Anonymous CVPR submission

Paper ID 38

## Abstract

001 Do the rich representations of multi-modal diffusion transformers (DiTs) exhibit unique properties that enhance their 002 interpretability? We introduce CONCEPTATTENTION, a 003 novel method that leverages the expressive power of DiT 004 005 attention layers to generate high-quality saliency maps that precisely locate textual concepts within images. Without re-006 quiring additional training, CONCEPTATTENTION repur-007 008 poses the parameters of DiT attention layers to produce highly contextualized concept embeddings, contributing the 009 major discovery that performing linear projections in the 010 output space of DiT attention layers yields significantly 011 012 sharper saliency maps compared to commonly used cross-013 attention mechanisms. Remarkably, CONCEPTATTEN-014 TION even achieves state-of-the-art performance on zeroshot image segmentation benchmarks, outperforming 15 015 other zero-shot interpretability methods on the ImageNet-016 Segmentation dataset and on a single-class subset of Pas-017 018 calVOC. Our work contributes the first evidence that the 019 representations of multi-modal DiT models like Flux are highly transferable to vision tasks like segmentation, even 020 outperforming multi-modal foundation models like CLIP. 021

#### **022 1. Introduction**

Diffusion models have recently gained widespread popular-023 ity, emerging as the state-of-the-art approach for a variety of 024 generative tasks, particularly text-to-image synthesis [21]. 025 These models transform random noise into photorealistic 026 027 images guided by textual descriptions, achieving unprecedented fidelity and detail. Despite the impressive generative 028 capabilities of diffusion models, our understanding of their 029 internal mechanisms remains limited. Diffusion models op-030 031 erate as black boxes, where the relationships between input 032 prompts and generated outputs are visible, but the decisionmaking processes that connect them are hidden from human 033 understanding. 034

Existing work on interpreting T2I models has predominantly focused on UNet-based architectures [19, 21], which utilize shallow cross-attention mechanisms between prompt



Figure 1. CONCEPTATTENTION interprets the representations of multi-modal diffusion transformers by producing highquality saliency maps of textual concepts. We compare to the activations of the cross attention mechanisms of the DiT.

embeddings and image patch representations. UNet cross 038 attention maps can produce high-fidelity saliency maps 039 that predict the location of textual concepts [27] and have 040 found numerous applications in tasks like image editing 041 [5, 12]. However, the interpretability of more recent 042 multi-modal diffusion transformers (DiTs) remains under-043 explored. DiT-based models have recently replaced UNets 044 [22] as the state-of-the-art architecture for image genera-045 tion, with models such as Flux [13] and SD3 [8] achieving 046 breakthroughs in text-to-image generation. The rapid ad-047 vancement and enhanced capabilities of DiT-based models 048 highlight the critical importance of methods that improve 049 their interpretability, transparency, and safety. 050

In this work, we propose CONCEPTATTENTION, a novel 051 method that leverages the representations of multi-modal 052 DiTs to produce high-fidelity saliency maps that localize 053 textual concepts within images. Our method provides insight into the rich semantics of DiT representations. CON-CEPTATTENTION is lightweight and requires no additional 056



Figure 2. CONCEPTATTENTION augments multi-modal DiTs with a sequence of concept embeddings that can be used to produce saliency maps. (Left) An unmodified multi-modal attention (MMATTN) layer processes both prompt and image tokens. (Right) CON-CEPTATTENTION augments these layers without impacting the image appearance to create a set of contextualized concept tokens.

057 training, instead it repurposes the existing parameters of 058 DiT attention layers. CONCEPTATTENTION works by producing a set of rich contextualized text embeddings each 059 corresponding to visual concepts (e.g. "dragon", "sun"). By 060 linearly projecting these concept embeddings and the im-061 062 age we can produce rich saliency maps that are even higher 063 quality than commonly used cross attention maps.

We evaluate the efficacy of CONCEPTATTENTION in a 064 zero-shot semantic segmentation task on real world images. 065 066 We compare our interpretative maps against annotated segmentations to measure the accuracy and relevance of the 067 068 attributions generated by our method. Our experiments and 069 extensive comparisons demonstrate that CONCEPTATTEN-TION provides valuable insights into the inner workings of 070 071 these otherwise complex black-box models. By explaining 072 the meaning of the representations of generative models our method paves the way for advancements in interpretability, 073 074 controllability, and trust in generative AI systems. 075

In summary, we contribute:

- CONCEPTATTENTION, a method for interpreting 076 077 text-to-image diffusion transformers. Our method requires no additional training, by leveraging the repre-078 sentations of multi-modal DiTs to generate highly inter-079 pretable saliency maps that depict the presence of arbi-080 trary textual concepts (e.g. "dragon", "sky", etc.) in im-081 082 ages (as shown in Figure 1).
- · The novel discovery that the output vectors of atten-083 tion operations produce higher-quality saliency maps 084 than cross attentions. CONCEPTATTENTION repurposes 085 086 the parameters of DiT attention layers to produce a set of rich textual embeddings corresponding to different con-087 cepts, something that is uniquely enabled by multi-modal 088 DiT architectures. By performing linear projections be-089 tween these concept embeddings and image patch repre-090 sentations in the attention output space we can produce 091 092 high quality saliency maps.

