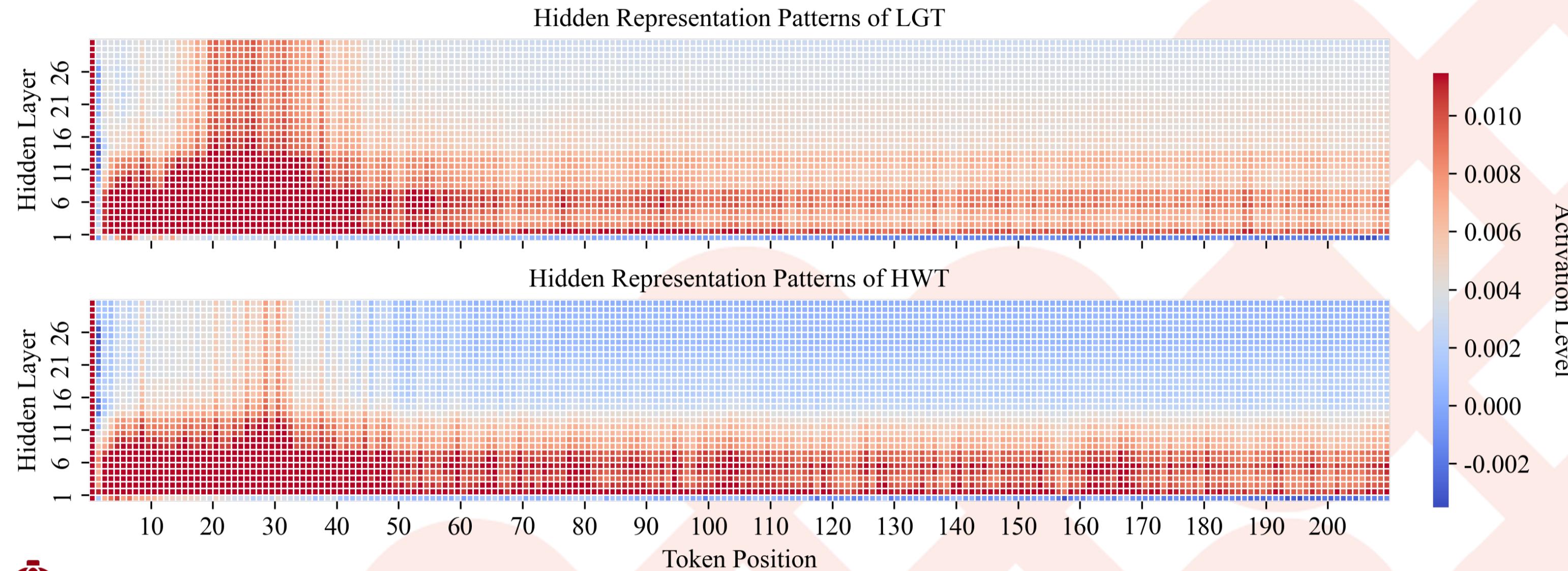


Background and Motivation

Background: Academic dishonesty, the spread of fake news, user mistrust in digital content, data pollution in web-sourced data...

Motivation: Current methods rely on large amounts of training data or single dimension with poor robustness across models and domains.



Observation: LLMs (Large Language Models) exhibit distinguishable hidden representation patterns when processing human-written texts (HWT) and LLM-generated texts (LGT), offering a promising method for achieving more robust LGT detection.

Main Result

Best Average AUROC (AUR.) And TPR@0.01 (TPR.) Performance on In-domain and Out-of-domain¹

