

WHITE-BOX PROMPT TRANSFORMERS: VARIATIONALLY GROUNDED PROMPT-ATTENTION COUPLING FOR UNIFIED IMAGE RESTORATION

Anonymous authors

Paper under double-blind review

ABSTRACT

Can soft prompts in vision Transformers be made explainable? Prompt-based models have achieved remarkable success in image restoration, yet they remain largely opaque: the underlying Transformer operations and the mechanism by which prompts modulate attention are poorly understood. This work revisits *guided image restoration*, where an auxiliary modality A assists in restoring a target modality B . We interpret A as a prompt and formulate a tailored structure-tensor total variation (STV) model, whose gradient suggests a white-box correspondence to prompt-attention interactions. This provides a principled bridge between prompts and attention. In scenarios where A is unavailable, we abstract its role into learnable soft prompts, enabling end-to-end training within standard Transformer pipelines. By unrolling the gradient flow of the STV variational problem, we derive the *White-Box Prompt Transformer (WBPT)*, a cascaded architecture that embeds interpretability directly into attention operations. Extensive experiments on multiple benchmarks demonstrate that WBPT achieves state-of-the-art restoration performance while offering interpretable, controllable, and robust prompt-attention dynamics.

1 INTRODUCTION

Prompt-based Transformers have recently reshaped unified image restoration, enabling a single model to tackle diverse degradations through learnable soft prompts (Potlapalli et al., 2023). These prompts condition the restoration process by modulating attention mechanisms and consistently deliver strong empirical results (e.g., Jia et al., 2022; Kong et al., 2025). However, despite their success, prompt-based designs remain fundamentally opaque: the inner workings of the Transformer and the interaction between prompts and attention lack interpretability, limiting both theoretical understanding and practical controllability (Chefer et al., 2021; Jain & Wallace, 2019). This opacity impedes reliable deployment in trust-sensitive applications (Rudin, 2019).

This motivates our central question:

Can prompt-driven attention be explained from first principles, providing a theoretically grounded interpretation of the black box?

We draw inspiration from *guided image restoration*, where an auxiliary modality A (e.g., T_1 -weighted MRI) provides structural guidance for restoring a target modality B (He et al., 2012; Li et al., 2016; Ehrhardt & Betcke, 2016a). In this setting, A acts as a prior, naturally analogous to a prompt guiding the restoration of B (Jia et al., 2022; Potlapalli et al., 2023). Since explicit auxiliary data are often unavailable (Havaei et al., 2016), we abstract the role of A into learnable soft prompts—trainable tokens that emulate auxiliary guidance through end-to-end optimization. This perspective reinterprets prompts not as heuristic inputs but as principled surrogates for classical guidance (Li & Liang, 2021; Zhou et al., 2022).

Building on this analogy, we introduce a variational perspective on prompt-based restoration. Specifically, we cast guided restoration as a *structure-tensor total variation (STV)* problem (Chambolle & Pock, 2011; Lefkimmiatis et al., 2015). Through gradient analysis, we show that the optimization dynamics naturally align with a white-box attention mechanism, suggesting a formal link

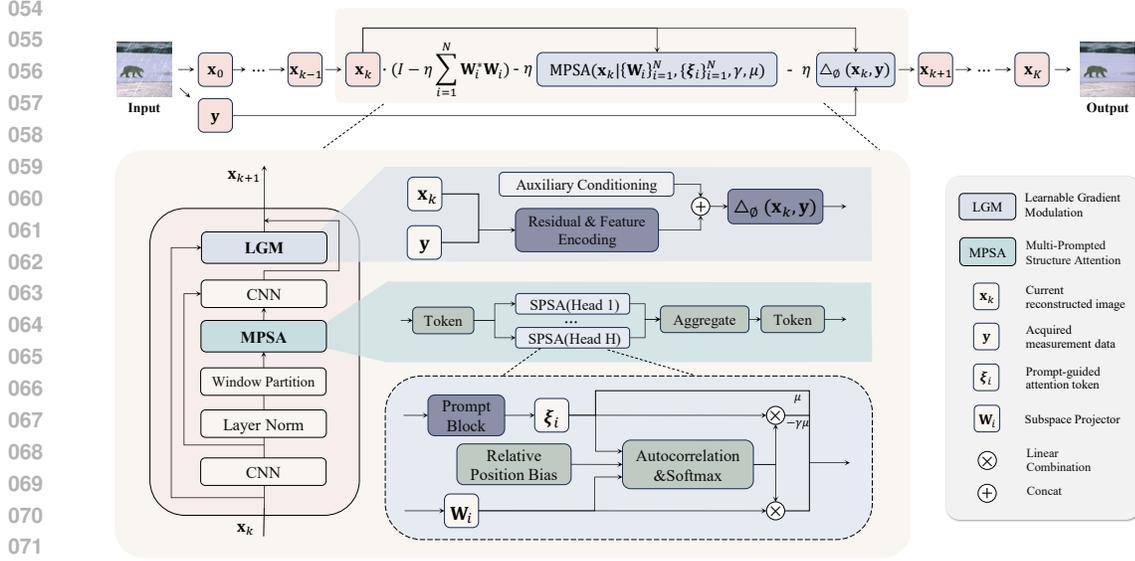


Figure 1: Overview of the White-box Prompt Transformer (WBPT). Image restoration is achieved by unrolling K gradient-flow steps. At iteration k (Eq. 8), the update combines a Multi-Prompted Structure Attention (MPSA) block with a learnable gradient-modulation (LGM) data-consistency term. MPSA consists of Single-Prompted Structure Attention (SPSA) heads: tokenized features interact with prompt tokens ξ_i and learnable projectors W_i , are aggregated, and mapped back to image domain. Across stages, prompts interact with features at multiple levels, enriching structural context while preserving fidelity to the measurements.

between guidance priors (or their prompt surrogates) and Transformer attention (Yu et al., 2023; Wang et al., 2018; Meng et al., 2024). While the derivation involves approximations, it provides a principled foundation for understanding and designing interpretable prompt-driven restoration.

Crucially, this formulation not only offers interpretability but also suggests a concrete architectural design. By unrolling the gradient flow of the *STV* variational problem, we obtain the *White-Box Prompt Transformer (WBPT)* (Chen et al., 2015; Monga et al., 2021). Each Transformer layer corresponds to an optimization step, with every attention operation tied intuitively to terms in the underlying energy functional. WBPT thus unifies variational analysis with deep learning, embedding interpretability into the model without compromising performance.

Contributions.

- *Variational Perspective on Prompt-based Restoration:* Guided restoration is formulated as a tailored *STV* problem. Its optimization dynamics reveal a white-box attention mechanism, offering a principled explanation of how prompts influence Transformer attention.
- *White-Box Prompt Transformer:* A cascaded Transformer derived by unrolling the *STV* gradient flow, where each layer corresponds to an optimization step and attention operations align with terms in the underlying energy functional.
- *Bridging Classical and Modern AI:* The framework connects variational principles with deep prompt-based models, providing a foundation for interpretable image restoration and controllable attention mechanisms with clear theoretical grounding.
- *Empirical Validation:* WBPT achieves state-of-the-art results on multiple image restoration benchmarks while enabling transparent analysis of prompt-attention dynamics via rigorous, comprehensive visualization and controlled perturbation studies.

2 METHODS

In this section, we introduce WBPT, a variationally inspired framework for guided image restoration that interprets the guidance modality as *soft prompts* in a principled manner. WBPT integrates a tailored *STV* prior with learnable transformations W_i and soft prompt tokens ξ_i , achieving restoration by unrolling K gradient-flow steps. The overall pipeline and information flow across the K cascaded stages are illustrated in Fig. 1, as depicted schematically.

2.1 STRUCTURED MODELING FRAMEWORK FOR GUIDED IMAGE RESTORATION

In guided image restoration, structural consistency across image modalities—where one modality (e.g., modality A) provides complementary information to enhance the restoration of another modality (e.g., modality B)—can be effectively exploited in practice. Rapidly acquired or higher-quality modality A images can serve as informative priors to guide the restoration of modality B (Ehrhardt & Betcke, 2016b). To systematically and rigorously model this guidance, we treat the modality A image as a *prompt*, explicitly encoding its anatomical information through dedicated operators to assist in restoring modality B (Potlapalli et al., 2023; Jia et al., 2022).

Formally, guided image restoration can be expressed as the following optimization problem:

$$\min_{\mathbf{x} \in \mathcal{X}} \mathcal{R}(\mathbf{x}; \boldsymbol{\xi}) := \mathcal{R}_1(\mathbf{x}) - \mu \mathcal{R}_2(\mathbf{x}; \boldsymbol{\xi}), \quad (1)$$

where $\mathbf{x} : \Omega \rightarrow \mathbb{R}^d$ denotes the target image (modality B) to be restored, and $\boldsymbol{\xi}$ represents the feature embedding of the guidance image (modality A). The function space \mathcal{X} is chosen appropriately (e.g., Sobolev space $H^1(\Omega; \mathbb{R}^d)$ or $L^2(\Omega; \mathbb{R}^d)$) to ensure well-definedness of the optimization and its variational derivatives (Evans, 2022). The functional $\mathcal{R}_1(\mathbf{x})$ encodes intrinsic priors of \mathbf{x} , while minimizing $-\mathcal{R}_2(\mathbf{x}; \boldsymbol{\xi})$ enforces consistency between \mathbf{x} and $\boldsymbol{\xi}$ in the transformed domain. The parameter $\mu > 0$ balances this trade-off.

Specifically, we redesign an enhanced *STV* prior to represent \mathcal{R} , which not only characterizes the structural priors of the target image but also enforces consistency with guidance features $\boldsymbol{\xi}_i$ in the domain defined by \mathbf{W}_i . This motivates our proposed weighted prompt formulation:

$$\begin{aligned} \mathcal{R}(\mathbf{x}; \{\mathbf{W}_i\}_{i=1}^N, \{\boldsymbol{\xi}_i\}_{i=1}^N) := & \frac{1}{2} \sum_{i=1}^N \int_{\Omega} \text{Tr}(\psi(\mathbf{W}_i \mathbf{x}(s)(\mathbf{W}_i \mathbf{x}(s))^\top)) ds \\ & - \mu \sum_{i=1}^N \int_{\Omega} \text{Tr}(\psi(\boldsymbol{\xi}_i(s)(\mathbf{W}_i \mathbf{x}(s))^\top)) ds. \end{aligned} \quad (2)$$

Classical *STV* instantiates \mathbf{W}_i as gradient operators capturing local structures, whereas nonlocal *STV* incorporates global interactions for superior performance. Motivated by this, we parameterize \mathbf{W}_i as learnable global transformations via fully connected layers rather than local convolutional kernels (Wang et al., 2018). In parallel, $\{\boldsymbol{\xi}_i\}_{i=1}^N$ serve as prompts in the transformed domain. When explicit guidance images are unavailable, these prompts are relaxed into learnable soft tokens (Jia et al., 2022; Potlapalli et al., 2023). Finally, $\psi(\cdot)$ is a sparsity-inducing penalty, for which nonconvex forms such as $\psi(u) = \ln(1 + u)$ are effective.

2.2 WHITE-BOX PROMPT TRANSFORMER VIA VARIATIONAL DERIVATION

The gradient of the energy functional (2) is derived via variational calculus, resulting in an interpretable form:

$$\begin{aligned} \frac{\delta \mathcal{R}(\mathbf{x}; \mathbf{W}_i, \boldsymbol{\xi}_i)}{\delta \mathbf{x}} \approx & \mathbf{W}_i^* \mathbf{W}_i \mathbf{x} + \gamma \mathbf{W}_i^* \mathbf{W}_i \mathbf{x} \cdot \text{softmax}((\mathbf{W}_i \mathbf{x})^\top \mathbf{W}_i \mathbf{x}) \\ & - \mu \mathbf{W}_i^* \boldsymbol{\xi}_i - \gamma \mu \mathbf{W}_i^* \boldsymbol{\xi}_i \cdot \text{softmax}((\mathbf{W}_i \mathbf{x})^\top \boldsymbol{\xi}_i). \end{aligned} \quad (3)$$

Here, $\mathcal{R}(\mathbf{x}; \mathbf{W}_i, \boldsymbol{\xi}_i)$ denotes the i -th component of

$$\mathcal{R}(\mathbf{x}; \{\mathbf{W}_i\}_{i=1}^N, \{\boldsymbol{\xi}_i\}_{i=1}^N) = \sum_{i=1}^N \mathcal{R}(\mathbf{x}; \mathbf{W}_i, \boldsymbol{\xi}_i).$$

Although Eq.3 replaces linear inner-product weights by column-wise softmax weights, Appendix C shows that, after optimal column-wise scaling, the resulting attention-based form admits a closed-form error expression with an explicit condition under which this approximation is highly accurate.

This gradient inspires the *Single-Prompted Structure Attention (SPSA)* module:

$$\text{SPSA}(\mathbf{x} \mid \mathbf{W}_i, \boldsymbol{\xi}_i, \gamma, \mu) := \mathbf{W}_i \mathbf{x} \cdot \text{softmax}((\mathbf{W}_i \mathbf{x})^\top \mathbf{W}_i \mathbf{x}) - \gamma \mu \boldsymbol{\xi}_i \cdot \text{softmax}((\mathbf{W}_i \mathbf{x})^\top \boldsymbol{\xi}_i) + \mu \boldsymbol{\xi}_i. \quad (4)$$

Eq.4 can thus be viewed as an operator-level implementation of this gradient-to-attention mapping, and inherits the same error bound and validity condition.

For multiple prompts, we define the *Multi-Prompted Structure Attention (MPSA)* module:

$$\text{MPSA}(\mathbf{x} \mid \{\mathbf{W}_i\}, \{\boldsymbol{\xi}_i\}, \gamma, \mu) := [\mathbf{W}_1^* \quad \cdots \quad \mathbf{W}_N^*] \begin{bmatrix} \text{SPSA}(\mathbf{x} \mid \mathbf{W}_1, \boldsymbol{\xi}_1, \gamma, \mu) \\ \vdots \\ \text{SPSA}(\mathbf{x} \mid \mathbf{W}_N, \boldsymbol{\xi}_N, \gamma, \mu) \end{bmatrix}. \quad (5)$$

This formulation provides an explicitly controllable, prompt-driven white-box attention mechanism with three functional components:

- *Self-Reconstruction Term:* $\mathbf{W}_i \mathbf{x} \cdot \text{softmax}((\mathbf{W}_i \mathbf{x})^\top \mathbf{W}_i \mathbf{x})$, enhancing intrinsic feature coherence via self-expression.
- *Prompt-Alignment Term:* $-\lambda_1 \boldsymbol{\xi}_i \cdot \text{softmax}((\mathbf{W}_i \mathbf{x})^\top \boldsymbol{\xi}_i)$, introducing a repulsive force to prevent trivial imitation while enabling structural adaptation.
- *Prompt-Bias Term:* $+\lambda_2 \boldsymbol{\xi}_i$, injecting prior knowledge as a static inductive bias to ensure faithful restoration.

