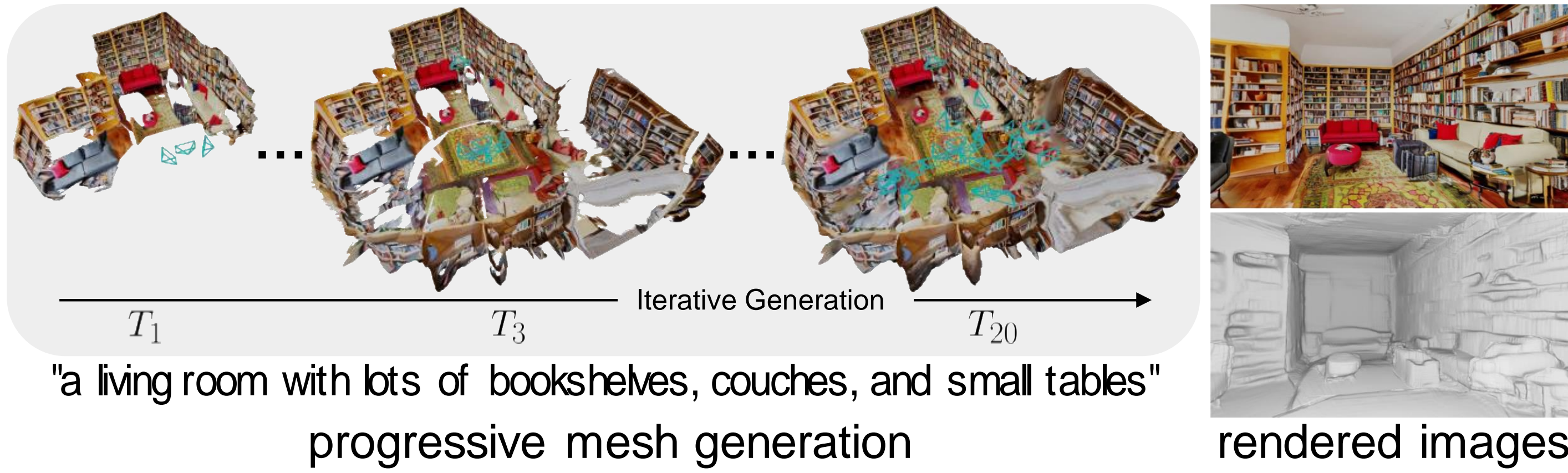


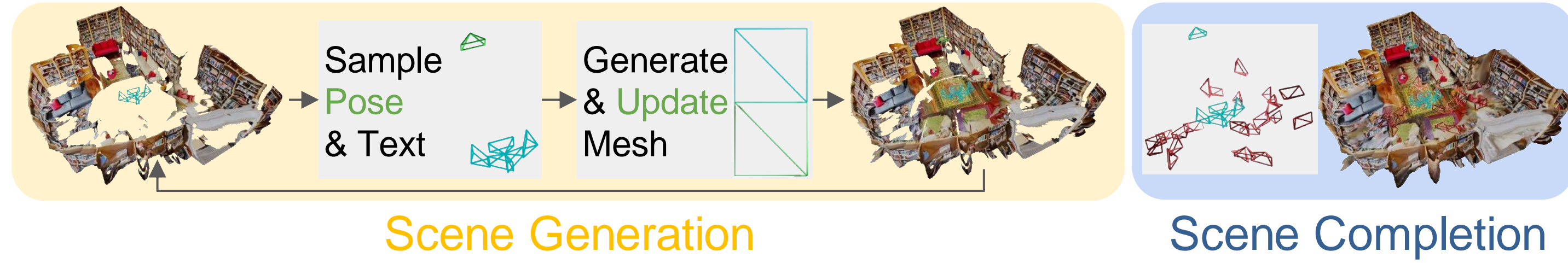


Introduction

Text2Room generates textured scene meshes from a given text prompt:

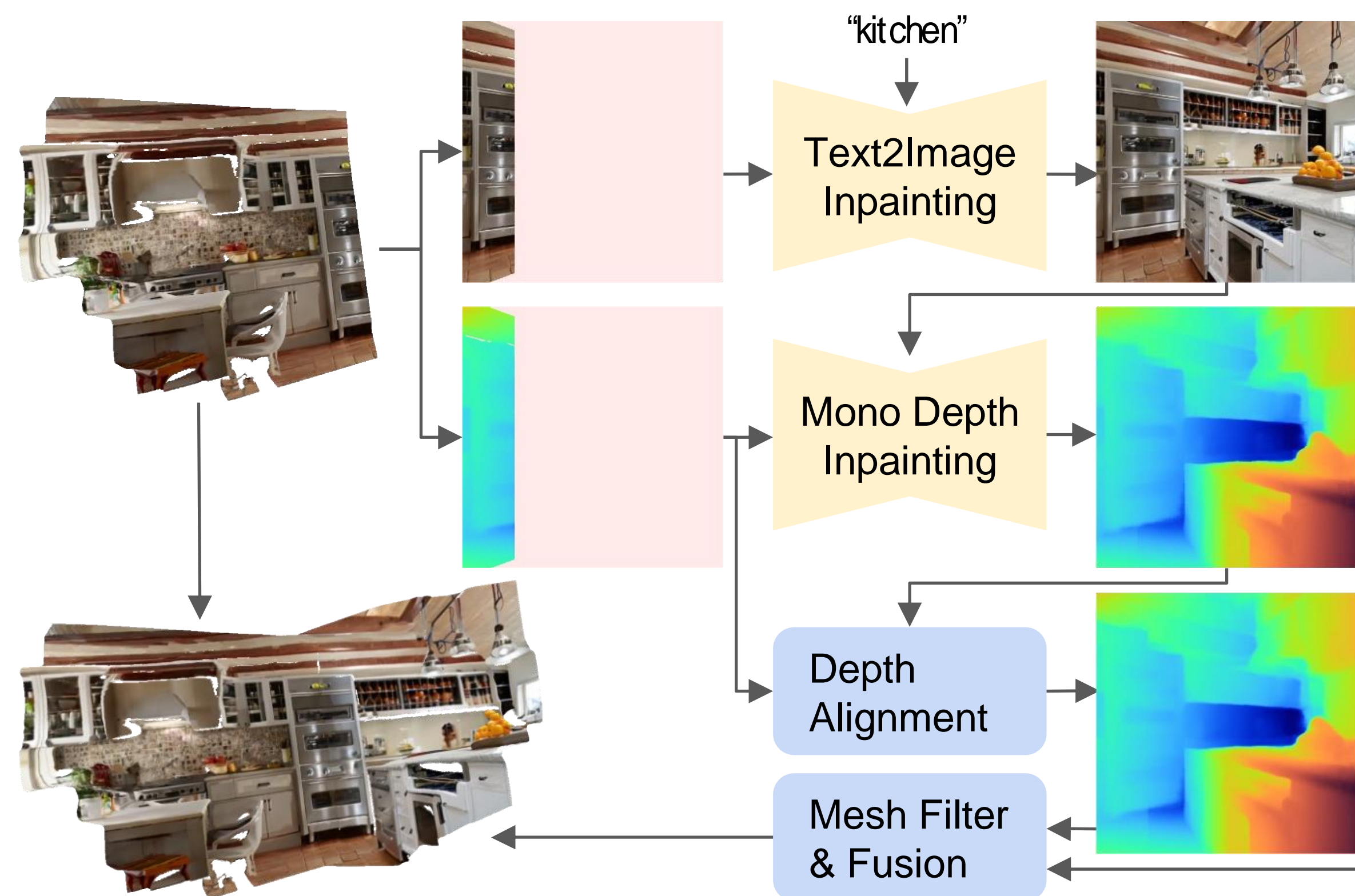


We iteratively generate meshes from (pose, text) input in two stages:



Iterative Scene Generation

For each pose, we use inpainting models to complete both rendered RGBs and depths. Then, we perform depth alignment and mesh filtering to get a next mesh patch, that is finally fused with the existing geometry.



Generated 3D Scenes



Editorial Style Photo, Coastal Bathroom, Clawfoot Tub, Seashell, Wicker, Mosaic Tile, Blue and White



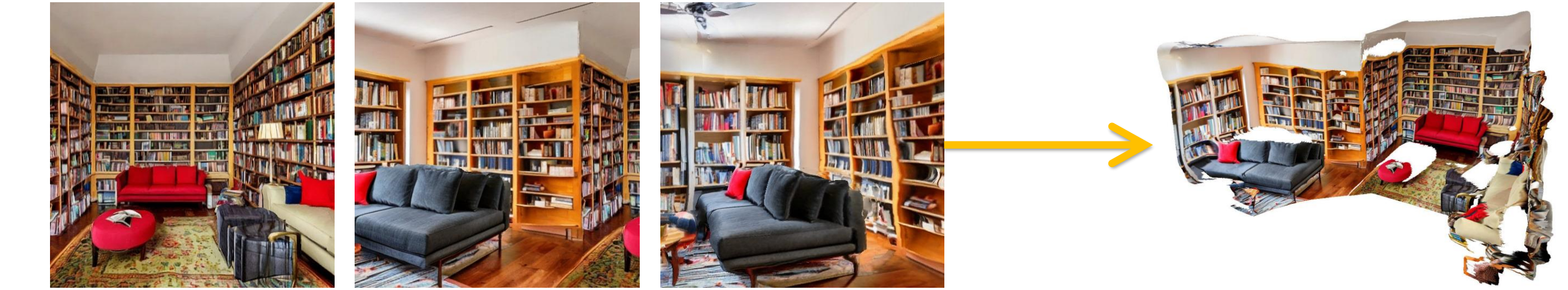
Editorial Style Photo, Modern Living Room, Large Window, Leather, Glass, Wood Paneling, Apartment



