

---

# Uni-ControlNet: All-in-One Control to Text-to-Image Diffusion Models

---

Shihao Zhao<sup>1†</sup> Dongdong Chen<sup>2\*</sup> Yen-Chun Chen<sup>2</sup> Jianmin Bao<sup>2</sup> Shaozhe Hao<sup>1</sup>  
Lu Yuan<sup>2</sup> Kwan-Yee K. Wong<sup>1\*</sup>  
<sup>1</sup> University of Hong Kong <sup>2</sup> Microsoft

## Abstract

Text-to-Image diffusion models have made tremendous progress over the past two years, enabling the generation of highly realistic images based on open-domain text descriptions. However, despite their success, text descriptions often struggle to adequately convey detailed controls, even when composed of long and complex texts. Moreover, recent studies have also shown that these models face challenges in understanding such complex texts and generating the corresponding images. Therefore, there is a growing need to enable more control modes beyond text description. In this paper, we introduce Uni-ControlNet, a novel approach that allows for the simultaneous utilization of different local controls (e.g., edge maps, depth map, segmentation masks) and global controls (e.g., CLIP image embeddings) in a flexible and composable manner within one model. Unlike existing methods, Uni-ControlNet only requires the fine-tuning of two additional adapters upon frozen pre-trained text-to-image diffusion models, eliminating the huge cost of training from scratch. Moreover, thanks to some dedicated adapter designs, Uni-ControlNet only necessitates a constant number (i.e., 2) of adapters, regardless of the number of local or global controls used. This not only reduces the fine-tuning costs and model size, making it more suitable for real-world deployment, but also facilitate composability of different conditions. Through both quantitative and qualitative comparisons, Uni-ControlNet demonstrates its superiority over existing methods in terms of controllability, generation quality and composability. Code is available at <https://github.com/ShihaoZhaoZSH/Uni-ControlNet>.

## 1 Introduction

In recent two years, diffusion models [1–3] have gained significant attention due to their remarkable performance in image synthesis tasks. Therefore, text-to-image (T2I) diffusion models [4–8] have emerged as a popular choice for synthesizing high-quality images based on textual inputs. By training on large-scale datasets with large models, these T2I diffusion models demonstrate exceptional ability in creating images that closely resemble the content described in text descriptions, and facilitate the connection between textual and visual domains. The substantially improved generation quality in capturing intricate texture details and complex relationships between objects, makes them highly suitable for various real-world applications, including but not limited to content creation, fashion design, and interior decoration.

However, text descriptions often prove to be either inefficient or insufficient to accurately convey detailed controls upon the final generation results, e.g., control the fine-grained semantic layout of multiple objects, not to mention the challenge in understanding complex text descriptions for such

---

\*Corresponding Author, † Intern at Microsoft

Table 1: Comparisons of different controllable diffusion models.  $N$  is the number of conditions. We define the fine-tuning cost as the number of times the model needs to be fine-tuned on  $N$  conditions. As Composer is trained from scratch, both fine-tuning cost and adapter number are not applicable. For T2I-Adapter, (+1) indicates that further joint fine-tuning is required on the  $N$ -based adapters along with an additional fuser to achieve composable conditions.

	Fine-tuning	Composable Control	Fine-tuning Cost	Adapter Number
Composer	✗	✓	-	-
ControlNet	✓	✗	$N$	$N$
T2I-Adapter	✓	✓	$N(+1)$	$N(+1)$
Uni-ControlNet (Ours)	✓	✓	2	2

models. As a result, there is a growing need to incorporate more additional control modes (e.g., user-drawn sketch, semantic mask) alongside the text description into such T2I diffusion models. This necessity has sparked considerable interest from both academia and industry, as it broadens the scope of T2I generation from a singular function to a comprehensive system.

Very recently, there are some attempts [9–11] studying controllable T2I diffusion models. One representative work, Composer [9], explores the integration of multiple different control signals together with the text descriptions and train the model from scratch on billion-scale datasets. While the results are promising, it requires massive GPU resources and incurs huge training cost, making it unaffordable for many researchers in this field. Considering there are powerful pretrained T2I diffusion models (e.g., Stable Diffusion [7]) publicly available, ControlNet [10] and T2I-Adapter [11] directly incorporate lightweight adapters into frozen T2I diffusion models to enable additional condition signals. This makes fine-tuning more affordable. However, one drawback is that they need one independent adapter for each single condition, resulting in a linear increase in fine-tuning cost and model size along as the number of the control conditions grows, even though many conditions share similar characteristics. Additionally, this also makes composability among different conditions remains a formidable challenge.

In this paper, we propose Uni-ControlNet, a new framework that leverages lightweight adapters to enable precise controls over pre-trained T2I diffusion models. As shown in Table 1, unlike previous methods, Uni-ControlNet categorizes various conditions into two distinct groups: local conditions and global conditions. Accordingly, we only add two additional adapters, regardless of the number of local and global controls involved. This design choice not only significantly reduces both the whole fine-tuning cost and the model size, making it highly efficient for deployment, but also facilitates the composability of different conditions. To achieve this, we dedicatedly design the adapters for local and global controls. Specifically, for local controls, we introduce a multi-scale condition injection strategy that uses a shared local condition encoder adapter. This adapter first converts the local control signals into modulation signals, which are then used to modulate the incoming noise features. And for global controls, we employ another shared global condition encoder to convert them into conditional tokens, which are concatenated with text tokens to form the extended prompt and interacted with the incoming features via cross-attention mechanism. Interestingly, we find these two adapters can be separately trained without the need of additional joint training, while still supporting the composition of multiple control signals. This finding adds to the flexibility and ease of use provided by Uni-ControlNet.

By only training on 10 million text-image pairs with 1 epoch, our Uni-ControlNet demonstrates highly promising results in terms of fidelity and controllability. Figure 1 provides visual examples showcasing the effectiveness of Uni-ControlNet when using either one or multiple conditions. To gain further insights, we perform in-depth ablation analysis and compare our newly proposed adapter designs with those of ControlNet [10] and T2I-Adapter [11]. The analysis results reveal the superiority of our adapter designs, emphasizing their enhanced performance over the counterparts offered by ControlNet and T2I-Adapter.

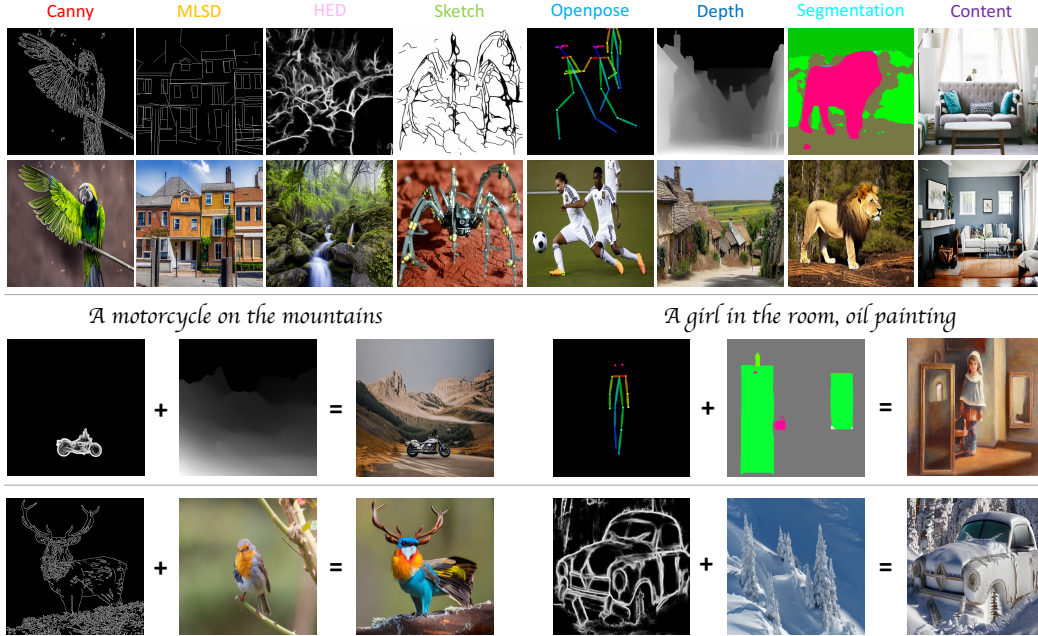


Figure 1: Visual results of our proposed Uni-ControlNet. The top and bottom two rows are results for single condition and multi-conditions respectively.

