# Supplmentary Material:
# L-CAD: Language-based Colorization with Any-level Descriptions using Diffusion Priors

**Zheng Chang**[#1]    **Shuchen Weng**[#2,3]    **Peixuan Zhang**[1]    **Yu Li**[4]    **Si Li**[*1]    **Boxin Shi**[2,3]

[1] School of Artificial Intelligence, Beijing University of Posts and Telecommunications
[2] National Key Laboratory for Multimedia Information Processing
School of Computer Science, Peking University
[3] National Engineering Research Center of Visual Technology
School of Computer Science, Peking University
[4] International Digital Economy Academy
{zhengchang98,pxzhang,lisi}@bupt.edu.cn
{shuchenweng, shiboxin}@pku.edu.cn    liyu@idea.edu.cn

## 6  Appendix

### 6.1  Robustness for contour estimation

We leverage a referring segmentation model to roughly estimate object contours mentioned in the description, which enables us to perform the instance-aware sampling strategy. To further demonstrate the robustness of our model, we manually annotate a sequence of contours ranging from coarse to fine and visualize the corresponding colorization results. As shown in Figure 8, our model presents a remarkable ability to produce condition-consistent colorization results even using imprecise contours. This is because the sampling is performed in the latent space using downsampled contours and the compression decoder in the pixel space could adaptively fix color bleeding issues.
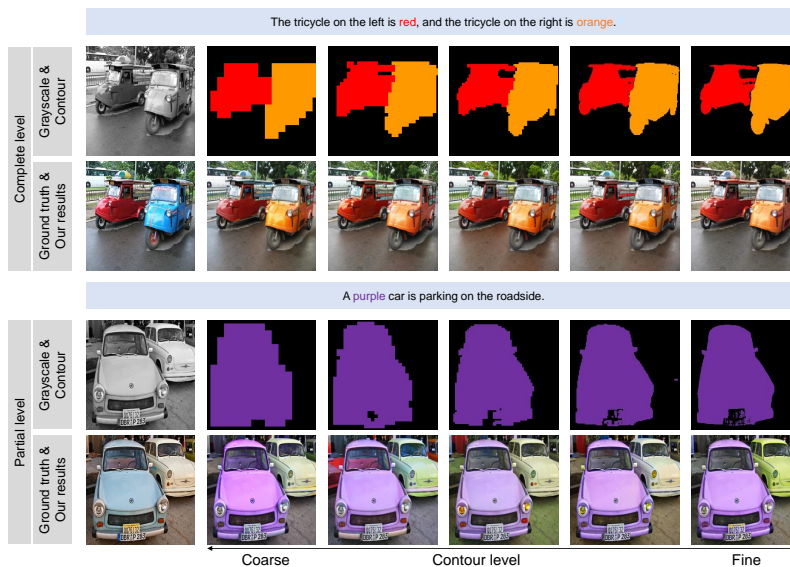


Figure 8: Visualization of colorization results by applying contours from coarse to fine.

# Equal contribution. * Corresponding author

## 6.2 Additional comparison results

As presented in Sec. 4.1 of the main paper, we comprehensively evaluate our method on language-based colorization datsets, where we make comparisons with language-based colorization methods (*e.g.*, LBIE [3], ML2018 [6], Xie2018 [14], L-CoDe [12], L-CoDer [1], and L-CoIns [2]) using complete-level and partial level descriptions, and comparisons with automatic colorization methods (*e.g.*, CIC [15], InstColor [9], ChromaGAN [10], BigColor [5], DISCO [13], and $CT^2$ [11]) using scarce-level descriptions.

Following the evaluation protocol on ImageNet dataset [7], we evaluate colorization results at the more common resolution of $256 \times 256$, instead of $224 \times 224$ resolution in previous works [1, 2, 12]. This higher resolution increases the difficulty of the colorization, resulting in slightly lower scores for the quantitative metrics (see Tab. 1 of the main paper), compared to those reported in previous works [1, 2, 12]. Additionally, we provide more qualitative comparison results with language-based colorization methods and automatic colorization methods in Fig. 9 and Fig. 10, respectively.

## 6.3 Additional ablation results

To demonstrate the effectiveness of our proposed luminance-guided image compression, semantic-aligned latent representation, and instance-aware sampling strategy (details in Sec. 4.3 of the main paper), we create three baselines by disabling corresponding modules. Additional qualitative ablation study results are shown in Fig. 11.