• CONCEPTATTENTION generalizes to achieve state-093 of-the-art performance in zero-shot segmentation on 094 benchmarks like ImageNet Segmentation and Pascal 095 VOC. We achieve superior performance to a diverse set of 096 zero-shot interpretability methods based on various foun-097 dation models like CLIP, DINO, and UNet-based diffu-098 sion models; this highlights the potential for the repre-099 sentations of DiTs to transfer to important downstream 100 vision tasks like segmentation. 101

# 2. Preliminaries

## 2.1. The Anatomy of a Multi-modal DiT Layer

Multi-modal DiTs like Flux and Stable Diffusion 3 leverage multi-modal attention layers (MMATTN) that process a combination of textual tokens and image patches. There are two key classes of layers: one that keeps separate residual streams for each modality and one that uses a single stream. In this work, we take advantage of the properties of these dual stream layers, which we refer to as multi-modal attention layers (MMATTNs).

The input to a given layer is a sequence of image patch representations  $x \in \mathbb{R}^{h \times w \times d}$  and prompt token embeddings  $p \in \mathbb{R}^{l \times d}$ . The initial prompt embeddings at the beginning of the network are formed by taking the T5 [20] embeddings of the prompt tokens.

Following [18], each MMATTN layer leverages a set of adaptive layer norm modulation layers [28], conditioned on the time-step and global CLIP vector. An adaptive layernorm operation is applied to the input image and text embeddings. The final modulated outputs are then residually added back to the original input. Notably, the image and text modalities are kept in separate residual streams. The exact details of this operation are omitted for brevity.

The key workhorse in MMATTN layers is the familiar 125 multi-head self attention operation. The prompt and im-126

103 104 105

106

107

108

109

110

111

112

113

114

115

116

117

118

119

120

121

122

123

124

102

133

171

173

175

187

188

200

201

202



Figure 3. CONCEPTATTENTION generates saliency maps for multiple concepts simultaneously, even concepts not in the prompt.

127 age embeddings have separate learned key, value, and query 128 projection matrices which we refer to as  $K_x, Q_x, V_x$  for im-129 ages and  $K_p, Q_p, V_p$  for text. The keys, queries, and values 130 for both modalities are collectively denoted  $q_{xp}, k_{xp}$ , and 131  $v_{xp}$ , where for example  $k_{xp} = [K_x x_1, \ldots, K_p p_1 \ldots]$ . A 132 self attention operation is then performed

$$o_x, o_p = \operatorname{softmax}(q_{xp}k_{xp}^T)v_{xp} \tag{1}$$

Here we refer to  $o_x$  and  $o_p$  as the *attention output* vectors. Another linear layer is then applied to these outputs and added to a separate residual streams weighted according to the output of the modulation layer. This gives us updated embeddings  $x^{L+1}$  and  $p^{L+1}$  which are given as input to the next layer.

## **3. Proposed Method: CONCEPTATTENTION**

We introduce CONCEPTATTENTION, a method for gener-141 ating high quality saliency maps depicting the location of 142 143 textual concepts in images. CONCEPTATTENTION works 144 by creating a set of contextualized concept embeddings for simple textual concepts (e.g. "cat", "sky", etc.). These 145 concept embeddings are sequentially updated alongside the 146 text and image embeddings, so they match the structure 147 that each MMATTN layer expects. However, unlike the 148 149 text prompt, concept embeddings do not impact the appear-150 ance of the image. We can produce high-fidelity saliency maps by projecting image patch representations onto the 151 concept embeddings. CONCEPTATTENTION requires no 152 additional training and has minimal impact on model la-153 tency and memory footprint. A high level depiction of our 154 155 methodology is shown in Figure 2.

Using CONCEPTATTENTION. The user specifies a set 156 of r single token concepts, like "cat", "sky", etc. which 157 are passed through a T5 encoder to produce an initial em-158 bedding  $c^0$  for each concept. For each MMATTN layer (in-159 dexed by L) we layer-normalize the input concept embed-160 161 dings  $c^{L}$  and repurpose the text prompt's projection matrices (i.e.  $K_p, Q_p, V_p$ ), to produce a set of keys, values, and 162 queries (i.e.  $k_c = [K_p c_1, \dots, K_p c_k]$ ). 163

164 **One-directional Attention Operation.** We would like 165 to update our concept embeddings so they are compatible 166 with subsequent layers, but also prevent them from impact-167 ing the image tokens. Let  $k_x$  and  $v_x$  be the keys and values 168 of the image patches x respectively. We can concatenate the 169 image and concept keys to get

170 
$$k_{xc} = [K_x x_1 \dots, K_x x_n, K_p c_1 \dots, K_p c_r]$$
 (2)

and the image and concept values to get

$$v_{xc} = [V_x x_1 \dots, V_x x_n, V_p c_1 \dots, V_p c_r]$$
 (3) 172

We can then perform the following attention operation

$$o_c = \operatorname{softmax}(q_c k_{xc}^T) v_{xc} \tag{4}$$

which produces a set of concept output embeddings.

