



Editable Scene Simulation for Autonomous Driving via Collaborative LLM-Agents



Yuxi Wei^{1*} Zi Wang^{3*} Yifan Lu^{1*} Chenxin Xu^{1*} Changxing Liu¹ Hao Zhao⁴ Siheng Chen^{1,2} Yanfeng Wang^{1,2}

¹ Shanghai Jiao Tong University ² Shanghai AI Laboratory ³ Carnegie Mellon University ⁴ Tsinghua University

* Equal Contribution



video here

yifanlu0227.github.io/
ChatSim

Overview

ChatSim is the first system achieves **language-controlled photorealistic** driving scene simulation. **Easy to use**

Control the simulation with language

1. **Multi-view consistent photorealistic rendering**

McNeRF for background scene reconstruction

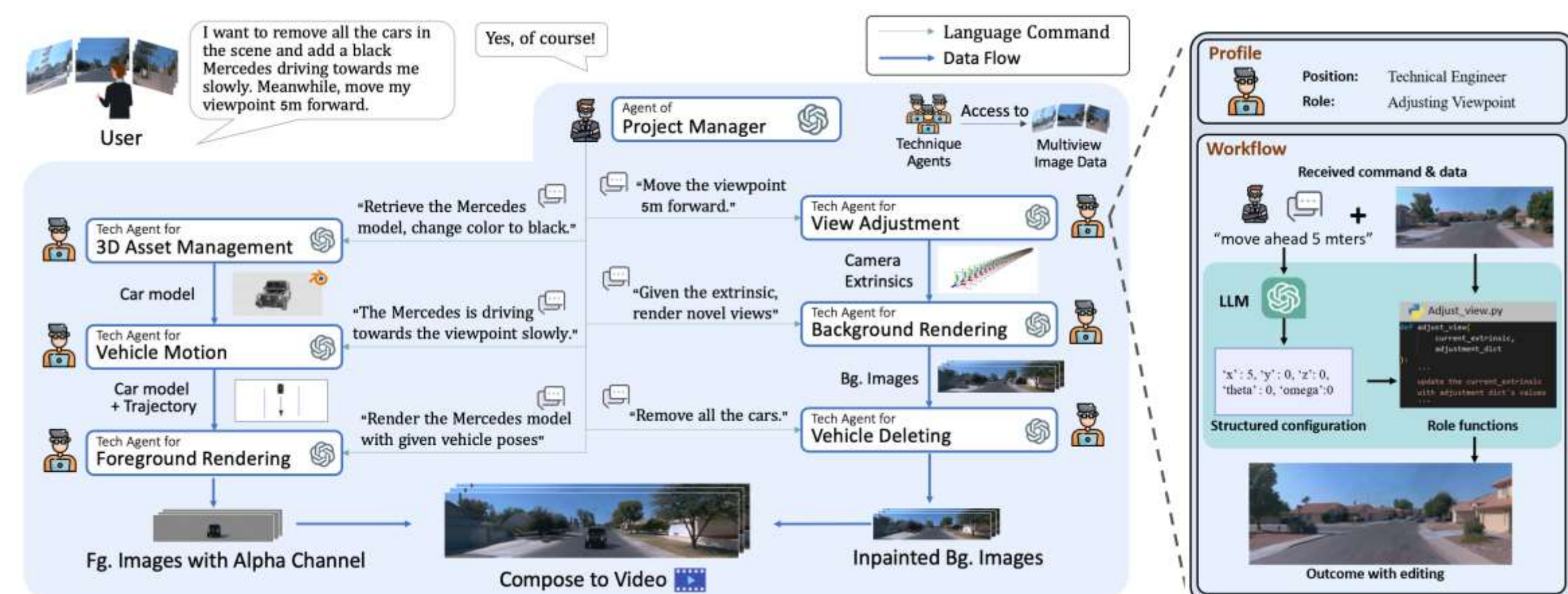
2. **Flexible simulation with external 3D assets**

