

540 **Reproducibility Statement.** We provide experimental details regarding the datasets, baselines,
541 evaluation metrics, and implementation in Sec. 4. Additional details on the dataset are provided
542 in Appendix B. In Appendix C, we provide additional details about the baselines, the implementation
543 of our method, and the computation resources used.

544 REFERENCES

545
546
547 Shahira Abousamra, Minh Hoai, Dimitris Samaras, and Chao Chen. Localization in the crowd
548 with topological constraints. In *Proceedings of the AAAI Conference on Artificial Intelligence*,
549 volume 35, pp. 872–881, 2021.

550
551 Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Ale-
552 man, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical
553 report. *arXiv preprint arXiv:2303.08774*, 2023.

554
555 Omri Avrahami, Dani Lischinski, and Ohad Fried. Blended diffusion for text-driven editing of
556 natural images. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern
Recognition*, pp. 18208–18218, 2022.

557
558 Samik Banerjee, Lucas Magee, Dingkan Wang, Xu Li, Bing-Xing Huo, Jaikishan Jayakumar,
559 Katherine Matho, Meng-Kuan Lin, Keerthi Ram, Mohanasankar Sivaprakasam, et al. Sema-
560 ntic segmentation of microscopic neuroanatomical data by combining topological priors with
561 encoder–decoder deep networks. *Nature machine intelligence*, 2(10):585–594, 2020.

562
563 Omer Bar-Tal, Lior Yariv, Yaron Lipman, and Tali Dekel. Multidiffusion: Fusing diffusion paths for
controlled image generation. 2023.

564
565 Dina Bashkirova, José Lezama, Kihyuk Sohn, Kate Saenko, and Irfan Essa. Masksketch: Unpaired
566 structure-guided masked image generation. In *Proceedings of the IEEE/CVF Conference on Com-
puter Vision and Pattern Recognition*, pp. 1879–1889, 2023.

567
568 Rickard Brüel-Gabrielsson, Bradley J Nelson, Anjan Dwaraknath, Primož Skraba, Leonidas J
569 Guibas, and Gunnar Carlsson. A topology layer for machine learning. *arXiv preprint
arXiv:1905.12200*, 2019.

570
571 Nick Byrne, James R Clough, Israel Valverde, Giovanni Montana, and Andrew P King. A persistent
572 homology-based topological loss for cnn-based multiclass segmentation of cmr. *IEEE transac-
573 tions on medical imaging*, 42(1):3–14, 2022.

574
575 Holger Caesar, Jasper Uijlings, and Vittorio Ferrari. Coco-stuff: Thing and stuff classes in context.
576 In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 1209–
577 1218, 2018.

578
579 Gunnar Carlsson. Topology and data. *Bulletin of the American Mathematical Society*, 46(2):255–
308, 2009.

580
581 Chao Chen, Xiuyan Ni, Qinxun Bai, and Yusu Wang. A topological regularizer for classifiers via
582 persistent homology. In *The 22nd International Conference on Artificial Intelligence and Statis-
583 tics*, pp. 2573–2582. PMLR, 2019.

584
585 Minghao Chen, Iro Laina, and Andrea Vedaldi. Training-free layout control with cross-attention
586 guidance. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer
Vision*, pp. 5343–5353, 2024.

587
588 James R Clough, Ilkay Oksuz, Nicholas Byrne, Julia A Schnabel, and Andrew P King. Explicit
589 topological priors for deep-learning based image segmentation using persistent homology. In
590 *International Conference on Information Processing in Medical Imaging*, pp. 16–28. Springer,
591 2019.

592
593 James R Clough, Nicholas Byrne, Ilkay Oksuz, Veronika A Zimmer, Julia A Schnabel, and An-
drew P King. A topological loss function for deep-learning based image segmentation using
persistent homology. *TPAMI*, 2020.

- 594 Tamal K Dey, Jiayuan Wang, and Yusu Wang. Road network reconstruction from satellite im-
595 ages with machine learning supported by topological methods. In *Proceedings of the 27th ACM*
596 *SIGSPATIAL International Conference on Advances in Geographic Information Systems*, pp.
597 520–523, 2019.
- 598 Prafulla Dhariwal and Alexander Nichol. Diffusion models beat gans on image synthesis. *Advances*
599 *in neural information processing systems*, 34:8780–8794, 2021.
- 600
- 601 Shiyi Du, Qicheng Lao, Qingbo Kang, Yiyue Li, Zekun Jiang, Yanfeng Zhao, and Kang Li. Distill-
602 ing knowledge from topological representations for pathological complete response prediction.
603 In *International Conference on Medical Image Computing and Computer-Assisted Intervention*,
604 pp. 56–65. Springer, 2022.
- 605 Edelsbrunner, Letscher, and Zomorodian. Topological persistence and simplification. *Discrete &*
606 *Computational Geometry*, 28:511–533, 2002.
- 607
- 608 Herbert Edelsbrunner and John Harer. *Computational Topology: An Introduction*. American Math-
609 ematical Soc., 2010.
- 610 Azade Farshad, Yousef Yeganeh, Yu Chi, Chengzhi Shen, Böjrn Ommer, and Nassir Navab. Scene-
611 genie: Scene graph guided diffusion models for image synthesis. In *Proceedings of the IEEE/CVF*
612 *International Conference on Computer Vision (ICCV) Workshops*, pp. 88–98, 2023.
- 613
- 614 J Funke, S Saalfeld, DD Bock, SC Turaga, and E Perlman. Miccai challenge on circuit reconstruction
615 from electron microscopy images, 2016.
- 616 Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair,
617 Aaron Courville, and Yoshua Bengio. Generative adversarial nets. *Advances in neural information*
618 *processing systems*, 27, 2014.
- 619
- 620 Saumya Gupta, Xiaoling Hu, James Kaan, Michael Jin, Mutshipay Mpoy, Katherine Chung, Gagan-
621 deep Singh, Mary Saltz, Tahsin Kurc, Joel Saltz, et al. Learning topological interactions for
622 multi-class medical image segmentation. In *ECCV*, 2022.
- 623 Saumya Gupta, Yikai Zhang, Xiaoling Hu, Prateek Prasanna, and Chao Chen. Topology-aware
624 uncertainty for image segmentation. *Advances in Neural Information Processing Systems*, 36,
625 2024.
- 626
- 627 Tiankai Hang, Shuyang Gu, Chen Li, Jianmin Bao, Dong Chen, Han Hu, Xin Geng, and
628 Baining Guo. Efficient diffusion training via min-snr weighting strategy. *arXiv preprint*
629 *arXiv:2303.09556*, 2023.
- 630 Hongliang He, Jun Wang, Pengxu Wei, Fan Xu, Xiangyang Ji, Chang Liu, and Jie Chen. Toposeg:
631 Topology-aware nuclear instance segmentation. In *Proceedings of the IEEE/CVF International*
632 *Conference on Computer Vision*, pp. 21307–21316, 2023.
- 633
- 634 Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter.
635 Gans trained by a two time-scale update rule converge to a local nash equilibrium. *Advances in*
636 *neural information processing systems*, 30, 2017.
- 637 Jonathan Ho and Tim Salimans. Classifier-free diffusion guidance. *arXiv preprint*
638 *arXiv:2207.12598*, 2022.
- 639
- 640 Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in*
641 *neural information processing systems*, 33:6840–6851, 2020.
- 642 Christoph Hofer, Roland Kwitt, Marc Niethammer, and Andreas Uhl. Deep learning with topologi-
643 cal signatures. *Advances in neural information processing systems*, 30, 2017.
- 644 Xiaoling Hu, Fuxin Li, Dimitris Samaras, and Chao Chen. Topology-preserving deep image seg-
645 mentation. In *NeurIPS*, 2019.
- 646
- 647 Xiaoling Hu, Yusu Wang, Li Fuxin, Dimitris Samaras, and Chao Chen. Topology-aware segmenta-
tion using discrete morse theory. In *ICLR*, 2021.