2.3 CASCADED TRANSFORMER ARCHITECTURE VIA GRADIENT FLOW UNROLLING

To optimize (2) while enforcing consistency with measurements \mathbf{y} , we consider the continuous-time gradient flow:

$$\frac{\partial \mathbf{x}(t)}{\partial t} = - \left(\frac{\delta \mathcal{R}}{\delta \mathbf{x}}(\mathbf{x}(t)) + \Delta(\mathbf{x}(t), \mathbf{y}) \right), \quad (6)$$

where $\Delta(\mathbf{x}(t), \mathbf{y})$ denotes the gradient of the data fidelity term.

Discretizing via explicit Euler with step size η gives:

$$\mathbf{x}_{k+1} = \mathbf{x}_k - \eta \left(\frac{\delta \mathcal{R}}{\delta \mathbf{x}}(\mathbf{x}_k) + \Delta(\mathbf{x}_k, \mathbf{y}) \right), \quad (7)$$

where k indexes the iteration. Substituting the MPSA module (5), the update becomes:

$$\mathbf{x}_{k+1} = \left(\mathbf{I} - \eta \sum_{i=1}^N \mathbf{W}_i^* \mathbf{W}_i \right) \mathbf{x}_k - \eta \text{MPSA}(\mathbf{x}_k \mid \{\mathbf{W}_i\}, \{\boldsymbol{\xi}_i\}, \gamma, \mu) - \eta \Delta_\phi(\mathbf{x}_k, \mathbf{y}). \quad (8)$$

In practice, the forward degradation model is often unknown, so Δ cannot be computed explicitly. We replace it with a learnable data-consistency term $\Delta_\phi(\mathbf{x}_k, \mathbf{y})$, parameterized by ϕ . This module captures the discrepancy between \mathbf{x}_k and \mathbf{y} while interacting with $\{\boldsymbol{\xi}_i\}$ and $\{\mathbf{W}_i\}$.

Unrolling K iterations of this process yields a deep cascaded network alternating between prompt-driven attention blocks and learnable data-consistency modules.

3 EXPERIMENTS

We evaluate our WBPT against both general-purpose restoration methods and specialized All-in-One approaches (Table 1). Averaged across tasks, WBPT raises the mean PSNR from 30.16 to 31.02,dB, narrowing the gap to PromptIR while maintaining transparency and controllability.

To explicitly assess the effect of multi-scale processing, we further introduce WBPT^\dagger as a ‘‘plus’’ variant: it augments the fully white-box single-scale WBPT with an additional pyramid pathway. Within this pathway, the feature interactions are still governed by our STV-consistent white-box attention blocks, while the down/up-sampling operators for cross-scale aggregation are implemented as a learned, black-box module. Thus, WBPT^\dagger retains a white-box core but uses a partially black-box multi-scale extension. A fully differentiable white-box pyramid is an active direction of ongoing work. Empirically, WBPT^\dagger matches PromptIR on average and exceeds it on selected tasks.

Task-level results further highlight the advantages of the proposed framework. On Rain100L (de-raining), WBPT^\dagger surpasses PromptIR, and Fig. 2 confirms removal of rain streaks with diverse

Table 1: Comparison in the All-in-One restoration setting. Results are reported as PSNR/SSIM. Within each block (single-scale vs. multi-scale), the best and second-best are **boldfaced** and underlined, and gray shading indicates **models whose core attention and reconstruction dynamics are white-box**. WBPT^\dagger shares the same white-box core as WBPT but incorporates a learned, black-box pyramid pathway for multi-scale aggregation.. WBPF yields a marked improvement over WBT, while WBPF^\dagger achieves performance comparable to PromptIR and surpasses it on several tasks.

Method	Dehazing SOTS	Deraining Rain100L	Denoising (BSD68)			Avg.
			$\sigma=15$	$\sigma=25$	$\sigma=50$	
<i>Single-scale methods</i>						
BRDNet	23.23/0.895	27.42/0.895	32.26/0.898	29.76/0.836	26.34/ 0.836	27.80/0.843
FDGAN	24.71/0.924	29.89/0.933	30.25/0.910	28.81/0.868	26.43/0.776	28.02/0.883
AirNet	27.94/0.962	34.90/0.967	33.92/0.933	31.26/0.888	28.00/0.797	31.20/0.910
WBT	27.40/0.958	32.13/0.940	33.17/0.923	30.68/0.875	27.41/0.770	30.16/0.893
WBPT	29.31/0.972	35.93/0.971	<u>33.66/0.929</u>	<u>31.01/0.881</u>	<u>27.72/0.781</u>	<u>31.02/0.907</u>
<i>Multi-scale methods</i>						
LPNet	20.84/0.828	24.88/0.784	26.47/0.778	24.77/0.748	21.26/0.552	23.64/0.738
MPRNet	25.28/0.954	33.57/0.954	33.54/0.927	30.89/0.880	27.56/0.779	30.17/0.899
DL	26.92/0.391	32.62/0.931	33.05/0.914	30.41/0.861	26.90/0.740	29.98/0.875
PromptIR	30.58/0.974	<u>36.37/0.972</u>	<u>33.98/0.933</u>	<u>31.31/0.888</u>	<u>28.06/0.799</u>	32.06/0.913
WBPT^\dagger	<u>29.94/0.970</u>	37.08/0.974	33.86/0.934	31.28/0.890	28.08/0.801	<u>32.05/0.914</u>

orientations more thoroughly than WBPT, producing cleaner outputs. On SOTS (dehazing), while WBPT^\dagger falls short of PromptIR in PSNR, it demonstrates the benefit of multi-scale modeling; as shown in Fig. 3, haze removal is clearer and scene details are better preserved. On BSD68 (denoising), WBPT^\dagger achieves PSNR comparable to PromptIR and yields consistently higher SSIM.

Datasets. For denoising, training is conducted on BSD400 (Arbelaez et al., 2010) and WED (Ma et al., 2016) with Gaussian noise levels $\sigma \in \{15, 25, 50\}$, and evaluation is performed on BSD68 (Martin et al., 2001) and Urban100 (Huang et al., 2015). For deraining, we use Rain100L (200 training / 100 test images) (Yang et al., 2020). For dehazing, training is on SOTS (72,135 images) and evaluation on SOTS (500 images) (Li et al., 2018). In the All-in-One setting, these datasets are combined to train a unified model, following the protocol of (Potlapalli et al., 2023).

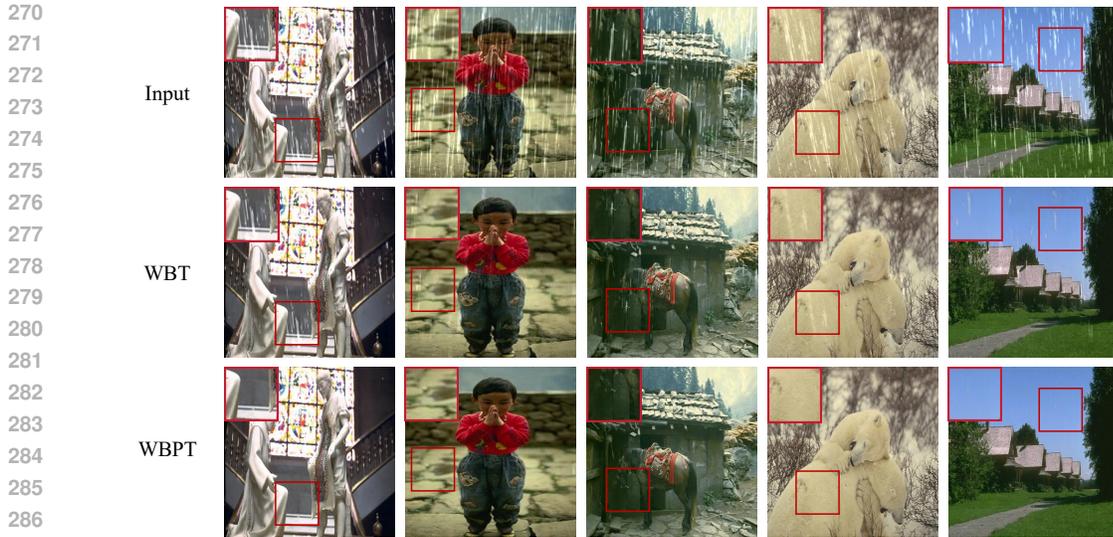
Model and Training. WBPT is trained end-to-end using a 10-iteration white-box framework, where each iteration integrates a learnable gradient update with a Transformer-based prompt branch. Prompts are injected specifically at the sixth Transformer block in each iteration. Training is performed with the standard Adam optimizer ($\beta_1=0.9$, $\beta_2=0.999$) at a fixed learning rate of 1×10^{-4} for 120 epochs in total. We use random 128×128 crops with rotations and flips, optimizing with L2 (MSE) loss. The best checkpoint is selected based on validation performance.

3.1 MULTIPLE DEGRADATION ALL-IN-ONE RESULTS

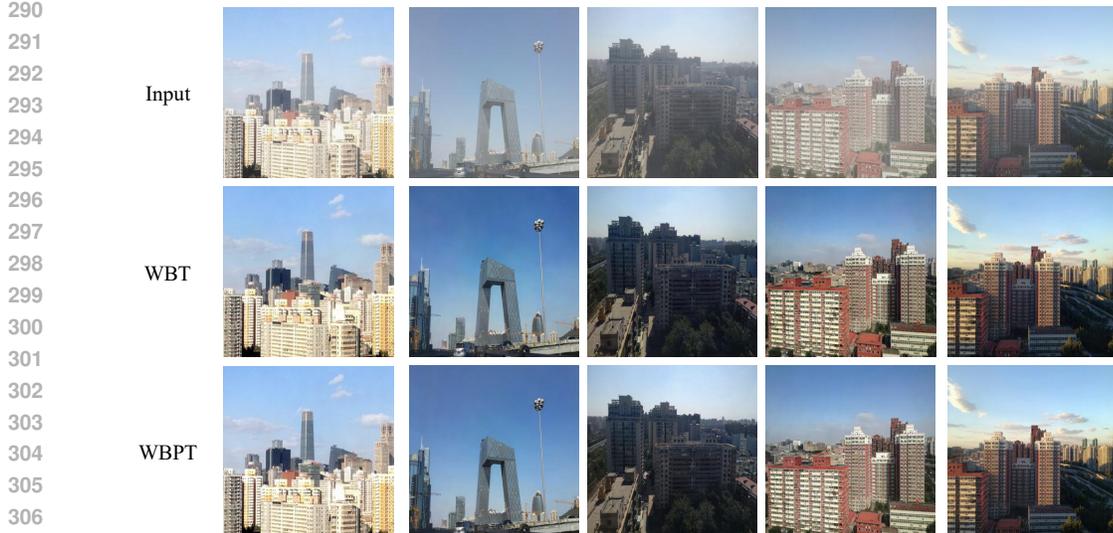
We compare our white-box models with general-purpose restoration approaches and specialized All-in-One methods (Table 1). Averaged across tasks, WBPT raises the mean PSNR from 30.16 to 31.02 dB, narrowing the gap to PromptIR while preserving transparency and controllability. Moreover, our multi-scale extension WBPT^\dagger , which augments the white-box WBPT core with a learned pyramid aggregator, performs on par with PromptIR and even surpasses it on certain tasks. On Rain100L for deraining, WBPT^\dagger outperforms PromptIR; visual comparisons in Fig. 2 show that, relative to WBT, WBPT more effectively removes rain streaks of diverse orientations, yielding cleaner rain-free results. On SOTS for dehazing, although WBPT^\dagger does not surpass PromptIR, it confirms the benefits of the multi-scale design; examples in Fig. 3 indicate clearer haze removal and more faithful scene restoration. On BSD68 for denoising, WBPT^\dagger attains PSNR comparable to PromptIR while overall delivering higher SSIM.

3.2 SINGLE DEGRADATION ONE-BY-ONE RESULTS

We evaluate PromptIR under the single-task setting, where a separate model is trained for each restoration task. This setting is intended to empirically verify that content-adaptive prompting via the prompt block is also effective for single-task networks. Table 2 reports the deraining results on standard datasets: our single-scale white-box WBPT consistently achieves the best performance,



287 Figure 2: Deraining results under the All-in-One setting. Compared with WBT, our WBPT re-
288 moves numerous residual rain streaks that WBT fails to eliminate, yielding cleaner backgrounds
289 and sharper details (see red zoom-in boxes).



307 Figure 3: Dehazing results for all-in-one methods. Compared with WBT, our WBPT recovers clearer
308 sky regions and sharper building edges, suppresses veiling glare and color cast, and yields more
309 natural contrast and details across diverse urban scenes.

310 surpassing PromptIR and [the multi-scale extension WBPT[†] with a learned pyramid aggregator](#); re-
311 lative to WBT (without prompts), WBPT delivers a 1.93 dB gain in PSNR. For dehazing and denois-
312 ing, although WBPT[†] does not surpass PromptIR, it achieves comparable performance; see Tables 3
313 and 4.

315 3.3 ATTENTION VISUALIZATION

316 To examine differences in model focus, we visualize the last-layer multi-head attention of the white-
317 box reconstruction model (WBPT) and the black-box method PromptIR. For each model, we extract
318 the last-layer attention tensor $A \in \mathbb{R}^{H \times N \times N}$, average across heads, and further aggregate along the
319 query dimension to obtain a single-channel response for each window. The window-wise responses
320 are then reassembled into a full-image heatmap via reverse window stitching. To ensure a fair
321 comparison, both models use identical inputs and visualization settings. As shown in Fig. 5, WBPT
322 exhibits strong responses on object boundaries and structural regions, indicating a preference for
323 image geometry and semantic content rather than directly following degradation textures; notably,
this boundary-centric attention aligns with our STV objective (Sec. 2.1). In contrast, PromptIR's

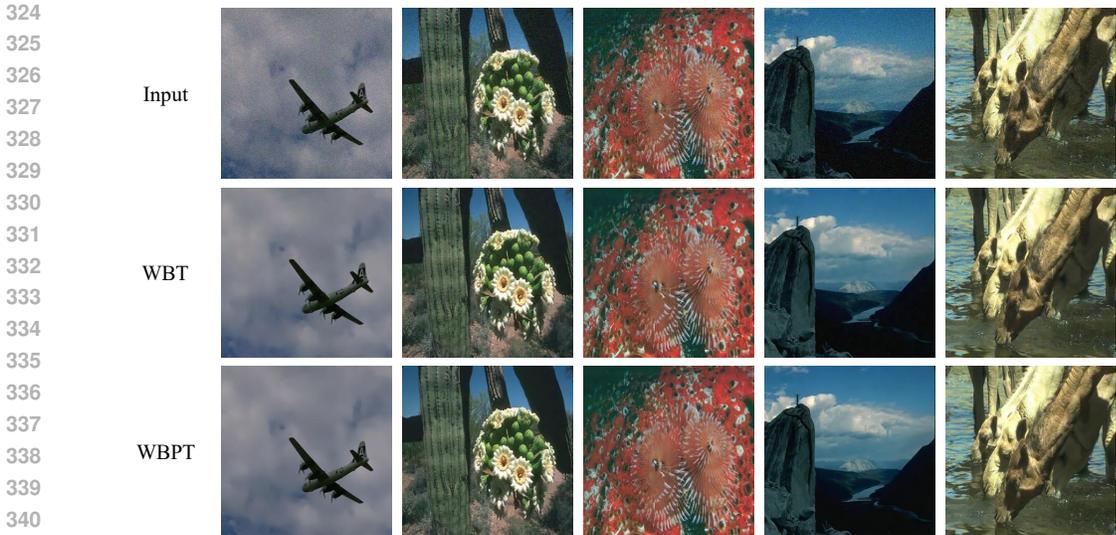


Figure 4: Denoising results for all-in-one methods.

