

MultiFusion: Fusing Pre-Trained Models for Multi-Lingual, Multi-Modal Image Generation

Marco Bellagente⁴⁺, Manuel Brack^{2,3}, Hannah Teufel^{1*}, Felix Friedrich^{3,6}, Björn Deiseroth^{1,3,6}, Constantin Eichenberg⁴, Andrew Dai¹, Robert J.N. Baldock², Souradeep Nanda⁵, Koen Oostermeijer¹, Andres Felipe Cruz-Salinas⁴, Patrick Schramowski^{2,3,6,8}, Kristian Kersting^{2,3,6,7*}, Samuel Weinbach^{1#}
 marco.bellagente@gmail.com, brack@cs.tu-darmstadt.de, hannah.teufel@aleph-alpha.com

¹Aleph Alpha, ²German Research Center for Artificial Intelligence (DFKI), ³Computer Science Department, TU Darmstadt, ⁴Stability AI, ⁵University of Texas, ⁶Hessian AI, ⁷Center for Cognitive Science, TU Darmstadt, ⁸LAION
 Work done while at Aleph Alpha. [#]Equal Contribution. ^{}Equal supervision



NEURAL INFORMATION
PROCESSING SYSTEMS



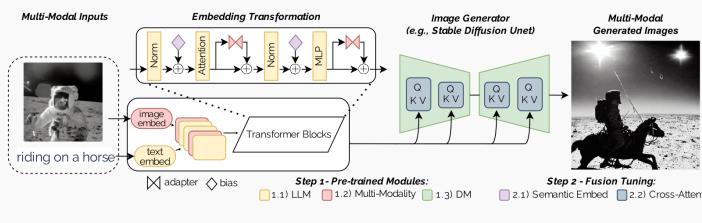
I. Method

How to enable the **intuitive, multilingual expression of complex and nuanced concepts** in **Diffusion Models** both **compute efficiently** and **without multilingual downstream training data**?

Interleaved multimodal, multilingual prompting



Fusion of pretrained models to a modular encoder, which transfers capabilities to the downstream diffusion model



To enable multimodal, multilingual prompting we replace Stable Diffusions Clip encoder by a custom modular one.

1.) An autoregressive **LLM pre-trained** on 5 languages enables **multilingualism**.

2.) Extending the LLM by an **image prefix** as well as adapters enables **multimodality**.

2.1) **Finetuning the biases** of the LLM provides embeddings, which capture the **semantic meaning** of the text prompt, thus simplifying the learning of mapping from embeddings to image outputs.

2.2) To **align the pre-trained Stable Diffusion** model (t4) with the embeddings of our modular encoder, we retrain the conditioning by finetuning the **cross-attention** weights.

II. Application

Image Composition:

MultiFusion increases **expressiveness in composition** through arbitrary and flexible prompting of image and text sequences.



Style Modification:

MultiFusion enables **simple style transfer** through one **reference** image capturing all the facets of a unique style such as color palette, composition, contrast, etc., making elaborate prompts obsolete. Additionally, MultiFusion enables **highly individual** prompting such as "In the style of a picture I drew".



Attention Manipulation:

Attention Manipulation allows us to **weight image and text** tokens at inference time and thus guide their **influence** on the resulting generation.



III. Evaluation

Image fidelity and image-text alignment:

We measure **image fidelity** and **image text-alignment** using the standard metrics FID-30k and Clip Scores. We find that MultiFusion prompted with text only performs **on par with stable diffusion** despite extension of the Encoder to support multiple languages and modalities.

Guidance Scale	FID-30k		CLIP Score (Text-to-Image)			
	SD	MFT (text)	MFT (image)	MFT (image)		
0.0	34.06	14.49	11.59	0.91	0.30	
0.5	37.73	12.35	10.29	1.18	0.31	0.29
1.0	30.93	9.90	8.63	6.03	0.30	0.29
2.0	0.94	12.21	8.61	1.15	0.20	0.28
3.0	46.69	31.82	24.42	9.83	0.27	0.25

Compositional Robustness through Multimodality:

Image composition is a known **limitation of Diffusion Models**. Evaluating on our new **benchmark MCC-250** we show that interleaved **multimodal prompting** leads to more **compositional robustness** as judged by humans.

Model	Two objects w/ correct colors (%)		Attribute Leakage	Interchanged Attributes	Missing Objects
	Stable Diffusion	MFT (text)			
Stable Diffusion	29.99				
Composable Diffusion	25.59				
MultiFusion(text)	21.66				
MultiFusion (multimodal)	58.35				

Stable Diffusion: a green apple and a red car; a blue book and a red cup.

Multilinguality:



Evaluating the **alignment** of prompt embeddings as well as generated images across **multiple languages** we show that **good embedding alignment** enables the **transfer** of multilingualism to downstream tasks even for task-specific monolingual training data.

