



Visual-Aware CoT: Achieving High-Fidelity Visual Consistency in Unified Models

Zixuan Ye¹, Quande Liu^{2,†}, Cong Wei³, Xintao Wang², Yuanxing Zhang², Pengfei Wan², Kun Gai², Wenhan Luo^{1,†}

¹Hong Kong University of Science and Technology, ²Kling Team, KuaiShou Technology, ³University of Waterloo, [†]Corresponding author

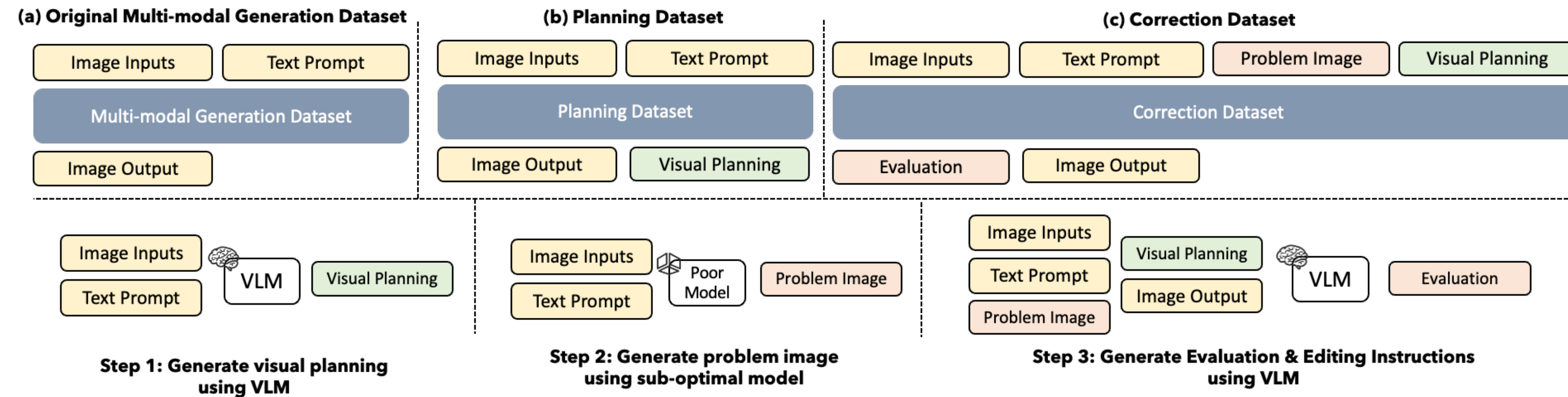


Project Page

TL;DR a paradigm shift from “text-following” to “visually-aware” reasoning in generation-understanding unified models.

