

GMM Diffusion: Gaussian Mixture Masks for Diffusion Models

Chang Liu, Karim Habashy, Peter Yuchen Pan, Hadi Sepanj

University of Waterloo

Motivation

- Diffusion models require significant computation and training time
- Our method proposes 2 novel modifications to allow small diffusion models generate similar quality images to their larger counterparts

Previous Work

- FreeU: Reweight U-Net's skip connections and backbone feature maps for diffusion models to improve image generation quality [1]
- GMM: Gaussian Mixture Masks applied to Retentive Networks [3]
- U-ViT: ViT-based architecture for diffusion models [2]

Proposed Methods

To the best of our knowledge, our work is the first to apply both Gaussian Mixture Mask and U-Net scale factors to diffusion models.

Specifically, our work is novel in 2 following ways.

1. Gaussian mixture mask

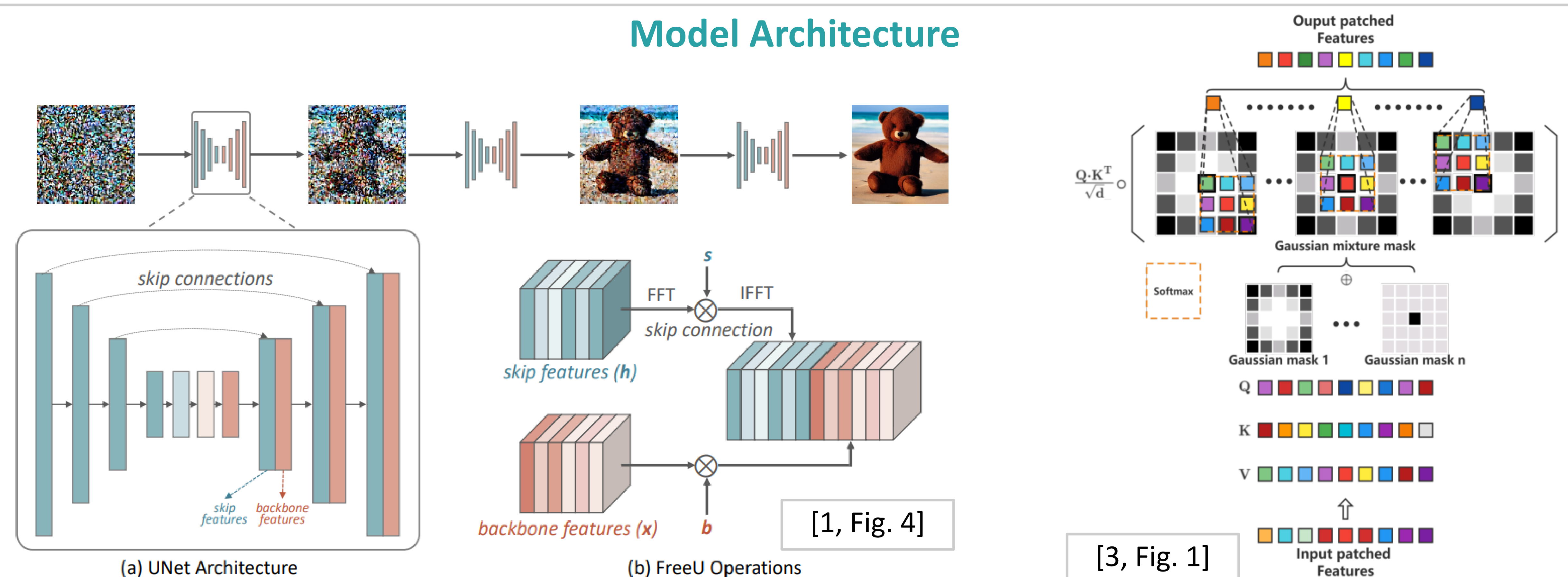
- Proposes Gaussian mixture mask to boost image generation capabilities for small dataset with almost zero additional parameters and computational cost
- Learn 2 parameters to implicitly generate Gaussian mixture mask on the attention heads

$$\text{GMMAttention}(Q, K, V) = \text{Softmax}\left(\frac{QK^T}{\sqrt{d_k}} \circ M\right) V$$