- 648 Xiaoling Hu, Dimitris Samaras, and Chao Chen. Learning probabilistic topological representations
649 using discrete morse theory. In *ICLR*, 2023.
- 650
- 651 Lianghua Huang, Di Chen, Yu Liu, Yujun Shen, Deli Zhao, and Jingren Zhou. Composer: Creative
652 and controllable image synthesis with composable conditions. *arXiv preprint arXiv:2302.09778*,
653 2023.
- 654 Phillip Isola, Jun-Yan Zhu, Tinghui Zhou, and Alexei A Efros. Image-to-image translation with
655 conditional adversarial networks. In *Proceedings of the IEEE conference on computer vision and
656 pattern recognition*, pp. 1125–1134, 2017.
- 657 Shizuo Kaji, Takeki Sudo, and Kazushi Ahara. Cubical ripser: Software for computing persistent
658 homology of image and volume data. *arXiv preprint arXiv:2005.12692*, 2020.
- 659
- 660 Valentin Khruikov and Ivan Oseledets. Geometry score: A method for comparing generative ad-
661 versarial networks. In *International conference on machine learning*, pp. 2621–2629. PMLR,
662 2018.
- 663 Gwanghyun Kim, Taesung Kwon, and Jong Chul Ye. Diffusionclip: Text-guided diffusion models
664 for robust image manipulation. In *Proceedings of the IEEE/CVF Conference on Computer Vision
665 and Pattern Recognition*, pp. 2426–2435, 2022.
- 666
- 667 Yuheng Li, Haotian Liu, Qingyang Wu, Fangzhou Mu, Jianwei Yang, Jianfeng Gao, Chunyuan Li,
668 and Yong Jae Lee. Gligen: Open-set grounded text-to-image generation. In *Proceedings of the
669 IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 22511–22521, 2023.
- 670 Maneesh Agrawala Lvmin Zhang. Controlnet-seg model. ControlNet-Seg model card, <https://huggingface.co/lllyasviel/sd-controlnet-seg>. 2023.
- 671
- 672 Inc Midjourney. Midjourney. <https://www.midjourney.com/>. 2024.
- 673
- 674 Chong Mou, Xintao Wang, Liangbin Xie, Yanze Wu, Jian Zhang, Zhongang Qi, Ying Shan, and
675 Xiaohu Qie. T2i-adapter: Learning adapters to dig out more controllable ability for text-to-image
676 diffusion models. *arXiv preprint arXiv:2302.08453*, 2023.
- 677
- 678 James R Munkres. *Elements of algebraic topology*. CRC press, 2018.
- 679
- 680 Xiuyan Ni, Novi Quadrianto, Yusu Wang, and Chao Chen. Composing tree graphical models with
681 persistent homology features for clustering mixed-type data. In *International Conference on
682 Machine Learning*, pp. 2622–2631. PMLR, 2017.
- 683 Alex Nichol, Prafulla Dhariwal, Aditya Ramesh, Pranav Shyam, Pamela Mishkin, Bob McGrew,
684 Ilya Sutskever, and Mark Chen. Glide: Towards photorealistic image generation and editing with
685 text-guided diffusion models. *arXiv preprint arXiv:2112.10741*, 2021.
- 686 Alexander Quinn Nichol and Prafulla Dhariwal. Improved denoising diffusion probabilistic models.
687 In *International Conference on Machine Learning*, pp. 8162–8171. PMLR, 2021.
- 688
- 689 OpenAI. Dall-e 3. <https://openai.com/product/dall-e-3>, a. 2023.
- 690
- 691 OpenAI. Lsun bedroom model. [https://openaipublic.blob.core.windows.net/
692 diffusion/march-2021/lsun_uncond_100M_2400K_bs64.pt](https://openaipublic.blob.core.windows.net/diffusion/march-2021/lsun_uncond_100M_2400K_bs64.pt), b. 2021.
- 693
- 694 Roni Paiss, Ariel Ephrat, Omer Tov, Shiran Zada, Inbar Mosseri, Michal Irani, and Tali Dekel.
695 Teaching clip to count to ten. In *Proceedings of the IEEE/CVF International Conference on
696 Computer Vision*, pp. 3170–3180, 2023.
- 697
- 698 Yaopeng Peng, Hongxiao Wang, Milan Sonka, and Danny Z Chen. Phg-net: Persistent homology
699 guided medical image classification. In *Proceedings of the IEEE/CVF Winter Conference on
700 Applications of Computer Vision*, pp. 7583–7592, 2024.
- 701
- 702 Quynh Phung, Songwei Ge, and Jia-Bin Huang. Grounded text-to-image synthesis with attention
703 refocusing. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recog-
704 nition*, pp. 7932–7942, 2024.

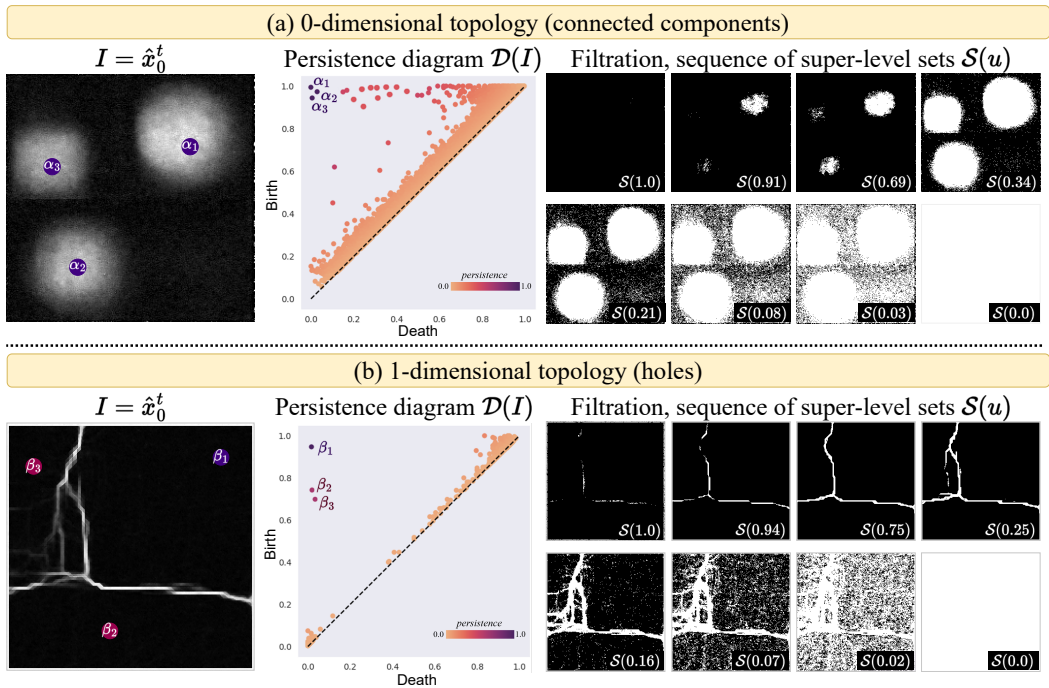
- 702 Can Qin, Shu Zhang, Ning Yu, Yihao Feng, Xinyi Yang, Yingbo Zhou, Huan Wang, Juan Car-
703 los Niebles, Caiming Xiong, Silvio Savarese, et al. Unicontrol: A unified diffusion model for
704 controllable visual generation in the wild. *arXiv preprint arXiv:2305.11147*, 2023.
- 705 Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal,
706 Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual
707 models from natural language supervision. In *International conference on machine learning*, pp.
708 8748–8763. PMLR, 2021.
- 709 Aditya Ramesh, Mikhail Pavlov, Gabriel Goh, Scott Gray, Chelsea Voss, Alec Radford, Mark Chen,
710 and Ilya Sutskever. Zero-shot text-to-image generation. In *International Conference on Machine*
711 *Learning*, pp. 8821–8831. PMLR, 2021.
- 712 Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-
713 resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF confer-*
714 *ence on computer vision and pattern recognition*, pp. 10684–10695, 2022.
- 715 Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomed-
716 ical image segmentation. In *Medical Image Computing and Computer-Assisted Intervention-*
717 *MICCAI 2015: 18th International Conference, Munich, Germany, October 5-9, 2015, Proceed-*
718 *ings, Part III 18*, pp. 234–241. Springer, 2015.
- 719 Nataniel Ruiz, Yuanzhen Li, Varun Jampani, Yael Pritch, Michael Rubinstein, and Kfir Aberman.
720 Dreambooth: Fine tuning text-to-image diffusion models for subject-driven generation. In *Pro-*
721 *ceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 22500–
722 22510, 2023.
- 723 Runway. Stability ai. Stable diffusion v1.5 model card, [https://huggingface.co/](https://huggingface.co/runwayml/stable-diffusion-v1-5)
724 [runwayml/stable-diffusion-v1-5](https://huggingface.co/runwayml/stable-diffusion-v1-5). 2022.
- 725 Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily Denton, Seyed Kam-
726 yar Seyed Ghasemipour, Burcu Karagol Ayan, S Sara Mahdavi, Rapha Gontijo Lopes, et al. Pho-
727 torealistic text-to-image diffusion models with deep language understanding. *URL https://arxiv.*
728 *org/abs/2205.11487*, 4, 2022a.
- 729 Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily L Denton, Kamyar
730 Ghasemipour, Raphael Gontijo Lopes, Burcu Karagol Ayan, Tim Salimans, et al. Photorealistic
731 text-to-image diffusion models with deep language understanding. *Advances in neural informa-*
732 *tion processing systems*, 35:36479–36494, 2022b.
- 733 Axel Sauer, Dominik Lorenz, Andreas Blattmann, and Robin Rombach. Adversarial diffusion dis-
734 tillation. In *European Conference on Computer Vision*, pp. 87–103. Springer, 2025.
- 735 Suprosanna Shit, Johannes C Paetzold, Anjany Sekuboyina, Ivan Ezhov, Alexander Unger, Andrey
736 Zhylka, Josien PW Pluim, Ulrich Bauer, and Bjoern H Menze. cldice-a novel topology-preserving
737 loss function for tubular structure segmentation. In *CVPR*, 2021.
- 738 Jascha Sohl-Dickstein, Eric Weiss, Niru Maheswaranathan, and Surya Ganguli. Deep unsupervised
739 learning using nonequilibrium thermodynamics. In *International conference on machine learn-*
740 *ing*, pp. 2256–2265. PMLR, 2015.
- 741 Kihyuk Sohn, Honglak Lee, and Xinchen Yan. Learning structured output representation using deep
742 conditional generative models. *Advances in neural information processing systems*, 28, 2015.
- 743 Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. *arXiv*
744 *preprint arXiv:2010.02502*, 2020a.
- 745 Yang Song and Stefano Ermon. Generative modeling by estimating gradients of the data distribution.
746 *Advances in neural information processing systems*, 32, 2019.
- 747 Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon, and Ben
748 Poole. Score-based generative modeling through stochastic differential equations. *arXiv preprint*
749 *arXiv:2011.13456*, 2020b.