## 2 Related Work

**Text-to-Image Generation** is an emerging field that aims to generate realistic images from text descriptions. To address this challenging task, various approaches have been proposed in the past years. Early works [12–14] primarily adopted Generative Adversarial Networks (GANs) and were often trained on specific domains. However, they faced two key challenges, i.e., training instability and poor generalization ability to open-domain scenarios. Motivated by the success of GPT models [15–18], recent works [19–22] have explored the use of autoregressive models for text-to-image generation and train on web-scale image-text pairs, which start to show strong generation capability under the zero-shot setting for open-domain scenarios. Another approach is the diffusion models, originally proposed by [1, 2]. Diffusion models comprise a forward process that gradually adds noise to natural images and a backward process that learns to denoise them back to generate clean output. They demonstrate stronger capability in modeling fine-grained structures and texture details compared to autoregressive models. Recently, vast variants of diffusion models have been developed, such as DALLÉ-2 [5], which uses one prior model and one decoder model to generate images from CLIP latent embeddings. Another phenomenal T2I diffusion model is Stable Diffusion (SD), which scaled up the latent diffusion model [7] with larger model and data scales, and made the pre-trained models publicly available. In this paper, we use SD as a base model and explore how to enable more control signals beyond the text description for pre-trained T2I diffusion models in an efficient and composable way.

**Controllable Diffusion Models** are designed to enable T2I diffusion models to accept more user controls for guiding the generation results. They have garnered increasing attention very recently. Broadly speaking, there are two strategies for implementing controllable diffusion models: training from scratch [9] and fine-tuning lightweight adapters [10, 11] on frozen pretrained T2I diffusion models. In the case of training from scratch, Composer [9] trains one big diffusion model from scratch to achieve great controllability for both single and multi-conditions. It obtains remarkable generation quality but comes with huge training cost. In contrast, ControlNet [10] and T2I-Adapter [11] propose to introduce lightweight adapters into publicly available SD models. By only fine-tuning the adapters while keeping original SD models frozen, they significantly reduce the training cost and make it affordable for the research community. However, both ControlNet and T2I-Adapter utilize independent adapters for each condition, resulting in increased fine-tuning cost and model size when handling increased number of conditions. Moreover, different adapters in ControlNet are isolated

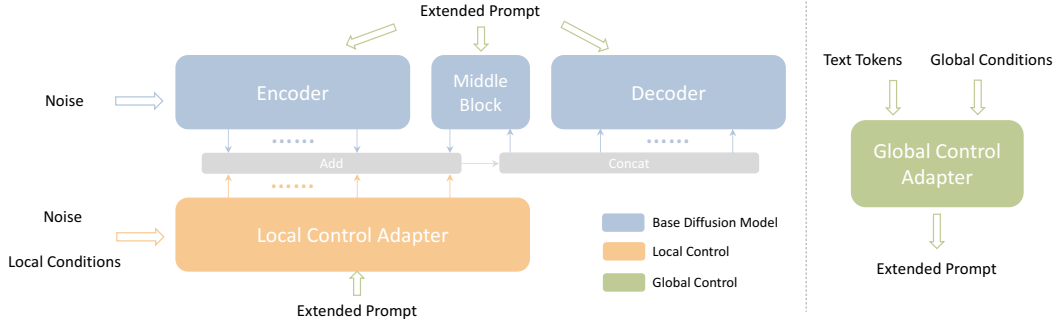


Figure 2: The overall framework of our proposed Uni-ControlNet.

from one another, limiting their composability. By testing CoAdapter [11], which is jointly trained using different T2I-Adapters, we find that it also exhibits inadequate performance in generating composable conditions. Our proposed Uni-ControlNet follows the second line of fine-tuning adapters and is much less expensive than Composer, while addressing the above limitations of ControlNet and T2I-Adapter. It groups conditions into two groups, i.e., local controls and global controls, and only requires two additional adapters accordingly. Thanks to our newly designed adapter structure, Uni-ControlNet is not only efficient in terms of training cost and model sizes, but also surpasses ControlNet and T2I-Adapter in controllability and quality.

### 3 Method

#### 3.1 Preliminary

A typical diffusion model involves two processes: a forward process which gradually adds small amounts of Gaussian noise onto the sample in  $T$  steps, and a corresponding backward process containing learnable parameters to recover input images by estimating and eliminating the noise. In this paper, we use SD as our example base model to illustrate how to enable diverse controls with our Uni-ControlNet. SD incorporates the UNet-like structure [23] as its denoising model, which consists of an encoder, a middle block, and a decoder, with 12 corresponding blocks in each of the encoder and decoder modules. For brevity, we denote the encoder as  $F$ , the middle block as  $M$ , and the decoder as  $G$ , with  $f_i$  and  $g_i$  denoting the output of the  $i$ -th block in the encoder and decoder, and  $m$  denoting the output of the middle block, respectively. It is important to note that, due to the adoption of skip connections in UNet, the input for the  $i$ -th block in the decoder is given by:

$$\begin{cases} \text{concat}(m, f_j) & \text{where } i = 1, \quad i + j = 13. \\ \text{concat}(g_{i-1}, f_j) & \text{where } 2 \leq i \leq 12, \quad i + j = 13. \end{cases} \quad (1)$$

Skip connections allow the decoder to directly utilize features from the encoder and thereby help minimize the information loss. In SD, cross-attention layers are employed to capture semantic information from the input text description. Here we use  $Z$  to denote the incoming noise features and  $y$  to denote text token embeddings encoded by the language encoder. The  $Q, K, V$  in cross-attention can be expressed as:

$$Q = W_q(Z), K = W_k(y), V = W_v(y), \quad (2)$$

where  $W_q, W_k$  and  $W_v$  are projection matrices.

#### 3.2 Control Adapter

In this paper, we consider seven example local conditions, including Canny edge [24], MLSD edge [25], HED boundary [26], sketch [27, 28], Openpose [29], Midas depth [30], and segmentation mask [31]. We also consider one example global condition, i.e., global image embedding of one reference content image that is extracted from the CLIP image encoder [32]. This global condition goes beyond simple image features and provides a more nuanced understanding of the semantic content of the condition image. By employing both local and global conditions, we aim to provide a comprehensive control over the generation process. We show the overview of our pipeline in Figure 2, and the details of local control adapter and global control adapter are given in Figure 3.



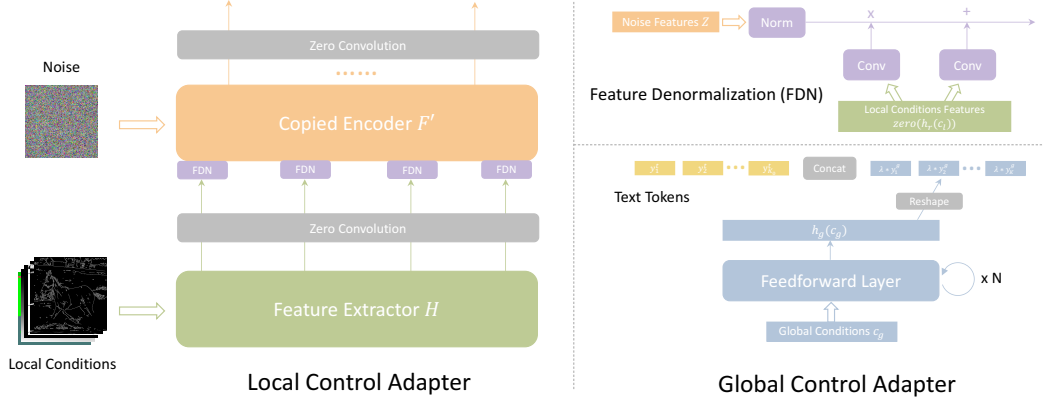


Figure 3: Details of the local and global control adapters.

**Local Control Adapter:** For our local control adapter, we have taken inspiration from ControlNet. Specifically, we fix the weights of SD and copy the structures and weights of the encoder and middle block, designated as  $F'$  and  $M'$  respectively. Thereafter, we incorporate the information from the local control adapter during the decoding process. To achieve it, we ensure that all other elements remain unchanged while modifying the input of the  $i$ -th block of the decoder as

$$\begin{cases} \text{concat}(m + m', f_j + \text{zero}(f'_j)) & \text{where } i = 1, \quad i + j = 13. \\ \text{concat}(g_{i-1}, f_j + \text{zero}(f'_j)) & \text{where } 2 \leq i \leq 12, \quad i + j = 13. \end{cases} \quad (3)$$

where  $\text{zero}$  represents one zero convolutional layer whose weights increase from zero to gradually integrate control information into the main SD model. In contrast to ControlNet that adds the conditions directly to the input noise and sends them to the copied encoder, we opt for a multi-scale condition injection strategy. Our approach involves injecting the condition information at all resolutions. In detail, we first concatenate different local conditions along the channel dimension and then use a feature extractor  $H$  (stacked convolutional layers) to extract condition features at different resolutions. Subsequently, we select the first block of each resolution (i.e.,  $64 \times 64, 32 \times 32, 16 \times 16, 8 \times 8$ ) in the copied encoder (i.e., the Copied Encoder in Figure 3) for condition injection. For the injection module, we take the inspiration from SPADE [33] and implement Feature Denormalization (FDN) that uses the condition features to modulate the normalized (i.e.,  $\text{norm}(\cdot)$ ) input noise features:

$$\text{FDN}_r(Z_r, c_l) = \text{norm}(Z_r) \cdot (1 + \text{conv}_\gamma(\text{zero}(h_r(c_l)))) + \text{conv}_\beta(\text{zero}(h_r(c_l))), \quad (4)$$

where  $Z_r$  denotes noise features at resolution  $r$ ,  $c_l$  is the concatenated local conditions,  $h_r$  represents the output of the feature extractor  $H$  at resolution  $r$ , and  $\text{conv}_\gamma$  and  $\text{conv}_\beta$  refer to learnable convolutional layers that convert condition features into spatial-sensitive scale and shift modulation coefficients. We will ablate different local feature injection strategies in following sections.