Table 2: Deraining results on Rain100L in the single-task setting. Within each scale group (Single or Multi), best results are **boldfaced** and second-best are underlined; Gray indicates white-box models.

Scale	<i>Single-scale methods</i>				<i>Multi-scale methods</i>				
Method	SIRR	AirNet	WBT	WBPT	MSPFN	LPNet	Restormer	PromptIR	WBPT [†]
PSNR	32.37	34.90	<u>36.77</u>	38.70	33.50	33.61	36.74	<u>37.04</u>	38.54
SSIM	0.926	0.977	<u>0.977</u>	0.983	0.948	0.958	0.978	<u>0.979</u>	0.984

responses align with rain streaks and are more tightly coupled to the degradation pattern, suggesting a greater reliance on degradation-pattern detection.

3.4 T-SNE ANALYSIS OF INTERMEDIATE REPRESENTATIONS

To analyze the intermediate representations of an all-in-one model across different degradations, we tap the *input* to the final convolution layer of the Transformer backbone during the forward pass. For each image, we apply global average pooling over the spatial dimensions to obtain a channel-wise embedding vector. We collect embeddings from three standard test sets—BSD68 denoising ($\sigma = 25$), Rain100L deraining, and SOTS-Outdoor dehazing—and project them to 2D using t-SNE.

Figure 6 compares the black-box PromptIR with our white-box WBPT under identical preprocessing and t-SNE settings: PromptIR (left) yields highly entangled embeddings with substantial cross-task overlap, whereas WBPT (right) forms well-separated clusters for the Noisy, Hazy, and Rainy samples, exhibiting tighter intra-cluster compactness and clearer inter-cluster margins. These results indicate that WBPT learns more discriminative, task-aware representations in the all-in-one setting.

To further examine whether the learned representations remain degradation-aware beyond single degradations, we additionally conduct t-SNE analysis on the CDD11 mixed-degradation dataset. CDD11 contains several two-way combinations of degradations; we select three representative protocols (low+haze, low+snow, haze+snow) and, following exactly the same feature-extraction protocol as above, compute embeddings for PromptIR and WBPT (Figure 7). While the three protocols share degradation components pairwise, the PromptIR embeddings still exhibit substantial cross-class overlap, whereas WBPT forms more compact clusters with clearer boundaries between mixed degradations. This suggests that our white-box design continues to organize intermediate features primarily according to degradation combinations, rather than image content, even in more complex mixed-degradation scenarios.

3.5 STABILITY UNDER PROMPT-PARAMETER PERTURBATIONS

To verify the stability of WBPT under prompt-parameter perturbations, we conduct a perturbation-sensitivity study in a controlled experimental setting and compare it with PromptIR. The test datasets are BSD68 (denoising), Rain100L (deraining), and SOTS-Outdoor (dehazing). We inject additive Gaussian noise *only* into the prompt parameters ($\sigma \in [0.001, 0.1]$), while keeping all other settings

Table 3: Dehazing results on SOTS dataset in the single-task setting. Within each block (single-scale vs. multi-scale), the best and second-best are **boldfaced** and underlined, and gray shading indicates white-box models. PSNR/SSIM reported; higher is better.

Scale	Single-scale methods					Multi-scale methods			
Method	AODNet	FDGAN	AirNet	WBT	WBPT	EPDN	Restormer	PromptIR	WBPT [†]
PSNR	20.29	23.15	23.18	28.72	<u>28.33</u>	22.57	<u>30.87</u>	31.31	30.47
SSIM	0.877	0.921	0.900	<u>0.961</u>	0.967	0.863	0.969	0.973	<u>0.972</u>

Table 4: Denoising comparisons in the single-task setting on BSD68 and Urban100. Results are reported as PSNR/SSIM. Within each block (single-scale vs. multi-scale), the best and second-best are **boldfaced** and underlined, respectively. gray shading indicates white-box models. At the challenging noise level of $\sigma = 50$ on Urban100, our WBPT achieves a 0.39 dB improvement over WBT. Meanwhile, WBPT[†] attains performance comparable to PromptIR.

Method	BSD68			Urban100		
	$\sigma=15$	$\sigma=25$	$\sigma=50$	$\sigma=15$	$\sigma=25$	$\sigma=50$
<i>Single-scale methods</i>						
CBM3D	33.50/0.922	30.69/0.868	27.36/0.763	33.93/0.941	31.36/0.909	27.93/0.840
DnCNN	33.89/0.930	31.23/0.883	27.92/0.789	32.98/0.931	30.81/0.902	27.59/0.833
IRCNN	33.87/0.929	31.18/0.882	27.88/0.790	27.59/0.833	31.20/0.909	27.70/0.840
FFDNet	33.87/0.929	31.21/0.882	27.96/0.789	33.83/0.942	31.40/0.912	28.05/0.848
BRDNet	34.10/0.929	31.43/0.885	28.16/0.794	34.42/0.946	31.99/0.919	28.56/0.858
AirNet	34.14/0.936	31.48/0.893	28.23/0.806	34.40/0.949	<u>32.10/0.924</u>	28.88/0.871
WBT	33.59/0.930	30.92/0.882	27.85/0.793	<u>33.43/0.956</u>	30.42/0.924	27.18/0.858
WBPT	34.02/ <u>0.935</u>	31.35/ <u>0.891</u>	28.03/ <u>0.797</u>	34.15/ 0.963	31.38/ 0.937	27.57/ <u>0.870</u>
<i>Multi-scale methods</i>						
Restormer	<u>34.29/0.937</u>	31.64/0.895	<u>28.41/0.810</u>	34.67/ 0.969	<u>32.41/0.927</u>	<u>29.31/0.878</u>
PromptIR	34.34/0.938	31.71/0.897	28.49/0.813	34.77/0.952	32.49/0.929	29.39/0.881
WBPT [†]	<u>34.31/0.938</u>	<u>31.60/0.895</u>	<u>28.36/0.811</u>	<u>34.76/0.952</u>	<u>32.27/0.927</u>	29.08/0.877

Table 5: Evaluation of stability under prompt-parameter perturbations, reported as relative drops in PSNR and SSIM (lower is better). Gaussian noise with $\sigma \in [0.001, 0.1]$ is injected exclusively into the prompt parameters. Results are averaged over multiple severity levels on BSD68, Rain100L, and SOTS-Outdoor. WBPT exhibits smaller drops than PromptIR, indicating greater stability.

Model	Denoising		Deraining		Dehazing		Average	
	PSNR	SSIM	PSNR	SSIM	PSNR	SSIM	PSNR	SSIM
PromptIR	-10.2%	-13.0%	-13.8%	-12.4%	-13.0%	-12.2%	-12.3%	-12.5%
WBPT	-3.05%	-0.35%	-1.02%	-0.64%	-2.18%	-0.05%	-2.08%	-0.31%

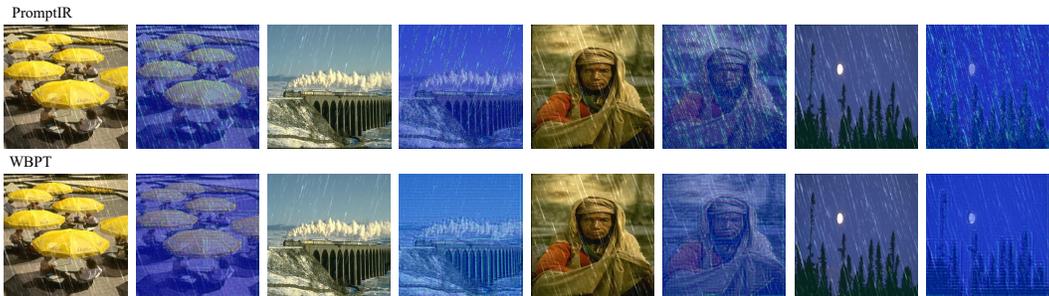
(e.g., the prompt insertion layer) identical to the previous configuration. The evaluation metric is the *average performance drop percentage* (lower indicates higher stability), averaged over multiple perturbation severities and the three datasets. The corresponding averages are summarized in Table 5; a representative qualitative comparison at $\sigma = 0.1$ is shown in Fig. 8.

From the visual results, PromptIR consistently exhibits *systematic contrast and color shifts* after perturbing the prompt, suggesting an undesirable coupling between the prompt representation and global imaging attributes. Under prompt-only perturbations, such global tone/contrast changes are not what degradation awareness is expected to primarily induce. In contrast, WBPT with $\sigma = 0.1$ still removes rain effectively while maintaining remarkably stable contrast and colors, indicating stronger robustness and better degradation–prompt decoupling.

3.6 GUIDANCE MODALITY VALIDATION: SOFT PROMPT VS REAL GUIDANCE

We compare a learnable *soft prompt* with a proxy of real guidance (*hard*: image gradients \rightarrow edge map plus Gaussian noise, $\sigma \in \{0.01, 0.02\}$). To control compute and isolate the modality effect, the *backbone is frozen* and only the prompt and fusion parameters are finetuned. For each test image, we report the paired difference $\Delta = \text{metric}_{\text{soft}} - \text{metric}_{\text{hard}}$; our goal is to assess the *relative gap be-*

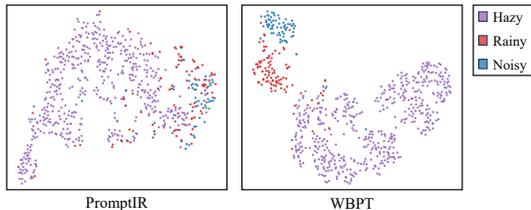
432
433
434
435
436
437
438
439
440
441



442
443
444
445
446

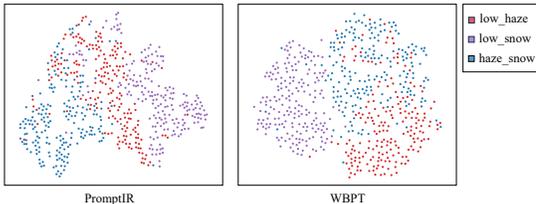
Figure 5: Attention-map visualizations for WBPT and PromptIR. Odd-numbered columns show input images; even-numbered columns show the corresponding attention maps. Top row: PromptIR; bottom row: WBPT. Attention heads and queries from the final layer are aggregated, and full-image heatmaps are reconstructed via reverse window stitching. WBPT focuses on object boundaries and main structures, whereas PromptIR emphasizes rain streaks.

447
448
449
450
451
452



453
454
455
456
457
458
459
460
461

Figure 6: t-SNE visualization of degradation embeddings from PromptIR and WBPT. Colors denote degradation types. With identical inputs and t-SNE settings, WBPT yields clearly separable clusters by degradation type, whereas PromptIR is more entangled.



462
463
464
465
466
467
468
469
470

tween soft and hard rather than absolute gains. Equivalence margins are pre-registered as ± 0.02 dB (PSNR) and ± 0.002 (SSIM), within which soft and hard are deemed practically equivalent.

Under the finetune-only setting (backbone frozen), soft and hard guidance behave nearly identically across denoising, deraining, and dehazing in our controlled evaluations. The paired differences $\Delta = \text{soft} - \text{hard}$ are consistently tiny and remain within the pre-registered equivalence margins (± 0.02 dB PSNR / ± 0.002 SSIM) for both $\sigma = 0.01$ and $\sigma = 0.02$. While dehazing yields lower absolute scores—reflecting its higher difficulty—the relative gap between soft and hard stays stable, which is precisely the comparison this experiment aims to isolate.

471
472
473
474
475
476
477

Table 6: Soft vs. hard guidance under the finetune-only setting (backbone frozen). Metrics are PSNR and SSIM in separate columns; parentheses denote the change relative to the baseline in the same column. Gray denote the paired difference $\Delta = \text{soft} - \text{hard}$; values within the pre-registered equivalence margins (± 0.10 dB PSNR, ± 0.002 SSIM) indicate practical equivalence. Results are reported for $\sigma \in \{0.01, 0.02\}$. While absolute scores for dehazing are lower due to its higher difficulty, the soft–hard gap remains small and stable across σ .

σ	Task	Denoise		Derain		Dehaze	
		PSNR	SSIM	PSNR	SSIM	PSNR	SSIM
0.01	WBT	33.59	0.930	36.77	0.977	28.72	0.961
	WBT+soft	33.61 (+0.02)	0.929 (−0.1%)	36.90 (+0.13)	0.978 (+0.1%)	27.83 (−0.89)	0.956 (−0.5%)
	WBT+hard	33.63 (+0.04)	0.929 (−0.1%)	36.89 (+0.12)	0.978 (+0.1%)	27.87 (−0.85)	0.958 (−0.3%)
	Δ (soft–hard)	−0.02	0.000	+0.01	0.000	−0.04	−0.002
0.02	WBT	33.59	0.930	36.77	0.977	28.72	0.961
	WBT+soft	33.62 (+0.03)	0.930 (+0.0%)	36.89 (+0.12)	0.978 (+0.1%)	27.83 (−0.89)	0.956 (−0.5%)
	WBT+hard	33.60 (+0.01)	0.929 (−0.1%)	36.89 (+0.12)	0.978 (+0.1%)	27.86 (−0.86)	0.958 (−0.3%)
	Δ (soft–hard)	+0.02	0.000	+0.00	0.000	−0.03	−0.002

485



496 Figure 8: Qualitative comparison under Gaussian perturbation of prompt parameters ($\sigma = 0.1$).
 497 Odd-numbered columns show input images; even-numbered columns show restored results. Top
 498 row: PromptIR; bottom row: WBPT. WBPT preserves deraining quality, contrast, and colors,
 499 whereas PromptIR exhibits noticeable shifts, indicating lower robustness.

500 4 CONCLUSION

502 In this work, we revisited prompt-based Transformers from a variational perspective and established
 503 a principled connection between prompts and attention. By casting guided image restoration as a
 504 *STV* problem, we derived a white-box attention mechanism that offers an interpretable foundation
 505 for prompt–attention coupling. Building on this formulation, we unrolled the gradient flow into
 506 WBPT, a cascaded architecture that integrates variational principles with modern prompt learning.
 507 Extensive experiments on the classic three-degradation benchmark (denoising, deraining, dehazing)
 508 demonstrate that WBPT delivers competitive performance while maintaining transparent and robust
 509 prompt–attention dynamics. Beyond empirical performance, WBPT is complementary to recent all-
 510 in-one restorers such as Perceive-IR, DA-RCOT, MoCE-IR, AdaIR and DFPIR, which push state-
 511 of-the-art results via stronger backbones and degradation-aware modules. Our *STV*-based atten-
 512 tion block can serve as an interpretable replacement or constraint for prompt/attention submodules
 513 inside these powerful architectures, or be combined with their semantic-, frequency- and feature-
 514 perturbation designs, suggesting “white-box mechanism + strong backbone” hybrids as a promising
 515 direction for future all-in-one restoration. Extending WBPT from three-degradation settings to more
 516 diverse protocols, such as five-degradation, mixed-degradation and real-world benchmarks, is an im-
 517 portant next step toward fully exploiting this connection in practice.

