic Segmentation. We first discuss the results for **different upsampling operations**. The remaining architecture is kept identical, with ResNet-style building blocks in the decoder, throughout these experiments. The clean test accuracies are shown in Table 2. We see that as we increase the kernel size of the transposed convolution layers, there is a slight increase across all three evaluation metrics. Moreover, Figure 6 visually demonstrates that, as we increase the size of the kernels in transposed convolution from 2×2 (baseline) to 11×11 , the segmentations of the thin end and protrusions, for example, in the wing of the aircraft sample image are improving. The baseline model with small transposed convolution kernels could not predict these details. As hypothesized in H1, we observe that increasing the context can reduce spectral artifacts caused when representation and images are upsampled using LCTC.

Further, in Table 2, we evaluate the performance of the segmentation models against FGSM [28] and the multi-step attack SegPGD [34] adversarial attacks for the indicated ϵ values. As expected, with the increasing intensity of the attack, the performance of all models drops. Yet, even at high attack intensities, the larger kernels perform better than

Table 3: Adversarially trained models using FGSM ($\epsilon \approx \frac{8}{255}$) from Table 2 tested against SegPGD adversarial attacks ($\epsilon \approx \frac{8}{255}$) on UNet with ConvNeXt encoder and decoder with different sized kernels in the **transposed convolution for upsampling**, while keeping rest of the architecture identical. See Tab. 10 in Appendix B.6 for more evaluations including PGD training.

Transposed Convolution Kernels	Clean Test Data			SegPGD attack iterations					
	mIoU	mAcc	allAcc	3			20		
				mIoU	mAcc	allAcc	mIoU	mAcc	allAcc
2×2	78.57	86.68	95.23	26.59	48.99	67.71	7.6	24.06	31.37
LCTC: 7×7 (Ours)	78.41	86.22	95.20	28.11	53.39	66.30	8.36	28.54	28.13
LCTC: 11×11 + 3×3 (Ours)	79.57	88.1	95.3	30.37	55.54	68.3	9.4	29.79	32.37

Table 4: Empirical evaluations for H2 using a UNet with ConvNeXt encoder. We observe that across different-sized kernels in **transposed convolution**, for a fixed kernel size, increasing the context in the **decoder building blocks** by using larger kernels causes performance deterioration. These observations for image decoding contrast the findings on image encoding by [17, 51, 54].

Transposed Convolution Kernels	Decoder Building Block Style	Test Accuracy			FGSM attack epsilon			SegPGD ($\epsilon \approx \frac{8}{255}$) attack iterations		
		mIoU / mAcc / allAcc	$\frac{1}{255}$ mIoU / mAcc / allAcc	$\frac{8}{255}$ mIoU / mAcc / allAcc	$\frac{1}{255}$ mIoU / mAcc / allAcc	$\frac{8}{255}$ mIoU / mAcc / allAcc	$\frac{8}{255}$ mIoU / mAcc / allAcc	3 mIoU / mAcc / allAcc	20 mIoU / mAcc / allAcc	
2×2	ResNet Style 3×3	78.34 / 86.89 / 95.15			53.54 / 70.96 / 86.08	47.02 / 65.41 / 82.78		23.06 / 46.51 / 60.04	5.54 / 18.79 / 23.72	
	ConvNeXt style 7×7	77.17 / 86.86 / 94.81			49.98 / 72.22 / 83.93	42.04 / 64.86 / 79.08		17.94 / 44.81 / 47.96	3.20 / 14.73 / 9.81	
	ConvNeXt style 11×11 + 3×3	77.17 / 86.86 / 94.81			47.34 / 67.72 / 83.34	37.91 / 57.79 / 78.21		13.97 / 35.82 / 45.68	2.21 / 10.75 / 5.29	
LCTC: 7×7 (Ours)	ResNet Style 3×3	78.92 / 88.06 / 95.23			56.02 / 74.13 / 86.45	49.24 / 68.89 / 82.87		26.53 / 53.05 / 61.16	7.17 / 23.05 / 27.52	
	ConvNeXt style 7×7	77.57 / 87.04 / 94.92			52.93 / 72.18 / 85.51	44.89 / 63.71 / 80.74		17.64 / 43.32 / 47.80	1.86 / 7.18 / 3.55	
	ConvNeXt style 11×11 + 3×3	77.99 / 87.86 / 94.96			51.61 / 73.01 / 84.85	43.93 / 66.22 / 80.73		17.07 / 42.30 / 48.78	1.80 / 7.11 / 3.04	
LCTC: 11×11 + 3×3 (Ours)	ResNet Style 3×3	79.33 / 87.81 / 95.41			58.04 / 74.93 / 87.80	51.25 / 69.31 / 84.64		27.49 / 53.08 / 64.13	7.08 / 23.30 / 26.82	
	ConvNeXt style 7×7	78.32 / 86.98 / 95.09			53.31 / 72.45 / 86.16	44.89 / 65.18 / 82.03		16.14 / 40.65 / 50.39	1.93 / 9.35 / 3.90	
	ConvNeXt style 11×11 + 3×3	77.42 / 86.24 / 94.94			54.48 / 72.53 / 86.25	46.67 / 66.59 / 82.29		18.76 / 44.60 / 51.49	2.31 / 8.70 / 3.50	

the small ones, and we see a trend of improvement in performance as we increase the kernel size, providing more evidence for Hypothesis 1.

Ablation Study. In the following, we first consider the effects of additional adversarial training, then ablate on the impact of other decoder building blocks and the filter size. Variations of the model encoder are ablated in the Appendix B.2, the impact of using small parallel kernels in addition to large kernels is ablated and discussed in Appendix B.3, and competing upsampling techniques are ablated in Appendix B.5.

Adversarial Training. In Table 3, we report results for FGSM adversarially trained models under SegPGD attack, with attacks as in Table 2. While the overall performance under attack is improved as expected, the trend of LCTC providing better results persists. More results for FGSM attack and SegPGD attacks with different numbers of iterations are given in Tab. 7 and Tab. 8 in the Appendix. In Table 15, we additionally evaluate image restoration models under adversarial training.

Change in the decoder backbone architecture. While all previous experiments focused on the upsampling using transposed convolutions in the decoder, we now evaluate the influence of the convolutional kernel size within the decoder which does not upsample (see Section 4). For these experiments, we use a UNet-like architecture with a ConvNeXt backbone in the encoder and the PASCAL VOC 2012 dataset.

In Table 4 we observe, for a **fixed transposed convolution** kernel size, as **we increase the size of the convolution kernel in the decoder building blocks**, the performance of the model *decreases*. This phenomenon extends to the performance of the architectures

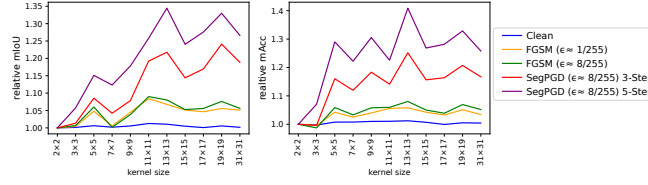


Fig. 7: Performance comparison on PASCAL VOC2012 using UNet with ConvNeXt encoder for different LCTC sizes from 2×2 (small) to 31×31 (LCTC) kernels. All, besides the baseline with 2×2 and 3×3 , have a parallel 3×3 kernel, as shown in Figure 4 (bottom left). For the decoder building block backbone, a ResNet Style 3×3 style is used. See Tab. 9 for the values.

Table 5: Comparison of performances of different upsampling methods in the UNet-like architecture. All architectures use the baseline (ConvNeXt) encoder and 3×3 convolution kernels in the decoder block. Please refer to Table 13 in Appendix B.5 for more evaluations and discussion, including those with ConvNeXt style $7 \times 7 + 3 \times 3$ Convolution kernels in the decoder blocks.

Upsampling Method	Test Accuracy			FGSM attack epsilon						SegPGD ($\epsilon \approx \frac{8}{255}$) attack iterations					
	mIoU	mAcc	allAcc	$\frac{1}{255}$ mIoU	$\frac{1}{255}$ mAcc	$\frac{1}{255}$ allAcc	$\frac{8}{255}$ mIoU	$\frac{8}{255}$ mAcc	$\frac{8}{255}$ allAcc	5 mIoU	5 mAcc	5 allAcc	20 mIoU	20 mAcc	20 allAcc
Pixel Shuffle	78.54	87.32	95.18	53.82	71.58	85.88	46.67	65.03	81.71	15.06	38.85	41.71	6.69	23.43	24.05
Nearest Neighbour Interpolation	78.40	88.16	95.09	52.68	73.51	84.55	46.08	67.96	80.22	15.34	44.53	36.21	7.65	27.89	20.48
Transposed Convolution 2×2	78.45	86.66	95.20	53.76	70.62	86.32	47.33	64.58	83.16	14.43	35.50	45.30	5.54	18.79	23.72
LCTC: $11 \times 11 + 3 \times 3$ (Ours)	79.33	87.81	95.41	58.04	74.93	87.80	51.25	69.31	84.64	18.15	43.51	49.36	7.08	23.30	26.82

under adversarial attacks, showing that a mere increase in parameters in the model decoder does not have a positive effect on model performance or on its stability. This proves the validity of hypothesis H2. An explanation for this phenomenon could be that we only need to increase context during the **actual upsampling step**, increasing context in the consequent **decoder building blocks** has a negligible effect on the quality of representations learned. However, the increase in the number of parameters makes the architecture more susceptible to adversarial attacks.

Ablation on filter size saturation. After proving H1 one could argue that networks will consistently improve with increased **kernel size for Large Context Transposed Convolutions**. Hence, we test larger kernel sizes of 15×15 , 17×17 , 19×19 and 31×31 kernels. Yet, as seen in Figure 7, the effect of the kernel size appears to saturate: the performance after 13×13 and the performance of 31×31 kernels is not better than for 11×11 kernels. Yet, they are significantly better than the baseline’s performance.

Ablation on different Upsampling Methods. Following, we compare different upsampling techniques thus justifying our advocacy for using LCTC instead of other upsampling techniques like interpolation and pixel shuffle in the real world.

We report the comparison in Table 5 and observe that both Pixel shuffle and Nearest Neighbor interpolation perform better than the usually used Transposed Convolution with a 2×2 kernel size. However, as we increase the kernel size for Transposed Convolution to 11×11 with a 3×3 small kernel in parallel, we observe that LCTC is strictly outperforming Pixel Shuffle, on both clean unperturbed images and under adversarial attacks, across all metrics used. Large Context Transposed Convolutions are either outperforming or performing at par with Nearest Neighbor interpolation. Thus we demonstrate the superior clean and adversarial performance of Large Context Transposed Convolutions operation over other commonly used techniques.

5.3 Disparity Estimation

To show that the observations extend from image restorations and segmentation to other tasks, we conduct additional experiments for disparity estimation. We consider the STTR-light [50] architecture, built from STTR, which is a recent state-of-the-art vision-transformer based model for disparity estimation and occlusion detection. To implement the proposed modification, we alter the kernel sizes in the transposed convolution layers used for pixel-wise upsampling in the “feature extractor” module of the architecture from 3×3 kernels to larger kernels. We conduct evaluations on FlyingThings3D [60] and keep all other details as implemented in [50].

In Table 6, we report the improvements in performance due to our architecture modification of increasing the size of the **transposed convolution kernels used for upsampling**, from the 3×3 in the baseline model to 7×7 (LCTC). Similar to previous applications, the increased kernel sizes with parallel 3×3 kernels further facilitate

Table 6: Comparison of performance of STTR-light architecture with different sized kernels in **transposed convolution for upsampling** the feature maps in the feature extractor (lower is better). The entire set of results is provided as Tab. 17 in Appendix D.1.

Transposed Convolution Kernels	Test Accuracy		3-Step PGD attack	
	epe↓	3px error↓	epe↓	3px error↓
STTR-light [50] reported	0.5	1.54	-	-
3×3 [50] reproduced	0.4927	1.54	4.05	18.5
LCTC: $7\times 7 + 3\times 3$ (Ours)	0.4788	1.50	4.02	18.3

to stabilize the model when attacked, as evaluated here for 3 attack iterations using PGD with $\epsilon \approx \frac{8}{255}$ on the disparity loss. Indicating that larger kernels in the transposed convolutions can better decode learned representations from the encoder regardless of the specific downstream task. We provide visual results in Appendix D.1.

6 Conclusion

We provide conclusive reasoning and empirical evidence for our hypotheses on the importance of context during upsampling. While increasing the **size of convolutions during upsampling (LCTC)** increases prediction stability, increasing the size of those **convolution layers without upsampling** does not benefit the network. This indicates that observations made for increased context during encoding do not translate to decoding. Further, we show that our simple LCTC can be directly incorporated into recent models, yielding better stability even in ViT-based architectures like Restormer, NAFNet, and STTR-light as well as in classical CNNs. Our observations are consistent across several architectures and downstream tasks.