## 6.4 Additional application results

We demonstrate our generalization capability by showing more colorization results on legacy black-and-white photos in Fig. 12, where results are presented sequentially from left to right using descriptions at the complete, partial, and scarce levels.

## 6.5 Diverse colorization results

By leveraging the inherent stochasticity of diffusion models [4, 8], which sample noise from a Gaussian distribution at each step of the denoising process, our method could effectively generate diverse colorization results for unmentioned objects in descriptions. We show our diverse colorization results with partial-level and scarce-level descriptions in Fig. 13.

Furthermore, we present more challenging results of our L-CAD using complete-level, partial-level, and scarce-level descriptions in Fig. 14, Fig. 15, and Fig. 16, respectively. These demonstrate that our method could produce high-quality colorization results for diverse and complex scenarios.
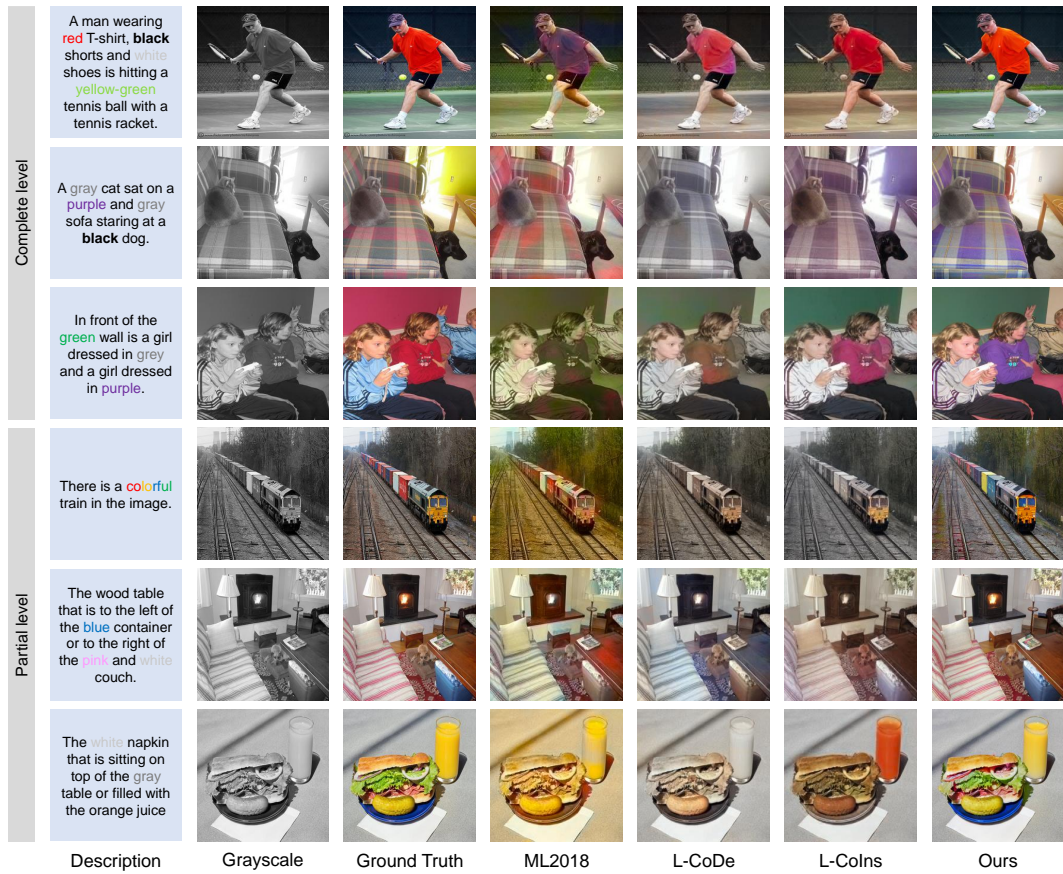
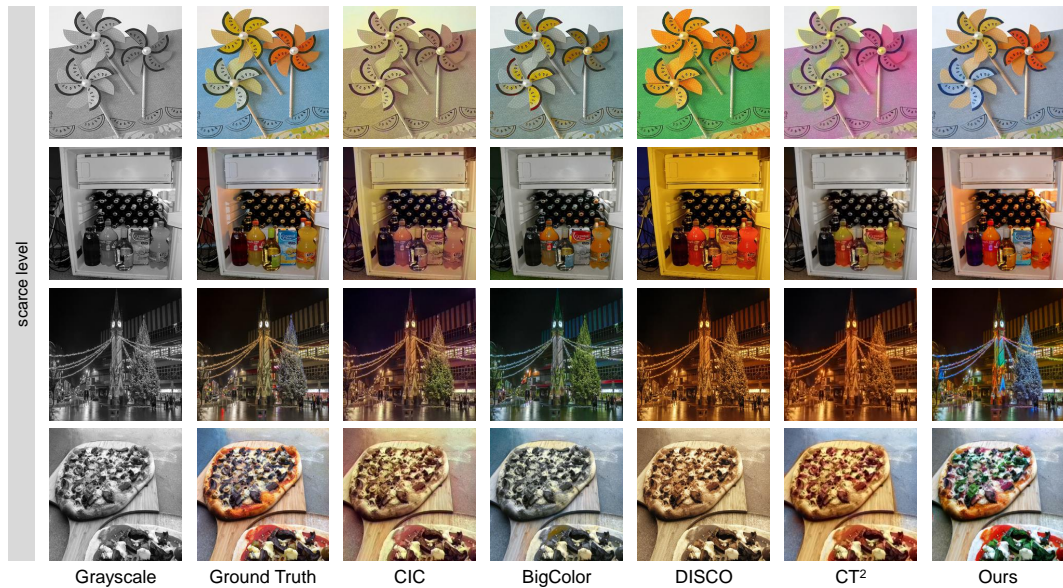Figure 9: More comparison results with language-based colorization methods.



