

# Diffusion Features in Stable Diffusion 3

Ilir Hajrullahu  
LMU Munich

I.Hajrullahu@campus.lmu.de

Ludwig Degenhardt  
LMU Munich

ludwig.degenhardt@lmu.de

## Abstract

This practical project explores the adaptation of DIFT from Stable Diffusion 2 to Stable Diffusion 3 and evaluates its performance in semantic correspondence tasks. DIFT extracts image correspondences from diffusion model features, and we analyze how SD3’s architectural shift from U-Net to transformer-based stacks affects its effectiveness. Through systematic experiments, we observe that SD3 performs worse, raising questions about transformer-based diffusion architectures for semantic correspondence tasks.

## 1. Introduction

The following is our project report for the practical "Image and Video Synthesis". Our topic, "Diffusion Features in recent T2I models," revolves around exploring semantic features in state-of-the-art text-to-image diffusion models. Our more concrete goal was to adapt the DIFT method for Stable Diffusion 3 and to investigate the adapted DIFT’s performance in semantic correspondence tasks. This presents us with an interesting assignment, as the original version of DIFT had excellent results with Stable Diffusion 2-1, outperforming OPENCLIP and DINO in a PCK (Percentage of Correct Keypoints) benchmark test on the SPair-71k dataset [7]. Our task raises various questions: How does the new model Stable Diffusion 3 and its changed architecture influence DIFT’s performance? Does it get better, stay the same, or even get worse? Can we learn something about diffusion models and their architectures in the process? Furthermore, adapting DIFT is interesting because of its possible applications, such as object recognition and edit propagation. Regarding the latter, DIFT can propagate edits across different images based on learned semantic correspondences, for example, adding a sticker to one image and propagating it accurately to another, even in different domains or categories [7].

The structure of this report is as follows: First we explain both DIFT and semantic correspondences, as well as

the architectural changes between SD2 and SD3 in section 038  
2. Section 039  
3 presents how we proceeded over the course of 040  
our project and how we approached our experimental 041  
tasks, followed by section 042 demonstrating the results of 043  
these experiments. Section 043 then provides a discussion of 044  
our results and suggestions for future work, and lastly, section 044  
concisely summarizes our project report.

## 2. Related work

### 2.1. DIFT and Semantic Correspondence

DIFT is a method introduced in the "Emergent Correspondence from Image Diffusion" paper, which explores how diffusion models inherently learn correspondences without explicit supervision [7]. Diffusion models, primarily designed for image generation, capture rich feature representations that we can leverage to establish correspondences between images. The name "DIFT" refers both to the method used to extract image features from pre-trained diffusion models and to the extracted features themselves, which encode semantic, geometric, and temporal correspondences. In the original paper, the researchers implemented DIFT with two different open-source diffusion models, Stable Diffusion 2-1 [6], and Ablated Diffusion Model [1], with Stable Diffusion showing better results for correspondence extraction [7].

The DIFT method works by extracting correspondence-rich image features. It does so by simulating the forward diffusion process and passing noisy images through a pre-trained diffusion model to obtain feature maps. These feature maps encode correspondences across images, and DIFT uses nearest-neighbor lookups using cosine similarity to find corresponding pixels.

For our project, we focused on exploring semantic correspondences. Semantic correspondence refers to finding corresponding pixel locations across different images where objects share similar semantic meanings. This means that we can use DIFT to identify, for example, the ears of cats in different images, even if they are in different perspectives, contexts, or even styles, such as a photograph and a drawing.

## 077 2.2. Stable Diffusion 2 and 3

078 Stable Diffusion 3 (SD3) marks a departure from the traditional latent diffusion approach used in Stable Diffusion 2  
 079 (SD2). SD3 adopts MMDiT (Multimodal Diffusion Trans-  
 080 former), a more flexible architecture that processes multiple  
 081 data modalities while enabling larger context windows. The  
 082 idea of MMDiT is based on Diffusion Transformers that use  
 083 separate weights for text and image and allow the flow of  
 084 information between the two modalities [5]. In contrast to  
 085 SD2’s reliance on a single text encoder (often CLIP), SD3  
 086 incorporates three different text encoders to capture nu-  
 087 ancied semantic information and accommodate more com-  
 088 plex or lengthy prompts.  
 089

