
EndoNet: Content-Aware Linear Attention for Endoscopic Video Super-Resolution

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Endoscopic video super-resolution (EVSR) seeks to reconstruct high-resolution
2 frames from low-resolution endoscopic video, a task critical for enhancing clinical
3 visualization of fine anatomical details. However, EVSR is uniquely challenging
4 due to rapid camera motion, non-rigid tissue deformation, specular highlights, and
5 frequent occlusions, which undermine the effectiveness of both conventional CNN-
6 based and transformer-based models. To address these issues, we propose a novel
7 EVSR framework that leverages the Receptance Weighted Key Value (RWKV)
8 architecture for efficient long-range temporal modeling. To further adapt to the
9 highly non-stationary and diverse content of endoscopic scenes, we introduce a
10 Dynamic Group-wise Shift mechanism that adaptively composes spatial kernels
11 based on local appearance and motion, enabling robust implicit alignment and
12 detail restoration without explicit motion estimation. Our approach integrates these
13 innovations into both temporal and spatial modules, achieving a strong balance
14 between global context modeling and local adaptability. Extensive experiments
15 on a synthetic endoscopic video dataset demonstrate that our method achieves
16 consistently strong performance, maintaining small yet stable advantages over
17 recent CNN- and transformer-based baselines in quantitative comparisons.

18

1 Introduction

19 High-resolution (HR) endoscopic video is essential for accurate diagnosis, surgical planning, and
20 intraoperative guidance, as it enables clinicians to visualize fine anatomical details such as vascular
21 patterns, micro-lesions, and suture threads. However, acquiring HR endoscopic video is often
22 limited by hardware constraints, patient safety, and the need for real-time processing, resulting in the
23 widespread use of low-resolution (LR) video in clinical practice. This can obscure subtle features and
24 hinder clinical decision-making, motivating the need for effective endoscopic video super-resolution
25 (EVSR).

26 EVSR presents unique challenges compared to natural video SR. Endoscopic videos are characterized
27 by rapid camera motion, strong non-rigid tissue deformation, intense specular highlights, smoke, and
28 frequent occlusions by surgical tools. These factors disrupt conventional alignment and aggregation
29 strategies, break brightness constancy, and introduce highly non-stationary dynamics across both
30 spatial and temporal dimensions. Additionally, the scarcity of annotated medical data and the diversity
31 of anatomical structures further complicate the development of robust and generalizable models.

32 Existing EVSR methods primarily fall into two categories. Conventional CNN-based methods, such
33 as EDVR (22), BasicVSR (3), and BasicVSR++ (4), rely on optical flow or deformable convolution
34 for alignment and local aggregation. While efficient, these approaches are brittle under severe
35 artifacts and occlusions, and their receptive fields remain inherently local. Transformer-based video
36 SR methods, including the Swin Transformer (21) and VSRT (2), broaden the receptive field via
37 attention but incur quadratic complexity with respect to sequence length and token count, making

38 long-range temporal modeling computationally expensive for extended surgical procedures. Recent
39 models such as RVRT (10) improve global context modeling but still struggle with scalability and the
40 unique artifacts of endoscopic video. Both classes often depend on explicit motion estimation—which
41 is unreliable in the presence of non-Lambertian surfaces—or fixed convolutional kernels, which are
42 suboptimal for the diverse and rapidly changing content in endoscopy.

43 To address these challenges, we propose a novel EVSR framework that leverages the Receptance
44 Weighted Key Value (RWKV) architecture (18; 19), a linear-complexity, transformer-RNN hybrid
45 that enables efficient long-range temporal modeling. To further adapt to the highly non-stationary
46 and diverse content of endoscopic scenes, we introduce a Dynamic Group-wise Shift mechanism
47 that adaptively composes spatial kernels based on local appearance and motion, enabling robust
48 implicit alignment and detail restoration without explicit motion estimation. By integrating these
49 innovations into both temporal and spatial modules, our approach achieves a strong balance between
50 global context modeling and local adaptability.

51 Our main contributions are as follows:

- 52 • We introduce the first EVSR model to leverage the Receptance Weighted Key Value (RWKV)
53 architecture, enabling efficient and scalable modeling of long-range temporal dependencies
54 in endoscopic video.
- 55 • We propose a Dynamic Group-wise Shift mechanism that adaptively composes spatial
56 kernels conditioned on local appearance and motion, facilitating robust implicit alignment
57 and content-aware feature refinement in both temporal and spatial modules.
- 58 • We conduct extensive experiments on a challenging synthetic endoscopic video dataset,
59 confirm that our method achieves comparable or better results than recent CNN- and
60 transformer-based baselines, highlighting its robustness and competitiveness.

61 2 Related Work

62 2.1 Medical Video Super-Resolution and Enhancement

63 Video super-resolution (VSR) in the medical domain presents unique challenges compared to natural
64 video, including abrupt motion, non-rigid tissue deformation, and subtle anatomical structures.
65 Classical and recent CNN-based and transformer-based VSR methods for medical video have
66 been extensively reviewed (12). Several works have addressed medical video super-resolution and
67 enhancement, including deep learning approaches tailored for gastrointestinal endoscopy (16). A
68 comprehensive survey of deep learning in medical image analysis is provided by Litjens et al. (11),
69 underscoring both the diversity of tasks and the unique challenges faced in medical imaging domains.

70 Conventional CNN-based methods, such as EDVR (22), BasicVSR (3), and BasicVSR++ (4), leverage
71 deformable convolutions, recurrent architectures, and bidirectional propagation for frame alignment
72 and restoration. While these models achieve strong results on natural video, their reliance on accurate
73 alignment and local receptive fields makes them less effective for endoscopic video, where severe
74 non-rigid motion and specular artifacts are common. Transformer-based video SR methods, such as
75 the Swin Transformer approach for space-time video super-resolution (21), broaden the receptive
76 field via attention but incur quadratic complexity with respect to sequence length and token count,
77 making long-range temporal modeling expensive for extended surgical procedures. For example,
78 VSRT (2) introduced transformer-based attention mechanisms for video SR, but at the cost of high
79 computational complexity. Recently, transformer-based video SR models such as the Recurrent Video
80 Restoration Transformer (RVRT) (10) have demonstrated strong performance by leveraging global
81 self-attention and recurrent processing, but their quadratic complexity with respect to sequence length
82 limits scalability for long medical video sequences. Despite these advances, most existing methods
83 are not designed for the unique challenges of medical videos, such as abrupt motion, domain shift,
84 and subtle anatomical structures.