Figure 10: More comparison results with automatic colorization methods.

Figure 11: More ablation study results.



Figure 12: More colorization results of legacy black-and-white photos.



Figure 13: Diverse colorization results.

There are two boats in the picture. The left is a green boat, and the right is a purple boat.

Two white teddy bears are dressed in pink and red.

There is chocolate cake on a white plate.

There is a colorful bus in front of the beige building.

The umbrella on the left is blue and the one on the right is red.

The car on the top is golden. The car on the bottom is blue.

The person on the left is wearing a blue T-shirt, while the person on the right is wearing a gray T-shirt.

There are two cows on the green grass. The one in front is white and the one behind is brown.

The left person is wearing black clothes, the middle person is wearing pink clothes, and the right person is wearing white clothes.

An adult and a child wearing a light pink T-shirt and a yellow helmet.

The left is orange pepper, the right is yellow pepper, the middle is red pepper.

There is a beige dog on the green grass, trying to grab a yellow Frisbee.

The beige dog is holding a green frisbee on the green grass.

The right is pale green cup, the middle is brown cup, the left is light pink cup.

The skateboards on the green grass are blue, green, white, orange and red from left to right.

Brown teddy bear lying on green grass.

The child on the left is wearing black clothes, while the child on the right is wearing pink striped clothes.

The top half of the hydrant is red and the bottom half is yellow.

Figure 14: More results of our L-CAD using complete-level descriptions.

A red bus is driving down the street.

A yellow school bus is traveling down a street.

A green mountain over a river.

Some vegetables on a white plate ready to be cut.

There is a cat on the blue sofa.

The wall in the living room is red.

Cat and dog are lying on the orange sofa.

There are red sofas in the room.

A slice of cake and a fork by sliced orange and a **black** Guinness.

The kitten is wearing a red Christmas hat.

Articulated bus is red on the street.

A fat orange cat sitting on the ground.

There is a **black** leather sofa in the living room.

The old woman on the bench is dressed in purple.

There is a **black** chair at the computer desk.

A dog lies on a orange blanket.

