

CVPR  
JUNE 3-7, 2026



DENVER  
COLORADO

# GeneVAR: Causal MeanFlow for Autoregressive Gene-to-WSI Tile Synthesis

Jianwei Zhao<sup>1</sup>, Fan Yang<sup>2</sup>, Xin Li<sup>1,#</sup>, Qiang Zhai<sup>3</sup>, Ao Luo<sup>4</sup>, Ziqi Ren<sup>5</sup>, Zhicheng Jiao<sup>6</sup>, and Hong Cheng<sup>1</sup>

UESTC<sup>1</sup>, ALPACA AI LAB<sup>2</sup>, SICAU<sup>3</sup>, SWJTU<sup>4</sup>, XDU<sup>5</sup>, Brown<sup>6</sup>



# 1. Introduction

## Why WSI:

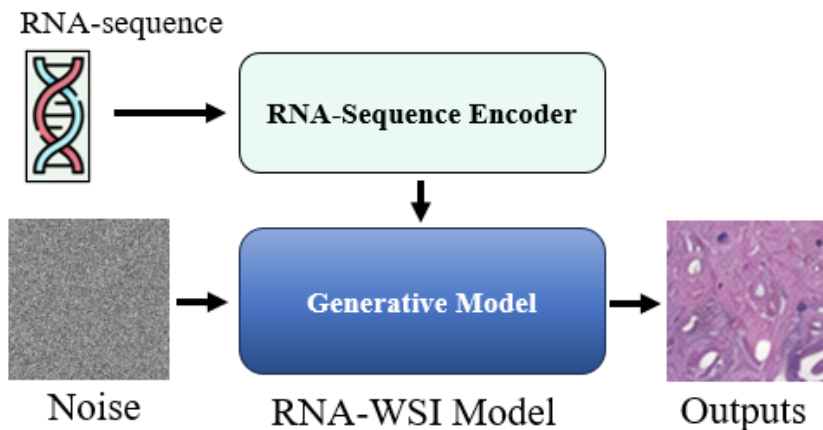
Histological Whole Slide Image (WSI) classification is essential in digital pathology, automating diagnosis and sub-typing while extracting insights from high-resolution scans to support prognosis and treatment planning

## Why Gene-to-WSI:

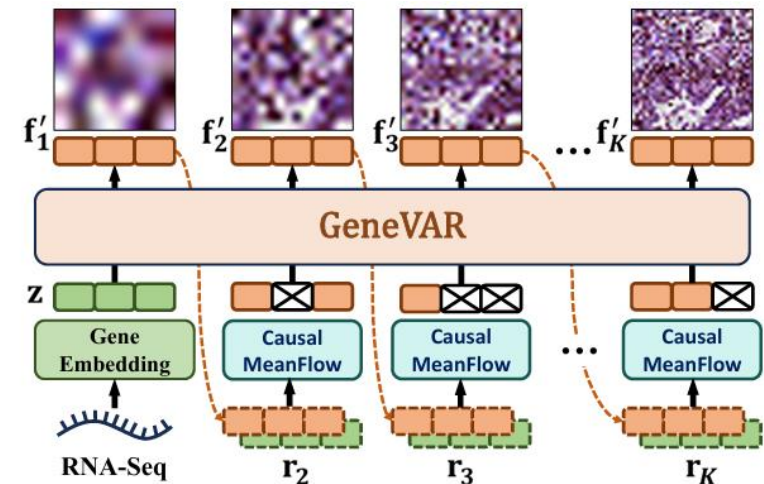
Understanding how transcriptomic programs shape tissue morphology remains a central challenge in computational pathology. Gene-to-WSI tile synthesis offers a principled generative framework to translate molecular profiles into histological images.

## Strategy:

Most existing methods compress RNA-Seq into a single global embedding injected once at initialization, an oversimplified design that weakens transcriptomic signals. We present GeneVAR, an Autoregressive Gene-to-WSI model that reformulates synthesis as an iterative, coarse-to-fine generative process.



Traditional solutions



Ours GeneVAR

# 1. Introduction

## Motivation

Despite recent progress, Gene-to-WSI tile synthesis remains underdeveloped. Existing methods typically compress RNA-Seq profiles into low-dimensional embeddings to condition GANs or diffusion models. However, this design has three key limitations: (1) one-time conditioning leads to signal decay, causing generated images to drift from gene-driven morphology; (2) fixed-resolution synthesis imposes scale rigidity, reducing cross-scale consistency.



Embeddings are learned in a purely correlational manner, leaving models vulnerable to confounders. An iterative paradigm alone does not immunize the generative trajectory against confounders, such as tumor purity, and staining variability, whose correlational footprints can divert decoding from gene-driven semantics.



