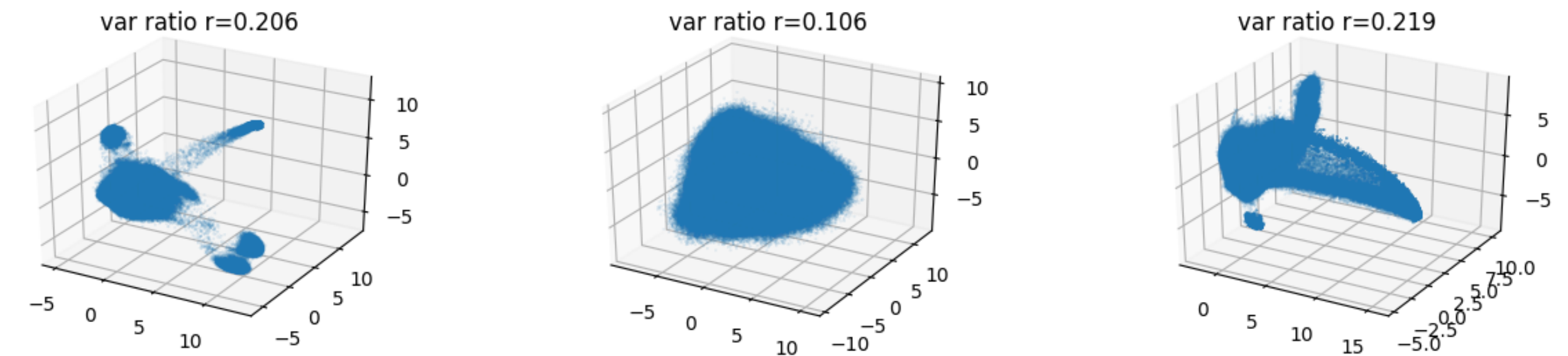


Introduction to the representation space of LLMs

- **High-dimensional data.** Embeddings form curved, twisted hypersurfaces with loops, pockets and bottlenecks.
- **Narrow-Cone Hypothesis** (Ethayarajh, 2019). Embeddings occupy a tight cone—highly anisotropic, not uniformly spread.
- **Manifold Hypothesis** (Clayton, 2015). Although embedded in high D , data lie on a much lower-dimensional manifold.