A car with red and white stripes was parked on the grass.
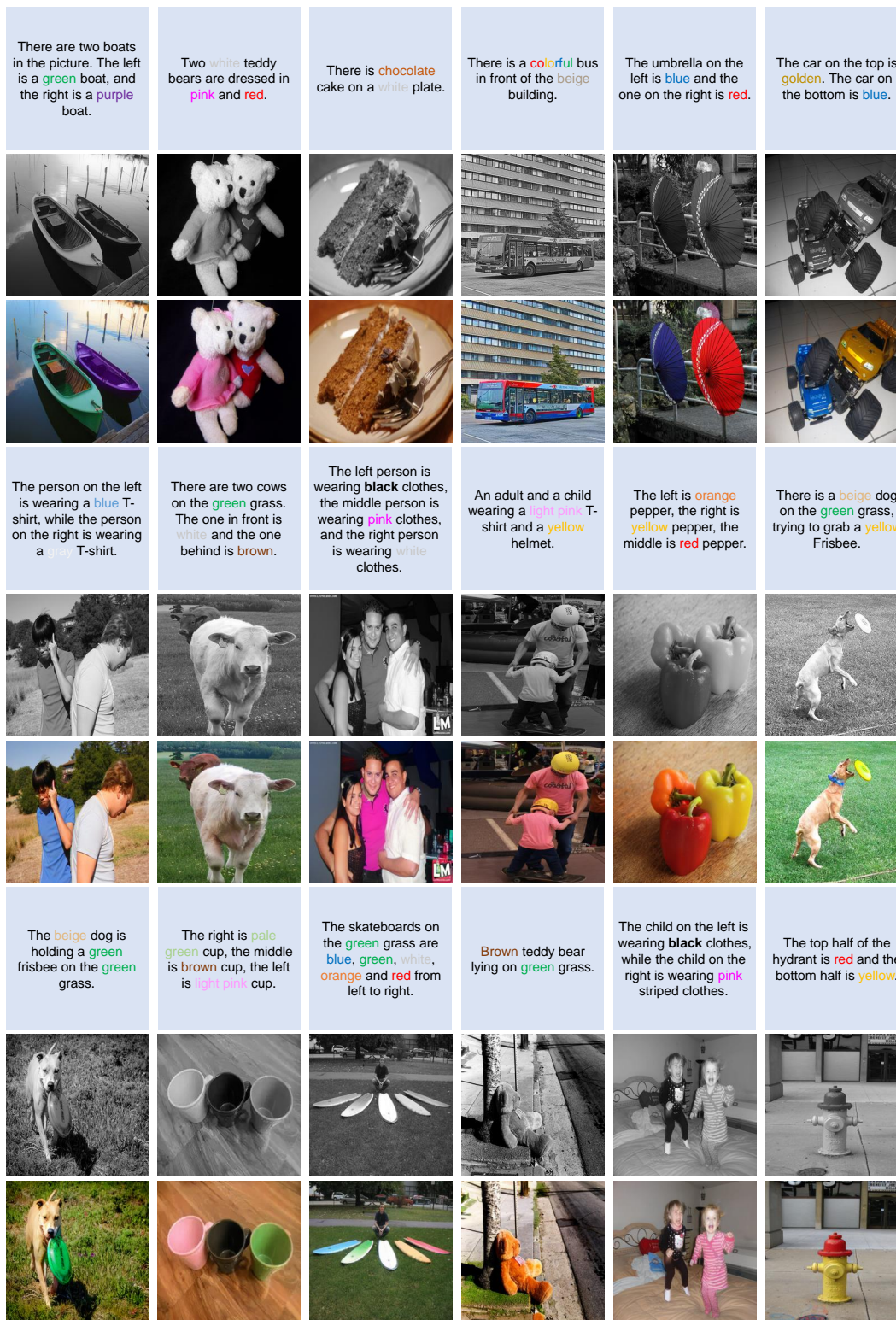
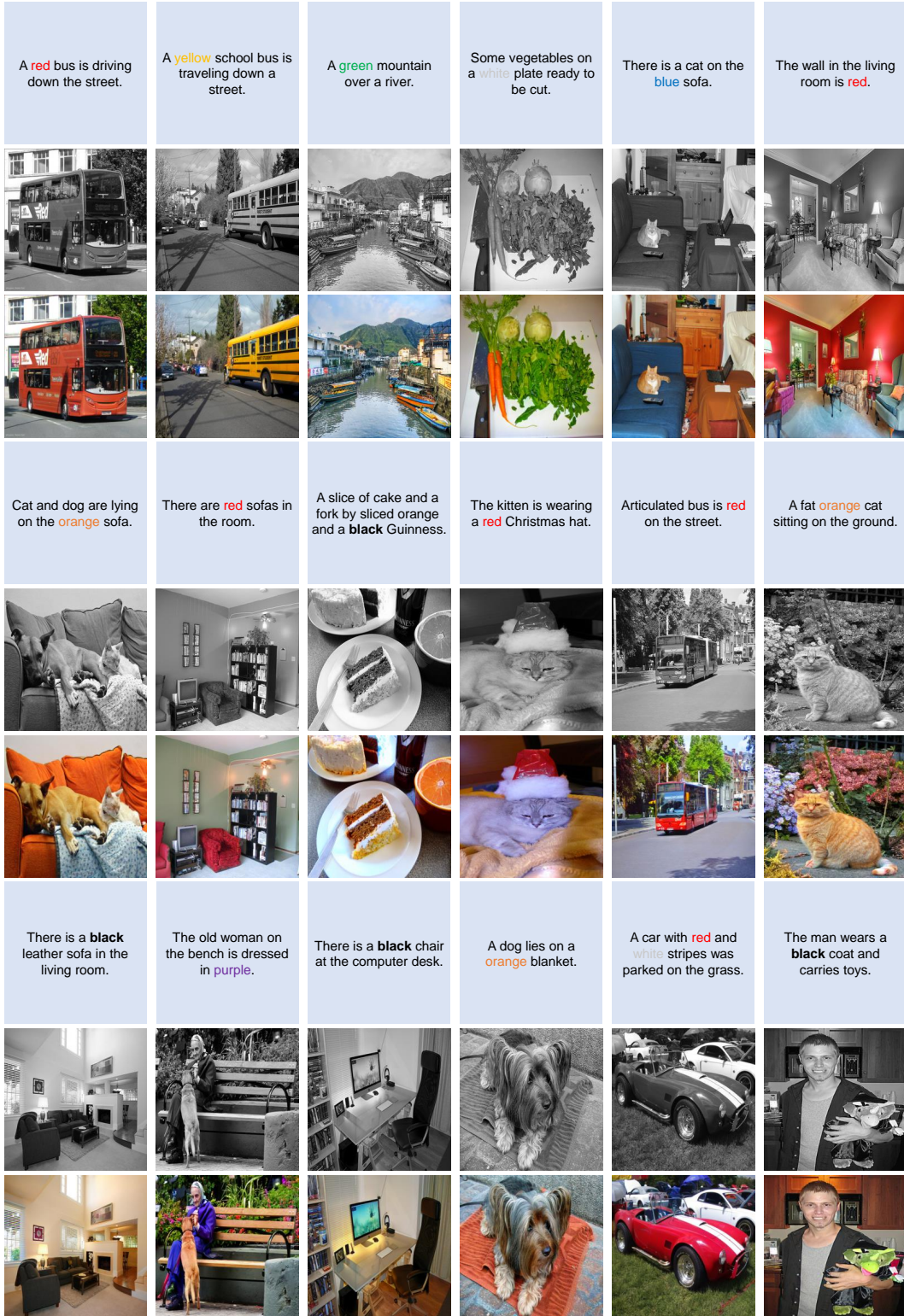The man wears a **black** coat and carries toys.