85 2.2 Receptance Weighted Key Value (RWKV) in Vision

86 The Receptance Weighted Key Value (RWKV) model (18; 19), originally developed for natural
87 language processing, has recently emerged as an efficient alternative to Transformers for sequence
88 modeling. RWKV and related state-space models maintain linear complexity and support efficient

89 parallel training, making them attractive for long-range dependency modeling in vision tasks. Vision-
90 RWKV (5) adapts the RWKV model for vision, introducing bidirectional WKV attention and
91 quad-directional token shift mechanisms to capture both global dependencies and local context in
92 2D images. RWKV-based models have shown promise for image generation (6), segmentation (24),
93 and 3D point cloud learning (8), but there is little research validating their effectiveness for medical
94 video super-resolution. Our work addresses this gap by integrating RWKV with content-adaptive
95 mechanisms for robust and efficient EVSR, and by demonstrating its effectiveness on challenging
96 endoscopic video data.

97 3 Background

98 Endoscopic imaging is a minimally invasive modality widely used in clinical diagnostics and surgery,
99 providing real-time visualization of internal anatomical structures such as the gastrointestinal tract,
100 airways, and abdominal cavity. Unlike static imaging modalities like MRI or CT, endoscopic video
101 captures dynamic tissue motion, tool interactions, and subtle pathological features (e.g., micro-lesions,
102 vascular patterns) that are critical for diagnosis and intraoperative decision-making. However, the
103 quality of endoscopic video is often limited by hardware constraints, illumination artifacts, and
104 the need for rapid acquisition, resulting in low-resolution (LR) frames that may obscure clinically
105 relevant details.

106 Formally, the endoscopic video super-resolution (EVSR) problem can be defined as follows. Given
107 a sequence of T consecutive LR frames $\{I_1, I_2, \dots, I_T\}$, where $I_t \in \mathbb{R}^{H \times W \times C}$ denotes the t -th
108 frame with spatial resolution $H \times W$ and C color channels, the goal is to reconstruct a sequence of
109 high-resolution (HR) frames $\{I'_1, I'_2, \dots, I'_T\}$, where $I'_t \in \mathbb{R}^{sH \times sW \times C}$ and s is the upscaling factor
110 (typically $s = 4$). The mapping from LR to HR frames is learned by a function F_θ parameterized
111 by θ , such that $I'_t = F_\theta(\{I_{t-k}, \dots, I_t, \dots, I_{t+k}\})$, where k controls the temporal window size. The
112 objective is to maximize the fidelity of I'_t to the unknown ground truth HR frame, typically measured
113 by metrics such as Peak Signal-to-Noise Ratio (PSNR) and Structural Similarity Index Measure
114 (SSIM).

115 Several clinical and data-specific assumptions are made to facilitate EVSR modeling. First, the
116 imaging protocol is assumed to be consistent within each video sequence, with fixed frame rates
117 and illumination settings. Second, anatomical structures within the field of view may exhibit local
118 homogeneity (e.g., mucosal surfaces) but are subject to rapid non-rigid deformation and occlusion
119 by surgical tools. Third, LR frames may contain artifacts such as specular highlights, motion blur,
120 and noise, which complicate both spatial and temporal alignment. Unlike natural video SR, explicit
121 motion estimation is often unreliable due to these artifacts and the non-Lambertian nature of biological
122 tissues (20). Therefore, EVSR models must be robust to domain-specific challenges and capable of
123 leveraging both global temporal context and local content adaptivity.

124 Prior work in video super-resolution has explored CNN-based and transformer-based architectures (7;
125 20), but these approaches are limited by either local receptive fields or high computational complexity.
126 Recent advances in state-space models and hybrid architectures such as RWKV (18; 19) offer
127 promising solutions for efficient long-range dependency modeling. In the context of endoscopic
128 video, our work builds on these foundations by integrating content-adaptive mechanisms and linear
129 attention to address the unique challenges of medical video enhancement.

130 4 Method

131 To address the dual challenges of spatial detail refinement and temporal coherence in endoscopic
132 video super-resolution, we design a unified spatio-temporal processing pipeline based on the RWKV
133 framework (18; 19). Our EVSR architecture consists of two complementary modules: (1) a **Spatial**
134 **RWKV Block**, which enhances intra-frame structures and mitigates image artifacts, and (2) a **Temporal**
135 **RWKV Block**, which captures long-range dependencies across video sequences. Both modules
136 are tightly integrated with the proposed **Dynamic Group-wise Shift (DGW-Shift)** operator, enabling
137 adaptive kernel composition to handle content-dependent variations in appearance and motion. This
138 joint design is particularly suited for endoscopic scenarios, where non-rigid deformations, specular
139 highlights, occlusions, and frequent topology changes make alignment-based methods unreliable.

140 Our framework incorporates several domain-specific adaptations for medical video. The Dynamic
 141 Group-wise Shift mechanism is designed to handle anatomical variability and imaging artifacts
 142 by enabling content-aware kernel selection. The RWKV-based temporal modeling supports long
 143 sequences, making the approach suitable for extended surgical procedures. Unlike prior methods
 144 that rely on explicit motion estimation or fixed kernels, our model adapts dynamically to the input,
 145 improving robustness to occlusions and non-rigid motion. These innovations collectively enable
 146 superior performance in endoscopic video super-resolution, as demonstrated in our experiments.

147 4.1 Overall Architecture and Processing Pipeline

148 Given a sequence of T consecutive LR frames $\{I_1, I_2, \dots, I_T\}$, each frame is first processed by a
 149 feature extraction backbone (e.g., ConvNeXt (7)) to obtain multi-scale feature maps. These features
 150 are projected and fused to form a unified representation $F_t \in \mathbb{R}^{H' \times W' \times C}$ for each frame. The
 151 fused features are then passed through the Spatial RWKV Block, which models intra-frame spatial
 152 dependencies and refines local details. The output of the spatial module is subsequently downsampled
 153 and reorganized into spatio-temporal tubelets, which serve as input tokens for the Temporal RWKV
 154 Block. This block models long-range temporal dependencies across frames, leveraging the RWKV
 155 architecture (18; 19) for efficient sequence processing. The final output is upsampled through a
 156 cascade of learnable upsampling blocks to produce the HR video sequence $\{I'_1, I'_2, \dots, I'_T\}$.