518
519
520
521
522
523
524
525
526
527
528
529
530
531
532
533
534
535
536
537
538
539

REFERENCES

- 540
541
542 Jonas Adler and Ozan Öktem. Learned primal-dual reconstruction. *IEEE Transactions on Medical*
543 *Imaging*, 37(6):1322–1332, 2018.
- 544 Hemant K Aggarwal, Merry P Mani, and Mathews Jacob. Modl: model-based deep learning archi-
545 tecture for inverse problems. *IEEE Transactions on Medical Imaging*, 38(2):394–405, 2019.
- 546 Pablo Arbelaez, Michael Maire, Charless Fowlkes, and Jitendra Malik. Contour detection and hier-
547 archical image segmentation. *IEEE Transactions on Pattern Analysis and Machine Intelligence*,
548 33(5):898–916, 2010.
- 549 Antonin Chambolle and Thomas Pock. A first-order primal-dual algorithm for convex problems with
550 applications to imaging. *Journal of Mathematical Imaging and Vision*, 40(1):120–145, 2011.
- 551 Hila Chefer, Shir Gur, and Lior Wolf. Transformer interpretability beyond attention visualization.
552 In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp.
553 782–791, 2021.
- 554 Hanjing Chen, Yunhe Wang, Tianyu Guo, Chang Xu, Yiping Deng, Zhenhua Liu, Siwei Ma, Chun-
555 jing Xu, Chao Xu, and Wen Gao. Pre-trained image processing transformer. In *Proceedings of*
556 *the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 12299–12310, 2021.
- 557 Liangyu Chen, Xiaojie Chu, Xiangyu Zhang, and Jian Sun. Simple baselines for image restoration.
558 In *European Conference on Computer Vision*, pp. 17–33. Springer, 2022.
- 559 Yunjin Chen, Wei Yu, and Thomas Pock. On learning optimized reaction diffusion processes for
560 effective image restoration. In *Proceedings of the IEEE Conference on Computer Vision and*
561 *Pattern Recognition*, pp. 5261–5269, 2015.
- 562 Hyungjin Chung, Byeonghu Sim, Minyong Ryu, and Jong Chul Ye. Diffusion posterior sampling
563 for general noisy inverse problems. In *International Conference on Learning Representations*,
564 2023.
- 565 Yuning Cui, Syed Waqas Zamir, Salman Khan, Alois Knoll, Mubarak Shah, and Fahad Shahbaz
566 Khan. AdaIR: Adaptive all-in-one image restoration via frequency mining and modulation. In
567 *International Conference on Learning Representations (ICLR)*, 2025.
- 568 Matthias J Ehrhardt and Marta M Betcke. Multicontrast MRI reconstruction with structure-guided
569 total variation. *SIAM Journal on Imaging Sciences*, 9(3):1084–1106, 2016a.
- 570 Matthias J Ehrhardt and Marta M Betcke. Multicontrast MRI reconstruction with structure-guided
571 total variation. *SIAM Journal on Imaging Sciences*, 9(3):1084–1106, 2016b.
- 572 Lawrence C Evans. *Partial differential equations*, volume 19. American Mathematical Society,
573 2022.
- 574 Chun-Mei Feng, Huazhu Fu, Tianfei Zhou, Yong Xu, Ling Shao, and David Zhang. Deep multi-
575 modal aggregation network for MR image reconstruction with auxiliary modality. *arXiv Preprint*
576 *arXiv:2110.08080*, 2021.
- 577 Kerstin Hammernik, Teresa Klatzer, Erich Kobler, Michael P Recht, Daniel K Sodickson, Thomas
578 Pock, and Florian Knoll. Learning a variational network for reconstruction of accelerated MRI
579 data. *Magnetic Resonance in Medicine*, 79(6):3055–3071, 2018.
- 580 Mohammad Havaei, Nicolas Guizard, Nicolas Chapados, and Yoshua Bengio. Hemis: hetero-modal
581 image segmentation. In *International Conference on Medical Image Computing and Computer-*
582 *Assisted Intervention*, pp. 469–477. Springer, 2016.
- 583 Kaiming He, Jian Sun, and Xiaoou Tang. Guided image filtering. *IEEE Transactions on Pattern*
584 *Analysis and Machine Intelligence*, 35(6):1397–1409, 2012.
- 585 Jia-Bin Huang, Abhishek Singh, and Narendra Ahuja. Single image super-resolution from trans-
586 formed self-exemplars. In *Proceedings of the IEEE Conference on Computer Vision and Pattern*
587 *Recognition*, pp. 5197–5206, 2015.
- 588
589
590
591
592
593

- 594 Sarthak Jain and Byron C Wallace. Attention is not explanation. *arXiv Preprint arXiv:1902.10186*,
595 2019.
- 596
- 597 Menglin Jia, Luming Tang, Bor-Chun Chen, Claire Cardie, Serge Belongie, Bharath Hariharan, and
598 Ser-Nam Lim. Visual prompt tuning. In *European Conference on Computer Vision*, pp. 709–727.
599 Springer, 2022.
- 600
- 601 Junjun Jiang, Zengyuan Zuo, Gang Wu, Kui Jiang, and Xianming Liu. A survey on all-in-one image
602 restoration: Taxonomy, evaluation and future trends. *IEEE Transactions on Pattern Analysis and*
603 *Machine Intelligence*, 2025.
- 604
- 605 Ulugbek S Kamilov, Charles A Bouman, Gregory T Buzzard, and Brendt Wohlberg. Plug-and-play
606 methods for integrating physical and learned models in computational imaging: theory, algo-
607 rithms, and applications. *IEEE Signal Processing Magazine*, 40(1):85–97, 2023.
- 608
- 609 Yonatan Kawar, Michael Elad, Tomer Michaeli, and Stefano Ermon. Denoising diffusion restoration
610 models. In *Advances in Neural Information Processing Systems*, 2022.
- 611
- 612 Dehong Kong, Fan Li, Zhixin Wang, Jiaqi Xu, Renjing Pei, Wenbo Li, and Wenqi Ren. Dual
613 prompting image restoration with diffusion transformers. In *Proceedings of the IEEE/CVF Con-*
614 *ference on Computer Vision and Pattern Recognition*, pp. 12809–12819, 2025.
- 615
- 616 Stamatios Lefkimmiatis, Anastasios Roussos, Petros Maragos, and Michael Unser. Structure tensor
617 total variation. *SIAM Journal on Imaging Sciences*, 8(2):1090–1122, 2015.
- 618
- 619 Boyi Li, Wenqi Ren, Dengpan Fu, Dacheng Tao, Dan Feng, Wenjun Zeng, and Zhangyang Wang.
620 Benchmarking single-image dehazing and beyond. *IEEE Transactions on Image Processing*, 28
621 (1):492–505, 2018.
- 622
- 623 Boyun Li, Xiao Liu, Peng Hu, Zhongqin Wu, Jiancheng Lv, and Xi Peng. All-in-one image restora-
624 tion for unknown corruption. In *Proceedings of the IEEE/CVF Conference on Computer Vision*
625 *and Pattern Recognition*, pp. 17452–17462, 2022.
- 626
- 627 Xiang Lisa Li and Percy Liang. Prefix-tuning: optimizing continuous prompts for generation. *arXiv*
628 *Preprint arXiv:2101.00190*, 2021.
- 629
- 630 Yijun Li, Jia-Bin Huang, Narendra Ahuja, and Ming-Hsuan Yang. Deep joint image filtering. In
631 *European Conference on Computer Vision*, pp. 154–169. Springer, 2016.
- 632
- 633 Jingyun Liang, Jiezhong Cao, Guolei Sun, Kai Zhang, Luc Van Gool, and Radu Timofte. Swinir:
634 image restoration using swin transformer. In *Proceedings of the IEEE/CVF International Confer-*
635 *ence on Computer Vision*, pp. 1833–1844, 2021.
- 636
- 637 Wenyang Luo, Haina Qin, Zewen Chen, Libin Wang, Dandan Zheng, Yuming Li, Yufan Liu, Bing
638 Li, and Weiming Hu. Visual-instructed degradation diffusion for all-in-one image restoration. In
639 *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*,
640 2025.
- 641
- 642 Michael Lustig, David Donoho, and John M. Pauly. Sparse MRI: The application of compressed
643 sensing for rapid MR imaging. *Magnetic Resonance in Medicine*, 2007.
- 644
- 645 Kede Ma, Zhengfang Duanmu, Qingbo Wu, Zhou Wang, Hongwei Yong, Hongliang Li, and Lei
646 Zhang. Waterloo exploration database: new challenges for image quality assessment models.
647 *IEEE Transactions on Image Processing*, 26(2):1004–1016, 2016.
- 648
- 649 David R. Martin, Charless C. Fowlkes, Doron Tal, and Jitendra Malik. A database of human seg-
650 mented natural images and its application to evaluating segmentation algorithms and measuring
651 ecological statistics. In *Proceedings of the Eighth IEEE International Conference on Computer*
652 *Vision (ICCV 2001)*, volume 2, pp. 416–423. IEEE, 2001.
- 653
- 654 Junying Meng, Faqiang Wang, and Jun Liu. Learnable nonlocal self-similarity of deep features for
655 image denoising. *SIAM Journal on Imaging Sciences*, 17(1):441–475, 2024.

- 648 Vishal Monga, Yuelong Li, and Yonina C Eldar. Algorithm unrolling: interpretable, efficient deep
649 learning for signal and image processing. *IEEE Signal Processing Magazine*, 38(2):18–44, 2021.
650
- 651 Vaishnav Potlapalli, Syed Waqas Zamir, Salman H Khan, and Fahad Shahbaz Khan. Promptir:
652 prompting for all-in-one image restoration. In *Advances in Neural Information Processing Sys-*
653 *tems*, volume 36, pp. 71275–71293, 2023.
- 654 Klaas P. Pruessmann, Markus Weiger, Michael B. Scheidegger, and Peter Boesiger. Sense: Sensi-
655 tivity encoding for fast MRI. *Magnetic Resonance in Medicine*, 1999.
656
- 657 Cynthia Rudin. Stop explaining black box machine learning models for high stakes decisions and
658 use interpretable models instead. *Nature Machine Intelligence*, 1(5):206–215, 2019.
659
- 660 Chitwan Saharia, William Chan, Huiwen Chang, Chris Lee, Jonathan Ho, Tim Salimans, David J.
661 Fleet, and Mohammad Norouzi. Palette: Image-to-image diffusion models. In *Proceedings of the*
662 *ACM SIGGRAPH*, 2022.
- 663 Jo Schlemper, Jose Caballero, Joseph V. Hajnal, Anthony Price, and Daniel Rueckert. A deep cas-
664 cade of convolutional neural networks for dynamic MR image reconstruction. *IEEE Transactions*
665 *on Medical Imaging*, 2018.
666
- 667 Jian Sun, Huibin Li, Zongben Xu, et al. Deep ADMM-net for compressive sensing MRI. In *Ad-*
668 *vances in Neural Information Processing Systems*, volume 29, 2016.
- 669 Xiaole Tang, Xiang Gu, Xiaoyi He, Xin Hu, and Jian Sun. Degradation-aware residual-conditioned
670 optimal transport for unified image restoration. *IEEE Transactions on Pattern Analysis and Ma-*
671 *chine Intelligence*, 2025.
672
- 673 Xiangpeng Tian, Xiangyu Liao, Xiao Liu, Meng Li, and Chao Ren. Degradation-aware feature
674 perturbation for all-in-one image restoration. In *Proceedings of the IEEE/CVF Conference on*
675 *Computer Vision and Pattern Recognition (CVPR)*, 2025.
- 676 Yuchuan Tian, Jianhong Han, Hanting Chen, Yuanyuan Xi, Ning Ding, Jie Hu, Chao Xu, and
677 Yunhe Wang. Instruct-IPT: all-in-one image processing transformer via weight modulation. *arXiv*
678 *Preprint arXiv:2407.00676*, 2024.
679
- 680 Dmitry Ulyanov, Andrea Vedaldi, and Victor Lempitsky. Deep image prior. In *Proceedings of the*
681 *IEEE Conference on Computer Vision and Pattern Recognition*, pp. 9446–9454, 2018.
- 682 Singanallur V Venkatakrishnan, Charles A Bouman, and Brendt Wohlberg. Plug-and-play priors
683 for model-based reconstruction. In *Proceedings of the IEEE Global Conference on Signal and*
684 *Information Processing*, pp. 945–948. IEEE, 2013.
685
- 686 Xiaolong Wang, Ross Girshick, Abhinav Gupta, and Kaiming He. Non-local neural networks. In
687 *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 7794–
688 7803, 2018.
- 689 Zhendong Wang, Xiaodong Cun, Jianmin Bao, Wengang Zhou, Jianzhuang Liu, and Houqiang Li.
690 Uformer: a general U-shaped transformer for image restoration. In *Proceedings of the IEEE/CVF*
691 *Conference on Computer Vision and Pattern Recognition*, pp. 17683–17693, 2022.
692
- 693 Fuzhi Yang, Huan Yang, Jianlong Fu, Hongtao Lu, and Baining Guo. Learning texture transformer
694 network for image super-resolution. In *Proceedings of the IEEE/CVF Conference on Computer*
695 *Vision and Pattern Recognition*, pp. 5791–5800, 2020.
- 696 Haidi Yang, Chuang Zhang, Chun-Mei Feng, Jianfu Zhang, and Huazhu Fu. Multi-modal guidance-
697 based deep unfolding network for MRI reconstruction. In *Proceedings of the 30th ACM Interna-*
698 *tional Conference on Multimedia (ACM MM)*, 2022.
699
- 700 Yaodong Yu, Sam Buchanan, Druv Pai, Tianzhe Chu, Ziyang Wu, Shengbang Tong, Benjamin Ha-
701 effele, and Yi Ma. White-box transformers via sparse rate reduction. In *Advances in Neural*
Information Processing Systems, volume 36, pp. 9422–9457, 2023.

702 Eduard Zamfir, Zongwei Wu, Nancy Mehta, Yuedong Tan, Danda Pani Paudel, Yulun Zhang, and
703 Radu Timofte. Complexity experts are task-discriminative learners for any image restoration. In
704 *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*,
705 2025.

706 Syed Waqas Zamir, Aditya Arora, Salman Khan, Munawar Hayat, Fahad Shahbaz Khan, and Ming-
707 Hsuan Yang. Restormer: efficient transformer for high-resolution image restoration. In *Proceed-*
708 *ings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 5728–5739,
709 2022.

710 Xu Zhang, Jiaqi Ma, Guoli Wang, Qian Zhang, Huan Zhang, and Lefei Zhang. Perceive-IR: Learn-
711 ing to perceive degradation better for all-in-one image restoration. *IEEE Transactions on Image*
712 *Processing*, 2025.

713 Kaiyang Zhou, Jingkang Yang, Chen Change Loy, and Ziwei Liu. Learning to prompt for vision-
714 language models. *International Journal of Computer Vision*, 130(9):2337–2348, 2022.

715 Wei Zhou, Xiaoyu Wang, Qi Zhang, Xiang Li, and Ke Chen. Multi-level modality fusion network
716 for multi-contrast MRI reconstruction. *Computerized Medical Imaging and Graphics*, 2024.

