

## Introduction

### Background

Recent advances in Text-to-Image (T2I) generation have profoundly revolutionized the vision landscape, facilitating the synthesis of highly authentic assets from textual prompts, e.g., text-driven Image-to-Image translation and video generation. Nevertheless, designing comprehensive prompts to meticulously control every aspect of an image can be both labor-intensive and time-consuming, posing challenges for efficient generation workflows.

### Motivation

However, existing methods still suffer from:

- Imprecise localization;
- Unrealistic artifacts.

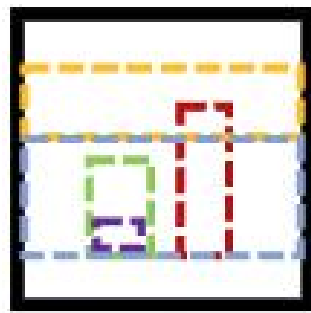
These issues primarily stem from:

- Commonly used attention energy function introduces inherent spatial distribution biases, hindering objects from being uniformly aligned with layout instructions.
- Vanilla backpropagation update rule can cause deviations from the pre-trained domain, leading to out-of-distribution artifacts.

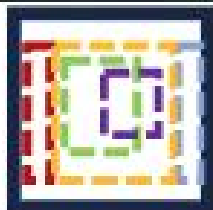
### How to solve these issues:

- Non-local attention prior is explored to redistribute attention scores, facilitating objects to better conform to the specified spatial conditions.
- Langevin dynamicsbased adaptive update scheme as a remedy that promotes in-domain updating while respecting layout constraints.

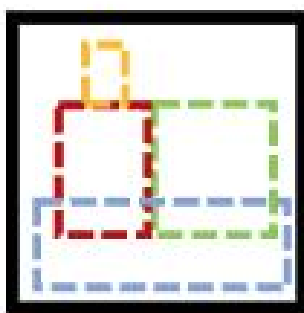
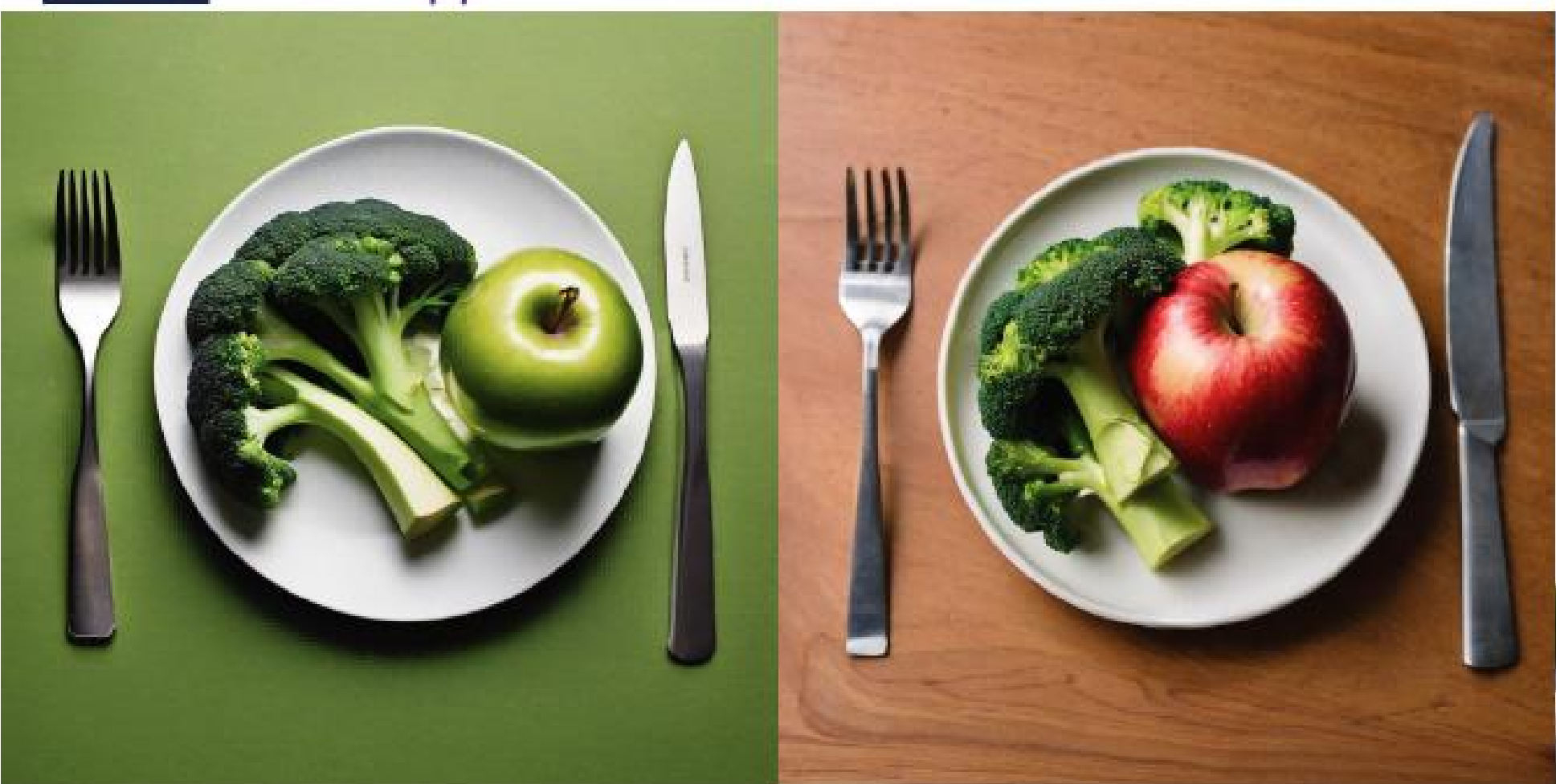
## Method