Limitations. Current metrics for measuring performance do not completely account for spectral artifacts. Spectral artifacts begin affecting these metrics only when they become pronounced such as under adversarial attacks, and here LCTC consistently performs better across tasks and architectures. Ideally, we would want infinitely large kernels, however, with increasing kernel size and task complexity, training extremely large kernels can be challenging. Thus, in this work, while having ablated over kernels as large as 31×31 , we propose using kernels only as large as 7×7 to 11×11 for good practical trade-offs. Further improvements might be possible when jointly optimizing the encoder *and* decoder. Moreover, there might exist other factors that contribute to the introduction and existence of spectral artifacts such as spatial bias.

Acknowledgements.

Margret Keuper acknowledges funding by the DFG Research Unit 5336 - Learning to Sense. The OMNI cluster of the University of Siegen was used for some of the initial computations. Additionally, Shashank Agnihotri would like to thank Dr. Bin Zhao for his help in translating [75].

References

1. Agnihotri, S., Gandikota, K.V., Grabinski, J., Chandramouli, P., Keuper, M.: On the unreasonable vulnerability of transformers for image restoration-and an easy fix. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 3707–3717 (2023)
2. Agnihotri, S., Grabinski, J., Keuper, J., Keuper, M.: Beware of aliases–signal preservation is crucial for robust image restoration. arXiv preprint arXiv:2406.07435 (2024)
3. Agnihotri, S., Jung, S., Keuper, M.: CosPGD: A unified white-box adversarial attack for pixel-wise prediction tasks. In: International Conference on Machine Learning (2024)
4. Aitken, A., Ledig, C., Theis, L., Caballero, J., Wang, Z., Shi, W.: Checkerboard artifact free sub-pixel convolution: A note on sub-pixel convolution, resize convolution and convolution resize. arXiv preprint arXiv:1707.02937 (2017)
5. Ba, J.L., Kiros, J.R., Hinton, G.E.: Layer normalization. arXiv preprint arXiv:1607.06450 (2016)
6. Badki, A., Troccoli, A., Kim, K., Kautz, J., Sen, P., Gallo, O.: Bi3D: Stereo depth estimation via binary classifications. In: The IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (2020)
7. Badrinarayanan, V., Kendall, A., Cipolla, R.: Segnet: A deep convolutional encoder-decoder architecture for image segmentation. *IEEE transactions on pattern analysis and machine intelligence* **39**(12), 2481–2495 (2017)
8. Baranchuk, D., Rubachev, I., Voynov, A., Khrulkov, V., Babenko, A.: Label-efficient semantic segmentation with diffusion models (2021)
9. Brigham, E.O., Morrow, R.: The fast fourier transform. *IEEE spectrum* **4**(12), 63–70 (1967)
10. Brunton, S.L., Kutz, J.N.: Data-driven science and engineering: Machine learning, dynamical systems, and control. Cambridge University Press (2022)
11. Carion, N., Massa, F., Synnaeve, G., Usunier, N., Kirillov, A., Zagoruyko, S.: End-to-end object detection with transformers. In: European conference on computer vision. pp. 213–229. Springer (2020)
12. Carlini, N., Wagner, D.: Towards evaluating the robustness of neural networks. In: 2017 IEEE symposium on security and privacy (SP). pp. 39–57. IEEE (2017)
13. Caron, M., Misra, I., Mairal, J., Goyal, P., Bojanowski, P., Joulin, A.: Unsupervised learning of visual features by contrasting cluster assignments. *Advances in neural information processing systems* **33**, 9912–9924 (2020)
14. Chandrasegaran, K., Tran, N.T., Cheung, N.M.: A closer look at fourier spectrum discrepancies for cnn-generated images detection. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR). pp. 7200–7209 (June 2021)
15. Chang, J.R., Chen, Y.S.: Pyramid stereo matching network. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 5410–5418 (2018)
16. Chen, L., Chu, X., Zhang, X., Sun, J.: Simple baselines for image restoration. In: European Conference on Computer Vision. pp. 17–33. Springer (2022)
17. Ding, X., Zhang, X., Han, J., Ding, G.: Scaling up your kernels to 31x31: Revisiting large kernel design in cnns. In: Proceedings of the IEEE/CVF conference on computer vision and pattern recognition. pp. 11963–11975 (2022)

18. Dosovitskiy, A., Brox, T.: Generating images with perceptual similarity metrics based on deep networks. *Advances in neural information processing systems* **29** (2016)
19. Dosovitskiy, A., Brox, T.: Inverting visual representations with convolutional networks. In: *Proceedings of the IEEE conference on computer vision and pattern recognition*. pp. 4829–4837 (2016)
20. Dosovitskiy, A., Fischer, P., Ilg, E., Hausser, P., Hazirbas, C., Golkov, V., Van Der Smagt, P., Cremers, D., Brox, T.: FlowNet: Learning optical flow with convolutional networks. In: *Proceedings of the IEEE international conference on computer vision*. pp. 2758–2766 (2015)
21. Dosovitskiy, A., Springenberg, J.T., Tatarchenko, M., Brox, T.: Learning to generate chairs, tables and cars with convolutional networks. *IEEE Transactions on Pattern Analysis and Machine Intelligence* **39**(4), 692–705 (2017)
22. Dumoulin, V., Visin, F.: A guide to convolution arithmetic for deep learning. *arXiv preprint arXiv:1603.07285* (2016)
23. Durall, R., Keuper, M., Keuper, J.: Watch your up-convolution: Cnn based generative deep neural networks are failing to reproduce spectral distributions. In: *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. pp. 7890–7899 (2020)
24. Everingham, M., Van Gool, L., Williams, C.K.I., Winn, J., Zisserman, A.: The PASCAL Visual Object Classes Challenge 2012 (VOC2012) Results. <http://www.pascal-network.org/challenges/VOC/voc2012/workshop/index.html> (2012)
25. Forsyth, D.A., Ponce, J.: *Computer vision: a modern approach*. Prentice Hall professional technical reference (2002)
26. Gal, R., Hochberg, D.C., Bermanno, A., Cohen-Or, D.: SWAGAN: A style-based wavelet-driven generative model. *ACM Transactions on Graphics (TOG)* **40**(4), 1–11 (2021)
27. Goodfellow, I., Pouget-Abadie, J., Mirza, M., Xu, B., Warde-Farley, D., Ozair, S., Courville, A., Bengio, Y.: Generative adversarial networks. *Communications of the ACM* **63**(11), 139–144 (2020)
28. Goodfellow, I.J., Shlens, J., Szegedy, C.: Explaining and harnessing adversarial examples. *arXiv preprint arXiv:1412.6572* (2014)
29. Grabinski, J., Jung, S., Keuper, J., Keuper, M.: Frequencylowcut pooling-plugin and play against catastrophic overfitting. In: *European Conference on Computer Vision*. pp. 36–57. Springer (2022)
30. Grabinski, J., Keuper, J., Keuper, M.: Aliasing and adversarial robust generalization of cnns. *Machine Learning* pp. 1–27 (2022)
31. Grabinski, J., Keuper, J., Keuper, M.: Aliasing coincides with cnns vulnerability towards adversarial attacks. In: *The AAAI-22 Workshop on Adversarial Machine Learning and Beyond*. pp. 1–5 (2022)
32. Grabinski, J., Keuper, J., Keuper, M.: Fix your downsampling asap! be natively more robust via aliasing and spectral artifact free pooling (2023)
33. Grabinski, J., Keuper, J., Keuper, M.: As large as it gets – studying infinitely large convolutions via neural implicit frequency filters. *Transactions on Machine Learning Research* (2024), <https://openreview.net/forum?id=xRyLYRcHWj>, featured Certification
34. Gu, J., Zhao, H., Tresp, V., Torr, P.H.: Segpgd: An effective and efficient adversarial attack for evaluating and boosting segmentation robustness. In: *European Conference on Computer Vision*. pp. 308–325. Springer (2022)
35. Guo, M.H., Lu, C.Z., Hou, Q., Liu, Z., Cheng, M.M., Hu, S.M.: Segnext: Rethinking convolutional attention design for semantic segmentation. *Advances in Neural Information Processing Systems* **35**, 1140–1156 (2022)
36. Hariharan, B., Arbeláez, P., Bourdev, L., Maji, S., Malik, J.: Semantic contours from inverse detectors. In: *2011 international conference on computer vision*. pp. 991–998. IEEE (2011)

37. He, K., Zhang, X., Ren, S., Sun, J.: Deep residual learning for image recognition. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 770–778 (2016)
38. He, Y., Yu, N., Keuper, M., Fritz, M.: Beyond the spectrum: Detecting deepfakes via re-synthesis. In: Zhou, Z.H. (ed.) Proceedings of the Thirtieth International Joint Conference on Artificial Intelligence, IJCAI-21. pp. 2534–2541. International Joint Conferences on Artificial Intelligence Organization (8 2021), main Track
39. Hendrycks, D., Dietterich, T.: Benchmarking neural network robustness to common corruptions and perturbations. Proceedings of the International Conference on Learning Representations (2019)
40. Hendrycks, D., Gimpel, K.: Gaussian error linear units (gelus). arXiv preprint arXiv:1606.08415 (2016)
41. Hoffmann, J., Agnihotri, S., Saikia, T., Brox, T.: Towards improving robustness of compressed cnns. In: ICML Workshop on Uncertainty and Robustness in Deep Learning (UDL) (2021)
42. Hossain, M.T., Teng, S.W., Lu, G., Rahman, M.A., Sohel, F.: Anti-aliasing deep image classifiers using novel depth adaptive blurring and activation function. *Neurocomputing* **536**, 164–174 (2023)
43. Ilg, E., Mayer, N., Saikia, T., Keuper, M., Dosovitskiy, A., Brox, T.: FlowNet 2.0: Evolution of optical flow estimation with deep networks. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 2462–2470 (2017)
44. Jung, S., Keuper, M.: Spectral distribution aware image generation. In: Proceedings of the AAAI conference on artificial intelligence. vol. 35, pp. 1734–1742 (2021)
45. Jung, S., Lukasik, J., Keuper, M.: Neural architecture design and robustness: A dataset. In: Eleventh International Conference on Learning Representations. OpenReview.net (2023)
46. Karras, T., Aittala, M., Laine, S., Härkönen, E., Hellsten, J., Lehtinen, J., Aila, T.: Alias-free generative adversarial networks. *Advances in Neural Information Processing Systems* **34**, 852–863 (2021)
47. Khayatkhoei, M., Elgammal, A.: Spatial frequency bias in convolutional generative adversarial networks. In: Proceedings of the AAAI Conference on Artificial Intelligence. vol. 36, pp. 7152–7159 (2022)
48. Krizhevsky, A., Sutskever, I., Hinton, G.E.: Imagenet classification with deep convolutional neural networks. *Advances in neural information processing systems* **25** (2012)
49. Kurakin, A., Goodfellow, I.J., Bengio, S.: Adversarial machine learning at scale. In: International Conference on Learning Representations (2017), <https://openreview.net/forum?id=BJm4T4Kgxx>
50. Li, Z., Liu, X., Drenkow, N., Ding, A., Creighton, F.X., Taylor, R.H., Unberath, M.: Revisiting stereo depth estimation from a sequence-to-sequence perspective with transformers. In: Proceedings of the IEEE/CVF international conference on computer vision. pp. 6197–6206 (2021)
51. Liu, S., Chen, T., Chen, X., Chen, X., Xiao, Q., Wu, B., Kärkkäinen, T., Pechenizkiy, M., Mocanu, D., Wang, Z.: More convnets in the 2020s: Scaling up kernels beyond 51x51 using sparsity. arXiv preprint arXiv:2207.03620 (2022)
52. Liu, S., Deng, W.: Very deep convolutional neural network based image classification using small training sample size. In: 2015 3rd IAPR Asian conference on pattern recognition (ACPR). pp. 730–734. IEEE (2015)
53. Liu, Z., Lin, Y., Cao, Y., Hu, H., Wei, Y., Zhang, Z., Lin, S., Guo, B.: Swin transformer: Hierarchical vision transformer using shifted windows. In: Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV) (2021)

