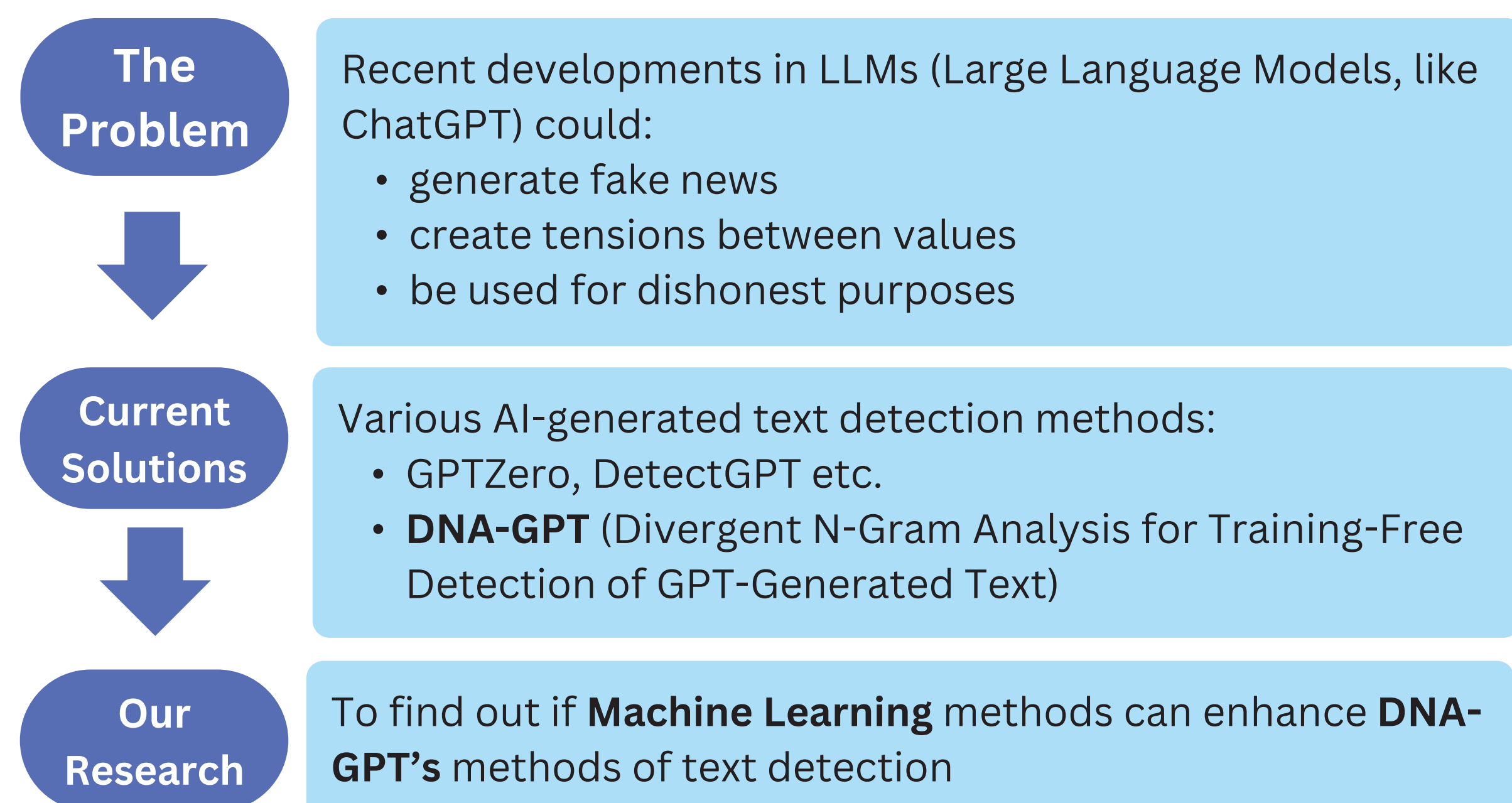


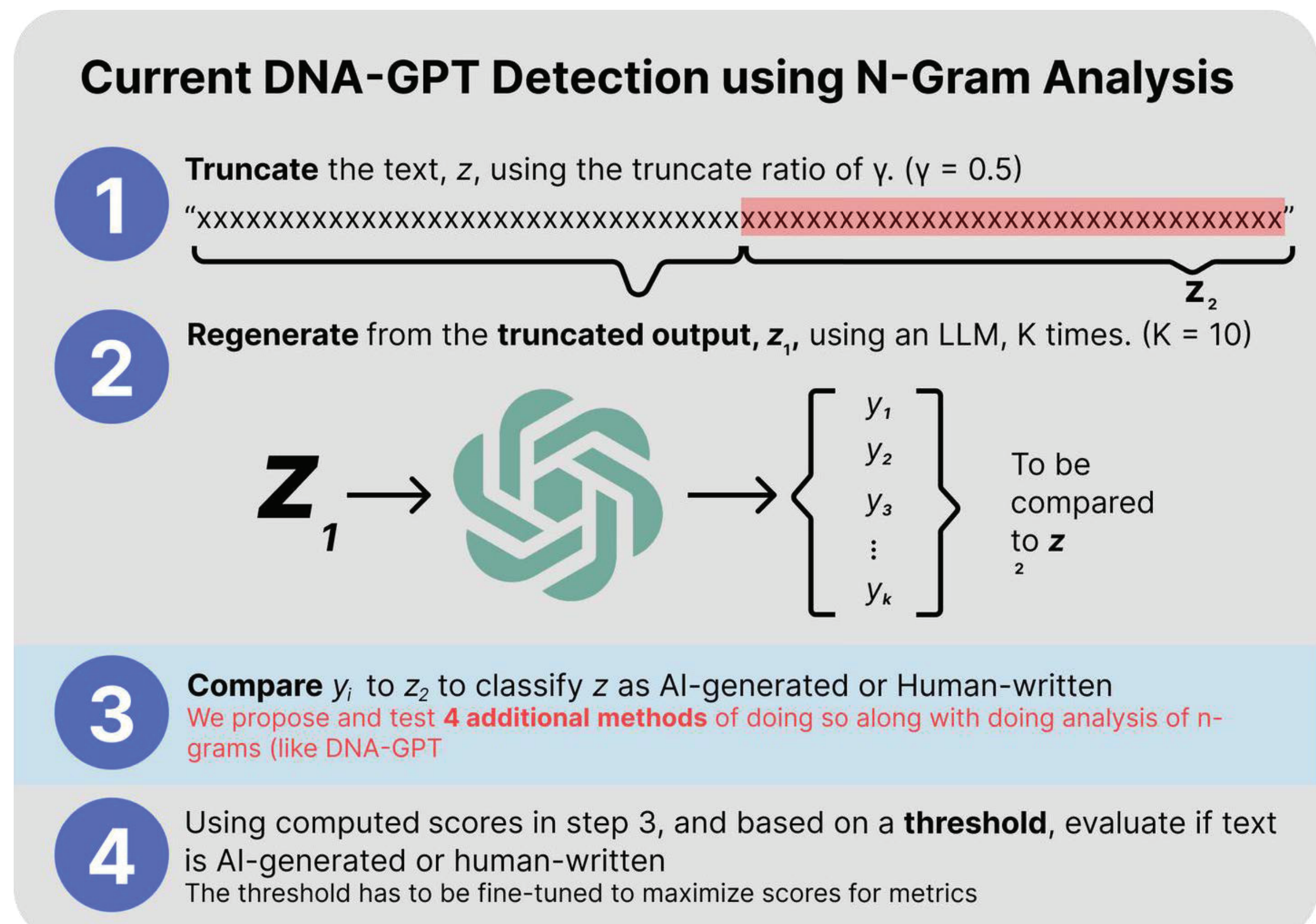