756 Nico Stucki, Johannes C Paetzold, Suprosanna Shit, Bjoern Menze, and Ulrich Bauer. Topologically
757 faithful image segmentation via induced matching of persistence barcodes. In *ICML*, 2023.
758
759 Student. The probable error of a mean. *Biometrika*, 6(1):1–25, 1908.
760
761 Robert Endre Tarjan. Efficiency of a good but not linear set union algorithm. *Journal of the ACM*
762 (*JACM*), 22(2):215–225, 1975.
763
764 Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez,
765 Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural informa-*
766 *tion processing systems*, 30, 2017.
767
768 Fan Wang, Huidong Liu, Dimitris Samaras, and Chao Chen. Topogan: A topology-aware generative
769 adversarial network. In *Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK,*
770 *August 23–28, 2020, Proceedings, Part III 16*, pp. 118–136. Springer, 2020.
771
772 Fan Wang, Saarthak Kapse, Steven Liu, Prateek Prasanna, and Chao Chen. Topotxr: a topological
773 biomarker for predicting treatment response in breast cancer. In *International Conference on*
774 *Information Processing in Medical Imaging*, pp. 386–397. Springer, 2021.
775
776 Haotian Wang, Min Xian, and Aleksandar Vakanski. Ta-net: Topology-aware network for gland
777 segmentation. In *WACV*, 2022a.
778
779 Tengfei Wang, Ting Zhang, Bo Zhang, Hao Ouyang, Dong Chen, Qifeng Chen, and Fang Wen.
780 Pretraining is all you need for image-to-image translation. *arXiv preprint arXiv:2205.12952*,
781 2022b.
782
783 Lvmin Zhang, Anyi Rao, and Maneesh Agrawala. Adding conditional control to text-to-image
784 diffusion models. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*,
785 pp. 3836–3847, 2023.
786
787 Shihao Zhao, Dongdong Chen, Yen-Chun Chen, Jianmin Bao, Shaozhe Hao, Lu Yuan, and Kwan-
788 Yee K Wong. Uni-controlnet: All-in-one control to text-to-image diffusion models. *Advances in*
789 *Neural Information Processing Systems*, 36, 2024.
790
791
792
793
794
795
796
797
798
799
800
801
802
803
804
805
806
807
808
809

810 The appendix is organized as follows.
 811 Appendix A provides the illustration of super-level sets for 1-dim topology.
 812 Appendix B provides additional details of the datasets.
 813 Appendix C provides additional baseline and implementation details.
 814 Appendix D contains 0-dim constraint-wise results on the COCO dataset.
 815 Appendix E provides qualitative results of Stable Diffusion (Rombach et al., 2022; Runway),
 816 DALL-E 3 (OpenAI, a), and AR (Phung et al., 2024) for 1-dim topological constraints.
 817 Appendix F presents experiments on 1-dim topology where there are non-boundary (standalone)
 818 holes.
 819 Appendix G includes an ablation study on the ‘Encoding Network’ for the topological constraint c .
 820 Appendix H provides a discussion on the limitations of our method.

825 **A PERSISTENT HOMOLOGY**

826 From Sec. 3.2, the equivalent of Fig. 4 for 1-dimensional topology is shown in Fig. 9.



852 Figure 9: Illustration of persistent homology and persistence diagrams of both types of topological
 853 structures, 0-dim connected components and 1-dim holes. (a) Despite the noise, we can visually see
 854 three prominent structures $\alpha_1, \alpha_2, \alpha_3$ in I . In the topological space, $\alpha_1, \alpha_2, \alpha_3$ thus appear in
 855 the top-left corner of the persistence diagram $\mathcal{D}(I)$, persisting through most of the filtration \mathcal{S} . Similarly
 856 in (b), $\beta_1, \beta_2, \beta_3$ denote the prominent holes. All the remaining connected components and holes
 857 are noisy, persisting over a short threshold in \mathcal{S} , thus appearing closer to the diagonal in $\mathcal{D}(I)$.
 858 Persistence diagrams are useful to distinguish between salient and noisy structures in an image.

860 **B DATASET DETAILS**

861 We provide the number of images per topological constraint c used for training on each dataset
 862 in Tab. 5. For COCO (Caesar et al., 2018), since we also consider the animal class, each animal is

Table 5: Dataset composition

Dataset	Topological Constraint (Betti Number)									
	1	2	3	4	5	6	7	8	9	10
Shapes	2K	2K	2K	2K	2K	2K	2K	2K	2K	2K
COCO (Animals)	520	517	503	297	176	85	64	38	30	27
Google Maps	549	669	1099	1220	1343	1806	602	1470	1054	662
CREMI	2160	1992	3726	3505	1644	580	187	207	170	112

distributed unequally across the different constraint values. For example, there were more images for ‘birds’ having $c = 10$ than compared to, say, ‘elephant.’ All the 10 animal classes are present in the dataset; they are bear, bird, cat, cow, dog, elephant, giraffe, horse, sheep, and zebra.

When curating the dataset for CREMI (Funke et al., 2016) and Google Maps (Isola et al., 2017), we manually added a (white) border to all the images. By definition of 1-dim topology, a hole is completely surrounded by a boundary. Hence, we needed to add a border to obtain the correct number of holes/regions.

C IMPLEMENTATION DETAILS

All ADM-T and TDN experiments were conducted on 1 NVIDIA RTX A6000 GPU, with a batch size of 16 and a learning rate of 2×10^{-5} . As mentioned in Sec. 3.3, our diffusion model is parameterized to predict in noise space, and we use Eq. (2) to get an estimate of the noiseless image. Although diffusion models can be parameterized to predict the noiseless state directly, we find from existing works (Hang et al., 2023; Wang et al., 2022b) that their performance is poorer compared to predicting the noise. Hence we stick to the configuration of predicting the noise. This also allows us to load pretrained weights from OpenAI (b) for our experiments instead of training from scratch.

For training ADM-T and TDN, we use the PyTorch codebase from Dhariwal & Nichol (2021)⁴ and use the LSUN Bedrooms pretrained model checkpoint (OpenAI, b) to fine-tune from. To compute the birth death pairs of each topological structure, we use the Cubical Ripser (Kaji et al., 2020) library. We will publicly release the code upon acceptance of the paper.

In Fig. 1, Fig. 6, Fig. 7, Fig. 11, and Fig. 12, the Stable Diffusion (Rombach et al., 2022) results are generated using the Diffusers⁵ library with pretrained checkpoint from Runway. For DALL-E 3, we generate images using the OpenAI API⁶. For Attention Refocusing (AR) (Phung et al., 2024), we use their publicly available codebase⁷ alongwith GPT-4 (Achiam et al., 2023) from the OpenAI API to generate the layout maps. For rendering images from masks via ControlNet (Zhang et al., 2023), we use the Diffusers library with pretrained checkpoint from Lvmin Zhang. For Fig. 7, however, we fine-tune ControlNet on the CREMI dataset (Funke et al., 2016) so as to generate appropriate results for the corresponding text prompt.

D CONSTRAINT-WISE RESULTS

In Fig. 10, we plot the accuracy of the different methods on the COCO dataset. For smaller object counts, that is, $c \leq 5$, the performance of each method is better than $c > 5$. At higher object counts, although the accuracy of each method reduces, TDN still significantly outperforms the baselines.

E ADDITIONAL 1-DIM QUALITATIVE RESULTS

In Fig. 11, we show qualitative results of pretrained Stable Diffusion (Rombach et al., 2022), DALL-E 3 (OpenAI, a) and Attention Refocusing (AR) (Phung et al., 2024) for the same 1-dim constraints shown in Fig. 7 of the main paper. As 1-dim topology is hard to describe in words, we tried a lot of variations for the text prompts, and show the results from the best ones in the figure.

⁴<https://github.com/openai/guided-diffusion>

⁵<https://huggingface.co/docs/diffusers/en/index>

⁶<https://platform.openai.com/docs/api-reference/introduction>

⁷<https://github.com/Attention-Refocusing/attention-refocusing>

918
919
920
921
922
923
924
925
926
927
928
929
930
931
932
933
934
935
936
937
938
939
940
941
942
943
944
945
946
947
948
949
950
951
952
953
954
955
956
957
958
959
960
961
962
963
964
965
966
967
968
969
970
971

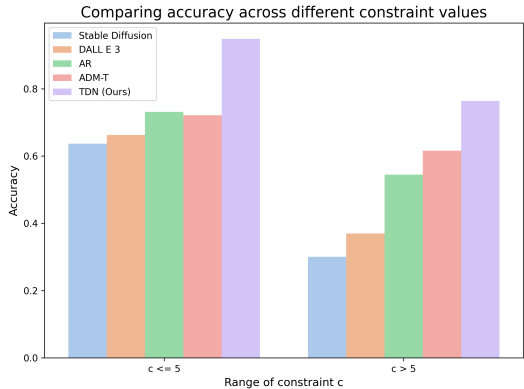


Figure 10: Accuracy results on low and high object counts (COCO dataset).

As the pretrained Stable Diffusion, DALL-E 3, and AR are not trained on Electron Microscopy images of cell neurons, their inaccurate results are understandable. For generating roads, however, these methods do generate visually appropriate images but struggle to maintain the correct number of holes/regions. AR additionally tends to generate images that appear to be divided into separate, unconnected sections. This is due to the use of layout maps in the reverse process. All these methods are limited to generating images with 0-dim topology (i.e., distinct objects), and do not extend to 1-dim topology. This shortcoming motivates our work on TopoDiffusionNet.