54. Liu, Z., Mao, H., Wu, C.Y., Feichtenhofer, C., Darrell, T., Xie, S.: A convnet for the 2020s. In: Proceedings of the IEEE/CVF conference on computer vision and pattern recognition. pp. 11976–11986 (2022)
55. Long, J., Shelhamer, E., Darrell, T.: Fully convolutional networks for semantic segmentation. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 3431–3440 (2015)
56. Loshchilov, I., Hutter, F.: SGDR: Stochastic gradient descent with warm restarts. In: International Conference on Learning Representations (2017), <https://openreview.net/forum?id=Skq89Scxx>
57. Loshchilov, I., Hutter, F.: Decoupled weight decay regularization. In: International Conference on Learning Representations (2019), <https://openreview.net/forum?id=Bkg6RiCqY7>
58. Maiya, S.R., Ehrlich, M., Agarwal, V., Lim, S.N., Goldstein, T., Shrivastava, A.: A frequency perspective of adversarial robustness (2021)
59. Mathew, A., Patra, A., Mathew, J.: Monocular depth estimators: Vulnerabilities and attacks. ArXiv **abs/2005.14302** (2020)
60. Mayer, N., Ilg, E., Hausser, P., Fischer, P., Cremers, D., Dosovitskiy, A., Brox, T.: A large dataset to train convolutional networks for disparity, optical flow, and scene flow estimation. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 4040–4048 (2016)
61. Moosavi-Dezfooli, S.M., Fawzi, A., Frossard, P.: Deepfool: a simple and accurate method to fool deep neural networks. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 2574–2582 (2016)
62. Mosleh, A., Langlois, J.M.P., Green, P.: Image deconvolution ringing artifact detection and removal via psf frequency analysis. In: Fleet, D., Pajdla, T., Schiele, B., Tuytelaars, T. (eds.) Computer Vision – ECCV 2014. pp. 247–262. Springer International Publishing, Cham (2014)
63. Nah, S., Kim, T.H., Lee, K.M.: Deep multi-scale convolutional neural network for dynamic scene deblurring. In: CVPR (July 2017)
64. Noh, H., Hong, S., Han, B.: Learning deconvolution network for semantic segmentation. In: Proceedings of the IEEE international conference on computer vision. pp. 1520–1528 (2015)
65. Odena, A., Dumoulin, V., Olah, C.: Deconvolution and checkerboard artifacts. Distill **1**(10), e3 (2016)
66. Peng, C., Zhang, X., Yu, G., Luo, G., Sun, J.: Large kernel matters – improve semantic segmentation by global convolutional network. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (July 2017)
67. Pervin, M., Tao, L., Huq, A., He, Z., Huo, L., et al.: Adversarial attack driven data augmentation for accurate and robust medical image segmentation. arXiv preprint arXiv:2105.12106 (2021)
68. Radford, A., Metz, L., Chintala, S.: Unsupervised representation learning with deep convolutional generative adversarial networks. arXiv preprint arXiv:1511.06434 (2015)
69. Radosavovic, I., Kosaraju, R.P., Girshick, R., He, K., Dollar, P.: Designing network design spaces. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) (June 2020)
70. Ronneberger, O., Fischer, P., Brox, T.: U-net: Convolutional networks for biomedical image segmentation. In: Medical Image Computing and Computer-Assisted Intervention–MICCAI 2015: 18th International Conference, Munich, Germany, October 5–9, 2015, Proceedings, Part III 18. pp. 234–241. Springer (2015)
71. Scheurer, E., Schmalz, J., Lis, A., Bruhn, A.: Detection defenses: An empty promise against adversarial patch attacks on optical flow. arXiv preprint arXiv:2310.17403 (2023)

72. Schmalfluss, J., Mehl, L., Bruhn, A.: Attacking motion estimation with adversarial snow. arXiv preprint arXiv:2210.11242 (2022)
73. Schmalfluss, J., Mehl, L., Bruhn, A.: Distracting downpour: Adversarial weather attacks for motion estimation (2023)
74. Schmalfluss, J., Scholze, P., Bruhn, A.: A perturbation-constrained adversarial attack for evaluating the robustness of optical flow (2022)
75. segcv: segcv/pspnet. <https://github.com/segcv/PSPNet/blob/master/Train.md> (2021)
76. Shannon, C.E.: Communication in the presence of noise. *Proceedings of the IRE* **37**(1), 10–21 (1949)
77. Shi, W., Caballero, J., Huszár, F., Totz, J., Aitken, A.P., Bishop, R., Rueckert, D., Wang, Z.: Real-time single image and video super-resolution using an efficient sub-pixel convolutional neural network (2016)
78. Si-Yao, L., Ren, D., Yin, Q.: Understanding kernel size in blind deconvolution. In: 2019 IEEE Winter Conference on Applications of Computer Vision (WACV). pp. 2068–2076. IEEE (2019)
79. Smith III, J.O.: Physical audio signal processing: For virtual musical instruments and audio effects. (No Title) (2010)
80. Sommerhoff, H., Agnihotri, S., Saleh, M., Moeller, M., Keuper, M., Kolb, A.: Differentiable sensor layouts for end-to-end learning of task-specific camera parameters. arXiv preprint arXiv:2304.14736 (2023)
81. Szegedy, C., Zaremba, W., Sutskever, I., Bruna, J., Erhan, D., Goodfellow, I., Fergus, R.: Intriguing properties of neural networks. arXiv preprint arXiv:1312.6199 (2013)
82. Tan, M., Le, Q.: Efficientnetv2: Smaller models and faster training. In: International conference on machine learning. pp. 10096–10106. PMLR (2021)
83. Teed, Z., Deng, J.: Raft: Recurrent all-pairs field transforms for optical flow. In: Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part II 16. pp. 402–419. Springer (2020)
84. Touvron, H., Cord, M., Douze, M., Massa, F., Sablayrolles, A., Jégou, H.: Training data-efficient image transformers & distillation through attention. In: International conference on machine learning. pp. 10347–10357. PMLR (2021)
85. Tsipras, D., Santurkar, S., Engstrom, L., Turner, A., Madry, A.: Robustness may be at odds with accuracy. In: International Conference on Learning Representations (2019), <https://openreview.net/forum?id=SyxAb30cY7>
86. Wang, Z., Bovik, A.C., Sheikh, H.R., Simoncelli, E.P.: Image quality assessment: from error visibility to structural similarity. *IEEE transactions on image processing* **13**(4), 600–612 (2004)
87. Xu, L., Ren, J.S., Liu, C., Jia, J.: Deep convolutional neural network for image deconvolution. *Advances in neural information processing systems* **27** (2014)
88. Yamanaka, K., Matsumoto, R., Takahashi, K., Fujii, T.: Adversarial patch attacks on monocular depth estimation networks. *IEEE Access* **8**, 179094–179104 (2020)
89. Zamir, S.W., Arora, A., Khan, S., Hayat, M., Khan, F.S., Yang, M.H.: Restormer: Efficient transformer for high-resolution image restoration. In: Proceedings of the IEEE/CVF conference on computer vision and pattern recognition. pp. 5728–5739 (2022)
90. Zhang, H., Yu, Y., Jiao, J., Xing, E., El Ghaoui, L., Jordan, M.: Theoretically principled trade-off between robustness and accuracy. In: ICML (2019)
91. Zhang, R.: Making convolutional networks shift-invariant again. In: ICML (2019)
92. Zhao, H.: semseg. <https://github.com/hszhao/semseg> (2019)
93. Zhao, H., Shi, J., Qi, X., Wang, X., Jia, J.: Pyramid scene parsing network. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 2881–2890 (2017)

94. Zhao, H., Zhang, Y., Liu, S., Shi, J., Loy, C.C., Lin, D., Jia, J.: PSANet: Point-wise spatial attention network for scene parsing. In: ECCV (2018)
95. Zou, X., Xiao, F., Yu, Z., Lee, Y.J.: Delving deeper into anti-aliasing in convnets. In: BMVC (2020)

Improving Feature Stability during Upsampling – Spectral Artifacts and the Importance of Spatial Context

Supplementary Material

In the following, we present results and figures to support our statements in the main paper and provide additional information. The following has been covered in the appendix:

- Appendix A: **Detailed experimental setup** for all downstream tasks.
 - Appendix A.1: Image Restoration experimental setup
 - Appendix A.2: Semantic Segmentation experimental setup
 - Appendix A.3: Disparity Estimation experimental setup
 - Appendix A.4: Detailed setup of Adversarial attacks for all downstream tasks.
 - Appendix A.5: Detailed setup of adversarial training for semantic segmentation and image restoration.

- Appendix B: **Semantic Segmentation:** Additional Experiments and Ablations. In detail:
 - Appendix B.1: Detailed results from Sec. 5.2 and Sec. 5.2.
 - Appendix B.1: Discussion on saturation of kernel size for upsampling.
 - Appendix B.2: An ablation on the impact of the capacity of the encoder block for standard options such as ResNet or ConvNeXt blocks.
 - Appendix B.3: Ablation about including or excluding a small parallel kernel during upsampling using transposed convolution.
 - Appendix B.4: Short study on drawbacks of using interpolation for pixel-wise upsampling.
 - Appendix B.5: A comparison to different kinds of upsampling Operations on Segmentation Models.
 - Appendix B.6: A comparison of the performance of different sized kernels in the transposed convolution operations of UNet-like models adversarially trained using FGSM attack and 3-step PGD attack on 50% of the mini-batches during training.

- Appendix C: **Image Restoration** : Additional Results:
 - Appendix C.1: Here we report the number of parameters and latency study of LCTC.
 - Appendix C.2: Adversarial training evaluation for Restormer and NAFNet for Image deblurring task.
 - Appendix C.3: Qualitative results for image reconstruction models using Restormer and NAFNet and evaluated on clean data, PDG and CosPGD attack with varying numbers of attack iterations.
 - Appendix C.4: **Visualizing Kernel Weights**: Here we visualize kernel weights from a random channel for models from Figure 5 to show the how different kernels handle uneven contributions of pixels that leads to spectral artifacts.
 - Appendix C.5: Out-Of-Distribution and Real World Generalization.
- Appendix D: **Disparity Estimation** : We provide additional results for Section 5.3: including performance against adversarial attacks.
 - Appendix D.1 Additional discussion on the results and importance of a parallel 3×3 kernel with large kernels for transposed convolution operation.
- Appendix E: **Nomenclature- What are “Large Context Transposed Convolutions?”**: We discuss the nomenclature used in this work and describe what comprises a LCTC.
- Appendix F: **Additional visualizations of Upsampling Artifacts and their Frequency Spectra**: Here we extend Figure 1 with more examples showing failure of upsampling operations used in prior work and superiority of LCTC both in the spatial and frequency domain.
- Appendix G: **Limitations**: Here we discuss the limitations of our work in detail.

A Experimental Setup

All the experiments were done using NVIDIA V100 16GB GPUs or NVIDIA Tesla A100 40GB GPUs. For image restoration, models were trained on 1 NVIDIA Tesla A100 40GB GPU. For the semantic segmentation downstream task, UNet [70] was trained using 1 GPU. For the disparity estimation task, STTR-light [50] was trained using 4 NVIDIA V100 GPUs in parallel.

A.1 Image Restoration

Architectures. We consider the recently proposed state-of-the-art transformer-based Image Restoration architectures Restormer [89] and NAFNet [16]. Both architectures as proposed use Pixel Shuffle [77] to upsample feature maps. We use these as our baseline models. We replace this pixel shuffle operation with a transposed convolution operation.

Dataset. For the Image Restoration task, we focus on Image Deblurring. For this, we use the GoPro image deblurring dataset [63]. This dataset consists of 3214 real-world images with realistic blur and their corresponding ground truth (deblurred images) captured using a high-speed camera. The dataset is split into 2103 training images and 1111 test images.

Training Regime. For Restormer we follow the same training regime of progressive training as that used by [89]. Similarly, for NAFNet we use the same training regime as that used by [16].

Evaluation Metrics. Following common practice [1, 16, 89], We report the PSNR and SSIM scores of the reconstructed images w.r.t. to the ground truth images, averaged over all images. PSNR stands for Peak Signal-to-Noise ratio, a higher PSNR indicates a better quality image or an image closer to the image to which it is being compared. SSIM stands for Structural similarity [86]. A higher SSIM score corresponds to better higher similarity between the reconstruction and the ground-truth image.

A.2 Semantic Segmentation

Here we describe the experimental setup for the segmentation task, the architectures considered, the dataset considered and the training regime.

Architectures. We considered UNet [70] with encoder layers from ConvNeXt [54]. For the decoder, the baseline comparison is done with 2×2 kernels in the transposed convolution layers and the commonly used ResNet [37] BasicBlock style layers for the convolution layers in the decoder building blocks. In our experiments, we used larger sized kernels, e.g. 7×7 and 11×11 in the transposed convolution while keeping the rest of the architecture, including the convolution blocks in the decoder identical to Sec. 5.2. When using kernels larger than 7×7 for transposed convolution we follow the work of [17, 51] and additionally include a parallel 3×3 kernel to keep the local context. Usage of this parallel kernel is denoted by “+ 3×3 ” Further, we analyze the behavior of a different block of convolution layers in the decoder, as explained in Sec. 4 and replace the ResNet-style layers with ConvNeXt-style layers in Sec. 5.2.

Dataset. We considered the PASCAL VOC 2012 dataset [24] for the semantic segmentation task. We follow the implementation of [92–94] and augment the training examples with semantic contours from [36] as instructed by [75].

Training Regime. We follow a similar training regime as [92, 93], and train for 50 epochs, with an AdamW optimizer [57] and the learning rate was scheduled using Cosine-Annealing [56]. In the implementation of [93], the authors slide over the images using a window of size 473×473 , however for computation reasons and for symmetry we use a window of size 256×256 . We use a starting learning rate of 10^{-4} and a weight decay of 5×10^{-2} .

Evaluation Metrics. We report the mean Intersection over Union (mIoU) of the predicted and the ground truth segmentation mask, the mean accuracy over all pixels (mAcc) and the mean accuracy over all classes (allAcc).

A.3 Disparity Estimation

Following, we describe the experimental setup for disparity estimation and occlusion detection tasks.

