

MOTIVATION

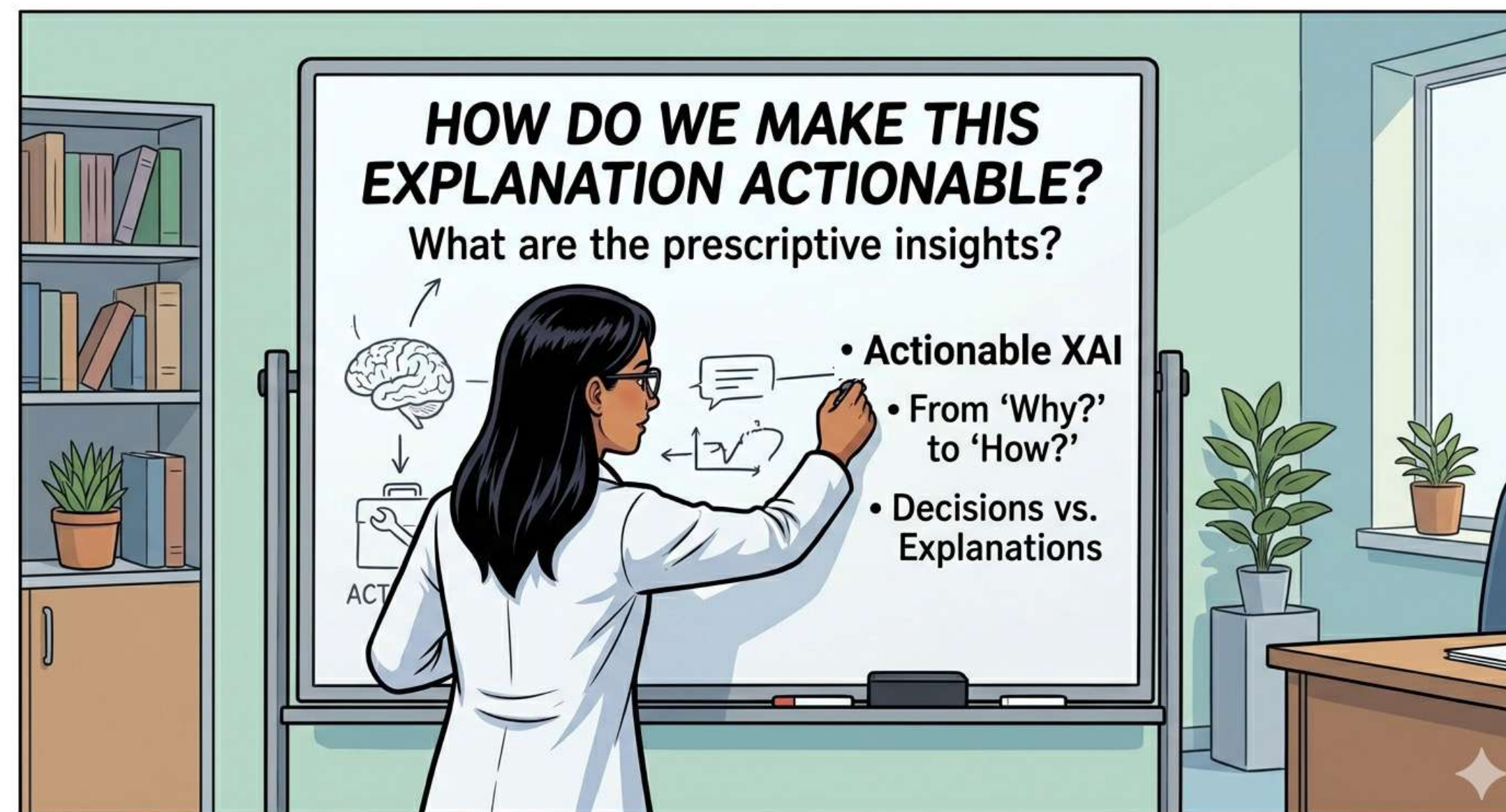


Figure generated with Google Nano Banana 2

XAI reveals feature influence on predictions ... but what's next?