717
718
719
720
721
722
723
724
725
726
727
728
729
730
731
732
733
734
735
736
737
738
739
740
741
742
743
744
745
746
747
748
749
750
751
752
753
754
755

756	APPENDIX	
757		
758	A Related Work	15
759		
760	B Ablation Experiment	16
761		
762	B.1 SPSA Component Ablations	16
763	B.2 Position of prompt blocks.	16
764	B.3 Complementarity between Prompt and Data-Consistency Modules.	17
765	B.4 Accuracy–overhead trade-off for prompt design	17
766		
767		
768	C Variational Derivation of the Energy Functional	17
769		
770	C.1 Variation of the First Term	18
771	C.2 Variation of the Second Term	19
772	C.3 Combining Both Terms	19
773	C.4 Variational Derivative Approximation	20
774		
775		
776	D Limitation	23
777		
778	E Qualitative results	23
779		
780	F Algorithm	23
781		
782	G Transformer and Prompt Blocks in WBPT	24
783		
784	G.1 Prompt Block in WBPT Framework	25
785	G.2 Transformer Block in WBPT Framework	27
786		
787		
788	H Reproducibility Statement	28
789		
790	I The Use of Large Language Models (LLMs)	28
791		
792		

A RELATED WORK

Transformer-based image restoration. Transformer architectures have advanced image restoration by modeling long-range dependencies, with strong single-task systems such as SwinIR, Restormer, and successors (Chen et al., 2021; Liang et al., 2021; Zamir et al., 2022; Chen et al., 2022; Wang et al., 2022). To curb task specialization, all-in-one models handle multiple degradations with a unified backbone (e.g., AirNet (Li et al., 2022)). More recent all-in-one restoration works further expand this line by systematizing unified benchmarks and architectures, including a comprehensive AiOIR survey (Jiang et al., 2025) and degradation-aware backbones such as DA-RCOT (Tang et al., 2025), MoCE-IR (Zamfir et al., 2025), and AdaIR (Cui et al., 2025). We likewise pursue unified restoration but differ in how conditioning is defined and used within the network. Parallel to Transformer priors, diffusion-based restoration has recently shown competitive performance and broad applicability (Kawar et al., 2022; Chung et al., 2023; Saharia et al., 2022), and our formulation is complementary: it can inject structure-aware conditioning regardless of the underlying prior family.

Prompt-based conditioning for restoration. Prompt-based Transformers inject task or condition cues (e.g., degradation type, maps, or control signals) into a shared backbone to enable multi-task

810 adaptation (Potlapalli et al., 2023; Li et al., 2022; Tian et al., 2024; Kong et al., 2025). On top of
 811 these designs, Perceive-IR (Zhang et al., 2025) and DFPIR (Tian et al., 2025) further improve all-
 812 in-one performance via quality-aware prompt learning and degradation-aware feature perturbations,
 813 while Defusion (Luo et al., 2025) leverages diffusion priors with degradation instructions for unified
 814 restoration. Prevailing designs treat prompts as black-box tokens or channel-wise modulations (e.g.,
 815 visual prompt tuning/adapters) that are concatenated to features and learned end-to-end, which limits
 816 interpretability and spatial control (Jia et al., 2022). In contrast, we treat prompts as **structural**
 817 **priors** and derive structure-aware prompt attention from a variational formulation via **learnable**
 818 **regularization gradients (LRG)**, providing mechanistic interpretability and spatial controllability
 819 (Yu et al., 2023).

820
 821 **Model-based deep learning and unrolled networks.** Model-based approaches explicitly couple
 822 data fidelity with learned priors through algorithm unrolling (e.g., ADMM-Net, MoDL, VarNet)
 823 (Sun et al., 2016; Aggarwal et al., 2019; Hammernik et al., 2018; Adler & Öktem, 2018). Plug-
 824 and-play and RED families further connect optimization with learned denoisers (Venkatakrisnan
 825 et al., 2013; Yang et al., 2022; Ulyanov et al., 2018; Kamilov et al., 2023). Our design follows this
 826 lineage: we unfold a variational objective, keep a physics-consistency step, and introduce a **learn-**
 827 **able data-consistency (LDC)** module that complements the physics-consistency step to suppress
 828 residual artifacts. A key difference is that our conditioning instead arises from a regularizer-driven
 829 decomposition and directly parameterizes attention via LRG.

830
 831 **Structure-guided and cross-modal reconstruction (MRI).** Guided reconstruction leverages side
 832 information to align the target with structures visible in a guide modality; in MRI, T1 can guide T2
 833 reconstruction under multi-contrast or cross-modal priors (Ehrhardt & Betcke, 2016b; Feng et al.,
 834 2021; Yang et al., 2022; Zhou et al., 2024). This line builds upon classical physics-consistent MRI,
 835 including parallel imaging and compressed sensing (Pruessmann et al., 1999; Lustig et al., 2007),
 836 and deep cascades that interleave learned priors with data consistency (Schlemper et al., 2018). We
 837 reinterpret guidance as prompting within a white-box formulation: the prompt enters the variational
 838 gradient as a structure-aligned term that controls attention, linking classical guided reconstruction
 839 with modern prompt-based Transformers.

840
 841 **Summary.** Compared to black-box prompt injection, our **White-Box Prompt Transformer**
 842 **(WBPT)** provides a variationally grounded route to structure-aware attention (via LRG) while main-
 843 taining faithful reconstruction through an LDC-augmented fidelity step, unifying multi-task restora-
 844 tion and guided MRI within the same unfolded architecture.

845 846 B ABLATION EXPERIMENT

847 848 B.1 SPSA COMPONENT ABLATIONS

849
 850 We ablate the **SPSA** operator defined in Eq. 4 by removing its last two terms: $A (-\gamma\mu \xi_i \cdot$
 851 $\text{softmax}((\mathbf{W}_i \mathbf{x})^\top \xi_i))$ and $B (+\mu \xi_i)$. All other implementation details, training protocol, and hy-
 852 perparameters (including γ, μ) follow the main setup. On BSD68 at $\sigma = 50$, removing either com-
 853 ponent degrades performance, while the full model (Eq. 4) attains the best average results, indicating
 854 that the two terms play complementary roles.
 855

856 857 B.2 POSITION OF PROMPT BLOCKS.

858
 859 In the hierarchical decoder, we ablate where to inject the prompt blocks. Table 8 compares placing
 860 prompts at blocks 1&2, at block 6, and at all decoder blocks on the denoising task with $\sigma = 15$.
 861 While placing prompts at all blocks yields only marginal gains over block 6 (up to 0.04 dB in PSNR
 862 and 0.001 in SSIM), it introduces nontrivial computational and latency overhead. We therefore adopt
 863 the single-block design at block 6 as the default, which closely matches the all-block variant while
 reducing wall-clock time and memory footprint.

B.3 COMPLEMENTARITY BETWEEN PROMPT AND DATA-CONSISTENCY MODULES.

To avoid attributing the overall improvement to a single component, we conduct a systematic ablation on the deraining task, comparing four configurations: removing the Prompt (w/o Prompt), removing the data-consistency module (w/o DC), removing both components as the cascaded baseline (w/o Prompt & DC), and enabling both components (Prompt+DC). To ensure fairness, all other settings—the number of unrolled steps K , step size η , training schedule, and parameter budget—are kept identical across variants. As summarized in Table 9, relative to the baseline without both modules (w/o Prompt DC), introducing either module alone yields consistent gains in PSNR/SSIM; enabling both simultaneously (Prompt+DC) produces the largest improvement, confirming the strong complementarity between the structural prior provided by the Prompt and the observation-consistency constraint enforced by the DC module.

Table 7: Ablation study of **SPSA** components in Eq. 4 on BSD68 with $\sigma = 50$. The complete formulation (Eq. 4) achieves the best overall performance. Here, A corresponds to $-\gamma\mu \xi_i \cdot \text{softmax}((\mathbf{W}_i \mathbf{x})^\top \xi_i)$, and B corresponds to $+\mu \xi_i$.

Variant	PSNR	SSIM
w/o A	27.57	0.753
w/o B	27.97	0.798
Full	28.33	0.967

Table 8: Ablation of prompt-injection positions on BSD68 at $\sigma = 15$. Injecting at block 6 matches all-block injection while substantially reducing computation. Results are shown for blocks 1&2, block 6, and all blocks.

Placement	PSNR	SSIM
block 6	34.02	0.963
blocks 1&2	33.97	0.962
all blocks	34.06	0.964

Table 9: Ablation study of **Prompt** and **Data-Consistency (DC)** components within the cascaded Transformer unrolled from Eq. 8, evaluated on the Rain100L dataset.

Variant	PSNR	SSIM
w/o DC&prompt	36.77	0.977
w/o DC	37.07	0.978
w/o prompt	37.46	0.979
Prompt&DC	38.70	0.983

B.4 ACCURACY-OVERHEAD TRADE-OFF FOR PROMPT DESIGN

To further quantify the efficiency of the prompt design, we plot parameter-overhead vs. PSNR curves on Rain100L under the single-task deraining setting. Figure 9 summarizes three sets of experiments: (a) different multi-insertion strategies, (b) the number of prompt tokens N with insertion fixed at the 6-th decoder block, and (c) the projector rank R , again with insertion fixed at the 6-th block. For all panels, the horizontal axis denotes the relative parameter overhead (%) with respect to the “No Prompt” baseline, and the vertical axis reports PSNR (dB). We also provide the exact numerical values in the accompanying appendix tables.

C VARIATIONAL DERIVATION OF THE ENERGY FUNCTIONAL

For notational simplicity, the subscript i is omitted throughout this section. To develop the optimization algorithm for the proposed model, we consider the variational derivative of the energy functional $\mathcal{R}(\mathbf{x})$, defined as:

$$\begin{aligned} \mathcal{R}(\mathbf{x}) = & \frac{1}{2} \int_{\Omega} \text{Tr} \left(\ln \left(\mathbf{I} + (\mathbf{W}\mathbf{x}(s))(\mathbf{W}\mathbf{x}(s))^\top \right) \right) ds \\ & - \mu \int_{\Omega} \text{Tr} \left(\ln \left(\mathbf{I} + (\mathbf{W}\mathbf{x}(s))(\boldsymbol{\xi}(s))^\top \right) \right) ds \end{aligned} \quad (9)$$

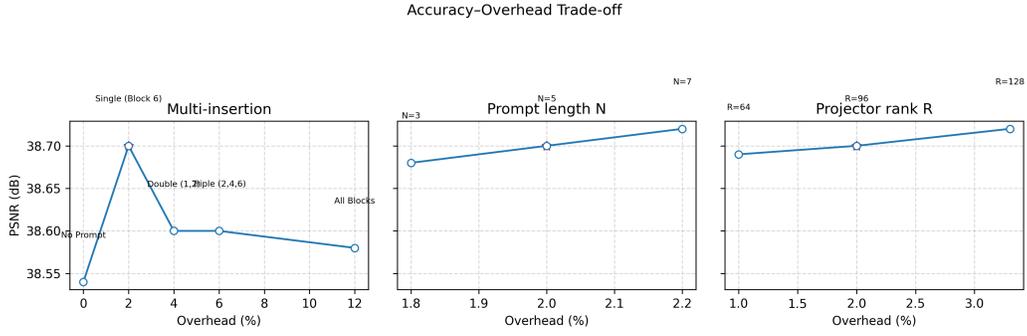


Figure 9: Accuracy–overhead trade-off for prompt design on Rain100L (single-task deraining). (a) Multi-insertion strategies (No Prompt, Single block 6, Double, Triple, All Blocks). (b) Number of prompt tokens N with insertion fixed at block 6. (c) Projector rank R of W_i with insertion fixed at block 6. The horizontal axis shows parameter overhead (%) relative to the “No Prompt” baseline; the vertical axis reports PSNR (dB). Star markers indicate the default configuration used in the main paper (single insertion at block 6 with $N = 5$ and $R = 96$), which lies on or very close to the Pareto front of the accuracy–overhead trade-off.

where $\mathbf{x}(s)$ is the optimization variable, $\boldsymbol{\xi}(s)$ denotes the structural prompt, \mathbf{W} is a linear transformation operator, and μ is a regularization weight.

Let us define $\mathbf{y}(s) = \mathbf{W}\mathbf{x}(s)$ and $\tilde{\mathbf{y}}(s) = \boldsymbol{\xi}(s)$.

C.1 VARIATION OF THE FIRST TERM

Consider the first term:

$$\mathcal{R}_1(\mathbf{x}) = \frac{1}{2} \int_{\Omega} \text{Tr} \left(\ln \left(\mathbf{I} + \mathbf{y}(s)\mathbf{y}(s)^\top \right) \right) ds \quad (10)$$

Using the matrix differential identity:

$$d \text{Tr} \left(\ln \left(\mathbf{I} + \mathbf{y}\mathbf{y}^\top \right) \right) = \text{Tr} \left(\left(\mathbf{I} + \mathbf{y}\mathbf{y}^\top \right)^{-1} d(\mathbf{y}\mathbf{y}^\top) \right) \quad (11)$$

Since $d(\mathbf{y}\mathbf{y}^\top) = d\mathbf{y} \cdot \mathbf{y}^\top + \mathbf{y} \cdot d\mathbf{y}^\top$, we have:

$$\begin{aligned} d \text{Tr} \left(\ln \left(\mathbf{I} + \mathbf{y}\mathbf{y}^\top \right) \right) &= \text{Tr} \left(\left(\mathbf{I} + \mathbf{y}\mathbf{y}^\top \right)^{-1} (d\mathbf{y} \cdot \mathbf{y}^\top + \mathbf{y} \cdot d\mathbf{y}^\top) \right) \\ &= 2 \text{Tr} \left(\left(\mathbf{I} + \mathbf{y}\mathbf{y}^\top \right)^{-1} \mathbf{y} d\mathbf{y}^\top \right) \\ &= 2 \left(\left(\mathbf{I} + \mathbf{y}\mathbf{y}^\top \right)^{-1} \mathbf{y} \right)^\top d\mathbf{y} \end{aligned} \quad (12)$$

where the last equality follows from the identity $\text{Tr}(A^\top) = \text{Tr}(A)$ and $\text{Tr}(A^\top B) = \langle A, B \rangle_F$.

