

# An In-Depth Exploration of ControlNet Training Strategy Ablations for Conditional Generation

Anonymous CVPR submission

Paper ID

## Abstract

In recent years, diffusion models represented by Stable Diffusion [6] have achieved great success in the field of text-to-image generation. The emergence of ControlNet [8] has greatly enhanced the controllability of the generation process by introducing external structural conditions such as attitude maps and edge maps. However, most of the existing studies follow the standard paradigm of freezing the backbone of pre-trained models, and there is insufficient exploration of the optimal training strategies in different fine-tuning scenarios. This paper aims to systematically study the training strategies of ControlNet, especially the impact of the parameter update (*sdlocked*) of the pre-trained U-Net backbone and the injection range of the control module (*onlymidcontrol*) on the model performance. We conducted a series of ablation experiments on the Fill50K dataset under the conditions of pose and edge maps. The experimental results reveal a key finding: Compared with the standard practice of completely freezing the backbone, adopting a strategy of "unlocking the backbone and simplifying the control" (that is, unfreezing the Stable Diffusion backbone while only applying the ControlNet module to the middle layer of U-Net) can achieve a better balance on small datasets. This strategy not only accelerates the "emergent convergence" of the model, improves the quality of the generated images and the accuracy of conditional alignment, but also effectively alleviates the severe overfitting and catastrophic forgetting problems that occur when all parameters are completely thawed. Our research indicates that this fine-tuning method similar to transfer learning offers a new and efficient paradigm for training ControlNet on specific styles or small-scale datasets, providing valuable practical guidance for the design and application in the field of controllable image generation.

## 1. Introduction

034

In recent years, diffusion models [2] have achieved significant success in the field of text-to-image generation. Latent diffusion models, represented by Stable Diffusion [6], can generate high-quality images based on natural language prompts. However, relying solely on text prompts makes it difficult to precisely control the spatial layout and structural details of the image; for example, it is challenging to accurately describe the complex poses of characters or the detailed arrangement of a scene through words alone. This often leads to the need for repeated experimentation and adjustment of the prompt to gradually approach the desired composition.

The introduction of ControlNet [8] provides a breakthrough for this challenge. ControlNet is an architecture that adds conditional control to pre-trained text-to-image diffusion models. Without altering the original weights of the Stable Diffusion model, it incorporates external structural information (such as edge detection maps, human key-point pose diagrams, etc.) as additional conditional inputs, thereby providing fine-grained spatial guidance for the generation process. Specifically, ControlNet duplicates and freezes the weights of the original Stable Diffusion model, using "zero-convolution" layers to connect a new conditional branch to the feature layers of the diffusion model. Since these convolutional layers start with zero initial weights, they have minimal impact on the original model at the beginning of training, gradually learning the conditional control signals from scratch. This design ensures that adding the control branch does not disrupt the original generation capabilities while effectively aligning the generated images with the input conditions. Experiments show that ControlNet can be stably trained under various conditions (regardless of the size of the dataset) and supports flexible combinations of single or multiple conditions such as edges, depth, segmentation, and human poses. By allowing users to provide additional structural images as conditions, ControlNet achieves finer control over image generation compared to pure text, enhancing the match be-

073      tween the generated results and user intentions, and expanding  
074      the application prospects of diffusion models in creative  
075      design, film and animation, content editing, and more.

076      Based on this background, we observe that current re-  
077      search and applications of ControlNet mainly focus on con-  
078      trolling the poses of static images of people, with insuf-  
079      ficient exploration of its broader potential uses and areas  
080      for improvement. For instance, in tasks such as video syn-  
081      thesis, medical imaging, biological feature simulation, and  
082      complex scene layout, there is great potential for structural  
083      condition control. Similarly, there are possibilities for fur-  
084      ther optimization in the ControlNet model architecture and  
085      training strategies to improve the quality of generation and  
086      the fidelity to the conditions. Therefore, this paper takes  
087      image generation under the conditions of pose maps (such  
088      as human keypoints) and edge maps (such as HED edge de-  
089      tection maps [7]) as a starting point, delving into how to en-  
090      hance the generation quality and conditional control capa-  
091      bilities of ControlNet. This research is of great significance:  
092      it can enrich the application of ControlNet in more fields,  
093      making the generation of diffusion models more "what you  
094      see is what you get," and also provide new ideas for the de-  
095      sign of the next generation of controllable image generation  
096      models.

