# Feature-Level Insights into Artificial Text Detection with Sparse Autoencoders

Kristian Kuznetsov[1,2], Laida Kushnareva[2], Polina Druzhinina[1,5], Anton Razzhigaev[1,5], Anastasia Voznyuk[3], Irina Piontkovskaya[2] Evgeny Burnaev[1,5] Serguei Barannikov[1,4]

[1]Skolkovo Institute of Science and Technology, [2]AI Foundation and Algorithm Lab, [3]Advacheck OÜ, Estonia, [4]CNRS, Université Paris Cité, France, [5]Artificial Intelligence Research Institute (AIRI)

**Skoltech**
Skolkovo Institute of Science and Technology

## Motivation

The increasing realism of LLM-generated text poses a major challenge for **artificial text detection (ATD)**. Many current ATD methods lack interpretability and robustness, leaving the specific linguistic and structural features of machine-generated text underexplored. To address this, we use **Sparse Autoencoders (SAEs) as a source of interpretable and generalizable features** for understanding and detecting AI-generated text.

## Methods

**Feature Extraction with SAEs:** we use sparse autoencoders to extract interpretable features from the residual stream of Gemma-2-2B. Given hidden activations $\mathbf{x} \in \mathbb{R}^d$, the SAE performs:
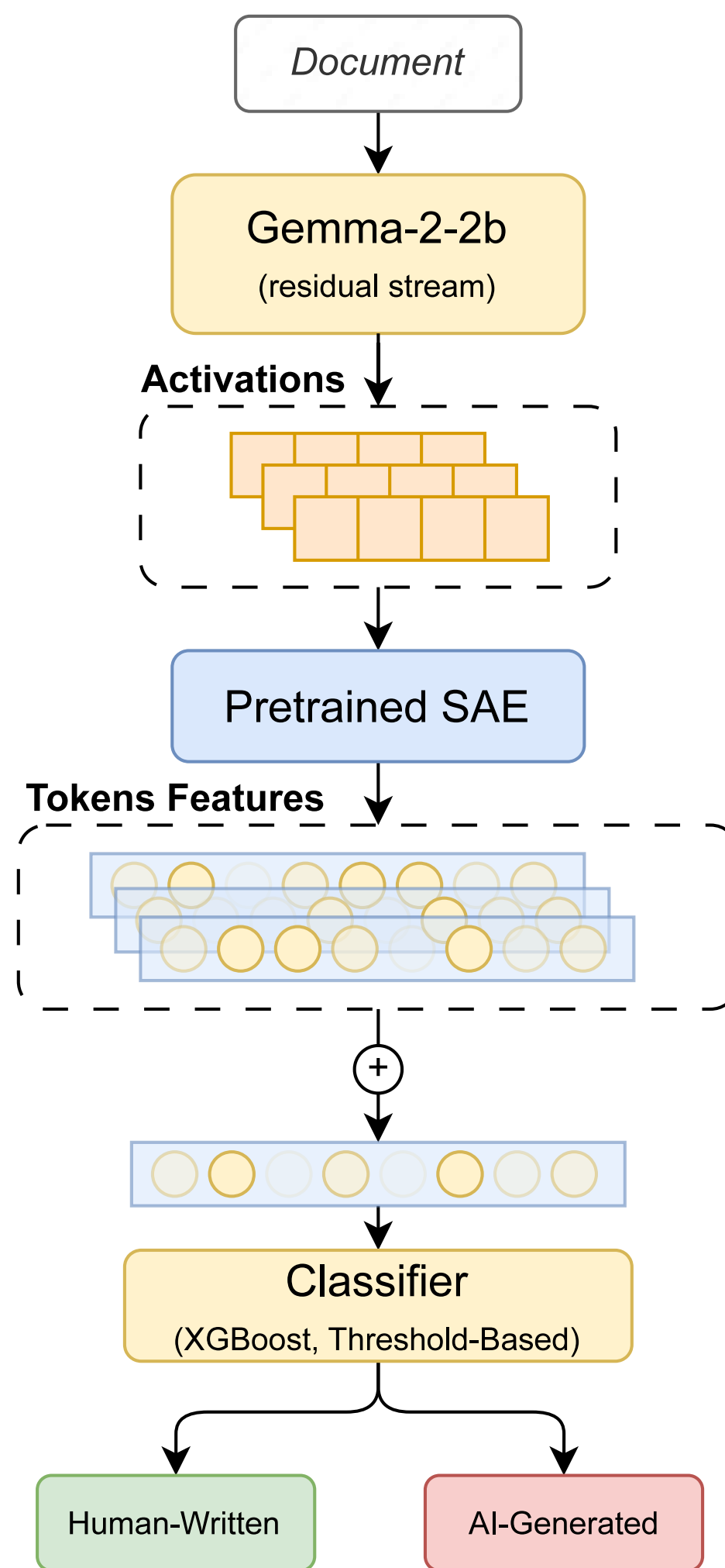
$$f(\mathbf{x}) = \sigma(\mathbf{W}_{\text{enc}}\mathbf{x} + \mathbf{b}_{\text{enc}}), \quad \hat{x}(f) = \mathbf{W}_{\text{dec}}f(\mathbf{x}) + \mathbf{b}_{\text{dec}}$$