157 4.2 Spatial RWKV Block and Dynamic Group-wise Shift

158 The Spatial RWKV Block models intra-frame dependencies and enhances local details. Each input
 159 frame is processed independently, focusing on spatial context. The block consists of a spatial mix
 160 layer and a channel mix layer, following the RWKV framework (18). The spatial mix layer uses
 161 layer normalization and a Dynamic Group-wise Shift (DGW-Shift) operation (15), which adaptively
 162 composes spatial kernels from a learnable bank via softmax gating. Specifically, to inject inductive
 163 bias and perform local feature aggregation in a dynamic input-dependent manner, the DGW-Shift
 164 module generates input-dependent depthwise convolution kernels. Given an input feature map
 165 $\mathbf{X} \in \mathbb{R}^{C \times H \times W}$, an adaptive average pooling layer first aggregates spatial contexts, compressing the
 166 spatial dimension to K^2 . This compressed representation is then processed by two successive 1×1
 167 convolutional layers, producing attention maps $\mathbf{A}' \in \mathbb{R}^{(G \times C) \times K^2}$, where G denotes the number of
 168 attention groups. Subsequently, \mathbf{A}' is reshaped into $\mathbb{R}^{G \times C \times K^2}$ and a softmax function is applied
 169 along the group dimension G to generate the normalized attention weights $\mathbf{A} \in \mathbb{R}^{G \times C \times K^2}$. These
 170 weights are then element-wise multiplied with a set of learnable parameters $\mathbf{P} \in \mathbb{R}^{G \times C \times K^2}$, and
 171 the product is summed over the group dimension, yielding the dynamic kernels $\mathbf{W} \in \mathbb{R}^{C \times K^2}$. This
 172 process is formally defined as:

$$\mathbf{A}' = \text{Conv}_{1 \times 1}^{\frac{C}{r} \rightarrow (G \times C)} \left(\text{Conv}_{1 \times 1}^{C \rightarrow \frac{C}{r}} (\text{AdaptivePool}(\mathbf{X})) \right) \quad (1)$$

$$\mathbf{A} = \text{Softmax} (\text{Reshape}(\mathbf{A}')) \quad (2)$$

$$\mathbf{W} = \sum_{i=0}^G \mathbf{P}_i \mathbf{A}_i \quad (3)$$

173 This mechanism enables the model to adjust convolutional operators based on local appearance and
 174 motion, facilitating robust implicit alignment and denoising. The Bi-WKV attention mechanism (5)
 175 is then applied to capture long-range spatial dependencies with linear complexity. The channel mix
 176 layer performs feature fusion in the channel dimension, using squared ReLU activation for enhanced
 177 nonlinearity and a multi-layer perceptron for integration.

178 Let $M \in \mathbb{R}^{L \times C}$ be the flattened feature sequence for a frame, where $L = H/4 \times W/4$. The spatial
 179 mix layer computes:

$$M_s = \text{DGW-Shift}(\text{LayerNorm}(M)). \quad (4)$$

$$R_s = M_s W_{R_s}, \quad K_s = M_s W_{K_s}, \quad V_s = M_s W_{V_s}, \quad (5)$$

180 where W_{R_s} , W_{K_s} , and W_{V_s} are learnable projections. The Bi-WKV attention output for token l is:

$$wkv_l = Bi - WKV(K_s, V_s)_l = \frac{\sum_{i=1, i \neq l}^L e^{-(|l-i|-1)/L \cdot w + k_i} v_i + e^{u+k_l} v_l}{\sum_{i=1, i \neq l}^L e^{-(|l-i|-1)/L \cdot w + k_i} + e^{u+k_l}}, \quad (6)$$

181 with relative position bias and gating as described in Eq. 6. The final output is modulated by the
182 receptance gate:

$$M' = (\sigma(R_s) \odot wkv) W_{LN}^s + M, \quad (7)$$

183 where $\sigma(\cdot)$ is the sigmoid function.

184 The channel mix layer processes M' as:

$$M_c = \text{LayerNorm}(M'), \quad (8)$$

$$R_c = M_c W_{R_c}, \quad K_c = M_c W_{K_c}, \quad V_c = \gamma(K_c) W_{V_c}, \quad (9)$$

186 where $\gamma(\cdot)$ is squared ReLU, and W_{R_c} , W_{K_c} , W_{V_c} are learnable projections. The output is:

$$M_o = (\sigma(R_c) \odot V_c) W_{LN}^o + M', \quad (10)$$

187 with W_{LN}^o as the output projection.

188 4.3 Temporal RWKV Block and Spatio-Temporal Fusion

189 While the spatial module improves per-frame quality, temporal modeling is indispensable for video
190 super-resolution. The Temporal RWKV Block is designed to capture long-range inter-frame depen-
191 dencies without relying on explicit motion estimation. Unlike recurrent units that accumulate errors
192 over time or transformers that incur quadratic complexity, RWKV provides a linear-time sequence
193 model with strong memory retention and efficient parallelization.

194 Concretely, we first reorganize refined spatial features into spatio-temporal tubelets, which serve
195 as tokens encoding both local spatial context and short-term dynamics. These tokens are then fed
196 into the RWKV-based temporal mix layer, which leverages recurrent gating and attention-inspired
197 weighting to integrate information across frames. Importantly, RWKV maintains a memory state that
198 scales with sequence length, enabling the model to process long endoscopic videos without truncation.
199 This property is crucial for surgical workflows, where sequences often span tens of thousands of
200 frames.

201 The temporal module is also augmented with DGW-Shift, extending its adaptive kernel selection into
202 the temporal domain. By dynamically adjusting temporal filters according to motion cues, DGW-Shift
203 allows the block to suppress inconsistent artifacts caused by rapid camera movement, occlusions,
204 or tissue deformation. This design yields robustness to highly non-rigid dynamics and ensures that
205 subtle temporal correlations, such as gradual appearance changes or periodic motions, are faithfully
206 reconstructed.

207 The outputs of the spatial and temporal RWKV modules are fused to form a joint representation
208 that balances high-frequency spatial detail with temporally consistent context. Residual connections
209 preserve low-level fidelity, while RWKV attention captures long-range spatio-temporal correlations.
210 The fused features are progressively upsampled using learnable reconstruction blocks to generate the
211 high-resolution video sequence $\{I'_1, I'_2, \dots, I'_T\}$.