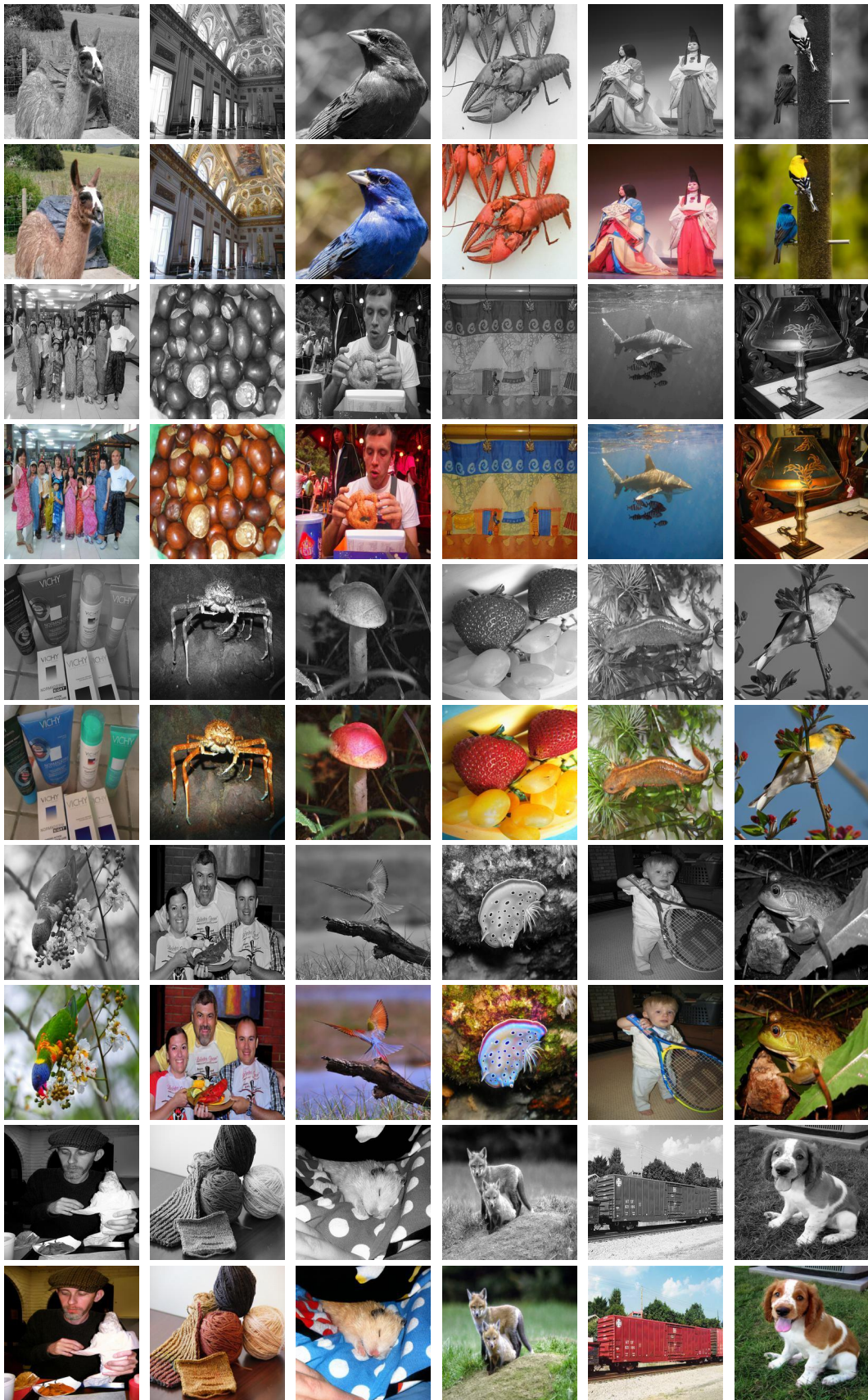Figure 15: More results of our L-CAD using partial-level descriptions.

