

# ANALYZING THE TRAINING DYNAMICS OF IMAGE RESTORATION TRANSFORMERS: A REVISIT TO LAYER NORMALIZATION

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

This work analyzes the training dynamics of Image Restoration (IR) Transformers and uncovers a critical yet overlooked issue: conventional LayerNorm (LN) drives feature magnitudes to diverge to a *million scale* and collapses channel-wise entropy. We analyze this in the perspective of networks attempting to bypass LayerNorm’s constraints, which conflict with IR tasks. Accordingly, we address two misalignments: 1) *per-token normalization* that disrupts spatial correlations, and 2) *input-independent scaling* that discards input-specific statistics. To address this, we propose Image Restoration Transformer Tailored Layer Normalization (*i*-LN), a simple drop-in replacement that normalizes features holistically and adaptively rescales them per input. We provide theoretical insights and empirical evidence that this design effectively captures important spatial correlations and better preserves input-specific statistics throughout the network. Experimental results verify that the proposed *i*-LN consistently outperforms vanilla LN on various IR tasks.

## 1 INTRODUCTION

Image restoration (IR) aims to reconstruct high-quality images from degraded inputs. With the success of Vision Transformers (Dosovitskiy et al., 2020), Transformer-based architectures have been actively adopted for IR tasks and are now a common standard for high-performance IR backbone (Liang et al., 2021; Chen et al., 2023a; Hsu et al., 2024). However, despite recent architectural advances, the underlying training dynamics of IR Transformers remain underexplored.

This inspires us to take a closer look at their internal behavior, leading us to uncover a critical yet overlooked phenomenon: feature magnitudes diverge dramatically, reaching scales up to a *million*, while channel-wise feature entropy drops sharply (Fig.1). Interestingly, this phenomenon aligns with previous studies (Karras et al., 2020; Wang et al., 2022a), which similarly observed visual artifacts and abnormal features when coupled with specific normalization layers. However, discussions specific to the unique requirements of IR tasks and IR Transformers were not made.

Building on these insights, we hypothesize that the observed feature divergence in IR Transformers arises from networks attempting to circumvent LayerNorm (LN), due to constraints of LN that do not align with the unique requirements of IR tasks. Accordingly, we identify two key mismatches between LayerNorm and IR tasks; supported by both theoretical insights and extensive empirical analysis. First, LayerNorm operates in a per-token manner, without considering inter-pixel (token) relationships. This disrupts the spatial correlations in input features, an aspect crucial for high-fidelity image restoration. Second, it maps intermediate features into a unified normalized space, limiting the range flexibility of internal representations. This thereby disregards the input-dependent statistical variability (Lim et al., 2017b) that is inherent in IR tasks. Together, these mismatches significantly hinder IR Transformer’s ability to accurately preserve low-level features throughout the network, which is necessary for faithful image restoration. While one intuitive solution could be the complete removal of normalization layers as prior works have done (Lim et al., 2017b; Wang et al., 2018; Karras et al., 2020; 2024), our experimental observations highlight significant training instability when normalization is entirely omitted from IR Transformers (Tab.1); the network fails to converge.

In this work, we show that these issues can be addressed in a surprisingly simple manner; leading to significant stability and substantial performance gain. We propose the *Image Restoration Transformer*

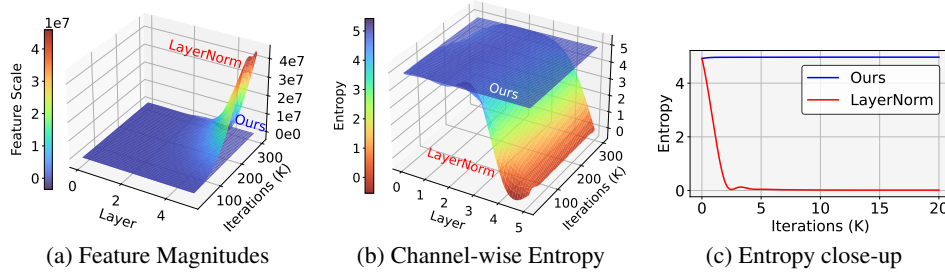


Figure 1: **Visualization of feature magnitudes and channel-wise entropy during training of an Image Restoration (IR) Transformer using conventional LayerNorm (LN) and *i*-LN (Ours).** (a) Evolution of feature magnitudes across layers and training iterations, highlighting the dramatic divergence (to million-scale) under conventional LN. (b-c) Channel-wise entropy with LN drops sharply at the very early stage of training, indicating the emergence of acute peaks hidden in specific channels. Ours *i*-LN exhibits well-distributed activation across channels and significant stability.

*Tailored Layer Normalization (i-LN)*, which acts as a drop-in replacement to conventional (vanilla) LayerNorm by better aligning with the unique requirements of IR tasks. Instead of normalizing each token independently, we propose to apply normalization across the entire spatio-channel dimension within IR Transformers (Fig.3), effectively preserving spatial correlations among tokens, contrary to vanilla per-token LayerNorm. Furthermore, We rescale features with the normalization parameters after each attention and feed-forward layer, explicitly enabling range flexibility and accounting for input-dependent variations in internal feature statistics. Together, these modifications effectively preserve low-level feature statistics throughout the network, better aligning with the requirements of IR tasks. Extensive experiments show that *i*-LN leads to both stable training dynamics with improved performance across various IR benchmarks. Additionally, we observe cues suggesting robustness under reduced-precision configurations and improved spatial correlation modeling.

## 2 METHOD

### 2.1 REVISITING LAYER NORMALIZATION

**Observation (Abnormal Feature Statistics).** Our initial analysis focuses on tracking the trajectory of internal features during the training of IR Transformers. We visualize the squared mean of intermediate features at each basic building block of the network, following (Karras et al., 2024). We select the x4 SR task using the HAT (Chen et al., 2023a) model as the representative IR task (Fig.1).

The analysis reveals that feature statistics diverge dramatically, reaching values up to a *million* scale. To pinpoint the origins of this feature divergence, we analyze the feature entropy across the channel-axis. Analysis demonstrates a sharp decrease in feature entropy, which indicates the presence of channels with extreme values that dominate the statistics. Since these extreme values are unusual, this motivates us to further investigate. Accordingly, we analyze the training dynamics across configurations by varying the network scale (Fig.2a2b), varying the IR tasks (Fig.2c), and varying the normalization scheme (Fig.4); and observe that this phenomenon occurs across all configurations utilizing standard IR Transformers. While this type of hidden abnormal behavior aligns with the observations in prior studies (Karras et al., 2020; Wang et al., 2018; 2022a), further discussion did not gain much attention, especially regarding the unique properties and requirements of IR Transformers.

In the following, we provide further insights into this phenomenon by examining the characteristics of LayerNorm (LN), the de facto normalization in IR Transformers. We start by defining the spatial relationship between pixels (i.e., inter-pixel structure), and further show that conventional LayerNorm cannot preserve this. For simplicity, we neglect the affine parameters for theoretical analysis.

**Definition 1** (Inter-pixel Structure and Preservation). *Let  $x \in \mathbb{R}^{L \times C}$  be a feature map with  $L$  tokens. We write the  $\ell$ -th token as  $x_\ell \in \mathbb{R}^C$  and the  $c$ -th element of it as  $x_{\ell,c} \in \mathbb{R}$ . The inter-pixel structure of a feature map is given by the set of relative differences  $\Delta x := \{x_\ell - x_k : 1 \leq \ell, k \leq L\}$ .*

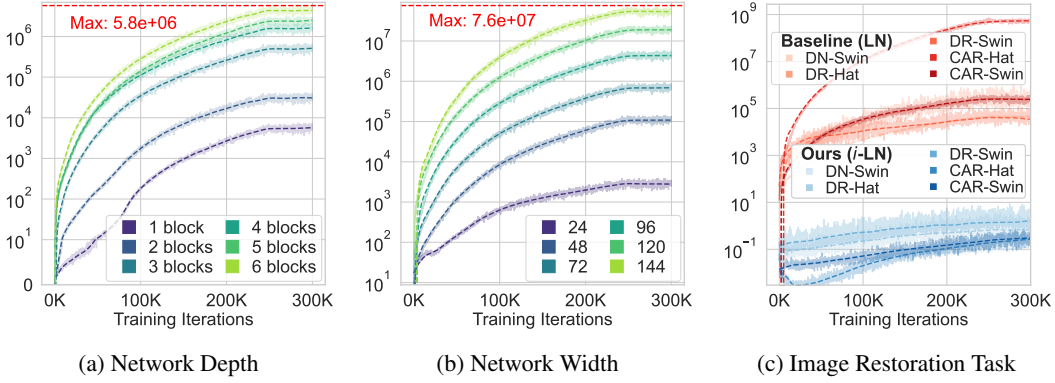


Figure 2: **Feature magnitude evolution in IR Transformers across different settings.** (a-b) Feature divergence signifies as the network scales. (c) Feature divergence appears across various transformer backbones and IR tasks: super-resolution (SR), denoising (DN), deraining (DR), JPEG compression artifact removal (CAR), demonstrating that this phenomenon is widespread. It can be effectively mitigated by simply replacing conventional LayerNorm with the proposed *i*-LN.

**Definition 2** (Structure Preserving Transformation). A transformation  $T$  is said to preserve inter-pixel structure up to scale if there exists a homothety  $H(x) = ax + b$ , with  $a > 0$  and  $b \in \mathbb{R}^C$ , such that

$$T(x_\ell) - T(x_k) = H(x_\ell - x_k) = a(x_\ell - x_k) \quad \text{for all } \ell, k.$$

Such maps preserve all angles and pairwise distance ratios, and correspond to a single global shift and uniform scaling across all tokens. For  $a = 1$ ,  $T$  is said to preserve structure absolutely.

Intuitively, consider  $x$  as a point cloud in  $\mathbb{R}^C$ , where each point represents a token. A structure-preserving transformation may only uniformly scale and shift the entire cloud. That is, the overall shape of the point cloud should be preserved up to a single global scaling factor and translation.

**Vanilla Per-token LayerNorm (Baseline).** Conventional Transformer architectures utilize the per-token LayerNorm (LN) as the de facto normalization scheme which operates as follows:

$$\text{LN}(x_\ell) = \gamma \frac{1}{\sqrt{\sigma_\ell^2 + \epsilon}} (x_\ell - \mu_\ell) + \beta, \quad \mu_\ell = \mathbb{E}_c[x_{\ell,c}], \quad \sigma_\ell^2 = \mathbb{E}_c[(x_{\ell,c} - \mu_\ell)^2], \quad (1)$$

where  $\mathbb{E}_c[\cdot]$  is taken over the channel dimension  $c$ , and  $\gamma, \beta \in \mathbb{R}^c$  are each affine parameters applied after the normalization step, and LN operates for each token  $x_\ell$  given the entire input feature  $x$ .

**Proposition 1. (Vanilla LayerNorm fails to preserve structure).** Let  $T_{\text{LN}}$  be the normalization in vanilla per-token LN. Then, in general, there do not exist  $a > 0$  and an orthogonal  $Q$  such that

$$T_{\text{LN}}(x_\ell) - T_{\text{LN}}(x_k) = aQ(x_\ell - x_k) \quad \text{for all } x_\ell, x_k,$$

Thus  $T_{\text{LN}}$  is not even conformal on the token set. Since homotheties are strict subclasses of conformal maps,  $T_{\text{LN}}$  is not a homothety and therefore it does not preserve inter-pixel structure in general.

**Remark.** The exception arises in degenerate cases where all tokens share identical per-token mean and variance, in which case a similarity map can exist (i.e., inter-pixel structure is preserved). Such cases are extremely rare in practice. Our intuition is that since LN cannot naturally preserve inter-pixel structure, networks learn to generate large magnitude features regardless of the input, thereby, manipulate the overall feature statistic to behave similarly to this exceptional degenerate scenario.

Inspired by prior observations, we hypothesize that feature divergence arises from a fundamental mismatch between the requirements of IR tasks and the constraints imposed by LayerNorm, leading us to propose a tailored normalization scheme that aligns with the unique requirements of IR tasks.

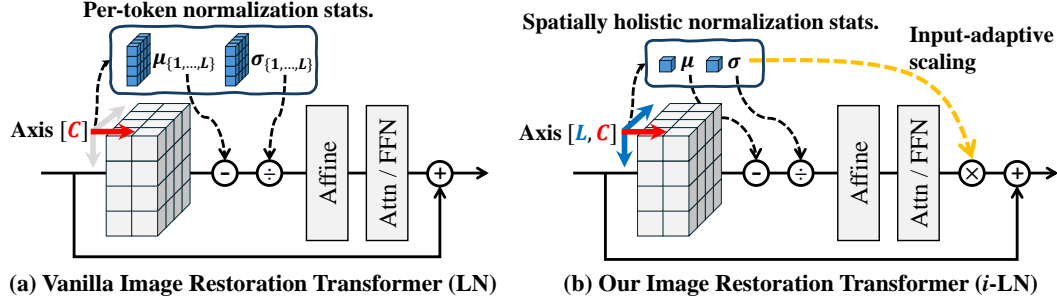


Figure 3: Comparison between IR Transformer blocks using conventional per-token LayerNorm (LN) and our proposed *i*-LN. Contrary to conventional LN, which normalizes each token independently, our *i*-LN applies holistic normalization across the entire spatio-channel dimension, preserving essential spatial correlations between tokens. Additionally, *i*-LN input-adaptively rescales features after the attention (Attn) and feedforward (FFN) layers, thereby better preserving input statistics and allowing feature range flexibility. These together enable IR Transformers to preserve low-level characteristics of input throughout the network, aligning with the unique requirements of IR.