Introduction



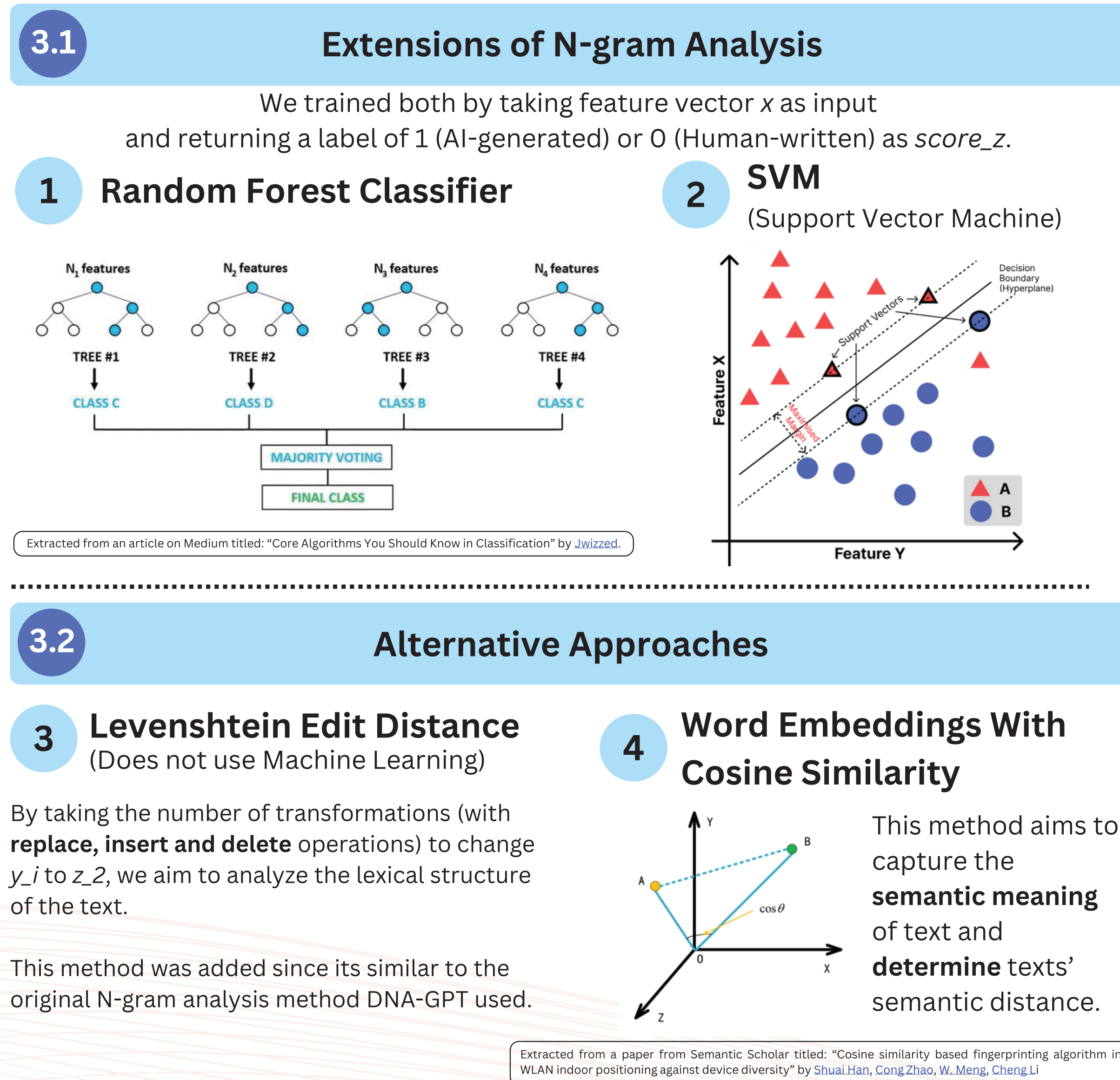
Scope & Methodology

4 Methods

Stemming from the DNA-GPT's original methods, we propose 4 methods to enhance their text-detection:

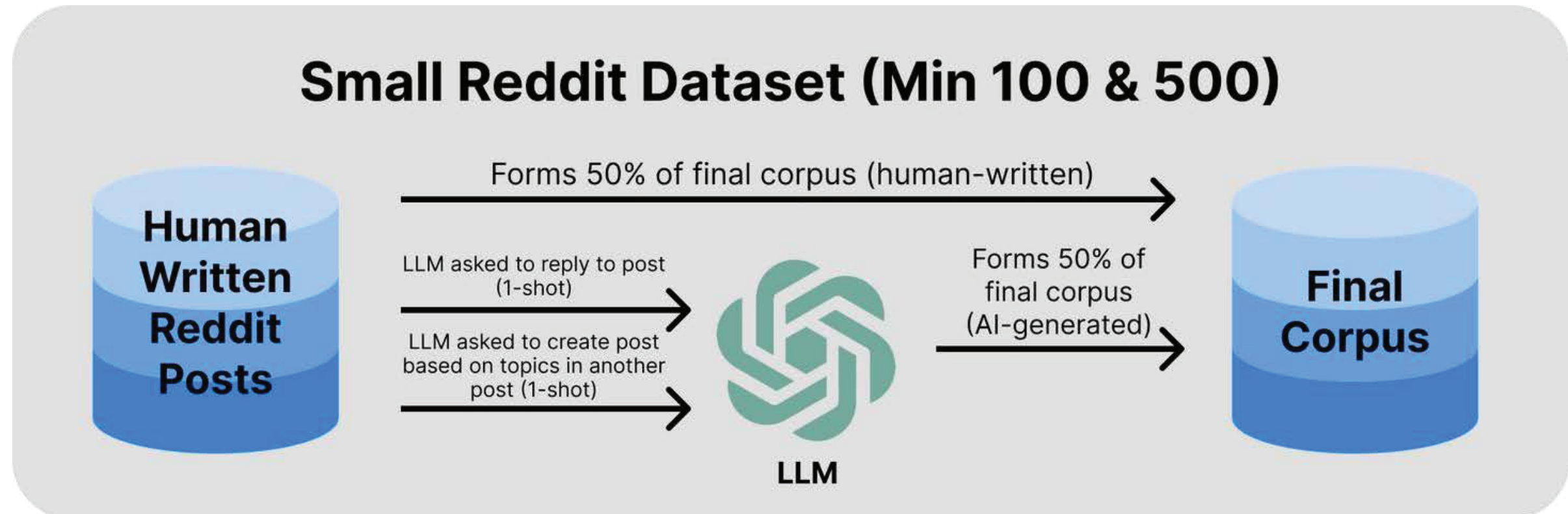


We grouped the 4 methods to improve step 3 into 2 groups:



3 Datasets

2 datasets with differing minimum word counts were pruned and generated as follows:



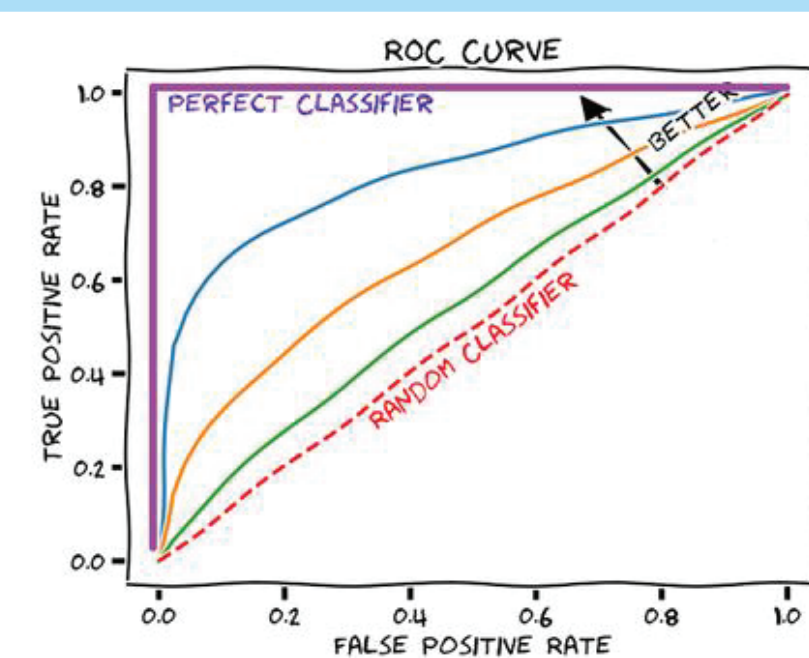
The final dataset, **ELI5 (Min 500)**, was constructed as follows:

- Human Section:** Replies to questions from the Explain Like I'm 5 community
- AI Section:** Generated replies to the same questions

It was also used by DNA-GPT, allowing us to compare our results with DNA-GPT's.