**Global Control Adapter:** For global controls, we use the image embedding of one condition image extracted from CLIP image encoder as the example. Inspired by the fact that the text description in T2I diffusion models can be also viewed as one kind of global control without explicit spatial guidance, we project the global control signals into condition embeddings by using a condition encoder  $h_g$ . The condition encoder consists of stacked feedforward layers, which aligns the global control signals with the text embeddings in SD. Next, we reshape the projected condition embeddings into  $K$  global tokens ( $K = 4$  by default) and concatenate them with the original  $K_0$  text tokens to create an extended prompt  $y_{ext}$  (total token number is  $K + K_0$ ) which serves as the input to all cross-attention layers in both main SD model and control adapters:

$$y_{ext} = [y_1^t, y_2^t, \dots, y_{K_0}^t, \lambda * y_1^g, \lambda * y_2^g, \dots, \lambda * y_K^g], \text{ where } y_i^g = h_g(c_g)[(i-1) \cdot d \sim i \cdot d], i \in [1, K] \quad (5)$$

where  $y^t$  and  $y^g$  represent the original text tokens and global condition tokens respectively, and  $\lambda$  is a hyper-parameter that controls the weight of the global condition.  $c_g$  denotes the global condition and  $d$  is the dimension of text token embedding.  $h_g(\cdot)[i_s \sim i_e]$  represents the sub-tensor of  $h_g(\cdot)$  that contains elements from the  $i_s$ -th to the  $i_e$ -th positions. Finally, the  $Q, K, V$  cross-attention operation in all cross-attention layers is changed to:

$$Q = W_q(Z), K = W_k(y_{ext}), V = W_v(y_{ext}), \quad (6)$$

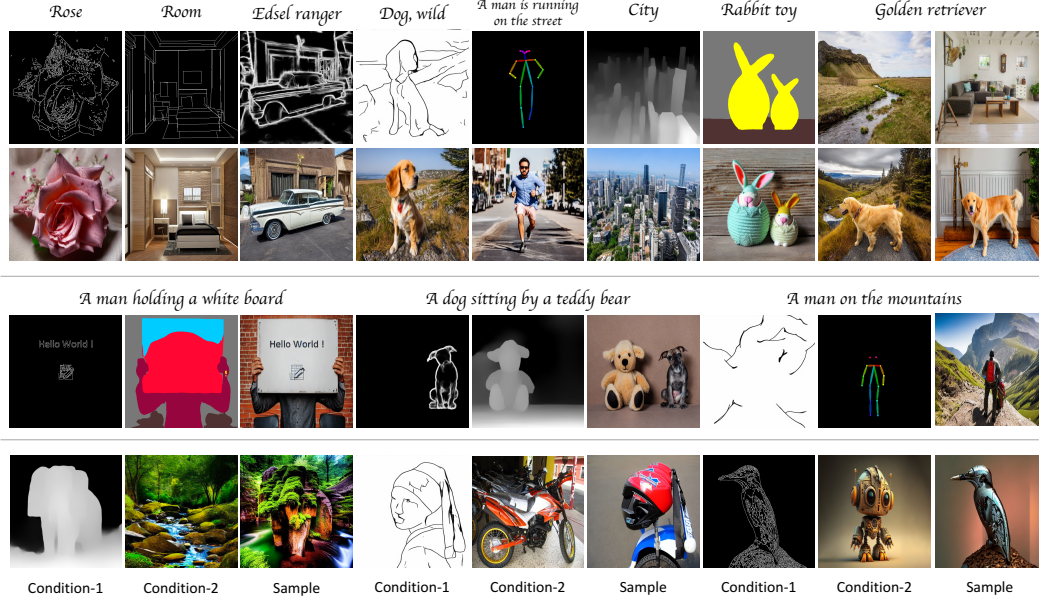


Figure 4: More visual results of Uni-ControlNet. The top two rows show results of a single condition, with columns 1-7 for local conditions and columns 8-9 for global condition. 3rd row shows the results of combining two local conditions, while row 4-th shows the results of integrating a local condition with a global condition. There is no text prompt for the examples in 4-th row.

### 3.3 Training Strategy

As the local control signals and global control signals often contain different amounts of condition information, we empirically find that directly joint fine-tuning these two types of adapters will produce poor controllable generation performance. Therefore, we opt to fine-tune these two types of adapters separately so that both of them can be sufficiently trained and contribute effectively to the final generation results. When fine-tuning each adapter, we employ a predefined probability to randomly dropout each condition, along with an additional probability to deliberately keep or drop all conditions. This can facilitate the model to learn generating the results based on one or multiple conditions simultaneously. Interestingly, by directly integrating these two separately trained adapters during inference, our Uni-ControlNet can already well combine global and local conditions together in a composable way, without the need of further joint fine-tuning. In Section 4.3, we will provide more detailed analysis about different fine-tuning strategies.

## 4 Experiments

**Implementation Details.** To fine-tune our model, we randomly sample 10 million text-image pairs from the LAION dataset [34] and fine-tune Uni-ControlNet for 1 epoch. We use the AdamW optimizer [35] with a learning rate of  $1 \times 10^{-5}$  and resize the input images and local condition maps to  $512 \times 512$ . As described, the local and global control adapters are fine-tuned separately by default. During inference, we merge the two adapters and adopt DDIM [36] for sampling, with the number of time steps set to 50 and the classifier free guidance scale [37] set to 7.5. During training, the hyper-parameter  $\lambda$  in Equation 6 is with a fixed value 1. At inference time, when there is no text prompt,  $\lambda$  remains at 1, while when there is a text prompt, the value is adjusted to around 0.75, depending on the intended weight between the text and global condition. As explained in Section 3.2, we employ 7 local conditions (Canny edge, MLSD edge, HED boundary, sketch, Openpose, Midas depth, and segmentation mask) and 1 global control condition (CLIP image embeddings) for control. Detailed structures of global and local condition adapters can be found in supplementary material.

### 4.1 Controllable Generation Results

In Figure 4, we provide more controllable generation results of Uni-ControlNet in both single and multi-condition setups. Notably, for visualization purposes, we use the original condition images

Table 2: FID on different controllable diffusion models. The best results are in **bold**.

	Canny	MLSD	HED	Sketch	Openpose	Depth	Segmentation	Style\Content
ControlNet	18.90	31.36	26.59	22.19	27.84	21.25	<b>23.08</b>	31.17
T2I-Adapter	18.98	-	-	<b>18.83</b>	29.57	21.35	23.84	28.86
Ours	<b>17.79</b>	<b>26.18</b>	<b>17.86</b>	20.11	<b>26.61</b>	<b>21.20</b>	23.40	<b>23.98</b>

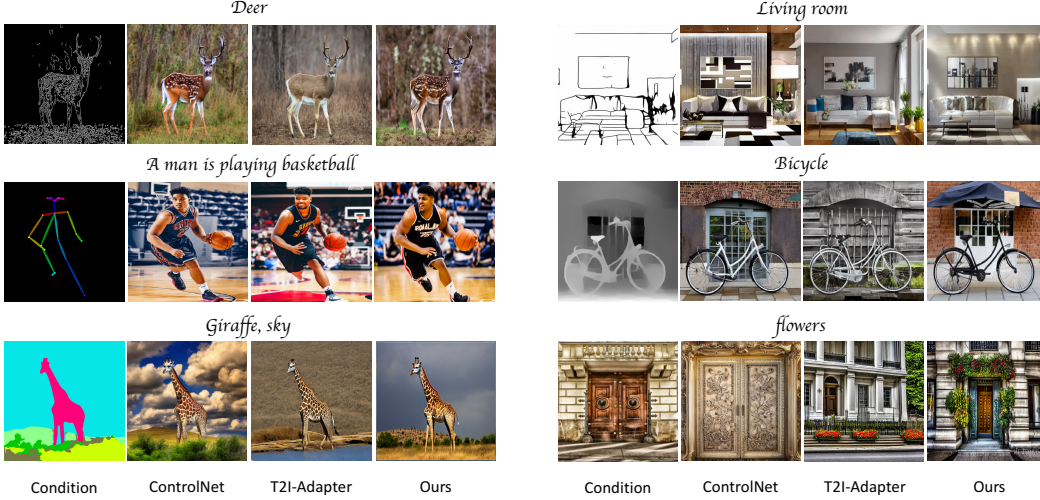


Figure 5: Comparison of existing controllable diffusion models on different single conditions.

to denote their CLIP image embeddings. It can be seen that our Uni-ControlNet can produce very promising results in terms of both controllability and generation fidelity. For example, in the case of a single sketch condition with the text prompt "Dog, wild" (rows 1-2, column 4), the resulting image accurately depicts a vivid dog and a background of grass and trees that align well with the given sketch condition. Similarly, when presented with the global CLIP image embedding conditions with the prompt "Golden retriever" (rows 1-2, columns 8-9), our model can seamlessly change the background of the dog from the wild to a room. Moreover, our model also handles multi-condition settings well, as demonstrated in the example of "A man on the mountains" (row 3, columns 7-9), where the combination of a sketch and a pose produces a cohesive and detailed image of a man on a mountainside. When presented with a local depth map and global CLIP image embeddings without any prompt (row 4, columns 1-3), our model produces an image of a forest, taking the contour of an elephant, which harmonizes with both the depth map and the content of the source global image.