Notice, that instead of just performing a cross attention 176 (i.e. softmax $(q_c k_x^T)v_x$ ) our approach leverages both cross 177 attention from the image patches to the concepts and self 178 attention among the concepts. We found that performing 179 both operations improves performance on downstream eval-180 uation tasks like segmentation (See Table 5). We hypoth-181 esize this is because it allows the concept embeddings to 182 repel from each other, avoiding redundancy between con-183 cepts. Meanwhile, the image patch and prompt tokens ig-184 nore the concept tokens and attend only to each other as in 185

$$o_x, o_p = \operatorname{softmax}(q_{xp}k_{xp}^T)v_{xp}.$$
(5) 186

A diagram of these operations is shown in Fig. 9 (b) in the Appendix.

A Concept Residual Stream. The above operations cre-189 ate a residual stream of concept embeddings distinct from 190 the image and patch embeddings. Following the pretrained 191 transformer's design, after the MMATTN we apply another 192 projection matrix P and MLP, adding the result residually 193 to  $c^{L}$ . We apply an adaptive layernorm at the end of the 194 attention operation which outputs several values: a scale  $\gamma$ , 195 shift  $\beta$ , and some gating values  $\alpha_1$  and  $\alpha_2$ . The residual 196 stream is then updated as 197

$$c^{L+1} \leftarrow c^L + \alpha_1(Po_c) \tag{6}$$

$$c^{L+1} \leftarrow c^{L+1} + \alpha_2 \operatorname{MLP}\left((1+\gamma)\operatorname{horm}(c^{L+1}) + \beta\right)$$
 199
(7)

where  $\leftarrow$  denotes assignment. The parameters from each of our modulation, projection, and MLP layers are the same as those used to process the text prompt.

Saliency Maps in the Attention Output Space. These203concept embeddings can be combined with the image patch204embeddings to produce saliency maps for each layer L.205Specifically, we found that taking a simple dot-product similarity between the image output vectors  $o_x$  and concept output vectors  $o_c$  produces high-quality saliency maps208208

209

250

251

252

253

254

255

256

257

258

259

260

261

262

263

264

265

266

267

268

269

270

271

272

286

287

288

289

290

291

|                    | ImageNet-Seg |       |       | PascalVOC |       |       |
|--------------------|--------------|-------|-------|-----------|-------|-------|
| Method             | Acc          | mIoU  | mAP   | Acc       | mIoU  | mAP   |
| CLIP SA [7]        | 67.84        | 46.37 | 80.24 | 68.51     | 44.81 | 83.63 |
| TextSpan [10]      | 75.21        | 54.50 | 81.61 | 75.00     | 56.24 | 84.79 |
| TransInterp [4]    | 79.70        | 61.95 | 86.03 | 76.90     | 57.08 | 86.74 |
| CLIPasRNN [26]     | 74.05        | 58.80 | 84.80 | 61.76     | 41.48 | 76.57 |
| DINOv2 SA [16]     | 77.39        | 63.12 | 84.19 | 79.65     | 57.61 | 87.26 |
| DAAM [27]          | 78.47        | 64.56 | 88.79 | 72.76     | 55.95 | 88.34 |
| Cross Attn (SD3.5) | 77.80        | 63.67 | 83.50 | 80.22     | 61.46 | 86.97 |
| Cross Attn (Flux)  | 74.92        | 59.90 | 87.23 | 80.37     | 54.77 | 89.08 |
| Ours (SD3.5)       | 81.92        | 67.47 | 90.79 | 83.90     | 69.93 | 90.02 |
| Ours (Flux)        | 83.07        | 71.04 | 90.45 | 87.85     | 76.45 | 90.19 |

Table 1. CONCEPTATTENTION outperforms a variety of Diffusion, DINO, and CLIP interpretability methods on ImageNet-Segmentation and PascalVOC (Single Class). An extended version of this table with additional baselines is shown in Table 2 of the Appendix.

$$\phi(o_x, o_c) = \operatorname{softmax}(o_x o_c^T). \tag{8}$$

This is in contrast to cross attention maps which are between the image keys  $k_x$  and prompt queries  $q_p$ .

We can aggregate the information from multiple layers by averaging them  $\frac{1}{|L|} \sum_{L=1}^{|L|} \phi(o_x^L, o_c^L)$  where |L| denotes the number of MMATTN layers (Flux has |L| = 19). These attention output space maps are unique to MM-DiT models as they leverage *concept embeddings* corresponding to textual concepts which fundamentally can not be produced by UNet-based models.

## **4.** Experiments

We are interested in investigating (1) the efficacy of CON-CEPTATTENTION to generate highly localized and semantically meaningful saliency maps, and (2) understand the transferability of multi-modal DiT representations to important downstream vision tasks. Zero-shot image segmentation is a natural choice for achieving these goals.