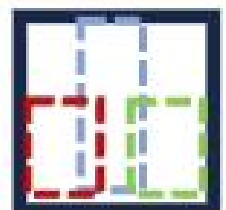
A **beggar** sitting on the **chair** under a **glowing lamp** in a utopia. The beautiful scenery scattered behind him : **Mountains** and **lakes**, studio ghibli style, art by Hayao Miyazaki, vivid colors, sharp angles, playful, 8K.



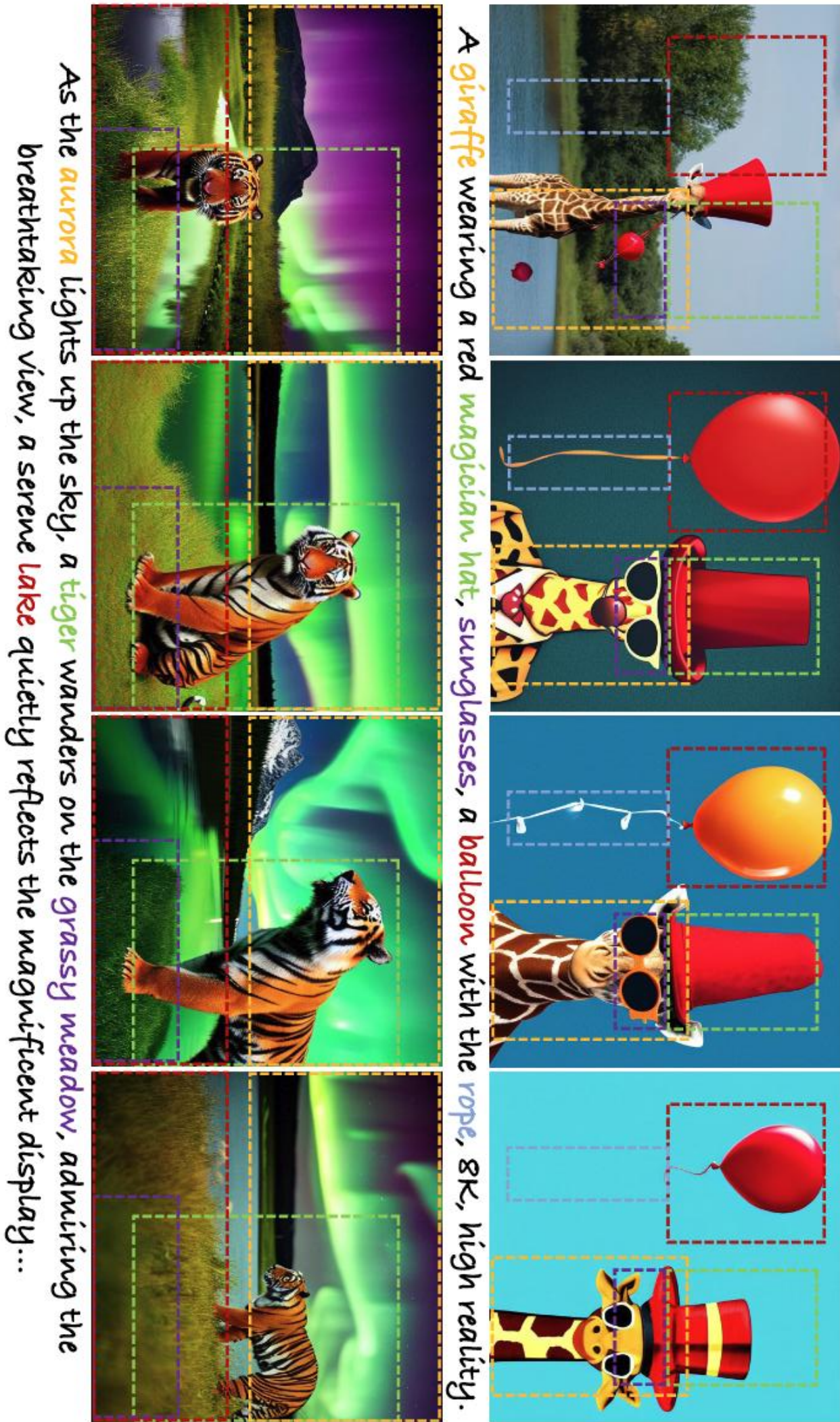
A **fork**, a **knife** and a **plate** with the **broccoli** and **apple**.



A **Pikachu** and a **dog** wearing a **pineapple hat** are playing **skateboard**, hyper realistic, highly detailed, best quality, high resolution, 8K, HD.



A **teddy bear** and a **hello kitty** sit in front of the **Eiffel Tower**.



Model	COCO2014			
	AP <sub>s</sub> ↑	AP↑	CLIP-s↑	
Att.Eng. Fun. + Back Upd.	25.8	8.4	0.310	
Non-local Att.Eng. Fun. + Back Upd.	44.1	17.4	0.318	
Att.Eng. Fun. + Ada Upd.	36.7	14.9	0.324	
Non-local Att.Eng. Fun. + Ada Upd.	49.2	19.7	0.327	