## 097      2. Objectives

098      This study aims to systematically enhance the image gener-  
099      ation effects based on pose maps and edge maps under con-  
100      ditional control using ControlNet. The specific objectives  
101      include:

102      Expanding Application Scenarios: To explore the ap-  
103      plicability of ControlNet in a wider range of tasks, such  
104      as character action transfer, medical image synthesis, com-  
105      plex scene layout generation, and video keyframe guidance,  
106      thereby validating the value of structured conditions across  
107      different fields.

108      Architectural Innovation and Optimization: To investi-  
109      gate improvements in the architectural design of the Con-  
110      trolNet model, including strategies for inserting conditions  
111      at intermediate layers of different diffusion models, opti-  
112      mizing the connection methods between the control branch  
113      and the backbone network, and developing efficient mech-  
114      anisms for fusing conditional features with generated fea-  
115      tures.

116      Improving Training Strategies: To develop effective  
117      training schemes that enhance the model's control capa-  
118      bilities and stability, such as reasonable freezing/thawing  
119      strategies for the main model, methods to improve the ex-  
120      pression and alignment of condition signals during training,  
121      and regularization techniques to boost training stability.

122      We aim to significantly improve the quality of images  
123      generated under target pose or edge conditions, ensuring  
124      that the output results more accurately match the given con-

ditions in terms of details and structure, while maintaining  
125      the diversity and visual realism of the images.

### 126      2.1. Research Questions

127      To achieve the research objectives, this paper will focus on  
128      exploring the following key issues:

- 129      1. Expansion of ControlNet Applications: Besides gen-  
130      erating static character images, under the condition of  
131      using pose keypoint maps and edge detection maps,  
132      in what other generation tasks can ControlNet be ex-  
133      panded? How does it play a role in motion transfer in  
134      dynamic videos, synthesis of medical images, genera-  
135      tion from scene layout to images, and continuous guid-  
136      ance for video keyframes, and what challenges does it  
137      face?
- 138      2. Optimization of Architecture Design: Currently, Con-  
139      trolNet introduces control signals by adding a zero-  
140      convolution conditional branch to a pre-trained diffusion  
141      model. Is there a more optimal network architecture  
142      or fusion mechanism that could improve the efficiency  
143      of condition utilization? For example, at which levels  
144      (shallow vs. deep, encoder vs. decoder) of the diffusion  
145      UNet should the conditional features be injected, and  
146      what kind of connection and fusion methods (addition,  
147      concatenation, or attention mechanisms) would yield the  
148      best generation results and precise condition alignment?
- 149      3. Improvement of Training Strategies: Without compro-  
150      mising the original generative capabilities, how can ef-  
151      fective training strategies be designed to enhance Con-  
152      trolNet's sensitivity and control accuracy towards con-  
153      ditions? For instance, in what situations should the  
154      weights of the pre-trained diffusion model be frozen, or  
155      gradually unfrozen to adapt to new domains? Can the  
156      introduction of additional loss functions (such as cycle  
157      consistency loss or intermediate feature alignment) im-  
158      prove the match between conditions and outputs, thereby  
159      enhancing training stability and generation quality?

### 160      2.2. Related Work

161      **T2I-Adapter: A Lightweight Control Module.** Follow-  
162      ing ControlNet [8], researchers have explored alternative ef-  
163      ficient methods for incorporating structural conditions. T2I-  
164      Adapter [5] introduces a lightweight, plug-and-play adapter  
165      module that injects external control signals without modi-  
166      fying the original Stable Diffusion architecture. This ap-  
167      proach trains small convolutional networks for each con-  
168      dition type (e.g., sketch, edge, depth, keypoints) to en-  
169      code multi-scale features, which are then added to the fea-  
170      ture maps at different layers of Stable Diffusion's U-Net  
171      encoder. While the base model remains frozen, only the  
172      adapter parameters are updated. With substantially fewer  
173      parameters (e.g., 77M compared to ControlNet's 567M),  
174      T2I-Adapter enables faster inference while achieving com-

parable or superior structural and textual alignment on datasets like COCO [4]. Additionally, it supports flexible module composition: adapters for different conditions can be weighted and combined during inference to enable multi-condition control without retraining. For instance, merging outputs from sketch and color palette adapters allows simultaneous control over shape and color. The method also employs non-uniform timestep sampling during training to enhance guidance for low-level visual features (e.g., edges, colors). Overall, T2I-Adapter demonstrates an effective pathway for leveraging large models’ implicit capabilities through compact modules, offering a valuable complement to the ControlNet framework.