Here, $f(\mathbf{x}) \in \mathbb{R}^M$ (with $M \gg d$) is a sparse, non-negative feature vector. For sequence-level representation: $\mathbf{f} = \sum_{i=1}^{n} f^{(l)}(\mathbf{x}_i^{(l)})$, where $\mathbf{x}_i^{(l)}$ is the residual stream at token $i$ in layer $l$.

**Classification:** we classify machine-generated and human-written texts from COLING dataset based on several stratagies: 1) **XGBoost** is trained on full feature vectors to measure global importance; 2) **Threshold classifiers** $\mathbb{I}[\mathbf{f}_j > \tau^*]$ and $\mathbb{I}[\mathbf{f}_j > 0]$ are used for analyzing individual feature activations.

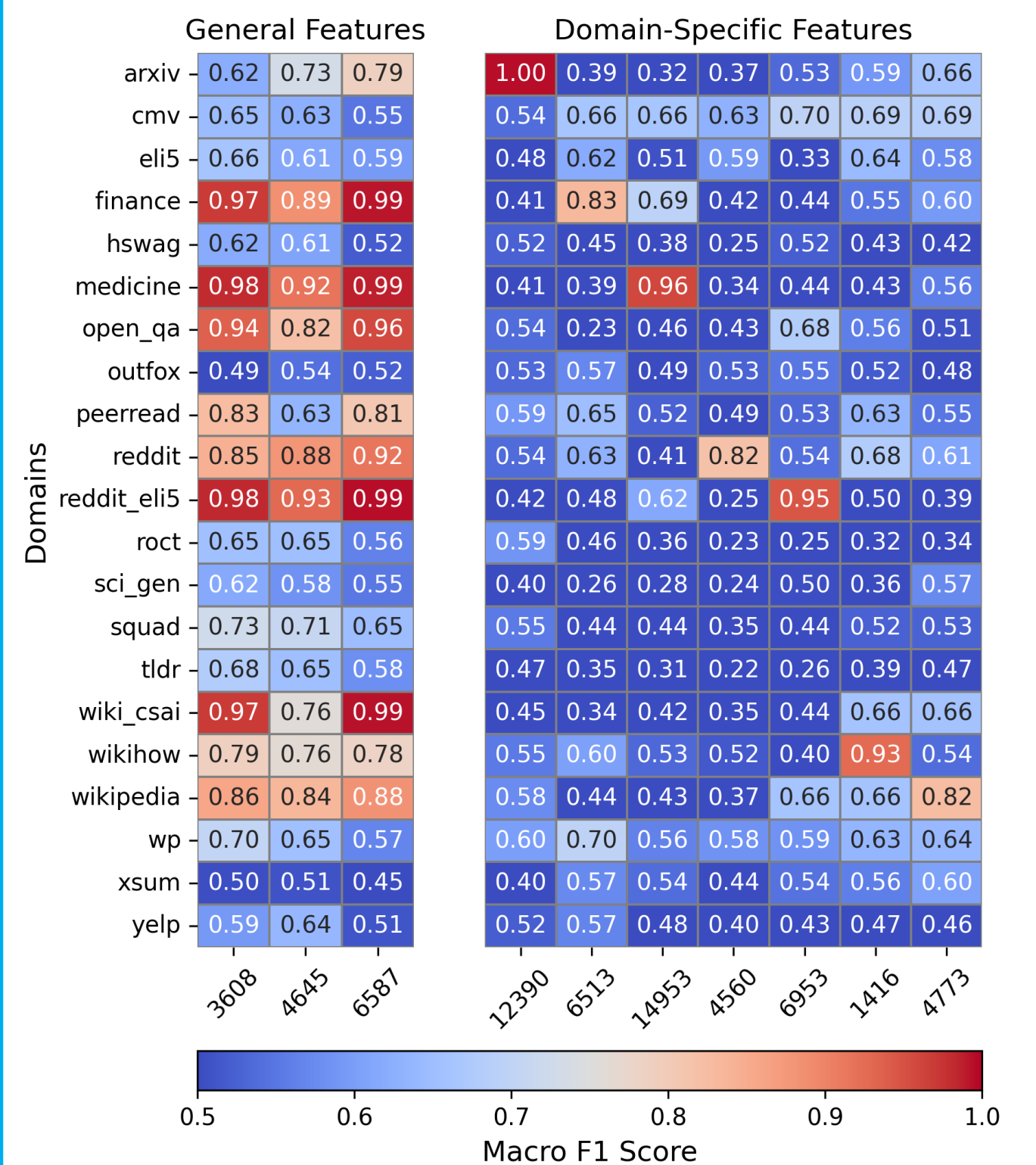**Feature Interpretation:** we interpret the top-performing features using three complementary methods:

1. **Automatic Interpretation:** use GPT-4o to describe the semantic patterns in Top-N examples.
2. **Manual Inspection:** human analysis of Top-N examples to validate common linguistic traits.
3. **Steering + Interpretation:** Use feature steering $\mathbf{x}' = \mathbf{x} + \lambda A_{\max}\mathbf{d}_i$, where $\mathbf{d}_i$ is the $i$-th column of $\mathbf{W}_{\text{dec}}$, and analyze these generations for interpretability via GPT-4o.

*Document* → Gemma-2-2b (residual stream) → **Activations** → Pretrained SAE → **Tokens Features** → ⊕ → Classifier (XGBoost, Threshold-Based) → Human-Written / AI-Generated

## Single Feature Classifiers

**General Features**: some features (e.g., 3608, 4645) generalize across domains and model families (e.g., Reddit, Wikipedia, Medicine).
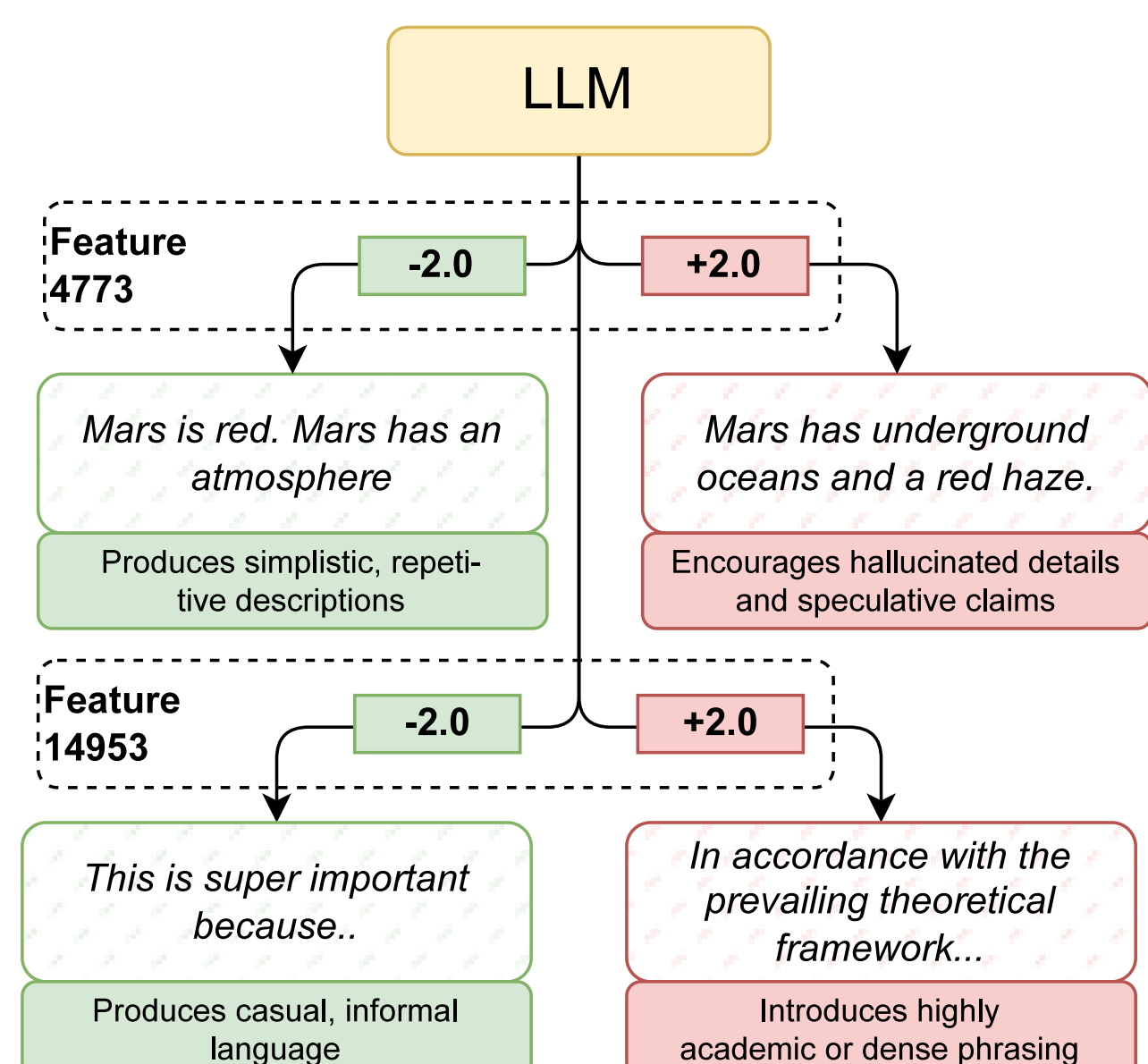
**Domain-Specific Features**: others are highly domain-specific, capturing traits like hallucinated facts in Wikipedia.