## 4.2 Comparison with Existing Methods

Here we compare our Uni-ControlNet with ControlNet [10] and T2I-Adapter (CoAdapter) [11]. Since Composer [9] is not open-sourced and trained from scratch, we do not include it in comparisons.

**Quantitative Comparison:** For quantitative evaluation, we use the validation set of COCO2017 [38] at a resolution of  $512 \times 512$ . Since this set contains 5k images, and each image has multiple captions, we randomly select one caption per image resulting in 5k generated images for our evaluation. We report the FID [39] in Table 2. It is important to note that for quantitative comparison, we limit our testing to different single conditions only. Additionally, we use Style\Content to represent the global condition as there are different settings in the ControlNet and T2I-Adapter. As T2I-Adapter does not take the MLSD and HED conditions into account, it has no results for MLSD and HED. Our model reveals superior performance across most conditions quantitatively compared to existing approaches.

**Qualitative Comparison:** We further provide qualitative comparison of single and composed multi-conditions in Figure 5 and Figure 6 respectively. For single conditions, we find that our Uni-ControlNet, ControlNet and T2I-Adapter can all perform overall well, and our results show slightly better alignments with input conditions. Notably, we only fine-tune 2 adapters for all conditions, whereas ControlNet and T2I-Adapter fine-tune eight adapters for eight different single conditions. Since ControlNet does not support composed multi-conditions, we only compare Uni-ControlNet with CoAdapter, which conducts further joint fine-tuning based on the base T2I-Adapters and an

Table 3: FID on condition injection methods and training strategies. The best results are in **bold**.

	Canny	MLSD	HED	Sketch	Openpose	Depth	Segmentation	Content
Injection-S1	25.89	27.22	22.48	23.51	27.89	24.71	26.25	-
Injection-S2	22.22	27.08	21.94	22.74	<b>26.56</b>	24.21	24.42	-
Injection-S3	-	-	-	-	-	-	-	27.06
Training-S1	21.21	27.20	20.78	23.22	27.83	25.01	24.99	28.51
Training-S2	18.80	<b>26.40</b>	19.12	20.91	27.17	<b>21.59</b>	23.93	<b>24.84</b>
Ours	<b>18.24</b>	26.91	<b>18.61</b>	<b>20.32</b>	27.76	21.97	<b>23.51</b>	24.86

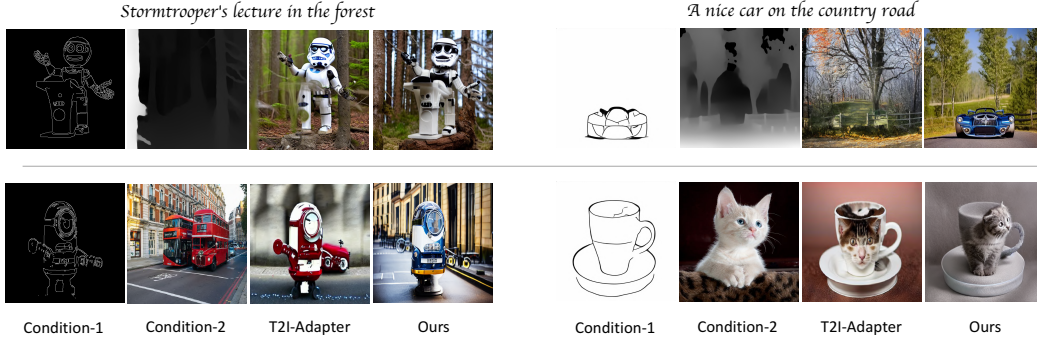


Figure 6: Comparison of different controllable diffusion models on composable multi-conditions.

extra Fuser module. As shown in Figure 6, CoAdapter shows poorer composability when dealing with two local conditions, e.g., missing the podium in the first example and no car in the second example. In contrast, our model can fuse the two conditions much better. As for composing a local condition with a global condition, CoAdapter performs okay, like the case of combining a sketch of a cup and a global condition of a cat. However, when the two conditions are not that related, e.g., the example where there is a Canny edge of a Minion and a global condition of a bus in London, the image generated by CoAdapter appears to be unrealistic and the two elements are not well integrated. And our model effectively creates a Minion-shaped bus with car windows and vivid background.

### 4.3 Ablation Analysis

For ablation study, we fine-tune our model using a smaller dataset for resource consideration. In detail, we utilize the 1 million subset of the 10 million dataset and fine-tune a single epoch, while keeping all other settings unchanged.

**Condition Injection Strategy:** For local conditions, we compare our proposed injection method with two other strategies. The first strategy is to directly use SPADE to inject the conditions, which involves resizing the conditions to the corresponding resolutions using interpolation. We call this Injection-S1. The second strategy is similar to Composer, ControlNet and T2I-Adapter, where the conditions are only sent to the adapter or the main model at the input layer, which we refer to as Injection-S2. When using these two strategies, all other parts of our method will remain unchanged. We follow the setting in Section 4.2 and evaluate the FID on different local condition injection strategies. The quantitative and qualitative results are presented in Table 3 and the upper part of Figure 7. The quantitative results show our proposed condition injection strategy performs better under most settings. For the visual results, we observe that for Injection-S1, the alignment with the conditions is poor. This may be because direct interpolation significantly destroys condition information. As for Injection-S2, it renders unsatisfactory results for composite control. For instance, in the "An elephant in the temple" case, the lanterns on the top of the image are not accurately aligned with depth condition. Moreover, the composite results are not as harmonious as ours. For example, in the "Gorilla wearing glasses" case, the gorilla's eyes and glasses are not well-merged. For the global condition, we compare our method to one way in which we only add the global condition into control adapter but not the main SD model, and we denote this injection strategy as Injection-S3. As shown in the lower part of Figure 7, without using the extended prompt in the main SD model, this method cannot inject the global condition into the final generated results.



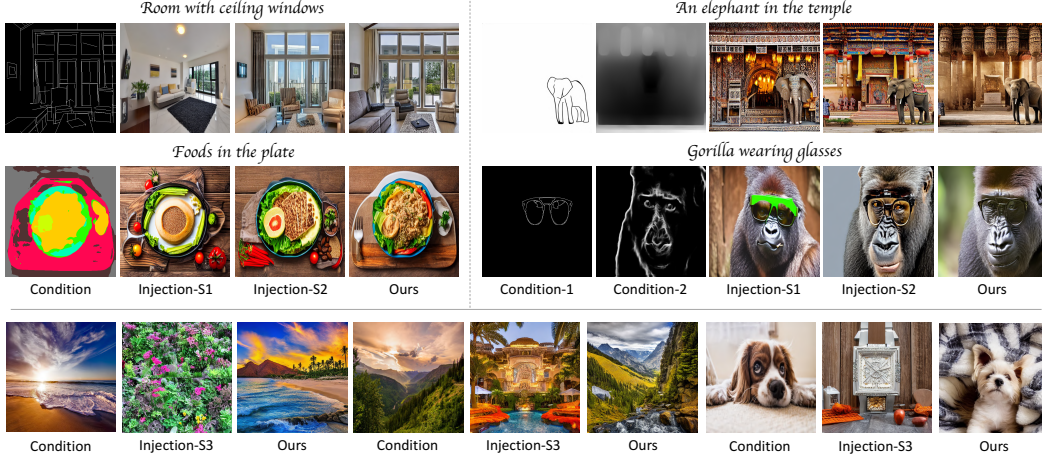


Figure 7: Ablation results on different condition injection strategies.

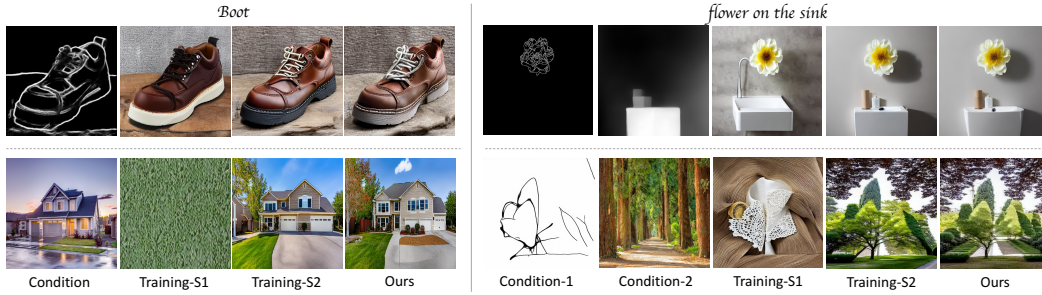


Figure 8: Ablation results on different training strategies.

