

BACKGROUND & MOTIVATION

A score-changing edit is not automatically evidence.

Generative editing can shift perceived scores. But the change may come from intended cue - or by confounds introduced during editing, such as changed lighting, layout, or other visual cues.



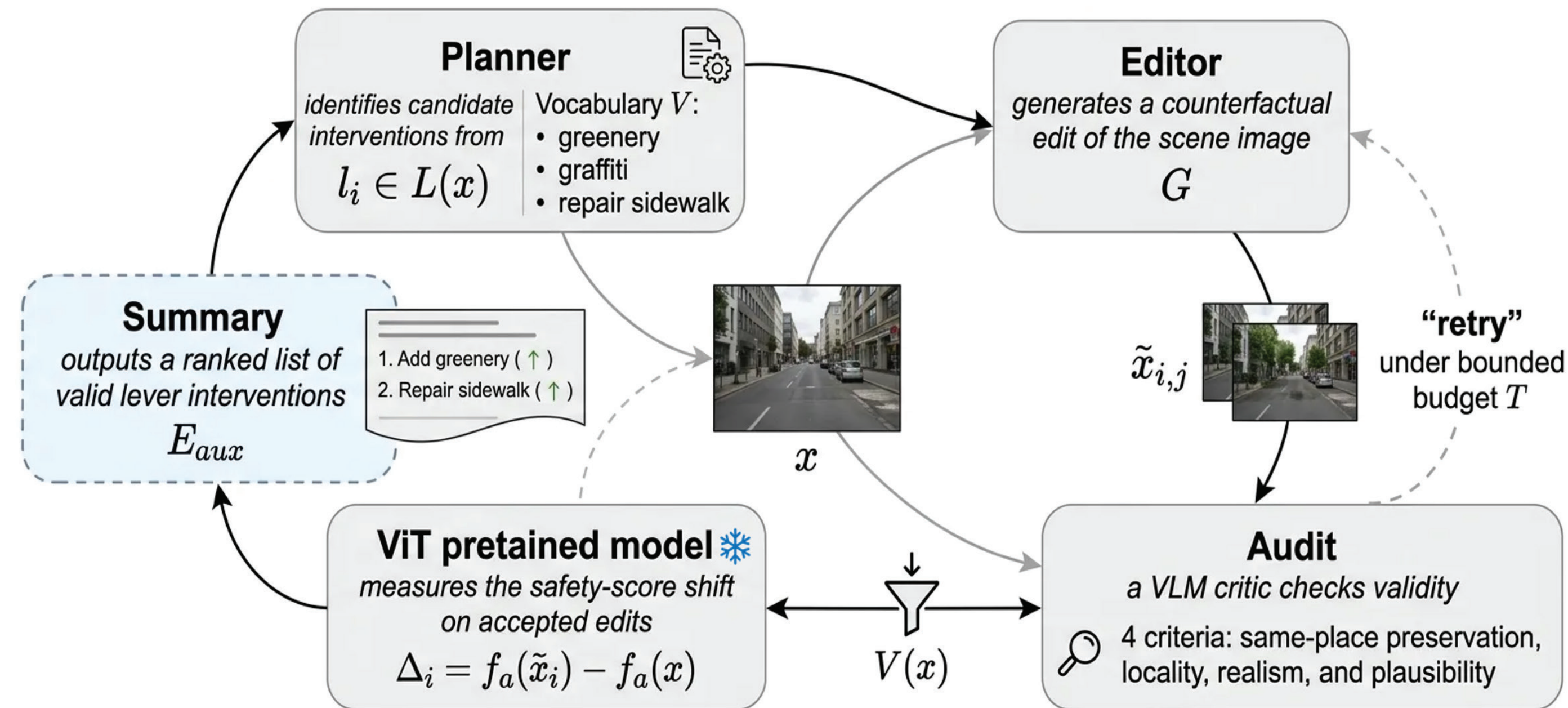
What should change?

Where should it change?

Would the edited scene still be the same place?