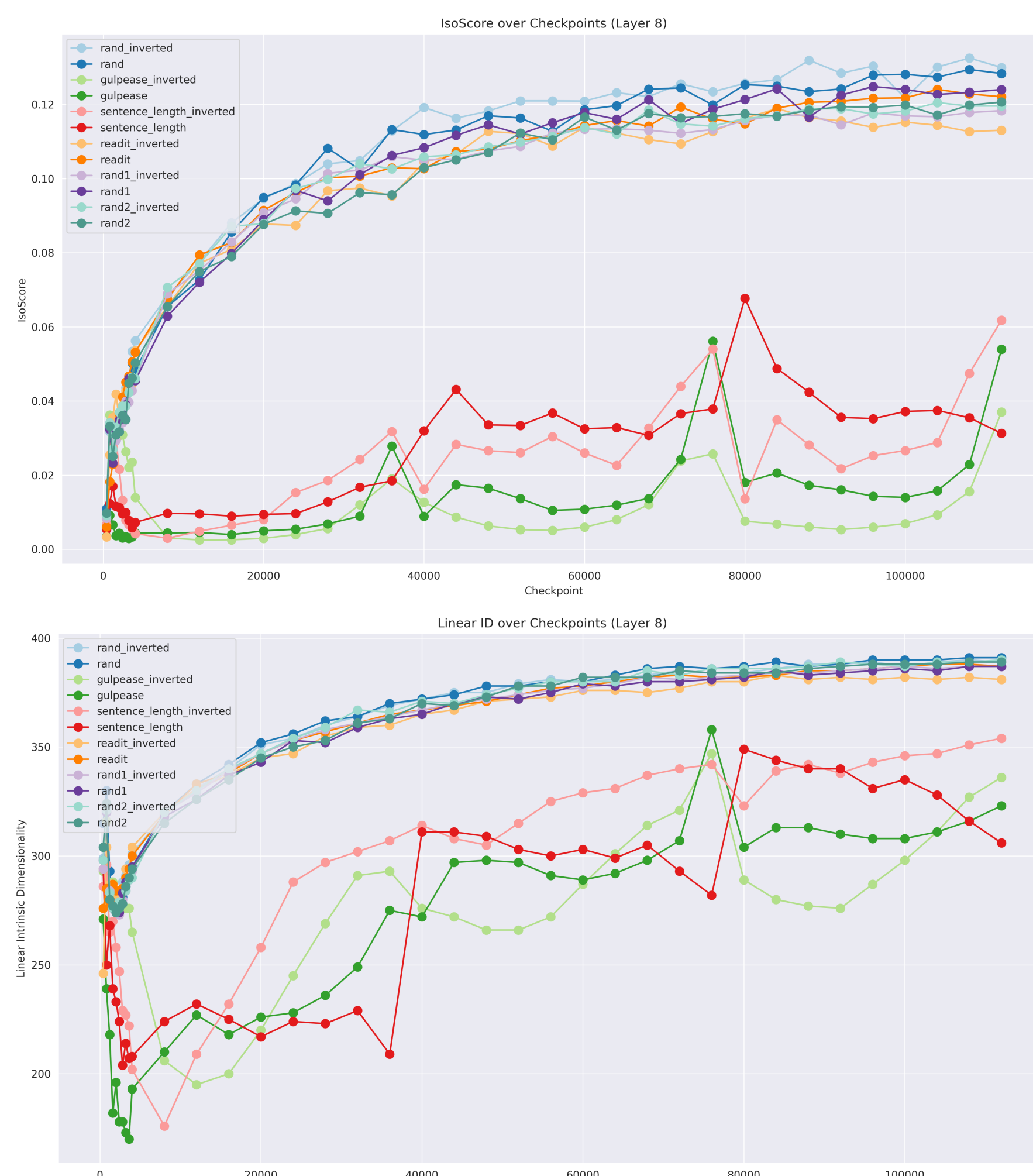
Research questions

- (i) **Linear-ID (99%):** $\text{LinearID}_{0.99} = \min\{k : \sum_{j=1}^k \frac{\lambda_j}{\lambda_D} \geq 0.99\}$
- (ii) **IsoScore:** $\tilde{\lambda}_i = \lambda_i \frac{\sqrt{D}}{\|\lambda\|_2}$, $\text{IsoScore} = 1 - \frac{\|\tilde{\lambda} - \mathbf{1}\|_2}{\sqrt{2(D - \sqrt{D})}}$

Eigenvalues $\lambda_1 \geq \dots \geq \lambda_D$ are taken from the covariance of centred embeddings.