## 2.2 TAILORING LAYER NORM FOR IMAGE RESTORATION TRANSFORMERS

**Spatially Holistic Normalization (LN\*).** We propose a simple variant of LN that improves in preserving inter-token spatial relationships of input features, which we refer to as LN\*. Instead of normalizing each token individually as LN, we derive normalization statistics from the entire spatio-channel dimension of the input feature as follows:

$$\text{LN}^*(x) = \gamma \frac{1}{\sqrt{\sigma^2 + \epsilon}}(x - \mu) + \beta, \quad \mu = \mathbb{E}_{\ell, c}[x_{\ell, c}], \quad \sigma^2 = \mathbb{E}_{\ell, c}[(x_{\ell, c} - \mu)^2], \quad (2)$$

where the expectation  $\mathbb{E}_{\ell, c}[\cdot]$  is taken over both spatial ( $\ell$ ) and channel dimensions ( $c$ ). This straightforward modification effectively mitigates the issue raised by the per-token operation in vanilla per-token LayerNorm. While normalization methods in CNNs already inherently work in a spatially holistic manner, the implications of such holistichness in normalization and the corresponding spatial structure corruption without it have received little attention, particularly in the context of IR Transformers. With this point, the following section aims to provide further intuition and establish connections between holistichness and spatial structure (i.e., inter-pixel structure) preservation.

**Proposition 2. (LN\* preserves structure).** Let  $T_{\text{LN}^*}$  be the normalization defined by LN\*, with global mean  $\mu$  and std.  $\sigma > 0$  computed over all tokens and channels. Then for any two tokens  $x_\ell, x_k$ ,

$$T_{\text{LN}^*}(x_\ell) - T_{\text{LN}^*}(x_k) = (1/\sigma)(x_\ell - x_k).$$

Thus,  $T_{\text{LN}^*}$  is a homothety, and accordingly, preserves spatial structure up to a global scale.

**Remark.** In short, LN\* is structure-preserving up to one missing scalar (i.e., the global scale). We handle this loss of information by explicitly reintroducing it later, as described below.

**Preserving Input Dependent Statistics.** We further tailor the normalization operator to better suit the requirements of IR tasks. Specifically, we address the issue of input-blind normalization. While IR tasks require the preservation of input-dependent feature statistics for faithful reconstruction, both conventional LayerNorm and even the holistic LN\* overlooks this aspect by mapping features into a unified normalized space. Although normalization improves training stability, it also causes the model to lose critical input-dependent information (i.e., the missing global scale term of inter-pixel structure) by restricting the range flexibility of internal representations (Lim et al., 2017b).

Accordingly, we propose a simple input-adaptive rescaling strategy that effectively tackles this issue. We rescale the output of Attention and FFN by their standard deviation computed in the preceding normalization process as the yellow line in Fig.3b, which we refer to as *i*-LN. Accordingly, a typical Attention or FFN block  $B$  could be further improved by coupling with *i*-LN as follows:

$$B(x; f, i\text{-LN}) = x + \sqrt{\sigma^2 + \epsilon} \cdot f(\text{LN}^*(x)), \quad (3)$$



Idx	Method	SH	Set14		BSD100		Urban100		Manga109	
			PSNR	SSIM	PSNR	SSIM	PSNR	SSIM	PSNR	SSIM
1	LayerNorm	✗	28.79	.7876	27.68	.7411	26.55	.8015	31.01	.9150
2	LayerScale	✗	28.89	.7887	27.76	.7426	26.75	.8058	31.37	.9178
3	RMSNorm	✗	28.88	.7879	27.74	.7417	26.67	.8037	31.24	.9165
4	ReZero	✓	28.81	.7861	27.70	.7406	26.41	.7964	31.05	.9147
5	None	✓	-	-	-	-	-	-	-	-
6	InstanceNorm	✓	28.98	.7907	27.80	.7445	27.02	.8136	31.46	.9199
7	BatchNorm <sup>†</sup>	✓	28.95	.7901	27.80	.7442	26.70	.8123	31.39	.9186
8	<i>i</i> -LN (Ours)	✓	<b>29.01</b>	<b>.7915</b>	<b>27.84</b>	<b>.7456</b>	<b>27.17</b>	<b>.8167</b>	<b>31.82</b>	<b>.9228</b>

Table 1: **Comparison between various normalization schemes.** <sup>†</sup> indicates that BatchNorm is evaluated in train-mode. SH indicates the spatial holistictness of the normalization scheme, including the setting without any normalization (None). Experiments are performed for  $\times 4$  SR with HAT<sub>1</sub>. The best result for each setting is highlighted in **bold**.

where  $f$  is either the according Attention or FFN operation of block  $B$ . Overall, this reintroduces the original feature statistic lost due to normalization. This simple strategy enables IR Transformers to better preserve the per-instance statistics throughout the network and allows range flexibility to intermediate features. We later show that this leads to an order of magnitude more stable feature distribution (i.e., higher entropy) and overall improved IR performance.

**Remark.** This simple input-adaptive rescaling strategy explicitly reintroduces the missing global scaling term that LN\* could not preserve (which leads to restricted range flexibility).

### 3 EXPERIMENTS

**Training Settings.** Since recent works have discrepancies in their detailed training settings (Chen et al., 2024), we reimplement baseline methods and our method under identical settings for fair comparison. Networks for deraining (DR) were trained on Rain13K (Jiang et al., 2020), while DF2K (DIV2K (Agustsson & Timofte, 2017) + Flickr2K (Lim et al., 2017a)) was used for other tasks. Only basic augmentations (random flips, rotations, crops) were applied, without mixing augmentations, progressive patch sizing, or warm-start. In order to provide thorough experimental results under various settings, the overall training budget was reduced as specified in Appendix.5. The representative SwinIR (Liang et al., 2021), HAT (Chen et al., 2023a), and DRCT (Hsu et al., 2024) were used.

**Evaluation Settings.** Standard benchmarks are employed including: Set5 (Bevilacqua et al., 2012), Set14 (Zeyde et al., 2010), BSD100 (Martin et al., 2001), Urban100 (Huang et al., 2015), Manga109 (Matsui et al., 2017) for SR; CBSD68 (Martin et al., 2001), Kodak (Franzen, 1999), McMaster (Zhang et al., 2011), Urban100 for DN; LIVE1 (Sheikh, 2005), Classic5 (Foi et al., 2007), Urban100 (Huang et al., 2015) for CAR; Test100 (Zhang et al., 2019) and Rain100L (Yang et al., 2017) for DR. We crop Urban100 into non-overlapping 256 $\times$ 256 patches due to memory limits for CAR and DN. We report PSNR and SSIM indices. Experiments were performed on NVIDIA A6000s.

#### 3.1 NORMALIZATION SCHEME VARIATION

We analyze the effects of various normalization techniques, including representative normalization schemes as vanilla LayerNorm (LN) (Ba et al., 2016), per-token RMSNorm (RMS) (Zhang & Sennrich, 2019), InstanceNorm (IN) (Ulyanov et al., 2016), BatchNorm (BN) (Ioffe & Szegedy, 2015), and our proposed *i*-LN. Considering previous studies where completely removing normalization from SR networks (Lim et al., 2017b; Wang et al., 2018) led to performance improvements, we additionally tested a similar configuration indicated as *None*, where normalizations are entirely removed.

Further, we investigate the empirical impacts of recent methods designed to stabilize Transformer training: ReZero (RZ) (Bachlechner et al., 2021) and LayerScale (LS) (Touvron et al., 2021). ReZero removes LayerNorm from Transformer blocks and multiplies a learnable zero-initialized scalar to the residual path. Similarly, LayerScale multiplies a near-zero-initialized learnable diagonal matrix to the residual path but reintroduces LayerNorm. Since both methods initially multiply a (near) zero-scale factor to the network output, we consider them as potential solutions to resolve the feature increasing issue in IR tasks. Notably, these methods also align with prior stud-

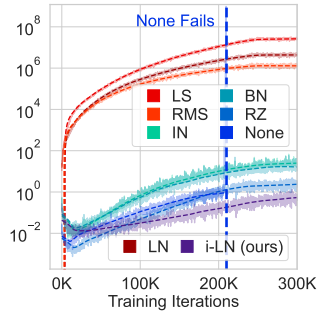


Figure 4: Feature divergence across various normalizations.

Table 2: Quantitative comparison between the conventional LayerNorm (LN) and our proposed *i*-LN across diverse IR tasks. The best result for each setting is highlighted in **bold**.

Backbone	Scale	Set5		Set14		BSD100		Urban100		Manga109		Backbone	Testset	Metric	
		PSNR	SSIM	PSNR	SSIM	PSNR	SSIM	PSNR	SSIM	PSNR	SSIM			PSNR	SSIM
HAT <sub>1</sub> + LN	×2	38.14	.9610	33.78	.9196	32.19	.9000	32.16	.9297	38.84	.9778	HAT <sub>1</sub> + LN	Rain100L	34.35	.9471
HAT <sub>1</sub> + <i>i</i> -LN	×2	<b>38.37</b>	<b>.9619</b>	<b>34.08</b>	<b>.9218</b>	<b>32.42</b>	<b>.9028</b>	<b>33.32</b>	<b>.9385</b>	<b>39.69</b>	<b>.9794</b>	HAT <sub>1</sub> + <i>i</i> -LN		<b>36.20</b>	<b>.9641</b>
DRCT <sub>1</sub> + LN	×2	38.19	.9613	33.28	.9197	32.28	.9010	32.60	.9323	39.23	.9785	SwinIR <sub>1</sub> + LN	Rain100L	33.00	.9434
DRCT <sub>1</sub> + <i>i</i> -LN	×2	<b>38.23</b>	<b>.9614</b>	<b>33.86</b>	<b>.9206</b>	<b>32.31</b>	<b>.9014</b>	<b>32.79</b>	<b>.9344</b>	<b>39.40</b>	<b>.9788</b>	SwinIR <sub>1</sub> + <i>i</i> -LN		<b>34.43</b>	<b>.9527</b>
HAT <sub>1</sub> + LN	×4	32.51	.8992	28.79	.7876	27.68	.7411	26.55	.8015	31.01	.9150	HAT <sub>1</sub> + LN	Test100	29.52	.8905
HAT <sub>1</sub> + <i>i</i> -LN	×4	<b>32.72</b>	<b>.9019</b>	<b>29.01</b>	<b>.7915</b>	<b>27.84</b>	<b>.7456</b>	<b>27.17</b>	<b>.8167</b>	<b>31.82</b>	<b>.9228</b>	HAT <sub>1</sub> + <i>i</i> -LN		<b>30.14</b>	<b>.9022</b>
DRCT <sub>1</sub> + LN	×4	32.50	.8989	28.85	.7871	27.73	.7414	26.63	.8021	31.24	.9169	SwinIR <sub>1</sub> + LN	Test100	27.45	.8766
DRCT <sub>1</sub> + <i>i</i> -LN	×4	<b>32.57</b>	<b>.8997</b>	<b>28.91</b>	<b>.7887</b>	<b>27.76</b>	<b>.7426</b>	<b>26.79</b>	<b>.8063</b>	<b>31.41</b>	<b>.9188</b>	SwinIR <sub>1</sub> + <i>i</i> -LN		<b>29.87</b>	<b>.8982</b>

(a) Single image super-resolution (SR)

(b) Image deraining (DR)

Backbone	$\sigma$	Urban100	CBSD68	Kodak24	McMaster	Backbone	q	Urban100	LIVE1	Classic5			
		PSNR	PSNR	PSNR	PSNR			PSNR	SSIM	PSNR	SSIM		
HAT <sub>1</sub> + LN	15	35.489	34.285	35.347	35.440	HAT <sub>1</sub> + LN	10	28.45	.8514	27.89	.8048	29.94	.8167
HAT <sub>1</sub> + <i>i</i> -LN	15	<b>35.558</b>	<b>34.296</b>	<b>35.366</b>	<b>35.477</b>	HAT <sub>1</sub> + <i>i</i> -LN	10	<b>28.52</b>	<b>.8530</b>	<b>27.90</b>	<b>.8057</b>	<b>29.96</b>	<b>.8178</b>
SwinIR <sub>1</sub> + LN	15	35.077	34.164	35.147	35.183	SwinIR <sub>1</sub> + LN	10	27.86	.8400	<b>27.65</b>	<b>.7995</b>	<b>29.72</b>	<b>.8111</b>
SwinIR <sub>1</sub> + <i>i</i> -LN	15	<b>35.138</b>	<b>34.181</b>	<b>35.177</b>	<b>35.223</b>	SwinIR <sub>1</sub> + <i>i</i> -LN	10	<b>27.92</b>	<b>.8410</b>	27.62	.7993	<b>29.72</b>	<b>.8111</b>
HAT <sub>1</sub> + LN	25	33.296	31.622	32.864	33.105	HAT <sub>1</sub> + LN	40	33.26	.9302	32.63	.9158	34.34	.9060
HAT <sub>1</sub> + <i>i</i> -LN	25	<b>33.384</b>	<b>31.632</b>	<b>32.887</b>	<b>33.139</b>	HAT <sub>1</sub> + <i>i</i> -LN	40	<b>33.36</b>	<b>.9312</b>	<b>32.67</b>	<b>.9162</b>	<b>34.39</b>	<b>.9066</b>
SwinIR <sub>1</sub> + LN	25	32.753	31.480	32.643	32.829	SwinIR <sub>1</sub> + LN	40	32.62	.9245	32.34	.9127	34.11	.9036
SwinIR <sub>1</sub> + <i>i</i> -LN	25	<b>32.803</b>	<b>31.489</b>	<b>32.660</b>	<b>32.848</b>	SwinIR <sub>1</sub> + <i>i</i> -LN	40	<b>32.68</b>	<b>.9252</b>	<b>32.35</b>	<b>.9129</b>	<b>34.12</b>	<b>.9038</b>

(c) Color image denoising (DN)

(d) Image JPEG compression artifact removal (CAR)

ies (Lim et al., 2017b; Wang et al., 2018), where multiplying a small scale factor to the residual path components helped the network to converge. Overall, this study aims to explore 1) the feature divergence tendency of per-token and holistic normalizations and 2) determine which normalization method yields the best performance.