Thus, the variation is:

$$d \text{Tr} \left(\ln \left(\mathbf{I} + \mathbf{y}\mathbf{y}^\top \right) \right) = \left\langle 2 \left(\mathbf{I} + \mathbf{y}\mathbf{y}^\top \right)^{-1} \mathbf{y}, d\mathbf{y} \right\rangle_F \quad (13)$$

Substituting $\mathbf{y}(s) = \mathbf{W}\mathbf{x}(s)$ and $d\mathbf{y}(s) = \mathbf{W}d\mathbf{x}(s)$, we obtain:

$$\begin{aligned} d\mathcal{R}_1 &= \frac{1}{2} \int_{\Omega} \left\langle 2 \left(\mathbf{I} + \mathbf{y}(s)\mathbf{y}(s)^\top \right)^{-1} \mathbf{y}(s), \mathbf{W}d\mathbf{x}(s) \right\rangle_F ds \\ &= \int_{\Omega} \left\langle \left(\mathbf{I} + \mathbf{y}(s)\mathbf{y}(s)^\top \right)^{-1} \mathbf{y}(s), \mathbf{W}d\mathbf{x}(s) \right\rangle_F ds \end{aligned} \quad (14)$$

972 Using the definition of the adjoint operator \mathbf{W}^* :

$$973 \quad d\mathcal{R}_1 = \int_{\Omega} \left\langle \mathbf{W}^* \left((\mathbf{I} + \mathbf{y}(s)\mathbf{y}(s)^\top)^{-1} \mathbf{y}(s) \right), d\mathbf{x}(s) \right\rangle_F ds \quad (15)$$

976 Therefore, the variational derivative is:

$$977 \quad \frac{\delta\mathcal{R}_1}{\delta\mathbf{x}}(s) = \mathbf{W}^* \left((\mathbf{I} + \mathbf{y}(s)\mathbf{y}(s)^\top)^{-1} \mathbf{y}(s) \right) \quad (16)$$

981 C.2 VARIATION OF THE SECOND TERM

982 Now consider the second term:

$$983 \quad \mathcal{R}_2(\mathbf{x}) = \int_{\Omega} \text{Tr} \left(\ln (\mathbf{I} + \mathbf{y}(s)\tilde{\mathbf{y}}(s)^\top) \right) ds \quad (17)$$

984 Using the matrix differential identity:

$$985 \quad d \text{Tr} \left(\ln (\mathbf{I} + \mathbf{y}\tilde{\mathbf{y}}^\top) \right) = \text{Tr} \left((\mathbf{I} + \mathbf{y}\tilde{\mathbf{y}}^\top)^{-1} d(\mathbf{y}\tilde{\mathbf{y}}^\top) \right) \quad (18)$$

$$986 \quad = \text{Tr} \left((\mathbf{I} + \mathbf{y}\tilde{\mathbf{y}}^\top)^{-1} d\mathbf{y}\tilde{\mathbf{y}}^\top \right) \quad (19)$$

987 where the second equality holds because $\tilde{\mathbf{y}}$ is independent of \mathbf{x} .

988 Using the cyclic property of the trace:

$$989 \quad \text{Tr} \left((\mathbf{I} + \mathbf{y}\tilde{\mathbf{y}}^\top)^{-1} d\mathbf{y}\tilde{\mathbf{y}}^\top \right) = \text{Tr} \left(\tilde{\mathbf{y}}^\top (\mathbf{I} + \mathbf{y}\tilde{\mathbf{y}}^\top)^{-1} d\mathbf{y} \right) \quad (20)$$

990 This can be written as an inner product:

$$991 \quad d \text{Tr} \left((\mathbf{I} + \mathbf{y}\tilde{\mathbf{y}}^\top)^{-1} d\mathbf{y}\tilde{\mathbf{y}}^\top \right) = \left\langle (\mathbf{I} + \mathbf{y}\tilde{\mathbf{y}}^\top)^{-\top} \tilde{\mathbf{y}}, d\mathbf{y} \right\rangle_F \quad (21)$$

992 Substituting $\mathbf{y}(s) = \mathbf{W}\mathbf{x}(s)$ and $d\mathbf{y}(s) = \mathbf{W}d\mathbf{x}(s)$, we obtain:

$$993 \quad d\mathcal{R}_2 = \int_{\Omega} \left\langle (\mathbf{I} + \mathbf{y}(s)\tilde{\mathbf{y}}(s)^\top)^{-\top} \tilde{\mathbf{y}}(s), \mathbf{W}d\mathbf{x}(s) \right\rangle_F ds \quad (22)$$

994 Using the adjoint operator \mathbf{W}^* :

$$995 \quad d\mathcal{R}_2 = \int_{\Omega} \left\langle \mathbf{W}^* \left((\mathbf{I} + \mathbf{y}(s)\tilde{\mathbf{y}}(s)^\top)^{-\top} \tilde{\mathbf{y}}(s) \right), d\mathbf{x}(s) \right\rangle_F ds \quad (23)$$

996 Therefore, the variational derivative is:

$$997 \quad \frac{\delta\mathcal{R}_2}{\delta\mathbf{x}}(s) = \mathbf{W}^* \left((\mathbf{I} + \mathbf{y}(s)\tilde{\mathbf{y}}(s)^\top)^{-\top} \tilde{\mathbf{y}}(s) \right) \quad (24)$$

1000 C.3 COMBINING BOTH TERMS

1001 Combining both components with their respective coefficients and regularization weight μ , we derive the complete variational derivative:

$$1002 \quad \frac{\delta\mathcal{R}}{\delta\mathbf{x}}(s) = \frac{\delta\mathcal{R}_1}{\delta\mathbf{x}}(s) - \mu \frac{\delta\mathcal{R}_2}{\delta\mathbf{x}}(s) \quad (25)$$

$$1003 \quad = \mathbf{W}^* \left((\mathbf{I} + \mathbf{y}(s)\mathbf{y}(s)^\top)^{-1} \mathbf{y}(s) \right) - \mu \mathbf{W}^* \left((\mathbf{I} + \mathbf{y}(s)\tilde{\mathbf{y}}(s)^\top)^{-\top} \tilde{\mathbf{y}}(s) \right)$$

1004 This gradient form supports the structure-aware optimization in the main algorithm and highlights the explicit interaction between the target variable \mathbf{x} and the structural prompt ξ within the proposed framework.

C.4 VARIATIONAL DERIVATIVE APPROXIMATION

Consider the variational derivative of the energy functional $\mathcal{R}[\mathbf{x}]$ with respect to $\mathbf{x}(s)$:

$$\frac{\delta \mathcal{R}}{\delta \mathbf{x}}(s) = \mathbf{W}^* \left((\mathbf{I} + \mathbf{y}(s)\mathbf{y}(s)^\top)^{-1} \mathbf{y}(s) \right) - \mu \mathbf{W}^* \left((\mathbf{I} + \mathbf{y}(s)\tilde{\mathbf{y}}(s)^\top)^{-\top} \tilde{\mathbf{y}}(s) \right), \quad (26)$$

where

$$\mathbf{y}(s) = \mathbf{W}\mathbf{x}(s), \quad \tilde{\mathbf{y}}(s) = \boldsymbol{\xi}(s). \quad (27)$$

Using the Neumann series expansion,

$$(\mathbf{I} + A)^{-1} = \mathbf{I} - A + A^2 - \dots \approx \mathbf{I} - A, \quad (28)$$

which holds for matrices with sufficiently small spectral norm, we obtain

$$(\mathbf{I} + \mathbf{y}\mathbf{y}^\top)^{-1} \approx \mathbf{I} - \mathbf{y}\mathbf{y}^\top, \quad (\mathbf{I} + \mathbf{y}\tilde{\mathbf{y}}^\top)^{-\top} \approx \mathbf{I} - \tilde{\mathbf{y}}\mathbf{y}^\top. \quad (29)$$

Substituting the above into the variational derivative yields

$$\frac{\delta \mathcal{R}}{\delta \mathbf{x}} \approx \mathbf{W}^* \left((\mathbf{I} - \mathbf{y}\mathbf{y}^\top) \mathbf{y} \right) - \mu \mathbf{W}^* \left((\mathbf{I} - \tilde{\mathbf{y}}\mathbf{y}^\top) \tilde{\mathbf{y}} \right), \quad (30)$$

and expanding the matrix products gives

$$\frac{\delta \mathcal{R}}{\delta \mathbf{x}} \approx \mathbf{W}^* \left(\mathbf{y} - \mathbf{y}(\mathbf{y}^\top \mathbf{y}) \right) - \mu \mathbf{W}^* \left(\tilde{\mathbf{y}} - \tilde{\mathbf{y}}(\mathbf{y}^\top \tilde{\mathbf{y}}) \right). \quad (31)$$

This derivation is rigorous and relies solely on the first-order Neumann expansion and standard matrix algebra, and it serves as the reference gradient when analyzing the approximation error of the attention-based form below.

To obtain a more intuitive subspace-membership interpretation, the inner-product terms $\mathbf{y}^\top \mathbf{y}$ and $\mathbf{y}^\top \tilde{\mathbf{y}}$ in equation 31 are heuristically replaced by softmax-normalized weights:

$$\frac{\delta \mathcal{R}}{\delta \mathbf{x}} \approx \mathbf{W}^* \left(\mathbf{y} - \gamma \mathbf{y} \cdot \text{softmax}(\mathbf{y}^\top \mathbf{y}) \right) - \mu \mathbf{W}^* \left(\tilde{\mathbf{y}} - \gamma \tilde{\mathbf{y}} \cdot \text{softmax}(\mathbf{y}^\top \tilde{\mathbf{y}}) \right), \quad (32)$$

where γ is a scaling factor. Although this softmax replacement is originally motivated by an intuitive subspace-membership view, we next provide a closed-form error expression and an explicit sufficient condition under which the approximation from the Neumann-expanded gradient equation 31 to the attention-like form equation 32 is accurate.

Approximation error analysis We denote $V := \tilde{\mathbf{y}} \in \mathbb{R}^{d \times n}$, $C = [c_1, \dots, c_n] := \mathbf{y}^\top \tilde{\mathbf{y}} \in \mathbb{R}^{n \times n}$, and $A = [a_1, \dots, a_n] := \gamma \text{softmax}(\mathbf{y}^\top \tilde{\mathbf{y}}) \in \mathbb{R}^{n \times n}$. We now analyze, for the optimal choice of γ , the approximation error incurred when Eq. equation 32 replaces $\tilde{\mathbf{y}}(\mathbf{y}^\top \tilde{\mathbf{y}})$ in Eq. equation 31 by $\gamma \tilde{\mathbf{y}} \cdot \text{softmax}(\mathbf{y}^\top \tilde{\mathbf{y}})$. The same analysis applies to the approximation error incurred when $\mathbf{y}(\mathbf{y}^\top \mathbf{y})$ is replaced by $\gamma \mathbf{y} \cdot \text{softmax}(\mathbf{y}^\top \mathbf{y})$. We further relax the scalar γ to a diagonal matrix $\Gamma = \text{diag}(\gamma_1, \dots, \gamma_n)$. The corresponding approximation-error minimization problem is then defined as

$$\min_{\Gamma = \text{diag}(\gamma_1, \dots, \gamma_n)} \|VC - V A \Gamma\|_F^2. \quad (33)$$

Because the Frobenius norm is a sum of column-wise Euclidean norms, problem equation 33 decouples into n independent one-dimensional problems:

$$\min_{\gamma_j \in \mathbb{R}} \|V(c_j - \gamma_j a_j)\|_2^2, \quad j = 1, \dots, n. \quad (34)$$

Thus it suffices to analyze the approximation error for a single column and then sum over j .

Error bound. Let $V \in \mathbb{R}^{d \times n}$ be fixed, and let $a, c \in \mathbb{R}^n$ denote a pair of weight vectors. Define

$$G := V^\top V \succeq 0, \quad \langle x, y \rangle_G := x^\top G y, \quad \|x\|_G := \sqrt{\langle x, x \rangle_G}, \quad (35)$$

and the G -cosine

$$\cos_G(a, c) := \frac{\langle a, c \rangle_G}{\|a\|_G \|c\|_G}, \quad \text{for } \|a\|_G > 0, \|c\|_G > 0. \quad (36)$$

1080 Consider the one-dimensional convex quadratic

$$1081 \quad f(\gamma) := \|V(c - \gamma a)\|_2^2 = (c - \gamma a)^\top G(c - \gamma a). \quad (37)$$

1083 Writing $G = H^\top H$ with a symmetric square root $H := G^{1/2}$ (defined on the closure of $\text{ran}(G)$),
1084 we have

$$1085 \quad f(\gamma) = \|Hc - \gamma Ha\|_2^2 = \|Hc\|_2^2 - 2\gamma \langle Ha, Hc \rangle + \gamma^2 \|Ha\|_2^2. \quad (38)$$

1086 The minimizer and minimum are therefore

$$1087 \quad \gamma^* = \frac{\langle Ha, Hc \rangle}{\|Ha\|_2^2} = \frac{\langle a, c \rangle_G}{\|a\|_G^2}, \quad (39)$$

1089 and

$$1091 \quad \min_{\gamma \in \mathbb{R}} f(\gamma) = \|Hc\|_2^2 - \frac{\langle Ha, Hc \rangle^2}{\|Ha\|_2^2} = \|c\|_G^2 \left(1 - \frac{\langle a, c \rangle_G^2}{\|a\|_G^2 \|c\|_G^2}\right) = \|Vc\|_2^2 (1 - \cos_G^2(a, c)). \quad (40)$$

1093 We thus obtain the following closed-form error identity:

$$1095 \quad \min_{\gamma \in \mathbb{R}} \|V(c - \gamma a)\|_2^2 = \|Vc\|_2^2 (1 - \cos_G^2(a, c)), \quad \gamma^* = \frac{\langle a, c \rangle_G}{\|a\|_G^2}. \quad (41)$$

1097 In particular, if $\cos_G^2(a, c) \geq 1 - \varepsilon$ for some $\varepsilon \in [0, 1)$, then

$$1099 \quad \min_{\gamma \in \mathbb{R}} \|V(c - \gamma a)\|_2^2 \leq \varepsilon \|Vc\|_2^2. \quad (42)$$

1101 Moreover, the minimum in equation 41 is zero if and only if $c \in \text{span}\{a\} \oplus \ker(G)$, i.e., there exist
1102 $\lambda \in \mathbb{R}$ and $z \in \ker(G)$ such that $c = \lambda a + z$; when $G \succ 0$, this reduces to strict collinearity $c = \lambda a$.

1103 Equations equation 41–equation 42 provide a closed-form error expression for replacing the linear
1104 weights c by attention weights a (plus optimal scaling). Summing over columns j gives the corre-
1105 sponding Frobenius-norm error for problem equation 33, and this identity will be used to justify the
1106 gradient-to-attention mapping in Eq. equation 3 and its operator realization in Eq. equation 4.

1108 **A sufficient condition for small approximation error.** We now derive an explicit sufficient condi-
1109 tion under which the factor $\cos_G(a, c)$ in equation 41 is close to 1, so that the error bound equa-
1110 tion 42 is very small.