| Domains | General Features | | | Domain-Specific Features | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | 3608 | 4645 | 6587 | 12390 | 6513 | 14953 | 4560 | 6953 | 1416 | 4773 |
| arxiv | 0.62 | 0.73 | 0.79 | 1.00 | 0.39 | 0.32 | 0.37 | 0.53 | 0.59 | 0.66 |
| cmv | 0.65 | 0.63 | 0.55 | 0.54 | 0.66 | 0.66 | 0.63 | 0.70 | 0.69 | 0.69 |
| eli5 | 0.66 | 0.61 | 0.59 | 0.48 | 0.62 | 0.51 | 0.59 | 0.33 | 0.64 | 0.58 |
| finance | 0.97 | 0.89 | 0.99 | 0.41 | 0.83 | 0.69 | 0.42 | 0.44 | 0.55 | 0.60 |
| hswag | 0.62 | 0.61 | 0.52 | 0.52 | 0.45 | 0.38 | 0.25 | 0.53 | 0.43 | 0.42 |
| medicine | 0.98 | 0.92 | 0.99 | 0.41 | 0.39 | 0.96 | 0.34 | 0.44 | 0.43 | 0.59 |
| open_qa | 0.94 | 0.82 | 0.96 | 0.54 | 0.23 | 0.46 | 0.43 | 0.68 | 0.56 | 0.51 |
| outfox | 0.49 | 0.54 | 0.52 | 0.53 | 0.57 | 0.46 | 0.53 | 0.55 | 0.52 | 0.48 |
| peerread | 0.83 | 0.63 | 0.81 | 0.59 | 0.65 | 0.52 | 0.49 | 0.53 | 0.63 | 0.57 |
| reddit | 0.85 | 0.58 | 0.88 | 0.54 | 0.63 | 0.41 | 0.82 | 0.54 | 0.68 | 0.61 |
| reddit_eli5 | 0.98 | 0.93 | 0.99 | 0.42 | 0.48 | 0.62 | 0.25 | 0.95 | 0.50 | 0.39 |
| roct | 0.65 | 0.65 | 0.56 | 0.59 | 0.46 | 0.36 | 0.23 | 0.25 | 0.32 | 0.34 |
| sci_gen | 0.62 | 0.58 | 0.55 | 0.40 | 0.26 | 0.28 | 0.24 | 0.50 | 0.36 | 0.57 |
| squad | 0.73 | 0.71 | 0.65 | 0.44 | 0.53 | 0.44 | 0.54 | 0.52 | 0.53 | |
| tldr | 0.68 | 0.65 | 0.58 | 0.47 | 0.35 | 0.31 | 0.22 | 0.26 | 0.39 | 0.47 |
| wiki_csai | 0.97 | 0.76 | 0.99 | 0.45 | 0.34 | 0.42 | 0.35 | 0.44 | 0.66 | 0.66 |
| wikihow | 0.79 | 0.76 | 0.78 | 0.55 | 0.60 | 0.53 | 0.52 | 0.40 | 0.93 | 0.61 |
| wikipedia | 0.86 | 0.84 | 0.88 | 0.58 | 0.44 | 0.43 | 0.37 | 0.66 | 0.66 | 0.82 |
| wp | 0.70 | 0.65 | 0.57 | 0.60 | 0.70 | 0.56 | 0.58 | 0.63 | 0.63 | 0.64 |
| xsum | 0.50 | 0.51 | 0.45 | 0.40 | 0.57 | 0.54 | 0.44 | 0.54 | 0.56 | 0.60 |
| yelp | 0.59 | 0.64 | 0.51 | 0.52 | 0.57 | 0.48 | 0.40 | 0.43 | 0.47 | 0.46 |