Figure 16: More results of our L-CAD using scarce-level descriptions.

# References

[1] Z. Chang, S. Weng, Y. Li, S. Li, and B. Shi. L-CoDer: Language-based colorization with color-object decoupling transformer. In *ECCV*, 2022.

[2] Z. Chang, S. Weng, P. Zhang, Y. Li, S. Li, and B. Shi. L-CoIns: Language-based colorization with instance awareness. In *CVPR*, 2023.

[3] J. Chen, Y. Shen, J. Gao, J. Liu, and X. Liu. Language-based image editing with recurrent attentive models. In *CVPR*, 2018.

[4] J. Ho, A. Jain, and P. Abbeel. Denoising diffusion probabilistic models. In *NIPS*, 2020.

[5] G. Kim, K. Kang, S. Kim, H. Lee, S. Kim, J. Kim, S.-H. Baek, and S. Cho. BigColor: Colorization using a generative color prior for natural images. In *ECCV*, 2022.

[6] V. Manjunatha, M. Iyyer, J. Boyd-Graber, and L. Davis. Learning to color from language. In *NAACL*, 2018.

[7] O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang, A. Karpathy, A. Khosla, M. Bernstein, et al. Imagenet large scale visual recognition challenge. *IJCV*, 2015.

[8] J. Song, C. Meng, and S. Ermon. Denoising diffusion implicit models. In *ICLR*, 2021.

[9] J.-W. Su, H.-K. Chu, and J.-B. Huang. Instance-aware image colorization. In *CVPR*, 2020.

[10] P. Vitoria, L. Raad, and C. Ballester. ChromaGAN: Adversarial picture colorization with semantic class distribution. In *WACV*, 2020.

[11] S. Weng, J. Sun, Y. Li, S. Li, and B. Shi. $CT^2$: Colorization transformer via color tokens. In *ECCV*, 2022.

[12] S. Weng, H. Wu, Z. C. Chang, J. Tang, S. Li, and B. Shi. L-CoDe: Language-based colorization using color-object decoupled conditions. In *AAAI*, 2022.

[13] M. Xia, W. Hu, T.-T. Wong, and J. Wang. Disentangled image colorization via global anchors. *TOG*, 2022.

[14] Y. Xie. Language-guided image colorization. Master's thesis, ETH Zurich, Departement of Computer Science, 2018.

[15] R. Zhang, P. Isola, and A. A. Efros. Colorful image colorization. In *ECCV*, 2016.