## Solution

**Coarse-to-Fine:** To address these challenges, we reformulate Gene-to-WSI synthesis as an *iterative, coarse-to-fine* generative process in which transcriptomic signals are injected at multiple stages during generation rather than used once at initialization. Building on this principle, we propose Gene VAR-the first autoregressive Gene-to-WSI model that delivers step-wise molecular reinforcement while preserving structural coherence and causal fidelity.

**Causal MeanFlow:** r-CM is introduced as a gene-driven causal module intrinsically embedded in the auto regressive trajectory. Within this framework, a Structural Causal Model (SCM) governs generative updates from coarse to fine semantics, ensuring that causal factors, *i.e.*, semantic information regulated by gene expression, remain consistent, while non-causal factors are suppressed.

# 1. Introduction

Overall, the contributions of this work can be summarized as:

- ◆ **Paradigm Shift.** GeneVAR reformulates Gene-to-WSI synthesis as an iterative, coarse-to-fine generative process in which transcriptomic signals are injected at multiple stages across scales, overcoming the signal decay and rigidity of static global embeddings.
- ◆ **Causal Fidelity.** By integrating the novel Causal MeanFlow module into the autoregressive trajectory, GeneVAR disentangles transcriptomic effects from confounders, ensuring that morphological synthesis remains biologically grounded and causally aligned.
- ◆ **State-of-the-Art Performance.** Comprehensive evaluation on five TCGA cohorts demonstrates that GeneVAR attains the lowest FID and the highest accuracy in downstream classification, setting a new benchmark for both generative fidelity and functional utility.

# 2. Method

- **Multi-Scale Vector Quantization:**

We employ multi-scale vector quantization (MSVQ) to discretize a WSI tile  $X \in R^{H \times W \times C}$  into  $K$  hierarchical token maps  $S = \{s_k\}_{k=1}^n$ , where each  $s_k \in R^{h_k \times w_k}$  represents a discrete map at scale  $k$ .

- **Problem Reformulation:**

Gene-to-WSI synthesis should preserve transcriptomic guidance dynamically across scales. However, existing methods collapse RNA-Seq into static embeddings injected once, leading to signal decay, scale rigidity, and spurious correlations. To overcome these issues, we reformulate Gene-to-WSI synthesis as a coarse-to-fine autoregressive process, where an image is quantized into hierarchical token maps  $S = \{s_k\}_{k=1}^n$  and generated sequentially under recurrent guidance from  $\mathbf{g}$ :

$$s_k = P_{\Theta}(s_{<k}, \mathbf{g})$$

Here,  $\mathbf{g}$  is injected at multiple stages rather than only once at initialization, ensuring that transcriptomic signals re-main persistently active throughout the generative trajectory while leveraging the distributional strength of autoregressive models.