Architectures. We consider the STTR-light [50] architecture for our work. To analyze the influence of implementing larger kernels in transposed convolution as described in Section 4 we alter the kernel sizes in the transposed convolution layers used for pixel-wise upsampling in the “feature extractor” module of the architecture. We consider the STTR-light architecture as proposed by [50] with 3×3 kernels in the transposed convolution layers as our baseline.

Dataset. Similar to [50] we train and test our models on FlyingThings3D dataset [60].

Training Regime. We follow the training regime as implemented in [50].

Evaluation Metrics. We report the end-point-error (epe) and the 3-pixel error (3px) for the disparity estimation w.r.t. the ground truth.

A.4 Adversarial Attacks

We consider the commonly used [34,59,67,88] FGSM attack [28] and a new segmentation-specific SegPGD attack [34] for testing the robustness of the models against adversarial attacks. For the **semantic segmentation downstream task**, each crop of the input was perturbed with FGSM and SegPGD, while for the disparity estimation downstream task, each of the left and right inputs were perturbed using FGSM.

For FGSM, we test our model against epsilons $\epsilon \in \{\frac{1}{255}, \frac{8}{255}\}$. Where, we follow common practice and use $\frac{1}{255} \approx 0.004$ and $\frac{8}{255} \approx 0.03$.

For SegPGD we follow the testing parameters as originally proposed in [34], with $\epsilon \approx \frac{8}{255}$, $\alpha=0.01$ and number of iterations $\in \{3, 5, 10, 20, 40, 100\}$. We use the same scheduling for loss balancing term λ as suggested by the authors. We use SegPGD for the semantic segmentation task as it is a stronger attack specifically designed for segmentation. Thus providing more accurate insights into the models’ performance and giving a better evaluation of the architectural design choices made.

For the **Image Restoration task**, we follow the evaluation method of [1], and evaluate against CosPGD [3] and PGD [49] adversarial attacks. For both attacks, we use $\epsilon \approx \frac{8}{255}$, $\alpha=0.01$ and test for number of attack iterations $\in \{5, 10, 20\}$.

For the **Depth Estimation task**, we use the PGD attack with $\epsilon \approx \frac{8}{255}$, $\alpha=0.01$ and test for number of attack iterations $\in \{5, 10, 20\}$.

A.5 Adversarial Training

Following, we describe the adversarial training setup employed in this work for adversarially training models for semantic segmentation and image restoration.

Semantic Segmentation. We follow the commonly used [34] procedure and split the batch into two 50%-50% mini-batches. One mini-batch is used to generate adversarial examples using FGSM attack with $\epsilon \approx \frac{8}{255}$ and PGD attack with 3 attack iterations and with $\epsilon \approx \frac{8}{255}$ and $\alpha=0.01$ during training.

Image Restoration. We follow the training procedure used by [1]. We split each training batch into two equal 50%-50% mini-batches. We use one of the mini-batches to generate adversarial samples using FGSM attack with $\epsilon \approx \frac{8}{255}$.

A.6 Frequency spectrum analysis

To analyze the images in the frequency domain, we use the Fast Fourier Transform [9] (FFT) $X_c = FFT(x_c)$ for all channels c of feature maps x and aggregate a 2D representation over frequencies w . We compute the mean over C channels of the FFT of the difference between the prediction and the ground truth.

$$\text{2D Frequency Spectra} = \frac{1}{C} \sum_{c \in C} FFT(x_c^{pred} - x_c^{gt}) \quad (4)$$

Here, x^{pred} are the predictions from the model, x^{gt} is the ground truth, and in Fig. 1 and Fig. 14 $C=3$ for the RGB channels. For better visualization, we plot the log of the magnitude of the Discrete Fourier Transform.

Next, we describe, from the literature, the process of performing a Discrete Fourier Transform.

Fast Fourier Transform (FFT) [9]. The discrete Fourier transform has been used in this work to convert the images from the spatial domain to the frequency domain.

“DFT is a linear operator (i.e. a matrix) that maps the data points in f to the frequency domain \hat{f} ” [10]

Equation 2.26 in [10] shows the formula to perform DFT is:

$$\hat{f}_k = \sum_{j=0}^{n-1} f_j \epsilon^{-i2\pi jk/n} \quad (5)$$

where \hat{f}_k from each sample n contains the amplitude and phase (of the sine and cosine components) information at frequency k . These are integer multiples of $\epsilon^{-2\pi j/n}$, the fundamental frequency, short-handed as ω_n [10]. Equation 2.29 in [10] shows the Discrete Fourier transform matrix (in terms of ω_n) that when multiplied by the samples in f , converts the information in those samples to frequency domain (a basis transformation). FFT is an algorithm by [9] to perform Discrete Fourier transform in an efficient manner. In Eq. (4), we use these frequencies w (referred to as k in Eq. (5)) from sample x_c obtained using an $FFT()$ function that uses the FFT algorithm.

B Additional Experiments and Ablation

Here we provide detailed results from Sec. 5 and Sec. 5.2 and additional results as mentioned in the main paper.

B.1 Semantic Segmentation

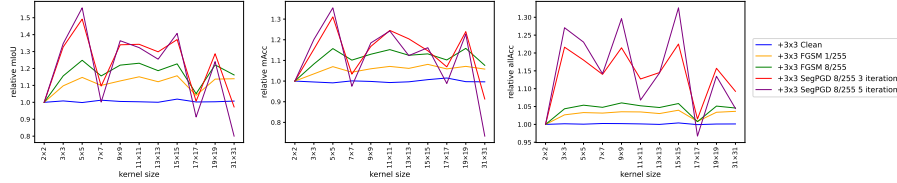


Fig. 8: Comparison of performance of different sizes of transpose convolutions from standard sizes like 2×2 as well as very large 31×31 kernels with **ConvNeXt style $11 \times 11 + 3 \times 3$ style in the decoder building blocks**. All have a parallel 3×3 kernel, as shown in Figure 4 (bottom left).

Table 7 and Table 8 provide all the results of empirical performance (across the considered upsampling blocks) on clean inputs images and input images perturbed by varying intensities of FGSM and SegPGD attacks respectively.

Limit of large kernels for Upsampling As discussed in Sec. 5.2, the performance of large kernels begins to saturate at a point. We report results from Figure 7 in tabular form in Table 9. In Table 9, we find that 13×13 appears to be the saturation point for this setting and 31×31 kernels are beyond this saturation point. While 31×31 performs worse or on-par with 17×17 , it still performs significantly better than the baseline of 2×2 . In Section 5.2 we explain the kernel size limit and that larger kernels are difficult to train. **We also find that these results further strengthen our Hypothesis 2.** For ease of understanding, we visualize the trends from Table 9 in Figure 8.

B.2 Choice of encoder

Following we aim to understand the importance of the encoder and its influence on the quality of representations later decoded during the upsampling. Consequently, we justify our choice of using ConvNeXt tiny encoder for the majority of our studies.

In Table 11 we compare different encoders: ResNet50, ConvNeXt tiny, and SLaK [51] while fixing the decoder to the baseline implementation. All encoders are pre-trained on the ImageNet-1k training dataset.

We observe that using ConvNeXt tiny and SLaK as the encoder backbone gives us significantly better performance than using ResNet50 as the encoder. This observation holds true for both clean and adversarially perturbed samples. We additionally observe that SLaK gives us marginally better performance than ConvNeXt. As shown by [51], SLaK is a significantly better encoder than ConvNeXt tiny as it provides significantly more context than ConvNeXt by using kernel sizes up to 51×51 in the convolution layers during encoding. This proves that better encoding can be harnessed during decoding which can lead to better upsampling.

However, in this work, we used the ConvNeXt tiny encoder since the SLaK encoder takes significantly longer to train for only a marginal gain in performance. We report

the performance results in Table 12. We observe that given our computation budget and the wall-clock time limit of 24 hours, we are unable to even compute the performance of the model with the SLaK encoder at 100 attack iterations.

B.3 Ablation over small parallel kernel

Following we ablate over the use of a small (3×3) kernel in parallel to a large ($\geq 7\times 7$) kernel for Large Context Transposed Convolutions. This concept is inspired by [17, 51] who use a small kernel in parallel with the large kernels to preserve local context when downsampling. Similar behavior is observed while upsampling. Table 7 compares the usage of this small parallel kernel. We observe, that while not using the small kernel results in marginal better performance on clean images (for a fixed backbone style), it lacks context and thus performs poorly (when compared to using a small parallel kernel) against adversarial attacks.

This is further highlighted in Tab. 8 when the performance is compared against strong adversarial attacks. Moreover, we observe that from medium-sized kernels i.e. , the upsampling seems to lose local context, and adding a kernel in parallel helps the model in getting this additional context. This effect can also be observed in the adversarial performances of the respective models.

B.4 Drawbacks of interpolation

As discussed in Section 3, architecture designs that use interpolation for pixel-wise upsampling suffer with over-smoothing of feature maps. This can be seen in the final predictions, as shown in Fig. 9b compared to the ground truth segmentation mask in Fig. 9a and prediction from a model with $11\times 11 + 3\times 3$ transposed convolution kernel in Fig. 9c.

In their work, [34] showed that PSPNet has considerably lower performance against adversarial attacks, similar to the analysis made in Section 5.2. This is explained by H2.

B.5 Different Upsampling Methods

Following we compare different upsampling techniques thus justifying our advocacy for using Transposed Convolution instead of other upsampling techniques like interpolation and pixel shuffle.

We report the comparison in Table 5 and observe that both Pixel shuffle and Nearest Neighbor interpolation perform better than the usually used Transposed Convolution with a 2×2 kernel size. However, as we increase the kernel size for Transposed Convolution to 11×11 with a 3×3 small kernel in parallel, we observe that Large Context Transposed Convolutions are strictly outperforming pixel shuffle, on both clean unperturbed images and under adversarial attacks, across all metrics used. Transposed Convolution with a large kernel is either outperforming or performing at par with Nearest Neighbor interpolation as well. Thus we demonstrate the superior clean and adversarial performance of large kernel-sized Transposed Convolution operation over other commonly used upsampling techniques.

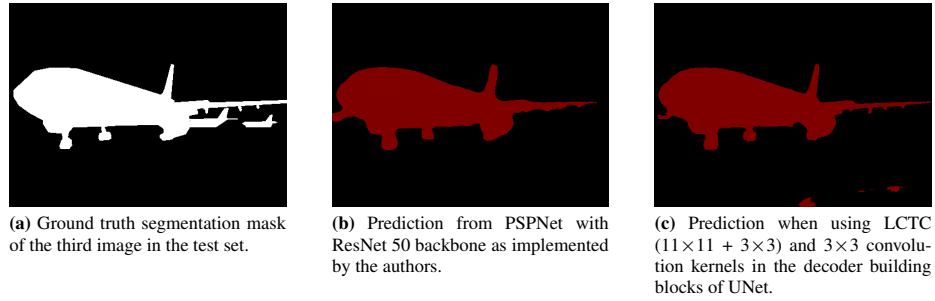


Fig. 9: A comparison of differences in the sharpness of final predictions due to different upsampling techniques. Fig. 9a is the ground truth segmentation mask with sharp and thin edges in the rear fin and wing with protrusions in the wing of the aircraft. We observe that PSPNet with a ResNet50 backbone as implemented by [93] is not able to accurately predict the thin edges and the protrusions, and is simply smoothening them out. This is due to the interpolation operation used in upsampling. However in comparison, as shown in Fig. 9c, when a transposed convolution operation is used for pixel-wise upsampling, the thin edges are sharper and protrusions are more accurately predicted.

There might be speculation if other downsampling techniques can utilize larger convolution kernels in the [decoder building blocks](#) better than transposed convolution. Thus, we additionally experiment using a ConvNeXt-like $7 \times 7 + 3 \times 3$ kernel in the Convolution operations in the decoder building blocks that follow the upsampling operation. We report these results in Table 13 and observe that similar to transposed convolution, other upsampling methods also do not benefit from an increase in the kernel size in the decoder building blocks.

B.6 Adversarial Training

Following, we present the results from adversarial training for semantic segmentation. In Table 10, we report the performance of different transposed convolution kernel-sized adversarially trained UNet on clean input and adversarially perturbed inputs. The observed performance improvement when increasing the transposed convolution kernel size during normal training also extends to adversarial training.

C Additional Results on Image Restoration

Following we provide additional results for the Image deblurring tasks, like the performance of models after adversarial training and some visual results of the deblurring for a better understanding of the impact of increased spatial context against different adversarial attack methods and strengths.

C.1 Latency Study

As PixelShuffle, when downsampling with a factor of 2, reduces the channel dims by a factor of 4, works [16, 89] use a 1×1 convolution layer before the PixelShuffle to increase the number of channels by a factor of 4. This added complexity is not needed for Transposed Convolution. Thus, in Table 14 we report the number of parameters in the models from Figure 5 and report latencies (mean over 1000 runs) of the upsampling operations, and show that these are comparable. In practice, these differences are negligible as other unchanged operations are more costly.

C.2 Adversarial Training

In Table 15 we provide additional results for adversarially training image restoration network NAFNet using FGSM attack on 50% of the training minibatch of the GoPro dataset each iteration. The state-of-the-art Image Restoration models are significantly larger w.r.t. the number of parameters, compared to the models considered for semantic segmentation. Thus, they are significantly more difficult to train adversarially. They require more training iterations. Due to the limited computing budget, we have only trained them for the same iterations as clean (non-adversarial) training iterations. We already observe the advantages of using a larger kernel for transposed convolution over pixel-shuffle in these experiments.