A living room with a lit furnace, couch, and cozy curtains, bright lamps

Two-Stage Viewpoint Selection

Generation: create the main parts of the scene following predefined camera trajectories and prompts, eventually covering the whole room.

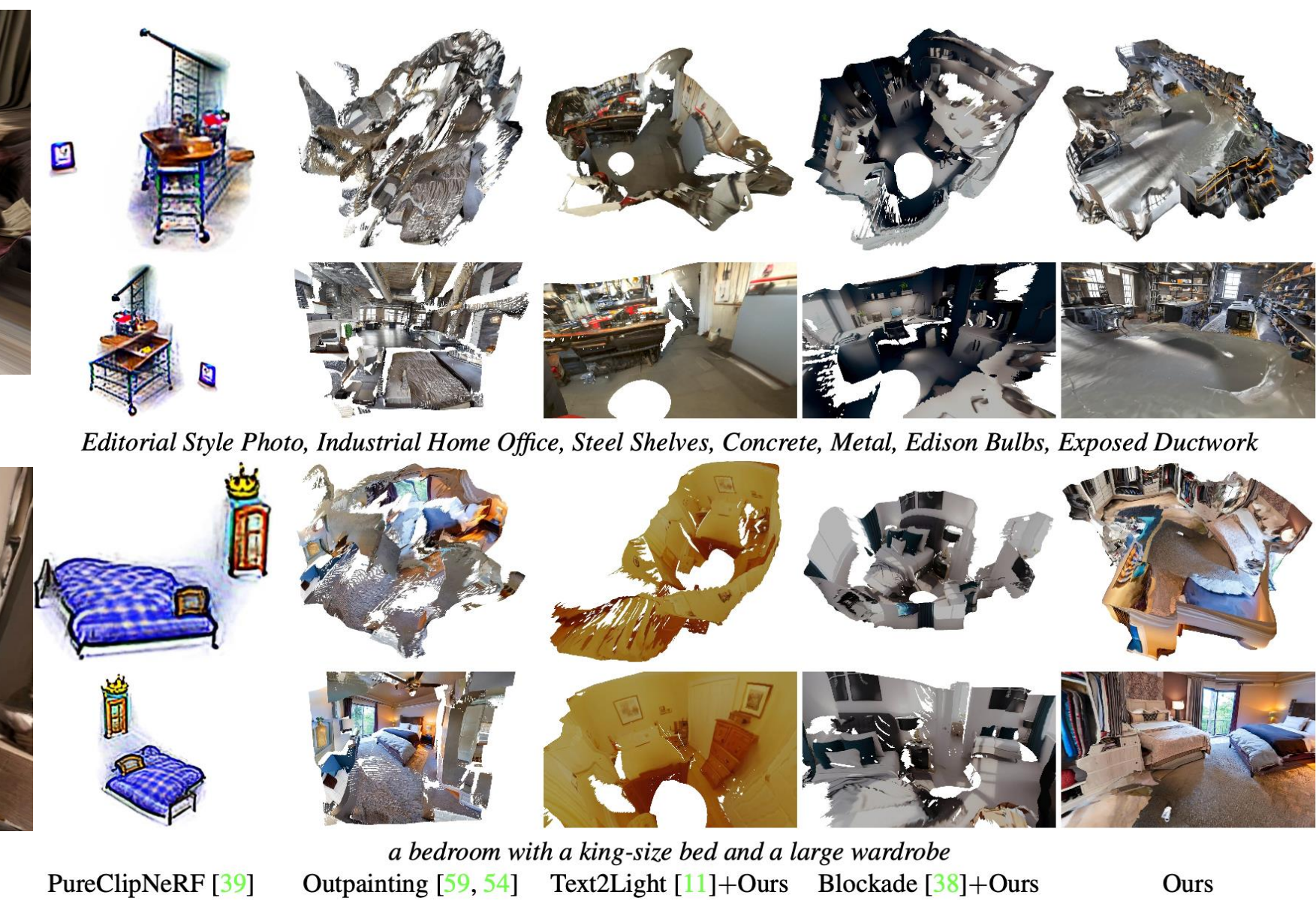


Completion: inpaint the remaining holes by sampling additional poses.