**Training Strategy:** As above mentioned, we fine-tune the local and global control adapters separately and merge them at inference without any further joint fine-tuning by default. Here, we also investigate two alternative training strategies: 1) joint fine-tuning together (“Train-S1”), where we fine-tune both adapters together from scratch; 2) further joint fine-tuning after separate fine-tuning (“Train-S2”), where we further fine-tune the adapters together after separate fine-tuning. The quantitative FID results are shown in Table 3. We find that our default strategy and Training-S2 perform much better consistently than Training-S1, but further joint fine-tuning in Train-S2 does not bring obvious performance gain in most cases. Some visual results are given in Figure 8. Note that, in order to better assess the controllability of the global condition, we do not provide text prompts for the cases with global condition. As we described before, the reason why Training-S1 gets poor controllability on the global condition is that the global control adapter does not learn as much as local adapter even equally treated during joint fine-tuning. One possible explanation is that the local conditions often contain more rich guidance information than global conditions, leading the model to pay less attention to the global condition.

## 5 Conclusion and Social Impact

In this paper, we propose Uni-ControlNet, a new solution that enhances the capabilities of text-to-image diffusion models by enabling efficient integration of diverse local and global controls. With better adapter designs, our Uni-ControlNet only requires two adapters for different conditions while existing methods often require independent adapters for each condition. The new design of Uni-ControlNet not only saves both fine-tuning cost and model size, but also facilitates composability, allowing for the simultaneous utilization of multiple conditions. Extensive experiments validate the effectiveness of Uni-ControlNet, showcasing its improved controllability, generation fidelity, and composability. While our system empowers artists, designers, and content creators to realize their creative visions with precise control, it is crucial to acknowledge the potential negative social impact that can arise from misuse or abuse, similar to other image generation and editing AI models. To address these concerns, responsible deployment practices, ethical regulations, and the inclusion of special flags in generated images to enhance transparency are vital steps towards responsible usage.

## References

- [1] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in Neural Information Processing Systems*, 33:6840–6851, 2020.
- [2] Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. *arXiv preprint arXiv:2010.02502*, 2020.
- [3] Florinel-Alin Croitoru, Vlad Hondru, Radu Tudor Ionescu, and Mubarak Shah. Diffusion models in vision: A survey. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2023.
- [4] Alex Nichol, Prafulla Dhariwal, Aditya Ramesh, Pranav Shyam, Pamela Mishkin, Bob McGrew, Ilya Sutskever, and Mark Chen. Glide: Towards photorealistic image generation and editing with text-guided diffusion models. *arXiv preprint arXiv:2112.10741*, 2021.
- [5] Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. Hierarchical text-conditional image generation with clip latents. *arXiv preprint arXiv:2204.06125*, 2022.
- [6] Shuyang Gu, Dong Chen, Jianmin Bao, Fang Wen, Bo Zhang, Dongdong Chen, Lu Yuan, and Baining Guo. Vector quantized diffusion model for text-to-image synthesis. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR 2022)*, 2022.
- [7] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10684–10695, 2022.
- [8] Wan-Cyuan Fan, Yen-Chun Chen, DongDong Chen, Yu Cheng, Lu Yuan, and Yu-Chiang Frank Wang. Frido: Feature pyramid diffusion for complex scene image synthesis. In *Thirty-Seventh AAAI Conference on Artificial Intelligence (AAAI)*, 2023.
- [9] Lianghua Huang, Di Chen, Yu Liu, Yujun Shen, Deli Zhao, and Jingren Zhou. Composer: Creative and controllable image synthesis with composable conditions. *arXiv preprint arXiv:2302.09778*, 2023.
- [10] Lvmin Zhang and Maneesh Agrawala. Adding conditional control to text-to-image diffusion models. *arXiv preprint arXiv:2302.05543*, 2023.
- [11] Chong Mou, Xintao Wang, Liangbin Xie, Jian Zhang, Zhongang Qi, Ying Shan, and Xiaohu Qie. T2i-adapter: Learning adapters to dig out more controllable ability for text-to-image diffusion models. *arXiv preprint arXiv:2302.08453*, 2023.
- [12] Scott Reed, Zeynep Akata, Xinchun Yan, Lajanugen Logeswaran, Bernt Schiele, and Honglak Lee. Generative adversarial text to image synthesis. In *ICML 2016*.
- [13] Han Zhang, Tao Xu, Hongsheng Li, Shaoting Zhang, Xiaogang Wang, Xiao lei Huang, and Dimitris N Metaxas. Stackgan: Text to photo-realistic image synthesis with stacked generative adversarial networks. In *ICCV 2017*.
- [14] Tao Xu, Pengchuan Zhang, Qiuyuan Huang, Han Zhang, Zhe Gan, Xiao lei Huang, and Xiaodong He. Attngan: Fine-grained text to image generation with attentional generative adversarial networks. In *CVPR 2018*.
- [15] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020.
- [16] Mark Chen, Alec Radford, Rewon Child, Jeffrey Wu, Heewoo Jun, David Luan, and Ilya Sutskever. Generative pretraining from pixels. In *International conference on machine learning*, pages 1691–1703. PMLR, 2020.
- [17] Ziyu Wan, Jingbo Zhang, Dongdong Chen, and Jing Liao. High-fidelity pluralistic image completion with transformers. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 4692–4701, 2021.

- [18] Qiankun Liu, Zhentao Tan, Dongdong Chen, Qi Chu, Xiyang Dai, Yinpeng Chen, Mengchen Liu, Lu Yuan, and Nenghai Yu. Reduce information loss in transformers for pluralistic image inpainting. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR 2022)*, 2022.
- [19] Aditya Ramesh, Mikhail Pavlov, Gabriel Goh, Scott Gray, Chelsea Voss, Alec Radford, Mark Chen, and Ilya Sutskever. Zero-shot text-to-image generation. In *ICML 2021*.
- [20] Ming Ding, Zhuoyi Yang, Wenyi Hong, Wendi Zheng, Chang Zhou, Da Yin, Junyang Lin, Xu Zou, Zhou Shao, Hongxia Yang, et al. Cogview: Mastering text-to-image generation via transformers. In *NeurIPS 2021*.
- [21] Junyang Lin, Rui Men, An Yang, Chang Zhou, Ming Ding, Yichang Zhang, Peng Wang, Ang Wang, Le Jiang, Xianyan Jia, et al. M6: A chinese multimodal pretrainer. *arXiv preprint arXiv:2103.00823*, 2021.
- [22] Jiahui Yu, Yuanzhong Xu, Jing Yu Koh, Thang Luong, Gunjan Baid, Zirui Wang, Vijay Vasudevan, Alexander Ku, Yinfei Yang, Burcu Karagol Ayan, et al. Scaling autoregressive models for content-rich text-to-image generation. *arXiv preprint arXiv:2206.10789*, 2022.
- [23] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical image segmentation. In *MICCAI 2015*.
- [24] John Canny. A computational approach to edge detection. *TPAMI 1986*.
- [25] Geonmo Gu, Byungsoo Ko, SeoungHyun Go, Sung-Hyun Lee, Jingeun Lee, and Minchul Shin. Towards light-weight and real-time line segment detection. In *AAAI 2022*.
- [26] Saining Xie and Zhuowen Tu. Holistically-nested edge detection. In *ICCV 2015*.
- [27] Edgar Simo-Serra, Satoshi Iizuka, Kazuma Sasaki, and Hiroshi Ishikawa. Learning to simplify: Fully convolutional networks for rough sketch cleanup. In *SIGGRAPH*, 2016.
- [28] Edgar Simo-Serra, Satoshi Iizuka, and Hiroshi Ishikawa. Mastering sketching: Adversarial augmentation for structured prediction. In *TOG 2018*.
- [29] Zhe Cao, Tomas Simon, Shih-En Wei, and Yaser Sheikh. Realtime multi-person 2d pose estimation using part affinity fields. In *CVPR 2017*.
- [30] René Ranftl, Katrin Lasinger, David Hafner, Konrad Schindler, and Vladlen Koltun. Towards robust monocular depth estimation: Mixing datasets for zero-shot cross-dataset transfer. *TPAMI 2020*.
- [31] Tete Xiao, Yingcheng Liu, Bolei Zhou, Yuning Jiang, and Jian Sun. Unified perceptual parsing for scene understanding. In *ECCV 2018*.
- [32] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *ICML 2021*.
- [33] Taesung Park, Ming-Yu Liu, Ting-Chun Wang, and Jun-Yan Zhu. Semantic image synthesis with spatially-adaptive normalization. In *CVPR 2019*.
- [34] Christoph Schuhmann, Richard Vencu, Romain Beaumont, Robert Kaczmarczyk, Clayton Mullis, Aarush Katta, Theo Coombes, Jenia Jitsev, and Aran Komatsuzaki. Laion-400m: Open dataset of clip-filtered 400 million image-text pairs. *arXiv preprint arXiv:2111.02114*, 2021.
- [35] Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *ICLR 2015*.
- [36] Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. In *ICLR 2021*.
- [37] Jonathan Ho and Tim Salimans. Classifier-free diffusion guidance. In *NeurIPS 2021 Workshop on Deep Generative Models and Downstream Applications*.