C.3 Visual Results

Figure 10 shows reconstruction under PGD attack for Restormer [89] and NAFNet [16]. Figure 11 shows reconstruction under CosPGD attack for Restormer [89] and NAFNet [16].

C.4 Visualizing Kernel Weights

An increase in kernel size leads to an increase in context and since the context is increased, the effect of uneven contributions of pixels is negated leading to reduced spectral artifacts. This can be seen in Figure 12. Here we observe that the weights for 3×3 are high at the edges, causing the described grid effect, whereas for 11×11 kernels there is a smooth fading towards the border of kernels, negating this effect.

C.5 Real World and Out-Of-Distribution (ODD) Generalization

Since LCTC leads to improved sampling that provides stability to feature maps learned by the network (not merely defense), inspired by observations from [29], we hypothesize that the trends on adversarial attacks should translate to Real-World noise. We show this in Table 16 by applying 2D common corruptions (CC) (severity=3) from [39] on images from the GoPro dataset and using NAFNet models from Figure 5. Since the task is deblurring, we consider all common corruptions but additional blurring and weather corruptions, as these would have to be captured before blurring.

D Additional Results Disparity Estimation

Following we report additional results for Disparity Estimation using STTR-light. In Table 17 we report the performance of STTR-light architecture on clean test images and under PGD attack. Whereas in Figure 13, we present a visual comparison of depth estimation predictions by a vanilla STTR-light as proposed by [50] and our proposed modification of **increasing the kernel size of the transposed convolution operation** in the “feature extractor” module of the architecture from 3×3 to Large Context Transposed Convolutions with kernel sizes $7\times 7+3\times 3$ and $11\times 11+3\times 3$.

D.1 Disparity Estimation Discussion

In Figure 13 as shown by the region in the **red circle**, both vanilla architecture and the architecture with our proposed change perform well compared to the ground truth on clean images. However, under a 10 iteration PGD adversarial attack, we observe small protrusion’s depth (shown by the **red arrow**) is incorrectly estimated by the vanilla architecture. The architecture with $7\times 7+3\times 3$ and $11\times 11+3\times 3$ transposed convolution kernels preserves the prediction of the disparity.

Additionally from Table 17, we observe the significance of the parallel 3×3 small kernel with the large 7×7 and 11×11 kernels. The stability of the performance of the large kernels without the small parallel kernel compared to the baseline is better. However, the stability of performance when only using larger kernels compared to larger kernels with small parallel kernels is marginally worse.

E Nomenclature: What are *Large Context Transposed Convolutions*?

In Section 4 we introduce the term “**Large Context Transposed Convolutions (LCTC)**”. In this work, we use this to describe the Transposed Convolution layers in the decoder with large kernel sizes and thus a large spatial context. However, terms like “large” are subjective, this in the following we discuss our interpretation of a “large” kernel size.

Most previous works use kernel sizes of 2×2 or 3×3 for any convolution operation, be it for downsampling [37,52] or be it for upsampling [70]. [54] introduced performing downsampling using convolution operations with a large kernel size which in their case was 7×7 . This “larger” kernel size for downsampling was further extended by other works like [17,35] to 31×31 and even up to 51×51 .

In Section 3, we show how increasing context during upsampling can reduce spectral artifacts from a theoretical perspective. Theoretically, we would want an infinite-sized kernel when performing upsampling. However, this is not practical, thus we used Transposed Convolution with kernel sizes sufficiently large to give a good trade-off between theorized context and practical trainability and compute requirements.

Thus, inspired by encoding literature [17,35,54] we use kernel sizes for upsampling that are larger than those used by previous works. Given that previous works used kernel sizes like 2×2 or 3×3 , anything bigger than this already provides more spatial context.

Thus, even a kernel size of 5×5 would be an interesting exploration and thus we explore this as well in Tab. 7 and Tab. 8.

However, given the theoretically ideal kernel size is infinity, a kernel size of 5×5 does not provide enough spatial context and thus we start calling transposed convolution operations as Large Context Transposed Convolution only when their kernel sizes are 7×7 or larger.

F Additional visualizations of Upsampling Artifacts and their Frequency Spectra

Following, we extend the example from Figure 1 to Figure 14 showing similar upsampling artifacts but on different input images to demonstrate that our findings are not limited to one example.

G Limitations

Current metrics for measuring performance do not completely account for spectral artifacts. Spectral artifacts begin affecting these metrics only when they become pronounced such as under adversarial attacks, and here Large Context Transposed Convolutions consistently perform better across tasks and architectures. Ideally, we would want infinitely large kernels, however, with increasing kernel size and task complexity, training extremely large kernels can be challenging. Thus, in this work, while having ablated over kernels as large as 31×31 , we propose using kernels only as large as 7×7 to 11×11 for good practical trade-offs. Further improvements might be possible when jointly optimizing the encoder *and* decoder of architectures.

In this work, we are focused on the reduction of spectral artifacts in upsampled images and features introduced due to the theoretical limitations of upsampling operations. However, there might exist other factors that contribute to the introduction and existence of spectral artifacts such as spatial bias. This might also present an interesting avenue to explore.

Table 7: Complete comparison of performances against FGSM attack, of UNet with ConvNeXt encoder and decoder with architectures along with different sized kernels in transposed convolution and different convolution blocks in the decoder for upscaling the feature maps.

Transposed Convolution Kernels	Backbone Style	Test Accuracy			FGSM attack ϵ_{mean}		
		mIoU	mAcc	allAcc	mIoU	mAcc	allAcc
2×2	ResNet Style 3×3	78.34	86.89	95.15	53.54	70.96	86.08
	ComNeXt style 7×7	77.17	86.86	94.81	77.42	86.24	94.94
	ComNeXt style 7×7 + 3×3	77.24	86.03	94.84	51.09	70.53	85.29
	ComNeXt style 11×11	77.68	86.42	94.97	50.73	69.78	84.88
	ComNeXt style 11×11 + 3×3	77.17	86.86	94.81	47.34	67.72	83.34
3×3	ResNet Style 3×3	78.45	86.66	95.20	53.76	70.62	86.32
	ComNeXt style 7×7	77.70	86.89	94.99	52.30	71.56	85.73
	ComNeXt style 7×7 + 3×3	77.33	87.53	94.79	50.90	72.77	83.78
	ComNeXt style 11×11	77.86	86.75	94.99	51.30	70.39	85.33
	ComNeXt style 11×11 + 3×3	77.81	86.48	94.98	51.95	70.08	85.57
5×5 (Ours)	ResNet Style 3×3	79.19	87.62	95.36	55.57	73.51	86.65
	ComNeXt style 7×7	78.52	87.39	95.13	54.4	72.48	86.29
	ComNeXt style 7×7 + 3×3	78.73	87.81	95.24	53.86	73.12	85.86
	ComNeXt style 11×11	77.83	86.99	94.91	53.76	72.8	85.96
	ComNeXt style 11×11 + 3×3	77.92	86.92	95.02	48.67	68.11	83.96
5×5 + 3×3 (Ours)	ResNet Style 3×3	78.83	87.56	95.28	56.11	73.97	86.91
	ComNeXt style 7×7	77.17	86.74	94.84	52.42	71.88	85.59
	ComNeXt style 7×7 + 3×3	77.87	86.65	94.92	52.93	72.18	85.53
	ComNeXt style 11×11	77.83	86.57	95.07	52.12	70.29	85.79
	ComNeXt style 11×11 + 3×3	77.07	86.11	94.87	54.31	72.45	86.1
LCtC: 7×7 (Ours)	ResNet Style 3×3	78.92	88.06	95.23	56.02	74.13	86.45
	ComNeXt style 7×7	77.17	86.74	94.84	52.42	71.88	85.59
	ComNeXt style 7×7 + 3×3	77.87	86.65	94.92	52.93	72.18	85.53
	ComNeXt style 11×11	77.9	87.35	94.94	53.47	72.61	85.79
	ComNeXt style 11×11 + 3×3	77.99	87.86	94.96	51.61	73.01	84.83
LCtC: 7×7 + 3×3 (Ours)	ResNet Style 3×3	78.5	87.57	95.13	53.85	72.75	85.87
	ComNeXt style 7×7	77.17	86.74	94.84	52.42	71.88	85.59
	ComNeXt style 7×7 + 3×3	77.87	86.65	94.92	52.93	72.18	85.53
	ComNeXt style 11×11	77.71	87.22	94.97	52.47	73.22	85.55
	ComNeXt style 11×11 + 3×3	77.84	86.94	95.05	52.08	70.63	85.98
LCtC: 9×9 (Ours)	ResNet Style 3×3	78.36	86.88	95.18	55.62	72.62	86.9
	ComNeXt style 7×7	77.17	86.74	94.84	52.42	71.88	85.59
	ComNeXt style 7×7 + 3×3	77.93	86.97	95.04	51.01	70.59	84.87
	ComNeXt style 11×11	77.80	86.80	94.99	52.42	72.22	85.39
	ComNeXt style 11×11 + 3×3	77.85	86.71	95.07	54.59	72.04	86.48
LCtC: 9×9 + 3×3 (Ours)	ResNet Style 3×3	78.77	87.77	95.24	55.94	73.79	86.67
	ComNeXt style 7×7	77.17	86.74	94.84	52.42	71.88	85.59
	ComNeXt style 7×7 + 3×3	77.96	87.24	94.98	51.21	70.01	85.24
	ComNeXt style 11×11	77.92	86.82	95.03	52.71	71.17	86.02
	ComNeXt style 11×11 + 3×3	77.57	86.71	95.02	53.32	71.75	86.29
LCtC: 11×11 (Ours)	ResNet Style 3×3	79.11	87.06	95.36	56.18	72.11	87.27
	ComNeXt style 7×7	77.17	86.74	94.84	52.42	71.88	85.59
	ComNeXt style 7×7 + 3×3	78.34	87.06	95.07	51.93	71.19	85.54
	ComNeXt style 11×11	77.42	86.68	94.94	53.11	71.43	86.03
	ComNeXt style 11×11 + 3×3	77.75	86.83	95.01	52.88	71.47	85.93
LCtC: 11×11 + 3×3 (Ours)	ResNet Style 3×3	79.33	87.81	95.41	58.04	74.93	87.8
	ComNeXt style 7×7	78.32	86.98	95.09	53.41	72.45	86.18
	ComNeXt style 7×7 + 3×3	78.64	86.78	95.17	54.32	71.27	86.63
	ComNeXt style 11×11	77.75	85.93	94.87	51.19	69.72	85.45
	ComNeXt style 11×11 + 3×3	77.42	86.24	94.94	54.48	72.53	86.25
LCtC: 13×13 (Ours)	ResNet Style 3×3	79.41	88.18	95.36	56.89	74.71	87.36
	ComNeXt style 7×7	77.99	87.11	95.06	54.96	73.32	86.69
	ComNeXt style 7×7 + 3×3	78.44	87.22	95.13	54.21	72.18	86.34
	ComNeXt style 11×11	77.57	85.99	95.00	53.51	70.31	86.67
	ComNeXt style 11×11 + 3×3	77.40	86.53	94.89	53.16	71.62	86.12
LCtC: 13×13 + 3×3 (Ours)	ResNet Style 3×3	79.17	87.96	95.38	57.17	75.08	87.44
	ComNeXt style 7×7	78.05	86.73	95.02	53.41	71.62	86.12
	ComNeXt style 7×7 + 3×3	77.76	86.14	95.06	54.09	72.11	86.29
	ComNeXt style 11×11	77.81	87.43	95.01	51.71	71.77	85.25
	ComNeXt style 11×11 + 3×3	77.20	86.55	94.81	53.1	71.88	85.87
LCtC: 15×15 (Ours)	ResNet Style 3×3	79.17	87.68	95.28	58.08	73.56	87.58
	ComNeXt style 7×7	78.34	87.14	95.03	53.86	72.77	86.11
	ComNeXt style 7×7 + 3×3	77.39	86.40	94.95	51.2	69.42	85.27
	ComNeXt style 11×11	77.14	86.36	94.82	50.14	69.32	84.49
	ComNeXt style 11×11 + 3×3	77.67	86.78	94.90	54.44	72.74	86.54
LCtC: 15×15 + 3×3 (Ours)	ResNet Style 3×3	78.72	87.50	95.25	56.28	73.97	87.15
	ComNeXt style 7×7	77.56	87.01	94.91	53.28	72.15	85.78
	ComNeXt style 7×7 + 3×3	77.09	86.27	94.76	52.25	70.01	85.41
	ComNeXt style 11×11	77.40	86.39	94.92	53.59	71.49	86.21
	ComNeXt style 11×11 + 3×3	77.64	87.46	95.20	54.77	73.2	86.65
LCtC: 17×17 (Ours)	ResNet Style 3×3	79.22	87.77	95.37	56.5	73.3	87.27
	ComNeXt style 7×7	77.36	87.64	94.99	54.06	73.88	85.84
	ComNeXt style 7×7 + 3×3	78.03	87.56	95.01	52.75	72.0	85.65
	ComNeXt style 11×11	77.82	87.40	94.92	51.43	70.57	85.22
	ComNeXt style 11×11 + 3×3	77.74	86.69	94.99	51.31	69.71	85.53
LCtC: 17×17 + 3×3 (Ours)	ResNet Style 3×3	78.41	86.84	95.26	56.03	73.28	87.16
	ComNeXt style 7×7	78.14	86.99	94.98	53.44	72.34	86.01
	ComNeXt style 7×7 + 3×3	78.62	87.64	95.14	55.54	73.87	86.83
	ComNeXt style 11×11	77.59	87.73	94.84	52.84	74.14	84.63
	ComNeXt style 11×11 + 3×3	77.33	88.15	94.75	49.29	71.71	84.04
LCtC: 19×19 (Ours)	ResNet Style 3×3	78.54	87.64	95.12	56.63	74.09	87.25
	ComNeXt style 7×7	78.74	87.66	95.15	56.28	73.79	87.11
	ComNeXt style 7×7 + 3×3	77.05	86.33	94.89	54.47	72.38	86.78
	ComNeXt style 11×11	77.66	86.61	95.00	51.58	71.51	84.83
	ComNeXt style 11×11 + 3×3	77.61	86.59	94.93	50.34	69.39	84.54
LCtC: 19×19 + 3×3 (Ours)	ResNet Style 3×3	78.78	87.34	95.28	56.53	74.59	86.97
	ComNeXt style 7×7	77.44	86.70	94.91	54.05	72.52	86.99
	ComNeXt style 7×7 + 3×3	78.14	87.14	95.02	55.82	74.54	86.96
	ComNeXt style 11×11	78.03	86.64	95.08	53.5	71.21	86.26
	ComNeXt style 11×11 + 3×3	77.42	86.61	94.91	53.83	72.54	86.17
LCtC: 31×31 (Ours)	ResNet Style 3×3	78.69	86.98	95.30	56.61	73.22	87.08
	ComNeXt style 7×7	77.54	87.30	94.84	52.36	72.27	85.14
	ComNeXt style 7×7 + 3×3	76.96	86.38	94.77	53.59	72.14	86.05
	ComNeXt style 11×11	76.84	86.72	94.71	50.74	70.53	84.61
	ComNeXt style 11×11 + 3×3	76.77	85.60	94.71	51.42	69.17	85.2
LCtC: 31×31 + 3×3 (Ours)	ResNet Style 3×3	78.47	87.26	95.16	56.27	73.39	87.22
	ComNeXt style 7×7	77.43	86.56	94.93	53.45	72.74	86.17
	ComNeXt style 7×7 + 3×3	78.43	87.07	95.17	56.72	73.65	87.6
	ComNeXt style 11×11	78.00	87.04	94.94	50.66	70.23	84.83
	ComNeXt style 11×11 + 3×3	77.73	86.54	94.93	53.94	71.65	86.39