226 Datasets. For our key evaluation we leverage two zeroshot-image segmentation datasets: ImageNet-Segmentation 227 [11] which was introduced for this task in [4, 10]. The 228 dataset contains 4,276 images from 455 classes. Each im-229 age depicts a single central object or subject, which makes 230 231 it a good method for comparing to a variety of interpretability methods that only generate a single saliency map, not a 232 233 class specific one. For the second dataset we leverage PascalVOC 2012 benchmark [9]. We investigate both a single 234 235 class (930 images) and multi-class split (1,449 images) of 236 this dataset. We also evaluate on a multi-class segmentation task in Appendix A.1. 237

Experimental Details. For our first task we closely follow the established evaluation framework from [10] and
[4]. We perform this evaluation setup on both ImageNet-Segmentation and a subset of 930 PascalVOC images con-

taining only a single class. For each method we assume the 242 class present in the image is known and use simplified de-243 scriptions of each ImageNet class (e.g. "Maltese dog"  $\rightarrow$ 244 "dog) this allows the concepts to be captured by a single 245 token. We construct a concept vocabulary for each image 246 composed of this target class and a set of fixed background 247 concepts for all examples (e.g. "background", "grass", 248 "sky"). 249

**Baselines.** We compare our approach to a variety of zero-shot interpretability methods which leverage several different multi-modal foundation models. We omit numerous models from Table 1 due to space constraints, but we have an extended version with additional baselines in Table 2 of the Appendix. We compare to numerous CLIP interpretability methods [1, 2, 4, 7, 10, 24, 26], the self-attentions of various DINO models [3, 6, 16], and approaches that leverage the cross attentions of UNet-based diffusion models [14, 27], and the cross attention maps of Flux and SD3.5 Turbo [23].

**Implementation Details.** For our experiments we implemented ConceptAttention for both the distilled Flux-Schnell model and Stable Diffusion 3.5 Turbo [23] in Py-Torch [17]. We encode real images with the DiT by first mapping them to the VAE latent space and then adding varying degrees of Gaussian noise before passing them through the DiT. We then cache all of the concept output  $o_c$  and image output vectors  $o_x$  from each MMATTN layer. At the end of generation we then construct our concept saliency maps for each layer and average them over all layers of interest. In our experiments we leverage the activations from the last 10 of the 19 MMATTN layers.

**Ouantitative Metrics.** Each method produces a set 273 of scalar raw scores for each image patch which we then 274 threshold based on the mean value to produce a binary 275 segmentation prediction. We compare each method using 276 standard segmentation evaluation metrics, namely: mean 277 Intersection over Union (mIoU), patch/pixelwise accuracy 278 (Acc), and mean Average Precision (mAP). Accuracy alone 279 is an insufficient metric as our dataset is highly imbalanced, 280 mIoU is significantly better, and mAP captures threshold 281 agnostic segmentation capability. We found that CONCEP-282 TATTENTION significantly out performs all of the baselines 283 we tested across all three metrics (Tab. 1). This is true for 284 diffusion, CLIP, and DINO based interpretability methods. 285

**Qualitative Evaluation.** We also show qualitative results comparing the segmentation performance to each baseline in Figures 1, 3 and in Appendix B.

**Multi Object Image Segmentation.** We also performed a quantitative evaluation for images with multiple objects, with details outlined in Appendix A.1.

Ablation Studies We performed various ar-292 chitectural ablation studies shown in Appendix 293 A.2 294