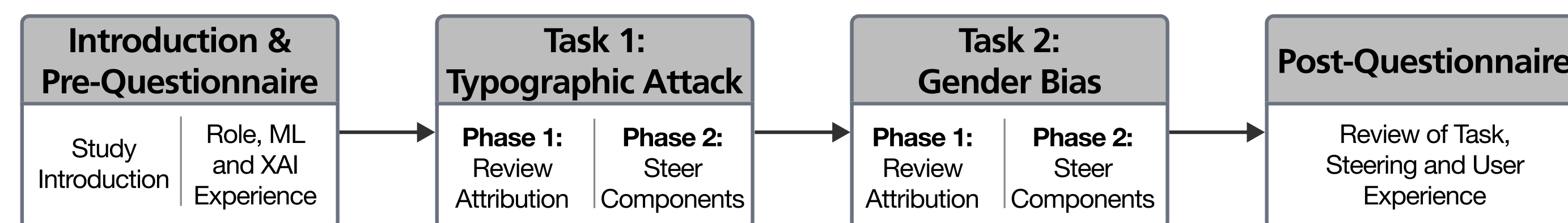
Attribution: Identify associated components, but cannot confirm causality.
 Steering: Modify model internals, but requires information to act upon.

What if we combine both, how will practitioners use it?

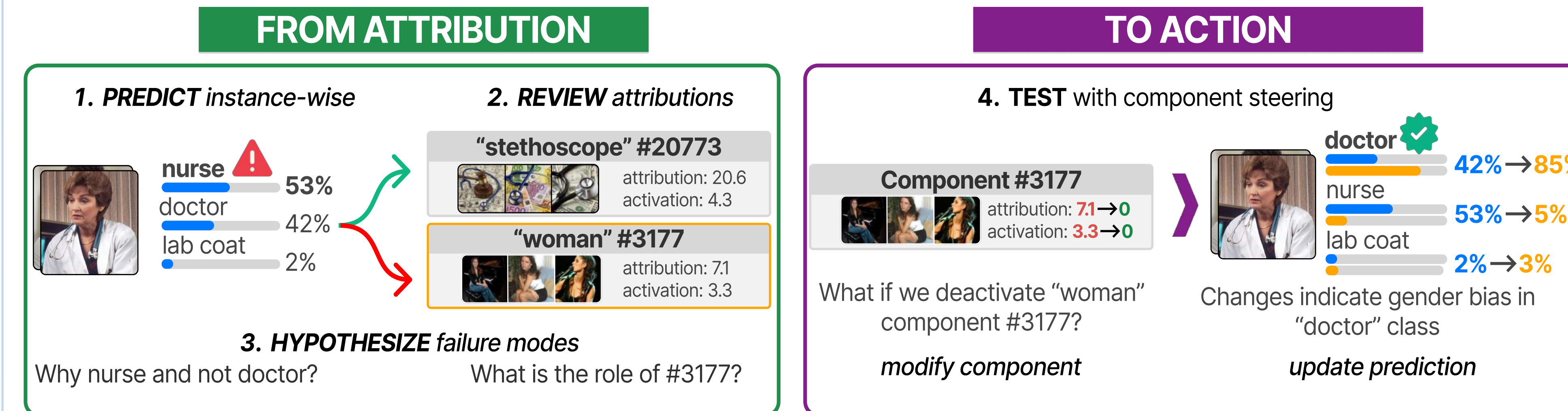
CONTRIBUTIONS

- Tool:** Web tool implementation of SemanticLens combining SAE-based attribution with activation steering for instance-level debugging of CLIP.
- Workflow:** Predict → Review → Hypothesize → Test.
- Expert Interviews (n=8):** First user study of practitioner reasoning, trust, and strategies for steering.