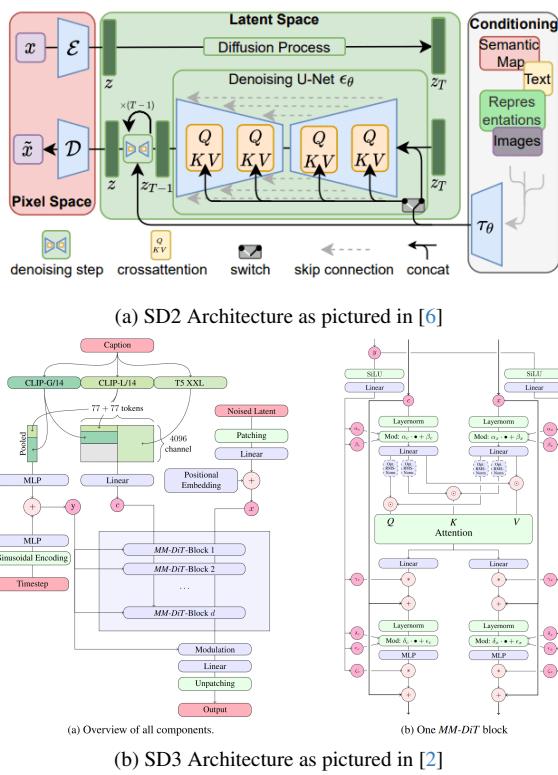


Figure 1. Stable Diffusion Architectures

090 Another prominent distinction is SD3’s transition away  
 091 from a U-Net-based noise prediction framework to a stack  
 092 of diffusion transformers. This change aligns the model’s  
 093 denoising operations more closely with Transformer archi-  
 094 tectures, making cross-attention more adaptive when  
 095 integrating textual features into the generation process. More-  
 096 over, SD3 employs rectified flow sampling, a refined tech-  
 097 nique for converting noisy intermediate representations into  
 098 clear output images with fewer inference steps and im-  
 099 proved coherence.

100 The underlying variational autoencoder (VAE) also re-  
 101 ceives a significant upgrade. Whereas SD2 typically used a

4-channel VAE, SD3 supports a 16-channel VAE, allowing  
 102 a richer latent representation of the image space and pro-  
 103 viding finer-grained detail in generated images. Taken to-  
 104 gether, these enhancements illustrate SD3’s push for larger  
 105 capacities, deeper representations, and more efficient sam-  
 106 pling, making it a noteworthy step forward in diffusion-  
 107 based image generation.

## 109 3. Methodology

This work is based on adapting the DIFT codebase, aligned  
 110 initially with SD2, to function seamlessly with the more re-  
 111 cent SD3 model. Our overall aim was to maintain DIFT’s  
 112 ability to measure and visualize semantic correspondences  
 113 between input prompts and generated images while taking  
 114 advantage of the updated architecture and textual condition-  
 115 ing mechanisms introduced in SD3. Initially, we cloned the  
 116 publicly available DIFT repository, which includes features  
 117 to extract intermediate representations of the diffusion pro-  
 118 cess and compare them across multiple images, prompts,  
 119 and layers in the U-Net.

To integrate SD3, we made several key modifications  
 121 to the original pipeline. One of the primary adjustments  
 122 involved updating the text encoder interface, since SD3  
 123 often relies on larger, more sophisticated text-encoder  
 124 backbones (such as T5). Our updated code handles the  
 125 expanded prompt-token capacity and the different attention  
 126 mechanisms these new encoders employ. Another signif-  
 127 icant change concerned the MMDiT structure. In some  
 128 variants of SD3, the network is equipped with additional  
 129 DiT-like blocks and an altered cross-attention mechanism.  
 130 Consequently, we reworked how DIFT retrieves and  
 131 processes these internal feature maps, which serve as  
 132 the basis for measuring semantic alignment—effectively  
 133 allowing us to see which parts of the image or latent rep-  
 134 resentation the model deems most relevant for a given prompt.

In order to systematically evaluate these changes, we set  
 137 up four distinct testing scenarios, all centered around the  
 138 concept of semantic correspondence. The first scenario in-  
 139 volved running DIFT on top of SD2.1 for baseline compari-  
 140 sons. By doing so, we wanted to get familiar with DIFT  
 141 and get a baseline for future comparisons with SD3. We  
 142 then moved on to the second scenario, where we tested  
 143 DIFT with the newly integrated SD3 pipeline. This step  
 144 aimed to verify how effectively DIFT captures semantic co-  
 145 herence with a bigger text encoder and transformer architec-  
 146 ture and whether the new features in SD3 lead to improved  
 147 or more detailed semantic maps.

