

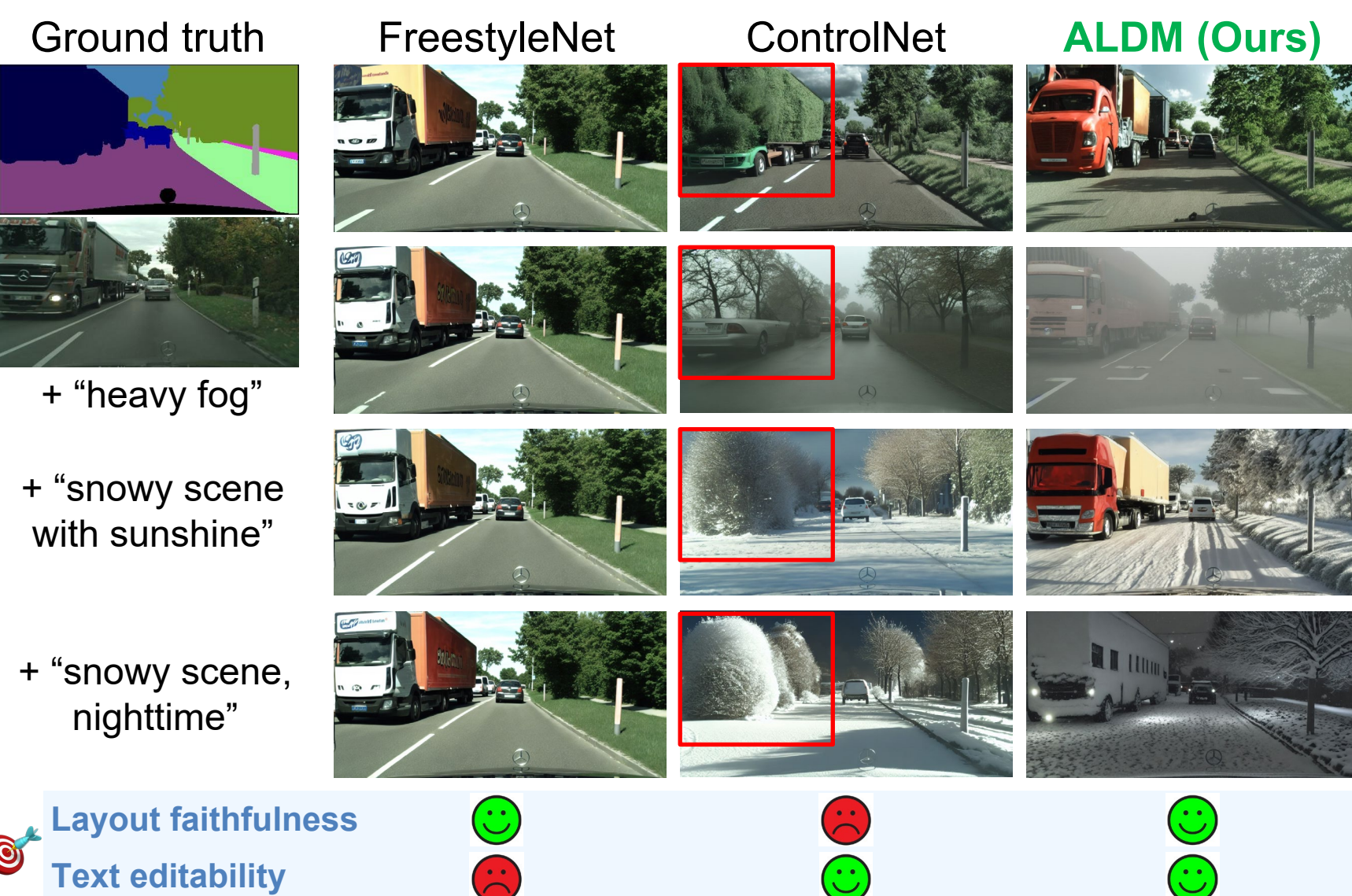


## 1. Motivation

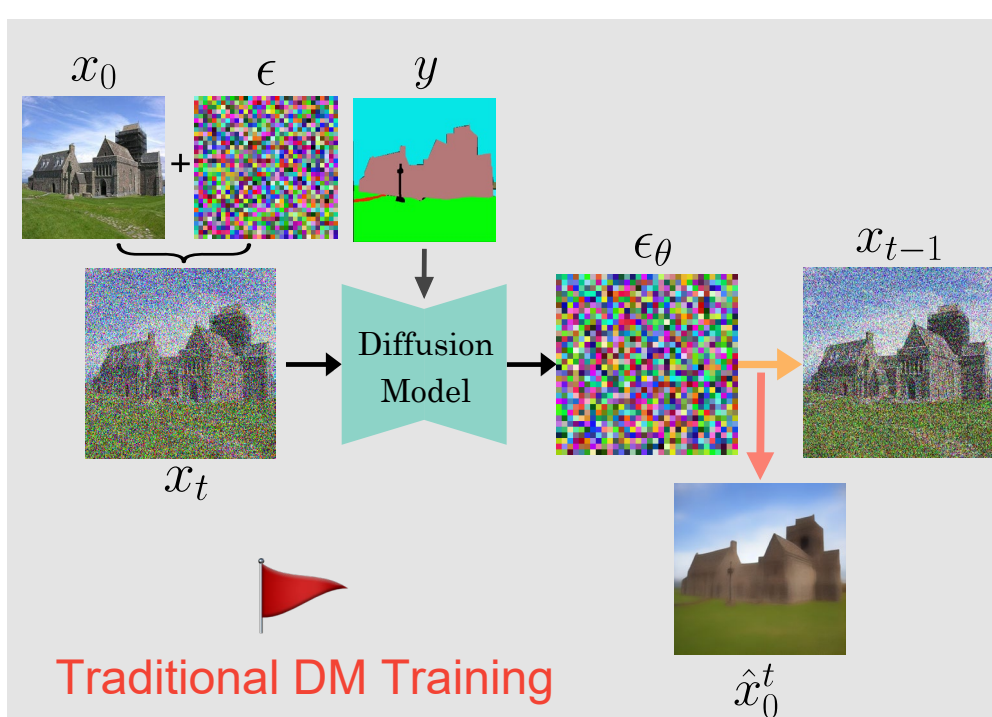
➤ Adopt **pretrained** large scale Text-to-Image diffusion models, e.g., Stable Diffusion, for Layout-to-Image task.

➤ We identify two challenges:

- Alignment** with the desired layout condition
- Editability** via text control



## 2. Traditional Training



Randomly sampled **single step**

- Sample the noisy latent  $x_t$  at a random timestep

$$x_t = \sqrt{\alpha_t}x_0 + \sqrt{1 - \alpha_t}\epsilon$$

MSE reconstruction loss

- Learn to denoise, i.e., predict the added noise  $\epsilon$

$$\mathcal{L}_{noise} = \mathbb{E}_{\epsilon \sim N(0, I), y, t} \left[ \left\| \epsilon - \epsilon_\theta(x_t, y, t) \right\|^2 \right]$$