# 2. Method

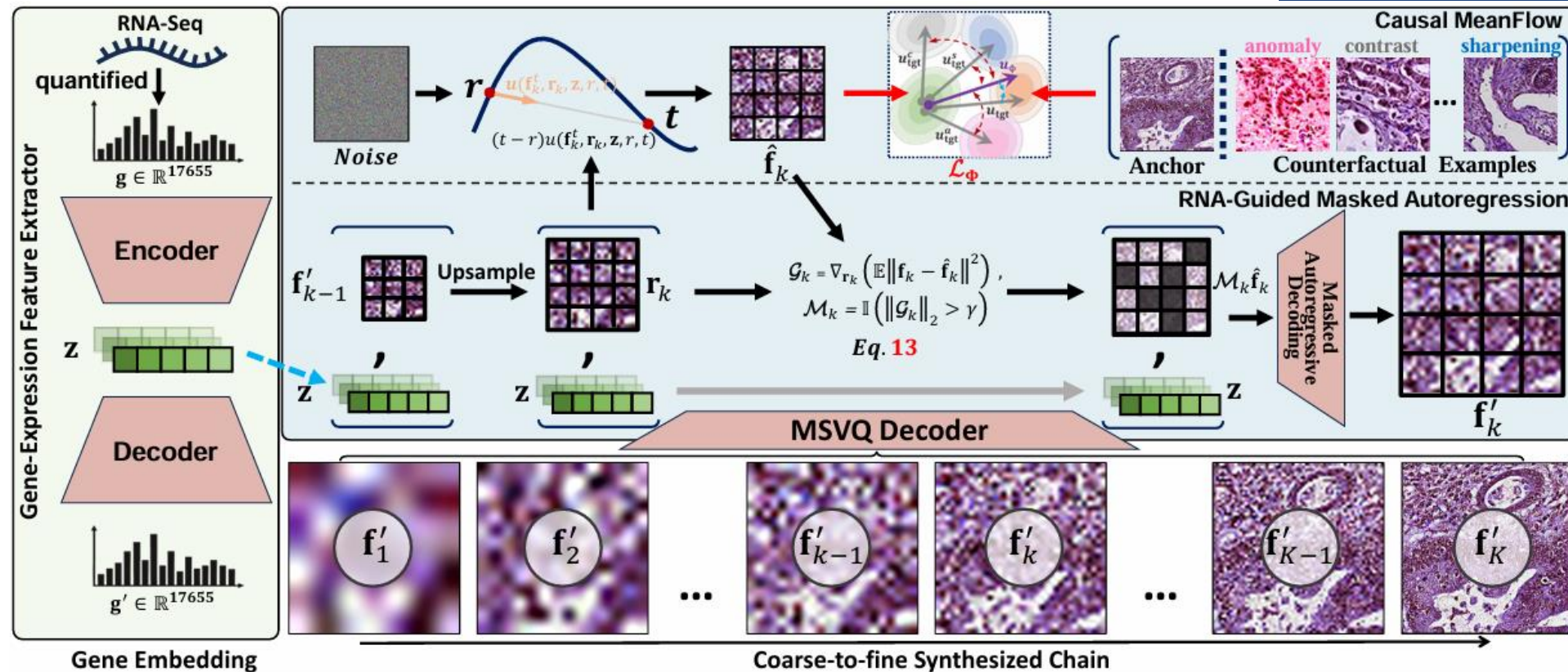


Figure 2. **Overview of GeneVAR.** RNA-Seq embeddings are injected into a multi-scale autoregressive pipeline, where morphology is reconstructed from  $f'_K$  via the MSVQ decoder and reinforced by Causal MeanFlow with counterfactual supervision (see Sec. 3 for details).

Fig. 2 presents an overview of our GeneVAR, an Autoregressive Gene-to-WSI model that reformulates synthesis as an iterative, coarse-to-fine generative process. At its core is a novel Causal MeanFlow module that reinforces transcriptome-informed guidance at multiple stages and mitigates non-causal factors through counterfactual-style interventions, thereby ensuring biological fidelity throughout the generative trajectory.

# 2. Method

## Causal MeanFlow:

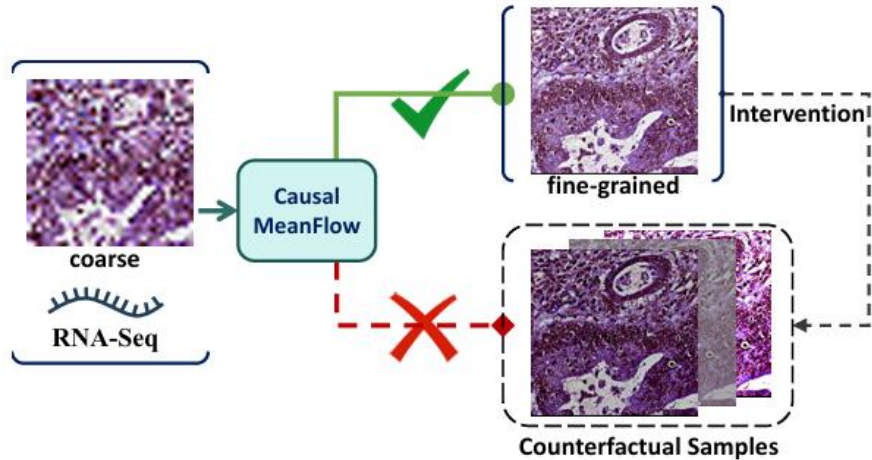
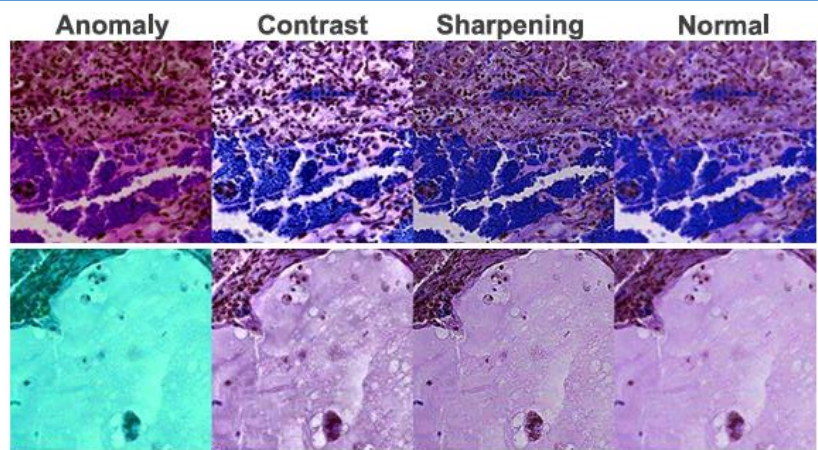


Figure 3. **Structural Causal Mechanism.** Our model transforms coarse gene-guided features into fine-grained morphology, while counterfactual interventions suppress non-causal variations.