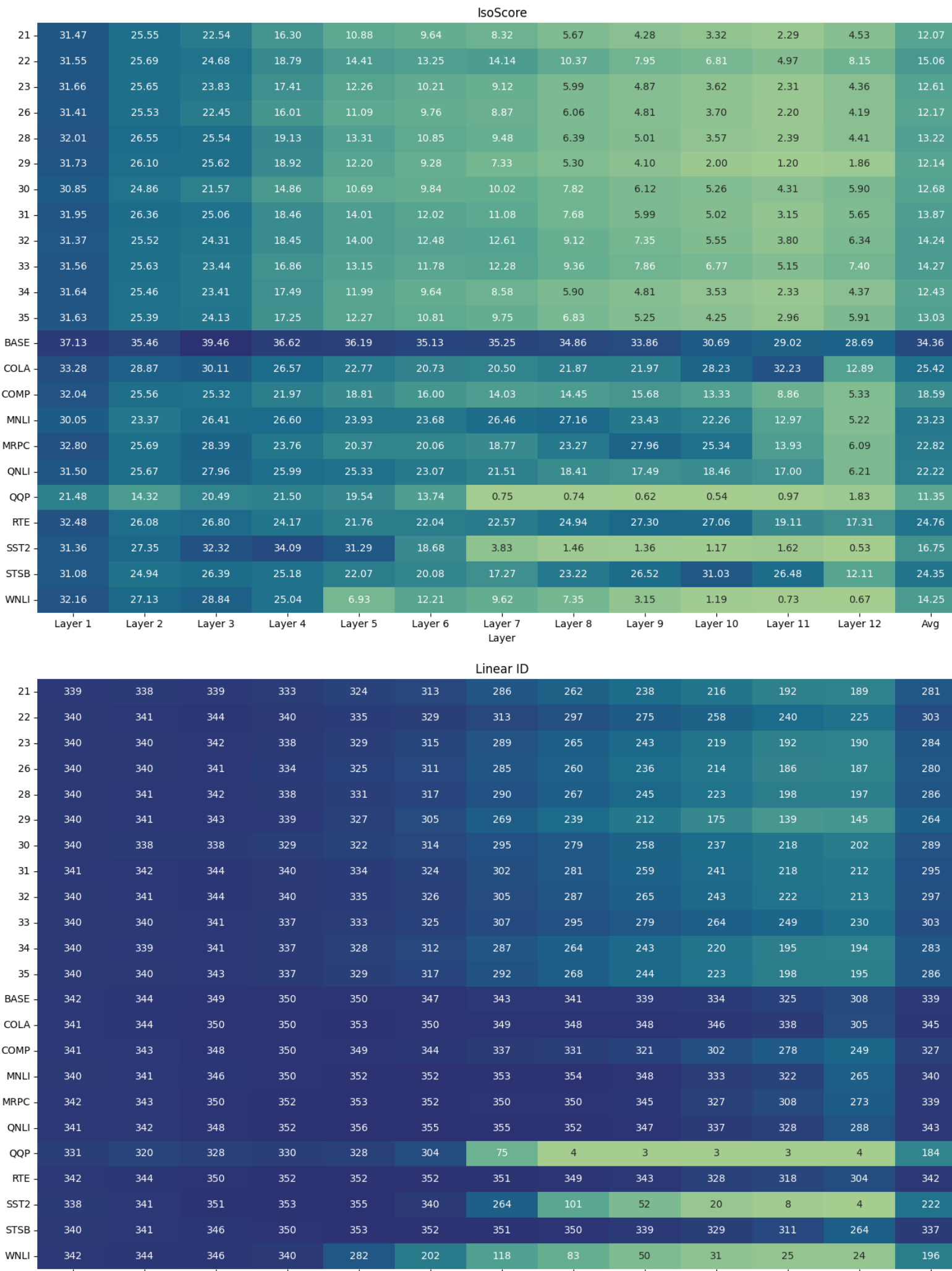
What happens during pre-training?

Curriculum learning Over checkpoints



What happens after fine-tuning?

Eye-tracking English



What happens to sub-spaces?