2 Metrics

AUROC (Area Under Receiver Operating Characteristic Curve)



TPR (True Positive Rate) at 1% FPR (False Positive Rate)

This metric was used to ensure the **reliability** of detection algorithms for real-life deployment since it is crucial to maintain a **high TPR** while **minimizing the FPR**.

Both metrics were used by DNA-GPT, allowing for effective comparison of results.

1 Model, Baseline and Detection Scenario

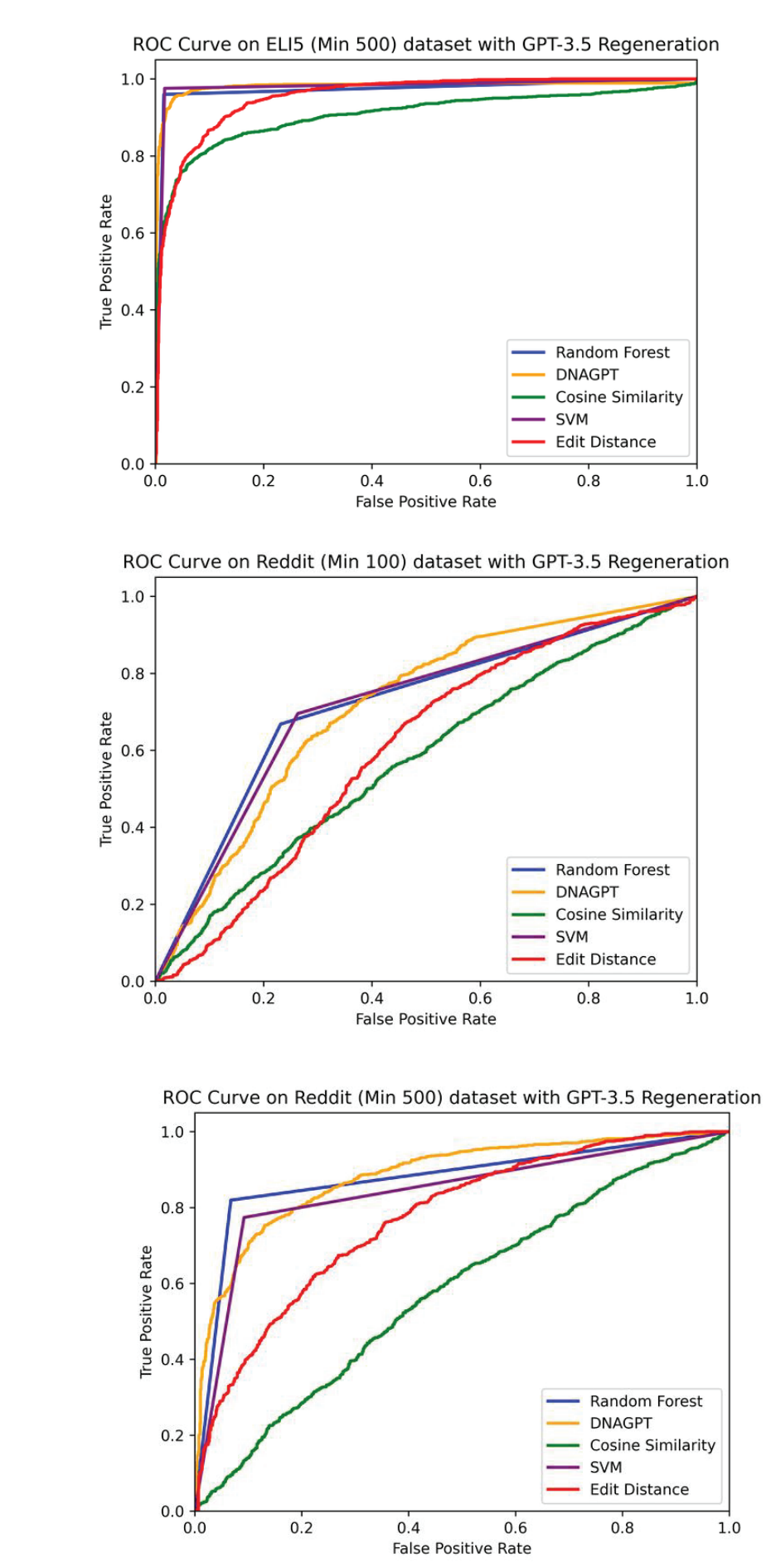
Due to time constraints, only one type of each were experimented with.

