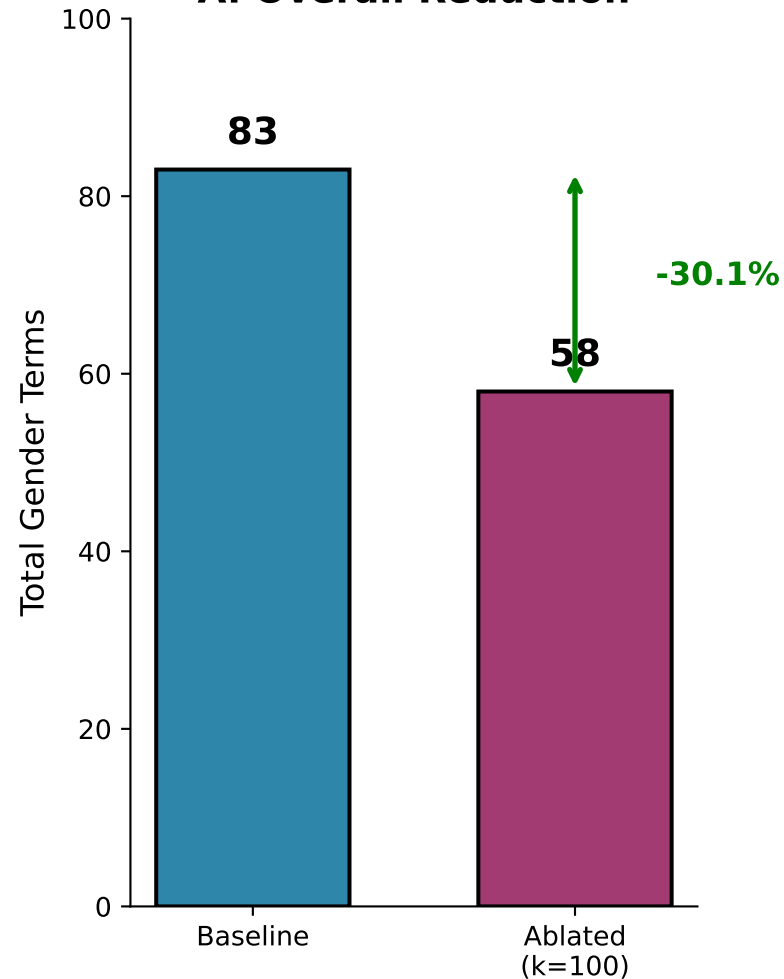
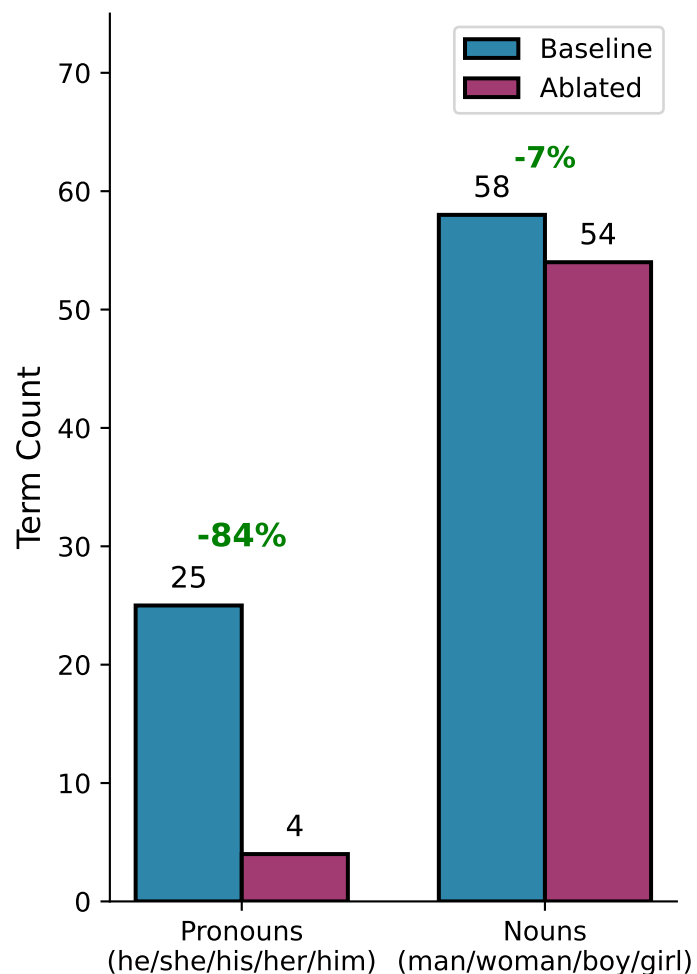


Causal Intervention: Gender Feature Ablation in PaLiGemma-3B (Layer 9)

A. Overall Reduction



B. Pronouns vs Nouns



C. Qualitative Examples

- Ex 1:** Baseline: two women standing and holding a kangaroo in their hands → Ablated: two persons are standing. Here we can see a kangaroo
women → persons
- Ex 2:** Baseline: holding a tennis racket in his hand → Ablated: holding a racket
his → [removed]
- Ex 3:** Baseline: a woman walks through a forest, her bare feet sinking → Ablated: A man is walking through a forest, crossing a river
woman/her → man

SAE feature ablation reduces pronouns by 84% while preserving semantic content