McLight for lighting estimation for virtual object insertion

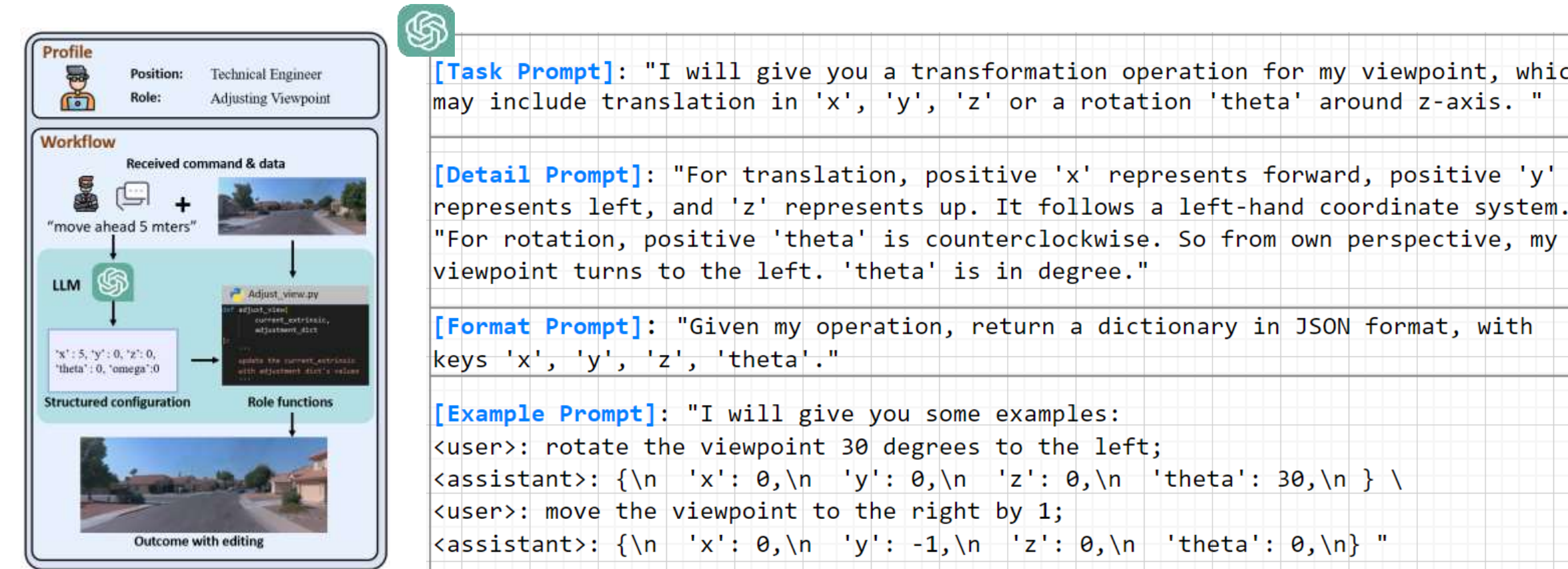


Framework

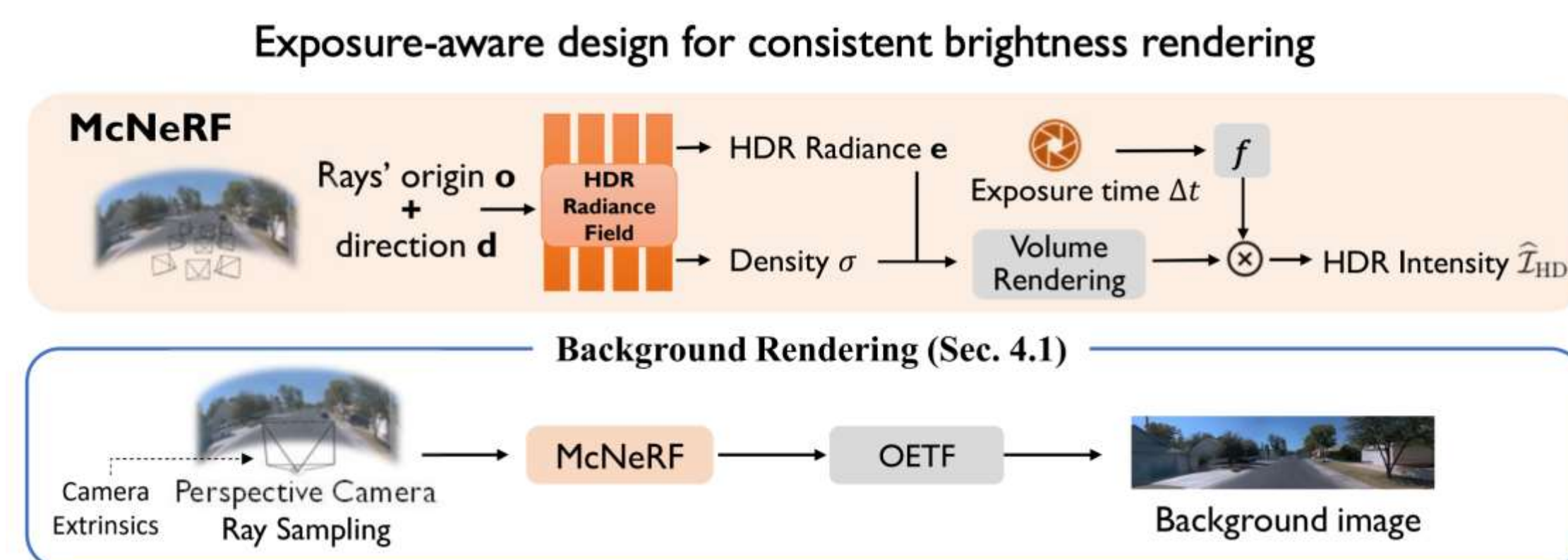
ChatSim adopts a large language model (LLM)-based multi-agent collaboration framework. The key idea is to exploit multiple LLM agents, each with a specialized role, .