## 3. Novel Training Strategy

Surpassing ControlNet

Model Agnostic Training Strategy

Our ALDM

Adversarial Supervision

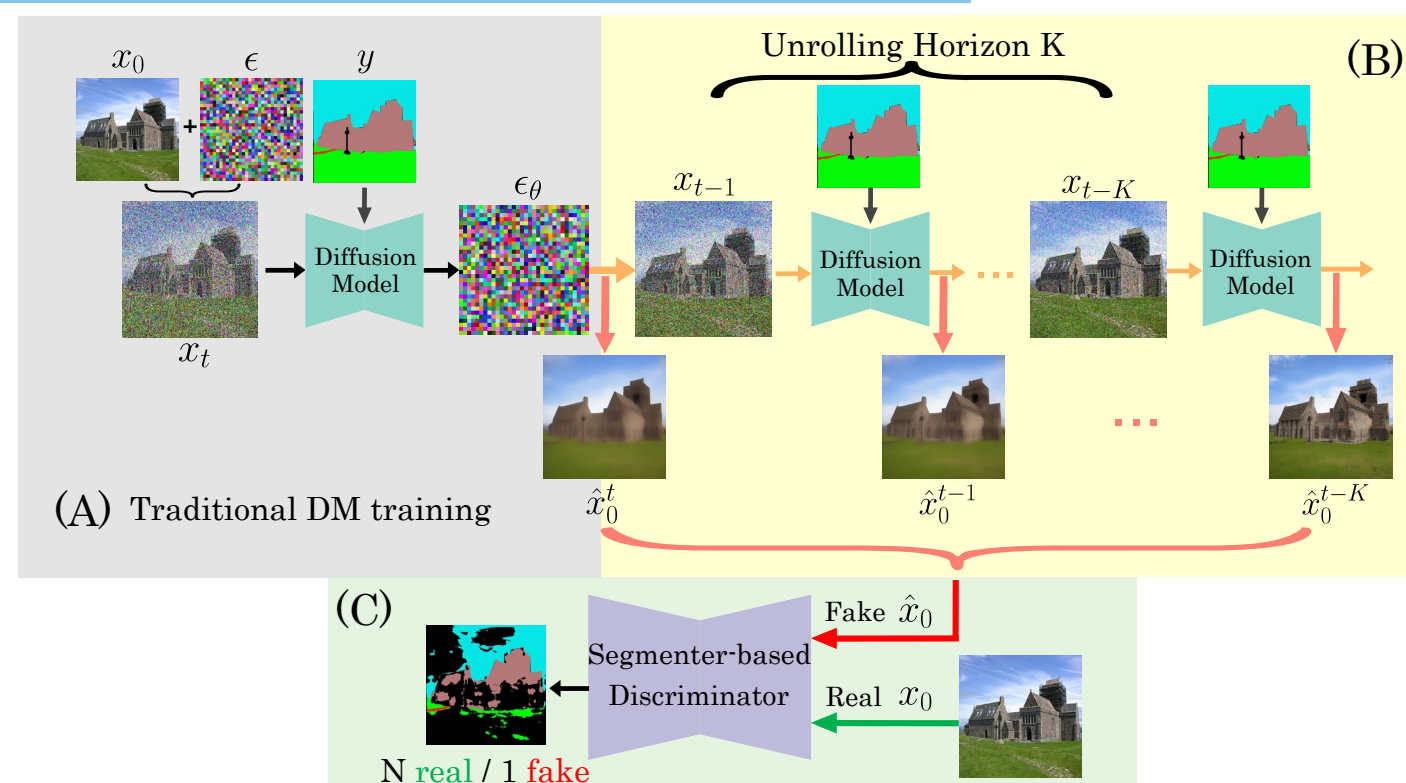
Explicitly leverage the label map condition for supervision

Multistep Unrolling

Better imitate the inference time sampling with more comprehensive supervision