329

330

331

332

333

363

364

365

366

367

368

369

370

371

372

373

374

375

376

377

378

379

380

381

382

383

384

385

386

387

388

389

390

391

392

393

394

395

396

397

398

399

400

401

402

403

404

405

406

### 295 References

- [1] Samira Abnar and Willem Zuidema. Quantifying Attention
   Flow in Transformers, 2020. arXiv:2005.00928 [cs]. 4, 7
- [2] Alexander Binder, Grégoire Montavon, Sebastian Bach, Klaus-Robert Müller, and Wojciech Samek. Layer-wise Relevance Propagation for Neural Networks with Local Renormalization Layers, 2016. arXiv:1604.00825 [cs]. 4, 7
- 302 [3] Mathilde Caron, Hugo Touvron, Ishan Misra, Hervé Jégou,
  303 Julien Mairal, Piotr Bojanowski, and Armand Joulin. Emerg304 ing Properties in Self-Supervised Vision Transformers, 2021.
  305 arXiv:2104.14294 [cs]. 4, 7
- 306 [4] Hila Chefer, Shir Gur, and Lior Wolf. Transformer
  307 Interpretability Beyond Attention Visualization, 2021.
  308 arXiv:2012.09838 [cs]. 4, 7
- [5] Hila Chefer, Yuval Alaluf, Yael Vinker, Lior Wolf, and
  Daniel Cohen-Or. Attend-and-Excite: Attention-Based Semantic Guidance for Text-to-Image Diffusion Models. *ACM Transactions on Graphics*, 42(4):148:1–148:10, 2023. 1
- [6] Timothée Darcet, Maxime Oquab, Julien Mairal, and Piotr
  Bojanowski. Vision Transformers Need Registers, 2024.
  arXiv:2309.16588 [cs]. 4, 7
- [7] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov,
  Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner,
  Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. An Image is
  Worth 16x16 Words: Transformers for Image Recognition at
  Scale, 2021. arXiv:2010.11929 [cs]. 4, 7
- [8] 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.
  arXiv:2403.03206. 1
  - [9] Mark Everingham, S. M. Ali Eslami, Luc Van Gool, Christopher K. I. Williams, John Winn, and Andrew Zisserman. The Pascal Visual Object Classes Challenge: A Retrospective. *International Journal of Computer Vision*, 111(1):98–136, 2015. 4
- [10] Yossi Gandelsman, Alexei A. Efros, and Jacob Steinhardt.
  Interpreting CLIP's Image Representation via Text-Based
  Decomposition, 2024. arXiv:2310.05916 [cs]. 4, 7
- [11] Matthieu Guillaumin, Daniel Küttel, and Vittorio Ferrari.
   ImageNet Auto-Annotation with Segmentation Propagation.
   *International Journal of Computer Vision*, 110(3):328–348,
   2014. 4
- [12] Amir Hertz, Ron Mokady, Jay Tenenbaum, Kfir Aberman, Yael Pritch, and Daniel Cohen-Or. Prompt-to-Prompt Image Editing with Cross Attention Control, 2022.
  arXiv:2208.01626 [cs]. 1
- 345 [13] Black Forest Labs. FLUX, 2023. 1
- [14] Ziyi Li, Qinye Zhou, Xiaoyun Zhang, Ya Zhang, Yanfeng
  Wang, and Weidi Xie. Open-vocabulary Object Segmentation with Diffusion Models, 2023. arXiv:2301.05221 [cs].
  4
- [15] Pablo Marcos-Manchón, Roberto Alcover-Couso, Juan C.
   SanMiguel, and Jose M. Martínez. Open-Vocabulary Atten-

tion Maps with Token Optimization for Semantic Segmentation in Diffusion Models, 2024. arXiv:2403.14291 [cs]. 7 353

- [16] Maxime Oquab, Timothée Darcet, Théo Moutakanni, Huy 354 Vo, Marc Szafraniec, Vasil Khalidov, Pierre Fernandez, 355 Daniel Haziza, Francisco Massa, Alaaeldin El-Nouby, Mah-356 moud Assran, Nicolas Ballas, Woiciech Galuba, Russell 357 Howes, Po-Yao Huang, Shang-Wen Li, Ishan Misra, Michael 358 Rabbat, Vasu Sharma, Gabriel Synnaeve, Hu Xu, Hervé Je-359 gou, Julien Mairal, Patrick Labatut, Armand Joulin, and Pi-360 otr Bojanowski. DINOv2: Learning Robust Visual Features 361 without Supervision, 2024. arXiv:2304.07193 [cs]. 4, 7 362
- [17] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Köpf, Edward Yang, Zach DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. PyTorch: An Imperative Style, High-Performance Deep Learning Library, 2019. arXiv:1912.01703 [cs]. 4
- [18] William Peebles and Saining Xie. Scalable Diffusion Models with Transformers, 2023. arXiv:2212.09748. 2
- [19] Dustin Podell, Zion English, Kyle Lacey, Andreas Blattmann, Tim Dockhorn, Jonas Müller, Joe Penna, and Robin Rombach. SDXL: Improving Latent Diffusion Models for High-Resolution Image Synthesis, 2023. arXiv:2307.01952 [cs]. 1
- [20] Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. Exploring the Limits of Transfer Learning with a Unified Text-to-Text Transformer, 2023. arXiv:1910.10683 [cs]. 2
- [21] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-Resolution Image Synthesis with Latent Diffusion Models, 2022. arXiv:2112.10752 [cs]. 1
- [22] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-Net: Convolutional Networks for Biomedical Image Segmentation, 2015. arXiv:1505.04597. 1
- [23] Axel Sauer, Frederic Boesel, Tim Dockhorn, Andreas Blattmann, Patrick Esser, and Robin Rombach. Fast High-Resolution Image Synthesis with Latent Adversarial Diffusion Distillation, 2024. arXiv:2403.12015 [cs]. 4
- [24] Ramprasaath R. Selvaraju, Michael Cogswell, Abhishek Das, Ramakrishna Vedantam, Devi Parikh, and Dhruv Batra. Grad-CAM: Visual Explanations from Deep Networks via Gradient-based Localization, 2019. arXiv:1610.02391. 4
- [25] Ramprasaath R. Selvaraju, Michael Cogswell, Abhishek Das, Ramakrishna Vedantam, Devi Parikh, and Dhruv Batra. Grad-CAM: Visual Explanations from Deep Networks via Gradient-based Localization. *International Journal of Computer Vision*, 128(2):336–359, 2020. arXiv:1610.02391 [cs]. 7
- [26] Shuyang Sun, Runjia Li, Philip Torr, Xiuye Gu, and Siyang Li. CLIP as RNN: Segment Countless Visual Concepts without Training Endeavor, 2024. arXiv:2312.07661 [cs]. 4, 7
- [27] Raphael Tang, Linqing Liu, Akshat Pandey, Zhiying Jiang,
   Gefei Yang, Karun Kumar, Pontus Stenetorp, Jimmy Lin,
   408