Single-Agent LLM Prompting



Background: An HDR radiance field



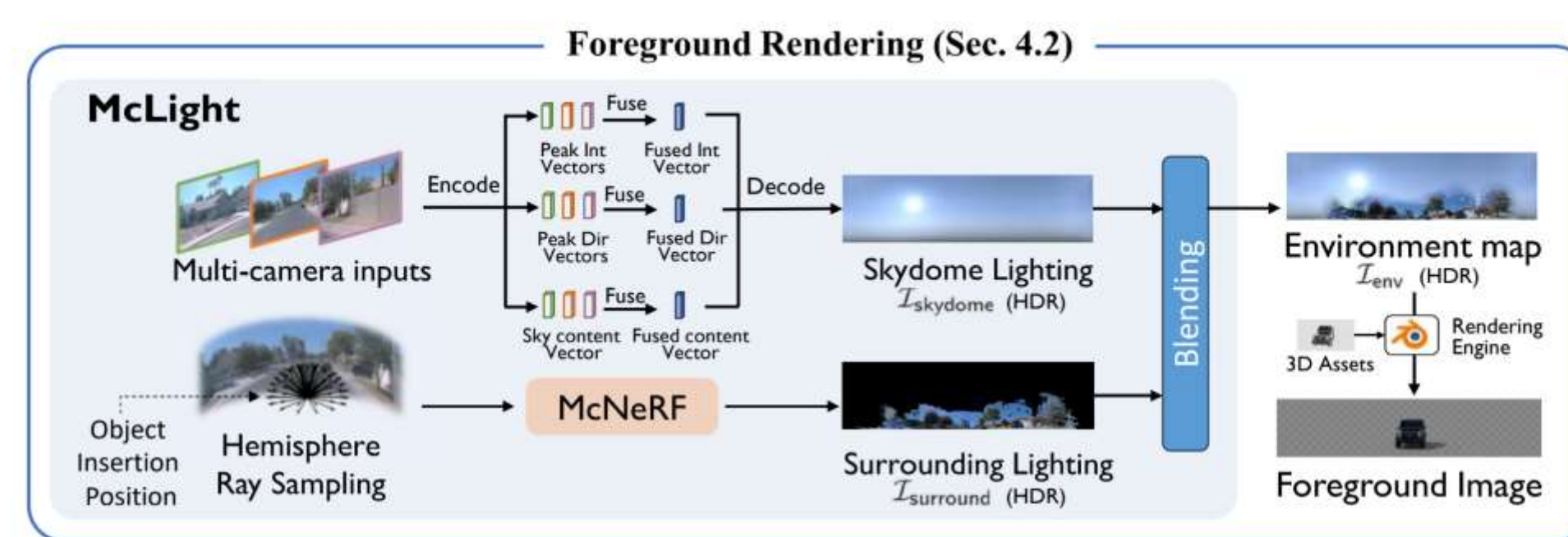
Predict radiance in HDR space (linear)

$$\hat{I}_{HDR}(\mathbf{r}) = f(\Delta t) \cdot \sum_{k=1}^K T_k \alpha_k \mathbf{e}_k$$

Calculate loss in LDR space (gamma-corrected)

$$\mathcal{L} = \frac{1}{|R|} \sum_{\mathbf{r} \in R} (\text{OETF}(\hat{I}_{HDR}(\mathbf{r})) - I(\mathbf{r}))^2$$

Foreground: Lighting Estimation + Blender Rendering



1. We predict peak intensity vectors, peak direction vectors and sky content vectors from multi-view image.
2. We leverage extrinsic-aware self-attention to fuse latent vectors to one, which will decode a skydome HDR.
3. We sample the rays directed at the hemispherical surface in McNeRF to obtain a surrounding HDR.
4. We blend two HDRs with alpha-blending and use Blender for virtual object rendering.

Results

Simulation with language command

Case 1 (highly abstract command)



Case 2 (complex command)



Virtual Object Insertion Comparison (foreground rendering)