Table 8: Comparison of performances against SegPGD attack, of UNet with ConvNeXt encoder and decoder with architectures along with different sized kernels in transposed convolution and different convolution blocks in the decoder for upscaling the feature maps.

Transposed Convolution Kernels	Backbone Style	SegPGD attack iterations									
		3		5		10		20		40	
		mIoU	mAcc	mIoU	mAcc	mIoU	mAcc	mIoU	mAcc	mIoU	mAcc
2×2	ResNet Style 3×3	23.06	46.51	60.04	14.43	35.50	45.30	08.12	24.67	29.88	05.54
	ComNeXt style 7×7	17.94	0.4481	47.96	10.64	33.63	30.64	05.47	21.74	15.8	03.2
	ComNeXt style 7×7 + 3×3	17.59	42.55	05.168	09.85	30.41	03.233	04.75	16.83	03.41	02.65
	ComNeXt style 11×11	16.39	0.4013	0.485	09.37	28.66	29.63	03.97	14.16	11.41	01.56
	ComNeXt style 11×11 + 3×3	13.97	35.82	45.68	07.61	25.07	28.33	03.4	14.38	12.04	02.21
3×3	ResNet Style 3×3	23.37	46.33	60.78	15.26	38.0	46.51	09.26	29.64	31.9	06.78
	ComNeXt style 7×7	17.59	43.57	51.41	09.9	30.84	33.14	04.74	18.3	14.55	02.23
	ComNeXt style 7×7 + 3×3	19.08	46.97	47.74	11.15	34.6	29.9	05.96	22.62	15.67	03.61
	ComNeXt style 11×11	16.2	39.11	50.93	09.52	29.32	32.61	04.93	20.31	14.82	02.86
	ComNeXt style 11×11 + 3×3	18.54	41.34	55.56	10.25	30.11	36.0	04.8	19.25	13.94	02.41
5×5 (Ours)	ResNet Style 3×3	24.23	51.8	57.82	16.16	42.98	43.29	10.11	32.79	30.3	07.32
	ComNeXt style 7×7	17.59	43.57	51.41	09.9	30.84	33.14	04.74	18.3	14.55	02.23
	ComNeXt style 7×7 + 3×3	18.7	43.18	52.74	10.56	31.41	33.32	04.87	18.5	14.78	02.49
	ComNeXt style 11×11	18.96	44.79	53.09	09.85	29.77	32.88	03.83	14.94	12.6	01.94
	ComNeXt style 11×11 + 3×3	13.38	33.61	45.0	06.84	20.94	25.99	02.51	08.85	08.5	01.18
5×5 + 3×3 (Ours)	ResNet Style 3×3	25.03	53.96	58.89	16.61	45.8	42.18	10.79	37.16	27.34	08.0
	ComNeXt style 7×7	17.65	44.79	48.41	09.79	31.78	28.51	04.62	18.37	11.12	02.58
	ComNeXt style 7×7 + 3×3	18.31	42.75	49.26	09.89	28.58	30.02	03.78	12.49	11.08	01.34
	ComNeXt style 11×11	17.87	40.62	52.77	09.74	27.94	34.21	04.65	14.98	14.34	02.0
	ComNeXt style 11×11 + 3×3	20.84	46.95	53.91	11.86	33.96	34.86	05.65	19.8	16.66	02.83
LTC: 7×7 (Ours)	ResNet Style 3×3	26.53	53.05	61.16	17.75	43.31	46.99	10.26	30.92	32.62	07.17
	ComNeXt style 7×7	17.64	43.32	47.8	09.95	30.43	28.02	04.21	15.08	10.07	01.86
	ComNeXt style 7×7 + 3×3	16.64	40.11	50.56	09.75	29.72	32.3	04.85	13.95	12.25	01.91
	ComNeXt style 11×11	17.37	45.07	47.32	08.86	30.03	26.48	03.47	14.22	07.94	01.53
	ComNeXt style 11×11 + 3×3	17.07	42.3	48.78	09.31	28.04	28.88	03.82	13.79	09.54	01.8
LTC: 7×7 + 3×3 (Ours)	ResNet Style 3×3	24.03	52.08	57.43	16.21	43.38	43.01	09.99	32.77	30.22	08.3
	ComNeXt style 7×7	16.19	43.4	48.59	09.02	32.38	29.17	04.23	19.63	10.47	02.46
	ComNeXt style 7×7 + 3×3	16.04	39.67	48.16	08.94	27.27	35.06	03.81	14.69	12.79	01.91
	ComNeXt style 11×11	18.08	46.24	50.64	10.18	33.17	31.35	04.49	18.33	12.04	02.01
	ComNeXt style 11×11 + 3×3	15.31	37.02	52.08	07.62	24.44	32.35	03.3	15.04	12.35	01.92
LTC: 9×9 (Ours)	ResNet Style 3×3	25.26	50.75	60.85	16.88	41.02	47.16	09.44	28.03	33.87	06.23
	ComNeXt style 7×7	16.56	36.5	53.58	08.74	23.95	35.67	04.01	13.92	16.64	02.13
	ComNeXt style 7×7 + 3×3	16.2	39.55	50.82	09.0	28.53	33.31	04.64	13.95	12.25	01.91
	ComNeXt style 11×11	17.02	43.01	48.45	08.92	28.35	28.13	03.64	14.36	10.06	01.17
	ComNeXt style 11×11 + 3×3	19.34	43.6	54.41	10.71	31.22	33.98	04.6	15.76	12.75	01.98
LTC: 9×9 + 3×3 (Ours)	ResNet Style 3×3	24.87	55.04	57.35	17.0	46.34	42.08	10.88	36.04	28.55	07.91
	ComNeXt style 7×7	16.56	36.5	53.58	08.74	23.95	35.67	04.01	13.92	16.64	02.13
	ComNeXt style 7×7 + 3×3	16.01	36.92	51.5	08.8	25.53	33.15	04.64	13.95	12.25	01.91
	ComNeXt style 11×11	16.42	39.19	51.71	08.32	26.64	31.11	03.66	15.7	11.61	01.94
	ComNeXt style 11×11 + 3×3	18.72	41.83	55.48	10.38	29.7	36.72	04.74	18.16	17.44	02.49
LTC: 11×11 (Ours)	ResNet Style 3×3	26.02	48.81	63.7	16.8	39.62	49.72	09.62	29.4	34.22	06.85
	ComNeXt style 7×7	19.04	45.39	52.63	10.17	32.3	32.46	04.58	20.16	13.36	02.44
	ComNeXt style 7×7 + 3×3	16.08	39.09	53.1	08.86	24.27	35.06	03.94	16.77	12.79	01.91
	ComNeXt style 11×11	18.09	40.72	53.7	09.93	29.6	34.68	04.55	18.22	14.17	02.21
	ComNeXt style 11×11 + 3×3	15.29	37.2	50.71	07.6	25.19	30.65	03.17	15.06	09.58	01.78
LTC: 11×11 + 3×3 (Ours)	ResNet Style 3×3	27.49	53.08	64.13	18.15	43.51	49.36	10.29	31.12	33.17	07.08
	ComNeXt style 7×7	18.14	40.58	50.39	08.98	27.2	31.4	03.34	15.36	12.29	01.93
	ComNeXt style 7×7 + 3×3	17.7	39.71	54.64	09.71	26.92	35.8	04.32	13.93	15.8	02.37
	ComNeXt style 11×11	14.62	34.73	49.37	07.26	22.21	29.37	02.76	12.24	10.69	01.23
	ComNeXt style 11×11 + 3×3	18.76	44.6	51.49	10.07	31.15	30.26	04.4	17.02	10.56	02.31
LTC: 13×13 (Ours)	ResNet Style 3×3	28.51	57.18	63.94	19.71	48.99	50.08	11.99	37.69	33.26	08.31
	ComNeXt style 7×7	19.29	46.62	55.13	12.32	34.21	35.91	04.14	21.39	16.39	02.16
	ComNeXt style 7×7 + 3×3	20.13	42.92	57.7	11.38	29.96	39.57	04.85	15.81	19.37	02.54
	ComNeXt style 11×11	18.65	39.48	56.4	10.02	27.46	38.02	04.69	17.27	19.03	02.47
	ComNeXt style 11×11 + 3×3	18.95	42.88	55.82	10.68	31.21	35.69	04.92	18.29	12.63	02.35
LTC: 13×13 + 3×3 (Ours)	ResNet Style 3×3	28.08	58.22	63.4	19.4	50.01	48.89	12.04	39.2	32.11	08.77
	ComNeXt style 7×7	18.62	41.52	51.26	10.23	30.56	30.5	04.37	16.47	15.67	02.16
	ComNeXt style 7×7 + 3×3	16.7	41.09	50.56	08.54	26.94	30.39	03.31	13.44	11.22	01.53
	ComNeXt style 11×11	14.1	36.4	47.79	07.01	23.32	27.54	02.87	12.05	10.26	01.51
	ComNeXt style 11×11 + 3×3	18.14	43.14	52.31	09.55	28.16	32.41	03.54	12.57	11.22	01.52
LTC: 15×15 (Ours)	ResNet Style 3×3	29.41	51.54	66.7	19.96	41.26	55.14	11.51	29.26	41.04	07.17
	ComNeXt style 7×7	18.62	44.42	51.51	10.55	32.47	32.54	04.69	18.66	12.4	02.64
	ComNeXt style 7×7 + 3×3	17.63	37.55	55.52	09.13	23.28	37.06	03.46	09.05	15.8	01.41
	ComNeXt style 11×11	15.24	36.62	49.45	08.05	24.89	31.93	03.68	14.62	13.68	02.03
	ComNeXt style 11×11 + 3×3	19.01	44.78	52.74	10.35	31.98	32.35	04.38	19.15	12.42	02.31
LTC: 15×15 + 3×3 (Ours)	ResNet Style 3×3	26.38	53.79	61.59	17.9	45.03	47.36	10.79	33.9	31.75	07.14
	ComNeXt style 7×7	19.81	42.18	53.73	11.14	28.25	37.03	04.72	13.45	17.41	01.92
	ComNeXt style 7×7 + 3×3	17.53	39.4	54.39	09.51	26.67	35.52	03.73	12.39	13.27	00.93
	ComNeXt style 11×11	16.69	39.29	52.18	08.78	27.55	32.8	03.72	17.08	12.37	02.06
	ComNeXt style 11×11 + 3×3	19.15	41.08	55.96	10.71	29.12	37.58	05.28	19.35	18.54	02.88
LTC: 17×17 (Ours)	ResNet Style 3×3	27.74	53.24	64.48	18.51	43.47	51.02	10.72	32.21	36.08	07.43
	ComNeXt style 7×7	19.82	40.01	54.08	11.09	32.71	36.79	03.35	19.1	17.49	02.67
	ComNeXt style 7×7 + 3×3	16.96	38.94	54.19	09.23	26.92	36.22	04.47	16.8	16.9	02.45
	ComNeXt style 11×11	13.72	34.03	48.09	06.57	20.94	26.93	02.32	09.22	07.84	01.08
	ComNeXt style 11×11 + 3×3	14.47	33.55	52.16	07.33	20.91	32.91	02.82	09.75	13.11	01.28
LTC: 17×17 + 3×3 (Ours)	ResNet Style 3×3	26.97	54.13	62.04	18.41	45.5	47.66	11.05	34.55	32.01	07.43
	ComNeXt style 7×7	17.96	41.81	54.93	09.89	27.73	35.7	03.85	15.35	15.1	01.95
	ComNeXt style 7×7 + 3×3	19.86	42.55	56.89	10.29	28.26	37.83	04.43	15.77	16.52	01.93
	ComNeXt style 11×11	16.84	44.91	45.35	09.41	31.25	26.04	03.8	14.79	09.08	01.4
	ComNeXt style 11×11 + 3×3	14.06	38.28	46.37	06.95	24.78	27.4	02.92	14.65	10.36	01.55
LTC: 19×19 (Ours)	ResNet Style 3×3	27.64	52.62	64.53	18.46	42.79	51.19	10.49	30.27	36.37	06.92
	ComNeXt style 7×7	20.28	46.96	56.75	10.06	30.29	36.2	03.5	13.07	14.56	01.51
	ComNeXt style 7×7 + 3×3	18.14	39.86	56.98	09.34	26.13	37.16	03.5	13.58	14.24	01.52
	ComNeXt style 11×11	15.85	40.55	46.13	08.33	28.21	25.47	03.27	15.94	07.52	01.8
	ComNeXt style 11×11 + 3×3	16.17	37.68	50.29	08.47	25.22	31.86	03.84	14.74	13.88	01.93
LTC: 19×19 + 3×3 (Ours)	ResNet Style 3×3	28.62	55.15	63.93	19.19	47.17	48.9	10.96	34.41	31.54	07.15
	ComNeXt style 7×7	17.45	40.44	51.24	09.13	27.29	31.41	03.56	13.76	10.77	01.61
	ComNeXt style 7×7 + 3×3	20.9	48.61	54.54	11.59	35.05	33.3	04.6	18.94	11.57	02.05
	ComNeXt style 11×11	19.01	41.41	55.17	10.56	28.49	37.54	05.23	19.35	17.13	02.92