The third testing scenario examined the effect of  
 149 different latent-space configurations on DIFT’s ability  
 150 to measure correspondences. Models in both SD2 and  
 151 SD3 typically apply a Variational Autoencoder (VAE) to  
 152 compress images into fewer channels, and this step can

154 influence how well certain semantic signals are preserved.  
 155 We, therefore, investigated 4-channel versus 16-channel  
 156 autoencoders, exploring how differences in compression  
 157 impact the detection of fine-grained features. For the  
 158 final scenario, we directed DIFT to extract features from  
 159 a specific transformer layer within the SD3 MMDiT  
 160 blocks, allowing us to pinpoint at which stage the strongest  
 161 semantic signals arise. This deeper look inside the network  
 162 is especially relevant, as DIFT was initially designed with a  
 163 focus on identifying how cross-attention layers in diffusion  
 164 models correlate image and text.

165  
 166 All four tests were conducted on a set of 14 images cho-  
 167 sen from the Spair-71k benchmark dataset to represent a va-  
 168 riety of real-world and synthetic scenes, lighting conditions,  
 169 and object categories [4]. These images enabled us to assess  
 170 the consistency of DIFT’s semantic alignment capabilities  
 171 under different circumstances. By comparing how well the  
 172 method performed in each scenario, we gained insight not  
 173 only into the robustness of DIFT itself but also into the ar-  
 174 chitectural and encoder-related changes in SD3. Overall,  
 175 these experiments thoroughly investigated how deeper tex-  
 176 tual contexts, larger model capacity, and varying levels of  
 177 latent compression affect the interpretation and visualiza-  
 178 tion of semantic features within diffusion-based generation  
 179 models.

## 180 4. Experiment results

181 Our experiments revealed a pronounced divergence in per-  
 182 formance between SD2 and SD3, as shown in Figure 2.  
 183 On the one hand, running DIFT with SD2 yielded crisp,  
 184 well-defined semantic correspondences. As illustrated in  
 185 the airplane example, the target image’s heatmap clearly  
 186 highlighted the areas of interest—such as the nose and es-  
 187 pecially the tail of the aircraft—corresponding to the cues  
 188 in the source image. The heatmap uses a Viridis scale, with  
 189 yellow demonstrating the highest semantic correspondence  
 190 and purple the lowest, normalized to values between 0 and  
 191 1. The resulting visualizations demonstrated that most se-  
 192 mantic information was captured in a concentrated and co-  
 193 herent manner, making it relatively easy to see how specific  
 194 regions in the target image map back to salient features in  
 195 the source.

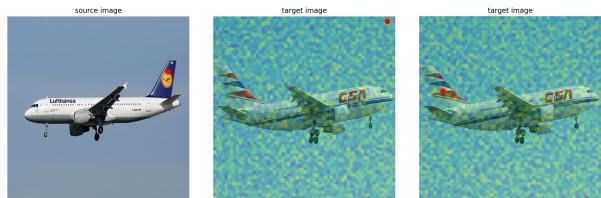
196 By contrast, DIFT performance degraded substantially  
 197 under SD3. The primary issue seemed to be that semantic  
 198 information became diffusely “spread out” across the entire  
 199 image rather than concentrated in a distinct region. This  
 200 diffuse pattern was apparent across all 14 test images.  
 201 In the case of the airplane, the target image heatmaps  
 202 essentially displayed broad noise patterns rather than  
 203 highlighting the plane’s fuselage or tail section. Similarly,  
 204 for the sheep example, neither the outline of the animal nor  
 205 its features stood out in any localized area of the target’s

latent representation; instead, the visualization indicated a  
 206 more uniform or random distribution of significance.

207  
 208



(a) Semantic Correspondence with SD2



(b) Semantic Correspondence with SD3



(c) Semantic Correspondence with SD2



(d) Semantic Correspondence with SD3

Figure 2. Comparison for semantic correspondence of SD2 and SD3

We also tested whether earlier layers of the SD3 model might contain cleaner, less “spread out” features. Specifically, we extracted and compared features at layer 2 and layer 6 out of the 12 available transformer blocks. Our initial hypothesis was that semantic structure might remain relatively intact at lower blocks before later transformer layers accumulate noise or disperse information.