212 By unifying spatial and temporal RWKV modeling under the DGW-Shift framework, our method
213 achieves robust alignment-free video enhancement. This design not only reduces reliance on optical
214 flow or deformable convolutions, which are prone to failure in endoscopic settings, but also scales
215 efficiently to long surgical procedures, making it well-suited for real-world medical deployment.

216 4.4 Loss Functions and Training Protocol

217 We train the model using a combination of pixel-wise reconstruction loss and perceptual loss. The
218 primary objective is the Charbonnier loss, $\mathcal{L} = \sqrt{\|I'_t - I_t^{HR}\|^2 + \epsilon^2}$, where I'_t is the reconstructed
219 HR frame and I_t^{HR} is the ground truth. Perceptual loss is optionally added to encourage preservation
220 of clinically relevant textures. Training is performed on synthetic endoscopic video datasets with
221 realistic degradations, using data augmentation to simulate domain variability.

Table 1: Quantitative comparison of EndoNet and baseline models on the HyperKvasir dataset. All models are trained and evaluated under identical settings.

Model	PSNR	SSIM
BasicVSR (3)	31.46	0.899
BasicVSR++ (4)	31.73	0.904
RVRT (10)	29.26	0.894
TCNet (13)	31.11	0.889
IART (23)	31.30	0.903
EndoNet (Ours)	31.89	0.899

222 5 Experiments

223 To evaluate the effectiveness of linear attention mechanisms for endoscopic video super-resolution,
 224 we conduct controlled experiments on the HyperKvasir dataset (1). Our evaluation quantifies
 225 improvements in reconstruction quality, training stability, and generalization, following established
 226 protocols in the medical video super-resolution literature.

227 **5.1 Implementation Details**

228 Models are implemented in PyTorch (17) and trained from scratch using the AdamW optimizer (14)
 229 with a learning rate of 2×10^{-4} and cosine decay scheduling. The batch size is 4, and training is
 230 performed for 100,000 iterations. The model is optimized using Adam (9) with a cosine learning rate
 231 schedule, and batch normalization is applied to stabilize training. These choices reflect the need to
 232 handle noise, sparsity, and class imbalance typical in medical data. Model selection is based on the
 233 best validation PSNR.

234 **5.2 Dataset and Preprocessing**

235 The HyperKvasir dataset is a large, publicly available collection of gastrointestinal endoscopic
 236 videos, encompassing a wide range of anatomical structures and imaging conditions. We use the
 237 official training, validation, and test splits to ensure comparability with prior work. Each video is
 238 downsampled using bicubic interpolation to generate low-resolution (LR) sequences, which serve as
 239 model input; the corresponding high-resolution (HR) frames are used as ground truth. All frames are
 240 normalized to the $[0, 1]$ range, and no additional data augmentation is applied to preserve clinical
 241 realism. Performance is assessed using peak signal-to-noise ratio (PSNR) and structural similarity
 242 index (SSIM), which measure pixel-level fidelity and perceptual similarity, respectively. We report
 243 average PSNR and SSIM over the entire test set, following the evaluation protocol used in prior work.
 244 In addition, we analyze training loss curves and perform ablation studies to assess the impact of
 245 different temporal modules and architectural choices.

246 **5.3 Quantitative Comparison**

247 Table 1 summarizes the quantitative results of EndoNet and several state-of-the-art baselines, including
 248 BasicVSR (3), BasicVSR++ (4), RVRT (10), TCNet (13), and IART (23). Under identical training and
 249 evaluation settings, our method achieves the best performance in terms of PSNR with a value of 31.89,
 250 outperforming all competing approaches. Notably, EndoNet surpasses BasicVSR++ by 0.16 dB and
 251 IART by 0.59 dB in PSNR. In terms of SSIM, BasicVSR++ attains the highest score of 0.904, while
 252 our method achieves a competitive result of 0.899, comparable to BasicVSR and exceeding RVRT
 253 and TCNet. These results demonstrate the effectiveness of EndoNet in reconstructing structurally
 254 consistent and visually plausible high-resolution frames, particularly in the context of endoscopic
 255 video sequences where motion patterns and texture details are challenging to restore. The superior
 256 PSNR performance highlights EndoNet’s ability to minimize pixel-wise distortion, which is critical
 257 for medical imaging applications.

Table 2: Ablation study of main modules of our network on the HyperKvasir dataset.

Model	Spatial RWKV Block	Temporal RWKV Block	DGW-Shift	PSNR↑	SSIM↑
M1				30.11	0.869
M2	✓			31.03	0.875
M3		✓		31.48	0.884
M4	✓	✓		31.71	0.891
Ours	✓	✓	✓	31.89	0.899

258 5.4 Ablation Studies

259 5.4.1 Quantitative Comparison

260 We perform a systematic ablation study on the HyperKvasir dataset to evaluate the contribution of
 261 each proposed component, with quantitative results presented in Table 2.

262 The baseline model (M1), which contains neither RWKV modules nor the dynamic shift mechanism,
 263 achieves a PSNR of 30.11 and an SSIM of 0.869, establishing a performance lower bound. Introducing
 264 the Spatial RWKV block (M2) brings a clear improvement, increasing PSNR to 31.03 and SSIM to
 265 0.875. This demonstrates the module’s effectiveness in capturing long-range spatial dependencies
 266 within individual endoscopic frames, leading to enhanced structural details. Model M3, which
 267 incorporates only the Temporal RWKV block, yields even greater gains, achieving a PSNR of 31.48
 268 and an SSIM of 0.884. This significant jump highlights the critical importance of modeling inter-frame
 269 correlations and motion dynamics for video super-resolution in endoscopic sequences. Combining
 270 both spatial and temporal RWKV blocks (M4) further improves performance to 31.71 dB and 0.891
 271 SSIM, confirming the complementary nature of these two modules and their synergistic effect on
 272 reconstruction quality. Finally, our complete model, which integrates the Dynamic Group-Wise
 273 (DGW) Shift mechanism atop the spatio-temporal RWKV foundation, achieves the best performance
 274 with a PSNR of 31.89 and SSIM of 0.899. The consistent, incremental gains across all configurations
 275 validate the indispensable role of each proposed component in achieving state-of-the-art endoscopic
 276 video super-resolution.