Table 9: Comparison of performance of Large Context Transposed Convolutions (LCTC) with very large: 31×31 kernels in transposed convolution to large (7×7 to 17×17) kernels. All have a parallel 3×3 kernel, as shown in Figure 4 (bottom left). Here we observe the saturation of performance for very large kernels for upsampling. This comparison is for the same encoder (ConvNeXt) and same ResNet-like building blocks in the decoder (our baseline). The complete table is provided in Appendix B.1.

Transposed Convolution Kernels	Test Accuracy			FGSM attack epsilon						SegPGD attack iterations		
	mIoU	mAcc	allAcc	$\frac{1}{255}$			$\frac{8}{255}$			20		
				mIoU	mAcc	allAcc	mIoU	mAcc	allAcc	mIoU	mAcc	allAcc
7×7	78.50	87.57	95.13	53.85	72.75	85.87	47.10	67.57	82.04	7.38	26.16	26.11
11×11	79.33	87.81	95.41	58.04	74.93	87.80	51.25	69.31	84.64	7.08	23.30	26.82
15×15	78.72	87.50	95.25	56.28	73.97	87.15	49.50	68.69	83.53	7.14	25.02	25.01
17×17	78.41	86.84	95.26	56.03	73.28	87.16	49.65	67.95	83.74	7.43	25.65	24.78
19×19	78.78	87.34	95.28	56.53	74.59	86.97	50.60	69.95	83.98	7.15	25.67	23.60
31×31	78.47	87.26	95.16	56.27	73.39	87.22	49.66	68.81	83.92	7.24	25.06	26.54

Table 10: Adversarially trained models using FGSM and PGD from Table 2 tested against adversarial attacks on UNet with ConvNeXt encoder and decoder with different sized kernels in the transposed convolution for upscaling, while keeping rest of the architecture identical.

Transposed Convolution Kernels	Clean Test Accuracy			FGSM attack epsilon						SegPGD attack iterations					
	mIoU	mAcc	allAcc	$\frac{1}{255}$			$\frac{8}{255}$			3			20		
				mIoU	mAcc	allAcc	mIoU	mAcc	allAcc	mIoU	mAcc	allAcc	mIoU	mAcc	allAcc
FGSM training															
2×2 (baseline)	78.57	86.68	95.23	54.28	70.80	86.91	52.45	68.38	86.26	26.59	48.99	67.71	7.6	24.06	31.37
LCTC: 7×7 (Ours)	78.41	86.22	95.20	56.87	72.92	87.70	51.31	68.4	85.17	28.11	53.39	66.30	8.36	28.54	28.13
LCTC: 11×11 + 3×3 (Ours)	79.57	88.1	95.3	57.90	74.64	87.61	52.15	70.23	84.96	30.37	55.54	68.3	9.4	29.79	32.37
PGD training with 3 attack iterations															
2×2 (baseline)	75.33	84.66	94.39	53.87	72.17	86.58	58.57	73.93	89.01	29.38	57.82	66.67	9.39	33.15	28.11
LCTC: 7×7 (Ours)	75.79	84.89	94.38	54.82	72.31	86.80	61.29	74.33	89.96	31.12	58.36	68.58	10.24	33.99	31.14
LCTC: 11×11 + 3×3 (Ours)	75.90	86.60	94.30	56.27	75.66	86.68	63.02	76.17	90.42	33.50	58.34	71.50	10.77	32.23	37.36

Table 11: Comparison of performances of different encoders in the UNet-like architecture. All architectures here have the baseline 2×2 transposed convolution kernel for upsampling followed by 3×3 convolution kernels in the decoder blocks. For more results please refer to Table 12.

Encoder	Test Accuracy			FGSM attack epsilon						SegPGD attack iterations		
	mIoU	mAcc	allAcc	$\frac{1}{255}$			$\frac{8}{255}$			20		
				mIoU	mAcc	allAcc	mIoU	mAcc	allAcc	mIoU	mAcc	allAcc
ResNet50	67.69	79.04	92.80	36.78	58.41	78.16	32.60	52.63	74.56	4.98	19.28	21.07
ConvNeXt tiny	78.45	86.66	95.20	53.76	70.62	86.32	47.33	64.58	83.16	5.54	18.79	23.72
SLaK tiny	78.82	87.01	95.17	55.22	71.72	86.97	48.69	66.45	83.57	8.45	25.42	32.37

Table 12: Comparison of performances of different encoders in the UNet-like architecture. All architectures here have the baseline 2×2 transposed convolution kernel followed by 3×3 convolution kernels in the decoder block.

Encoder	Test Accuracy			FGSM attack epsilon						SegPGD attack iterations					
	mIoU	mAcc	allAcc	$\frac{1}{255}$			3			5			10		
				mIoU	mAcc	allAcc	mIoU	mAcc	allAcc	mIoU	mAcc	allAcc	mIoU	mAcc	allAcc
ResNet50	67.69	79.04	92.80	36.78	58.41	78.16	32.60	52.63	74.56	4.98	19.28	21.07	3.95	16.49	18.35
ConvNeXt tiny	78.45	86.66	95.20	53.76	70.62	86.32	47.33	64.58	83.16	5.54	18.79	23.72	4.39	14.98	23.70
SLaK tiny	78.82	87.01	95.17	55.22	71.72	86.97	48.69	66.45	83.57	8.45	25.42	32.37	6.22	19.58	29.06

Table 13: Comparison of performances of different upsampling methods in the UNet-like architecture. All architectures here have the baseline i.e. ConvNeXt encoder and a ResNet style 3×3 or ConvNext style $7 \times 7 + 3 \times 3$ convolution kernels in the decoder block.

Upsampling Method	Convolution Kernel in Decoder blocks	Test Accuracy		FGSM attack epsilon						SupPGD attack iterations					
		mIoU	mAcc	mIoU	mAcc	mIoU	mAcc	mIoU	mAcc	mIoU	mAcc	mIoU	mAcc	mIoU	mAcc
Pixel Shuffle	ResNet Style 3×3	78.54	87.32	85.18	53.82	71.58	85.88	46.67	65.03	81.71	23.08	48.18	56.54	15.06	38.85
	ConvNeXt Style $7 \times 7 + 3 \times 3$	77.19	85.90	94.88	51.78	69.68	85.44	43.80	62.24	81.06	17.52	40.16	50.31	9.43	27.37
Nearest Neighbour Interpolation	ResNet Style 3×3	78.40	88.16	95.09	52.68	73.51	84.45	46.08	67.96	80.22	22.82	53.16	51.75	15.34	44.53
	ConvNeXt Style $7 \times 7 + 3 \times 3$	77.86	86.92	94.97	50.71	71.21	84.45	41.97	64.92	78.89	15.77	44.36	42.09	8.56	30.25
Transposed Convolution 2×2	ResNet Style 3×3	78.45	86.66	95.20	53.76	70.62	86.32	47.33	64.58	83.16	23.06	46.51	60.04	14.43	35.50
	ConvNeXt Style $7 \times 7 + 3 \times 3$	77.24	86.03	94.84	51.09	70.53	85.29	43.52	63.74	81.19	17.59	42.55	51.69	9.88	30.41
LCTC: $11 \times 11 + 3 \times 3$ (Ours)	ResNet Style 3×3	79.33	87.81	95.41	58.04	74.93	87.8	51.25	69.31	84.64	27.49	53.08	64.13	18.15	43.51
	ConvNeXt Style $7 \times 7 + 3 \times 3$	78.64	86.78	95.17	54.32	71.27	86.63	45.48	63.62	82.32	17.7	39.71	54.64	9.71	26.92

Table 14: Comparing latency and number of parameters for models from Figure 5.

Upsampling Method	Latency (ms)	No. of Params
Pixel Shuffle	0.26	17.11 M
Trans. Conv. 3×3	0.27	16.43 M
LCTC $11 \times 11 + 3 \times 3$	0.38	16.54 M

Table 15: Comparison of performances of adversarially trained *SotA* Image Restoration Networks. The considered architectures use Pixel Shuffle for Upsampling, we propose replacing the Pixel Shuffle with Transposed Convolution operations using the large filter. Testing for image deblurring on GoPro dataset.

Network	Upsampling Method	Test Accuracy		PGD attack iterations			
		PSNR	SSIM	5	10	20	
NAFNet + ADV	Pixel Shuffle	29.91	0.9291	15.76	0.5228	13.91	0.4445
	Transposed Conv 3×3	31.26	0.9448	15.89	0.5390	13.43	0.4627
	LCTC: $7 \times 7 + 3 \times 3$ (Ours)	31.21	0.9446	16.46	0.5061	14.55	0.4211
	LCTC: $11 \times 11 + 3 \times 3$ (Ours)	30.70	0.9390	13.68	0.4857	11.91	0.4085

Table 16: Performance of different upsampling methods in NAFNet in real-world ODD setting by applying 2D common corruption [39] (severity=3) on GoPro dataset. We use all common corruptions from [39] GitHub repository except weather conditions (ideally these should happen before the motion blurring) and blurring (since the images are already motion blurred). Here “Mean” is performance over all the considered corruptions:

Common Corruption	Upsampling Method					
	Pixel Shuffle		Trans. Conv. 3×3		LCTC $11 \times 11 + 3 \times 3$	
	PSNR	SSIM	PSNR	SSIM	PSNR	SSIM
Gaussian Noise	4.8501	0.0104	8.7346	0.1014	13.6475	0.1523
Shot Noise	4.8616	0.0127	8.9524	0.0984	13.2464	0.1564
Impulse Noise	5.0154	0.0214	9.2451	0.1065	14.8425	0.187
Brightness	32.3199	0.9576	30.676	0.9394	30.4098	0.9361
Contrast	26.5941	0.7759	25.9743	0.7561	25.8733	0.7525
Elastic Transform	17.944	0.6392	19.7686	0.703	19.7672	0.702
Pixelate	4.4977	0.246	4.4999	0.246	4.4958	0.246
JPEG Compression	25.2767	0.8095	25.1014	0.8032	25.3788	0.8104
Speckle Noise	4.8287	0.0158	9.2336	0.1044	14.6622	0.2473
Saturate	32.1969	0.958	30.5904	0.9399	30.3005	0.9365
Mean	15.8385	0.4447	17.2776	0.4798	19.262	0.5127

