Figure 7: ROC curves for 'detection' auto-interpretability for Pythia-70m over 100 SAE latents. These results demonstrate the similarity in performance between the SAE variants, as well as the overall degradation in performance as the layer index increases.
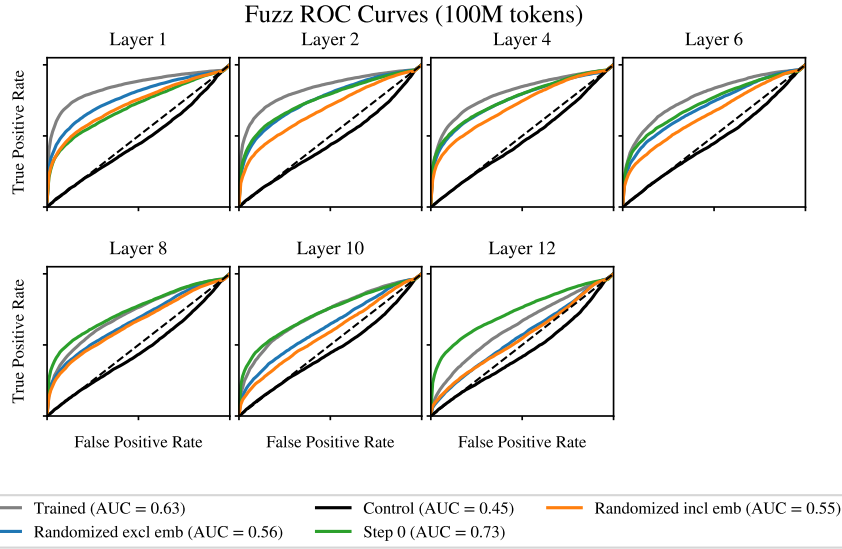
## B.2  PYTHIA 160M



Figure 8: ROC curves for 'fuzzing' auto-interpretability for Pythia-160m over 100 SAE latents. These results demonstrate the similarity in performance between the SAE variants, as well as the overall degradation in performance as the layer index increases.

Figure 9: ROC curves for 'detection' auto-interpretability for Pythia-160m over 100 SAE latents. These results demonstrate the similarity in performance between the SAE variants, as well as the overall degradation in performance as the layer index increases.
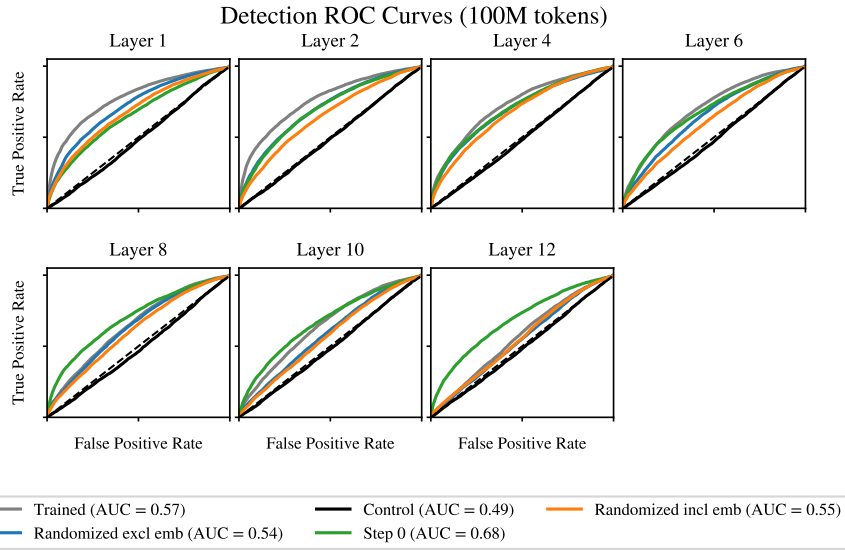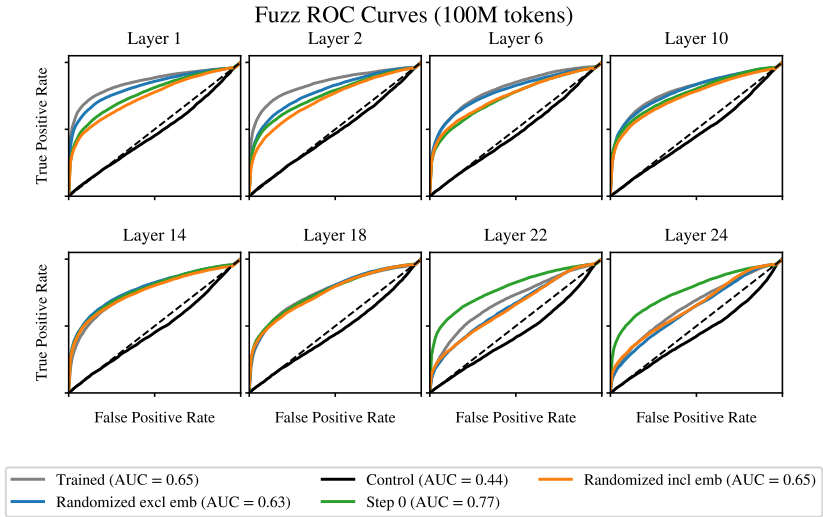
## B.3    PYTHIA 410M