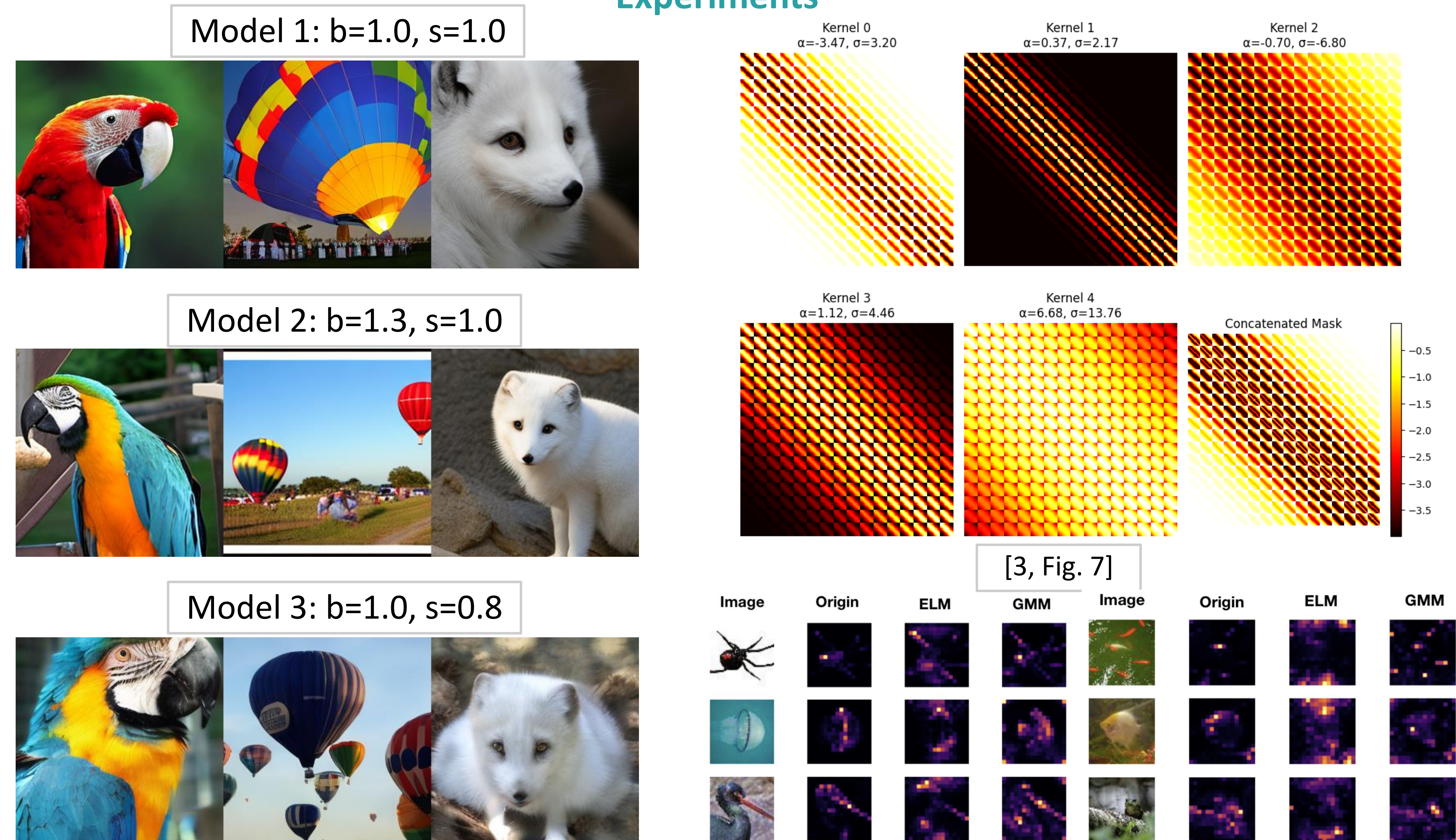
2. Skip connections during inference

- Re-scales U-Net's skip connection feature maps and backbone feature maps to improve image quality without additional training or finetuning
- Adjusts 2 scaling factors (scaling backbone based on averaged feature maps and scaling long skip connections across different decoder blocks).