## Counterfactual examples construction and velocity formulation:

**b.) Counterfactual Regularization.** Autoregressive modeling can be viewed as an iterative reconstruction of quantized features, producing latent representations  $\{\mathbf{f}'_k\}_{k=1}^K$ . Within our framework,  $\mathbf{f}'_k$  is reconstructed from  $\mathbf{r}_k$  under the guidance of gene embedding  $\mathbf{z}$ , progressively enriching semantics across scales. To prevent spurious enhancements from non-causal factors such as tumor purity and staining [1, 23, 28, 38], we introduce a causal regularization strategy based on counterfactual interventions. Following the Structural Causal Model (SCM) [34] and the principle that interventions must destruct non-causal cues [45], we apply three perturbations, *i.e.*, color anomaly and contrast adjustment (to model stain variation) and sharpening (to simulate tumor purity differences), to generate counterfactual variants  $X^a, X^c, X^s$  for each WSI tile  $X$ , as illustrated in Fig. 3. After quantization, these yield feature maps  $\{\mathbf{f}_k^a, \mathbf{f}_k^c, \mathbf{f}_k^s\}_{k=1}^K$ , which enforce  $\Phi$  to emphasize causal, scale-invariant morphology.

To endow  $r\text{-CM}$  with causal regularization, we enlarge the discrepancy between the predicted average velocity  $u_\phi$  and the target velocities  $\{u_{\text{tgt}}^a, u_{\text{tgt}}^c, u_{\text{tgt}}^s\}$  obtained from degraded counterfactuals, thereby stabilizing reconstruction. For example,  $u_{\text{tgt}}^a$  is not derived via costly partial derivatives in Eq. 6. Instead, we decompose average velocity into magnitude and direction, where the direction is computed by normalizing the flow between two sampled fields  $\mathbf{f}_k^{a,t}$

and  $\mathbf{f}_k^{a,r}$ , and the magnitude is estimated by scaling  $\|u_{\text{tgt}}\|$  with a stochastic factor  $\lambda$ . Targets  $\{u_{\text{tgt}}^c, u_{\text{tgt}}^s\}$  are constructed analogously. The procedure is formally expressed as:

$$u_{\text{tgt}}^a = \lambda_a \|u_{\text{tgt}}\| \cdot \frac{\mathbf{f}_k^{a,t} - \mathbf{f}_k^{a,r}}{\|\mathbf{f}_k^{a,t} - \mathbf{f}_k^{a,r}\|}. \quad (9)$$

**c.) Learning Objective.** We further define a causality-driven learning objective to disentangle scale-invariant morphological semantics (causal factors) from degradation-related cues (non-causal factors). For the  $k$ -th iteration,  $u_{\text{tgt}}$  from  $\mathbf{f}_k$  serves as the anchor, while  $\mathbf{f}_{k-1}$  is treated as an extreme counterfactual, yielding  $u_{\text{tgt}}^l$  via Eq. 9. Additional degradation samples  $\{u_{\text{tgt}}^a, u_{\text{tgt}}^c, u_{\text{tgt}}^s\}$  derived from independent WSIs, distinct from the source of  $u_{\text{tgt}}$ , forming the counterfactual set  $\mathcal{C}_u$ . This design blocks potential shortcuts that exploit non-causal signals and enforces reliance on scale-invariant morphology. The training objective is:

$$\mathcal{L}_\Phi = \mathbb{E}[\|u_\phi - \text{sg}(u_{\text{tgt}})\|^2] - \frac{\alpha}{N} \sum_{u_{\text{tgt}}^n \sim \mathcal{C}_u} \mathbb{E}[\|u_\phi - \text{sg}(u_{\text{tgt}}^n)\|^2], \quad (10)$$

where  $\text{sg}(\cdot)$  denotes stop-gradient and  $\alpha$  controls the strength of causal regularization.