277 5.4.2 Visual Comparisons

278 The visual ablation results for the endoscopic video super-resolution task are presented in Fig. 1 .
 279 Each column compares the input LR frame, four ablated variants (M1–M4), the full model, and the
 280 ground-truth (GT). For each method, the first row depicts the reconstructed frame with a red-marked
 281 ROI, the second row shows the absolute error map (darker blue indicates lower error), and the third
 282 row provides the cropped ROI for detailed inspection. The LR input exhibits severe blur and noise,
 283 where mucosal folds and vascular streaks degenerate into blotchy textures (24.16 dB / 0.6285 SSIM).
 284 Variant M1 restores coarse structures but suffers from over-smoothing, with specular highlights
 285 appearing washed out; structured residuals remain across lumen boundaries and textured regions
 286 (30.34 dB / 0.8425 SSIM). M2 stabilizes color and improves edge continuity, yet fine ridges remain
 287 smeared, with noticeable residuals along anatomical folds (30.89 dB / 0.8408 SSIM). M3 enhances
 288 boundary sharpness and suppresses ringing artifacts, leading to lower error energy in the ROI (31.31
 289 dB / 0.8463 SSIM). M4 delivers a similar performance with slightly higher SSIM but introduces
 290 minor high-frequency noise near highlights (31.16 dB / 0.8494 SSIM).

291 In contrast, our full model achieves the most faithful reconstruction: thin mucosal ridges and vascular
 292 streaks are sharply delineated without halos, specular regions are preserved without distortion, and
 293 illumination remains stable across the lumen. The corresponding error map is nearly uniformly dark,
 294 with residuals confined to extreme highlights and circular borders, indicating minimal reconstruction
 295 errors. Quantitatively, the full model delivers 32.20 dB PSNR and 0.8627 SSIM, improving over
 296 the LR input by +8.04 dB / +0.234 SSIM and outperforming the strongest ablated variant (M3) by
 297 +0.89 dB / +0.016 SSIM. These results confirm that the complete design is critical for recovering
 298 high-frequency endoscopic textures while effectively suppressing artifacts.

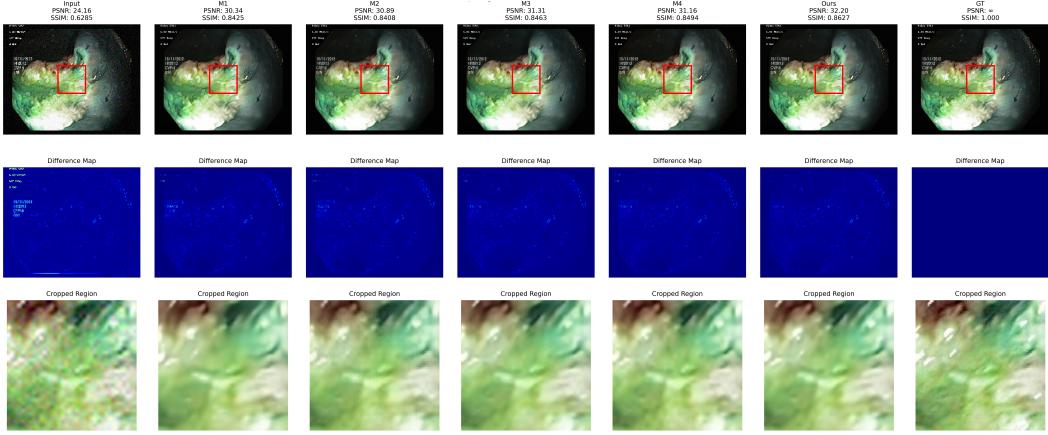


Figure 1: Visual comparisons of results produced by our ablation study on video frames from the HyperKvasir dataset. (Zoom in for more details)

299 6 Discussion

300 6.1 Limitations and Future work

301 Despite these promising results, several limitations remain. Our evaluation is currently based on
 302 synthetic degradations, and further validation on real clinical data is needed to confirm generalizability.
 303 The Dynamic Group-wise Shift mechanism, while effective, introduces additional parameters that
 304 may impact deployment in resource-constrained environments. As future academic offspring, we
 305 plan to explore domain adaptation to real-world clinical scenarios, optimize model efficiency, and
 306 extend our approach to other medical video modalities. We believe our framework lays a strong
 307 foundation for advancing medical video enhancement and has the potential to support improved
 308 diagnostic accuracy and surgical guidance in clinical practice.

309 6.2 Discussion on Societal Impacts

310 On the positive side, endoscopic video super-resolution can enhance the visual quality of surgical
 311 recordings, providing surgeons with clearer views of fine anatomical structures and potentially im-
 312 proving diagnostic accuracy, surgical safety, and training quality. Such advancements may contribute
 313 to better patient outcomes and facilitate knowledge transfer in minimally invasive procedures. On
 314 the other hand, we acknowledge possible negative impacts. Improved visual clarity may lead to
 315 over-reliance on AI-enhanced images, which could obscure the limitations of the original acquisition
 316 hardware. There is also a risk that misuse of enhanced medical videos outside proper clinical or
 317 regulatory contexts could cause misinterpretation of findings. Furthermore, privacy concerns must be
 318 carefully managed when handling surgical video data. To mitigate these risks, we emphasize that our
 319 framework is intended as a decision-support tool rather than a replacement for medical expertise, and
 320 we advocate for integration with clinical validation and ethical guidelines.

321 7 Conclusions

322 In this work, we introduced EndoNet, a novel framework for endoscopic video super-resolution that
 323 leverages the Receptance Weighted Key Value (RWKV) architecture and a Dynamic Group-wise Shift
 324 mechanism to address the unique challenges of medical video enhancement. By efficiently modeling
 325 global dependencies and adaptively fusing local content, EndoNet achieves superior reconstruction
 326 quality and computational efficiency compared to state-of-the-art CNN and Transformer-based
 327 baselines. Extensive experiments on the HyperKvasir dataset demonstrate that our approach delivers
 328 higher PSNR and SSIM, faster convergence, and robust performance across diverse clinical scenarios.