**Feature Divergence Behavior.** Fig.4 illustrates that feature divergence always emerges when using per-token normalizations: vanilla LN, RMSNorm, and LayerScale. In contrast, spatially consistent normalizations as our *i*-LN or BN, IN, ReZero do not exhibit the divergence trend. For the configuration without any normalization, we observe failure to converge due to unstable training. However, the feature magnitudes are well-bounded before this failure occurs, aligning with other normalization schemes without the per-token operation. This observation also aligns with our hypothesis that the feature divergence phenomena is closely related to the per-token normalization, and also reveals that any spatially consistent normalization could potentially reduce this effect.

**Performance Comparisons.** We further analyze the empirical performance for each normalization scheme in Tab.1. Conventional LN performs the worst since it neglects inter-token spatial relationships and maps features into a unified normalized space, disregarding the input-dependent feature statistics.

LayerScale and RMSNorm show improvement against vanilla LN, but perform worse than methods with spatially consistent normalization. Meanwhile, without any normalization (None), the network fails to converge potentially due to unstable gradients raised by the absence of normalization, similar to prior studies in RZ (Bachlechner et al., 2021). BN leads to a significant performance drop in eval-mode, despite being healthy in train-mode; consistent with prior studies (Lim et al., 2017b; Wang et al., 2022a). This signifies the necessity of per-image statistics within the normalization scheme for IR tasks. IN performs better than vanilla LN but worse than ours. Both IN and BN discard crucial channel-wise information necessary for representing deep features, resulting in limited performance. However, despite these limitations in current spatially holistic normalization schemes (IN, BN), they already outperform those with per-token schemes (LN, LS, RMS). Meanwhile, our *i*-LN achieves the best performance among all examined methods, demonstrating its effectiveness in preserving important inter-token spatial relationships and internal statistics, and ultimately the input low-level features throughout the network.

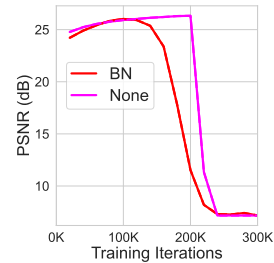


Figure 5: Eval-mode BN and removing all normalization (None) fails.

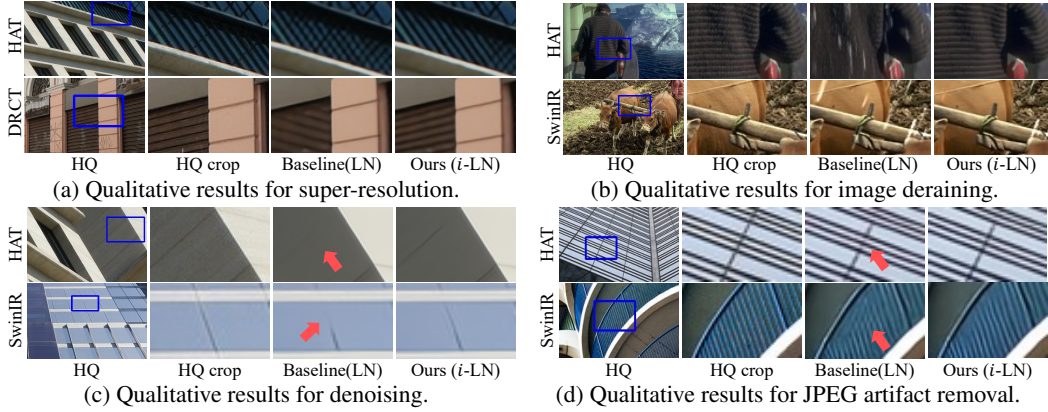


Figure 6: Qualitative comparison across four representative image restoration tasks.

### 3.2 ANALYSIS UNDER TASK VARIATION

**Feature Divergence Behavior.** Fig. 2c illustrates the evolution of feature magnitudes across various Image Restoration (IR) tasks, including Image Super-Resolution (SR), Image Denoising (DN), Image Deraining (DR), and JPEG Compression Artifact Removal (CAR). The figure clearly demonstrates that feature divergence consistently occurs across all restoration tasks under conventional LayerNorm. In contrast, integrating our *i*-LN effectively resolves this issue, maintaining stable and well-bounded feature scales throughout training. This consistent stabilization of internal feature magnitudes confirms the general applicability and robustness of our proposed method across diverse IR scenarios.

**Benchmark: Image Super-Resolution (SR).** Tab.2a and Fig.6a illustrate quantitative and qualitative results for SR. Compared to vanilla LayerNorm, we achieve significant improvements across benchmarks. Notably, SR benefits greatly from our method due to the inherent nature of SR: the input is entirely reliable, since it exactly aligns with the low-frequency information in the ground truth. By precisely preserving these input features, our method substantially enhances restored image quality. We additionally provide a comparison against the official public models under computationally extensive settings in Appendix.B.1.

**Benchmark: Image Deraining (DR).** Similarly, Tab. 2b and Fig. 6b demonstrate substantial improvements of our method in image deraining compared to conventional LayerNorm. This improvement is particularly pronounced because our method effectively preserves reliable input regions, specifically the local areas unaffected by rain streaks. By explicitly maintaining these local correspondences with the ground truth, our *i*-LN method achieves improved restoration accuracy.

**Benchmark: Image Denoising (DN)** Tab. 2c and Fig. 6c demonstrate that our method consistently outperforms conventional LayerNorm in image denoising tasks. However, the observed performance improvements are smaller compared to SR and Deraining. This relatively reduced benefit arises because denoising involves uniformly distributed corruptions across the entire image, limiting the advantage gained from explicitly preserving particular input features. Despite this, visual examples confirm meaningful improvements in recovering sharp edges.

**Benchmark: JPEG compression artifact removal (CAR).** Similarly, Tab. 2d and Fig. 6d demonstrate consistent improvements of our method over LayerNorm for JPEG compression artifact removal. However, these performance gains remain smaller than those achieved in SR and Deraining. Similar to denoising, JPEG artifacts affect images globally and irregularly, limiting the advantage of explicitly preserving specific input details. Still, visual examples illustrate consistent improvement in accurate artifact reductions, highlighting our method’s broad effectiveness across various IR tasks.

**Real-world Degradation Scenarios.** To further validate the effectiveness and robustness of the proposed method, we conduct further experiments under the challenging real-world degradation configurations. Here, we choose the representative Real-ESRGAN Wang et al. (2021) degradation pipeline and synthesize both the train and test images accordingly. Experiments are performed under the  $\times 4$  SR task with the HAT<sub>1</sub> model. In Fig. 9 and Tab. 15, we provide qualitative and quantitative results, respectively. As demonstrated, our *i*-LN shows significant improvements even under the complex real-world degradation settings, successfully reconstructing fine-details and sharp edges.

### 3.3 ABLATION STUDY

To analyze the contribution of each component in  $i$ -LN, we conduct an ablation study by selectively removing spatial holistness and rescaling. Compared to Tab.2, we increase the network capacity and training iterations (denoted as  $HAT_2$ ) to ensure that the observed benefits are not simply due to faster convergence. In Tab. 3, Fig.13 and Fig.12, we provide a quantitative and qualitative analysis results under the  $\times 4$  SR task. Removing either the rescaling strategy (Rs) or the spatial holistness (SH) consistently reduces restoration quality, confirming their complementary roles in improving IR performance.

We then examine feature statistics in Fig. 7. Starting from our full method, we remove the rescaling method (falling back to  $LN^*$ ) and subsequently the spatial holistic scheme (falling back to vanilla per-token  $LN$ ). Here, we observe that channel entropy collapses *exponentially*, indicating that each component contributes to maintaining well-distributed activations across channels.

Overall, using both components together achieves the best results in terms of both restoration quality and stable feature statistics. Spatial holistness ( $LN^*$ ) effectively preserves inter-token relationships, while the rescaling strategy further restores the missing global scale that  $LN^*$  alone cannot maintain.

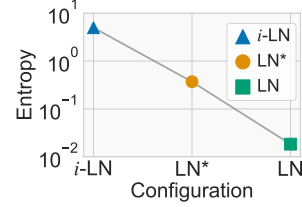


Figure 7: Channel entropy collapses *exponentially* as we remove each component (spatially holistness and rescaling) of our  $i$ -LN; falling back to vanilla LN.

### 3.4 INTRIGUING PROPERTIES OF FEATURE DIVERGENCE UNDER VANILLA LN

#### 3.4.1 HOW NETWORK SCALE IMPACTS FEATURE DIVERGENCE

We further investigate how the overall network size affects feature divergence by varying the depth and width of the IR Transformer individually. As shown in Fig. 2a–2b, larger models consistently diverge faster and to higher magnitudes. In particular, the emergence of an extreme valued feature appears to be a cumulative process: in order for a newly generated outlier channel to dominate the statistics, it must surpass the already abnormal activations propagated through the residual path, resulting in increasingly extreme values as the network scales. Taken together, our analysis reveals a potential vulnerability unique to low-level restoration at scale, where enlarging capacity does not merely amplify representational power but also exacerbates pathological feature growth.

#### 3.4.2 CHANNEL IMBALANCE AND BIAS ALIGNMENT

Earlier, we observed extreme feature norms and imbalances in channel entropy, indicating highly peaky feature distributions concentrated in specific channels. Interestingly, despite these severe imbalances, baseline IR Transformers manage to converge and produce outputs with well-bounded magnitudes. To gain insights to this paradox, we take a closer look at the final normalization layer (LN). In Fig. 8, we visualize the unnormalized feature magnitudes along the channel dimension before the final normalization layer and compare them with the learned affine bias parameter ( $\gamma$ ) across various IR tasks.

We observe sharp peaks in the bias parameters precisely aligning with the channels exhibiting high magnitudes. This exact alignment reveals a compensatory mechanism where the learnable affine terms ( $\gamma, \beta$ ) of LayerNorm counteract abnormal channel activations, allowing baselines to yield normal images. Additionally, this also indicates that the normalization operation ( $\mu, \sigma$ ) itself is incapable of directly removing these extreme peaks.

Moreover, although the observed bias–feature alignment allows baseline IR Transformers to maintain reasonable outputs, this mechanism should be regarded as a compensatory shortcut rather than a fundamental fix. The fact that networks must rely on such peaky biases to counteract extreme channel activations leaves the model fragile and prone to failures, including potential training instability and failure in practical scenarios such as reduced-precision inference as discussed in Sec.3.5.2.

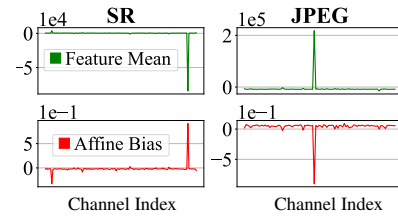


Figure 8: Alignment of affine bias parameters in the last LN and channel-wise magnitude of input feature; showing a compensatory mechanism.



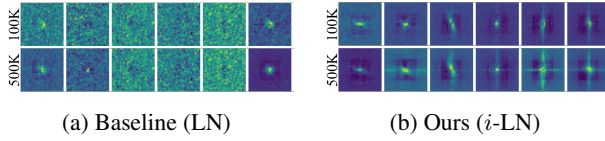


Figure 9: Visualization of Relative Position Embeddings (RPE) per head, for training iteration 100K and 500K. Ours exhibit well-structured RPEs, indicating the superiority in understanding the spatial relationship between pixels.

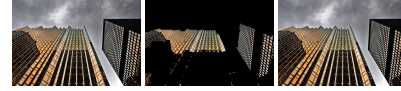


Figure 10: Half-precision inference results for  $\times 4$  SR. LN leads to artifacts while our *i*-LN achieves near-zero fidelity loss compared to full-precision.

Idx	Backbone	SH	Rs	BSD100	Urban100	Manga109
1	HAT <sub>2</sub> (LN)			27.7897	26.8779	31.5444
2	HAT <sub>2</sub>		✓	27.8615	27.3373	31.8888
3	HAT <sub>2</sub>	✓		27.9034	27.5335	32.0837
4	HAT <sub>2</sub> ( <i>i</i> -LN)	✓	✓	<b>27.9206</b>	<b>27.5849</b>	<b>32.1694</b>

Table 3: Ablation study. *SH* indicates introducing spatial holistiness (identical to LN\*) and *Rs* indicates our rescaling strategy. Idx 1 and 4 are each identical to vanilla LN and our *i*-LN, respectively. Experiments conducted for  $\times 4$  SR.

Idx	Backbone	Quantization	Urban100	Manga109
1	HAT <sub>2</sub> + LN	W int8	26.8711	31.5266
2	HAT <sub>2</sub> + <i>i</i> -LN	W int8	27.5818	32.1657
3	HAT <sub>2</sub> + LN	W int4	25.0242	28.0831
4	HAT <sub>2</sub> + <i>i</i> -LN	W int4	26.8292	30.6596
5	HAT <sub>2</sub> + LN	W+F fp16	7.4640	5.0736
6	HAT <sub>2</sub> + <i>i</i> -LN	W+F fp16	27.5849	32.1693

Table 4: Quantitative results under low-precision inference. *W* indicates weight-only quantization, *W+F* indicates weight and feature quantization.

### 3.5 INTRIGUING PROPERTIES UNDER *i*-LN

#### 3.5.1 ENHANCED SPATIAL CORRELATION VIA STRUCTURED RPE

Relative Position Embeddings (RPE) explicitly encodes relative spatial positions between tokens in an input-agnostic manner, similar to the convolution operation that inherently captures the spatial locality through their structured kernel patterns. Accordingly, we can consider well-structured RPEs as a strong indicator of enhanced spatial correlation understanding. In Fig.9, we analyze how our proposed normalization method influences spatial relationship modeling by visualizing the learned RPE of both the baseline IR Transformer and our proposed method. The baseline Transformer exhibits noisy, unstructured embedding patterns, suggesting a limited capability to effectively model spatial correlations. Conversely, our method produces RPEs that resemble well-structured convolutional filter patterns, clearly indicating superior capture of spatial relationships. This structured embedding aligns with our hypothesis that our spatially holistic normalization better preserves intrinsic spatial correlations, helping the network to learn spatial relations more effectively. In Fig. 12, we provide further visual examples of RPEs between LN, *i*-LN and also the ablated variants LN\* and LN+Rescaling, which shows aligning results with the discussion above.