Part-of-speech

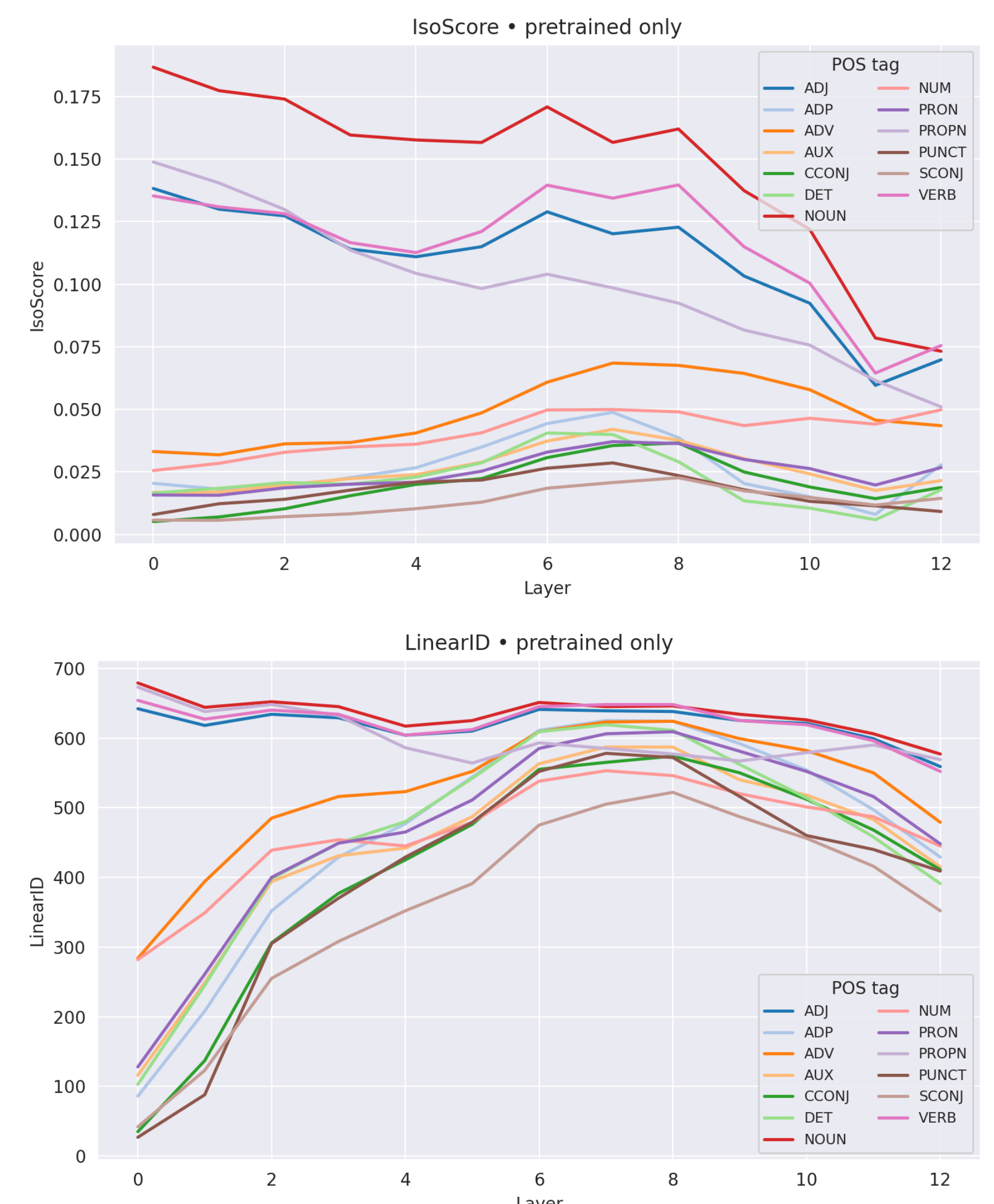


Figure 3: IsoScore and Linear-ID across layers on diverse fine-tuning tasks.