F ADDITIONAL EXPERIMENTS ON HOLES

In the main paper, we show experiments on 1-dim topology using the CREMI (Funke et al., 2016), and Google Maps (Isola et al., 2017) datasets. In these datasets, it makes sense that the holes are with respect to the image frame/boundary and span the whole image. However, TDN can handle 1-dim holes in general, not just those with respect to the boundary. To demonstrate this, we conduct experiments on standalone holes (not relative to the boundary), and present the results in Tab. 6 and Fig. 12. We generate a synthetic dataset of circular rings (similar to the Shapes dataset) to train TDN, and use the prompt ‘donuts’ in ControlNet to render the image. While generating c donuts could also be achieved using the 0-dim topological constraint, this experiment highlights the 1-dim generalizability of TDN.

Table 6: Standalone holes

Method	Accuracy \uparrow	F1 \uparrow
ADM-T	0.79 ± 0.12	0.81 ± 0.11
TDN (Ours)	0.95 ± 0.03	0.96 ± 0.02

G ADDITIONAL ABLATION STUDY

We conduct an additional ablation study apart from the ones presented in the main paper.

Ablation study of the encoding network. In TDN, we use the topological constraint c as a condition. We do this by first obtaining an embedding of c via an Encoding Network (Fig. 3), and then passing it to all the residual blocks in the denoising model.

In the results reported in the main paper, the Encoding Network is composed of a few linear layers. We use ‘LL’ to denote this configuration.

Another way to configure the Encoding Network is to use the Transformer sinusoidal position embedding (Vaswani et al., 2017). Let ‘PE’ denote this configuration. It uses the same code as the network generating the sinusoidal timestep embedding.

We show results comparing LL and PE in Tab. 7 when using our proposed objective function \mathcal{L}_{top} . We find that both configurations have comparable performance, with LL slightly outperforming PE.

972
973
974
975
976
977
978
979
980
981
982
983
984
985
986
987
988
989
990
991
992
993
994
995
996
997
998
999
1000
1001
1002
1003
1004
1005
1006
1007
1008
1009
1010
1011
1012
1013
1014
1015
1016
1017
1018
1019
1020
1021
1022
1023
1024
1025

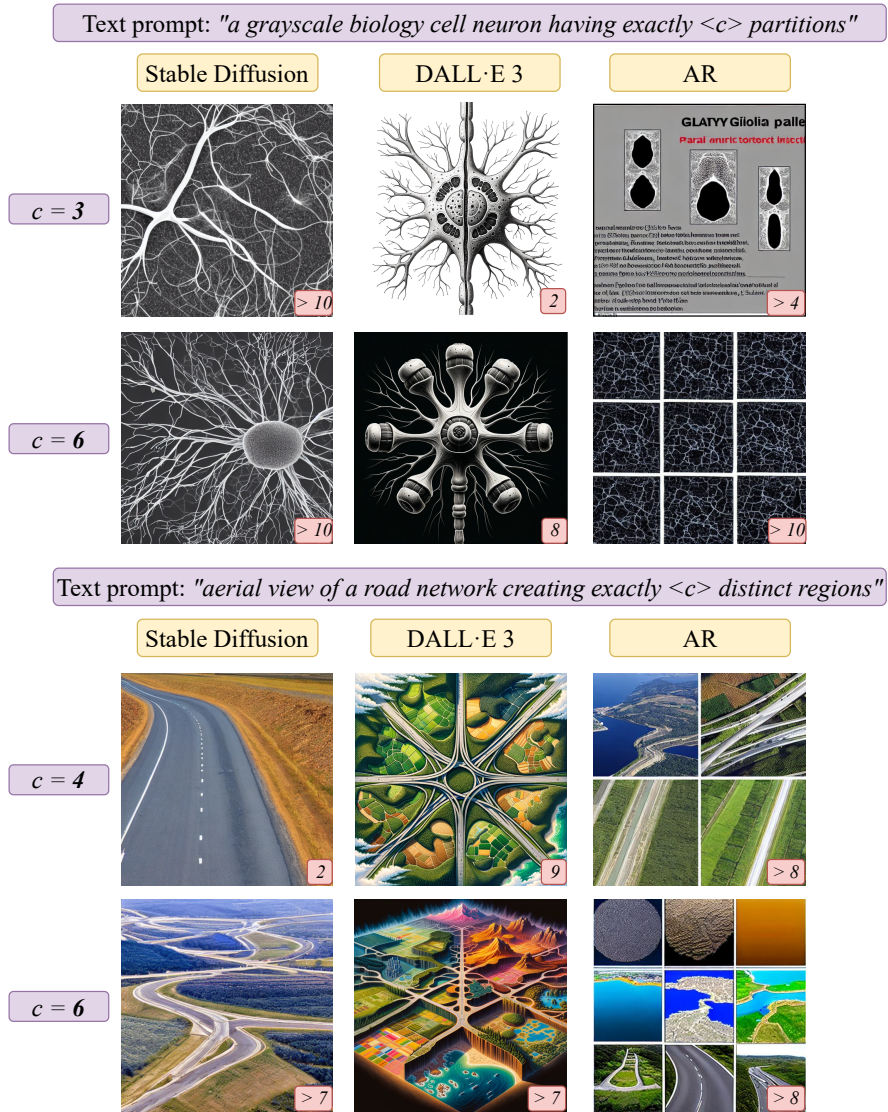


Figure 11: Qualitative results for 1-dim topological constraint. Stable Diffusion, DALL-E 3 and AR take text prompts as input (purple box). Rows 1-2: Results equivalent to CREMI. Rows 3-4: Results equivalent to Google Maps. Number of holes within each image is noted in its bottom-right inset.

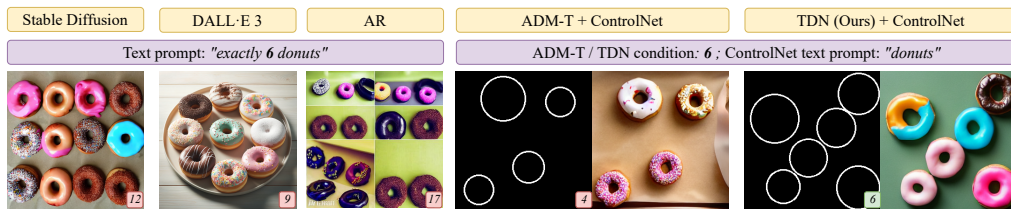


Figure 12: Results on standalone holes (holes not relative to the image frame/boundary). Number of donuts within each mask/image is noted in its bottom-right inset.

Table 7: Ablation study on Encoding Network for TDN

Dataset	Encoding	Accuracy \uparrow	Precision \uparrow	F1 \uparrow
Shapes	LL	0.9478 \pm 0.0420	0.9499 \pm 0.0492	0.9488 \pm 0.0370
	PE	0.9011 \pm 0.0730	0.9132 \pm 0.0683	0.9033 \pm 0.0385
COCO (Animals)	LL	0.8557 \pm 0.0805	0.8670 \pm 0.0636	0.8613 \pm 0.0970
	PE	0.8395 \pm 0.0952	0.8436 \pm 0.1152	0.8411 \pm 0.1014

H LIMITATIONS

Presently, TDN requires at least a few samples of each constraint in the training data and is not guaranteed to extrapolate to unseen constraints. This limitation stems from the broader challenges associated with state-of-the-art diffusion models, which require a large amount of training data. Exploring the numeric relationships between the different constraints, instead of treating them independently, has the potential to generalize to unseen constraints.

I DETAILS ON PERSISTENT HOMOLOGY COMPUTATION

We expand on Sec. 3.2 regarding the persistent homology and persistence diagram (PD) computation. As mentioned in the implementation details (Appendix C), we use the Cubical Ripser library Kaji et al. (2020) for the homology computation. Their manuscript has details about the optimized algorithms they implement for the computation. Here, we provide a summarized version of the details.

Consider a 2D image $I \in \mathbb{R}^2$. We consider the 0-dim case in which we track the birth and death of connected components (CC). As we construct a super-level filtration \mathcal{S} , we start thresholding an image I starting from the maximum value to the minimum value. For computational purposes, we do not need to consider all threshold values $u \in \mathbb{R}$, rather, since the image has size $H \times W$, the maximum number of unique values in the image is $H * W$. Hence we need to consider atmost $H * W$ values of threshold u to determine the birth and death times of all the CC in the image.

First, preprocessing is done to store all the intensity values in I along with their (x, y) location in decreasing order in a sorted data structure. This allows for faster identification of the next threshold value as well as the pixels that get newly included in the super-level set. Second, we utilize the standard Union-Find data structure Tarjan (1975), which maintains a collection of disjoint sets. This data structure supports the operation `FIND` $((x, y))$, which finds the highest-value representative of the CC containing the pixel (x, y) . It also supports the operation `UNITE` $((x, y), (v, w))$, which unites the CCs represented by root pixels (x, y) and (v, w) and, assuming $I(x, y) > I(v, w)$, making (x, y) the representative of the merged components.

At each threshold value u , we use the sorted data structure to identify which pixels are included in the set $\mathcal{S}(u)$. `FIND` is called on all these newly added pixels to determine if they belong to any existing CC; if not, they are considered roots of a new CC, with a new entry to the list of disjoint sets in the union-find data structure. Furthermore, this results in the creation of a dot $(u, -\infty)$ in the persistence diagram whose birth time is u , and a default death time of $-\infty$. Then, `UNITE` is called for every pair of roots in the disjoint set list, to determine if any of the sets are in fact merged. Consider if sets with roots (x, y) and (v, w) are merged, with $I(x, y) > I(v, w)$. In that case, CC with root (v, w) has now ‘died’, and so the persistence dot whose birth time was $I(v, w)$ in the diagram, will now have an updated death time of u . Hence, $(I(v, w), u)$ will replace the old dot $(I(v, w), -\infty)$ in the persistence diagram.