**Other Structured Conditional Generation Methods.** Beyond the above, several recent architectures have been proposed to enhance controllability by integrating structural information into diffusion models. For example, GLIGEN [3] focuses on controlling scene layout and object placement. By incorporating learnable gated units into Stable Diffusion’s cross-attention layers, GLIGEN accepts bounding box coordinates and corresponding object labels as additional inputs, enabling precise object positioning. Similar to other methods, it freezes most of the pre-trained weights and trains only a small set of gating and embedding parameters, thereby endowing the model with grounding capabilities. Experiments show that GLIGEN, built upon pre-trained diffusion models, achieves higher image quality and layout accuracy in layout-to-image tasks compared to models trained from scratch. Another related direction is multi-conditional diffusion, where models are trained to accept composite controls (e.g., text, segmentation, edges) in a single forward pass. However, such approaches typically require full retraining or large-scale fine-tuning, incurring high computational costs and potential limitations in generalizing to unseen condition combinations.

**Training Strategies and Enhanced Control Performance.** As structural conditional generation methods evolve, improving condition adherence and output quality has become a key research focus. Some studies introduce additional constraint losses during training to strengthen condition-image consistency. For example, a pixel-level cycle consistency loss can be used: pre-trained discriminative models (e.g., edge detectors, segmentation networks) are applied to generated images to extract condition signals, which are then compared with the original input conditions to directly optimize alignment. To mitigate the high computational cost of full diffusion sampling during loss calculation, an efficient single-step perturbation strategy can approximate the generated output for consistency evaluation.

Another line of work explores intermediate feature alignment training. At each denoising step, lightweight convolutional probes are trained to reconstruct input condition maps (e.g., edges or depth) from the U-Net’s intermediate fea-

tures. During training, a consistency loss is computed between the predicted “pseudo-condition” from noisy latents and the ground-truth condition, encouraging the model to maintain condition awareness throughout the diffusion process. This strategy embeds control signals more deeply into generation, enhancing structural fidelity and control precision.

In summary, recent advances—spanning model architectures (e.g., ControlNet, T2I-Adapter, GLIGEN) and training strategies (e.g., frozen fine-tuning, adapter composition, consistency constraints)—have continuously advanced the field of structurally conditioned controllable generation. Building upon these works, this study further extends applicable scenarios and proposes improvements for generation tasks under pose and edge map conditions, contributing to enhanced quality and precision in controllable image synthesis.

## 3. Methods

### 3.1. Merging Conditional Control with the Stable Diffusion Backbone

ControlNet [8] augments a pretrained text-to-image diffusion model by introducing a secondary, learnable pathway that injects spatial conditioning signals into the denoising network while preserving the capabilities of the original backbone (Stable Diffusion [6]). Let  $F(\cdot; \Theta)$  denote a pre-trained U-Net block that maps a feature map  $x \in \mathbb{R}^{h \times w \times c}$  to an output  $y = F(x; \Theta)$ . To incorporate structural conditions such as edges, poses, or depth maps, ControlNet constructs a *dual-path* architecture: the original block is kept intact as a frozen backbone, and a trainable copy  $F(\cdot; \Theta_c)$  is created to process condition-aware features. These two paths are fused through lightweight  $1 \times 1$  convolutions, denoted by  $Z(\cdot; \Theta_z)$ , which align feature dimensions and modulate the backbone activations with task-specific information.

Given a conditioning feature map  $c$ , the ControlNet block computes

$$y_c = F(x; \Theta) + Z(F(x + Z(c; \Theta_{z1}); \Theta_c); \Theta_{z2}). \quad (1)$$

This formulation allows the trainable copy to interpret the conditioning signal and contribute corrective residuals, while the frozen backbone ensures that the pretrained generative behavior is preserved. The connector convolutions  $Z(\cdot; \Theta_{z1})$  and  $Z(\cdot; \Theta_{z2})$  serve as projection layers that map the conditioning features into the appropriate channel space and then inject the processed representation back into the main feature stream.