Encourage the **consistent adherence** to the given label map **over a time horizon**

## 4. Method Overview



(B) **Multistep unrolling:**

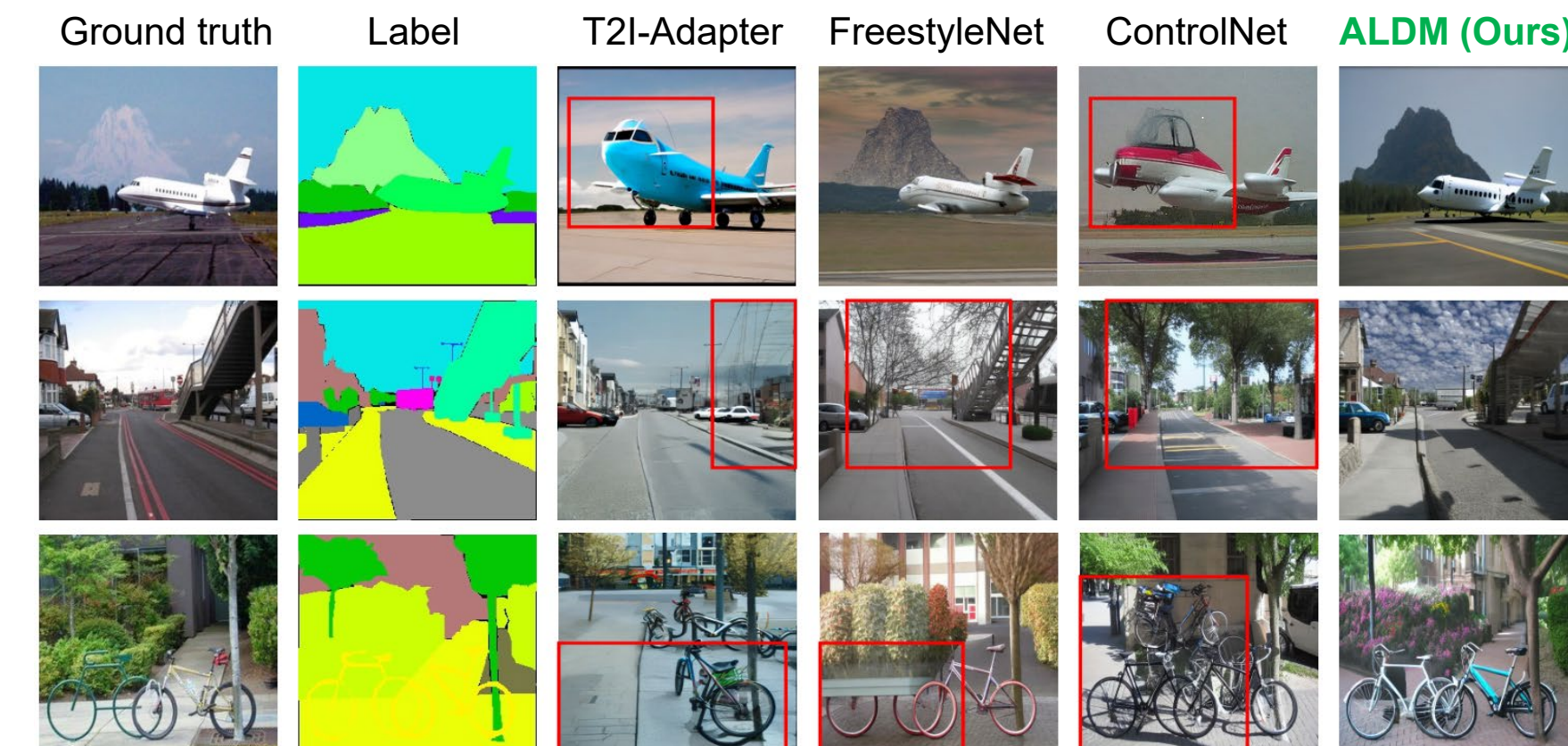
- Bridge the gap between inference time sampling and single timestep sampling during training
- Employ supervision over consecutive denoising steps → **Consistent alignment** with the given layout condition