1111 Let $c \in \mathbb{R}^n$ be a column-wise linear weight vector. Using the translation invariance of softmax, we
1112 define a nonnegative shifted version

$$1114 \quad \hat{c} := c - \min_i c_i \mathbf{1} \geq 0, \quad s := \mathbf{1}^\top \hat{c} > 0, \quad \pi := \hat{c}/s \in \Delta_{n-1}, \quad (43)$$

1115 and the corresponding attention weight

$$1117 \quad a := \text{softmax}(\beta \hat{c}) \in \Delta_{n-1}, \quad (44)$$

1118 where $\beta > 0$ is the temperature. Note that π and \hat{c} are strictly collinear:

$$1120 \quad \langle \pi, \hat{c} \rangle_G = \frac{\|\hat{c}\|_G^2}{s}, \quad \|\pi\|_G = \frac{\|\hat{c}\|_G}{s}, \quad \cos_G(\pi, \hat{c}) = 1. \quad (45)$$

1122 Denote the index of the main peak by

$$1124 \quad i^* = \arg \max_i \hat{c}_i, \quad (46)$$

1125 the main-peak ratio by

$$1127 \quad \rho := \frac{\max_i \hat{c}_i}{s} \in [1/n, 1], \quad (47)$$

1128 and the logit gap by

$$1129 \quad \Delta := \min_{i \neq i^*} (\hat{c}_{i^*} - \hat{c}_i) \geq 0. \quad (48)$$

1130 Writing $a = \pi + \delta$, we first bound $\cos_G(a, \hat{c})$ in terms of

$$1131 \quad \varepsilon := \frac{\|\delta\|_G}{\|\pi\|_G}. \quad (49)$$

1134 Using the Cauchy–Schwarz and triangle inequalities, we obtain

$$1135 \quad \langle a, \hat{c} \rangle_G = \langle \pi, \hat{c} \rangle_G + \langle \delta, \hat{c} \rangle_G \geq \|\pi\|_G \|\hat{c}\|_G - \|\delta\|_G \|\hat{c}\|_G = (1 - \varepsilon) \|\pi\|_G \|\hat{c}\|_G, \quad (50)$$

$$1137 \quad \|a\|_G \leq \|\pi\|_G + \|\delta\|_G = (1 + \varepsilon) \|\pi\|_G,$$

1138 hence

$$1139 \quad \cos_G(a, \hat{c}) = \frac{\langle a, \hat{c} \rangle_G}{\|a\|_G \|\hat{c}\|_G} \geq \frac{1 - \varepsilon}{1 + \varepsilon}, \quad \cos_G^2(a, \hat{c}) \geq \left(\frac{1 - \varepsilon}{1 + \varepsilon}\right)^2. \quad (51)$$

1142 We next bound ε by relating the G -norm to the Euclidean norm on the two-dimensional subspace
1143 $S := \text{span}\{\pi, \delta\}$. Let

$$1144 \quad \lambda_{\min} := \lambda_{\min}(G|_S), \quad \lambda_{\max} := \lambda_{\max}(G|_S), \quad \kappa_S := \frac{\lambda_{\max}}{\lambda_{\min}} \geq 1. \quad (52)$$

1147 Then for any $z \in S$,

$$1148 \quad \sqrt{\lambda_{\min}} \|z\|_2 \leq \|z\|_G \leq \sqrt{\lambda_{\max}} \|z\|_2. \quad (53)$$

1149 In particular,

$$1151 \quad \|\delta\|_G \leq \sqrt{\lambda_{\max}} \|\delta\|_2, \quad \|\pi\|_G \geq \sqrt{\lambda_{\min}} \|\pi\|_2 = \sqrt{\lambda_{\min}} \frac{\|\hat{c}\|_2}{s}, \quad (54)$$

1153 so that

$$1154 \quad \varepsilon = \frac{\|\delta\|_G}{\|\pi\|_G} \leq \kappa_S \frac{s}{\|\hat{c}\|_2} \|\delta\|_2. \quad (55)$$

1156 Using $s = \|\hat{c}\|_1 \leq \sqrt{n} \|\hat{c}\|_2$ (Cauchy–Schwarz), we further obtain

$$1158 \quad \varepsilon \leq \kappa_S \sqrt{n} \|\delta\|_2. \quad (56)$$

1159 It remains to bound $\|\delta\|_2 = \|a - \pi\|_2$ in terms of ρ and Δ . For the softmax weights, we have

$$1161 \quad a_{i^*} = \frac{1}{1 + \sum_{i \neq i^*} \exp[-\beta(\hat{c}_{i^*} - \hat{c}_i)]} \geq \frac{1}{1 + (n-1)e^{-\beta\Delta}}, \quad (57)$$

1164 so that

$$1165 \quad 1 - a_{i^*} \leq (n-1)e^{-\beta\Delta}. \quad (58)$$

1166 Moreover,

$$1167 \quad \sum_{i \neq i^*} a_i^2 \leq \sum_{i \neq i^*} a_i = 1 - a_{i^*}, \quad (59)$$

1169 which implies

$$1170 \quad \|a - e_{i^*}\|_2^2 = (1 - a_{i^*})^2 + \sum_{i \neq i^*} a_i^2 \leq 2(1 - a_{i^*}) \leq 2(n-1)e^{-\beta\Delta}. \quad (60)$$

1173 Similarly, for π we have

$$1174 \quad \|\pi - e_{i^*}\|_2^2 = (1 - \rho)^2 + \sum_{i \neq i^*} \pi_i^2 \leq 2(1 - \rho). \quad (61)$$

1177 By the triangle inequality,

$$1179 \quad \|\delta\|_2 = \|a - \pi\|_2 \leq \|a - e_{i^*}\|_2 + \|\pi - e_{i^*}\|_2 \leq \sqrt{2(1 - \rho)} + \sqrt{2(n-1)e^{-\beta\Delta}}. \quad (62)$$

1181 Absorbing numerical constants into a universal constant $C > 0$, we obtain the bound

$$1183 \quad \varepsilon \leq C \kappa_S \left(\sqrt{1 - \rho} + \sqrt{n-1} e^{-\beta\Delta/2} \right). \quad (63)$$

1185 Combining equation 51 and equation 63, we see that when the right-hand side of equation 63 is
1186 much smaller than 1 (i.e., the linear weights are sharply peaked with large ρ , the logit gap Δ is
1187 sufficiently large under the given temperature β , and G is not severely ill-conditioned on S), we
have $\cos_G(a, \hat{c}) \approx 1$ and hence $\cos_G^2(a, \hat{c}) \approx 1$. If, in addition, G is insensitive to constant shifts

(for instance, $G1 = 0$ or we apply a G -orthogonal projection to remove the mean component), then $\cos_G(a, c) = \cos_G(a, \hat{c})$ and equation 63 directly controls $\cos_G(a, c)$.

Substituting this into equation 41–equation 42 shows that, under equation 63, the approximation error between the softmax-based form equation 32 and the Neumann-expanded gradient equation 31 is negligible. In summary, the identity equation 41 together with the bound equation 42 provides a closed-form error bound for the gradient-to-attention mapping in Eq. equation 3, and the sufficient condition equation 63 characterizes when this error bound is very small. Since Eq. equation 4 is an operator-level instantiation of the same mapping (via SPSA/MPSA), the same error identity and condition equation 63 also justify the approximation steps in Eq. equation 4.

D LIMITATION

Despite achieving competitive results across restoration tasks and affording transparent prompt–attention dynamics, one aspect merits further investigation: *computational efficiency*. Specifically, unrolling the gradient flow and stacking Transformer blocks increase the backpropagation memory footprint and wall-clock training time; at inference, the near-quadratic complexity of attention with respect to sequence length (or spatial resolution) can exacerbate latency. We view this as an engineering trade-off rather than a fundamental limitation, and it does not compromise our core conclusion of a variationally anchored, interpretable prompt–attention coupling; nevertheless, further optimization is warranted for resource- and latency-constrained deployments.

A second limitation lies in our current treatment of multi-scale processing. While the single-scale WBPT is fully white-box, the multi-scale extension WBPT[†] introduces a learned pyramid aggregator for cross-scale fusion. This component remains black-box, even though the underlying attention and gradient-flow dynamics are white-box. From a variational perspective, a natural way to obtain a fully differentiable white-box pyramid is to use a multigrid V-cycle as the mathematical backbone: restriction and prolongation operators implement variational down/up sampling, while the STV-derived attention acts as the inter-grid operator on multi-channel features at each scale. Developing such a multigrid-style white-box pyramid, with STV-consistent cross-scale regularization explicitly encoded in these operators, is an important topic of ongoing work.

A third limitation concerns the scope of our experimental evaluation. In the present version, we systematically train and tune WBPT on the classic three-degradation all-in-one benchmark (denoising, deraining, dehazing), but do not yet provide full training and hyper-parameter search on larger unified protocols such as five-degradation or mixed-degradation settings, or on real-world benchmarks summarized in recent surveys. Extending WBPT to these more diverse settings and conducting a comprehensive study of its behavior under diverse degradation combinations is left for future work.

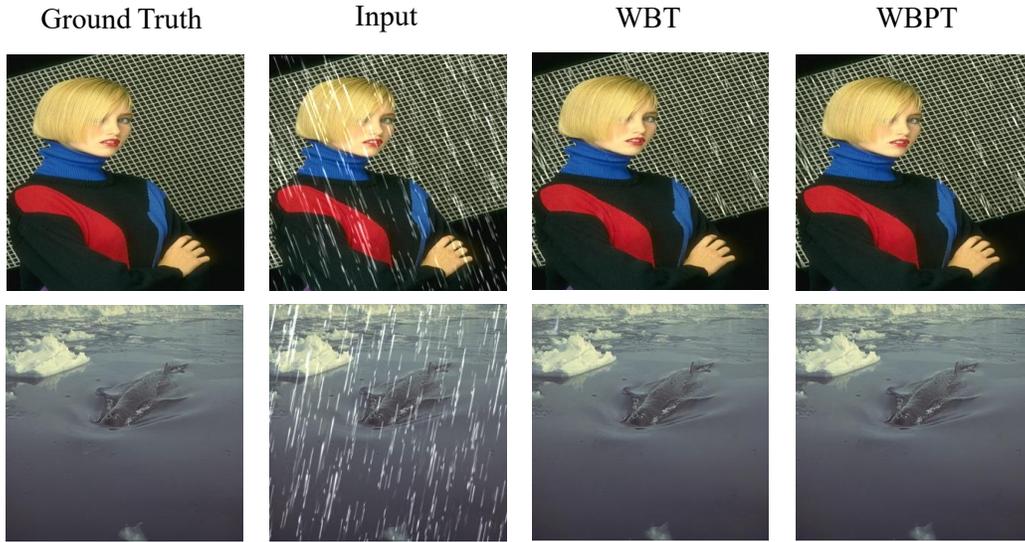
E QUALITATIVE RESULTS

We present additional qualitative results under the single-task setting to further demonstrate the effectiveness of *prompt-block*. The presented examples correspond one-to-one with the three quantitative single-task tables in the main text (Tables 2, 3, 4), serving to complement and visually illustrate the trends reported in the main paper.

F ALGORITHM

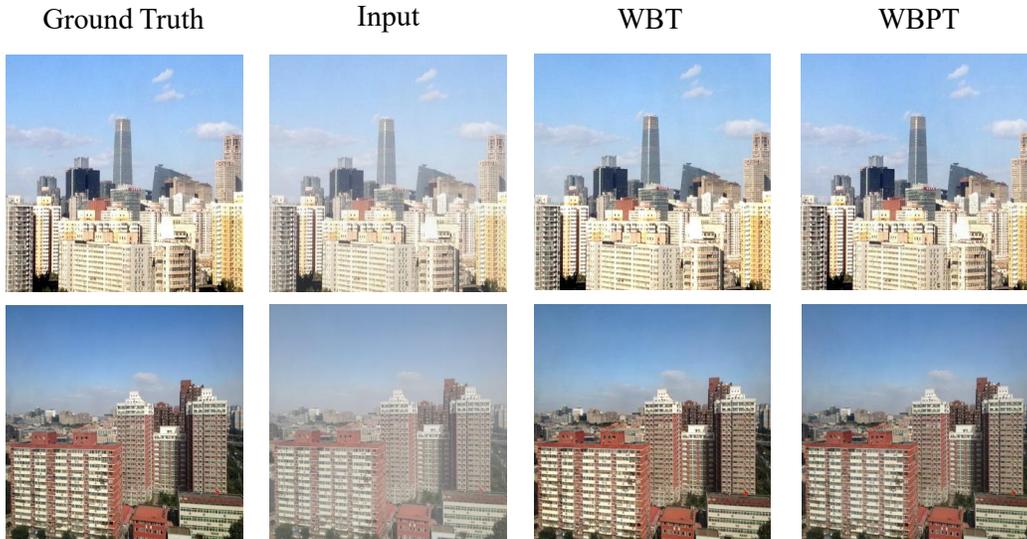
This part provides PyTorch-style pseudocode for the WBPT. Alg. 1 outlines the overall training loop with T -stage learnable gradient updates; Alg. 2 details one iteration stage with the SwinIR backbone and prompt injection; Alg. 3 specifies the prompted window attention and feedforward modules. Unless otherwise noted, we use the following defaults in the pseudocode: $epochs = 120$, $stages T = 10$, $embed_dim = 96$, $num_heads = [6, 6, 6]$, $window = 8$, $prompt_len = 5$, Adam with $(\beta_1=0.9, \beta_2=0.999)$, learning rate 1×10^{-4} , MSE reconstruction loss, and 128×128 random crops with rotation/flip augmentations. The pseudocode focuses on core computations; engineering aspects (I/O, multi-GPU, logging, checkpointing) are omitted for brevity. Complexity expressions report the dominant terms (attention and MLP). Notation: x_{noisy}/x_{clean} denote inputs/targets, z the

1242
1243
1244
1245
1246
1247
1248
1249
1250
1251
1252
1253
1254
1255
1256
1257
1258
1259



1260 Figure 10: Deraining results for all-in-one methods.

1261
1262
1263
1264
1265
1266
1267
1268
1269
1270
1271
1272
1273
1274
1275
1276
1277
1278
1279
1280



1281 Figure 11: Dehazing results for all-in-one methods.

1282
1283
1284
1285

current state, U_k the backbone output/prompt at stage k , and t the iteration index. Prompt injection is fixed at block indices $\{2, 4, 6\}$ within the SwinIR backbone.

1286
1287
1288

1289 G TRANSFORMER AND PROMPT BLOCKS IN WBPT

1290
1291
1292
1293
1294
1295

As stated in Section 2.2 of the main manuscript, we present in Fig. 13 the block diagram of the Prompt block corresponding to ξ_i , and further elaborate on the implementation details of the Transformer block used within this Prompt block in Fig. 14. The Prompt block and the Transformer block follow the design and hyper-parameter settings outlined in Potlapalli et al. (2023) and Zamir et al. (2022), respectively.

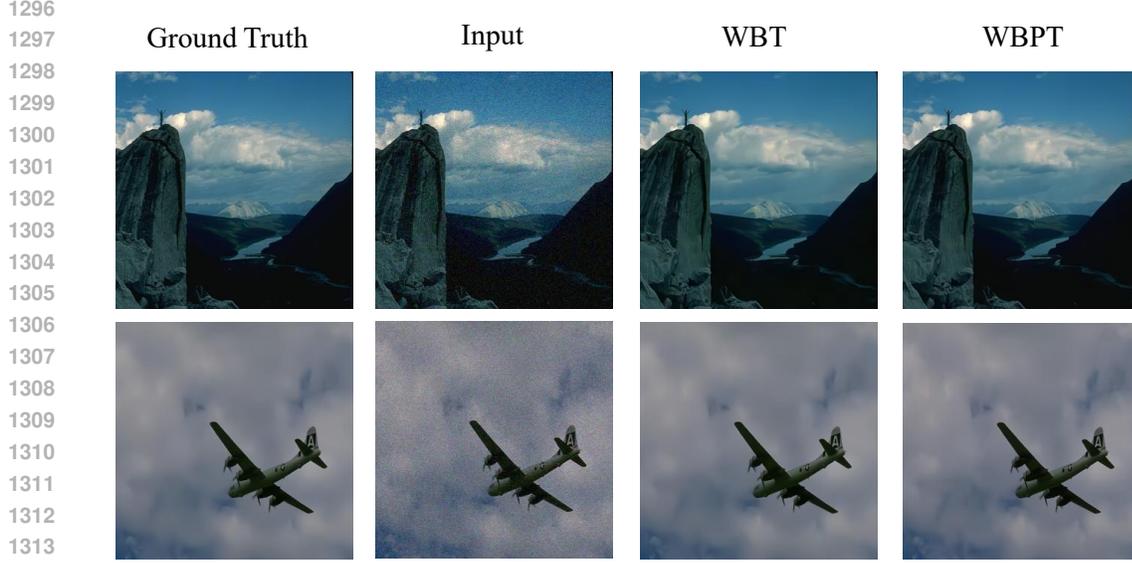


Figure 12: Denoising results for all-in-one methods.