# 2. Method

## RNA-Guided Masked Autoregression:

We formulate GeneVAR as an autoregressive model over tokenized sequences. This factorization is realized by a ViT-based transformer with a causal mask, ensuring that generation proceeds strictly left-to-right under continuous transcriptomic conditioning, akin to GPT-2, as shown in **Eq.11**.

$$p(s_1, s_2, \dots, s_K) = \prod_{k=1}^K p(s_k | s_{<k}, s_0 = \mathbf{z}), \quad (11)$$

**a.) Gradient-Guided Masking.** After  $r\text{-}\mathcal{CM}$  converges, it produces stable reconstructions  $\hat{\mathbf{f}}_k$  from  $\mathbf{r}_k$  before passing them to the transformer. We define the gradient as the derivative of the MSE loss between  $\mathbf{f}_k$  and  $\hat{\mathbf{f}}_k$  with respect to  $\mathbf{r}_k$ , and compute its  $\ell_2$  norm along the channel dimension:

$$\begin{aligned} \mathcal{G}_k &= \nabla_{\mathbf{r}_k} (\mathbb{E} \|\mathbf{f}_k - \hat{\mathbf{f}}_k\|^2), \\ \mathcal{M}_k &= \mathbb{I}(\|\mathcal{G}_k\|_2 > \gamma), \gamma \sim \mathcal{N}(0, 1)^{h_k \times w_k}. \end{aligned} \quad (12)$$

Here  $\mathbb{I}(\cdot)$  denotes the indicator function. Tokens with large gradient responses are deemed RNA-conditioned salient regions, as perturbations there strongly affect reconstruction. These positions are therefore preferentially retained, dynamically injecting transcriptomic cues at the  $k$ -th iteration.

**b.) Masked Autoregressive Decoding.** Following class/text-conditioned VAR models, the latent gene embedding  $\mathbf{z}$  is used as the initial condition token  $\mathbf{r}_1$ . For subsequent steps, tokens at positions with  $\mathcal{M}_k = 0$  are replaced by a learnable embedding  $e$ . The autoregressive generation and cross-entropy training objective are formulated as:

$$\begin{aligned} \{s'_k\}_{k=1}^K &= \mathcal{P}_{\Theta}(\{\mathbf{r}_1, \dots, \mathcal{M}_k \hat{\mathbf{f}}_k, \mathcal{M}_K \hat{\mathbf{f}}_K\}), \\ \mathcal{L}_{\Theta} &= \sum_{k=1}^K \mathbb{CE}(s'_k, s_k). \end{aligned} \quad (13)$$

For efficiency, masking is applied only when  $k \geq K_m$ , balancing accuracy and cost. Together, these advances endow GeneVAR with causality-aware, biologically grounded generation, which we validate in the next section through comprehensive experiments.

# 3. Quantitative Results

Table 1. **Quantitative comparison on TCGA benchmarks.** Fréchet Inception Distance (**FID, lower is better**) is reported along with model parameters (#Para) and generation steps (#Step). GeneVAR consistently achieves the best performance across all cohorts.

Method	Type	#Para	#Step	GBM	CESC	KIRP	COAD	LUAD	ALL
RNA-CDM [4]	Diffusion	1146M	2000	24.15	25.76	24.46	33.60	27.98	23.36
PathLDM [43]	Diffusion	400M	50	22.46	19.87	20.39	26.96	21.28	20.29
U-ViT [2]	Diffusion	297M	100	13.89	<u>15.74</u>	21.83	26.75	17.86	18.55
DiT [29]	Diffusion	305M	250	15.03	18.75	21.53	29.13	18.57	18.11
SiT [27]	Flow Matching	305M	25	14.63	19.48	22.10	29.47	19.52	18.84
LlamaGen [37]	Token-wise Autoregression	343M	256	15.48	18.34	19.89	27.91	17.52	17.43
VQGAN [11]	Token-wise Autoregression	227M	256	21.74	23.48	26.10	32.47	26.48	25.09
VAR [39]	Scale-wise Autoregressive	310M	10	<u>12.96</u>	17.21	<u>17.32</u>	<u>25.84</u>	<u>15.40</u>	<u>16.83</u>
ImageFolder [26]	Scale-wise Autoregressive	314M	10	23.75	23.97	25.81	33.25	25.46	24.81
<b>GeneVAR (Ours)</b>	Scale-wise Autoregressive	317M	10	<b>11.20</b>	<b>13.68</b>	<b>14.49</b>	<b>19.86</b>	<b>13.65</b>	<b>12.95</b>

1. GeneVAR achieves leading generative performance (FID) across all datasets, setting new benchmarks.

Table 2. **Cell distribution comparison.** Mean±std of cell proportions across all cohorts. Extended results are shown in Supplementary Material.

Dataset	Method	Neoplastic	Dead
GBM	VAR	9.96%±21.99	69.55%±34.59
	<b>Ours</b>	<b>5.04%±12.07</b>	<b>78.67%±26.91</b>
	Real	6.04%±16.29	77.49%±30.90
LUAD	VAR	25.06%±27.38	55.99%±30.00
	<b>Ours</b>	<b>20.20%±22.78</b>	<b>61.10%±26.56</b>
	Real	19.32%±25.21	63.12%±29.60

2. While FID captures visual idelity, biological plausibility is critical for clinical relevance. We confirm that GeneVAR not only preserves overall cell composition but also generates morphologies consistent with cancer type, demonstrating biological realism beyond mere visual fidelity.