(C) **Adversarial supervision:**

- (N+1)-Classes-Segmenter-Based **Discriminator**  
Real → N semantic classes, Fake → One extra "fake" class
- Generator** (diffusion model) should learn to fool the discriminator, i.e., synthesize samples that well comply with the label map

## 5. Results

ALDM can effectively enhance the layout faithfulness!



**Metric:**

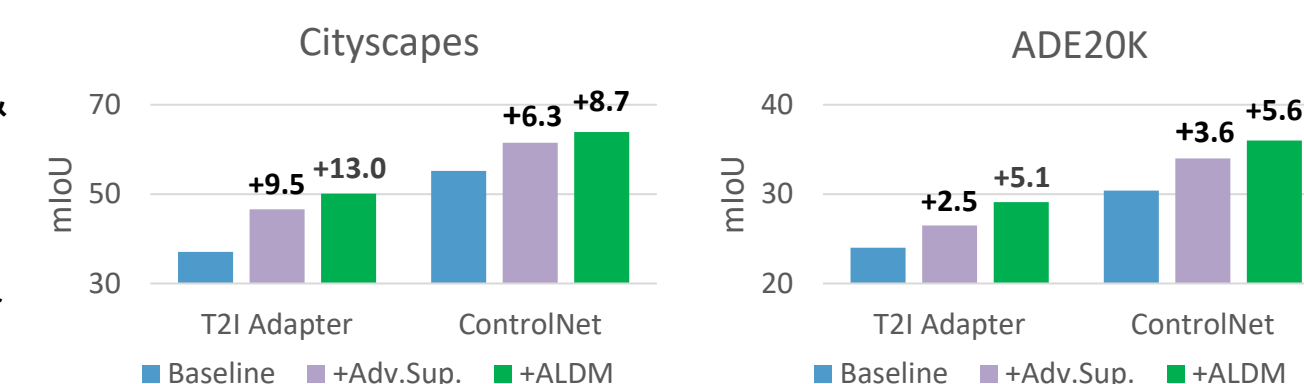
- mIoU**: measure alignment with the layout condition
- TIFA**: measure text editability

➤ By default, ALDM represents ControlNet + Adv. Supervision + Multistep unrolling.

Method	Cityscapes		ADE20K	
	mIoU ↑	TIFA ↑	mIoU ↑	TIFA ↑
FreestyleNet	68.8	0.300	36.1	0.740
T2I-Adapter	37.1	0.902	24.0	0.892
ControlNet	55.2	0.822	30.4	0.838
<b>ALDM (Ours)</b>	<b>63.9</b>	<b>0.856</b>	<b>36.0</b>	<b>0.888</b>

## 6. Ablation Study

The proposed Adversarial Supervision & Multistep Unrolling can effectively boost different layout-to-image diffusion models, e.g., T2I-Adapter and ControlNet.



## 7. Application: Improved Domain Generalization