- [38] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *ECCV 2014*.
- [39] Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter. Gans trained by a two time-scale update rule converge to a local nash equilibrium. In *NeurIPS 2017*.

## A The Weight of the Global Condition

In the global control module, we have implemented a hyper-parameter  $\lambda$ . This hyper-parameter plays a role in determining the influence of the global condition while concatenating the projected global condition to the text. Illustrating the effect of varying  $\lambda$ , we present two representative visualization results in Figure 9. It can be seen that the value of  $\lambda$  plays a significant role in displaying the elements of the global conditions in the generated images. As the value of  $\lambda$  increases, the global condition takes precedence over the original text content, leading to a decrease the influence of the text prompt on the results. For instance, in the first case, it appears that the forest has experienced a reduction in coverage area as compared to an increase in the number of houses with an increase in the value of  $\lambda$ . Similarly, in the second case, it is evident that the city is shrinking while the desert’s coverage is expanding with the rise of  $\lambda$  value. In the real-world applications, we can adjust the hyper-parameter  $\lambda$  to generate our desired results with flexibility.

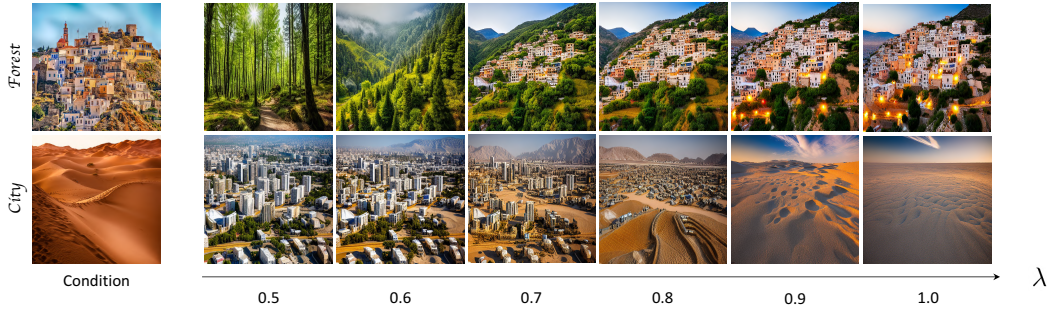


Figure 9: The effect of the hyper-parameter  $\lambda$ . On the left side are the textual prompts and global conditions provided. On the right side are the images generated under increased  $\lambda$  values.

## B Condition Conflicts

Since our Uni-ControlNet can support multiple conditions simultaneously, we are curious about its behavior when providing multiple conflicting conditions. Need to that, this is very rare in the real-world applications, and this experiment is just for the analysis purpose. For example, we consider the case of providing the model with two local conditions that are fundamentally incompatible, such as the conditions of two dogs shown in Figure 10. Through this experiment, we can possibly evaluate the relative importance of each condition and learn how Uni-ControlNet resolves conflicts, which may help us design more robust integration of conditions that can adequately handle diverse and ambiguous situations.

To provide a comprehensive analysis of different condition compositions, we have assigned each condition in the first column a number 1 and each condition in the first row a number 2. This allows us to refer to the dog in the first row as dog-1 and the dog in the first line as dog-2. The results depicted in Figure 10 demonstrate that HED is the most powerful condition, with generated images closely following the HED boundary when depicting text. Other conditions can only influence areas that do not overlap. For instance, when we combine the HED boundary of dog-2 with the Canny edge map of dog-1, the resulting image adopts the HED boundary of dog-2 but recognizes the head of dog-1 as a small element positioned near the head of dog-2. Similarly, when dog-2’s HED boundary is combined with the sketch of dog-1, the model fails to identify the head of dog-1 even though it does not conflict with the HED boundary. Among the other conditions, the Canny edge map is the next most powerful, followed by the sketch, depth, MLSD, and segmentation map. The Openpose condition is the weakest, whereby the model generally disregards it in the event of a conflict. Only when combined with the segmentation map, the Openpose condition produces recognizable human elements. For better visualization, we have reordered the conditions based on their strength, which implies that the upper and left conditions have greater influence.

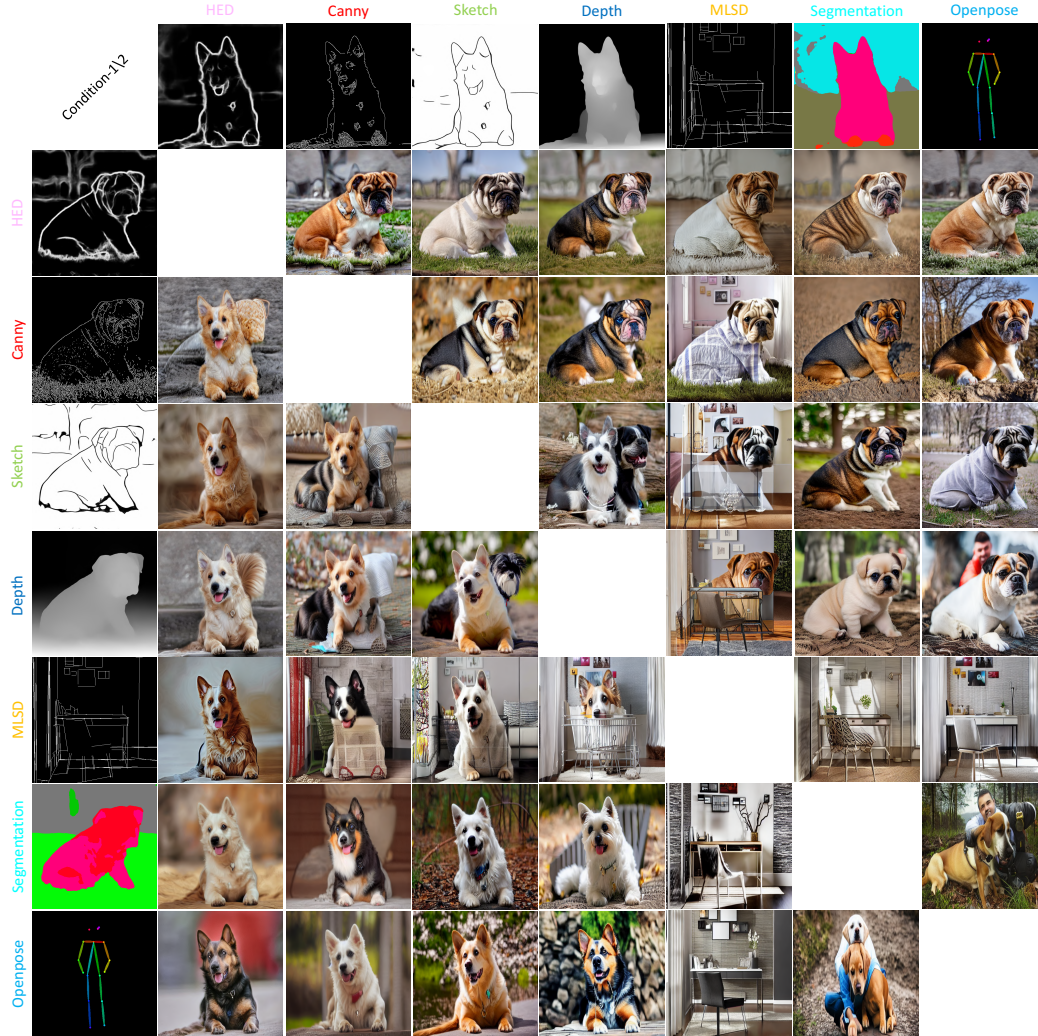


Figure 10: Study for the cases where composed conditions contradict each other. We choose the detection results of a room as the condition for MLSD and a man as the condition for Openpose. To test other conditions, we select two different dogs, which allows us to observe the model's output when given different dog-shaped conditions. We use "room" as the prompt for MLSD, "man" for Openpose, and "dog" for other conditions. When combining two types of conditions, we integrate their prompts, such as "dog", "dog and room", and "dog and man".

## C More Quantitative Results

**CLIP Score:** Besides FID, we also test CLIP scores for comparing different controllable diffusion models, and ablating condition injections methods and training strategies. We follow the settings in the Section 4.2 and Section 4.3 in the main paper. The results are shown in the Table 4 and Table C respectively. Our model demonstrates superior performance quantitatively across most conditions when compared to existing controllable diffusion models. Moreover, for different condition injection methods and training strategies, our method, along with Training-S2, consistently outperforms other strategies. However, joint fine-tuning in Training-S2 does not yield obvious performance gains in most cases. These finds are consistent with those presented in the Section 4.2 and Section 4.3 of the main paper.