Table 3. **Tile classification performance.** ACC/F1 under varying substitution ( $p$ ) and pretraining ( $q$ ) ratios. Full results are provided in Supplementary Material.

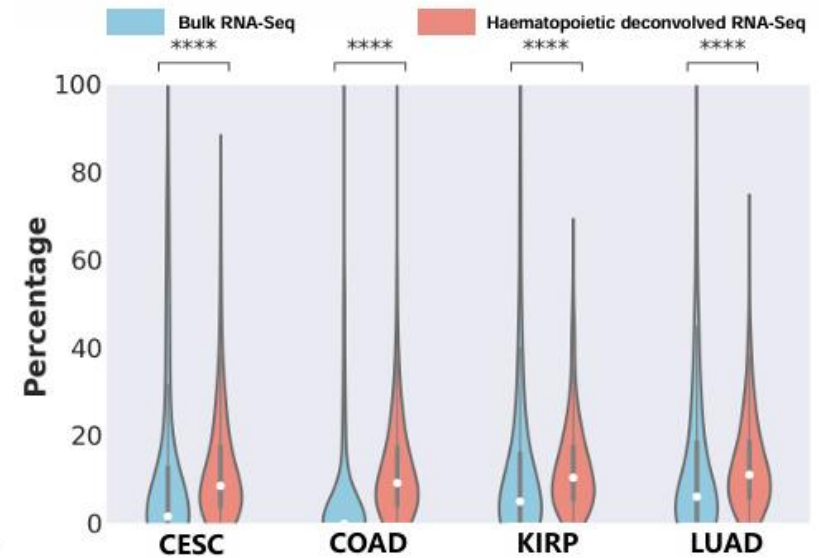
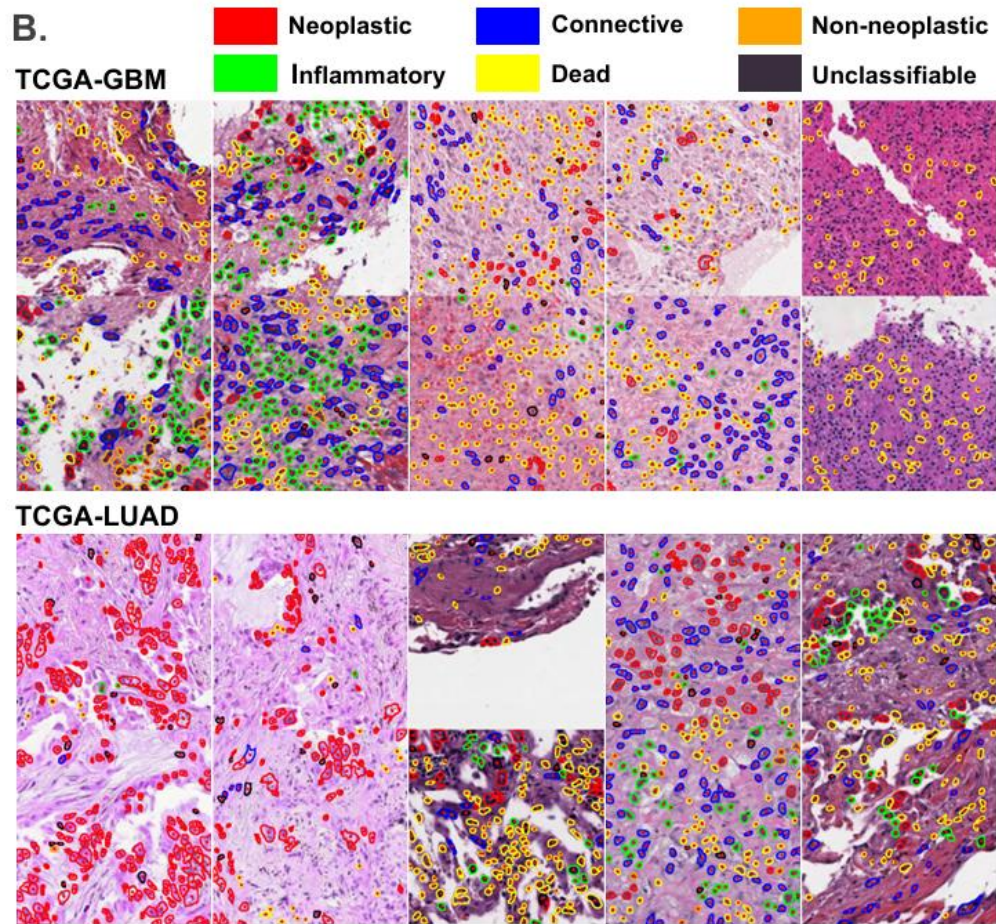
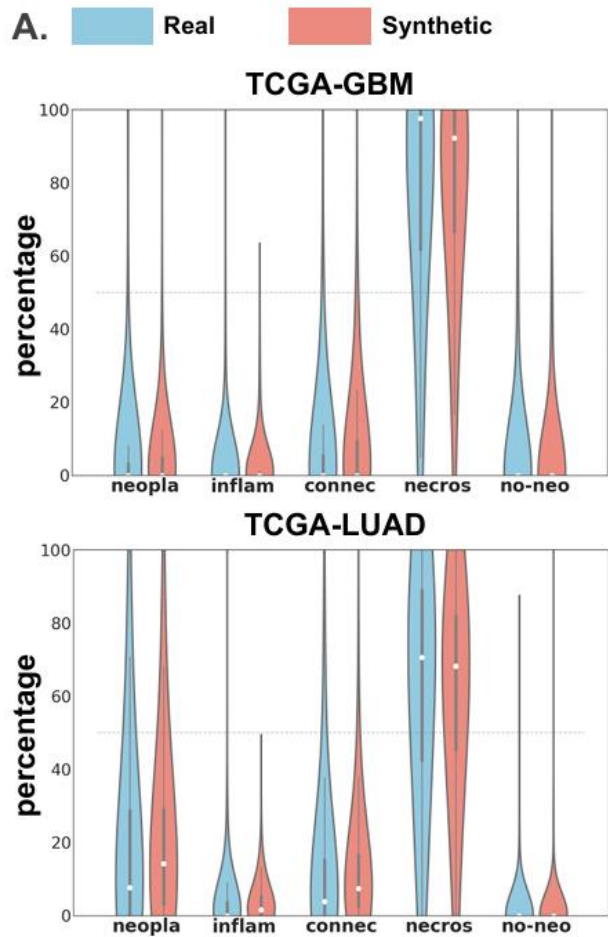
Method	$p=0.0$		$p=0.75$		
	ACC	F1	ACC	F1	
R-CDM [4]	0.573±0.020	0.556±0.022	0.492±0.016	0.472±0.046	
VAR [39]	0.573±0.034	0.556±0.015	0.521±0.030	0.510±0.035	
<b>Ours</b>	<b>0.579±0.032</b>	<b>0.570±0.039</b>	<b>0.592±0.029</b>	<b>0.588±0.031</b>	
		$q=0.5$		$q=1.0$	
R-CDM [4]	0.618±0.026	0.612±0.030	0.650±0.021	0.641±0.031	
VAR [39]	0.661±0.020	0.662±0.019	0.708±0.011	0.709±0.010	
<b>Ours</b>	<b>0.722±0.023</b>	<b>0.722±0.020</b>	<b>0.767±0.015</b>	<b>0.766±0.017</b>	