Over layers

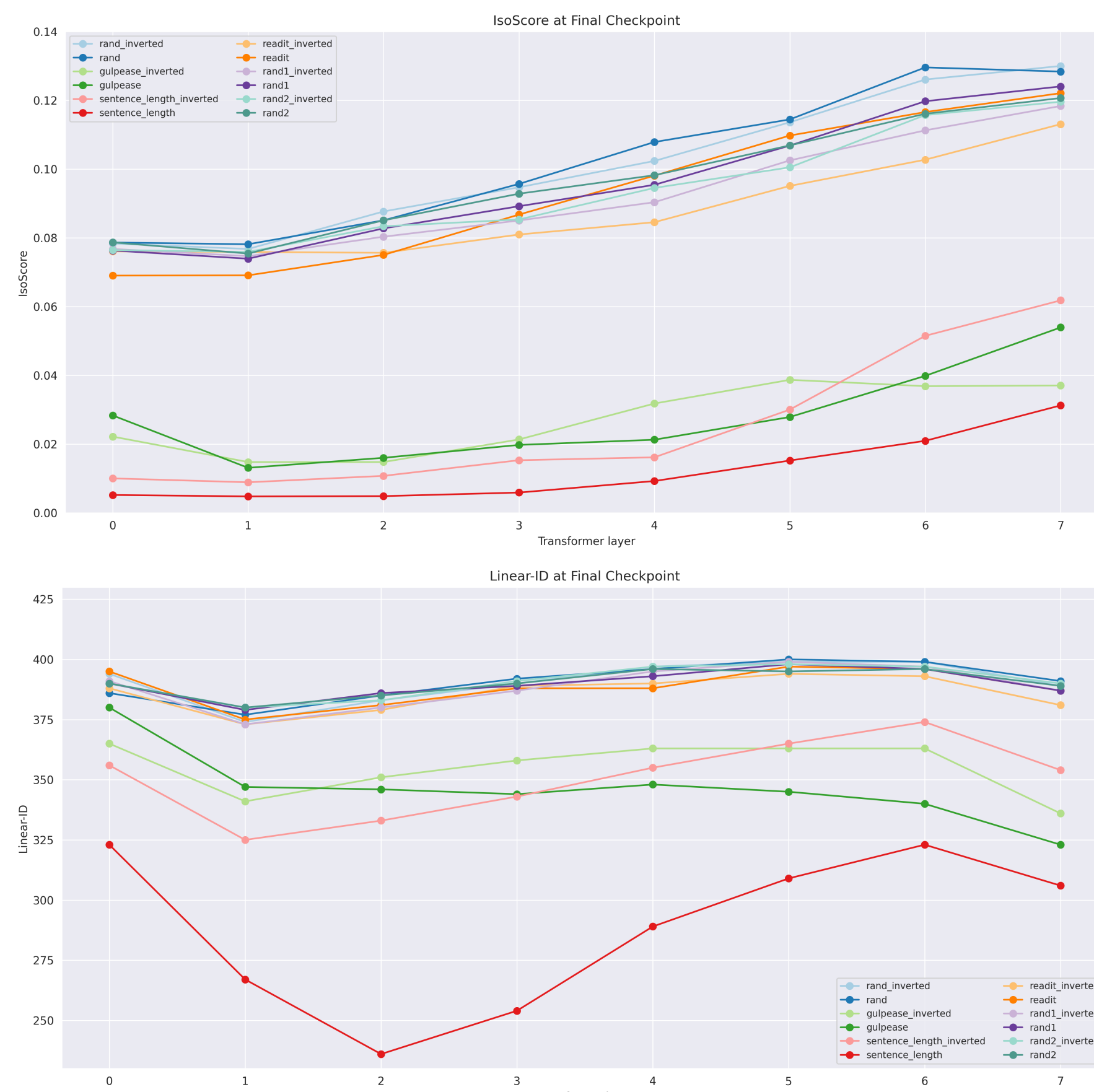


Figure 2: IsoScore and Linear-ID across checkpoint (top) and at the final checkpoint (bottom) for BERT-medium.