Table 4: CLIP score on different controllable diffusion models. The best results are in bold.

	Canny	MLSD	HED	Sketch	Openpose	Depth	Segmentation	Style\Content
ControlNet	0.2538	0.2481	0.2530	0.2499	0.2572	0.2558	0.2531	0.2352
T2I-Adapter	0.2513	-	-	<b>0.2584</b>	<b>0.2608</b>	0.2559	0.2478	0.2366
Ours	<b>0.2539</b>	<b>0.2485</b>	<b>0.2556</b>	0.2542	0.2514	<b>0.2561</b>	<b>0.2540</b>	<b>0.2402</b>

Table 5: CLIP score on condition injection methods and training strategies. The best results are in bold.

	Canny	MLSD	HED	Sketch	Openpose	Depth	Segmentation	Content
Injection-S1	0.2513	0.2497	0.2518	0.2507	0.2527	0.2525	0.2508	-
Injection-S2	0.2504	<b>0.2506</b>	0.2518	0.2523	0.2527	0.2544	0.2540	-
Injection-S3	-	-	-	-	-	-	-	<b>0.2502</b>
Training-S1	0.2506	0.2504	0.2511	0.2510	0.2529	0.2538	0.2526	0.2478
Training-S2	<b>0.2528</b>	0.2504	0.2530	0.2537	<b>0.2533</b>	0.2547	<b>0.2545</b>	0.2421
Ours	<b>0.2528</b>	0.2483	<b>0.2535</b>	<b>0.2539</b>	0.2503	<b>0.2549</b>	0.2522	0.2420

**User Study:** As FID and CLIP score may be not always consistent with human preference, we further conduct user study to quantitatively compare our approach with the baseline methods ControlNet [10] and T2I-Adapter [11]. More specifically, we carry out tests in both single and multi-condition settings, with 20 cases for each setting. Each case is evaluated based on three metrics: the quality of generated images, the match with the given text, and the alignment with the given conditions. Users should select the best one for each metric from the generated images of ControlNet, T2I-Adapter, and our Uni-ControlNet. We collect responses from 20 users and analyze the total number of votes for each metric under each setting.

The results are presented in Figures 11 and 12. It can be seen that, our approach outperforms both ControlNet and T2I-Adapter in the single condition setting, demonstrating a clear advantage. Additionally, in the multi-conditions setting, our approach performed significantly better than CoAdapter.

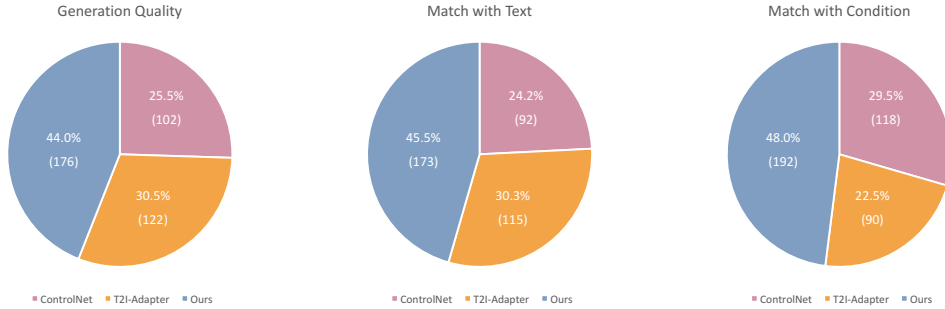


Figure 11: User study of the preference rate for the single condition setting.

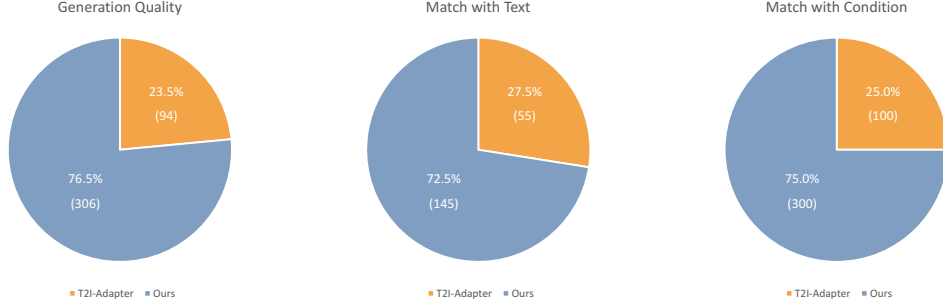


Figure 12: User study of the preference rate for the multi-conditions setting.

## D More Visualization Results

In this section, we present additional qualitative results. Figure 13 illustrates the results for the single-condition setting, while Figure 14 shows the results for the multi-conditions setting. Moreover, we demonstrate our performance on a more challenging case where there are four conditions, as seen in rows 7-8 of Figure 14.

## E Adapter Details

We provide the details of our proposed local control adapter and global control adapter in Figure 15, Figure 16 and Figure 17.



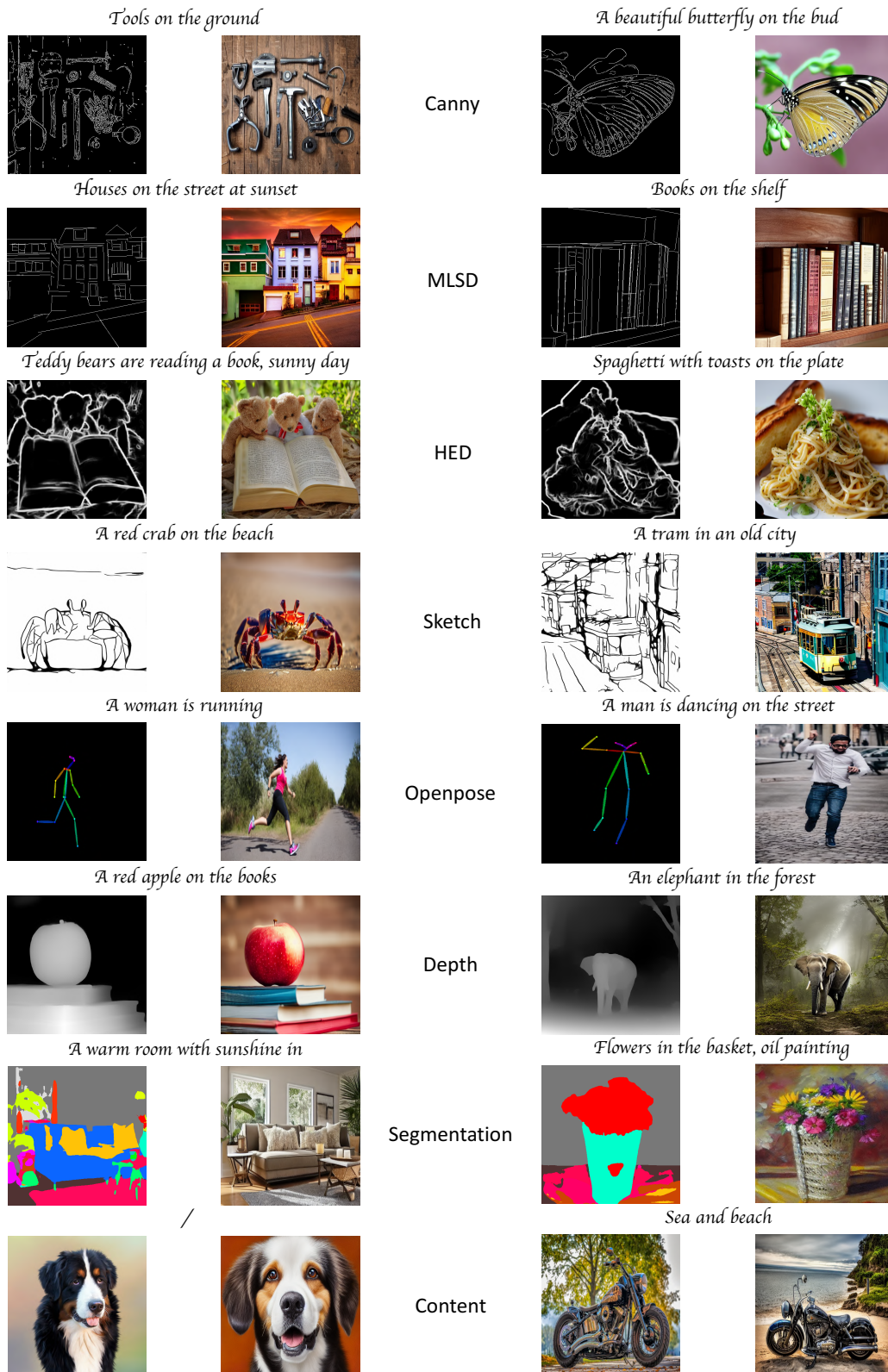


Figure 13: More visual results of Uni-ControlNet for single condition setting.