409and Ferhan Ture. What the DAAM: Interpreting Stable Dif-410fusion Using Cross Attention, 2022. arXiv:2210.04885 [cs].4111, 4, 7

- [28] Jingjing Xu, Xu Sun, Zhiyuan Zhang, Guangxiang Zhao, and
  Junyang Lin. Understanding and Improving Layer Normalization, 2019. arXiv:1911.07013 [cs]. 2
- 415[29] Zhuoyi Yang, Jiayan Teng, Wendi Zheng, Ming Ding, Shiyu416Huang, Jiazheng Xu, Yuanming Yang, Wenyi Hong, Xiao-417han Zhang, Guanyu Feng, Da Yin, Yuxuan Zhang, Weihan418Wang, Yean Cheng, Bin Xu, Xiaotao Gu, Yuxiao Dong, and419Jie Tang. CogVideoX: Text-to-Video Diffusion Models with420An Expert Transformer, 2025. arXiv:2408.06072 [cs]. 13,42114
- [30] Jinghao Zhou, Chen Wei, Huiyu Wang, Wei Shen, Cihang
  Xie, Alan Yuille, and Tao Kong. iBOT: Image BERT PreTraining with Online Tokenizer, 2022. arXiv:2111.07832
  [cs]. 7

|                   |              | ImageNet-Segmentation        |       | ntation | PascalVOC (Single Class) |       | Class) |
|-------------------|--------------|------------------------------|-------|---------|--------------------------|-------|--------|
| Method            | Architecture | $\operatorname{Acc}\uparrow$ | mIoU↑ | mAP↑    | Acc $\uparrow$           | mIoU↑ | mAP↑   |
| LRP [2]           | CLIP ViT     | 51.09                        | 32.89 | 55.68   | 48.77                    | 31.44 | 52.89  |
| Partial-LRP [2]   | CLIP ViT     | 76.31                        | 57.94 | 84.67   | 71.52                    | 51.39 | 84.86  |
| Rollout [1]       | CLIP ViT     | 73.54                        | 55.42 | 84.76   | 69.81                    | 51.26 | 85.34  |
| ViT Attention [7] | CLIP ViT     | 67.84                        | 46.37 | 80.24   | 68.51                    | 44.81 | 83.63  |
| GradCAM [25]      | CLIP ViT     | 64.44                        | 40.82 | 71.60   | 70.44                    | 44.90 | 76.80  |
| TextSpan [10]     | CLIP ViT     | 75.21                        | 54.50 | 81.61   | 75.00                    | 56.24 | 84.79  |
| TransInterp [4]   | CLIP ViT     | 79.70                        | 61.95 | 86.03   | 76.90                    | 57.08 | 86.74  |
| CLIPasRNN [26]    | CLIP ViT     | 74.05                        | 58.80 | 84.80   | 61.76                    | 41.48 | 76.57  |
| OVAM [15]         | SDXL UNet    | 79.41                        | 65.02 | 88.12   | 73.50                    | 58.12 | 87.91  |
| DINO SA [3]       | DINO ViT     | 81.97                        | 69.44 | 86.12   | 80.71                    | 64.33 | 88.90  |
| DINOv2 SA [16]    | DINOv2 ViT   | 77.39                        | 63.12 | 84.19   | 79.65                    | 57.61 | 87.26  |
| DINOv2 Reg SA [6] | DINOv2 Reg   | 72.04                        | 56.31 | 80.83   | 77.16                    | 56.60 | 86.35  |
| iBOT SA [30]      | iBOT ViT     | 76.34                        | 61.73 | 82.04   | 74.96                    | 55.80 | 85.26  |
| DAAM [27]         | SDXL UNet    | 78.47                        | 64.56 | 88.79   | 72.76                    | 55.95 | 88.34  |
| DAAM [27]         | SD2 UNet     | 64.52                        | 47.62 | 78.01   | 64.28                    | 45.01 | 83.04  |
| Cross Attention   | Flux DiT     | 74.92                        | 59.90 | 87.23   | 80.37                    | 54.77 | 89.08  |
| Cross Attention   | SD3.5 DiT    | 77.80                        | 63.67 | 83.50   | 80.22                    | 61.46 | 86.97  |
| CONCEPTATTENTION  | SD3.5 DiT    | 81.92                        | 67.47 | 90.79   | 83.90                    | 69.93 | 90.02  |
| CONCEPTATTENTION  | Flux DiT     | 83.07                        | 71.04 | 90.45   | 87.85                    | 76.45 | 90.19  |

Table 2. CONCEPTATTENTION outperforms a variety of Diffusion, DINO, and CLIP ViT interpretability methods on ImageNet-Segmentation and PascalVOC (Single Class). An extended version of this table with additional baselines is shown in Appendix ??.