We show a walkthrough of this in Fig. 13. We start with an empty union-find data structure. In (a), we start with $u = 5$, as 5 is the maximum value in the image. We call `FIND` and update the union-find data structure with a new entry. Additionally, the persistence diagram (PD) in (e) has an entry for $(5, -\infty)$. This CC is denoted as α_1 . Next, in (b) $u = 4$, all the blue and red pixels are newly added to this filtration. `FIND` is called on each of them. The blue pixels indicate that they were found to be a part of an existing CC. The red pixels indicate that no parent CC was found, and hence are considered as creating new CCs. The PD includes two new dots $\alpha_2(4, -\infty)$ and $\alpha_3(4, -\infty)$. Next, as intensity value 3 is not present in I , we do not need to compute $\mathcal{S}(3)$. The sorted data structure directly leads us to the next u value which is $u = 2$ in (c). `FIND` is called on all the blue pixels. `UNITE` is called on all pairs of root CC nodes, resulting `True` for the fact that components α_1 and α_2 have now merged. Since the birth time of α_1 is larger, the component α_2 ‘dies’ and becomes a part of α_1 . Thus, the persistent dot is updated from $\alpha_2(4, -\infty)$ to $\alpha_2(4, 2)$ as it was born at $u = 4$ has died at $u = 2$. Similarly in (d), at $u = 1$, `FIND` is called on all the blue pixels. Then `UNITE` (α_1, α_3) is called, showing that they are actually merged. Thus, the persistent dot is updated from $\alpha_3(4, -\infty)$ to $\alpha_3(4, 1)$ as α_3 was born at $u = 4$ and has now died at $u = 1$. Thus, the final persistent dots we end up with are $\alpha_1(5, -\infty)$, $\alpha_2(4, 2)$, $\alpha_3(4, 1)$ as shown in (e).

For 1-dimensional and higher features (e.g., loops and voids), while we do not go into details here, the Cubical Ripser library constructs a sparse coboundary matrix that encodes relationships between pixels. This matrix is reduced to upper-triangular form using column operations, optimizing the identification of the birth and death of cycles. Each non-zero pivot in the reduced matrix corresponds to a topological feature. This computational framework ensures the efficient generation of persistence diagrams, leveraging union-find for 0-dim and matrix operations for higher dimensions. More details are available in the Cubical Ripser Kaji et al. (2020) paper.

1134
 1135
 1136
 1137
 1138
 1139
 1140
 1141
 1142
 1143
 1144
 1145
 1146
 1147
 1148
 1149
 1150
 1151
 1152
 1153
 1154
 1155
 1156
 1157
 1158
 1159
 1160
 1161
 1162
 1163
 1164
 1165
 1166
 1167
 1168
 1169
 1170
 1171
 1172
 1173
 1174
 1175
 1176
 1177
 1178
 1179
 1180
 1181
 1182
 1183
 1184
 1185
 1186
 1187

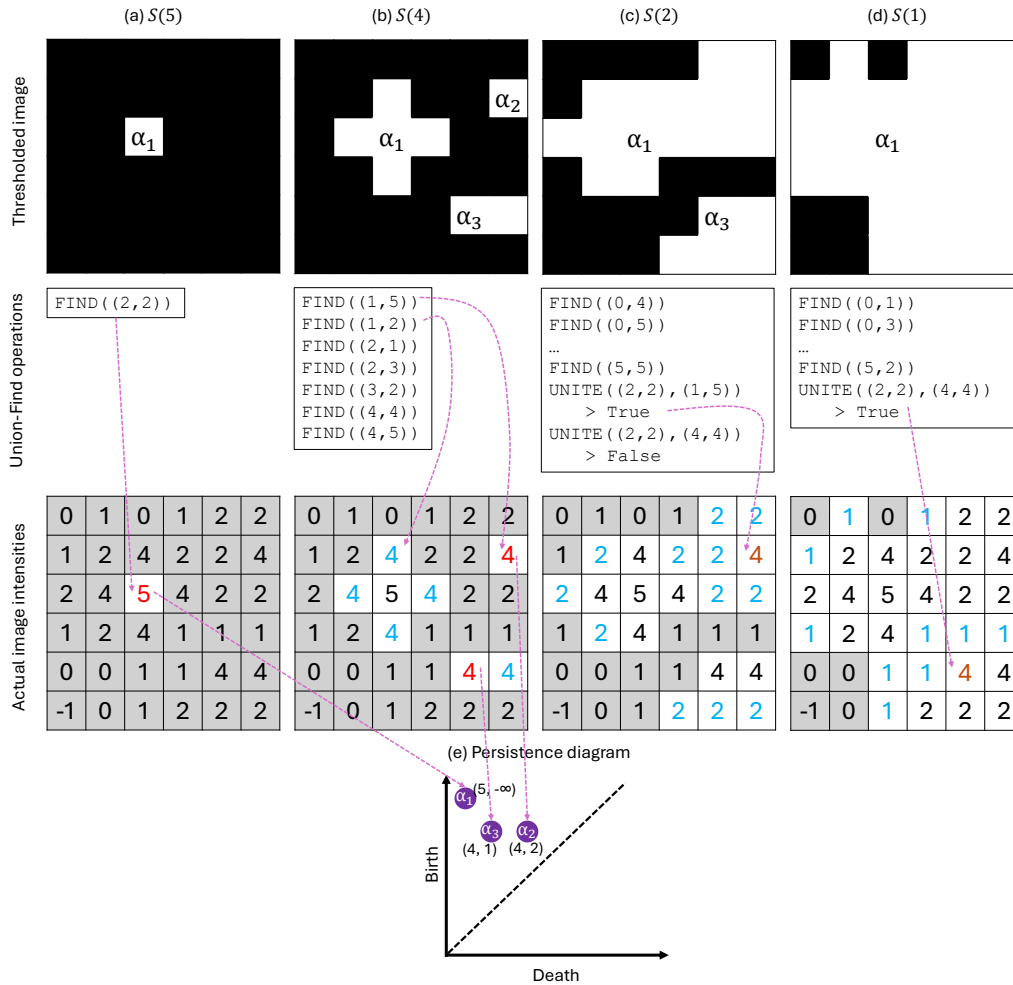


Figure 13: Figure inspired from Kaji et al. (2020). The figure shows a sample 2D image, thresholded at different values. (a)-(d) shows the sequence of super-level sets $S(u)$, which is nothing but thresholding the image at the value u . It shows how different connected components (CC) are created or destroyed as u changes. White-background cells indicate pixels included in the set $S(u)$. Red indicates root pixels causing the birth of a CC. Blue indicates pixels that are newly added to the set $S(u)$ but are absorbed in existing CCs. Brown indicates root pixels whose CC has died, having merged with another CC whose root had a higher birth time. (e) shows the corresponding 0-dim persistence diagram.

J QUALITY OF GENERATED MASKS

We provide FID scores across all the datasets in Tab. 8. Since our method TDN focuses on generating masks, we report the FID of the mask [FID (Mask)] and the FID of the images [FID (Image)] generated by ControlNet Zhang et al. (2023) when using these masks as condition.

Our analysis shows that the quality of the masks generated by TDN closely resembles the true masks, that is, the ground truth (GT) masks annotated by humans, as evidenced by the consistently low FID (Mask) scores across all datasets. Additionally, to evaluate the impact on final image quality, we compare the FID (Image) metric between two scenarios: ControlNet using GT masks (ControlNet + GT) versus ControlNet using TDN-generated masks. The results demonstrate comparable FID scores, indicating that using TDN-generated masks as conditioning does not degrade the image quality. In other words, TDN-generated masks are as good as GT masks for generating real images.

Results from Tab. 1 and Tab. 8 highlight that TDN achieves improved topological control while maintaining the overall visual quality of the generated images.

Table 8: **Comparison of FID (Mask) and FID (Image) against baselines across all datasets.** For Shapes and Google Maps datasets, we cannot report FID (Image) as there is no ground truth image dataset to compare to. For COCO and CREMI, we report FID (Image) against the respective dataset images. ControlNet + GT indicates that ControlNet is using the GT masks as the condition. ControlNet + FT indicates that it has been fine-tuned on the CREMI dataset to generate similarly textured images. **Bold** denotes the best results, while *italics* denotes the second best. Note that ControlNet + GT has been included as the upper bound performance; realistically, it cannot be used during inference as the GT masks are not available

Dataset	Method	FID (Mask) ↓	FID (Image) ↓
Shapes	ADM-T	0.092	-
	TDN (Ours)	0.068	-
COCO (Animals)	Stable Diffusion	-	29.41
	DALL-E 3	-	17.49
	AR	-	35.05
	ADM-T	0.267	21.72
	TDN (Ours)	0.222	21.28
	ControlNet + GT	-	20.94
Google Maps	ADM-T	0.198	-
	TDN (Ours)	0.156	-
CREMI	Stable Diffusion	-	48.18
	DALL-E 3	-	54.72
	AR	-	69.86
	ADM-T	0.518	3.322
	TDN (Ours)	0.467	3.286
	ControlNet + GT + FT	-	3.126

K USING 0-DIM AND 1-DIM SIMULTANEOUS TOPOLOGICAL CONSTRAINTS

In our work, we focused on using either 0-dim (number of objects) or 1-dim (number of holes) constraints individually, depending on the property of the dataset. In this section, we conduct experiments on synthetic data using both 0-dim and 1-dim constraints simultaneously. The constraint is of the form (a, b) where a denotes 0-dim while b denotes 1-dim. Both are first separately encoded and then added together to be fed as a condition. Sample masks generated by TDN along with ControlNet-rendered images are shown in Fig. 14. We report quantitative results in Tab. 9.