Dataset

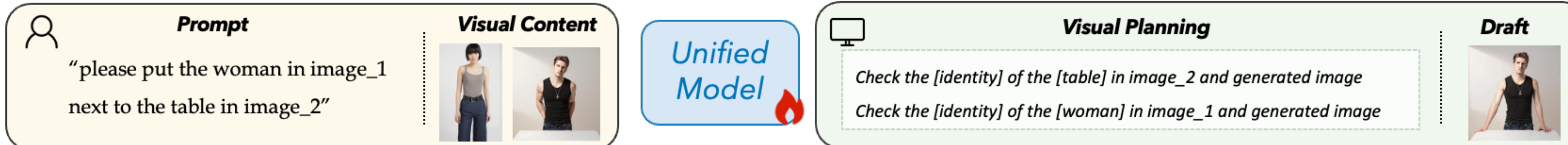
Planning Dataset & Correction Dataset



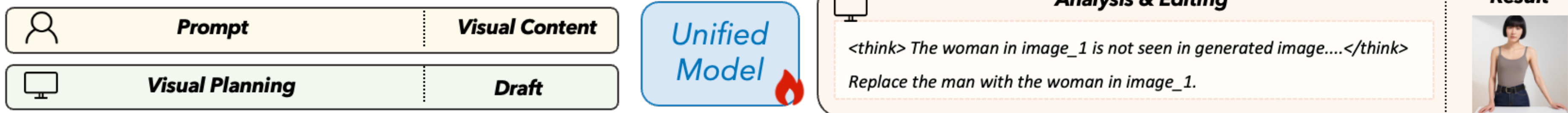
Training

SFT + GRPO

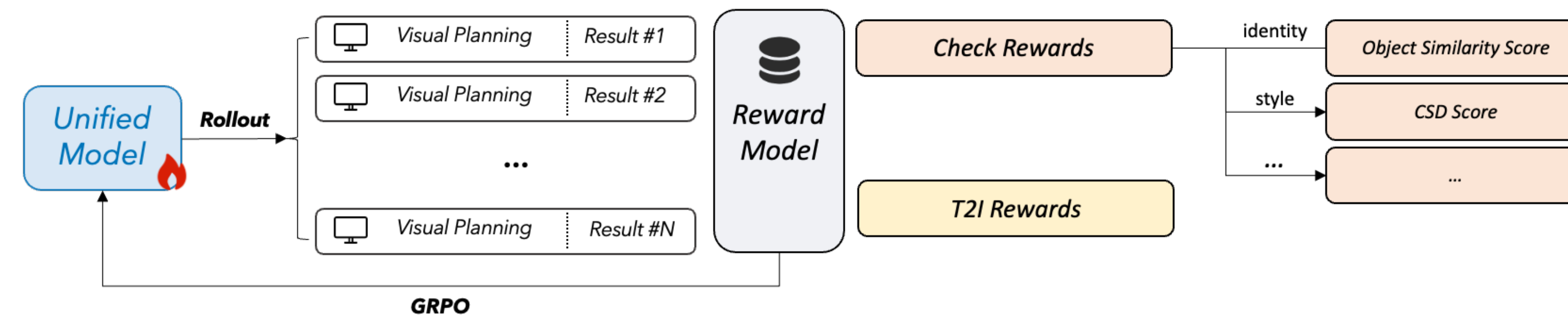
Adaptive Visual Planning



Iterative Visual Correction

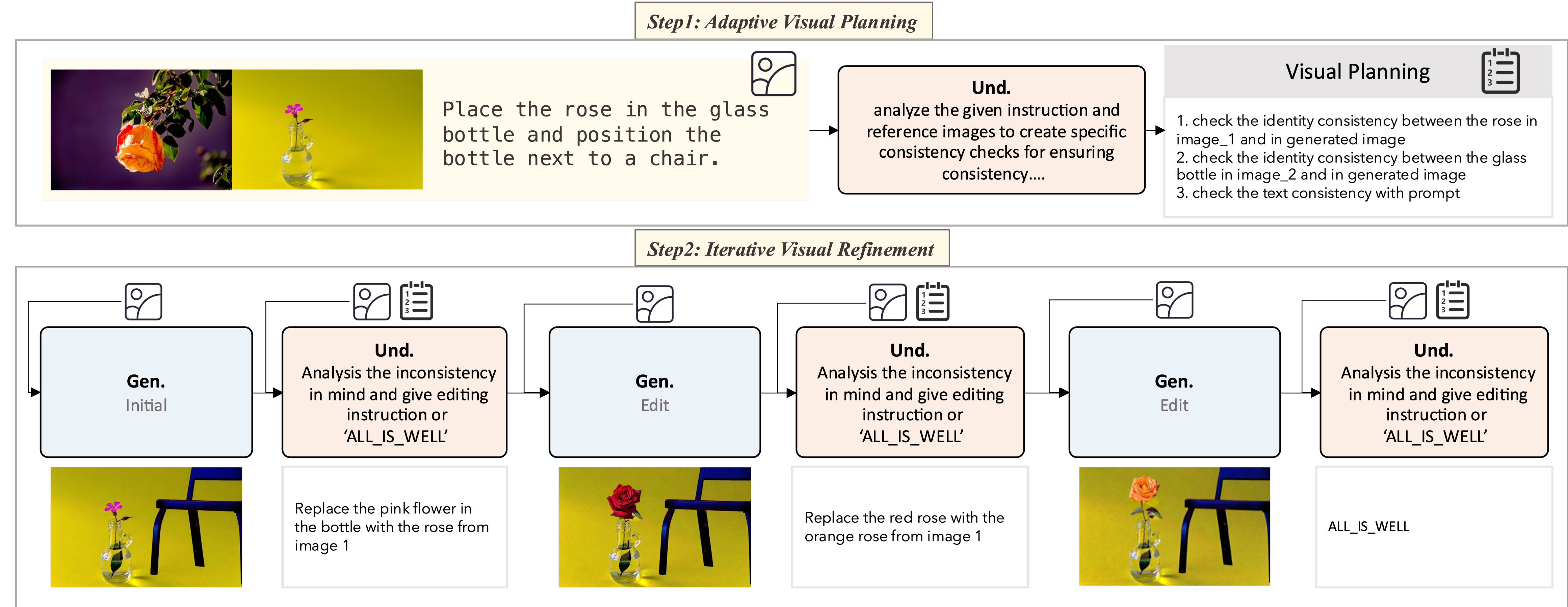


Visual Consistency GRPO



Method

Adaptive planning → Iterative Visual Refinement



Performance

Achieve High-Fidelity in Multi-reference Generation

