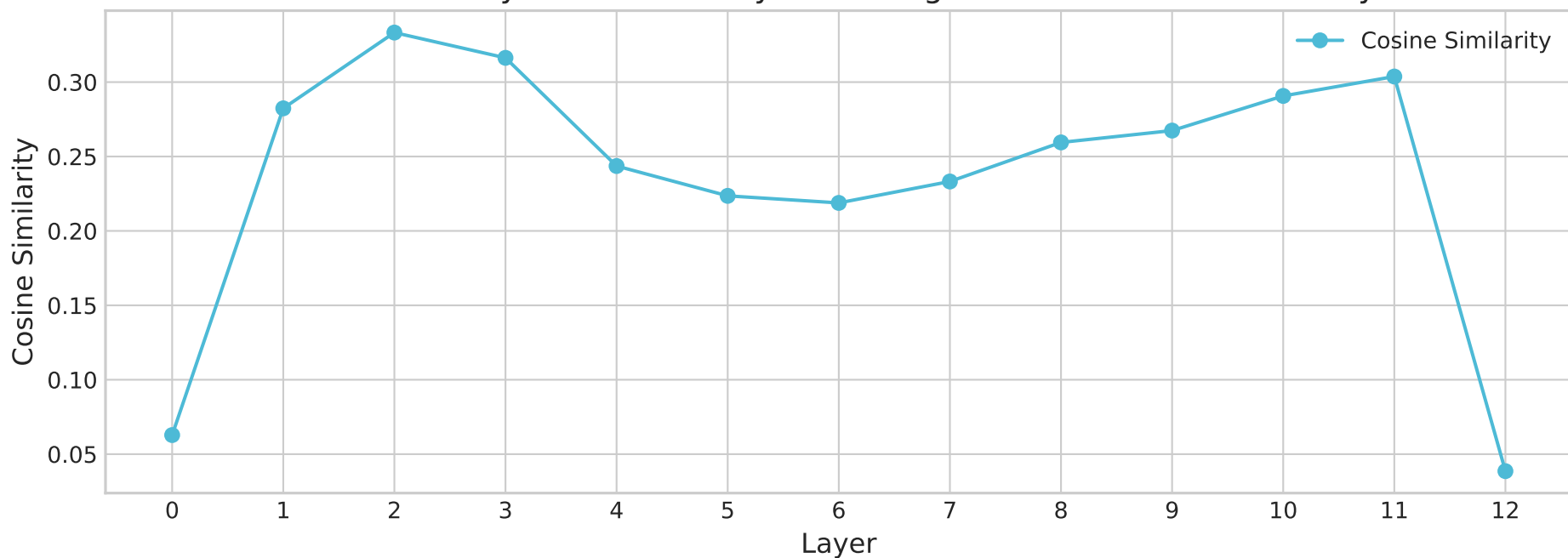Similarity Between Sticky and Benign Hidden States Across Layers

Cosine Similarity Between Sticky and Benign Hidden States Across Layers

Euclidean Distance Between Sticky and Benign Hidden States Across Layers