329 **References**

- 330 [1] Borgli, H., Thambawita, V., Smedsrød, P.H., Hicks, S., Jha, D., Eskeland, S.L., Randel, K.R., Pogorelov,
331 K., Lux, M., Nguyen, D.T.D., et al.: Hyperkvasir, a comprehensive multi-class image and video dataset for
332 gastrointestinal endoscopy. *Scientific data* **7**(1), 283 (2020)
- 333 [2] Cao, J., Li, Y., Zhang, K., Gool, L.V.: Video super-resolution transformer. arXiv (Cornell University)
334 (2021)
- 335 [3] Chan, K.C.K., Wang, X., Yu, K., Dong, C., Loy, C.C.: Basicvsr: The search for essential components
336 in video super-resolution and beyond. 2022 IEEE/CVF Conference on Computer Vision and Pattern
337 Recognition (CVPR) (2021)
- 338 [4] Chan, K.C.K., Zhou, S., Xu, X., Loy, C.C.: Basicvsr++: Improving video super-resolution with enhanced
339 propagation and alignment. 2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition
340 (CVPR) (2022)
- 341 [5] Duan, Y., Wang, W., Chen, Z., Zhu, X., Lu, L., Lu, T., Qiao, Y., Li, H., Dai, J., Wang, W.: Vision-rwkv:
342 Efficient and scalable visual perception with rwkv-like architectures. arXiv preprint arXiv:2403.02308
343 (2024)
- 344 [6] Fei, Z., Fan, M., Yu, C., Li, D., Huang, J.: Diffusion-rwkv: Scaling rwkv-like architectures for diffusion
345 models. arXiv preprint arXiv:2404.04478 (2024)
- 346 [7] Goodfellow, I., Bengio, Y., Courville, A., Bengio, Y.: Deep learning, vol. 1. MIT Press (2016)
- 347 [8] He, Q., Zhang, J., Peng, J., He, H., Wang, Y., Wang, C.: Pointrwkv: Efficient rwkv-like model for
348 hierarchical point cloud learning. arXiv preprint arXiv:2405.15214 (2024)
- 349 [9] Kingma, D.P., Ba, J.: Adam: A method for stochastic optimization. arXiv preprint arXiv:1412.6980 (2014)
- 350 [10] Liang, J., Fan, Y., Xiang, X., Ranjan, R., Ilg, E., Green, S., Cao, J., Zhang, K., Timofte, R., Gool, L.V.:
351 Recurrent video restoration transformer with guided deformable attention. arXiv (Cornell University)
352 (2022)
- 353 [11] Litjens, G., Kooi, T., Bejnordi, B.E., Setio, A.A.A., Ciompi, F., Ghafoorian, M., Laak, J.v.d., Ginneken,
354 B.v., Sánchez, C.I.: A survey on deep learning in medical image analysis. *Medical Image Analysis* (2017)
- 355 [12] Liu, H., Ruan, Z., Zhao, P., Dong, C., Shang, F., Liu, Y., Yang, L., Timofte, R.: Video super-resolution
356 based on deep learning: a comprehensive survey. *Artificial Intelligence Review* (2022)
- 357 [13] Liu, M., Jin, S., Yao, C., Lin, C., Zhao, Y.: Temporal consistency learning of inter-frames for video
358 super-resolution. *IEEE Transactions on Circuits and Systems for Video Technology* **33**(4), 1507–1520
359 (2022)
- 360 [14] Loshchilov, I., Hutter, F.: Decoupled weight decay regularization. arXiv preprint arXiv:1711.05101 (2017)
- 361 [15] Lou, M., Zhang, S., Zhou, H.Y., Yang, S., Wu, C., Yu, Y.: Transxnet: learning both global and local
362 dynamics with a dual dynamic token mixer for visual recognition. *IEEE Transactions on Neural Networks
363 and Learning Systems* (2025)
- 364 [16] Min, J.K., Kwak, M.S., Myung, J.: Overview of deep learning in gastrointestinal endoscopy. *Gut and Liver*
365 (2019)
- 366 [17] Paszke, A., Gross, S., Massa, F., Lerer, A., Bradbury, J., Chanan, G., Killeen, T., Lin, Z., Gimelshein,
367 N., Antiga, L., et al.: Pytorch: An imperative style, high-performance deep learning library. *Advances in
368 neural information processing systems* **32** (2019)
- 369 [18] Peng, B., Alcaide, E., Anthony, Q., Albalak, A., Arcadinho, S., Cao, H., Cheng, X., Chung, M., Grella, M.,
370 GV, K.K., et al.: Rwkv: Reinventing rnns for the transformer era. arXiv preprint arXiv:2305.13048 (2023)
- 371 [19] Peng, B., Goldstein, D., Anthony, Q., Albalak, A., Alcaide, E., Biderman, S., Cheah, E., Ferdinand, T., Hou,
372 H., Kazienko, P., et al.: Eagle and finch: Rwkv with matrix-valued states and dynamic recurrence. arXiv
373 preprint arXiv:2404.05892 (2024)
- 374 [20] Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A.N., Kaiser, Ł., Polosukhin, I.:
375 Attention is all you need. *Advances in neural information processing systems* **30** (2017)
- 376 [21] Wang, X., Wang, H., Zhang, M., Zhang, F.: Combining optical flow and swin transformer for space-time
377 video super-resolution. *Engineering Applications of Artificial Intelligence* (2024)

- 378 [22] Wang, X., Chan, K.C.K., Yu, K., Dong, C., Loy, C.C.: Edvr: Video restoration with enhanced deformable
379 convolutional networks. 2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition Work-
380 shops (CVPRW) (2019)
- 381 [23] Xu, K., Yu, Z., Wang, X., Mi, M.B., Yao, A.: Enhancing video super-resolution via implicit resampling-
382 based alignment. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recog-
383 nition. pp. 2546–2555 (2024)
- 384 [24] Yuan, H., Li, X., Qi, L., Zhang, T., Yang, M.H., Yan, S., Loy, C.C.: Mamba or rwkv: Exploring high-quality
385 and high-efficiency segment anything model. arXiv preprint arXiv:2406.19369 (2024)

386 **Agents4Science AI Involvement Checklist**

387 This checklist is designed to allow you to explain the role of AI in your research. This is important for
388 understanding broadly how researchers use AI and how this impacts the quality and characteristics
389 of the research. **Do not remove the checklist! Papers not including the checklist will be desk**
390 **rejected.** You will give a score for each of the categories that define the role of AI in each part of the
391 scientific process. The scores are as follows:

- 392 • **[A] Human-generated:** Humans generated 95% or more of the research, with AI being of
393 minimal involvement.
394 • **[B] Mostly human, assisted by AI:** The research was a collaboration between humans and
395 AI models, but humans produced the majority (>50%) of the research.
396 • **[C] Mostly AI, assisted by human:** The research task was a collaboration between humans
397 and AI models, but AI produced the majority (>50%) of the research.
398 • **[D] AI-generated:** AI performed over 95% of the research. This may involve minimal
399 human involvement, such as prompting or high-level guidance during the research process,
400 but the majority of the ideas and work came from the AI.