Model Architecture



Experiments



References

- [1] C. Si, Z. Huang, Y. Jiang, and Z. Liu, "FreeU: Free Lunch in Diffusion U-Net," <https://arxiv.org/pdf/2309.11497.pdf>, Sep. 2023.
- [2] F. Bao et al., "All are worth words: A ViT backbone for diffusion models," 2023 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), 2023. doi:10.1109/cvpr52729.2023.02171
- [3] C. Li and C. Zhang, "Toward a deeper understanding: Retnet viewed through convolution, 2023. doi:10.2139/ssrn.4637493

UNIVERSITY OF
WATERLOO



GMM Diffusion: Gaussian Mixture Masks for Diffusion Models

Chang Liu, Karim Habashy, Peter Yuchen Pan, Hadi Sepanj

University of Waterloo

Motivation

- Diffusion models require significant computation and training time
- Our method proposes 2 novel modifications to allow small diffusion models generate similar quality images to their larger counterparts

Previous Work

- FreeU: Reweight U-Net's skip connections and backbone feature maps for diffusion models to improve image generation quality [1]
- GMM: Gaussian Mixture Masks applied to Retentive Networks [3]
- U-ViT: ViT-based architecture for diffusion models [2]

Proposed Methods

To the best of our knowledge, our work is the first to apply both Gaussian Mixture Mask and U-Net scale factors to diffusion models.

Specifically, our work is novel in 2 following ways.

1. *Gaussian mixture mask*

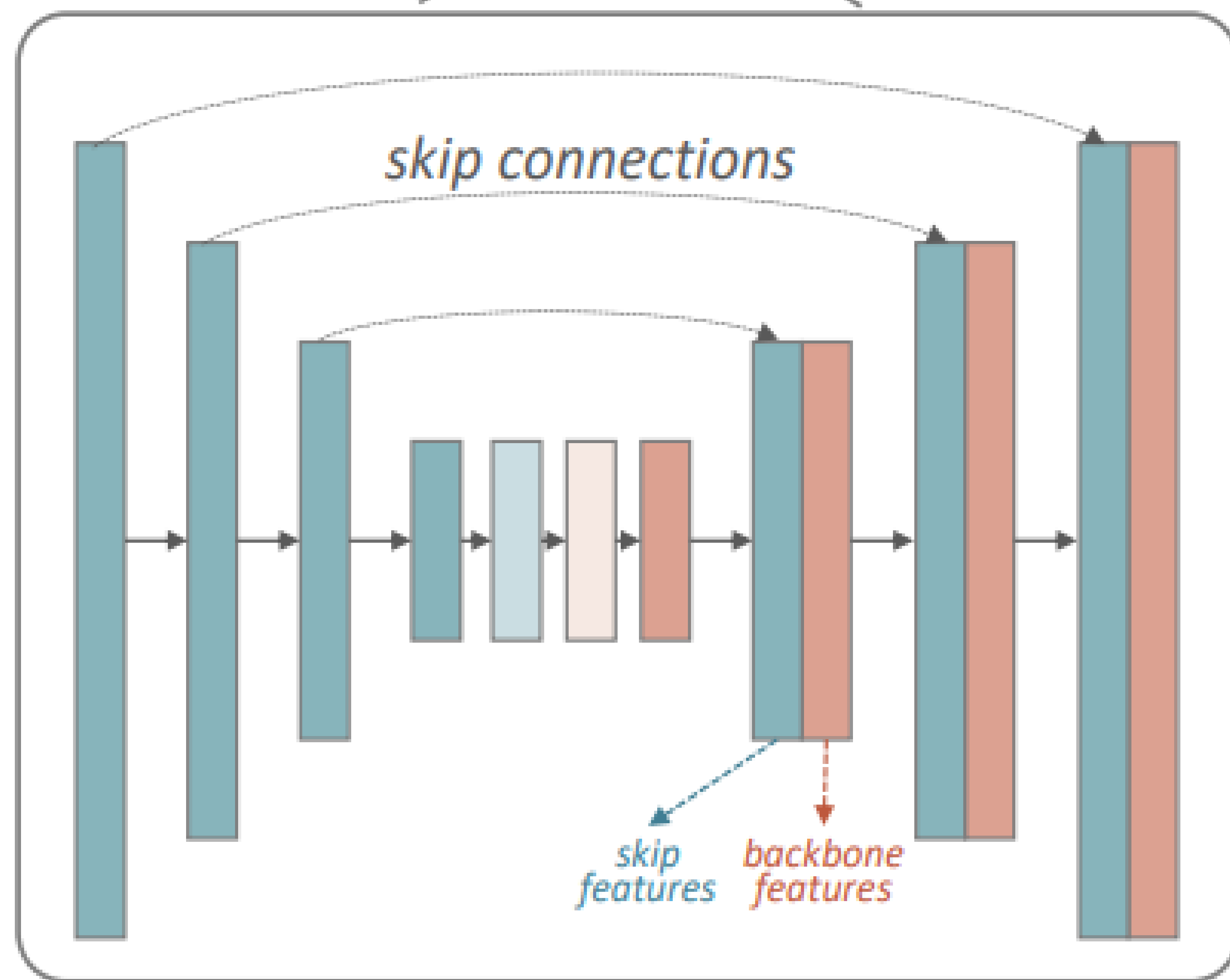
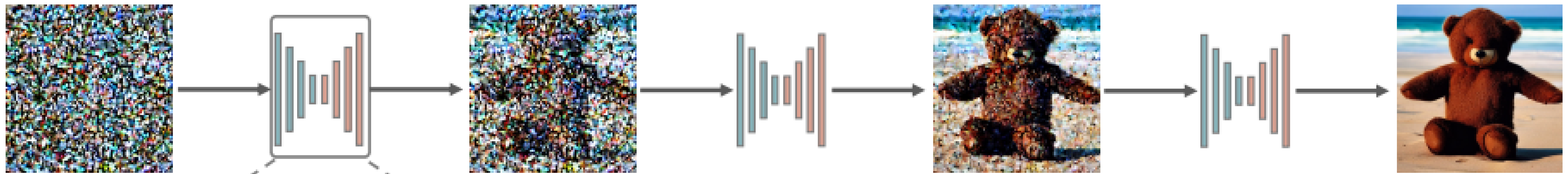
- Proposes Gaussian mixture mask to boost image generation capabilities for small dataset with almost zero additional parameters and computational cost
- Learn 2 parameters per attention block to implicitly generate Gaussian mixture mask on the attention heads