Table 4. **Performance Comparison** of MIL methods (w/ vs. w/o pretraining). Our GeneVAR generates 512 patch-wise tiles using each RNA-Seq.

Method	ACC		F1-score		AUC	
	w/o	w/	w/o	w/	w/o	w/
TransMIL [36]	0.849	<b>0.877</b>	0.847	<b>0.875</b>	0.894	<b>0.934</b>
ACMIL [46]	0.767	<b>0.863</b>	0.749	<b>0.858</b>	0.817	<b>0.938</b>
WiKG [25]	0.767	<b>0.822</b>	0.753	<b>0.818</b>	0.820	<b>0.917</b>
MambaMIL [42]	0.836	<b>0.918</b>	0.833	<b>0.917</b>	0.925	<b>0.941</b>

3. For downstream evaluation (tile-level and slide-level classification), we train tile-level and WSI-level classifiers on synthetic tiles and report F1-score, accuracy (ACC), and AUC.

# 3. Qualitative Interpretability



**Synthetic Vs. Real tiles.** Panel A. Cell-type distributions in real and synthetic tiles across TCGA cohorts. Panel B. HoverNet visualizations showing consistent cellular composition between synthetic and real samples.

**Lymphocyte Distribution Comparison** between WSIs generated from bulk RNA-Seq and those generated from hematopoietic-deconvolved expression. **Lymphocyte proportions show consistently increased** trends across TCGA-CESC, TCGA-COAD, TCGA-KIRP and TCGA-LUAD.