We find that TDN is capable of handling these joint constraints. While the performance is slightly lower compared to using each constraint individually, TDN still significantly outperforms ADM-T across all metrics. This shows that persistent homology (PH) can enhance the performance of the base diffusion model, even under multi-constraint scenarios. This finding opens up interesting possibilities for future work in handling even richer combinations of topological constraints.

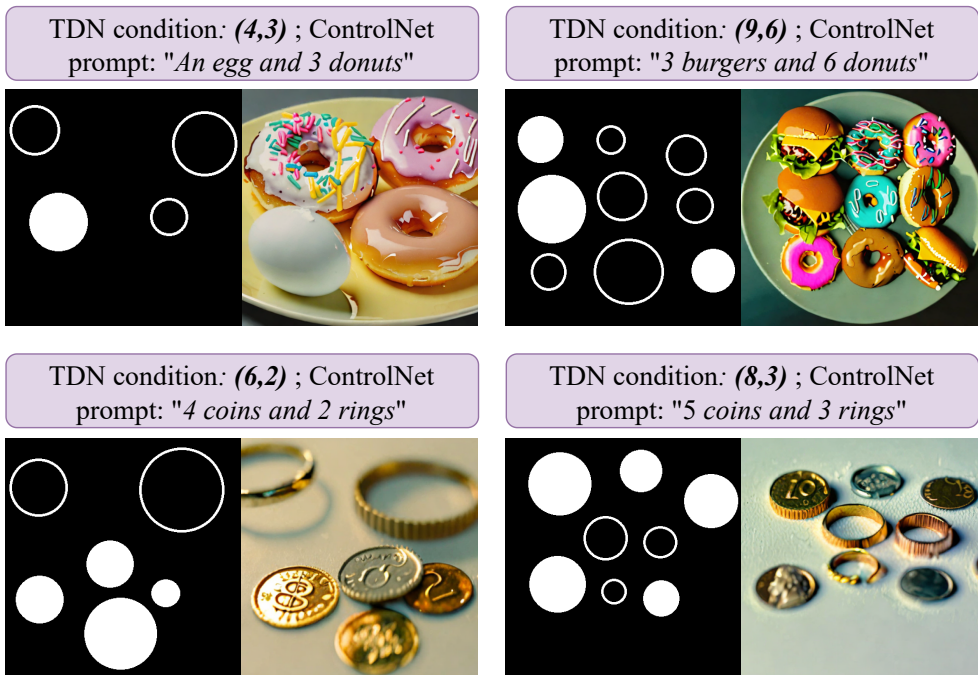


Figure 14: Masks generated by TDN using both 0-dim and 1-dim topological constraints simultaneously. Condition (a, b) denotes a for 0-dim (#objects) and b for 1-dim (#holes). Images are generated via ControlNet by using mask as a condition.

Table 9: Quantitative performance of using both 0-dim and 1-dim topological constraints simultaneously. The FID (Mask) is reported once for each method. Best results are highlighted in **bold**

TopoDim	Method	Accuracy \uparrow	Precision \uparrow	F1 \uparrow	FID (Mask) \downarrow
0-dim	ADM-T	0.7383 \pm 0.1305	0.7997 \pm 0.1268	0.7677 \pm 0.1229	0.1279
	TDN (Ours)	0.9183 \pm 0.0731	0.9338 \pm 0.0993	0.9261 \pm 0.0906	0.0982
1-dim	ADM-T	0.7616 \pm 0.0905	0.7892 \pm 0.1129	0.7752 \pm 0.1082	-
	TDN (Ours)	0.9233 \pm 0.0705	0.9492 \pm 0.0913	0.9360 \pm 0.0720	-

1296
1297
1298
1299
1300
1301
1302
1303
1304
1305
1306
1307
1308
1309
1310
1311
1312
1313
1314
1315
1316
1317
1318
1319
1320
1321
1322
1323
1324
1325
1326
1327
1328
1329
1330
1331
1332
1333
1334
1335
1336
1337
1338
1339
1340
1341
1342
1343
1344
1345
1346
1347
1348
1349

K.1 GENERATING SPECIFIC TOPOLOGY LIKE DOUBLE ANNULUS

As mentioned above, we represent the joint constraints as pairs (a, b) , where a specifies the number of connected components (0-dim) and b specifies the number of holes (1-dim). A double annulus consists of concentric circles, with $(a, b) = (2, 2)$. Our loss, based on homology, cannot distinguish whether two loops are positioned concentric or not, however, it can generate structures that are homotopy equivalent to the double annulus. As shown in Fig. 15, by sampling for the $(2, 2)$ condition, TDN can generate nested circles (visually the closest constraint to a double annulus) and other homotopy equivalent structures. This experiment demonstrates that our method can handle complex topological specifications through the combination of constraints, though geometric relationships (like equidistance and concentricity) remain an interesting direction for future work.

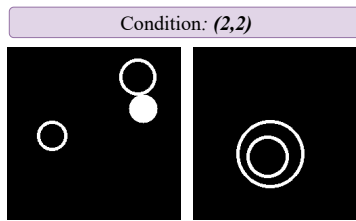


Figure 15: Sample images generated by TDN when using both 0-dim and 1-dim topological constraints simultaneously. With constraints $(a, b) = (2, 2)$, we can generate nested circles which are the closest to achieving a double annulus (two concentric circles) using our method.

L COMPARISON AGAINST CVAE

In this section, we justify the motivation behind using diffusion models instead of other generative models like Conditional Variational Autoencoders (CVAE) Sohn et al. (2015). We conduct experiments using CVAE as the generative model for the same task of generating masks. Although we trained on multiple image resolutions, CVAE performed poorly on larger image sizes. Hence we report results on 64×64 images, using a latent embedding dimension of 128. We show results in Tab. 10 and Fig. 16 for both standard CVAE and a version incorporating our proposed loss \mathcal{L}_{top} .

Our experimental results show significant limitations of the CVAE approach. From Fig. 16, we see that the quality of the masks generated by CVAE, with and without \mathcal{L}_{top} , is severely degraded. From Tab. 10, it has a significantly higher FID compared to our TDN masks, indicating a lower resemblance to the true masks. Additionally, CVAE also struggles to preserve topological constraints, as evidenced by the low performance on evaluation metrics like Accuracy and F1. While CVAE + \mathcal{L}_{top} slightly improves performance, it remains far weaker than diffusion models. Specifically, in the 0-dim, CVAE generates fragmented objects and fails to preserve their overall shapes. In 1-dim, it fails to generate connected structures. Although mask images seem easy because they are binary images, they are in fact a challenge to generate. This is because of the difference in the number of objects, their varied shapes and sizes, and unconstrained spatial locations (eg: not fixed to the center of the image). Similarly, the diverse orientation and connection patterns make the 1-dim datasets challenging. Such complexities make it difficult for the CVAE to generalize to these datasets.

Using diffusion models is critical to the success of TDN. First, the base diffusion model ADM-T effectively captures spatial arrangements and object shapes (in 0-dim), and connectivity patterns in 1-dim. However, it struggles with the number of objects/holes. This is where our persistent homology loss brings big improvements. Preserving the number of structures requires heavy global reasoning and detail. In contrast, CVAE’s bottleneck layer compresses information, losing significant information in the image space. This limits their ability to preserve object shape and enforce strict topological constraints. The significant gain of our loss in diffusion models is owing to the gradual denoising steps, as well as intact image resolution. These experiments highlight the necessity of diffusion models, making them critical to the success of our method in topological control.

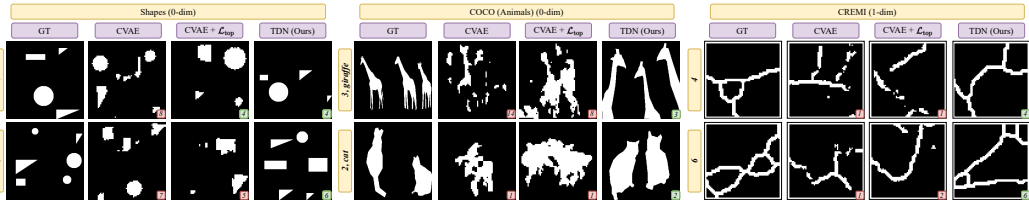


Figure 16: Qualitative comparison of using CVAE and diffusion models as the generative model. Ground truth (GT) denotes the true masks used to train the models

Table 10: Quantitative comparison of using CVAE and diffusion models as the generative model. We include ADM-T and TDN results from Tab. 1 for convenience. Best results are in **bold**

Dataset (TopoDim)	Method	Accuracy \uparrow	Precision \uparrow	F1 \uparrow	FID (Mask) \downarrow
Shapes (0-dim)	CVAE Sohn et al. (2015)	0.3133 \pm 0.1321	0.3235 \pm 0.1984	0.3078 \pm 0.1443	2.231
	CVAE + \mathcal{L}_{top}	0.3816 \pm 0.0919	0.3444 \pm 0.1276	0.3220 \pm 0.1365	1.981
	ADM-T	0.7500 \pm 0.1889	0.7809 \pm 0.1582	0.7651 \pm 0.1210	0.092
	TDN (Ours)	0.9478 \pm 0.0420	0.9499 \pm 0.0492	0.9488 \pm 0.0370	0.068
COCO (0-dim)	CVAE Sohn et al. (2015)	0.2442 \pm 0.0410	0.2696 \pm 0.0285	0.2562 \pm 0.0318	4.208
	CVAE + \mathcal{L}_{top}	0.3094 \pm 0.0967	0.3342 \pm 0.0729	0.3213 \pm 0.1165	4.083
	ADM-T	0.6685 \pm 0.1485	0.6917 \pm 0.1079	0.6799 \pm 0.1931	0.267
	TDN (Ours)	0.8557 \pm 0.0805	0.8670 \pm 0.0636	0.8613 \pm 0.0970	0.222
CREMI (1-dim)	CVAE Sohn et al. (2015)	0.2785 \pm 0.4296	0.2267 \pm 0.1156	0.2499 \pm 0.1391	5.971
	CVAE + \mathcal{L}_{top}	0.3267 \pm 0.2834	0.3091 \pm 0.1079	0.3176 \pm 0.1447	5.751
	ADM-T	0.5357 \pm 0.1879	0.4777 \pm 0.1797	0.4881 \pm 0.1571	0.518
	TDN (Ours)	0.7785 \pm 0.1901	0.8142 \pm 0.1925	0.7959 \pm 0.1659	0.467

M COMPARISON AGAINST BASIC MASKS

To justify the importance of our approach, we conduct experiments using simple shapes (circles and squares) as conditioning masks for ControlNet Zhang et al. (2023). Quantitative results are shown in Tab. 11, and qualitative examples are provided in Fig. 17.