#### 3.5.2 LOW PRECISION INFERENCE

Image restoration networks often require deployment on lightweight edge devices, creating significant demand for efficient inference in IR Transformers. A common approach to enhance inference efficiency is reducing precision during model deployment. Consequently, we conducted experiments under reduced-precision inference conditions to empirically evaluate the effects of *i*-LN. Initially, we applied linear weight quantization to the model weights. As shown in Tab.4, vanilla LayerNorm resulted in substantial performance degradation, while *i*-LN demonstrated remarkable stability. We further conducted half-precision inference experiments, casting both internal feature values and weights to half-precision floating-point numbers. Fig.10 illustrates that vanilla LayerNorm generated extensive regions of black dots, indicating network-generated infinity values due to extreme internal feature magnitudes inadequate for low-precision conditions. Notably, no substantial performance degradation was observed in regions where the network maintained finite feature values. This highlights the necessity of well-bounded feature values achieved by *i*-LN, emphasizing its critical role in enabling efficient inference for IR Transformers.

#### 3.5.3 EMPIRICAL EFFECTS OF IMPROVED STABILITY

**Multiple Runs.** To further probe training stability, we conduct extensive multi-run experiments across diverse random seeds (Fig.11). Here, we use a small batch size of 2 to induce training instability. Vanilla LN exhibits inconsistent optimization trajectories, with large fluctuations in feature statistic evolution patterns and substantial discrepancy in final PSNR results for each run. In contrast, our *i*-LN produces significantly lower variance between multiple runs, in terms of both training statistics and

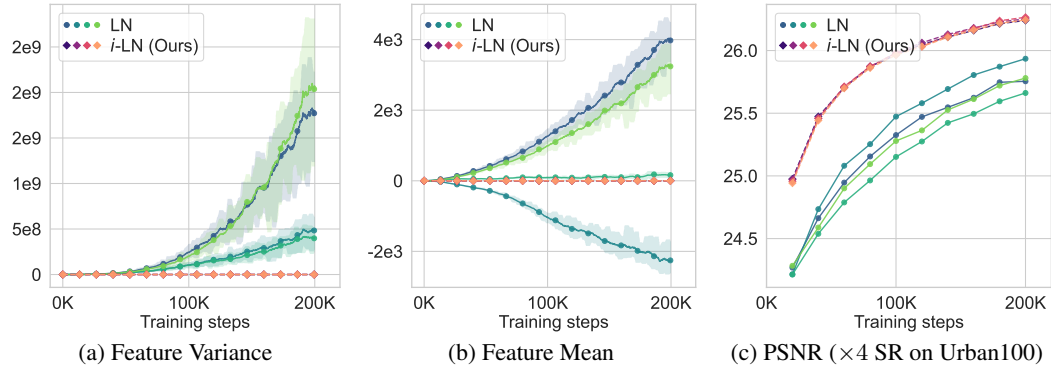


Figure 11: **Feature statistics and PSNR across multiple runs with different random seeds.** Trajectories show significant fluctuation under vanilla LN, while our *i*-LN maintains well-bounded and consistent results across all seeds ( $\times 4$  SR, Urban100).

the final reconstruction performance. These results demonstrate that *i*-LN provides a more reliable optimization landscape, reducing susceptibility to randomness in initialization or data ordering, which is an important practical advantage for training IR networks.

## 4 RELATED WORK

**Image Restoration Transformers.** Recent advances in Image Restoration (IR) transformers (Chen et al., 2024; Zamir et al., 2022; Wang et al., 2022c; Zhang et al., 2022) show superior performance over CNNs (Dong et al., 2015; Kim et al., 2016; Zhang et al., 2018) by leveraging attention mechanisms to effectively model long-range context. A pioneering work, SwinIR (Liang et al., 2021), adopted an efficient Swin-Transformer (Liu et al., 2021) based architecture in IR tasks, balancing computational cost and restoration quality. A notable method is HAT, which originated as a super-resolution model (Chen et al., 2023b) but expanded to general image restoration tasks (Chen et al., 2023a). By unifying spatial and channel attention within a hybrid attention framework, HAT surpasses existing IR Transformers in both restoration fidelity and robustness across various IR tasks.

**Abnormal Feature Behaviors.** Normalization is a key element in enhancing stability and performance in deep networks, but also can lead to unintended feature behavior. EDSR (Lim et al., 2017b), which is a foundational work in super-resolution pointed out that BatchNorm removes range flexibility of intermediate features, leading to a performance drop. Accordingly, normalization layers are removed in the most recent CNN-based SR architectures. Meanwhile, ESRGAN (Wang et al., 2018) and StyleGAN2 (Karras et al., 2020) observe that InstanceNorm and BatchNorm, respectively, cause water droplet-like artifacts. They suggest that the generator might learn to deceive feature statistics by sneaking abnormal values in internal features to reduce the effects of normalization. EDM2 (Karras et al., 2024) identifies feature magnitude divergence in diffusion models. Accordingly, they redesign the network architecture to preserve the magnitude based on statistical assumptions, leading to overall performance enhancement. DRCT (Hsu et al., 2024) notes that feature map intensities drop sharply at the end of SR networks, leading to information bottlenecks, and shows that dense residuals help.

## 5 CONCLUSION

We analyzed the training dynamics of Image Restoration (IR) Transformers and highlighted an overlooked phenomenon: divergence of feature magnitudes accompanied by collapses in channel-wise entropy. We interpret this as networks attempting to bypass the constraints of conventional LayerNorm, whose per-token normalization and input-independent scaling disrupt spatial correlations and restrict the flexibility needed for accurate restoration. To address this, we introduced Image Restoration Transformer Tailored Layer Normalization (*i*-LN), a simple drop-in replacement for LayerNorm. It is designed to better align with the unique characteristics of IR tasks and preserve important low-level features of the input throughout the network. *i*-LN normalizes jointly across spatial and channel dimensions and incorporates input-dependent rescaling, aligning normalization more closely with the demands of IR tasks. Extensive experiments show that *i*-LN prevents feature divergence, stabilizes channel entropy, improves robustness under low-precision inference, and significantly enhances IR performance across diverse tasks.



## 6 REPRODUCIBILITY STATEMENT

Experimental settings for both training and evaluation are described in Sec.3. Detailed hyperparameter settings and network configurations for each model variant are described in Appendix.B.1 and Tab.5. Detailed algorithm to calculate the channel entropy is in Appendix.A.2 We plan to release the code for further reproducibility.

## REFERENCES

- Eirikur Agustsson and Radu Timofte. Ntire 2017 challenge on single image super-resolution: Dataset and study. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops*, pp. 126–135, 2017.
- Jimmy Lei Ba, Jamie Ryan Kiros, and Geoffrey E Hinton. Layer normalization. *arXiv preprint arXiv:1607.06450*, 2016.
- Thomas Bachlechner, Bodhisattwa Prasad Majumder, Henry Mao, Gary Cottrell, and Julian McAuley. Rezero is all you need: Fast convergence at large depth. In *Uncertainty in Artificial Intelligence*, pp. 1352–1361. PMLR, 2021.
- Marco Bevilacqua, Aline Roumy, Christine Guillemot, and Marie Line Alberi-Morel. Low-complexity single-image super-resolution based on nonnegative neighbor embedding. 2012.
- Yochai Blau and Tomer Michaeli. Rethinking lossy compression: The rate-distortion-perception tradeoff. In *International Conference on Machine Learning*, pp. 675–685. PMLR, 2019.
- Rob Brekelmans, Daniel Moyer, Aram Galstyan, and Greg Ver Steeg. Exact rate-distortion in autoencoders via echo noise. *Advances in neural information processing systems*, 32, 2019.
- Xiangyu Chen, Xintao Wang, Wenlong Zhang, Xiangtao Kong, Yu Qiao, Jiantao Zhou, and Chao Dong. Hat: Hybrid attention transformer for image restoration. *arXiv preprint arXiv:2309.05239*, 2023a.
- Xiangyu Chen, Xintao Wang, Jiantao Zhou, Yu Qiao, and Chao Dong. Activating more pixels in image super-resolution transformer. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 22367–22377, 2023b.
- Xiangyu Chen, Zheyuan Li, Yuandong Pu, Yihao Liu, Jiantao Zhou, Yu Qiao, and Chao Dong. A comparative study of image restoration networks for general backbone network design. In *European Conference on Computer Vision*, pp. 74–91. Springer, 2024.
- Chao Dong, Chen Change Loy, Kaiming He, and Xiaoou Tang. Image super-resolution using deep convolutional networks. *IEEE transactions on pattern analysis and machine intelligence*, 38(2): 295–307, 2015.
- Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*, 2020.
- Alessandro Foi, Vladimir Katkovnik, and Karen Egiazarian. Pointwise shape-adaptive dct for high-quality denoising and deblocking of grayscale and color images. *IEEE transactions on image processing*, 16(5):1395–1411, 2007.
- R. Franzen. Kodak lossless true color image suite. <http://r0k.us/graphics/kodak>, 1999. Volume 4, no. 2.
- Irina Higgins, Loic Matthey, Arka Pal, Christopher Burgess, Xavier Glorot, Matthew Botvinick, Shakir Mohamed, and Alexander Lerchner. beta-vae: Learning basic visual concepts with a constrained variational framework. In *International conference on learning representations*, 2017.
- Chih-Chung Hsu, Chia-Ming Lee, and Yi-Shiuan Chou. Drct: Saving image super-resolution away from information bottleneck. *arXiv preprint arXiv:2404.00722*, 2024.

- Jia-Bin Huang, Abhishek Singh, and Narendra Ahuja. Single image super-resolution from transformed self-exemplars. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 5197–5206, 2015.
- Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. In *International conference on machine learning*, pp. 448–456. pmlr, 2015.
- Kui Jiang, Zhongyuan Wang, Peng Yi, Chen Chen, Baojin Huang, Yimin Luo, Jiayi Ma, and Junjun Jiang. Multi-scale progressive fusion network for single image deraining. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 8346–8355, 2020.
- Tero Karras, Samuli Laine, Miika Aittala, Janne Hellsten, Jaakko Lehtinen, and Timo Aila. Analyzing and improving the image quality of stylegan. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 8110–8119, 2020.
- Tero Karras, Miika Aittala, Jaakko Lehtinen, Janne Hellsten, Timo Aila, and Samuli Laine. Analyzing and improving the training dynamics of diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 24174–24184, 2024.
- Jiwon Kim, Jung Kwon Lee, and Kyoung Mu Lee. Accurate image super-resolution using very deep convolutional networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 1646–1654, 2016.
- Jingyun Liang, Jiezhang Cao, Guolei Sun, Kai Zhang, Luc Van Gool, and Radu Timofte. Swinir: Image restoration using swin transformer. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 1833–1844, 2021.
- Bee Lim, Sanghyun Son, Heewon Kim, Seungjun Nah, and Kyoung Mu Lee. Enhanced deep residual networks for single image super-resolution. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR) Workshops*, July 2017a.
- Bee Lim, Sanghyun Son, Heewon Kim, Seungjun Nah, and Kyoung Mu Lee. Enhanced deep residual networks for single image super-resolution. In *Proceedings of the IEEE conference on computer vision and pattern recognition workshops*, pp. 136–144, 2017b.
- Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining Guo. Swin transformer: Hierarchical vision transformer using shifted windows. In *Proceedings of the IEEE/CVF international conference on computer vision*, pp. 10012–10022, 2021.
- David Martin, Charless Fowlkes, Doron Tal, and Jitendra Malik. A database of human segmented natural images and its application to evaluating segmentation algorithms and measuring ecological statistics. In *Proceedings Eighth IEEE International Conference on Computer Vision. ICCV 2001*, volume 2, pp. 416–423. IEEE, 2001.
- Yusuke Matsui, Kota Ito, Yuji Aramaki, Azuma Fujimoto, Toru Ogawa, Toshihiko Yamasaki, and Kiyoharu Aizawa. Sketch-based manga retrieval using manga109 dataset. *Multimedia Tools and Applications*, 76:21811–21838, 2017.
- H Sheikh. Live image quality assessment database release 2. <http://live.ece.utexas.edu/research/quality>, 2005.
- Hugo Touvron, Matthieu Cord, Alexandre Sablayrolles, Gabriel Synnaeve, and Hervé Jégou. Going deeper with image transformers. In *Proceedings of the IEEE/CVF international conference on computer vision*, pp. 32–42, 2021.
- Dmitry Ulyanov, Andrea Vedaldi, and Victor Lempitsky. Instance normalization: The missing ingredient for fast stylization. *arXiv preprint arXiv:1607.08022*, 2016.
- Xintao Wang, Ke Yu, Shixiang Wu, Jinjin Gu, Yihao Liu, Chao Dong, Yu Qiao, and Chen Change Loy. Esrgan: Enhanced super-resolution generative adversarial networks. In *Proceedings of the European conference on computer vision (ECCV) workshops*, pp. 0–0, 2018.

