

Based on the text description T of the room's features, atomic edits will be generated. All of these atomic edits represent the edit space E . State S_t is made up of the embedding of the current image I and an ordered set of all edits completed of which each edit $e \in E$ of all atomic edits.

$$S_t = (I_t, E_t)$$

The action consists of the selection of a new edit to update state S_t . This process will either be random or based on a heuristic that can help decide what edits should be made first. After each new action a new image embedding I_t is generated and the list of completed edits E_t .

Actions = members of E not in E_t

$$\pi(e|S_t) = 1 / |\text{members of } E \text{ not in } E_t| \text{ or } \pi(e|S_t) \propto h(e)$$

No more actions are allowed for one search after all edits are completed so each run will have $|E|$ edits to make. Edit paths will be selected using naive backtracking and since the output of each step is a new image (or the embedded information in the image embedding) it is fully observable. The state S_t will transition upon the action of selecting a new edit to a new picture and edit list (I_t, E_t) using a stable diffusion model with img2img capabilities.

$$I_{t+1} = f(I_t, e_t) \text{ where } f \text{ represents the img2img process}$$

$$E_{t+1} = E_t + e_t$$

The base image for the img2img operation is the current image I_t and the text prompt e_t will be the selected edit, for example "the bed had blue pillows". An observation score r_t can be computed as the cosine similarity between the current image embedding and the embedding full text description T known as T_e .

$$r_{t+1} = \cos(I_{t+1}, T_e)$$

If the edit ever causes the cosine similarity between I_t and T_e to go down then that branch is not working well and can therefore be pruned. Each edit should push the similarity score higher as features start matching better and better.

$$\text{Prune if } r_t < r_{t-1}$$

A greedy selection policy can even be applied to always continue expanding the state node with the highest scoring image. This makes the most sense if there is a goal threshold in mind like cosine similarity of 0.8, for example and there are so many edits that expanding most branches is infeasible. This addition also changes the end state to $E_t = E$ and cosine similarity of 0.8 or higher.

$$t = |E|$$

$$r_t > 0.8$$