# 4. Conclusion

In this work, we introduced GeneVAR, an autoregressive Gene-to-WSI tile synthesis model that unites multi-stage transcriptomic conditioning with causality-aware modeling. By reformulating synthesis as an iterative, coarse-to-fine autoregressive process, GeneVAR continually reinforces transcriptomic signals, mitigating signal decay and ensuring cross-scale consistency. At its core lies the Causal Mean-Flow module, guided by RNA-Seq embeddings, which uses counterfactual interventions to suppress spurious variation and enforce causal fidelity in morphology. Beyond benchmarking, GeneVAR serves as a controllable tool for probing gene-morphology relationships under transcriptomic perturbations, opening new avenues for integrative studies in computational pathology.



## GeneVAR: Causal MeanFlow for Autoregressive Gene-to-WSI Tile Synthesis

Jianwei Zhao<sup>1</sup>, Fan Yang<sup>2</sup>, Xin Li<sup>1\*</sup>, Qiang Zhai<sup>3</sup>, Ao Luo<sup>4</sup>, Ziqi Ren<sup>5</sup>, Zhicheng Jiao<sup>6</sup>, and Hong Cheng<sup>1</sup>

UESTC<sup>1</sup>, ALPACA I LAB<sup>2</sup>, SICAU<sup>3</sup>, SWJTU<sup>4</sup>, XDU<sup>5</sup>, Brown<sup>6</sup>



### Motivation

- Gene-to-WSI tile synthesis offers a principled generative framework to translate molecular profiles into histological images.
- First, collapsing the transcriptome into a static embedding injected only once causes signal decay, with molecular guidance fading as generation unfolds and images drifting toward superficial correlations rather than gene-driven morphology.
- Models are vulnerable to confounders such as staining variability and imaging artifacts, thereby entangling true gene-driven morphology with non-biological factors.

Figure 1. Key Idea. GeneVAR integrates Causal MeanFlow into an autoregressive framework, iteratively injecting transcriptomic guidance while enforcing invariance to non-biological factors, enabling biologically faithful WSI synthesis.

### Contribution

- Paradigm Shift.** GeneVAR reformulates Gene-to-WSI synthesis as an iterative, coarse-to-fine generative process in which transcriptomic signals are injected at multiple stages across scales, overcoming the signal decay and rigidity of static global embeddings.
- Biological Fidelity.** By integrating the novel Causal MeanFlow module into the autoregressive trajectory, GeneVAR disentangles true transcriptomic signals from nonbiological confounders. This ensures that morphological synthesis remains biologically grounded and artifact-invariant, explicitly avoiding claims of gene-level causal discovery.
- State-of-the-Art Performance.** Comprehensive evaluation on five TCGA cohorts demonstrates that GeneVAR attains the lowest FID and the highest accuracy in downstream classification, setting a new benchmark for both generative fidelity and functional utility.

### Method

RNA-Seq embeddings are injected into a multi-scale autoregressive pipeline, where morphology is reconstructed from  $f_k^k$  via the MSQ decoder and reinforced by Causal Mean Flow with counterfactual supervision.

### Causal Analysis

Figure 1. Visualizations of the Causal Graph for generative updates from coarse to fine semantics.

Figure 3. Structural Causal Mechanism. Our model transforms coarse gene-guided features into fine-grained morphology, while counterfactual interventions suppress non-causal variations.

### Experimental Results

Method	Type	Wtata	thsp	GBI	ERCC	KIRP	COAD	LIAD	ALL
RNA-CBM [1]	Diffusion	14600	2700	24.85	27.76	24.66	11.66	27.98	21.26
SCALE [10]	Diffusion	4904	78	22.65	19.87	20.59	70.96	21.28	28.28
U-VIT [1]	Diffusion	29753	436	15.80	15.72	21.83	16.75	17.86	18.55
SEIT [8]	Diffusion	30581	286	18.80	18.75	21.53	19.12	18.57	18.11
SEIT-DL	Flow Matching	30581	25	14.83	19.48	22.33	19.47	19.52	18.84
LowGen [11]	Time-wise Autoregression	8134	276	15.44	14.34	19.89	17.91	17.52	17.43
VQGAN [12]	Time-wise Autoregression	22781	276	21.74	21.48	26.83	22.47	26.19	25.99
TOR [13]	Scale-wise Autoregression	33981	80	15.86	17.21	17.12	17.64	15.69	16.85
img2img [17]	Scale-wise Autoregression	31494	80	15.25	17.87	18.81	17.25	17.46	18.18
GeneVAR (Ours)	Scale-wise Autoregression	10751	50	14.23	12.88	14.49	18.06	14.88	12.98

#### 1. Quantitative comparison with SOTA methods

#### 2. Cell-type distributions in real and synthetic tiles