- Xintao Wang, Liangbin Xie, Chao Dong, and Ying Shan. Real-esrgan: Training real-world blind super-resolution with pure synthetic data. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 1905–1914, 2021.
- Xintao Wang, Chao Dong, and Ying Shan. Repsr: Training efficient vgg-style super-resolution networks with structural re-parameterization and batch normalization. In *Proceedings of the 30th acm international conference on multimedia*, pp. 2556–2564, 2022a.
- Xintao Wang, Liangbin Xie, Ke Yu, Kelvin C.K. Chan, Chen Change Loy, and Chao Dong. BasicSR: Open source image and video restoration toolbox. <https://github.com/XPixelGroup/BasicSR>, 2022b.
- Zhendong Wang, Xiaodong Cun, Jianmin Bao, Wengang Zhou, Jianzhuang Liu, and Houqiang Li. Uformer: A general u-shaped transformer for image restoration. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 17683–17693, 2022c.
- Wenhan Yang, Robby T Tan, Jiashi Feng, Jiaying Liu, Zongming Guo, and Shuicheng Yan. Deep joint rain detection and removal from a single image. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 1357–1366, 2017.
- Jingfeng Yao, Bin Yang, and Xinggang Wang. Reconstruction vs. generation: Taming optimization dilemma in latent diffusion models. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pp. 15703–15712, 2025.
- Syed Waqas Zamir, Aditya Arora, Salman Khan, Munawar Hayat, Fahad Shahbaz Khan, and Ming-Hsuan Yang. 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.
- Roman Zeyde, Michael Elad, and Matan Protter. On single image scale-up using sparse-representations. In *International conference on curves and surfaces*, pp. 711–730. Springer, 2010.
- Biao Zhang and Rico Sennrich. Root mean square layer normalization. *Advances in Neural Information Processing Systems*, 32, 2019.
- Dafeng Zhang, Feiyu Huang, Shizhuo Liu, Xiaobing Wang, and Zhezhu Jin. Swinfr: Revisiting the swinir with fast fourier convolution and improved training for image super-resolution. *arXiv preprint arXiv:2208.11247*, 2022.
- He Zhang, Vishwanath Sindagi, and Vishal M Patel. Image de-raining using a conditional generative adversarial network. *IEEE transactions on circuits and systems for video technology*, 30(11): 3943–3956, 2019.
- Lei Zhang, Xiaolin Wu, Antoni Buades, and Xin Li. Color demosaicking by local directional interpolation and nonlocal adaptive thresholding. *Journal of Electronic imaging*, 20(2):023016–023016, 2011.
- Yulun Zhang, Kunpeng Li, Kai Li, Lichen Wang, Bineng Zhong, and Yun Fu. Image super-resolution using very deep residual channel attention networks. In *Proceedings of the European conference on computer vision (ECCV)*, pp. 286–301, 2018.
- Yupeng Zhou, Zhen Li, Chun-Le Guo, Song Bai, Ming-Ming Cheng, and Qibin Hou. Srformer: Permuted self-attention for single image super-resolution. In *Proceedings of the IEEE/CVF international conference on computer vision*, pp. 12780–12791, 2023.

Table 5: Overview of model variants. ► indicates the group (type) of each model variant: from lightweight to computationally extensive settings. The according placements of the experiments are specified as *Main* (i.e., the main article) and *Appendix*. The placements of the detailed network hyperparameters and training configurations for each variant are highlighted in **bold**.

Variant	Description
<b>► Type1 (Lightweight Computational Configuration)</b>	
These configurations are used for most analyses. All models were trained from scratch without the Warm-start strategy (i.e., the $\times 4$ SR models are not finetuned from the $\times 2$ SR weights), Mixing Augmentations, Progressive Patch Sizing.	
<b>SwinIR<sub>1</sub> - Main</b> . . . . .	Details are provided in <b>Tab. 10</b> . This variant shares the same network architecture as the official SwinIR-light model implementation.
<b>DRCT<sub>1</sub> - Main</b> . . . . .	Details are provided in <b>Tab. 11</b> . This variant is a lightweight variant of the DRCT Hsu et al. (2024) model implementation. The embedding dimension is reduced in order align the network architecture with the HAT <sub>1</sub> model.
<b>HAT<sub>1</sub> - Main &amp; Appendix</b> . . . .	Details are provided in <b>Tab. 13</b> . This is a variant is a lighter version of the HAT-S Chen et al. (2023b) model, modified with a slightly reduced embedding dimension. This change was made since the standard HAT-S, despite being denoted as <i>small</i> , requires more Mult-Adds than the full-sized SwinIR model.
<b>SRFormer<sub>1</sub> - Appendix</b> . . . . .	Details are provided in <b>Tab. 12</b> . This variant is a lightweight variant of the SRFormer Zhou et al. (2023) model. The overall capacity is reduced to align with the networks specified above.
<b>► Type2 (Moderate Computational Configuration)</b>	
These configurations are used for the ablation study and low-precision inference analysis.	
<b>HAT<sub>2</sub> - Main</b> . . . . .	Details are provided in <b>Tab. 14</b> . This variant shares the same network capacity as the official HAT-S implementation, which is slightly heavier than HAT <sub>1</sub> . However, the training budget is reduced for computational efficiency compared to HAT-S <sup>†</sup> (the public model); the patch size and the batch size were halved each. Aligning with Type1 configurations, all models were trained from scratch without the warm-start strategy for 300K.
<b>► Type3 (Extensive Computational Configuration)</b>	
These configurations are used when comparing with official public models and validating the scalability of our method.	
<b>HAT<sup>†</sup> - Main &amp; Appendix</b> . . . .	Details are provided in <b>Tab. 15</b> . This variant shares the same network architecture as the official full-sized HAT implementation and precisely follows the training configuration of the public model. Quantitative results from this variant are copied from the original paper Chen et al. (2023b).

Table 6: Quantitative results for classical image super-resolution under **computationally extensive** setting.  $\dagger$  indicates that we have precisely followed the architecture and training settings of the **official public model**, as specified in Tab.15.  $\text{HAT}^\dagger$  requires 40 GPU days for  $\times 2$  SR (500K train iterations) and additional 20 GPU days for  $\times 4$  SR (250K finetuning iterations). The best results for each setting are highlighted in **bold**, respectively.

Backbone	Scale	Set5		Set14		BSD100		Urban100		Manga109	
		PSNR	SSIM	PSNR	SSIM	PSNR	SSIM	PSNR	SSIM	PSNR	SSIM
$\text{HAT}^\dagger + \text{LN}$	$\times 2$	38.63	.9630	34.86	.9274	32.62	.9053	34.45	.9466	40.26	.9809
$\text{HAT}^\dagger + i\text{-LN}$ (Ours)	$\times 2$	<b>38.65</b>	<b>.9631</b>	<b>34.92</b>	<b>.8276</b>	<b>32.63</b>	<b>.9053</b>	<b>34.60</b>	<b>.9476</b>	<b>40.38</b>	<b>.9811</b>
$\text{HAT}^\dagger + \text{LN}$	$\times 4$	33.04	.9056	29.23	.7973	<b>28.00</b>	.7517	27.97	.8368	32.48	.9292
$\text{HAT}^\dagger + i\text{-LN}$ (Ours)	$\times 4$	<b>33.12</b>	<b>.9064</b>	<b>29.26</b>	<b>.7981</b>	<b>28.00</b>	<b>.7520</b>	<b>28.04</b>	<b>.8388</b>	<b>32.56</b>	<b>.9299</b>

## A EXPERIMENTAL DETAILS

### A.1 MODEL IMPLEMENTATION DETAILS

Since this work provides extensive analysis for more than 60 configurations, analyses throughout this work are performed on various settings due to computational efficiency. We provide implementation details in terms of both network architectural hyperparameters and training configuration for each model variant in Tab.5.

- **Type1 (Lightweight Setting):** These models are the lightweight variants of the original implementations. These configurations are used in Tab.1 and Tab.2, where effects of different normalization schemes and task variations are analyzed.
- **Type2 (Moderate Setting):** These model variants indicate moderate computational budget settings. They are used for the ablation study and also for the analysis in low-precision settings (Tab.3, Tab.4, Fig.10).
- **Type3 (Computationally Extensive Setting):** These model variants indicate computationally extensive settings. This configuration is used to validate the scalability of our method (Tab.6), which aligns with the official implementation of the public models.

### A.2 CHANNEL ENTROPY

Algorithm 1 represents a simple pseudocode to calculate the channel-axis entropy used in our analysis. A sharp drop in channel-axis entropy indicates that feature activations are becoming concentrated in a few specific channels. Analysis throughout this work shows that this entropy collapse is intrinsically linked to the feature divergence problem that arises from conventional LayerNorm in Image Restoration (IR) Transformers.

---

#### Algorithm 1 Channel Entropy Calculation

---

**Require:** Activation tensor  $x$  of shape (C, H, W), a small constant  $\epsilon$  for numerical stability.

**Ensure:** A single scalar entropy value.

- ▶ Step 1: Average the total activation magnitude over spatial-dim.
- ▶ Step 2: Convert to a probability distribution.
- ▶ Step 3: Compute channel entropy.

```

1: function CHANNELENTROPY( $x, \epsilon$ )
2:    $x_{\text{avg}} \leftarrow \text{mean}(\text{abs}(x), \text{dims} = (H, W))$                                 ▷ Step 1
3:    $p \leftarrow \text{softmax}(x_{\text{avg}})$                                                   ▷ Step 2
4:    $\text{entropy} \leftarrow -1 \cdot \text{sum}(p \cdot \log(p + \epsilon))$                              ▷ Step 3
5:   return entropy
6: end function

```

---

Table 7: Quantitative results for  $\times 4$  super-resolution on the SRFormer (Zhou et al., 2023) network architecture. The network capacity and the training budget are adjusted as in Tab.12, which aligns with experimental settings in Tab.2. The best results for each setting are highlighted in **bold**, respectively.

Backbone	Scale	Set5		Set14		BSD100		Urban100		Manga109	
		PSNR	SSIM	PSNR	SSIM	PSNR	SSIM	PSNR	SSIM	PSNR	SSIM
SRFormer <sub>1</sub> + LN	$\times 4$	32.41	.8972	28.77	.7853	27.68	.7398	26.43	.7957	30.98	.9141
SRFormer <sub>1</sub> + <i>i</i> -LN (Ours)	$\times 4$	<b>32.45</b>	<b>.8979</b>	<b>28.81</b>	<b>.7862</b>	<b>27.70</b>	<b>.7407</b>	<b>26.49</b>	<b>.7997</b>	<b>31.10</b>	<b>.9152</b>

Table 8: Quantitative results for  $\times 4$  super-resolution with additional regularization methods. *GC* denotes Gradient Clipping, and *KLD* denotes an auxiliary KL-Divergence loss. Neither proved effective at addressing the instability caused by LayerNorm. GC slightly improves stability but still allows extreme feature magnitudes ( $5.6 \times 10^6$ ), comparable to the vanilla baseline ( $5.8 \times 10^6$ ). KLD regularization enforces smoother statistics but leads to a notable performance drop. In contrast, our proposed *i*-LN yields magnitudes close to  $\mathcal{N}(0, 1)$  (around 1.2) while consistently outperforming all alternatives. The best results for each setting are highlighted in **bold**.

Backbone	Scale	Set5		Set14		BSD100		Urban100		Manga109	
		PSNR	SSIM	PSNR	SSIM	PSNR	SSIM	PSNR	SSIM	PSNR	SSIM
HAT <sub>1</sub> + LN	$\times 4$	32.51	.8992	28.79	.7876	27.68	.7411	26.55	.8015	31.01	.9150
HAT <sub>1</sub> + LN + <i>GC</i>	$\times 4$	32.55	.8996	28.87	.7882	27.74	.7417	26.70	.8037	31.31	.9169
HAT <sub>1</sub> + LN + <i>KLD</i>	$\times 4$	32.36	.8974	28.65	.7853	27.64	.7402	26.34	.7972	30.41	.9105
HAT <sub>1</sub> + <i>i</i> -LN (Ours)	$\times 4$	<b>32.72</b>	<b>.9019</b>	<b>29.01</b>	<b>.7915</b>	<b>27.84</b>	<b>.7456</b>	<b>27.17</b>	<b>.8167</b>	<b>31.82</b>	<b>.9228</b>

## B ADDITIONAL BENCHMARK RESULTS

### B.1 SCALING MODELS AND COMAPRISON AGAINST PUBLIC MODELS

In Tab.6, we validate the scalability of the proposed *i*-LN under computationally extensive settings. Specifically, we train our models on top of the full-sized HAT architecture variant, with the exact training configurations of the public model as specified in Tab.15. The models are indicated as HAT<sup>†</sup>, where <sup>†</sup> means that we have precisely followed the exact network architecture hyperparameters and training configurations for fair comparison. HAT<sup>†</sup> for  $\times 2$  SR and  $\times 4$  SR variants requires 40 GPU days and 20 GPU days on wall-clock time each, with NVIDIA RTX A6000s under the representative BasicSR (Wang et al., 2022b) framework.

**Benchmark.** In Tab.6 we validate that replacing the conventional LayerNorm with the proposed *i*-LN leads to significant performance gain also in the computationally extensive setting where the networks have significantly larger capacity. Accordingly, we conclude that the proposed *i*-LN is effective in both 1) lightweight settings, as shown in our main article and 2) also in computationally extensive settings as in Tab.6, showing the scalability of our *i*-LN.

**Training Details.** Scores are from the original paper for the baselines. Here, we follow the original training scheme where  $\times 4$  SR models are trained under warm-start configuration (i.e., finetuned from  $\times 2$  SR model weight).

### B.2 ADAPTATION TO EFFICIENT SR NETWORK

In Tab.7, we further validate the effectiveness on top of the SRFormer Zhou et al. (2023) architecture, a representative efficient SR network. Similar to other Type1 model variants, the training configurations are adjusted. Refer to Tab.12 for the detailed configurations.

**Discussion.** SRFormer utilizes a Permuted Self-Attention (PSA) mechanism. Accordingly, features across multiple pixels are reshaped into a single feature-pixel (pixelshuffle-style). Thus, the per-token vanilla LN implicitly takes normalization parameters across multi-pixels. While the effect of permuted self-attention in the perspectives of normalization was not discussed in the original work, our work suggests insights that PSA induces (partially) spatial holistcness in normalization, a potential factor for the performance gain of SRFormer (i.e., potentially reducing the performance gap against ours). Seeking further improvements regarding the relationship between the reshaping operation and normalization may be a valuable direction for future work.