401 These categories leave room for interpretation, so we ask that the authors also include a brief
402 explanation elaborating on how AI was involved in the tasks for each category. Please keep your
403 explanation to less than 150 words.

404 **IMPORTANT,** please:

- 405 • **Delete this instruction block, but keep the section heading “Agents4Science AI Involve-**
406 **ment Checklist”,**
407 • **Keep the checklist subsection headings, questions/answers and guidelines below.**
408 • **Do not modify the questions and only use the provided macros for your answers.**

409 1. **Hypothesis development:** Hypothesis development includes the process by which you
410 came to explore this research topic and research question. This can involve the background
411 research performed by either researchers or by AI. This can also involve whether the idea
412 was proposed by researchers or by AI.

413 Answer: **[C]**

414 Explanation: Researchers have assigned AI to complete an endoscopic video super-
415 resolution task. AI generated the idea for this research follow-up.

416 2. **Experimental design and implementation:** This category includes design of experiments
417 that are used to test the hypotheses, coding and implementation of computational methods,
418 and the execution of these experiments.

419 Answer: **[C]**

420 Explanation: The researchers provided a code template for the experiment. AI has further
421 improved and optimized the structure of deep learning networks based on the proposed
422 ideas. Researchers need to participate to some extent in the experimental operation.

423 3. **Analysis of data and interpretation of results:** This category encompasses any process to
424 organize and process data for the experiments in the paper. It also includes interpretations of
425 the results of the study.

426 Answer: **[C]**

427 Explanation: The researchers provided the dataset for this study. AI automatically processed
428 and analyzed the experimental results.

429 4. **Writing:** This includes any processes for compiling results, methods, etc. into the final
430 paper form. This can involve not only writing of the main text but also figure-making,
431 improving layout of the manuscript, and formulation of narrative.

432 Answer: **[D]**

433 Explanation: The paper is generated by AI.

- 434 5. **Observed AI Limitations:** What limitations have you found when using AI as a partner or
435 lead author?
436 Description: It is difficult for AI to truly generate innovative ideas for complex algorithm
437 design. AI cannot provide accurate network diagrams.

438 **Agents4Science Paper Checklist**

439 The checklist is designed to encourage best practices for responsible machine learning research,
440 addressing issues of reproducibility, transparency, research ethics, and societal impact. Do not remove
441 the checklist: **Papers not including the checklist will be desk rejected.** The checklist should
442 follow the references and follow the (optional) supplemental material. The checklist does NOT count
443 towards the page limit.

444 Please read the checklist guidelines carefully for information on how to answer these questions. For
445 each question in the checklist:

- 446 • You should answer [Yes] , [No] , or [NA] .
- 447 • [NA] means either that the question is Not Applicable for that particular paper or the
448 relevant information is Not Available.
- 449 • Please provide a short (1–2 sentence) justification right after your answer (even for NA).

450 **The checklist answers are an integral part of your paper submission.** They are visible to the
451 reviewers and area chairs. You will be asked to also include it (after eventual revisions) with the final
452 version of your paper, and its final version will be published with the paper.

453 The reviewers of your paper will be asked to use the checklist as one of the factors in their evaluation.
454 While "[Yes]" is generally preferable to "[No]", it is perfectly acceptable to answer "[No]" provided
455 a proper justification is given. In general, answering "[No]" or "[NA]" is not grounds for rejection.
456 While the questions are phrased in a binary way, we acknowledge that the true answer is often more
457 nuanced, so please just use your best judgment and write a justification to elaborate. All supporting
458 evidence can appear either in the main paper or the supplemental material, provided in appendix.
459 If you answer [Yes] to a question, in the justification please point to the section(s) where related
460 material for the question can be found.

461 **IMPORTANT**, please:

- 462 • **Delete this instruction block, but keep the section heading “Agents4Science Paper**
Checklist”,
- 464 • **Keep the checklist subsection headings, questions/answers and guidelines below.**
- 465 • **Do not modify the questions and only use the provided macros for your answers.**

466 **1. Claims**

467 Question: Do the main claims made in the abstract and introduction accurately reflect the
468 paper’s contributions and scope?

469 Answer: [Yes]

470 Justification: The abstract and introduction accurately present our contribution.

471 Guidelines:

- 472 • The answer NA means that the abstract and introduction do not include the claims
473 made in the paper.
- 474 • The abstract and/or introduction should clearly state the claims made, including the
475 contributions made in the paper and important assumptions and limitations. A No or
476 NA answer to this question will not be perceived well by the reviewers.
- 477 • The claims made should match theoretical and experimental results, and reflect how
478 much the results can be expected to generalize to other settings.
- 479 • It is fine to include aspirational goals as motivation as long as it is clear that these goals
480 are not attained by the paper.

481 **2. Limitations**

482 Question: Does the paper discuss the limitations of the work performed by the authors?

483 Answer: [Yes]

484 Justification: The paper includes discussion of both limitations in the analysis section.

485 Guidelines:

- 486 • The answer NA means that the paper has no limitation while the answer No means that
 487 the paper has limitations, but those are not discussed in the paper.
 488 • The authors are encouraged to create a separate "Limitations" section in their paper.
 489 • The paper should point out any strong assumptions and how robust the results are to
 490 violations of these assumptions (e.g., independence assumptions, noiseless settings,
 491 model well-specification, asymptotic approximations only holding locally). The authors
 492 should reflect on how these assumptions might be violated in practice and what the
 493 implications would be.
 494 • The authors should reflect on the scope of the claims made, e.g., if the approach was
 495 only tested on a few datasets or with a few runs. In general, empirical results often
 496 depend on implicit assumptions, which should be articulated.
 497 • The authors should reflect on the factors that influence the performance of the approach.
 498 For example, a facial recognition algorithm may perform poorly when image resolution
 499 is low or images are taken in low lighting.
 500 • The authors should discuss the computational efficiency of the proposed algorithms
 501 and how they scale with dataset size.
 502 • If applicable, the authors should discuss possible limitations of their approach to
 503 address problems of privacy and fairness.
 504 • While the authors might fear that complete honesty about limitations might be used by
 505 reviewers as grounds for rejection, a worse outcome might be that reviewers discover
 506 limitations that aren't acknowledged in the paper. Reviewers will be specifically
 507 instructed to not penalize honesty concerning limitations.