Figure 10: ROC curves for 'fuzzing' auto-interpretability for Pythia-410m over 100 SAE latents. These results demonstrate the similarity in performance between the SAE variants, as well as the overall degradation in performance as the layer index increases. The auto-interpretability scores here fail to distinguish between trained and randomized models.
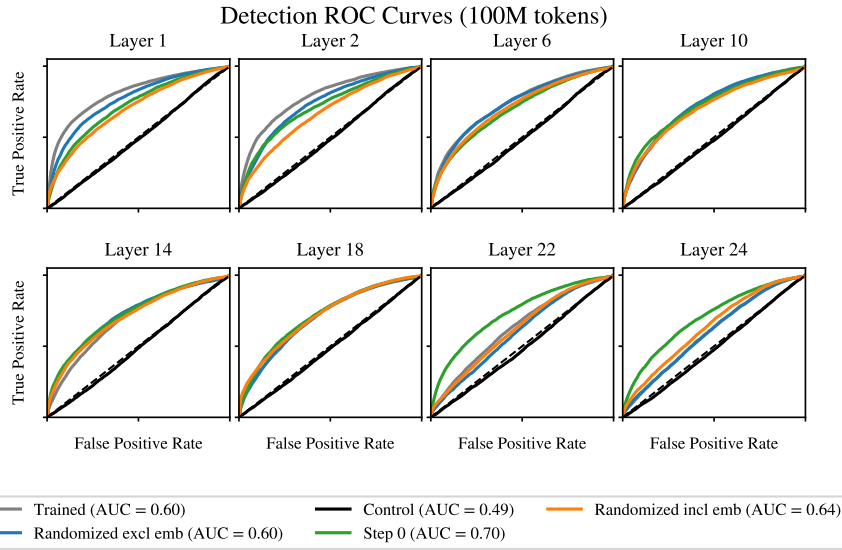
## Detection ROC Curves (100M tokens)

Layer 1 · Layer 2 · Layer 6 · Layer 10

Layer 14 · Layer 18 · Layer 22 · Layer 24

True Positive Rate / False Positive Rate

Trained (AUC = 0.60) · Control (AUC = 0.49) · Randomized incl emb (AUC = 0.64)
Randomized excl emb (AUC = 0.60) · Step 0 (AUC = 0.70)

Figure 11: ROC curves for 'detection' auto-interpretability for Pythia-410m over 100 SAE latents. These results demonstrate the similarity in performance between the SAE variants, as well as the overall degradation in performance as the layer index increases.

### B.4 PYTHIA-1B

## Fuzz ROC Curves (100M tokens)

Layer 1 · Layer 3 · Layer 5 · Layer 7

Layer 9 · Layer 11 · Layer 13 · Layer 15

True Positive Rate / False Positive Rate

Trained (AUC = 0.77) · Control (AUC = 0.49) · Randomized incl emb (AUC = 0.82)
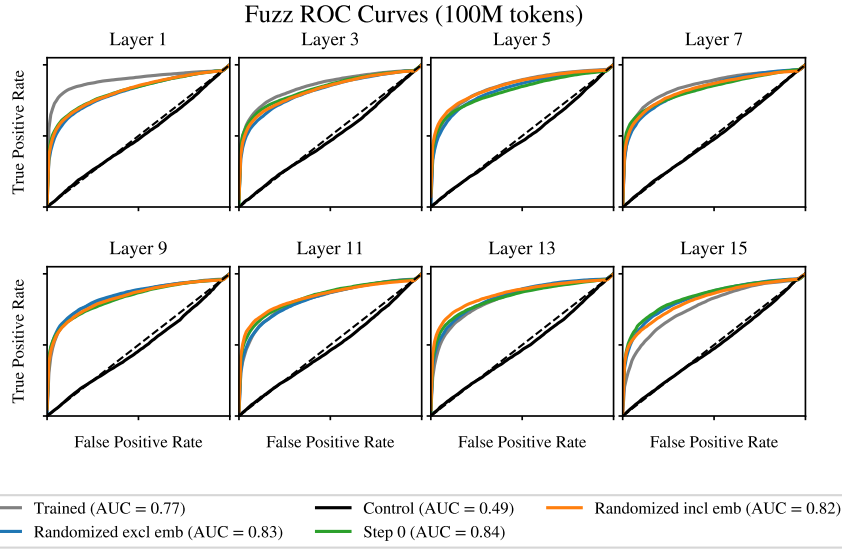Randomized excl emb (AUC = 0.83) · Step 0 (AUC = 0.84)

Figure 12: ROC curves for 'fuzzing' auto-interpretability for Pythia-1b over 100 SAE latents. These results demonstrate the similarity in performance between the SAE variants, although here we do not observe an overall degradation in quality.