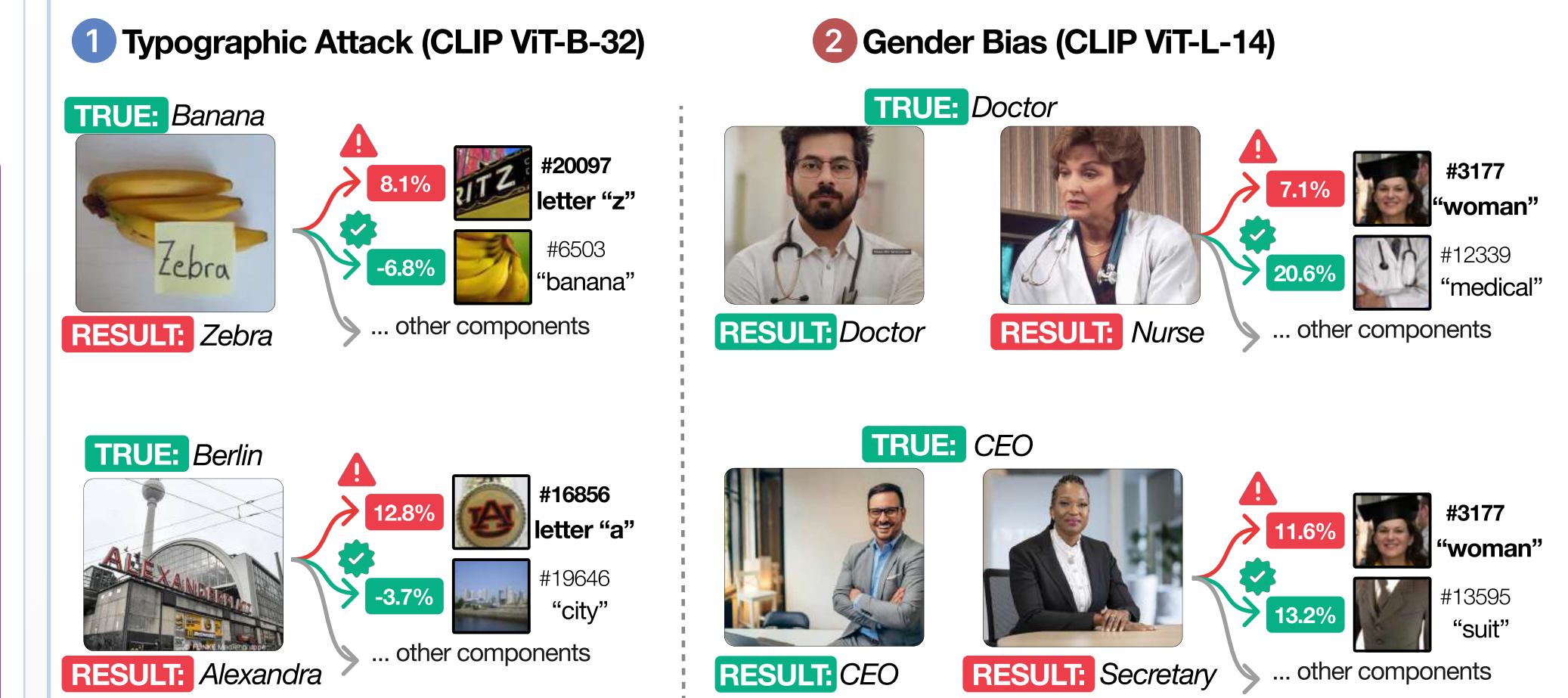
USER STUDY



WORKFLOW



STUDY TASKS



KEY FINDINGS

Reasoning shift: All 8 completed causal test cycles after steering. 6/8 formed explicit hypotheses before intervening. 5/8 spontaneously acknowledged limitations of causal claims.

Trust calibration: 6/8 shifted to evidence-based trust (grounded in observed model responses, not plausibility). 3/8 demanded external validation before deployment.

Strategies: Dominant strategy: component suppression (7/8, necessity test). Only 2/8 tested amplification. Interface attribution ranking implicitly biases toward suppression.

Perceived risks: Ripple effects / non-orthogonality (3/8); insufficient instance-level validation (3/8); over-steering (2/8).

CONCLUSION

Embedding steering within a structured workflow enables practitioners to transition from **correlational inspection** to **actionable investigation**.

TOOL

ID	Concept Examples	Attribution	Activation	Concept Description	Most relevant for Classes
19764		7.5	1.86	asymmetrical, indistinct bo...	regression structures, irregular b...
27752		6.9	1.06	melanoma, multicolored, ...	melanoma, regression structure...
28126		5.6	1.34	atypical dots or globules, ...	melanoma, blue tape, atypical d...
2020		5.4	2.31	shades of red, large size, r...	shades of red, shades of pink, r...
3860		4.9	1.36	black and brown color, re...	black and brown color, shades...