$$\text{GMMA}(\text{Attention}(Q, K, V)) = \text{Softmax}\left(\frac{QK^T}{\sqrt{d_k}} \circ M\right)V$$

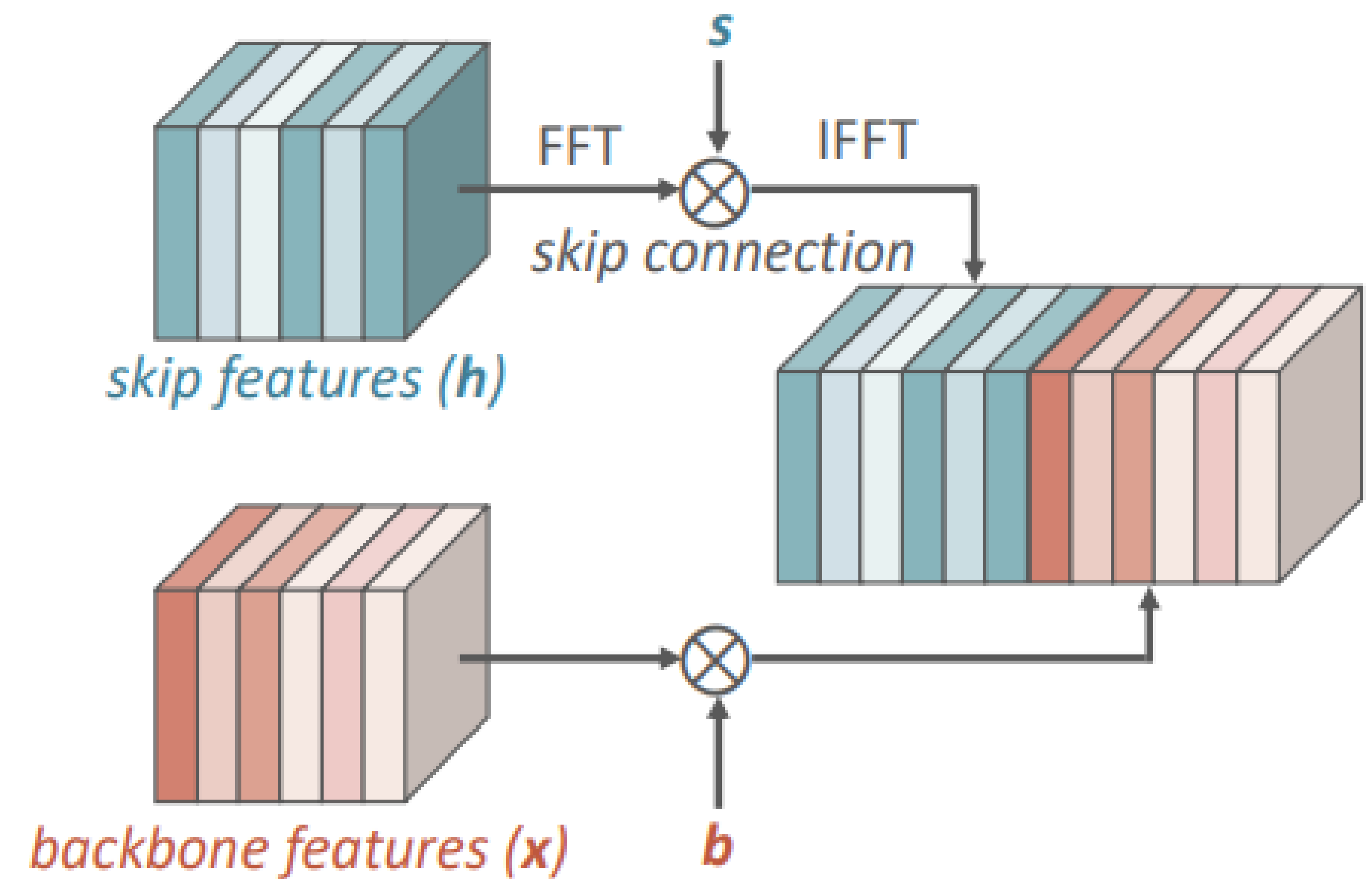
2. *Skip connections during inference*

- Re-scales U-Net's skip connection feature maps and backbone feature maps to adaptively improve image quality without additional training or finetuning
- Adjusts 2 scaling factors (scaling backbone based on averaged feature maps and scaling long skip connections across different decoder blocks).