The conditioning feature map  $c$  itself is derived from a raw conditioning image  $c_i \in \mathbb{R}^{512 \times 512 \times 3}$  through a small encoder  $E(\cdot)$ . This encoder consists of a few strided convolution layers, where the first convolution primarily aligns

279 the spatial resolution with the latent space of Stable Diffusion [6] (typically  $64 \times 64$ ). The goal of  $E(\cdot)$  is not to  
 280 perform deep semantic extraction but to convert the conditioning  
 281 image into a feature map compatible with the U-Net  
 282 resolution at which ControlNet operates.  
 283

### 284 285 3.2. Training Considerations: Frozen Backbone and Zero-Initialized Connectors

286 To stabilize learning and prevent catastrophic forgetting, all  
 287 pretrained parameters  $\Theta$  of Stable Diffusion [6] are frozen  
 288 throughout training. Since conditional datasets are often  
 289 much smaller than the large-scale corpora used to train  
 290 the original model, freezing the backbone prevents over-  
 291 specialization and ensures that the core generative abilities  
 292 remain intact. Only the parameters of the trainable copy  $\Theta_c$ ,  
 293 the conditioning encoder  $E(\cdot)$ , and the connector convolu-  
 294 tions  $\Theta_{z1}, \Theta_{z2}$  are updated.

295 A second key design choice is the zero-initialization  
 296 of the connector convolutions. At initialization, both  
 297  $Z(c; \Theta_{z1})$  and  $Z(\cdot; \Theta_{z2})$  output zero for any input, causing  
 298 the entire ControlNet block to reduce to

$$299 \quad y_c = F(x; \Theta), \quad (2)$$

300 which means the model behaves identically to the pretrained  
 301 Stable Diffusion at the start of training. This avoids injecting  
 302 untrained, potentially harmful noise into the backbone  
 303 and allows the influence of the conditional branch to “grow  
 304 in” smoothly as training progresses. Empirically, this strat-  
 305 egy leads to a stable optimization process and a character-  
 306 istic “sudden convergence” effect: after a period in which  
 307 the model behaves like the vanilla backbone, the connec-  
 308 tor weights learn meaningful transformations that allow the  
 309 network to abruptly acquire strong conditioning fidelity.

310 Together, the dual-path structure, frozen backbone, and  
 311 zero-initialized connector convolutions form a robust and  
 312 data-efficient architecture. ControlNet is thus able to learn  
 313 spatially localized, task-specific controls while preserving  
 314 the expressive power and visual quality of the underlying  
 315 diffusion model.

### 316 3.3. Conditioning Generation

317 In our experiments, we evaluated the *Human Pose* con-  
 318 ditioning. The human pose condition is generated using an  
 319 OpenPose [1]-based keypoint detector, which produces a  
 320 structured skeletal representation that captures the global  
 321 configuration and articulation of the subject. This con-  
 322 ditioning type is widely used in controllable generation  
 323 tasks because pose encodes coarse geometry and provides  
 324 a salient spatial scaffold for the generative model to follow.

325 For the condition-type experiments, we adopted a pub-  
 326 licly available pretrained ControlNet model from Hugging-  
 327 Face. This model was originally trained on large-scale

328 general-domain datasets and has demonstrated strong gen-  
 329 eralization to diverse image conditions. In our replication,  
 330 the human pose conditioning yielded high-quality genera-  
 331 tions, and the model consistently interpreted the skeletal  
 332 signals provided by the pose skeletons. The results indicate  
 333 that the pretrained ControlNet weights possess robust rep-  
 334 resentational capacity and that the conditioning mechanism  
 335 is sufficiently expressive to guide the generation process in  
 336 a semantically meaningful way.

### 337 3.4. Training Procedure

338 We train our models using the Fill50K dataset released with  
 339 the ControlNet training pipeline [8], a collection of 50,000  
 340 image–mask pairs used for conditioning-guided finetuning.  
 341 The dataset contains diverse structural patterns and  
 342 object boundaries, providing meaningful supervision for  
 343 conditioning-guided generation even though it is signifi-  
 344 cantly smaller than the large-scale datasets used for training  
 345 Stable Diffusion. As our generative backbone, we adopt  
 346 Stable Diffusion v1.5 [6], a latent diffusion model trained  
 347 on LAION-5B that produces high-quality natural images  
 348 through a U-Net denoiser operating in a compressed latent  
 349 space. Its strong visual prior makes it an ideal foundation  
 350 for studying conditional finetuning behavior.

351 To understand the effect of architectural choices dur-  
 352 ing finetuning, we performed an ablation study follow-  
 353 ing the configuration conventions in the ControlNet train-  
 354 ing documentation.<sup>1</sup> Specifically, we varied two key flags:  
 355 `sd_locked` and `only_mid_control`. These options de-  
 356 termine which portions of the Stable Diffusion U-Net are  
 357 updated during training.