Italian

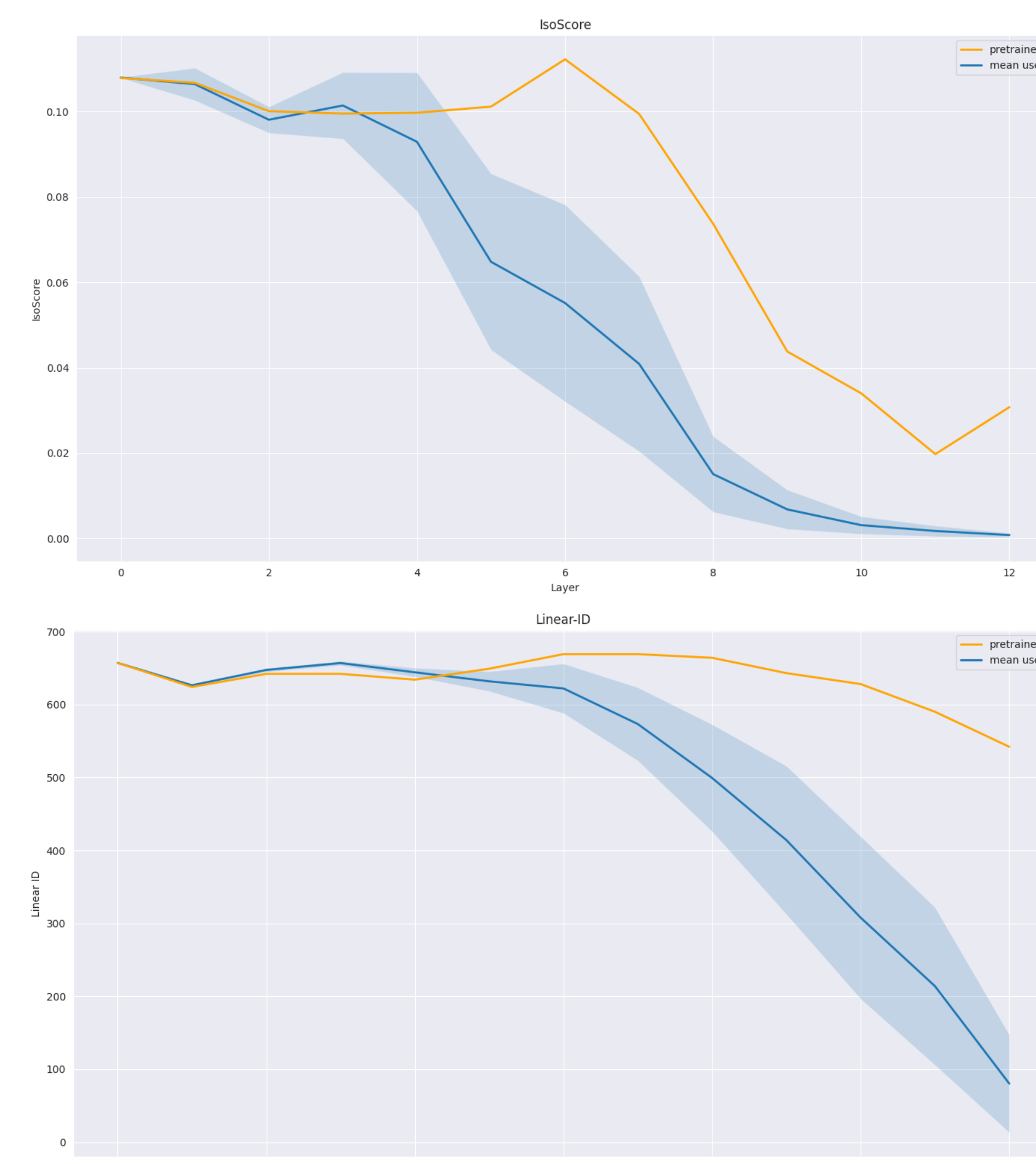


Figure 4: IsoScore and Linear-ID across layers for baseline XLM-RoBERTa and eye-tracking-fine-tuned model.