209  
 210  
 211  
 212  
 213  
 214  
 215  
 216  
 217  
 218

Unfortunately, the results at both of these depths remained poor, as shown in Figure 3. Rather than isolating any coherent semantic points, the heatmaps still showed

219 generalized patterns across nearly the entire image with  
 220 no distinct alignment to specific objects or regions of  
 221 interest. The pattern was consistent across all 14 images  
 222 we tested, suggesting that the problem is systematic rather  
 223 than image-specific.  
 224

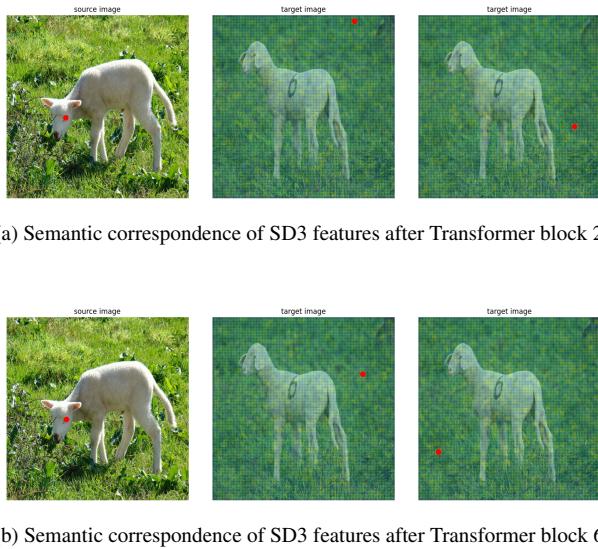


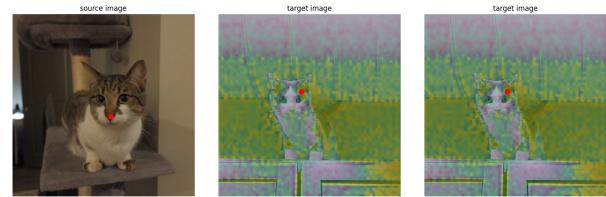
Figure 3. Semantic correspondence of SD3 after specific Transformer blocks

225 As for the results of our experiment using DIFT on vae-  
 226 encoded images without utilizing Stable Diffusion, in both  
 227 the 4-channel and 16-channel variants, DIFT could not ex-  
 228 tract meaningful semantic correspondences. However, we  
 229 could notice the richer latent representation of the encoded  
 230 image providing finer-grained details, making the target im-  
 231 age’s heatmap smoother, as shown in Figure 4.

232 Moreover, SD3 proved to be notably more resource-  
 233 intensive during inference compared to SD2 without offer-  
 234 ing a corresponding improvement in semantic correspon-  
 235 dence. GPU memory usage was higher, and inference  
 236 times were longer, yet the DIFT method could not exploit  
 237 any additional capacity that SD3 might have had to pro-  
 238 duce clearer correspondences. In effect, the newer diffusion  
 239 architecture’s improvements for general generation tasks  
 240 did not translate into benefits for DIFT’s targeted semantic  
 241 alignment metrics.

## 242 5. Discussion and Future Work

243 When using DIFT with SD3, the semantic information  
 244 we extract from the feature maps seems to be dispersed  
 245 throughout the images, as visible in the according figures  
 246 we provided in the section 4. To make the differences



(a) Semantic correspondence using 4-channel vae



(b) Semantic correspondence using 16-channel vae

Figure 4. Semantic correspondence with vae-encoded images

247 between the heatmaps in the SD2 and SD3 versions of  
 248 DIFT more comprehensible, we added a visualization of the  
 249 heatmap’s Viridis scale to the target images and plotted the  
 250 distribution of semantic correspondence values, both visible  
 251 in the example in Figure 5.

252 In the SD2 version of DIFT, most semantic correspon-  
 253 dence values lie between 0.0 and 0.4. The higher values are  
 254 located in semantically meaningful areas corresponding to  
 255 the input area chosen in the source image. In the case of  
 256 SD3, the heatmap’s values align in a pattern resembling a  
 257 standard distribution around medium-high values of 0.5 to  
 258 0.6, with only a few values around 0 and 1. This further  
 259 supports our assumption that the semantic information is  
 260 being spread through the image, with no indication of DIFT  
 261 extracting meaningful semantic correspondences from the  
 262 image features.

263 We assume that the reason for this behavior in the SD3  
 264 version of DIFT lies in its architecture. We hypothesize  
 265 that when the transformer projects the input to a linear  
 266 sequence to process it further, the semantic information is  
 267 dispersed or mixed throughout the image.

268 We have some propositions for future work to explore  
 269 the behavior of DIFT with state-of-the-art diffusion models,  
 270 especially diffusion models with transformer architectures,  
 271 more profoundly. First, it would be interesting to test the  
 272 SD3-adapted DIFT on geometric and temporal corre-  
 273 spondence tasks, which were outside the scope of our practical  
 274 project. Another possibility would be to adapt DIFT for  
 275 another different state-of-the-art text-to-image diffusion  
 276 model, such as FLUX [3].