## C OTHER REGULARIZATION TECHNIQUES FOR TRAINING STABILITY

Beyond our proposed  $i$ -LN, one may ask whether simpler regularization methods could mitigate the training instabilities of IR Transformers. We therefore examined common strategies such as gradient clipping (GC) and KL divergence (KLD) regularization in Tab.8.

While GC is widely used to bound exploding gradients, our experiments confirmed that it does not prevent the emergence of extreme feature magnitudes in IR Transformers. The maximum feature magnitude observed during training with GC was  $5.6 \times 10^6$ . As a reference, the maximum feature magnitude for the vanilla HAT<sub>1</sub>+LN was  $5.8 \times 10^6$ . In contrary, our HAT<sub>1</sub>+ $i$ -LN shows 1.2, very closely aligning with the expected magnitude of a random noise sampled from the normal distribution  $\mathcal{N}(0, 1)$ , which is 1. Additionally, while GC leads to slight performance improvement against the vanilla model, it consistently underperforms compared to  $i$ -LN.

Likewise, KLD regularization can stabilize feature statistics, but at the cost of substantial reconstruction performance degradation. Specifically, we observed that although KLD encourages well-behaved distributions, the resulting models suffered PSNR drops even below the baseline with vanilla LN. This is consistent with prior findings in rate-distortion theory (Brekelmans et al., 2019; Blau & Michaeli, 2019), and also to VAE literature (Higgins et al., 2017; Yao et al., 2025), where strong regularization penalties reduce reconstruction fidelity.

Overall, these results highlight that although general-purpose regularization may offer partial remedies, they are either ineffective (GC) or detrimental to reconstruction quality (KLD). This further emphasizes the necessity of normalization methods tailored to the unique requirements of IR Transformers.

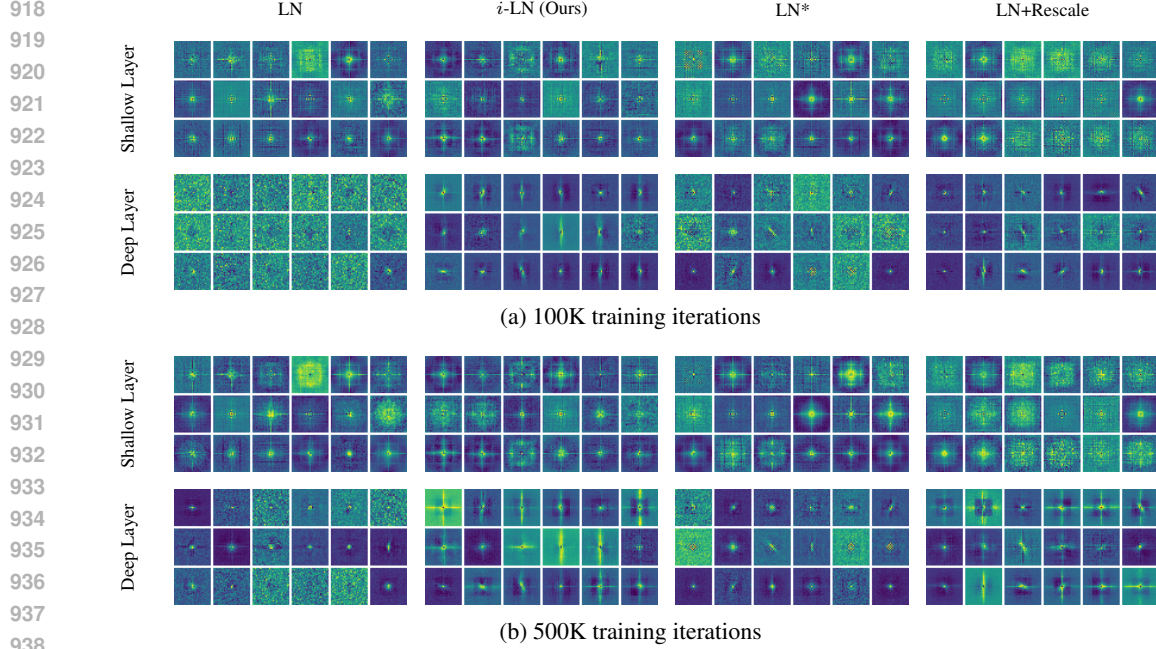


Figure 12: Visual comparison of Relative Position Embeddings (RPE) across attention heads. We show the RPEs from the last three attention layers of each a shallow (RHAG . 0) and a deep (RHAG . 5) building block of HAT Chen et al. (2023a) (i.e., the RHAG Block), as well as early training (100K iterations) versus fully converged models (500K iterations). Each corresponds to an experimental setting aligned with Tab. 3, where vanilla LN and our *i*-LN are the primary comparison, and LN\* and LN+Rescale represent ablation variants obtained by selectively removing components of our method. Our full method (*i*-LN) yields cleaner and more stable RPE patterns, with substantially reduced noise across variations in training iteration and layer depth.

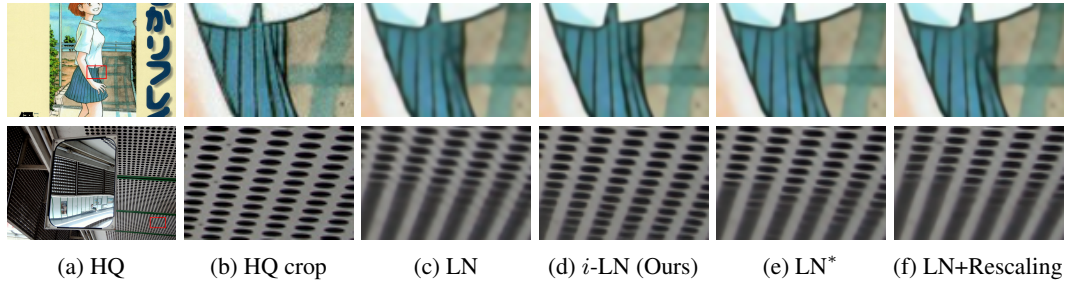


Figure 13: Visual comparison between LN and *i*-LN, along with the ablated variants LN\* and LN+Rescaling with the HAT<sub>2</sub> configuration. Experimental settings follow the ablation study in the main article (Tab. 3). The proposed *i*-LN more faithfully reconstructs fine details, producing sharper edges and clearer complex patterns than the LN baseline.

## D FURTHER ABLATION STUDY

Here, we further compare our full method *i*-LN against the baseline (vanilla) LN, and two variants: LN\* and LN+Rescaling, which each are ablations of our method without rescaling (LN\*) and without the spatial holistic type of normalization (LN+Rescaling). Below, we provide additional visualization of relative position embeddings (RPEs) and also further visual comparison of the according  $\times 4$  SR result. Quantitative comparison can be seen in Tab.3 of the main article. Each of these experimental setups is directly aligned with the ablation analysis presented in Tab. 3 of the main article. Accordingly, all experiments were performed with the HAT<sub>2</sub> configuration for the synthetic (bicubic)  $\times 4$  SR task.

## D.1 ABLATION STUDY: RELATIVE POSITION EMBEDDINGS (RPE) VISUALIZATION

In Fig. 12, we provide a comprehensive visual comparison of the learned Relative Position Embeddings (RPE) across different normalization strategies. The figure compares vanilla Layer Normalization (LN), our proposed *i*-LN, as well as two ablated variants, LN\* and LN+Rescaling.

To obtain deeper insight into the dynamics and stability of RPE formation, we visualize both early-stage training (100K iterations) and fully converged models (500K iterations), and further examine representations from a shallow block (RHAG . 0) and a deep block (RHAG . 5).

Across all settings, vanilla LN exhibits highly noisy RPE structures throughout the entire training process. Even after convergence, its embeddings fail to organize into meaningful spatial patterns, suggesting a limited ability to encode coherent spatial correlations. The LN\* variant, which adopts a spatially holistic normalization, occasionally reveals global structures, but these patterns remain weak and are overshadowed by considerable noise. The LN+Rescaling variant shows improved structure in deeper layers, yet its shallow-layer embeddings remain unstable and inconsistent. This indicating that rescaling alone is insufficient to guide early-layer RPE formation, reflected in low reconstruction scores in Tab.3.

In contrast, our proposed *i*-LN consistently produces substantially clearer and more structured RPE maps, with significantly reduced noise across both shallow and deep layers and across all training stages. The strong spatial coherence visible in the embeddings aligns with the quantitative improvements reported in Tab. 3, where *i*-LN achieves the highest performance among all evaluated variants. These visual results confirm that *i*-LN facilitates stable and meaningful spatial relational modeling throughout the entire network depth and training trajectory.

## D.2 ABLATION STUDY: $\times 4$ SR VISUALIZATION

In Fig.13, we provide additional visual comparisons between the SR outputs obtained by networks each employing LN, LN\*, LN+Rescale and our *i*-LN.

Across all cases, the model equipped with *i*-LN produces the sharpest and most faithful reconstruction of fine-grained structures, including thin edges, repetitive patterns, and high-frequency textures. The restored images exhibit not only improved clarity but also enhanced local contrast and reduced artificial smoothing, indicating that *i*-LN effectively preserves low-level feature statistics throughout the network.

By contrast, the baseline vanilla LN often yields blurry and overly smoothed outputs, where crucial high-frequency details are lost. This degradation is consistent with the unstable and noisy RPE behavior observed earlier, suggesting that vanilla LN struggles to maintain coherent spatial relations required for accurate detail reconstruction. The ablated variants as LN\* (spatial holistic normalization without rescaling) and LN+Rescaling (rescaling without spatial holistictness) show partial improvements over LN but still fall short of *i*-LN.

Overall, the qualitative comparisons provide visual evidence that both components of *i*-LN (spatial holistictness and input-adaptive rescaling) are essential.

## E ADDITIONAL EXPERIMENTS ON THE STABILITY AND ROBUSTNESS

In this section, we provide additional experiments that further validate the stability, robustness, and general applicability of the proposed *i*-LN across a wide range of practical and challenging training scenarios. Specifically, we evaluate: 1) real-world super-resolution settings, 2) robustness under multiple training batch sizes.

Across all settings, *i*-LN consistently outperforms and exhibits more stable behavior than the conventional per-token LayerNorm (LN). Unless otherwise specified, the base configuration follows the HAT<sub>1</sub> setting for the  $\times 4$  SR task. In all experiments, we compare against the vanilla LN baseline under identical training configurations.

### E.1 ROBUSTNESS ACROSS VARYING BATCH SIZES

Run	Backbone	Set5		Set14		BSD100		Urban100		Manga109	
		PSNR	SSIM	PSNR	SSIM	PSNR	SSIM	PSNR	SSIM	PSNR	SSIM
1	RealHAT <sub>1</sub> + LN	26.27	.7601	24.43	.6272	24.39	.5881	22.01	.6064	23.65	.7466
	RealHAT <sub>1</sub> + <i>i</i> -LN	<b>26.38</b>	<b>.7656</b>	<b>24.58</b>	<b>.6322</b>	<b>24.47</b>	<b>.5918</b>	<b>22.19</b>	<b>.6146</b>	<b>23.96</b>	<b>.7557</b>
2	RealHAT <sub>1</sub> + LN	26.77	.7739	24.27	.6201	24.29	.5839	22.00	.6056	23.61	.7503
	RealHAT <sub>1</sub> + <i>i</i> -LN	<b>26.90</b>	<b>.7777</b>	<b>24.46</b>	<b>.6272</b>	<b>24.37</b>	<b>.5876</b>	<b>22.19</b>	<b>.6141</b>	<b>23.84</b>	<b>.7587</b>
3	RealHAT <sub>1</sub> + LN	24.39	.6927	24.77	.6325	24.32	.5823	21.96	.5973	23.71	.7522
	RealHAT <sub>1</sub> + <i>i</i> -LN	<b>24.58</b>	<b>.6993</b>	<b>24.94</b>	<b>.6377</b>	<b>24.40</b>	<b>.5862</b>	<b>22.16</b>	<b>.6055</b>	<b>23.99</b>	<b>.7602</b>
4	RealHAT <sub>1</sub> + LN	25.86	.7472	24.72	.6389	24.35	.5878	21.84	.5954	23.22	.7365
	RealHAT <sub>1</sub> + <i>i</i> -LN	<b>26.05</b>	<b>.7531</b>	<b>24.93</b>	<b>.6447</b>	<b>24.43</b>	<b>.5913</b>	<b>22.04</b>	<b>.6046</b>	<b>23.48</b>	<b>.7450</b>
5	RealHAT <sub>1</sub> + LN	25.08	.7086	24.46	.6355	24.38	.5899	21.93	.6032	23.46	.7458
	RealHAT <sub>1</sub> + <i>i</i> -LN	<b>25.26</b>	<b>.7120</b>	<b>24.63</b>	<b>.6411</b>	<b>24.46</b>	<b>.5933</b>	<b>22.10</b>	<b>.6115</b>	<b>23.73</b>	<b>.7554</b>
Avg.	RealHAT <sub>1</sub> + LN	25.68	.7365	24.53	.6308	24.35	.5864	21.95	.6016	23.53	.7463
	RealHAT <sub>1</sub> + <i>i</i> -LN	<b>25.83</b>	<b>.7415</b>	<b>24.71</b>	<b>.6366</b>	<b>24.43</b>	<b>.5900</b>	<b>22.14</b>	<b>.6101</b>	<b>23.80</b>	<b>.7550</b>

Table 9: Quantitative results for real-world  $\times 4$  super-resolution across five random seeds. Both the training images and the test images were synthesized following the Real-ESRGAN (Wang et al., 2021) degradation pipeline. The best results for each setting are highlighted in **bold**.