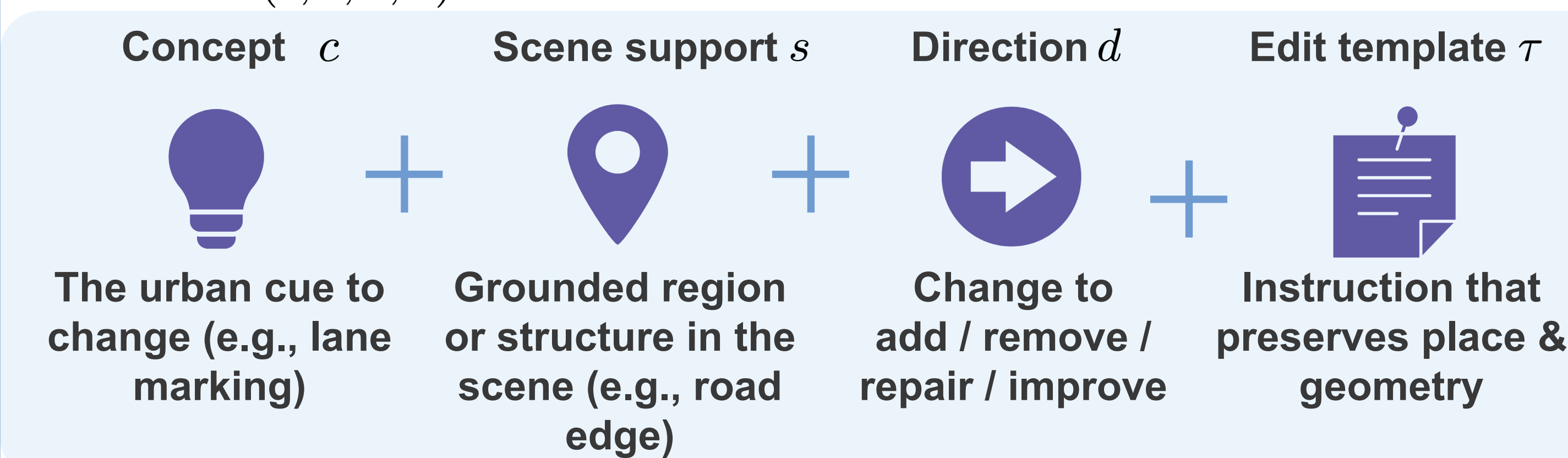
VALIDITY-GATED INTERVENTIONAL ATTRIBUTION

Protocol overview



Lever intervention (the unit of explanation)

A lever $l = (c, s, d, \tau)$ specifies:



Example lever: same intended lever, different evidence status

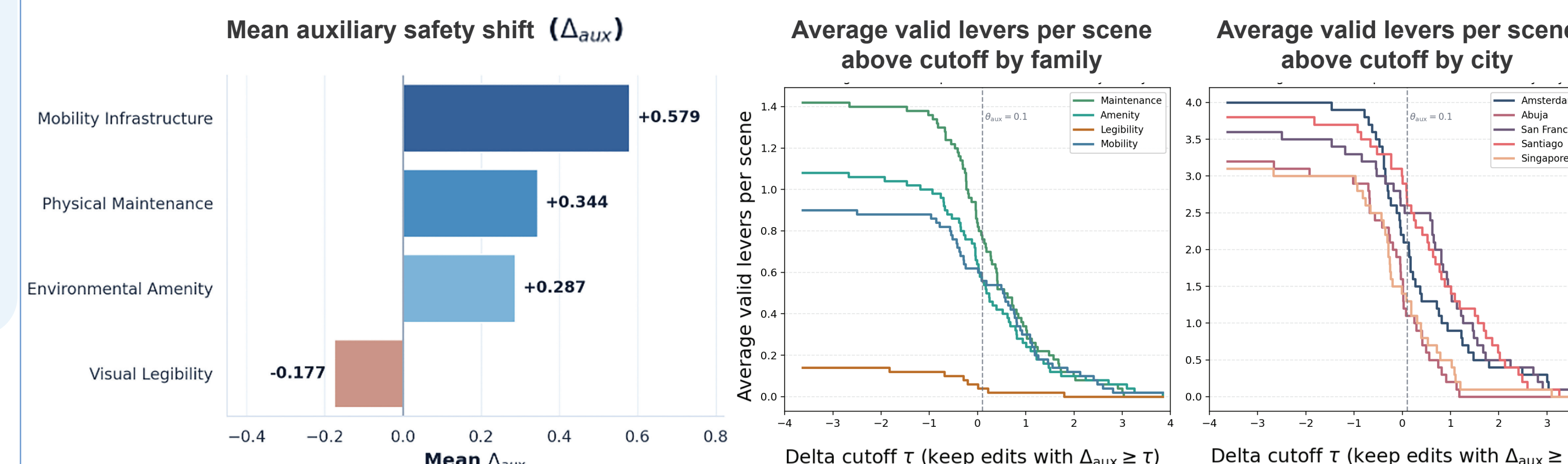
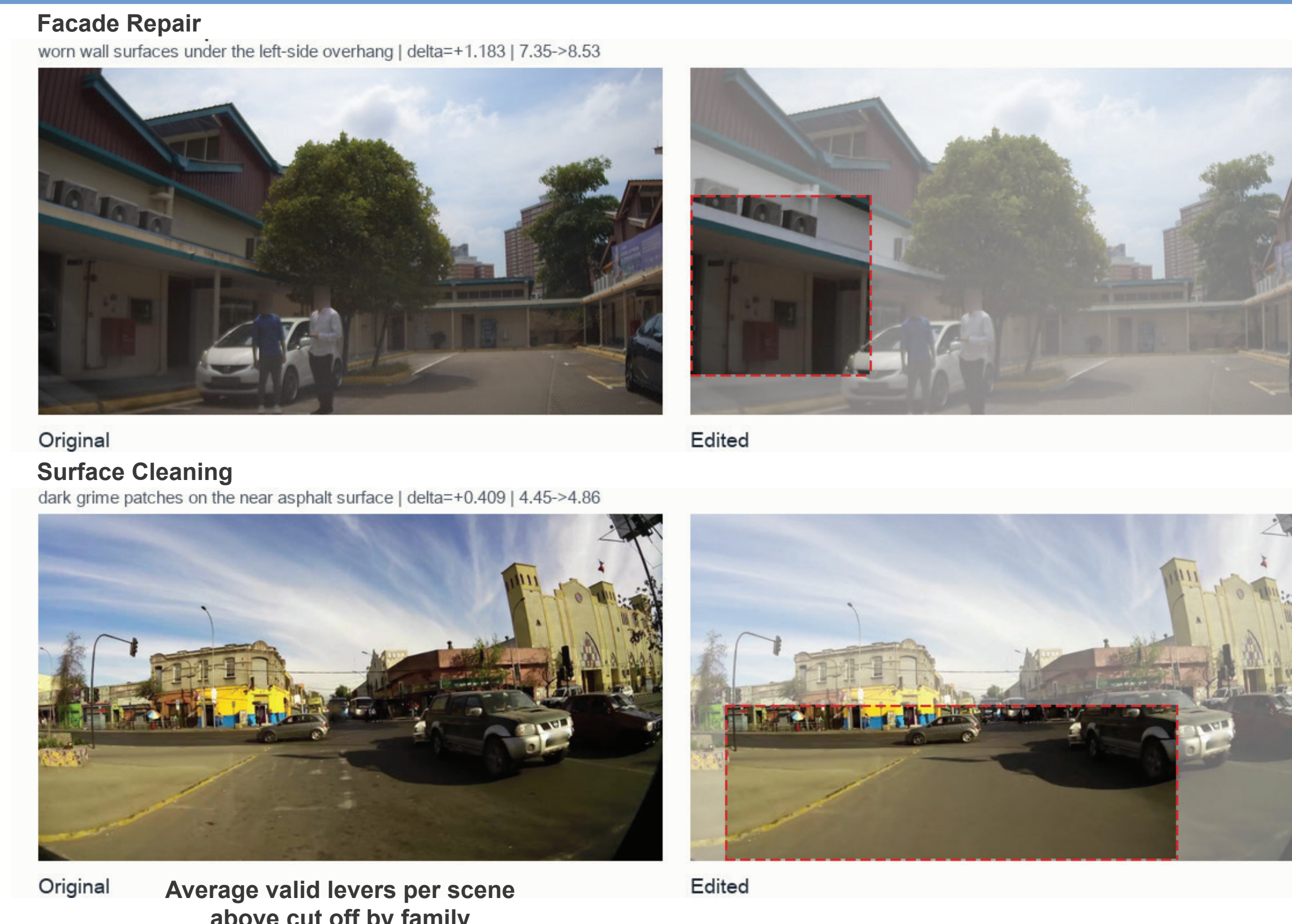


VALID EDIT: Same place • Local • Realistic • Plausible

INVALID EDIT: Confounded changes

An edit is scored and retained only if it passes all validity checks (same-place preservation, locality, realism, and plausibility).

PILOT STUDY & RESULTS



TAKEAWAY

- What we claim**: A protocol for ranking visual levers only after edits become admissible evidence.
- What we do NOT claim**: Proxy shifts are not human-validated causal effects.
- What this enables**: A cheaper and better-prioritized path toward human evaluation and real-world urban design interventions.

Why counterfactuals still need auditing

| | |
|--|---------------------------|
| Saliency Map Show where the model looks. | Not an edit. |
| SHAP / Concept attribution Identifies correlational visual cues. | Not intervention-tested |
| Full-scene counterfactuals Generates alternative scenes | Changes too much at once. |

None of these tell us whether a single, localised, valid edit can plausibly shift perception.

The motivation

| | | | |
|---|--|---|--|
| Localised Change only what matters. | Plausible Realistic & context-appropriate. | Same-place preserving Keep the scene identity | Auditable Checked before we interpret. |
|---|--|---|--|

We treat each edit as an auditable intervention hypothesis, the strongest lever is the largest score shift among edits that remain valid.