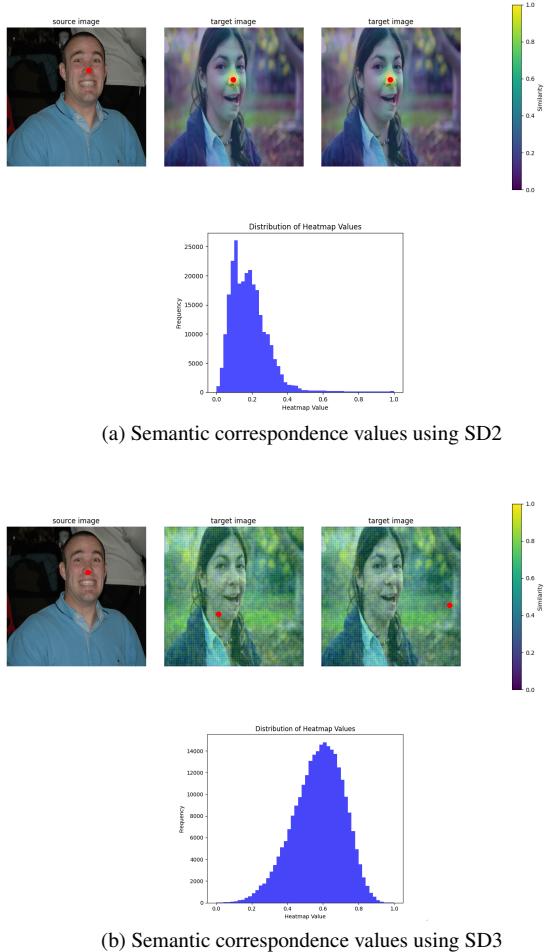


Figure 5. Visualization of semantic correspondence values of SD2 and SD3

278

279 By doing so, we could further explore if transformer-  
280 based architectures for text-to-image diffusion models per-  
281 form poorly in semantic correspondence tasks or if there  
282 is a problem inherent to SD3. Exploring such possibilities  
283 could even come to influence architectural decisions for fu-  
284 ture models.

285

## 6. Conclusion

286 This project explored the adaptation of DIFT from Sta-  
287 ble Diffusion 2 (SD2) to Stable Diffusion 3 (SD3), specif-  
288 ically evaluating its performance in semantic correspon-  
289 dence tasks. DIFT excelled in semantic correspondence  
290 tasks in its initial implementation using SD2. Our find-  
291 ings indicate a notable decline in DIFT’s effectiveness when  
292 applied to SD3, raising questions about the suitability of  
293 transformer-based diffusion architectures for such tasks.

The key observation was that semantic information in SD3’s feature maps appeared more diffusely spread across images rather than being localized in meaningful regions, as seen with SD2. This degradation in correspondence quality may be attributed to SD3’s architectural shift, particularly its use of transformer-based MMDiT instead of U-Net. Despite the model’s improved generative capabilities, it does not translate into better performance in semantic correspondence tasks for DIFT-based feature extraction.

Future research should investigate whether this limitation is inherent to all transformer-based diffusion models or specific to SD3. Testing alternative architectures such as FLUX and evaluating how DIFT with SD3 performs in geometric and temporal correspondence tasks could provide further insights.

## References

- [1] Prafulla Dhariwal and Alexander Nichol. Diffusion models beat gans on image synthesis. *Advances in neural information processing systems*, 34:8780–8794, 2021. 1
- [2] Patrick Esser, Sumith Kulal, Andreas Blattmann, Rahim Entezari, Jonas Müller, Harry Saini, Yam Levi, Dominik Lorenz, Axel Sauer, Frederic Boesel, Dustin Podell, Tim Dockhorn, Zion English, Kyle Lacey, Alex Goodwin, Yannik Marek, and Robin Rombach. Scaling rectified flow transformers for high-resolution image synthesis, 2024. 2
- [3] Black Forest Labs. Flux ai official github repository: <https://github.com/black-forest-labs/flux>. 4
- [4] Juhong Min, Jongmin Lee, Jean Ponce, and Minsu Cho. Spair-71k: A large-scale benchmark for semantic correspondence, 2019. 3
- [5] William Peebles and Saining Xie. Scalable diffusion models with transformers, 2023. 2
- [6] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models, 2022. 1, 2
- [7] Luming Tang, Menglin Jia, Qianqian Wang, Cheng Perng Phoo, and Bharath Hariharan. Emergent correspondence from image diffusion, 2023. 1

294  
295  
296  
297  
298  
299  
300  
301  
302  
303  
304  
305  
306  
307  
308

309  
310  
311  
312  
313  
314  
315  
316  
317  
318  
319  
320  
321  
322  
323  
324  
325  
326  
327  
328  
329  
330  
331