





Overview

Diffusion models are the State-of-the-Art for text-to-image generation, and increasing research effort has been dedicated to adapting the inference process of pretrained diffusion models to achieve zero-shot capabilities. An example is the generation of long images, which has been tackled in recent works by combining strided diffusions over overlapping latent features. This yield perceptually aligned but semantically incoherent panoramas. We propose the Merge-Attend-Diffuse operator (MAD), pluggable into different types of diffusion models featuring attention operations, and an inference-time strategy for generating perceptually and semantically coherent panoramas.

MAD is modular: use it where and when you want

MAD can be applied in different blocks of the backbone and for a variable portion of the diffusion denoising chain, obtaining the desired tradeoff between uniformity and variability. MAD applied in *different backbone blocks*: MAD applied for a *different fraction of the denoising chain*:





MAD is effective and efficient

We perform quantitative comparison on standard panorama generation literature prompts and compare with direct inference across multiple resolutions. For MAD, we sweep over several application tresholds, quantifying the tradeoff between uniformity and variability.

						mCLIP ↑	I-LPIPS	\downarrow FID \downarrow	$KID\downarrow$					
	mCLIP ↑ I-LPIPS ↓		, FID ↓	KID↓		2 Inference Steps								
SD [1]	31.63	0.74	28.31	<0.01	LCD [4]	30.57	0.53	26.82	<0.01		mCLIP	↑ I-LPIPS ↓	↓ FID↓	KID↓
SD-L [1] SD-L+Attn-S [5]	32.01 32.02	0.50 0.52	87.64 80.16	76.83 67.54	$\frac{100 - 1}{100}$	30.75	0.47	27.96	13.39	SD [1]	32.45	0.67	49.32	< 0.01
MD [2] SyncD [3]	31.77	0.69	33.52	9.04	MAD (τ =1) MAD (τ =2)	30.97 30.87	0.50 0.47	35.70 52.88	23.74 48.04	SD-L [1] SD-L+Attn-S [5]	31.89 32.03	0.52 0.53	68.03 65.08	6.61 5.47
MAD (τ=0)	31.65	0.64	34.51	9.19	4 Inference Steps					MD [2] SvncD [3]	32.46 32.34	0.63 0.55	49.41 53.52	0.42 1.11
MAD ($ au$ =5) MAD ($ au$ =15) MAD ($ au$ =25)	31.86 32.03 32.15	0.59 0.56 0.53	38.10 48.52 61.76	13.75 27.15 43 31	LCD [4] LCD-L [4]	31.37 31.30	0.55 0.50	29.06 55.52	<0.01 51.72	MAD MAD+Attn-S [5]	32.47 32.40	0.58	54.44 55.58	1.28 1.82
MAD (τ =40) MAD (τ =50)	32.16 32.14	0.50 0.49	86.20 98.01	76.15 91.51	MAD (τ =0) MAD (τ =2)	31.36 31.48	0.55 0.52	31.35 38.78	13.28 27.42					

Compared to direct inference methods MAD scales better by performing split convolutions for all timesteps and split attentions for the timesteps after the user-defined threshold.

MAD (τ =**4)** 31.41 0.48 62.82 61.70



MERGING AND SPLITTING DIFFUSION PATHS FOR SEMANTICALLY COHERENT PANORAMAS

We generate 1000 complex-scene prompts (GPT1k) and perform a quantitative comparison with different baselines.



MAD works just fine on different aspect ratios, on LDMs and LCMs

We generate horizontal panorama images at various aspect ratios, from 512x1024 to 512x4096. Top: we apply MAD to an LCM [4] model and compare with a baseline. Bottom: we apply MAD to an LDM [1] and show images generated by competitors, namely MultiDiffusion [2] and SyncDiffusion [3].



[3] Lee, Y., et al. "Syncdiffusion: Coherent montage via synchronized joint diffusions." NeurIPS 2024. [4] Luo, S., et al. "Latent consistency models: Synthesizing high-resolution images with few-step inference." arXiv 2023.

Fabio Quattrini, Vittorio Pippi, Silvia Cascianelli, Rita Cucchiara {fabio.quattrini, vittorio.pippi, silvia.cascianelli, rita.cucchiara}@unimore.it





[5] Jin, Z., Shen, X., Li, B., Xue, X.: Training-free Diffusion Model Adaptation for Variable-Sized Text-to-Image Synthesis. NeurIPS 2023. [6] Zhang, Lvmin, Anyi Rao, and Maneesh Agrawala. "Adding conditional control to text-to-image diffusion models." ICCV 2023.







You can use MAD off-the-shelf for controlled image generation

We apply the MAD operator in other Plug&Play tasks for controllable image generation, namely MultiDiffusion [2] and ControlNet [6], showing increased global coherence in the generated images



An ancient building in the style of an oil painting



Not all prompts are panorama-friendly...

Inference-time panorama generation approaches struggle with scenes or objects that do not fit well with the specified aspect ratio and prompts where the base model itself does not produce good-quality results. We provide examples with the prompt "A gothic cathedral nave", which does not fit well the horizontal aspect ratio. On a vertical canvas, results improve. Bottom: we provide examples with the prompt "A fancy living room".





... But some are!