Model: GPT-3.5-Turbo, **Baseline:** DNA-GPT's original method, **Detection Scenario:** Black-Box

Results & Discussion

Scores for **AUROC** and **TPR** metrics for all 3 datasets and 5 methods (including DNA-GPT's original score), with the best-scores (with a margin of 1%) **bolded red**

| Datasets | ELI5 (Min 500) | | Reddit Small (Min 500) | | Reddit Small (Min 100) | |
|--------------------|----------------|--------------|------------------------|--------------|------------------------|-------------|
| | AUROC | TPR | AUROC | TPR | AUROC | TPR |
| DNA-GPT (original) | 96.85 | 63.50 | - | - | - | - |
| DNA-GPT | 98.07 | 59.08 | 88.32 | 8.06 | 71.50 | 1.36 |
| Random Forest | 97.20 | 61.04 | 87.60 | 12.16 | 71.83 | 2.89 |
| SVM | 97.91 | 56.78 | 84.09 | 8.04 | 71.62 | 2.64 |
| Cosine Similarity | 90.75 | 38.41 | 58.07 | 1.11 | 57.58 | 1.13 |
| Edit Distance | 95.45 | 34.19 | 77.29 | 3.96 | 60.58 | 0.78 |



- DNA-GPT**
 - Able to **closely match** their results with our attempt at replicating DNA-GPT
 - Original N-Gram Analysis methods proved to be **extremely competitive**
- Random Forest Classifier and SVM**
 - Random forest classifier consistently **performs the best** among the 5 methods
 - SVM has inferior performance, possibly because Random forest classifier relies on **multiple models**
 - However, both Random forest classifier and SVM require training, thus **with a larger dataset size**, improved results can definitely be achieved
- Cosine similarity with Word Embeddings**
 - Performed **much worse** than other methods
 - Semantic meaning** of y_i is very likely to match z_2 , with z_1 as context, irrelevant of whether z is human-written or AI-generated
 - Scores **unable to properly differentiate** AI-generated from human-written
- Edit Distance vs N-Gram Analysis**
 - Lexical analysis of edit distance is **much less effective** than n-gram analysis
 - Edit distance is **unable to detect** words with similar spelling could have completely different meanings (eg. "Stationary" vs "Stationery")

Acknowledgements & References

[1] Yang, X. (2023, May 27). DNA-GPT: Divergent N-Gram Analysis for Training-Free Detection of GPT-Generated Text.
 [2] Angela Fan, Yacine Jernite, Ethan Perez, David Grangier, Jason Weston, and Michael Ruli. ELI5: Long form question answering. In Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics, pp. 3558-3567
 [3] Eric Mitchell, Yooho Lee, Alexander Khazatsky, Christopher D Manning, and Chelsea Finn. Detectgpt: Zero-shot machine-generated text detection using probability curvature. arXiv preprint arXiv:2301.12305, 2023.
 [4] Kalpesh Krishna, Yiyao Song, Marzena Karpinska, John Wieting, and Mohit Iyyer. Paraphrasing evades detectors of ai-generated text, but retrieval is an effective defense. arXiv preprint arXiv:2303.13408, 2023