1317
1318
1319
1320
1321
1322
1323
1324
1325
1326
1327
1328
1329
1330
1331
1332
1333
1334
1335

Algorithm 1: PyTorch-style pseudocode for WBPT training (overall)

Input: Dataset \mathcal{D} ; epochs E (default 120); stages K (default 10); hyperparameters and prompt config P

Output: Trained parameters θ^* ; best validation metrics (PSNR/SSIM)

Optimizer & Loss: Adam($\beta_1=0.9, \beta_2=0.999$), lr 1×10^{-4} ; reconstruction loss $\mathcal{L} = \text{MSE}$

for epoch $\leftarrow 1$ **to** E **do**

foreach $(\mathbf{x}_{\text{noisy}}, \mathbf{x}_{\text{clean}}) \in \text{DataLoader}(\mathcal{D})$ **do**

$\mathbf{x}_0 \leftarrow \mathbf{x}_{\text{noisy}}; \mathbf{W}_{\text{hist}} \leftarrow []; \boldsymbol{\xi}_{\text{hist}} \leftarrow []$

for $k \leftarrow 1$ **to** K **do**

$\{\mathbf{W}_i\}_{i=1}^N, \{\boldsymbol{\xi}_i\}_{i=1}^N \leftarrow \text{SwinIR}(\mathbf{x}_{k-1}; P, k)$ # with optional prompt injection

Append $\{\mathbf{W}_i\}_{i=1}^N$ to \mathbf{W}_{hist} and $\{\boldsymbol{\xi}_i\}_{i=1}^N$ to $\boldsymbol{\xi}_{\text{hist}}$; drop oldest if length > 5

$\mathbf{x}_k \leftarrow \text{MPSA}(\mathbf{x}_{k-1}, \mathbf{y}=\mathbf{x}_{\text{clean}}, \{\mathbf{W}_i\}, \{\boldsymbol{\xi}_i\}, k)$

$\mathcal{L} \leftarrow \text{MSE}(\mathbf{x}_K, \mathbf{x}_{\text{clean}})$

optimizer.zero_grad(); $\mathcal{L}.$ backward(); optimizer.step()

ValidateAndSaveIfBest(θ) # compute PSNR/SSIM on val, save best

Complexity: $\mathcal{O}(E \cdot N \cdot K \cdot (F_{\text{swin}} + F_{\text{mpsa}}))$ # Here N denotes the number of batches per epoch.

1336 1337 1338 G.1 PROMPT BLOCK IN WBPT FRAMEWORK

1339 1340 1341 G.1.1 PROMPT BLOCK OVERVIEW

1342
1343
1344
1345

Given prompt components $\mathbf{P}_c \in \mathbb{R}^{N \times \hat{H} \times \hat{W} \times \hat{C}}$ and input features $\mathbf{F}_1 \in \mathbb{R}^{\hat{H} \times \hat{W} \times \hat{C}}$, the prompt block refines the input via:

$$1346 \quad \hat{\mathbf{F}}_1 = \text{PIM}(\text{PGM}(\mathbf{P}_c, \mathbf{F}_1), \mathbf{F}_1) \quad (64)$$

1347
1348

1349
The block contains two modules: the Prompt Generation Module (PGM) and the Prompt Interaction Module (PIM), detailed below.

1350 **Algorithm 2:** One WBPT iteration stage (SwinIR with prompt injection)

1351 **Input:** Current state \mathbf{x}_{k-1} ; prompt config P ; SwinIR depths $[2, 2, 2]$, channels $C=96$, window

1352 $M=8$

1353 **Output:** $\{\mathbf{W}_i\}_{i=1}^N, \{\boldsymbol{\xi}_i\}_{i=1}^N$ (transformations and prompts for MPSA)

1354 $\mathbf{f} \leftarrow \text{Conv}_{3 \times 3}(\mathbf{x}_{k-1})$ # shallow feature

1355 **for** $\ell \leftarrow 1$ **to** 3 **do**

1356 $(\mathbf{f}, \text{size}) \leftarrow \text{PatchEmbed}(\mathbf{f})$

1357 **for** $b \leftarrow 1$ **to** 2 **do**

1358 **if** $\text{ShouldUsePrompt}(\ell, b, P)$ **then**

1359 $\{\boldsymbol{\xi}_i\}_{i=1}^N \leftarrow \text{PromptGenBlock}(\mathbf{f}); \mathbf{f} \leftarrow \mathbf{f} + \text{Inject}(\{\boldsymbol{\xi}_i\})$ # e.g., at

1360 fixed block indices $\{2, 4, 6\}$

1361 $\mathbf{f} \leftarrow \text{SwinTransformerBlock}(\mathbf{f}, \text{size})$

1362 **if** $\ell < 3$ **then**

1363 $\mathbf{f} \leftarrow \text{PatchMerging}(\mathbf{f}, \text{size})$

1364 $\{\mathbf{W}_i\}_{i=1}^N, \{\boldsymbol{\xi}_i\}_{i=1}^N \leftarrow \text{ExtractTransformationsAndPrompts}(\mathbf{f})$

1365 **return** $\{\mathbf{W}_i\}_{i=1}^N, \{\boldsymbol{\xi}_i\}_{i=1}^N$

1366 **Complexity:** $\sum_{\ell=1}^3 \sum_{b=1}^2 (F_{\text{W-MSA}} + F_{\text{MLP}}) \approx \mathcal{O}(n_W \cdot B \cdot M^2 C + LC^2)$

1370 **Algorithm 3:** Multi-Prompted Structure Attention (MPSA) with Learnable Data Consistency

1371 **Input:** $\mathbf{x}_k \in \mathbb{R}^{B \times H \times W \times C}$; transforms $\{\mathbf{W}_i\}_{i=1}^N$; prompts $\{\boldsymbol{\xi}_i\}_{i=1}^N$; measurements/targets \mathbf{y}

1372 **Output:** $\mathbf{x}_{k+1} \in \mathbb{R}^{B \times H \times W \times C}$

1373 **Multi-Prompted Structure Attention:**

1374 **for** $i \leftarrow 1$ **to** N **do**

1375 $\text{spsa}_i \leftarrow \text{SPSA}(\mathbf{x}_k \mid \mathbf{W}_i, \boldsymbol{\xi}_i, \gamma, \mu)$ # Eq. 4

1376 $\text{mps_out} \leftarrow [\mathbf{W}_1^* \ \cdots \ \mathbf{W}_N^*] \begin{bmatrix} \text{spsa}_1 \\ \vdots \\ \text{spsa}_N \end{bmatrix}$ # Eq. 5

1377

1378

1379

1380 **Learnable Data Consistency:**

1381 $\Delta_\phi(\mathbf{x}_k, \mathbf{y}) \leftarrow \text{LearntGradient}(\mathbf{x}_k, \mathbf{y})$

1382 **Gradient Flow Update:**

1383 $\mathbf{x}_{k+1} \leftarrow (\mathbf{I} - \eta \sum_{i=1}^N \mathbf{W}_i^* \mathbf{W}_i) \mathbf{x}_k - \eta \text{mps_out} - \eta \Delta_\phi(\mathbf{x}_k, \mathbf{y})$

1384 # Eq. 8

1385 **return** \mathbf{x}_{k+1}

1386 **Complexity:** MPSA $\mathcal{O}(N \cdot HWC^2)$; data consistency $\mathcal{O}(HWC^2)$

1389 G.1.2 PROMPT GENERATION MODULE (PGM)

1390 PGM dynamically generates input-conditioned prompts. First, global average pooling (GAP) is

1391 applied on \mathbf{F}_1 to produce a channel-wise descriptor $\mathbf{v} \in \mathbb{R}^{\hat{C}}$. A 1×1 convolution and softmax yield

1392 prompt weights $w \in \mathbb{R}^N$:

$$1393 w_i = \text{Softmax}(\text{Conv}_{1 \times 1}(\text{GAP}(\mathbf{F}_1))) \quad (65)$$

1394 These weights modulate the learned prompt components to form the final prompt \mathbf{P} :

$$1395 \mathbf{P} = \text{Conv}_{3 \times 3} \left(\sum_{c=1}^N w_i \mathbf{P}_c \right) \quad (66)$$

1400 To support variable-resolution inputs, prompt components are upsampled to match the spatial size

1401 of \mathbf{F}_1 via bilinear interpolation.

1404

1405

1406

1407

1408

1409

1410

1411

1412

1413

1414

1415

1416

1417

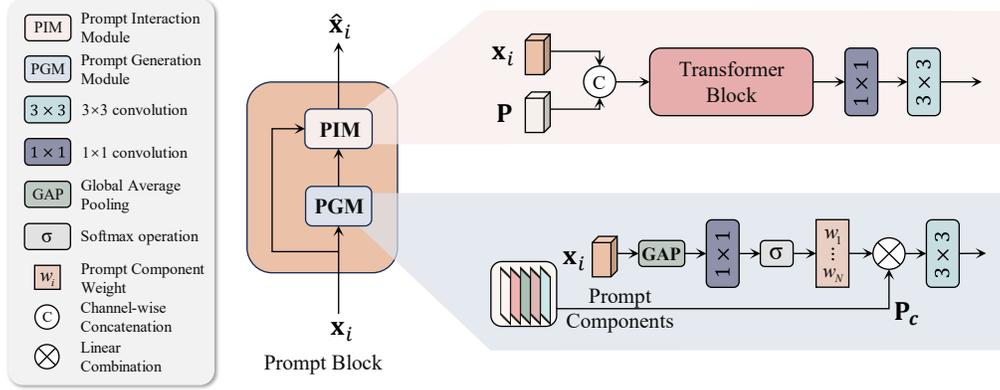


Figure 13: Overview of the Prompt block used in the WBPT framework. The Prompt block is composed of two sub modules, the Prompt Generation Module (PGM) and the Prompt Interaction Module (PIM).

1421

1422

G.1.3 PROMPT INTERACTION MODULE (PIM)

1423

1424

1425

1426

PIM fuses the generated prompt \mathbf{P} with features \mathbf{F}_1 via channel-wise concatenation, followed by a Transformer block:

1427

1428

1429

1430

1431

1432

1433

1434

1435

1436

1437

1438

1439

1440

1441

1442

1443

1444

1445

1446

1447

1448

1449

1450

1451

1452

1453

1454

1455

1456

1457

G.2 TRANSFORMER BLOCK IN WBPT FRAMEWORK

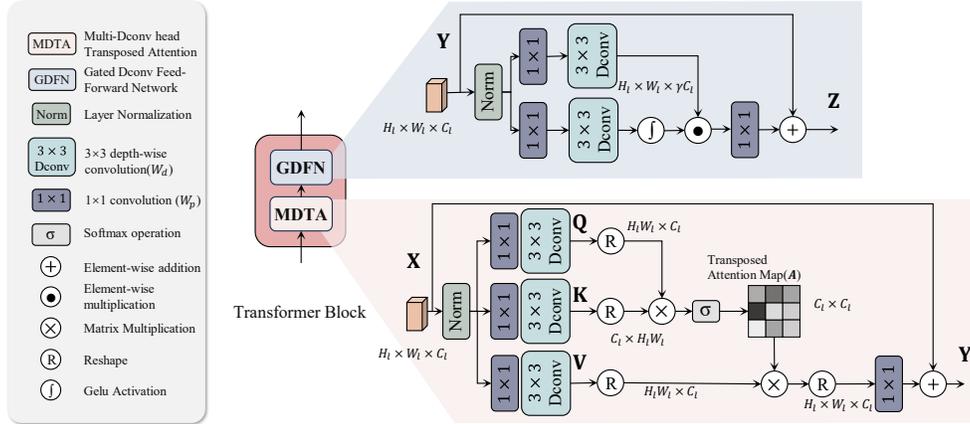


Figure 14: Overview of the Transformer block used in the Prompt block. The Transformer block is composed of two sub modules, the Multi Dconv head transposed attention module (MDTA) and the Gated Dconv feed-forward network (GDFN).

MDTA Module. Let the input feature map be $\mathbf{X} \in \mathbb{R}^{H_i \times W_i \times C_i}$. MDTA first applies Layer Normalization, then projects the input to query (\mathbf{Q}), key (\mathbf{K}), and value (\mathbf{V}) tensors using a sequence of 1×1 pointwise convolution followed by 3×3 depth-wise convolution, all bias-free. To enable channel-wise attention, \mathbf{Q} and \mathbf{K} are reshaped to $\mathbb{R}^{H_i W_i \times C_i}$ and $\mathbb{R}^{C_i \times H_i W_i}$ respectively, resulting in a transposed attention map of shape $C_i \times C_i$ via dot-product interaction. Multi-head computation is performed in parallel.

1458 **GDFN Module.** The GDFN submodule begins with a channel expansion using a 1×1 convolution
 1459 by a factor γ . The expanded features are split into two parallel branches, each followed by a 3×3
 1460 depth-wise convolution. One branch passes through a GeLU activation while the other remains
 1461 linear. The outputs are combined via element-wise multiplication, and finally projected back to
 1462 the original channel dimension through a 1×1 convolution. Residual connections are maintained
 1463 throughout the block.

1464 **MDTA** performs self-attention along channels:

$$1466 \mathbf{Y} = W_p \mathbf{V} \cdot \text{Softmax}(\mathbf{K} \cdot \mathbf{Q} / \alpha) + \mathbf{X}$$

1468 **GDFN** transforms the result as:

$$1470 \mathbf{Z} = W_p^0 (\phi(W_d^1 W_p^1(\text{LN}(\mathbf{Y}))) \odot W_d^2 W_p^2(\text{LN}(\mathbf{Y}))) + \mathbf{Y}$$

1472 Here, LN is layer normalization, ϕ denotes GELU activation, and \odot is element-wise multiplication.

1474 H REPRODUCIBILITY STATEMENT

1477 We provide all necessary details to support reproducibility. All experiments are conducted on publicly
 1478 available datasets, and the model architectures, hyperparameters, training protocols, and evaluation
 1479 metrics are specified in the paper. We will release our codebase, training scripts, and pretrained
 1480 checkpoints on GitHub upon acceptance.

1482 I THE USE OF LARGE LANGUAGE MODELS (LLMs)

1484 We used large language models only for light editorial assistance during manuscript preparation
 1485 (grammar and wording refinement, minor style/formatting suggestions). No LLMs were used for
 1486 research ideation, dataset curation, modeling, experiment design, analysis, or drafting substantive
 1487 sections.

1488
 1489
 1490
 1491
 1492
 1493
 1494
 1495
 1496
 1497
 1498
 1499
 1500
 1501
 1502
 1503
 1504
 1505
 1506
 1507
 1508
 1509
 1510
 1511