Our experiments reveal several limitations when using simplified shapes as conditioning masks, as they fail to provide sufficient spatial and structural guidance. This setup is comparable to AR Phung et al. (2024) (shown in Tab. 1 and Fig. 6 of the main paper), which uses rectangular bounding boxes with an additional attention step. From our observations, we detect multiple failure modes. First, from Fig. 17, we see that such conditions tend to introduce visual artifacts in the generated images. The images are often fragmented, with objects isolated within assigned areas, leading to a divided and visually disjointed image. Second, the basic shapes also fail to constrain the number of objects within them, often resulting in 0 objects or multiple objects per shape (see the ‘cats’ row). Finally, for large animals like zebras, ControlNet struggles to complete the image beyond the boundary of the shape—while the zebra texture is present within the shape, the overall image is not as desired. Similarly, for even larger animals like giraffes, the discrepancy between the basic shape and the natural object proportions leads to incomplete or distorted generations.

In contrast, the masks generated by TDN provide tighter spatial and structural control by offering sufficient detail to guide the network effectively. This results in correct object counts and better visual quality, as evidenced by superior FID scores and count metrics. TDN’s masks strike a balance between oversimplification and unnecessary complexity, as generative models are powerful enough to fill in finer textural details.

Our results highlight that relying solely on a text prompt and simplified masks is insufficient for the model to generate visually coherent images. Instead, meaningful masks like those generated by TDN are essential for ensuring reliable control over cardinality and topology in generative models.

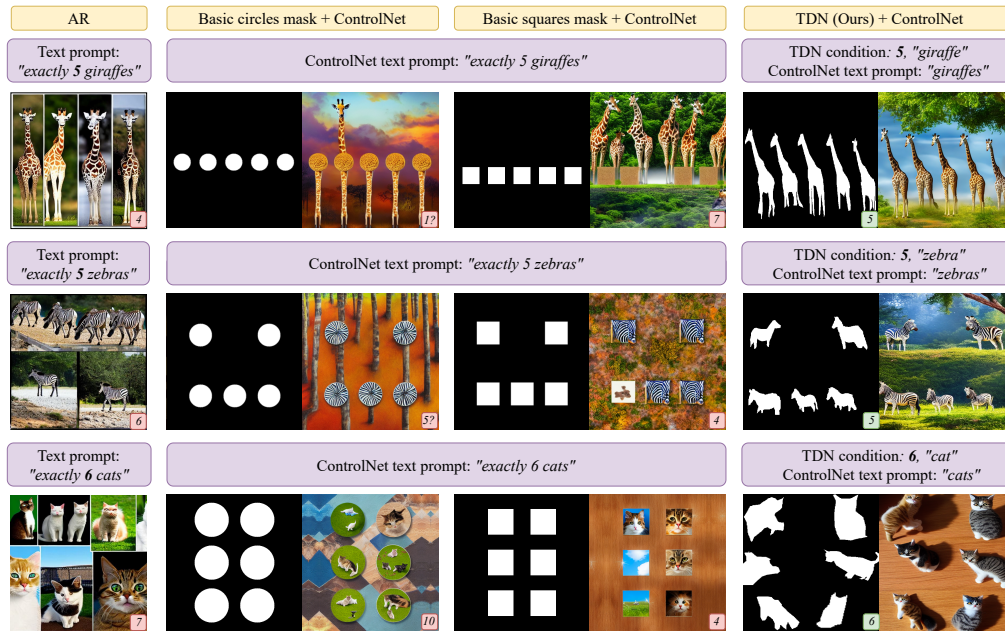


Figure 17: Qualitative comparison of using different conditioning types, particularly masks containing basic shapes like circles/squares to guide cardinality. ControlNet uses the SD1.5 backbone (Lvmin Zhang).

Table 11: Quantitative comparison of using different conditioning types for the COCO dataset. We include AR (Phung et al., 2024) and TDN results from Tab. 1 for convenience. Best results in **bold**

Method / Mask condition type	Accuracy \uparrow	Precision \uparrow	F1 \uparrow	FID (Image) \downarrow
Basic circles mask + ControlNet	0.2972 \pm 0.3592	0.3194 \pm 0.3046	0.3079 \pm 0.2208	49.47
Basic squares mask + ControlNet	0.3123 \pm 0.2114	0.3459 \pm 0.2209	0.3282 \pm 0.2618	47.16
AR (Bounding box w/ attention)	0.6379 \pm 0.2062	0.7360 \pm 0.1658	0.6611 \pm 0.1851	35.05
TDN (Ours)	0.8557 \pm 0.0805	0.8670 \pm 0.0636	0.8613 \pm 0.0970	21.28

M.1 USING SDXL-TURBO BACKBONE FOR CONTROLNET

The ControlNet results generated in this paper have used the SD1.5 backbone (Lvmin Zhang). In this sub-section, we generate ControlNet results using the newer SDXL-Turbo backbone⁸ (Sauer et al., 2025). The results are in Fig. 18, Fig. 19, and Fig. 20.

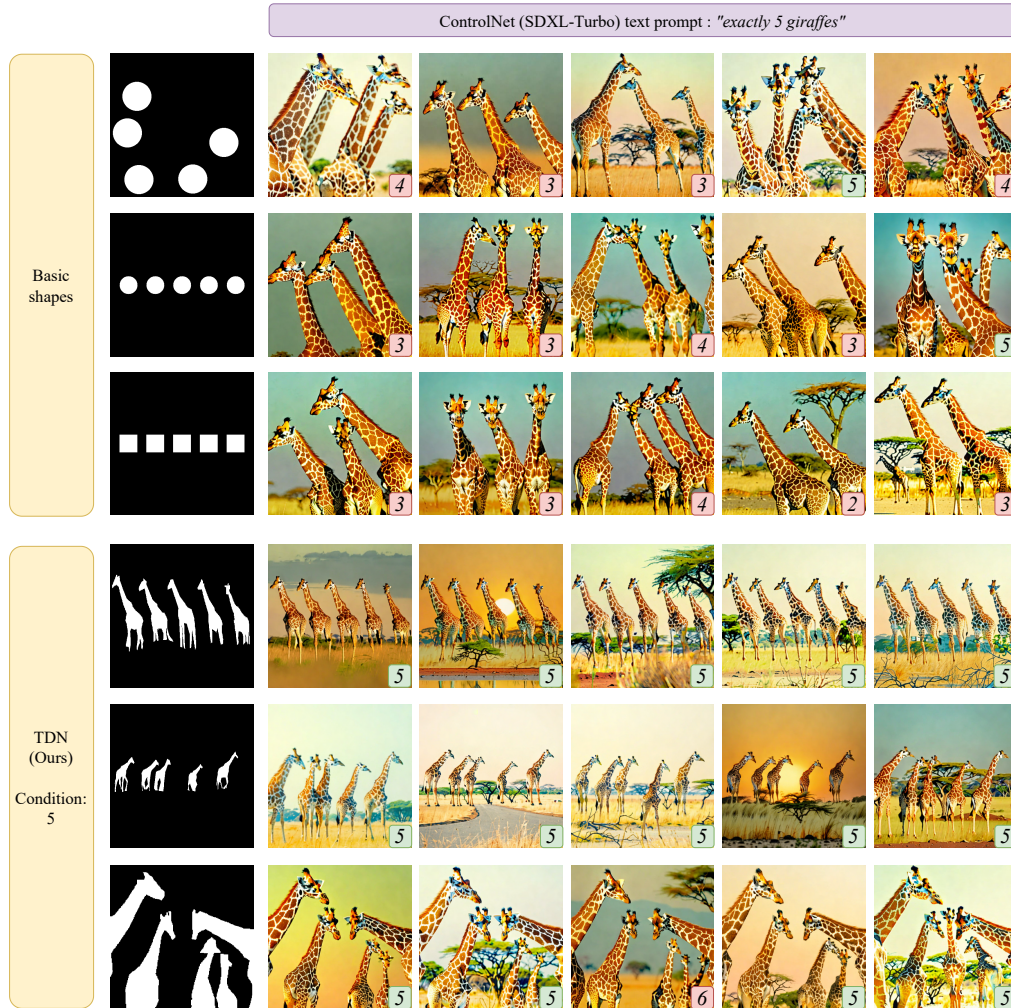


Figure 18: Generating exactly 5 giraffes. Qualitative comparison of using different conditioning types, such as masks containing basic shapes like circles/squares to guide cardinality. ControlNet uses the SDXL-Turbo backbone (Sauer et al., 2025).