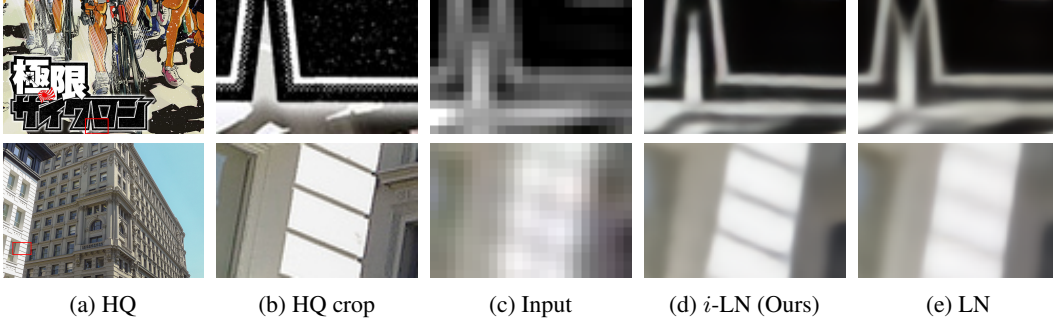


Figure 15: Visual comparison between LN and *i*-LN for the real-world  $\times 4$  super-resolution task with HAT<sub>1</sub>. Both the training images and the test images were synthesized following the Real-ESRGAN (Wang et al., 2021) degradation pipeline.

To evaluate the robustness of *i*-LN, we train HAT<sub>1</sub> models but with varying batch size from 2 to 8; while keeping all other hyperparameters fixed (Fig. 14). This experimental design is chosen in order to mimic unstable training configurations without heavy hyperparameter search. For each configuration, we report PSNR and SSIM for the  $\times 4$  SR task on Urban100 throughout training. As shown in Fig. 14, models with *i*-LN consistently achieve higher PSNR/SSIM than the baseline across all batch sizes. Notably, the performance gap remains stable as the batch size decreases. These results demonstrate that the benefit of *i*-LN is not tied to a particular training setup, and that its stability advantages persist even under extremely small-batch training (e.g., batch size 2). This property is especially valuable for memory-constrained environments where large batches are infeasible.

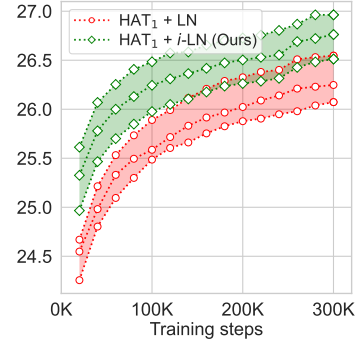


Figure 14: PSNR for  $\times 4$  SR with the HAT<sub>1</sub> model, but with batch size 2, 4, 8 (Urban100).

## E.2 REAL-WORLD SUPER-RESOLUTION

To evaluate the practical effectiveness of *i*-LN, we adopt a real-world degradation setup following the RealESRGAN Wang et al. (2021) pipeline. These analyses complement the main text by demonstrating that the advantages of *i*-LN are not restricted to controlled laboratory settings, but generalize to real-world usage and unstable training regimes.

Training is performed on synthetic DF2K pairs, and testing is conducted on synthetically degraded versions of standard SR benchmarks (Set5, Set14, BSD100, Urban100, Manga109). Degradation syn-

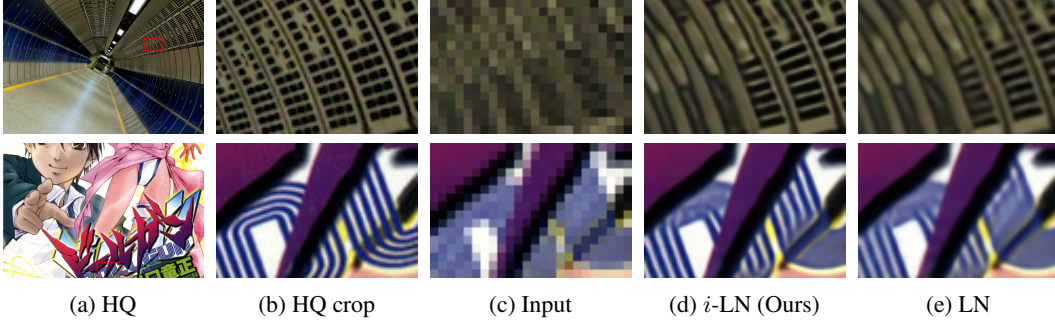
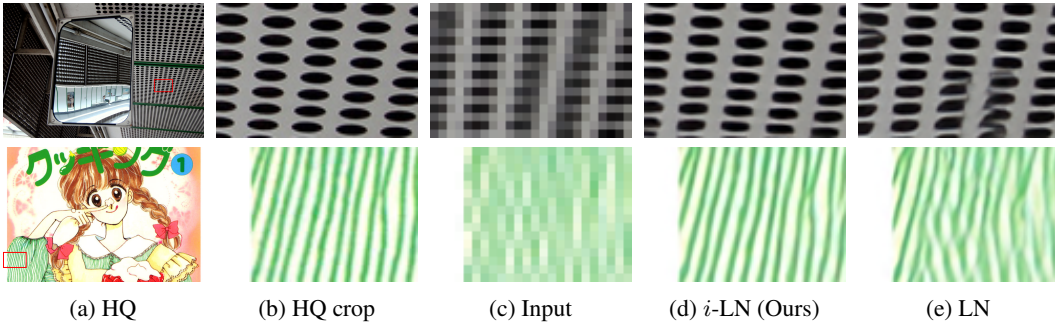
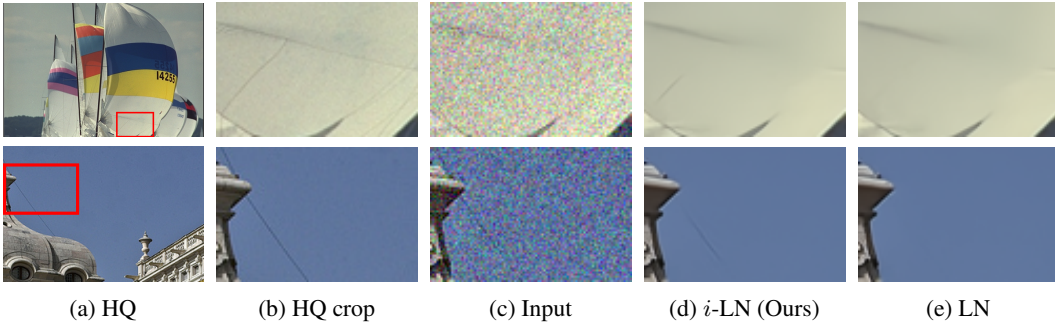
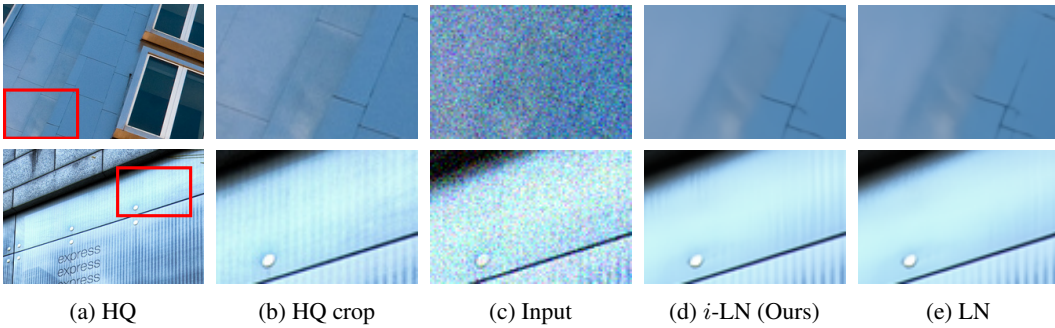
thesis strictly follows the RealESRGAN procedures. To assess robustness, we repeat each experiment across five different random seeds and report the average score of the resulting PSNR/SSIM scores.

As shown in Tab.9 and Fig.15, *i*-LN leads to consistent and significant improvements over vanilla LN across all benchmarks, despite the increased complexity of the real-world degradation pipeline. These results highlight that *i*-LN is not only theoretically grounded but also practically beneficial in more challenging real-world restoration scenarios.

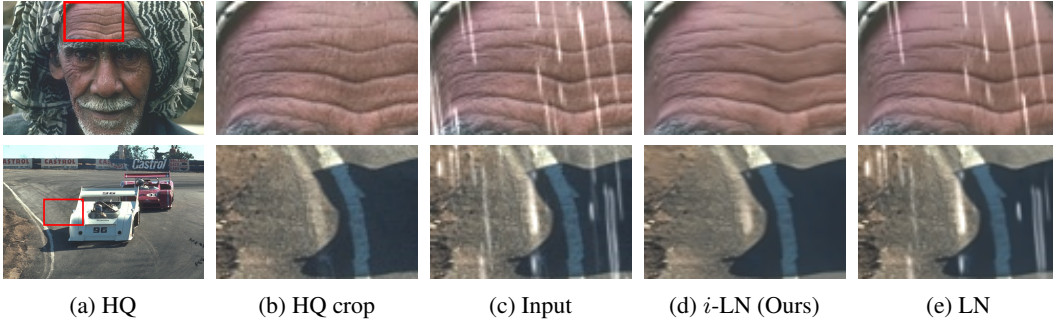
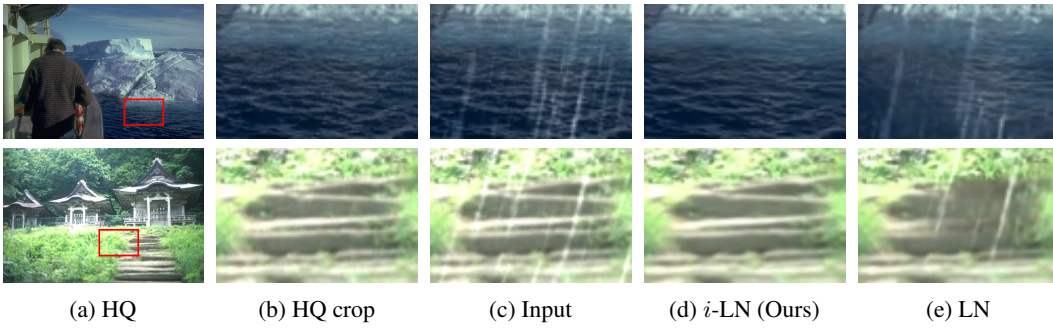
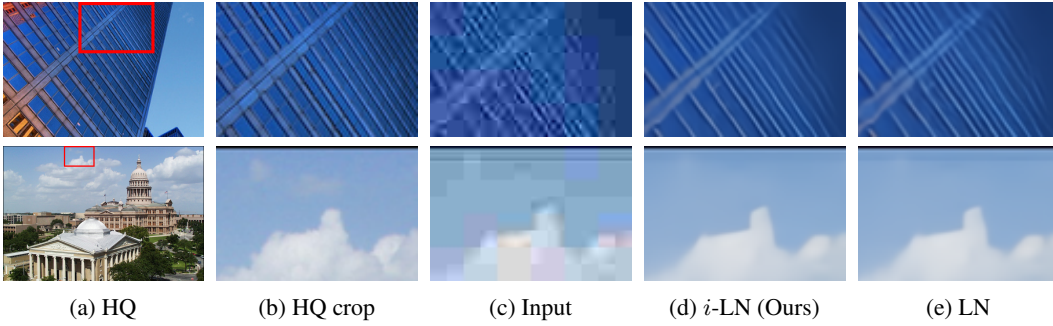
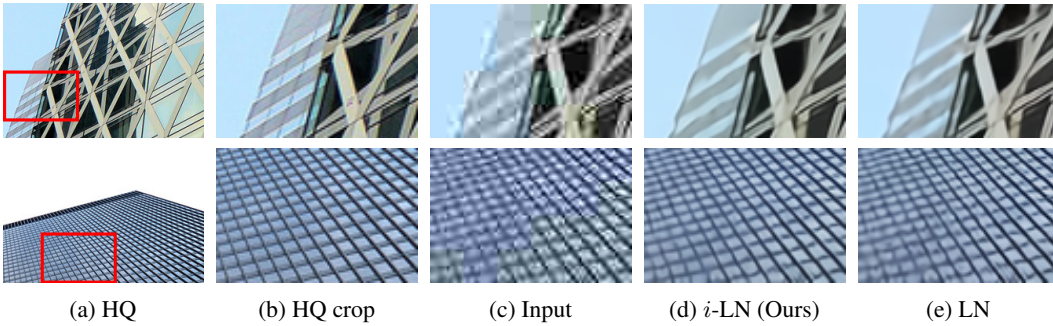
## F ADDITIONAL QUALITATIVE COMPARISON

We provide additional visual comparisons between IR Transformers with our proposed *i*-LN against their counterparts using vanilla Layer Normalization (LN). As shown in the following figures, *i*-LN consistently produces sharper structures, cleaner textures, and fewer artifacts across a range of low-level vision tasks. Qualitative results are provided for: (i) super-resolution (Figs. 16 and 17), (ii) image denoising (Figs. 18 and 19), (iii) JPEG compression artifact removal (Figs. 23 and 22), and (iv) image deraining (Figs. 20 and 21).



Figure 16: Visual comparison between LN and *i*-LN for the super-resolution task with  $HAT_1$ .Figure 17: Visual comparison between LN and *i*-LN for the super-resolution task with  $DRCT_1$ .Figure 18: Visual comparison between LN and *i*-LN for the denoising task with  $HAT_1$ .Figure 19: Visual comparison between LN and *i*-LN for the denoising task with  $SwinIR_1$ .



Figure 20: Visual comparison between LN and *i*-LN for the deraining task with HAT<sub>1</sub>.Figure 21: Visual comparison between LN and *i*-LN for the deraining task with SwinIR<sub>1</sub>.Figure 22: Visual comparison between LN and *i*-LN for the JPEG artifact removal task with HAT<sub>1</sub>.Figure 23: Visual comparison between LN and *i*-LN for the JPEG artifact removal task with SwinIR<sub>1</sub>.