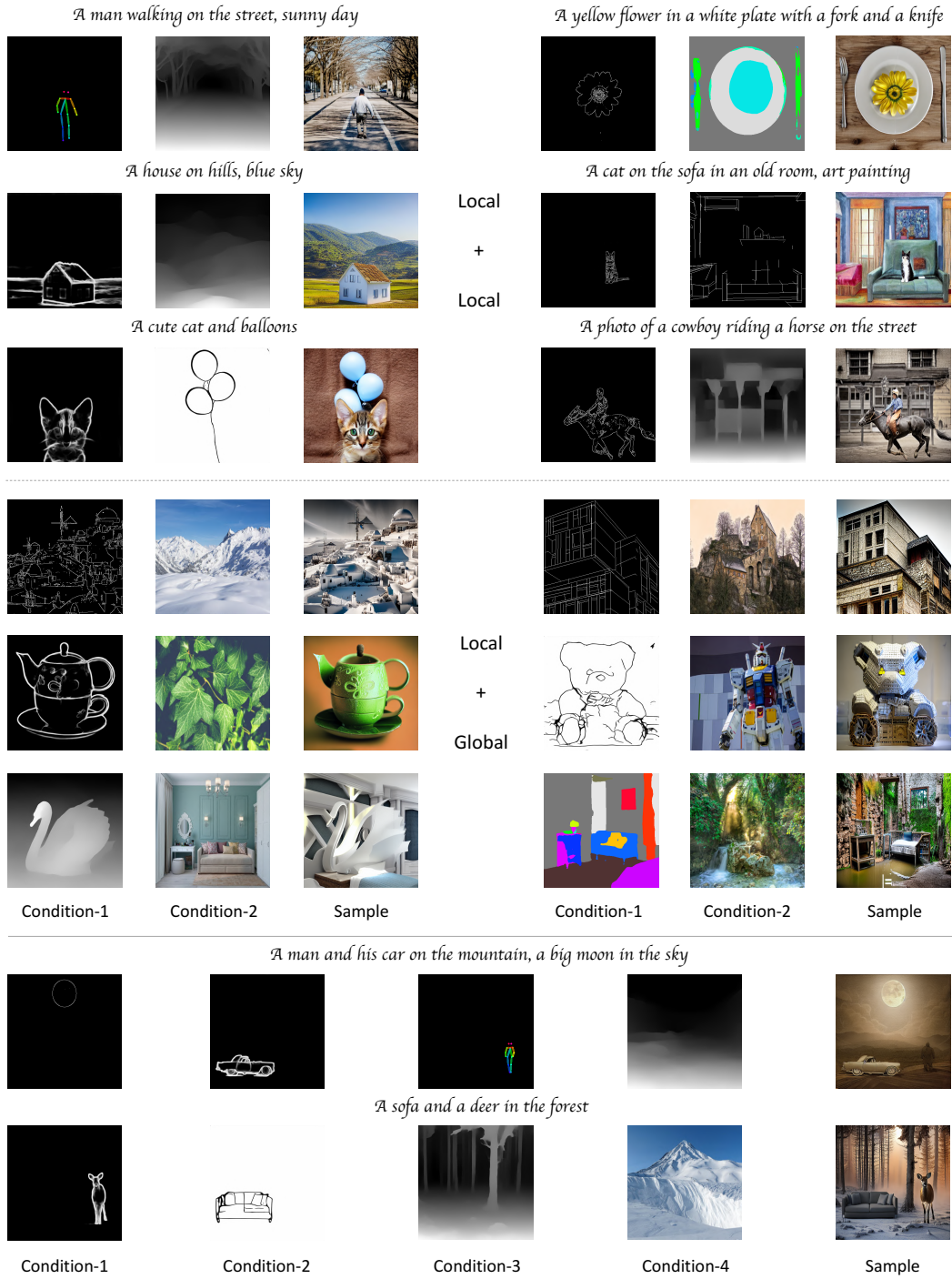


Figure 14: More visual results of Uni-ControlNet for multi-conditions setting.



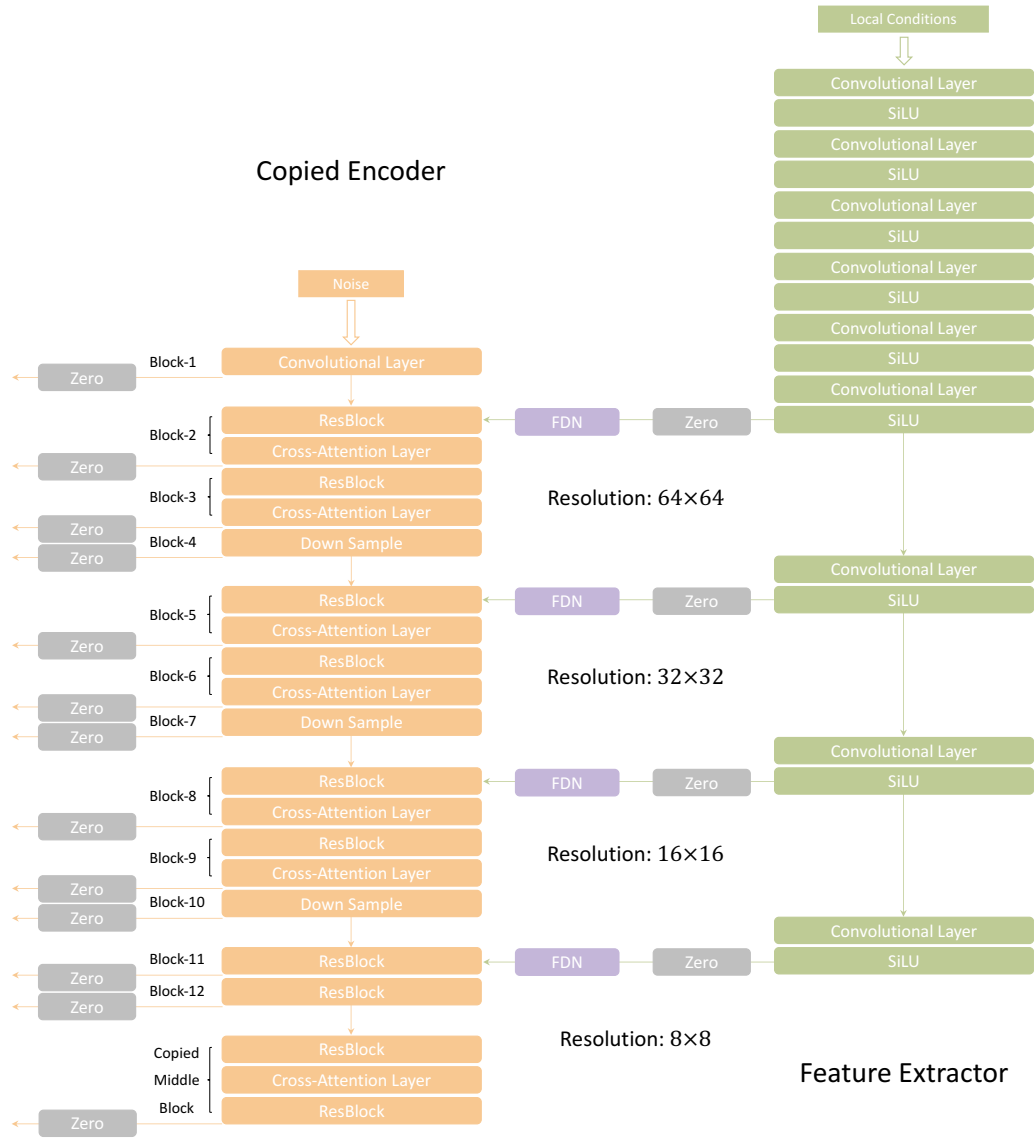


Figure 15: Details of the local control adapter.

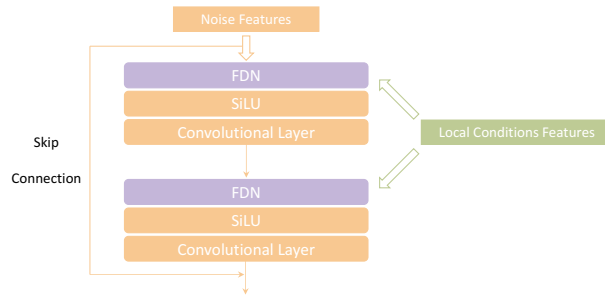


Figure 16: Details of the ResBlock with FDN in local control adapter.

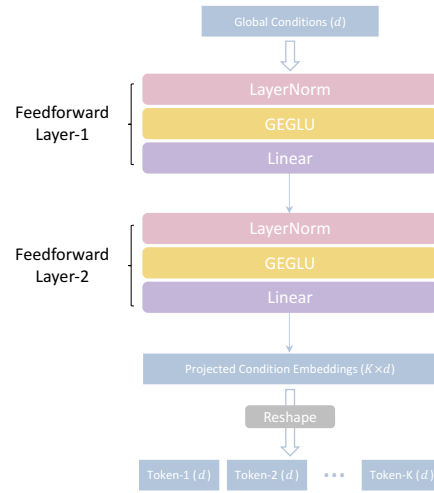


Figure 17: Details of the global control adapter.  $d$  is the dimension of text token embedding and  $K$  is the number of the global tokens.