Head distance

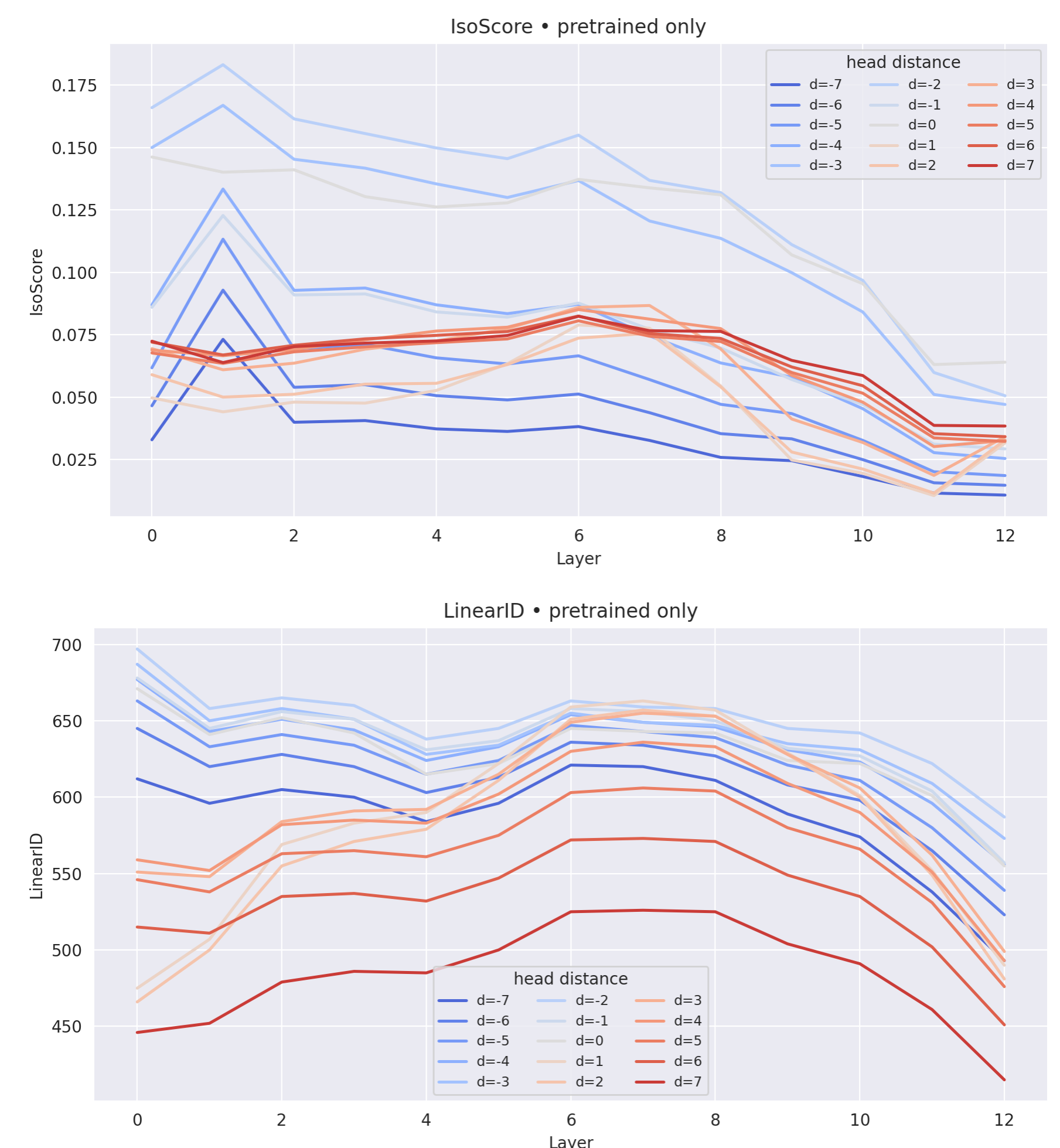


Figure 5: IsoScore and Linear-ID across layers for POS and head-dependent distance subspaces.