**Definition of `sd_locked`.** When the configuration  
 358 `sd_locked = True`, all parameters of the original Sta-  
 359 ble Diffusion model are frozen. Only the ControlNet branch  
 360 (trainable copy, connector convolutions, and condition  
 361 encoder) receives gradient updates. This configuration pre-  
 362 serves the pretrained generative behavior and prevents the  
 363 model from drifting away from the domain of natural im-  
 364 ages. When `sd_locked = False`, the Stable Diffusion  
 365 U-Net becomes trainable, greatly increasing the number of  
 366 parameters updated at each iteration. This allows strong  
 367 adaptation to the conditioning task, but also increases the  
 368 risk of overfitting and destabilizing the pretrained backbone.

**Definition of `only_mid_control`.** When the config-  
 370 uration `only_mid_control = True`, ControlNet at-  
 371 taches trainable modules only to the middle block of the  
 372 U-Net, rather than to all encoder blocks. This substantially  
 373 reduces the number of trainable parameters, as the skip-  
 374 connected encoder pathways remain unchanged. When

<sup>1</sup><https://github.com/lillyasviel/ControlNet/blob/main/docs/train.md>

376       only\_mid\_control = False, ControlNet is attached  
 377       to every encoder block, maximizing control capacity but  
 378       also enlarging the effective finetuning footprint.

379       These configurations produce four experimental vari-  
 380       ants. To ensure column compatibility, we present the results  
 381       in a compact table:

ExpID	sd_locked	mid_only	Sudden Conv.	Behavior
A	True	False	2700	Medium quality
B	True	True	N/A	Low quality
C	False	False	1500	Overfitting, forgetting
D	False	True	2400	High quality, mild overfit

Table 1. Ablation results of training configurations.

382       Across all configurations, we consistently observe the  
 383       “sudden convergence” phenomenon: for several training  
 384       steps, the model behaves similarly to the pretrained back-  
 385       bone, and then at a specific iteration, it abruptly begins  
 386       to follow conditioning signals and generate structurally  
 387       aligned images. However, the step at which this conver-  
 388       gence occurs varies significantly between training settings,  
 389       as shown in Table 1.

390       When both the Stable Diffusion model and all ControlNet  
 391       modules are trainable (`sd_locked = False`,  
 392       `only_mid_control = False`), the network receives  
 393       gradients across a very large parameter space. This acceler-  
 394       ates the sudden convergence (1500 steps) and yields high-  
 395       quality generations early in training. However, this flexi-  
 396       bility also leads to severe overfitting: as training continues,  
 397       images develop blurry or irregular boundaries, and catas-  
 398       trophic forgetting may occur, with some samples collapsing  
 399       into disordered patterns.

400       When the backbone is frozen but all ControlNet blocks  
 401       are active (`sd_locked = True`, `only_mid_control`  
 402       = `False`), the training process becomes more stable. This  
 403       is the recommended configuration in the ControlNet pa-  
 404       per [8]. The model maintains reasonable image quality and  
 405       avoids overfitting even at high training iterations. Notably,  
 406       the generated images tend to inherit aesthetic properties of  
 407       Stable Diffusion v1.5 itself—such as characteristic textures  
 408       and shading—indicating that the frozen backbone strongly  
 409       shapes the model output.

410       Restricting ControlNet to the middle block while  
 411       keeping the backbone frozen (`sd_locked = True`,  
 412       `only_mid_control = True`) significantly limits the  
 413       effective gradient pathways. As a result, we observe no sud-  
 414       den convergence even after 12k steps, and the model fails to  
 415       learn meaningful structural patterns from the dataset. The  
 416       limited backpropagation route through the single mid-block  
 417       likely reduces optimization efficiency.

418       Interestingly, when the backbone is unfrozen but Con-  
 419       trolNet is attached only to the middle block (`sd_locked`  
 420       = `False`, `only_mid_control = True`), the model  
 421       becomes more stable than in the fully trainable case. Over-

422       fitting still appears but is substantially mitigated, and im-  
 423       age quality improves. We hypothesize two contribut-  
 424       ing factors: (1) the number of trainable parameters under  
 425       this configuration becomes comparable to the standard  
 426       setting (`sd_locked = True`, `only_mid_control =`  
 427       `False`), whereas the fully trainable configuration involves  
 428       significantly more parameters; and (2) because only a sin-  
 429       gle decoding pathway is controllable, backpropagation be-  
 430       comes more direct, enabling more effective learning of  
 431       dataset-specific features without destabilizing the entire U-  
 432       Net.