# A. Additional Quantitative Results

## A.1. Multi-object Semantic Segmentation

We also wanted to evaluate the capabilities of our method at differentiating between multiple classes in an image. However, only a subset of methods produce distinct saliency maps for open ended classes. For this experiment we compare to DAAM using a SDXL backbone, TextSpan using a CLIP backbone, and the raw cross attentions of Flux. Instead of binarizing the image to produce segmentations, for each patch we predict the concept with the highest score. We used pixelwise accuracy and mIoU as our evaluation metrics and found that our method significantly outperformed the baselines (See Table 3). We also show qualitative results of our approach differentiating between multiple concepts in a single image in Figures 1, ?? and we show more results in Appendix B. 428 429 430 430 431 432 433

## A.2. Ablation Studies

We perform several ablation studies to investigate the impact of various architectural choices and hyperparameters on the performance of CONCEPTATTENTION.

Impact of Layer Depth on SegmentationWe hypothesized that deeper MMATTN layers in the DiT would have more438refined representations that better transfer to segmentation. This was confirmed by our evaluation (see Figure 4). We pull439features from each diffusion layer and evaluated the segmentation performance of these features on ImageNet Segmentation.440We also compare the performance of combining all layers simultaneously, which we found performs better than any individual441layer.442

Impact of Diffusion Timestep on SegmentationWe add Gaussian noise to encoded images before passing them to the443DiTs, this conforms with the expectations of the models. Intuitively one might expect the later timesteps (less noise) to have444much higher segmentation performance as less information about the original image is corrupted. However, we found that the445middle diffusion timesteps best (See ??). Throughout the rest of our experiments we use timestep 500 out of 1000 following446this result.447

CVPR #38

> 426 427

435

436

437



Figure 4. (Left) Later MMATTN layers encode richer features for zero-shot segmentation. We investigated the impact of using features from various MMATTN layers and found that deeper layers lead to better performance on segmentation metrics like pixelwise accuracy, mIoU, and mAP. We also found that combining the information from all layers further improves performance. (Right) Optimal segmentation performance requires some noise to be present in the image. We evaluated the performance of CONCEPTATTENTION by encoding samples from a variety of timesteps (determines the amount of noise). Interestingly, we found that the optimal amount of noise was not zero, but in the middle to later stages of the noise schedule.

| Method               | Acc↑  | mIoU↑ |
|----------------------|-------|-------|
| TextSpan             | 73.84 | 38.10 |
| DAAM                 | 62.89 | 10.97 |
| Flux Cross Attention | 79.52 | 27.04 |
| CONCEPTATTENTION     | 86.99 | 51.39 |

Table 3. **CONCEPTATTENTION outperforms alternative methods on images with multiple classes from PascalVOC.** Notably, the margin between CONCEPTATTENTION and other methods is even higher for this task than when a single class is in each image.

| Space  | Softmax      | Acc↑  | mIoU↑ | mAP↑  |
|--------|--------------|-------|-------|-------|
| CA     |              | 66.59 | 49.91 | 73.17 |
| CA     | $\checkmark$ | 74.92 | 59.90 | 87.23 |
| Value  |              | 45.93 | 29.81 | 65.79 |
| Value  | $\checkmark$ | 45.78 | 29.68 | 39.61 |
| Output |              | 78.75 | 64.95 | 88.39 |
| Output | $\checkmark$ | 83.07 | 71.04 | 90.45 |

Table 4. The output space of DiT attention layers produces more transferable representations than cross attentions. We explore the transferability of several representation spaces of a DiT: the cross attentions (CA), the value space, and the output space. We performed linear projections of the image patches and concept vectors in each of these spaces and evaluated their performance on ImageNet-Segmentation.

448 Concept Attention Operation Ablations We compared the performance on the ImageNet Segmentation benchmark of 449 performing (a) just cross attention from the image patches to the concept vectors, (b) just self attention, (c) no attention 450 operations, and (d) both cross and self attention. Our results seen in Table 5 indicate that using a combination of both cross 451 and self attention achieves the best performance.

We also investigated the impact of applying a pixelwise softmax operation over our predicted segmentation coefficients. We found that it slightly improves segmentation performance in the attention output space and significantly improves performance for the cross attention maps (see Table 4.