⁸<https://huggingface.co/stabilityai/sd-xl-turbo>

1512
 1513
 1514
 1515
 1516
 1517
 1518
 1519
 1520
 1521
 1522
 1523
 1524
 1525
 1526
 1527
 1528
 1529
 1530
 1531
 1532
 1533
 1534
 1535
 1536
 1537
 1538
 1539
 1540
 1541
 1542
 1543
 1544
 1545
 1546
 1547
 1548
 1549
 1550
 1551
 1552
 1553
 1554
 1555
 1556
 1557
 1558
 1559
 1560
 1561
 1562
 1563
 1564
 1565

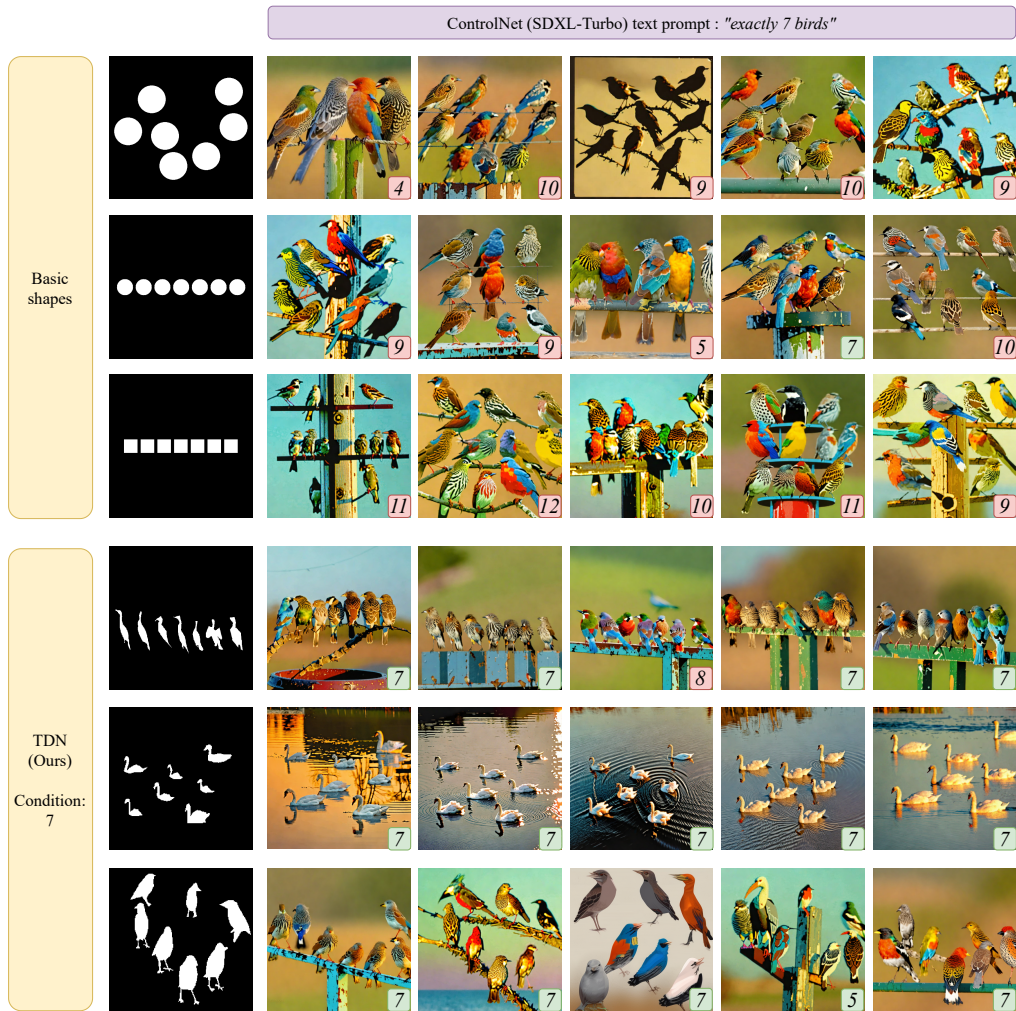


Figure 19: Generating exactly 7 birds. Qualitative comparison of using different conditioning types, such as masks containing basic shapes like circles/squares to guide cardinality. ControlNet uses the SDXL-Turbo backbone (Sauer et al., 2025).

1566
 1567
 1568
 1569
 1570
 1571
 1572
 1573
 1574
 1575
 1576
 1577
 1578
 1579
 1580
 1581
 1582
 1583
 1584
 1585
 1586
 1587
 1588
 1589
 1590
 1591
 1592
 1593
 1594
 1595
 1596
 1597
 1598
 1599
 1600
 1601
 1602
 1603
 1604
 1605
 1606
 1607
 1608
 1609
 1610
 1611
 1612
 1613
 1614
 1615
 1616
 1617
 1618
 1619

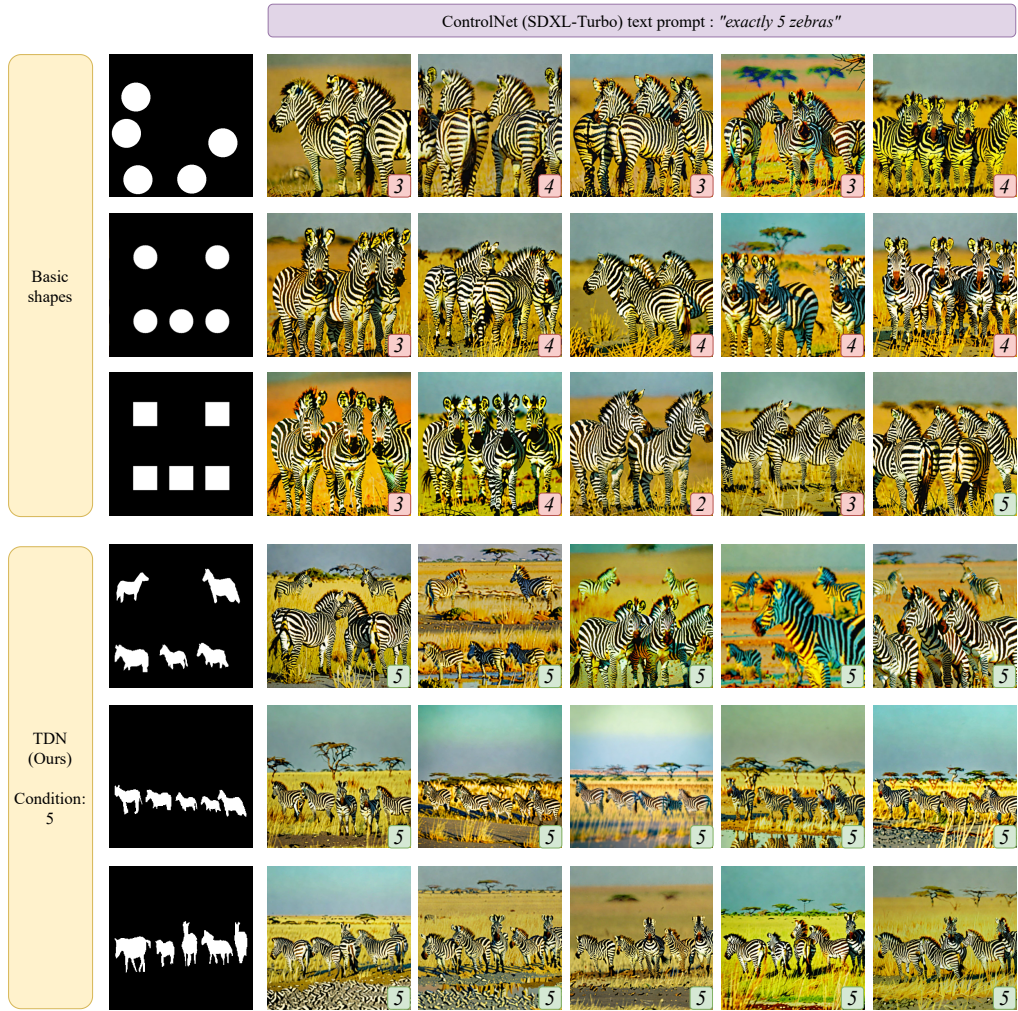


Figure 20: Generating exactly 5 zebras. Qualitative comparison of using different conditioning types, such as masks containing basic shapes like circles/squares to guide cardinality. ControlNet uses the SDXL-Turbo backbone (Sauer et al., 2025).

1620 N FUTURE WORK
1621

1622 Our proposed TDN currently uses persistent homology to control the number of objects (in 0-dim)
1623 and the number of holes (in 1-dim). However, persistent homology can theoretically be extended to
1624 higher dimensions, as persistence diagrams can capture topological features in arbitrary dimensions.
1625 For future work, we are looking into graph network generations, as well as 3D applications, where
1626 we can control not just connected components (0-dim) and holes (1-dim), but also voids (2-dim).
1627 3D point clouds and volumetric medical imaging data are important applications where maintaining
1628 topology is crucial in generating realistic synthetic data.
1629

1630
1631
1632
1633
1634
1635
1636
1637
1638
1639
1640
1641
1642
1643
1644
1645
1646
1647
1648
1649
1650
1651
1652
1653
1654
1655
1656
1657
1658
1659
1660
1661
1662
1663
1664
1665
1666
1667
1668
1669
1670
1671
1672
1673