433       Taken together, these results highlight the trade-offs be-  
 434       tween model stability, learning efficiency, and overfitting.  
 435       Allowing too many components to update accelerates con-  
 436       vergence but risks rapid degradation, whereas restricting  
 437       updates can improve robustness but slow down or even pre-  
 438       vent meaningful learning. Based on our observations, it  
 439       may be advantageous to unfreeze the backbone while at-  
 440       taching ControlNet only to the middle block (`sd_locked`  
 441       = `False`, `only_mid_control = True`) when fine-  
 442       tuning on a small dataset with a strong and distinctive style,  
 443       which behaves more like a transfer-learning procedure. In  
 444       contrast, when the pretrained generative properties of Sta-  
 445       ble Diffusion v1.5 are crucial—such as maintaining real-  
 446       ism or preserving its characteristic aesthetic—and the pri-  
 447       mary goal of finetuning is to teach the model how to in-  
 448       terpret a new conditioning modality rather than to change  
 449       its visual style, the configuration (`sd_locked = True`,  
 450       `only_mid_control = False`) is preferable.

## 4. Members and Contributions

**Hao Zheng:** In charge of model training and experimen-  
 452       tal work, setting up test platforms for different tasks (such  
 453       as action transfer, medical synthesis, etc.), comparing and  
 454       analyzing the performance differences between Baseline  
 455       (original ControlNet) and improved methods, ensuring the  
 456       rigor and reliability of the experiments. **Zhiyi Chen:** Re-  
 457       sponsible for preparing and processing training/ testing data  
 458       and conducting quantitative and qualitative analysis of ex-  
 459       perimental results, including metric evaluation, visualiza-  
 460       tion result comparison, identifying the strengths and weak-  
 461       nesses of the methods and suggesting improvement direc-  
 462       tions. **Rongpei Li:** Mainly responsible for proposing and  
 463       implementing architectural improvement plans, including  
 464       new ControlNet intermediate layer connection strategies,  
 465       conditional feature fusion module designs, etc., and collab-  
 466       orating on the formulation of training schemes.

## References

- [1] Zhe Cao, Tomas Simon, Shih-En Wei, and Yaser Sheikh. Realtime multi-person 2d pose estimation using part affinity fields. In *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 1302–1310, 2017.
- [2] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffu-

- 474 sion probabilistic models. In *Advances in Neural Information*  
475 *Processing Systems*, 2020.
- 476 [3] Yuheng Li, Haotian Liu, Qingyang Wu, Fangzhou Mu, Jian-  
477 wei Yang, Jianfeng Gao, Chunyuan Li, and Yong Jae Lee.  
478 Gligen: Open-set grounded text-to-image generation. In  
479 *2023 IEEE/CVF Conference on Computer Vision and Pattern*  
480 *Recognition (CVPR)*, pages 22511–22521, 2023.
- 481 [4] Tsung-Yi Lin, Michael Maire, Serge Belongie, Lubomir  
482 Bourdev, Ross Girshick, James Hays, Pietro Perona, Deva  
483 Ramanan, C. Lawrence Zitnick, and Piotr Dollár. Microsoft  
484 coco: Common objects in context. In *Computer Vision –*  
485 *ECCV 2014*, pages 740–755, 2014.
- 486 [5] Chong Mou, Xintao Wang, Liangbin Xie, Yanze Wu, Jian  
487 Zhang, Zhongang Qi, Ying Shan, and Xiaohu Qie. T2i-  
488 adapter: Learning adapters to dig out more controllable abil-  
489 ity for text-to-image diffusion models. *Proceedings of the*  
490 *AAAI Conference on Artificial Intelligence*, 38(5):4296–4304,  
491 2024.
- 492 [6] Robin Rombach, Andreas Blattmann, Dominik Lorenz,  
493 Patrick Esser, and Björn Ommer. High-resolution image syn-  
494 thesis with latent diffusion models. In *2022 IEEE/CVF Con-*  
495 *ference on Computer Vision and Pattern Recognition (CVPR)*,  
496 pages 10674–10685, 2022.
- 497 [7] Saining Xie and Zhuowen Tu. Holistically-nested edge de-  
498 tection. In *2015 IEEE International Conference on Computer*  
499 *Vision (ICCV)*, 2015.
- 500 [8] Lvmin Zhang, Anyi Rao, and Maneesh Agrawala. Adding  
501 conditional control to text-to-image diffusion models. In  
502 *2023 IEEE/CVF International Conference on Computer Vi-*  
503 *sion (ICCV)*, pages 3813–3824, 2023.