## 455 B. Additional Qualitative Results

456 Here we show a variety of qualitative results for our method that we could not fit into the original paper.

#### CVPR 2025 Submission #38. CONFIDENTIAL REVIEW COPY. DO NOT DISTRIBUTE.

| CA           | SA           | Acc↑  | mIoU↑ | mAP↑  |  |
|--------------|--------------|-------|-------|-------|--|
|              |              | 52.63 | 35.72 | 70.21 |  |
|              | $\checkmark$ | 51.68 | 34.85 | 69.36 |  |
| $\checkmark$ |              | 76.51 | 61.96 | 86.73 |  |
| $\checkmark$ | $\checkmark$ | 83.07 | 71.04 | 90.45 |  |

Table 5. CONCEPTATTENTION performs best when we utilize both cross and self attention. We tested the effectiveness of performing just a cross attention operation between the concepts and image tokens, just a self attention among the concepts, both cross and self attention, and neither. We found that doing both operations leads to the best results. Metrics are computed on the ImageNet Segmentation benchmark.



Figure 5. A qualitative comparison between our method and several others.

| C. Pseudo-code for the CONCEPTATTENTION Algorithm   | 457        |
|---|------------|
| We show pseudo-code depicting the difference between a vanilla multi-modal attention mechanism and a multi-modal atten-<br>tion mechanism with concept attention added to it. | 458<br>459 |
| D. Concept Attention on Video Generation Models   | 460        |



Figure 6. A qualitative comparison between our method and several others.



Figure 7. A qualitative comparison between our method and several others.



Figure 8. Qualitative comparisons between numerous baselines on ImageNet Segmentation Images. We show the soft predictions for each method and binarized segmentation predictions.



Figure 9. (a) MMATTN combines cross and self attention operations between the prompt and image tokens. (b) Our CONCEPTATTENTION allows the concept tokens to incorporate information from other concept tokens and the image tokens, but not the other way around.



Figure 10. CONCEPTATTENTION produces higher fidelity raw scores and saliency maps than alternative methods, sometimes surpassing in quality even the ground truth saliency map provided by the ImageNet-Segmentation task. Top row shows the soft predictions of each method and the bottom shows the binarized predictions.

| (a) Multi-Modal Attention (b  | b) Multi-modal Attention with Concept Attention  |
|---|--|
| <pre>def multi_modal_attn(img, txt):</pre>  | <pre>attn_with_concept_attn(img, txt, concepts):</pre>   |
| img_k, img_q, img_v = img_projection(img)     img_k, img_q       txt k, txt q, txt v = txt projection(txt)     txt k, txt q | $ \begin{array}{llllllllllllllllllllllllllllllllllll$  |
| + concept_k, c  | concept_q, concept_v = txt_projection(concepts)  |
| # Concat the image and text keys, queries, and vals # Concat the  | e image and text keys, queries, and vals   |
| <pre>img_txt_k = concat([img_k, txt_k])</pre>   | concat([img_k, txt_k])   |
| <pre>img_txt_q = concat([img_q, txt_q])</pre>   | <pre>concat([img_q, txt_q])</pre>  |
| <pre>img_txt_v = concat([img_v, txt_v])</pre>   | <pre>concat([img_v, txt_v])</pre>  |
| # Perform self attention on combined sequence # Perform se  | elf attention on combined sequence   |
| attn_out = self_attention(img_txt_k, img_txt_q, img_txt_v) attn_out = s   | <pre>self_attention(img_txt_k, img_txt_q, img_txt_v)</pre>   |
| # Unpack the attention outputs # Unpack the   | e attention outputs  |
| <pre>img = attn_out[:img.shape[0]], attn_out[img.shape[0]:] img, txt = a</pre>  | attn_out[:img.shape[0]], attn_out[img.shape[0]:]   |
| + # Concatenat  | e the image and concept keys and values  |
| + img_concept_  | _K = concat([img_K, concept_K])  |
| + Img_concept_  | _v = concat([img_v, concept_v])  |
|   | man = matmul(concent a jma concent k T)  |
|   | <pre>_ mathdf(concept_q, img_concept_k.i) _ map = coftmax(concept_attn map _ axic=1) * ccale</pre> |
| + concepts = m  | natmul(concept_attn_mapimg_concept_v)  |
| + concepts - m  | acina (concept_actn_map, _img_concept_v)   |
| return img, txt + return img,   | txt, concepts  |

Figure 11. **Pseudo-code depicting the (a) multi-modal attention operation used by Flux DiTs and (b) our CONCEPTATTENTION operation.** We leverage the parameters of a multi-modal attention layer to construct a set of contextualized concept embeddings. The concepts query the image tokens (cross-attention) and other concept tokens (self-attention) in an attention operation. The updated concept embeddings are returned in addition to the image and text embeddings.



Figure 12. **CONCEPTATTENTION generalizes seamlessly to video generation MMDiT models like CogVideoX.** We apply CON-CEPTATTENTION to a CogVideoX [29] video generation model. We take several frames from the video and compare the saliency maps generated by CONCEPTATTENTION to the model's internal cross attention maps.



Figure 13. **CONCEPTATTENTION generalizes seamlessly to video generation MMDiT models like CogVideoX.** We apply CON-CEPTATTENTION to a CogVideoX [29] video generation model. We take several frames from the video and compare the saliency maps generated by CONCEPTATTENTION to the model's internal cross attention maps.