Ablations



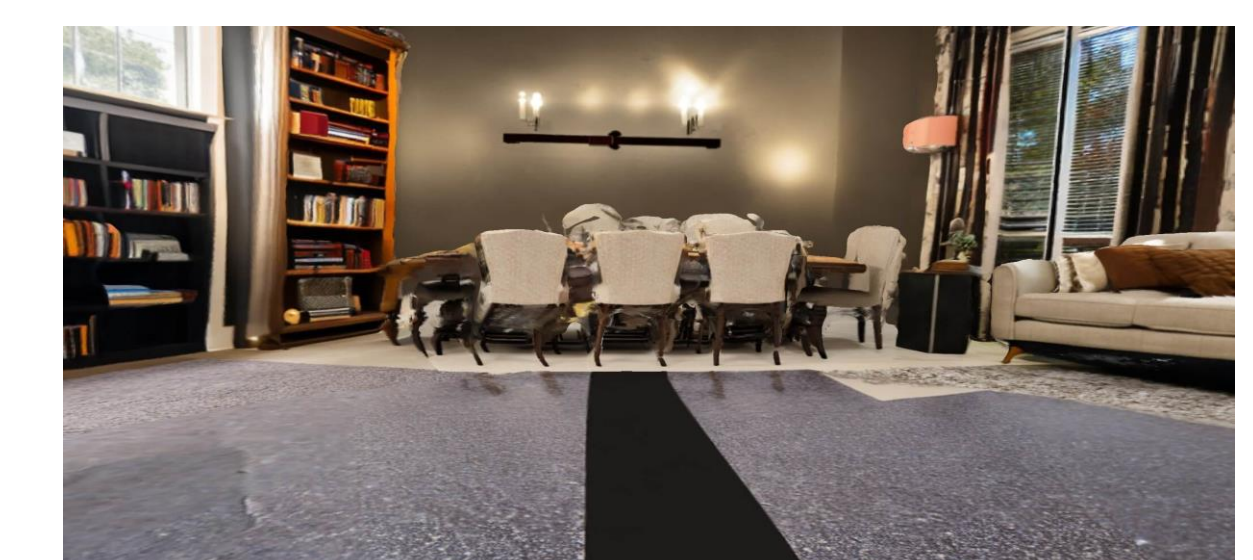
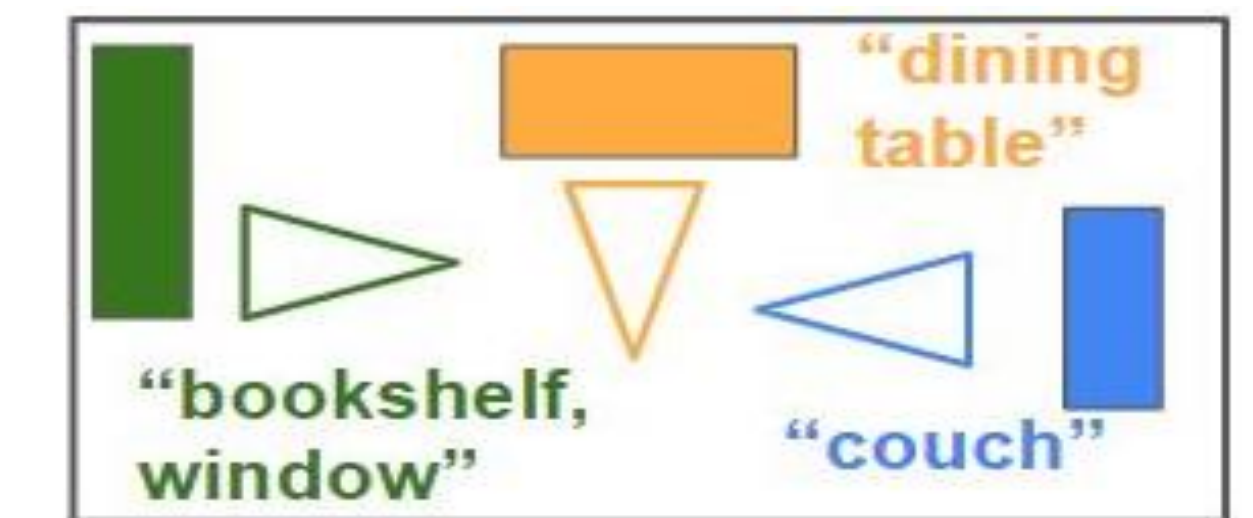
Baselines



Spatially Varying Scene Generation



Mixed Prompts



Layout Guidance