Macro F1 Score: 0.5 – 1.0
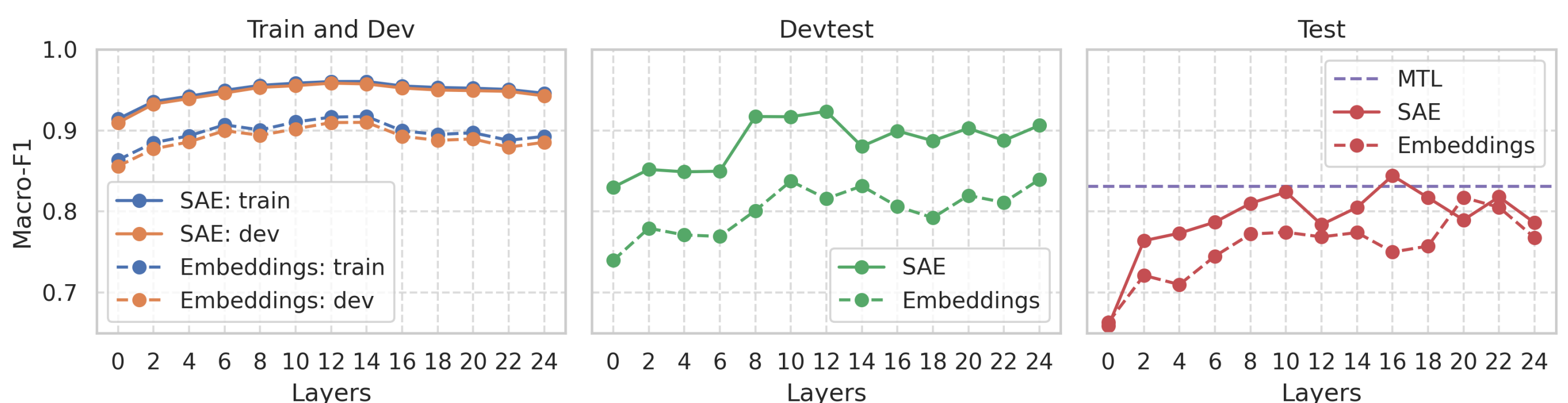
## Interpretation Insights

We identified several patterns in features and combined into groups:

1. **Common traits of AI-generated text**: *excessive complexity, assertive claims, wordy introductions, repetition, and over-formality*
2. **Domain-specific signals**: *overcomplicated syntax, hallucinated facts, speculative tone, and overly polite phrasing depending on the source*

LLM

Feature 4773: **-2.0** / **+2.0**
- *Mars is red. Mars has an atmosphere.* → Produces simplistic, repetitive descriptions
- *Mars has underground oceans and a red haze.* → Encourages hallucinated details and speculative claims

Feature 14953: **-2.0** / **+2.0**
- *This is super important because..* → Produces casual, informal language
- *In accordance with the prevailing theoretical framework...* → Introduces highly academic or dense phrasing

## Overall Detection Performance

SAE-derived features consistently outperform mean-pooled transformer embeddings on the COLING dataset across all splits (Train, Dev, DevTest, Test). On the 16th layer, SAE-based classifiers surpass even the state-of-the-art multitask learning (MTL) baseline.

Train and Dev — SAE: train, SAE: dev, Embeddings: train, Embeddings: dev

Devtest — SAE, Embeddings

Test — MTL, SAE, Embeddings

(Macro-F1 vs Layers)

## Conclusion

Sparse Autoencoders provide a powerful and interpretable alternative for detecting AI-generated text. A small number of learned features capture robust signals that **generalize across domains and models**. These features are not only effective for classification but also **human-interpretable through activation analysis, steering, and language-based explanations**. Our approach bridges performance and explainability, enabling more transparent and reliable AI-generated text detection.

## Paper Link

For more information you can check paper via QR-code or contact me:
**Telegram:** @pyashy
**Email:** kris@kuznetsov.su