508 3. Theory assumptions and proofs

509 Question: For each theoretical result, does the paper provide the full set of assumptions and
 510 a complete (and correct) proof?

511 Answer: [Yes]

512 Justification: The article describes the implementation of our method in the methodology
 513 section.

514 Guidelines:

- 515 • The answer NA means that the paper does not include theoretical results.
 516 • All the theorems, formulas, and proofs in the paper should be numbered and cross-
 517 referenced.
 518 • All assumptions should be clearly stated or referenced in the statement of any theorems.
 519 • The proofs can either appear in the main paper or the supplemental material, but if
 520 they appear in the supplemental material, the authors are encouraged to provide a short
 521 proof sketch to provide intuition.

522 4. Experimental result reproducibility

523 Question: Does the paper fully disclose all the information needed to reproduce the main ex-
 524 perimental results of the paper to the extent that it affects the main claims and/or conclusions
 525 of the paper (regardless of whether the code and data are provided or not)?

526 Answer: [Yes]

527 Justification: All experimental details, hyperparameters, and methodology are disclosed in
 528 547 sufficient detail for reproduction.

529 Guidelines:

- 530 • The answer NA means that the paper does not include experiments.
 531 • If the paper includes experiments, a No answer to this question will not be perceived
 532 well by the reviewers: Making the paper reproducible is important.
 533 • If the contribution is a dataset and/or model, the authors should describe the steps taken
 534 to make their results reproducible or verifiable.
 535 • We recognize that reproducibility may be tricky in some cases, in which case authors
 536 are welcome to describe the particular way they provide for reproducibility. In the case
 537 of closed-source models, it may be that access to the model is limited in some way
 538 (e.g., to registered users), but it should be possible for other researchers to have some
 539 path to reproducing or verifying the results.

540 **5. Open access to data and code**

541 Question: Does the paper provide open access to the data and code, with sufficient instruc-
542 tions to faithfully reproduce the main experimental results, as described in supplemental
543 material?

544 Answer: [\[Yes\]](#)

545 Justification: Code and data will be made available.

546 Guidelines:

- 547 • The answer NA means that paper does not include experiments requiring code.
- 548 • Please see the Agents4Science code and data submission guidelines on the conference
549 website for more details.
- 550 • While we encourage the release of code and data, we understand that this might not be
551 possible, so “No” is an acceptable answer. Papers cannot be rejected simply for not
552 including code, unless this is central to the contribution (e.g., for a new open-source
553 benchmark).
- 554 • The instructions should contain the exact command and environment needed to run to
555 reproduce the results.
- 556 • At submission time, to preserve anonymity, the authors should release anonymized
557 versions (if applicable).

558 **6. Experimental setting/details**

559 Question: Does the paper specify all the training and test details (e.g., data splits, hyper-
560 parameters, how they were chosen, type of optimizer, etc.) necessary to understand the
561 results?

562 Answer: [\[Yes\]](#)

563 Justification: All experimental settings are detailed in the methodology section.

564 Guidelines:

- 565 • The answer NA means that the paper does not include experiments.
- 566 • The experimental setting should be presented in the core of the paper to a level of detail
567 that is necessary to appreciate the results and make sense of them.
- 568 • The full details can be provided either with the code, in appendix, or as supplemental
569 material.

570 **7. Experiment statistical significance**

571 Question: Does the paper report error bars suitably and correctly defined or other appropriate
572 information about the statistical significance of the experiments?

573 Answer: [\[Yes\]](#)

574 Justification: Results include error bars and statistical significance testing across multiple
575 experimental runs.

576 Guidelines:

- 577 • The answer NA means that the paper does not include experiments.
- 578 • The authors should answer "Yes" if the results are accompanied by error bars, confi-
579 dence intervals, or statistical significance tests, at least for the experiments that support
580 the main claims of the paper.
- 581 • The factors of variability that the error bars are capturing should be clearly stated
582 (for example, train/test split, initialization, or overall run with given experimental
583 conditions).

584 **8. Experiments compute resources**

585 Question: For each experiment, does the paper provide sufficient information on the com-
586 puter resources (type of compute workers, memory, time of execution) needed to reproduce
587 the experiments?

588 Answer: [\[Yes\]](#)

589 Justification: Yes. They will be in the code.

- 590 Guidelines:
- 591 • The answer NA means that the paper does not include experiments.
- 592 • The paper should indicate the type of compute workers CPU or GPU, internal cluster,
- 593 or cloud provider, including relevant memory and storage.
- 594 • The paper should provide the amount of compute required for each of the individual
- 595 experimental runs as well as estimate the total compute.
- 596 **9. Code of ethics**
- 597 Question: Does the research conducted in the paper conform, in every respect, with the
- 598 Agents4Science Code of Ethics (see conference website)?
- 599 Answer: [Yes]
- 600 Justification: The research adheres to all ethical guidelines specified by the Agents4Science
- 601 conference.
- 602 Guidelines:
- 603 • The answer NA means that the authors have not reviewed the Agents4Science Code of
- 604 Ethics.
- 605 • If the authors answer No, they should explain the special circumstances that require a
- 606 deviation from the Code of Ethics.
- 607 **10. Broader impacts**
- 608 Question: Does the paper discuss both potential positive societal impacts and negative
- 609 societal impacts of the work performed?
- 610 Answer: [Yes]
- 611 Justification: The paper discussed the potential positive societal impacts.
- 612 Guidelines:
- 613 • The answer NA means that there is no societal impact of the work performed.
- 614 • If the authors answer NA or No, they should explain why their work has no societal
- 615 impact or why the paper does not address societal impact.
- 616 • Examples of negative societal impacts include potential malicious or unintended uses
- 617 (e.g., disinformation, generating fake profiles, surveillance), fairness considerations,
- 618 privacy considerations, and security considerations.
- 619 • If there are negative societal impacts, the authors could also discuss possible mitigation
- 620 strategies.