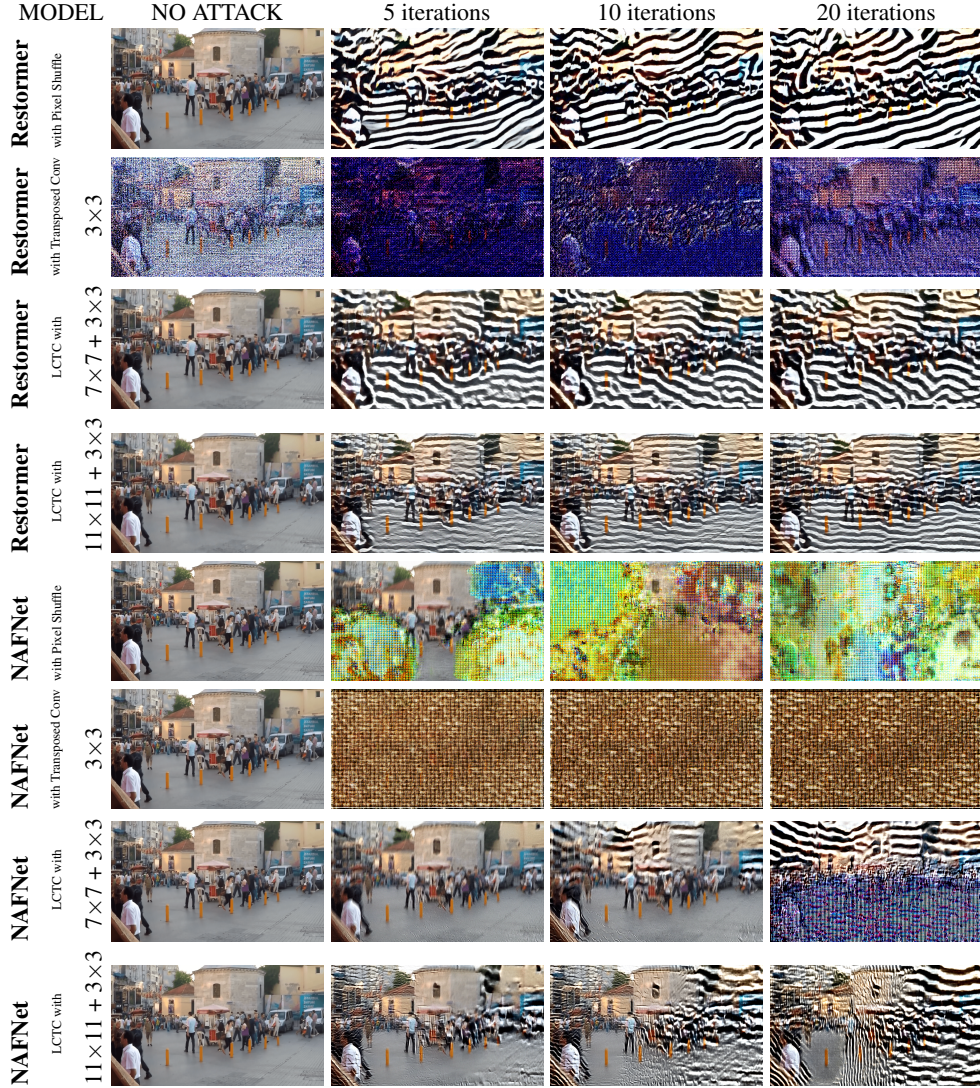


Fig. 10: Comparing images reconstructed by all models after **PGD attack** on variants of Upsampling.

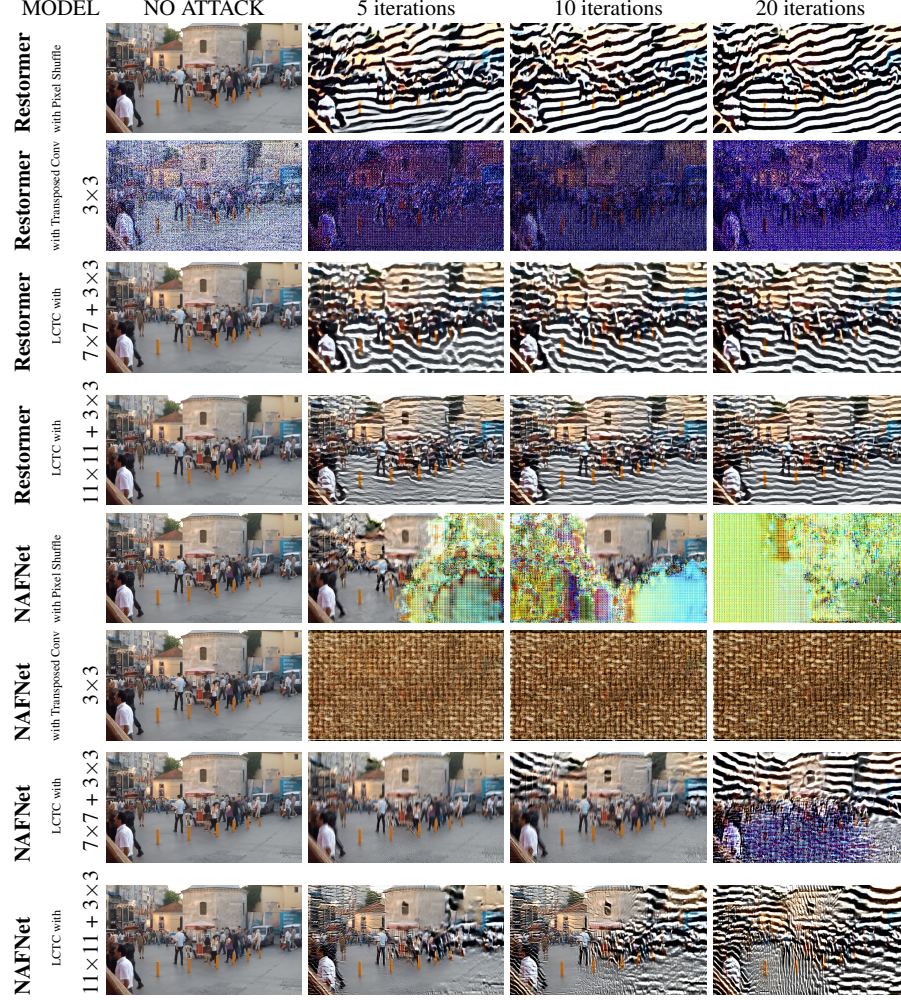


Fig. 11: Comparing images reconstructed by the considered variants of the SotA models after CosPGD attack [3]. We observe that the originally proposed Restormer and NAFNet architectures that use Pixel Shuffle for upsampling perform considerably well under no adversarial attack but even a small perturbation of $\epsilon = \frac{8}{255}$ causes ringing and other spectral artifacts to occur in the deblurred images to the extent that the images are unrecognizable. However, on replacing the Pixel Shuffle operation in these architectures with a Transposed Convolution operation with a large kernel ($11 \times 11 + 3 \times 3$), we observe a significant reduction in the spectral artifacts in the images restored under adversarial attack while the image restored under no attack are very comparable to those restored by the original architectures.

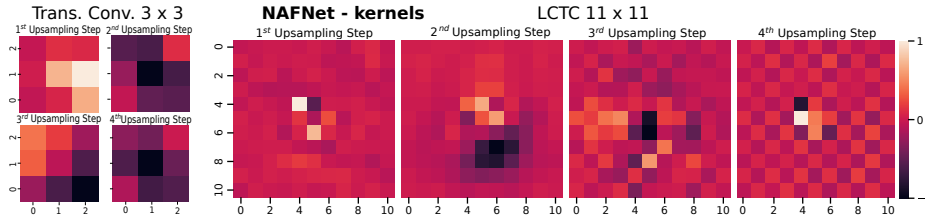


Fig. 12: Normalized kernel weights from a random channel each for the models from Figure 5.

Table 17: Comparison of performance of STTR-light architecture with different sized kernels in transposed convolution for upscaling the feature maps in the feature extractor.

Transposed Convolution Kernels	Test Accuracy		PGD Attack					
	epe↓	3px error↓	3 Iterations		5 Iterations		10 Iterations	
	epe↓	3px error↓	epe↓	3px error↓	epe↓	3px error↓	epe↓	3px error↓
STTR-light [50] reported	0.5	1.54						
3×3 [50] reproduced	0.4927	1.54	4.05	18.46	4.07	18.59	4.08	18.6
LCTC: 7×7 (Ours)	0.487	1.52	4.26	19.09	4.289	19.21	4.294	19.23
LCTC: 7×7 + 3×3 (Ours)	0.4788	1.50	4.02	18.3	4.0474	18.43	4.05	18.45
LCTC: 9×9 (Ours)	0.4983	1.50	4.36	18.02	4.386	18.14	4.39	18.16
LCTC: 11×11+3×3 (Ours)	0.5124	1.57	4.004	18.29	4.028	18.42	4.032	18.44

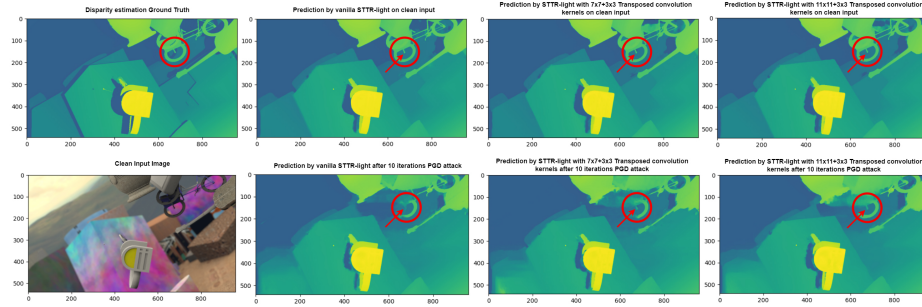


Fig. 13: Visual comparison of Disparity Estimation predictions by a vanilla STTR-light as proposed by [50] and our proposed modification of **increasing the kernel size of the transposed convolution operation** in the “feature extractor” module of the architecture from 3×3 to LCTC with 7×7+3×3 and 11×11+3×3 sized kernels. As shown by the region in the **red circle**, both vanilla architecture and the architecture with our proposed change perform well compared to the ground truth on clean images. However, under 10 iteration PGD adversarial attack, we observe small protrusion’s depth (shown by the **red arrow**) is incorrectly estimated by the vanilla architecture, however, the architectures with LCTC preserve the prediction of the disparity.

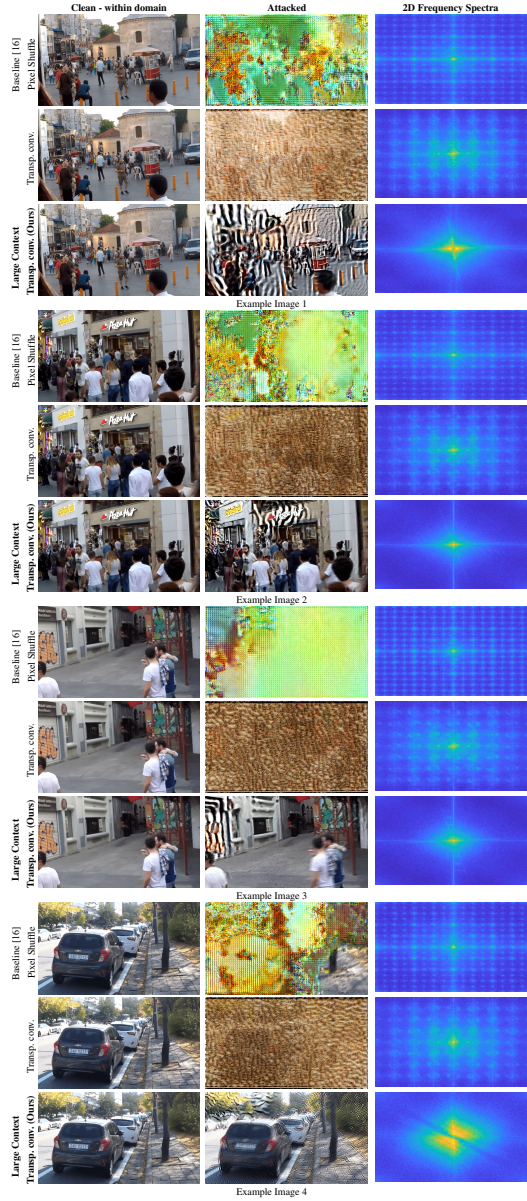


Fig. 14: This is extension to Fig. 1, here we observe the same artifacts both in the spatial and frequency domain as that observed in Fig. 1. Here we perform Image restoration using NAFNet [16] variants on GoPro [63]. Normal Transposed Convolution uses 3×3 sized kernels. Large Context Transposed Convolution uses kernels of size $7 \times 7 + 3 \times 3$ for upsampling. LCTC significantly increases the model's stability during upsampling, observable in the restored image under attack and the frequency spectrum. The procedure for obtaining the 2D Frequency Spectra has been explained in Appendix A.6.