## G DETAILED DERIVATIONS FOR STRUCTURE PRESERVATION

### G.1 NOTATION AND PRELIMINARIES

Let  $X \in \mathbb{R}^{L \times C}$  be the feature matrix with tokens  $x_\ell \in \mathbb{R}^C$  (row-vectors). Define the inter-pixel (inter-token) structure by the set of pairwise displacements

$$\Delta X := \{x_\ell - x_k : 1 \leq \ell, k \leq L\}.$$

A map  $T : \mathbb{R}^C \rightarrow \mathbb{R}^C$  *preserves inter-pixel structure up to scale* if there exists a homothety  $H(x) = ax + b$  with  $a > 0, b \in \mathbb{R}^C$  such that

$$T(x_\ell) - T(x_k) = a(x_\ell - x_k) \quad \text{for all } \ell, k.$$

Equivalently, all angles and pairwise distance ratios are preserved.

We analyze the *pure normalization maps* (i.e., the normalization before the affine  $(\gamma, \beta)$  is applied); any global translation by  $\beta$  does not affect  $\Delta X$ , and a scalar post-scale can be absorbed into the homothety factor  $a$ .

### G.2 PROPOSITION 1 (VANILLA LAYER NORM FAILS TO PRESERVE STRUCTURE)

Let  $T_{\text{LN}}$  denote the transformation defined by the normalization in vanilla *per-token* LayerNorm. In general there do not exist  $a > 0$  and an orthogonal  $Q$  such that for all tokens  $x_\ell, x_k$ ,

$$T_{\text{LN}}(x_\ell) - T_{\text{LN}}(x_k) = a Q(x_\ell - x_k).$$

Hence  $T_{\text{LN}}$  is not conformal on the token set. By the nested class relation Homothety  $\subset$  Similarity  $\subset$  Conformal, it follows that  $T_{\text{LN}}$  is neither a similarity nor a homothety, and thus does *not* preserve inter-pixel structure in general.

**Proof.** Write per-token means and standard deviations as

$$\mu_\ell = \frac{1}{C} \sum_{c=1}^C x_{\ell,c}, \quad \sigma_\ell = \left( \frac{1}{C} \sum_{c=1}^C (x_{\ell,c} - \mu_\ell)^2 \right)^{1/2}.$$

The pure LN map (no  $\gamma, \beta$ ) acts componentwise as

$$T_{\text{LN}}(x_\ell) = \frac{x_\ell - \mu_\ell \mathbf{1}}{\sigma_\ell},$$

so for two tokens  $\ell, k$ ,

$$\Delta_{\ell k} := T_{\text{LN}}(x_\ell) - T_{\text{LN}}(x_k) \tag{4}$$

$$= \frac{x_\ell}{\sigma_\ell} - \frac{x_k}{\sigma_k} - \left( \frac{\mu_\ell}{\sigma_\ell} - \frac{\mu_k}{\sigma_k} \right) \mathbf{1}. \tag{1}$$

Assume for contradiction there exist  $a > 0$  and orthogonal  $Q$  such that  $\Delta_{\ell k} = a Q(x_\ell - x_k)$  for all  $\ell, k$ . Compare the coefficients of  $x_\ell$  and  $x_k$  on both sides of (1). Because the equality must hold for arbitrary token values, we must have

$$\frac{1}{\sigma_\ell} I = aQ \quad \text{and} \quad \frac{1}{\sigma_k} I = aQ,$$

hence  $\sigma_\ell = \sigma_k$  for all  $\ell, k$  and  $Q$  must be proportional to the identity. With  $\sigma_\ell \equiv \sigma$ , the bias term in (1) reduces to  $\left( \frac{\mu_k - \mu_\ell}{\sigma} \right) \mathbf{1}$ , which must vanish for all  $\ell, k$ ; thus  $\mu_\ell = \mu_k$  for all  $\ell, k$ . Therefore the assumed similarity can hold only in the *degenerate* case where all tokens share identical per-token mean and variance.

For real features,  $\{\mu_\ell, \sigma_\ell\}$  are not constant across tokens, so the assumption leads to a contradiction. Hence no single similarity map exists;  $T_{\text{LN}}$  is not conformal and does not preserve spatial structure.

**Remark.** The degenerate equal-statistics case is precisely the rare exception noted in the main text.

### G.3 PROPOSITION 2 (LN\* PRESERVES STRUCTURE)

Let  $T_{\text{LN}^*}$  denote the transformation defined by the normalization in *spatially holistic* LayerNorm (LN\*) with global mean and standard deviation

$$\mu = \frac{1}{LC} \sum_{\ell,c} x_{\ell,c}, \quad \sigma = \left( \frac{1}{LC} \sum_{\ell,c} (x_{\ell,c} - \mu)^2 \right)^{1/2} > 0.$$

Then for any tokens  $x_\ell, x_k$ ,

$$T_{\text{LN}^*}(x_\ell) - T_{\text{LN}^*}(x_k) = \frac{1}{\sigma}(x_\ell - x_k),$$

so  $T_{\text{LN}^*}$  is a homothety and preserves inter-pixel structure up to a global scale.

**Proof.**  $T_{\text{LN}^*}$  (without  $\gamma, \beta$ ) is

$$T_{\text{LN}^*}(x) = \frac{x - \mu \mathbf{1}}{\sigma}.$$

Hence

$$T_{\text{LN}^*}(x_\ell) - T_{\text{LN}^*}(x_k) = \frac{x_\ell - \mu \mathbf{1}}{\sigma} - \frac{x_k - \mu \mathbf{1}}{\sigma} = \frac{1}{\sigma}(x_\ell - x_k).$$

This is exactly a homothety with scale factor  $a = \sigma^{-1}$ ; therefore angles and pairwise distance ratios of  $\Delta X$  are preserved and the spatial configuration is rigid up to a uniform scale.

Network Architecture Hyperparameters	
Embedding Dimension	60
Layer Depths	[6, 6, 6, 6]
Attention Heads	[6, 6, 6, 6]
Window Size	$8 \times 8$
MLP Ratio	2
Residual Connection	'1conv'
Dataset Configuration	
Training Dataset	DIV2K + Flickr2K
PatchSize   BatchSize	
- Denoising	$64 \times 64   16$
- Deraining	$128 \times 128   8$
- JPEG Artifact Removal	$64 \times 64   16$
Noise Degradation	torch.randn
JPEG Degradation	OpenCV
Optimizing Configuration	
Total Iterations	300K
Optimizer	Adam
Learning Rate (LR)	$2 \times 10^{-4}$
Adam Betas	(0.9, 0.99)
Weight Decay	0
Scheduler ( $\gamma = 0.5$ )	StepLR
Milestones (K)	250
Loss Function	L1 Loss

Table 10: **Hyperparameters and training configurations for the model variant SwinIR<sub>1</sub>**. Network architectural terminology is based on terminologies either in the official implementation of each work or the implementation in BasicSR Wang et al. (2022b).

Network Architecture Hyperparameters	
Embedding Dimension	96
Layer Depths	[6, 6, 6, 6, 6, 6]
Attention Heads	[6, 6, 6, 6, 6, 6]
Window Size	$16 \times 16$
MLP Ratio	2
Residual Connection	'1conv'
Dataset Configuration	
Training Dataset	DIV2K + Flickr2K
PatchSize   BatchSize	
- $\times 2$ Super Resolution	$64 \times 64   16$
- $\times 4$ Super Resolution	$128 \times 128   16$
SR Degradation	MATLAB
Optimizing Configuration	
Total Iterations	300K
Optimizer	Adam
Learning Rate (LR)	$2 \times 10^{-4}$
Adam Betas	(0.9, 0.99)
Weight Decay	0
Scheduler ( $\gamma = 0.5$ )	StepLR
Milestones (K)	250
Loss Function	L1 Loss

Table 11: **Hyperparameters and training configurations for the model variant DRCT<sub>1</sub>**. Network architectural terminology is based on terminologies either in the official implementation of each work or the implementation in BasicSR Wang et al. (2022b).

Network Architecture Hyperparameters	
Embedding Dimension	60
Layer Depths	[6, 6, 6, 6]
Attention Heads	[6, 6, 6, 6]
Window Size	$16 \times 16$
MLP Ratio	2
Residual Connection	'1conv'
Dataset Configuration	
Training Dataset	DIV2K + Flickr2K
PatchSize   BatchSize	
- $\times 4$ Super Resolution	$128 \times 128$   16
SR Degradation	MATLAB
Optimizing Configuration	
Total Iterations	300K
Optimizer	Adam
Learning Rate (LR)	$2 \times 10^{-4}$
Adam Betas	(0.9, 0.99)
Weight Decay	0
Scheduler ( $\gamma = 0.5$ )	StepLR
Milestones (K)	250
Loss Function	L1 Loss

Table 12: **Hyperparameters and training configurations for the model variant SRFormer<sub>1</sub>**. Network architectural terminology is based on terminologies either in the official implementation of each work or the implementation in BasicSR Wang et al. (2022b).

Network Architecture Hyperparameters	
Embedding Dimension	96
Layer Depths	[6, 6, 6, 6, 6, 6]
Attention Heads	[6, 6, 6, 6, 6, 6]
Window Size	$16 \times 16$
MLP Ratio	2
Compress Ratio	24
Squeeze Factor	24
Overlap Ratio	0.5
Conv Scale	0.01
Residual Connection	'1conv'
Dataset Configuration	
Training Dataset	DIV2K + Flickr2K
PatchSize   BatchSize	
- Denoising	$64 \times 64$   16
- Deraining	$128 \times 128$   8
- JPEG Artifact Removal	$64 \times 64$   16
- $\times 2$ Super Resolution	$64 \times 64$   16
- $\times 4$ Super Resolution	$128 \times 128$   16
Noise Degradation	torch.randn
JPEG Degradation	OpenCV
SR Degradation	MATLAB
Optimizing Configuration	
Total Iterations	300K
Optimizer	Adam
Learning Rate (LR)	$2 \times 10^{-4}$
Adam Betas	(0.9, 0.99)
Weight Decay	0
Scheduler ( $\gamma = 0.5$ )	StepLR
Milestones (K)	250
Loss Function	L1 Loss

Table 13: **Hyperparameters and training configurations for the model variant HAT<sub>1</sub>**. Network architectural terminology is based on terminologies either in the official implementation of each work or the implementation in BasicSR Wang et al. (2022b).

Network Architecture Hyperparameters	
<b>Embedding Dimension</b>	<b>144</b>
Layer Depths	[6, 6, 6, 6, 6, 6]
Attention Heads	[6, 6, 6, 6, 6, 6]
Window Size	$16 \times 16$
MLP Ratio	2
Compress Ratio	24
Squeeze Factor	24
Overlap Ratio	0.5
Conv Scale	0.01
Residual Connection	‘1conv’
Dataset Configuration	
Training Dataset	DIV2K + Flickr2K
PatchSize   BatchSize	
- $\times 4$ Super Resolution	$128 \times 128$   16
SR Degradation	MATLAB
Optimizing Configuration	
<b>Total Iterations</b>	<b>500K</b>
Optimizer	Adam
Learning Rate (LR)	$2 \times 10^{-4}$
Adam Betas	(0.9, 0.99)
Weight Decay	0
<b>Scheduler</b> ( $\gamma = 0.5$ )	<b>MultiStepLR</b>
<b>Milestones (K)</b>	<b>[250, 400, 450, 475]</b>
Loss Function	L1 Loss

Table 14: **Hyperparameters and training configurations for HAT<sub>2</sub>**. This variant uses the **HAT-S** architecture but is trained with a reduced budget. Differences from HAT<sub>1</sub> are highlighted in **bold**. Network architectural terminology is based on terminologies either in the official implementation of each work or the implementation in BasicSR Wang et al. (2022b).



Network Architecture Hyperparameters	
Embedding Dimension	180
Layer Depths	[6, 6, 6, 6, 6, 6]
Attention Heads	[6, 6, 6, 6, 6, 6]
Window Size	$16 \times 16$
MLP Ratio	2
Compress Ratio	3
Squeeze Factor	30
Overlap Ratio	0.5
Conv Scale	0.01
Residual Connection	'1conv'
Dataset Configuration	
Training Dataset	DIV2K + Flickr2K
PatchSize   BatchSize	
- $\times 2$ Super Resolution	$128 \times 128 \mid 32$
- $\times 4$ Super Resolution	$256 \times 256 \mid 32$
SR Degradation	MATLAB
Optimizing Configuration	
Optimizer	Adam
Adam Betas	(0.9, 0.99)
Weight Decay	0
Scheduler ( $\gamma = 0.5$ )	MultiStepLR
Loss Function	L1 Loss
$\times 2$ Super Resolution	
- Total Iterations	500K
- Learning Rate (LR)	$2 \times 10^{-4}$
- Scheduler Milestones (K)	[250, 400, 450, 475]
$\times 4$ Super Resolution	
- Total Iterations	250K
- Learning Rate (LR)	$1 \times 10^{-4}$
- Scheduler Milestones (K)	[125, 200, 225, 240]
- Pretrained	finetune from $\times 2$ SR weight

Table 15: **Hyperparameters and training configurations for HAT<sup>†</sup>**. This variant uses the full-sized HAT architecture and precisely follows the training settings of the public model. Network architectural terminology is based on terminologies either in the official implementation of each work or the implementation in BasicSR Wang et al. (2022b).

1566 H THE USE OF LARGE LANGUAGE MODELS (LLMs)  
1567

1568 In this study, LLMs were used for text editing, grammar correction, and coding assistance for graph  
1569 visualization.  
1570  
1571  
1572  
1573  
1574  
1575  
1576  
1577  
1578  
1579  
1580  
1581  
1582  
1583  
1584  
1585  
1586  
1587  
1588  
1589  
1590  
1591  
1592  
1593  
1594  
1595  
1596  
1597  
1598  
1599  
1600  
1601  
1602  
1603  
1604  
1605  
1606  
1607  
1608  
1609  
1610  
1611  
1612  
1613  
1614  
1615  
1616  
1617  
1618  
1619