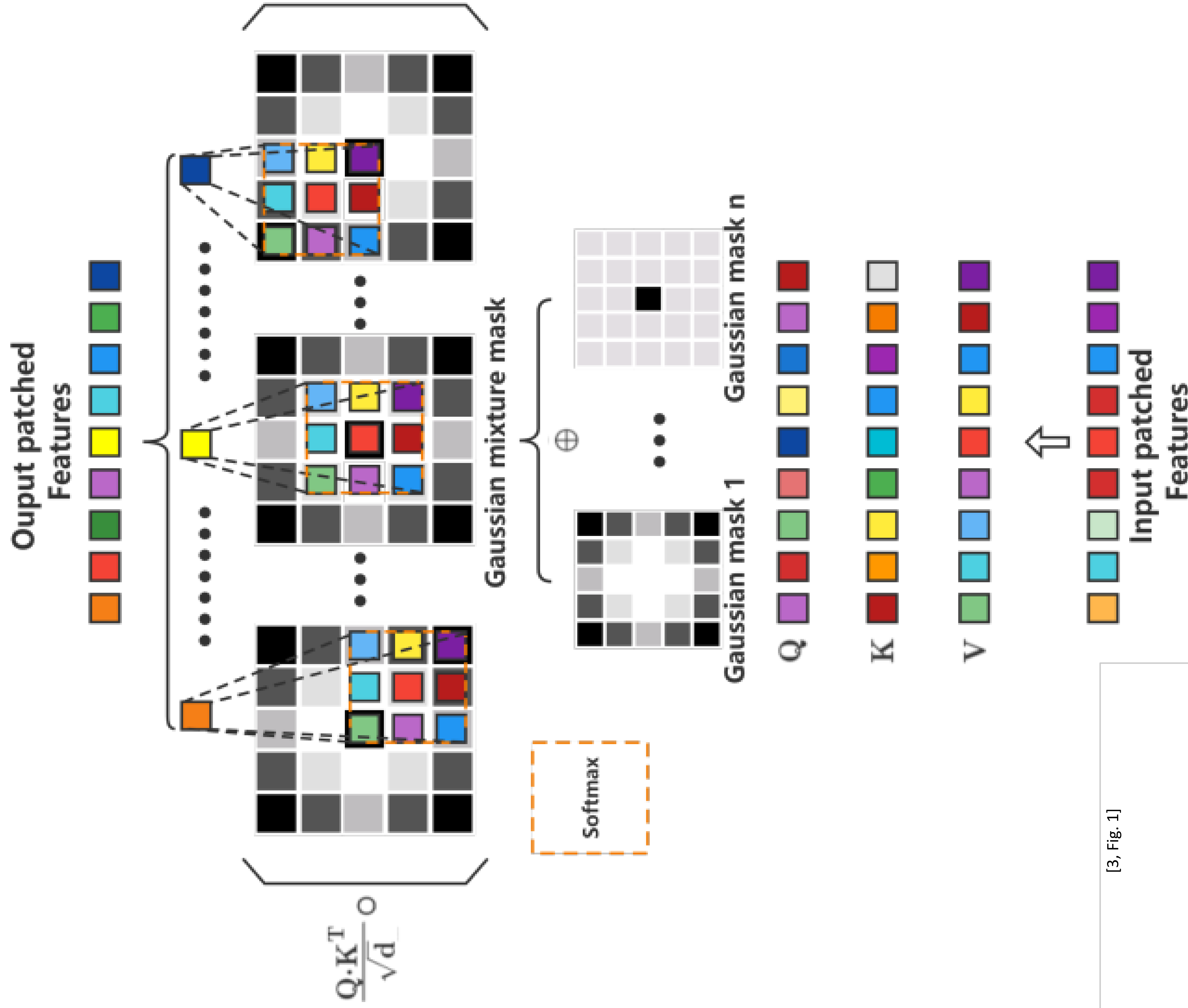
Model Architecture



(a) UNet Architecture



(b) FreeU Operations



[3, Fig. 1]

Model 1: $b=1.0$, $s=1.0$

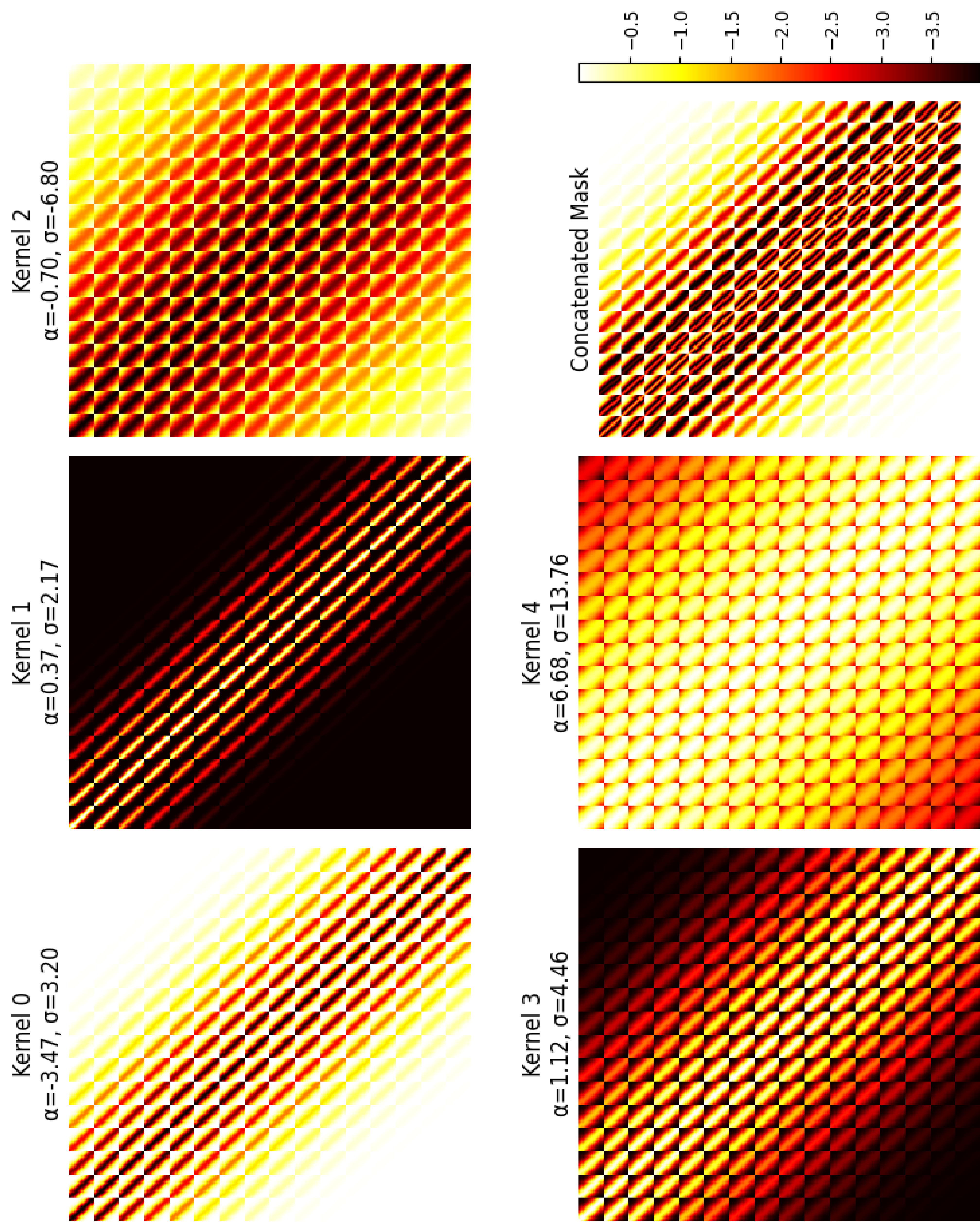


Model 2: $b=1.3$, $s=1.0$

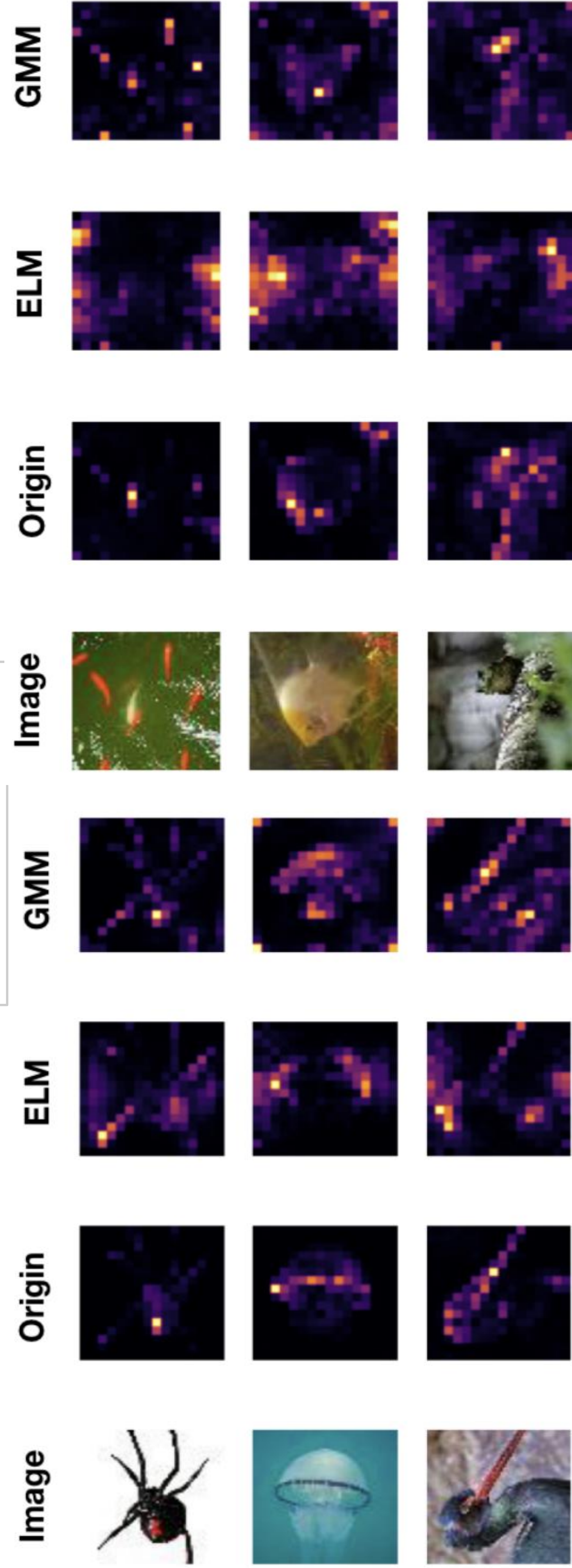


Model 3: $b=1.0$, $s=0.8$





[3, Fig. 7]



References

- [1] C. Si, Z. Huang, Y. Jiang, and Z. Liu, "FreeU: Free Lunch in Diffusion U-Net," <https://arxiv.org/pdf/2309.11497.pdf>, Sep. 2023.
- [2] F. Bao et al., "All are worth words: A ViT backbone for diffusion models," 2023 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), 2023. doi:10.1109/cvpr52729.2023.02171
- [3] C. Li and C. Zhang, *Toward a deeper understanding: Retnet viewed through convolution*, 2023. doi:10.2139/ssrn.4637493