| Test → Train ↓ | Detector ↓ Metrics → | ChatGPT | AUR. | TPR. | Llama-2-70b | AUR. | TPR. | Google-PaLM | AUR. | TPR. | Claude-instant | AUR. | TPR. | Avg. | TPR. |
|----------------|----------------------|------------------------------|------------------------------|------------------------------|------------------------------|------------------------------|------------------------------|------------------------------|------------------------------|------------------------------|------------------------------|------|------|------|------|
| ChatGPT | Roberta | 98.38 _{±0.32} | 90.52 _{±1.93} | 81.71 _{±2.27} | 54.64 _{±16.57} | 74.56 _{±0.89} | 20.20 _{±12.87} | 66.74 _{±1.49} | 22.36 _{±13.26} | 80.35 _{±1.24} | 46.93 _{±11.16} | | | | |
| | LRR | 92.61 _{±0.39} | 26.20 _{±0.00} | 95.84 _{±0.34} | 86.20 _{±0.00} | 81.98 _{±0.14} | 39.80 _{±0.00} | 57.82 _{±0.80} | 0.10 _{±0.00} | 82.06 _{±0.42} | 38.07 _{±0.00} | | | | |
| | DetectGPT | 54.87 _{±0.25} | 0.11 _{±0.14} | 59.21 _{±0.53} | 0.66 _{±1.77} | 55.29 _{±0.37} | 0.08 _{±0.00} | 55.92 _{±0.61} | 0.00 _{±0.00} | 56.32 _{±0.40} | 0.21 _{±0.49} | | | | |
| | Fast-Detect. | 75.65 _{±0.28} | 11.60 _{±0.00} | 85.49 _{±0.20} | 2.50 _{±0.00} | 80.36 _{±0.63} | 17.70 _{±0.00} | 47.29 _{±0.44} | 0.00 _{±0.00} | 72.20 _{±0.42} | 7.95 _{±0.00} | | | | |
| | Str-Detect. | 52.05 _{±0.00} | 0.01 _{±0.00} | 52.30 _{±0.00} | 0.01 _{±0.00} | 56.05 _{±0.00} | 0.01 _{±0.00} | 57.50 _{±0.00} | 0.01 _{±0.00} | 54.47 _{±0.40} | 0.01 _{±0.00} | | | | |
| | Binoculars | 97.36 _{±0.52} | 40.70 _{±0.00} | 99.45 _{±0.44} | 98.70 _{±0.00} | 98.03_{±0.34} | 89.80_{±0.00} | 61.86 _{±3.64} | 3.40 _{±0.00} | 89.18 _{±0.24} | 58.15 _{±0.00} | | | | |
| | RepreGuard | 99.84 _{±0.12} | 100.00 _{±0.00} | 99.55_{±0.12} | 99.26_{±0.11} | 88.67 _{±0.65} | 64.66 _{±1.28} | 85.00_{±1.40} | 56.12_{±2.20} | 93.26_{±0.57} | 80.01_{±0.89} | | | | |
| Llama-2-70b | Roberta | 95.49 _{±0.07} | 71.16 _{±4.59} | 94.21 _{±2.06} | 63.66 _{±10.40} | 86.00 _{±2.17} | 43.60 _{±2.56} | 76.98 _{±0.40} | 22.88 _{±5.13} | 88.17 _{±2.31} | 50.32 _{±6.67} | | | | |
| | LRR | 91.88 _{±1.05} | 26.20 _{±0.00} | 96.47 _{±0.52} | 86.20 _{±0.00} | 81.65 _{±0.55} | 39.80 _{±0.00} | 56.20 _{±1.77} | 0.10 _{±0.00} | 81.55 _{±0.97} | 38.07 _{±0.00} | | | | |
| | DetectGPT | 54.33 _{±0.81} | 0.51 _{±0.80} | 59.36 _{±0.86} | 1.28 _{±2.18} | 55.02 _{±0.98} | 0.02 _{±0.00} | 56.11 _{±0.24} | 0.00 _{±0.00} | 56.20 _{±0.72} | 0.45 _{±0.76} | | | | |
| | Fast-Detect. | 94.08 _{±1.26} | 71.90 _{±0.00} | 98.76 _{±0.05} | 94.20 _{±0.00} | 92.39 _{±0.28} | 81.80 _{±0.00} | 51.45 _{±0.62} | 0.40 _{±0.00} | 84.17 _{±0.55} | 62.08 _{±0.00} | | | | |
| | Str-Detect. | 52.05 _{±0.00} | 0.01 _{±0.00} | 52.30 _{±0.00} | 0.01 _{±0.00} | 56.05 _{±0.00} | 0.01 _{±0.00} | 57.50 _{±0.00} | 0.01 _{±0.00} | 54.47 _{±0.40} | 0.01 _{±0.00} | | | | |
| | Binoculars | 97.94 _{±0.74} | 85.90 _{±0.00} | 99.55_{±0.06} | 98.70 _{±0.00} | 97.23_{±0.93} | 89.80_{±0.00} | 57.49 _{±3.57} | 3.40 _{±0.00} | 88.07 _{±1.32} | 69.45 _{±0.00} | | | | |
| | RepreGuard | 99.54_{±0.00} | 99.38_{±0.14} | 99.38 _{±0.18} | 98.84_{±0.11} | 88.84 _{±1.28} | 77.08 _{±0.45} | 84.08_{±3.52} | 60.66_{±1.66} | 92.96_{±1.27} | 83.99_{±0.59} | | | | |
| Google-PaLM | Roberta | 88.72 _{±1.35} | 40.94 _{±5.18} | 85.36 _{±5.00} | 32.70 _{±6.15} | 82.09 _{±3.99} | 36.52 _{±4.64} | 72.98 _{±4.00} | 21.54 _{±8.64} | 82.29 _{±3.58} | 32.92 _{±6.15} | | | | |
| | LRR | 91.40 _{±2.01} | 26.20 _{±0.00} | 92.84 _{±2.78} | 86.20 _{±0.00} | 83.13 _{±1.53} | 39.80 _{±0.00} | 61.11 _{±2.05} | 0.10 _{±0.00} | 82.12 _{±2.00} | 38.07 _{±0.00} | | | | |
| | DetectGPT | 54.04 _{±0.32} | 0.00 _{±0.00} | 59.21 _{±0.22} | 0.00 _{±0.00} | 55.30 _{±0.35} | 0.02 _{±0.00} | 55.53 _{±4.46} | 0.96 _{±1.63} | 56.02 _{±0.34} | 0.24 _{±0.42} | | | | |
| | Fast-Detect. | 74.62 _{±0.99} | 11.60 _{±0.00} | 86.47 _{±0.50} | 2.50 _{±0.00} | 80.64 _{±0.22} | 17.70 _{±0.00} | 48.46 _{±4.47} | 0.00 _{±0.00} | 72.55 _{±0.54} | 7.95 _{±0.00} | | | | |
| | Str-Detect. | 52.05 _{±0.00} | 0.01 _{±0.00} | 52.30 _{±0.00} | 0.01 _{±0.00} | 56.05 _{±0.00} | 0.01 _{±0.00} | 57.50 _{±0.00} | 0.01 _{±0.00} | 54.47 _{±0.40} | 0.01 _{±0.00} | | | | |
| | Binoculars | 98.56_{±0.35} | 85.90 _{±0.00} | 99.55_{±0.20} | 98.70 _{±0.00} | 97.98_{±0.38} | 89.80_{±0.00} | 61.15 _{±2.49} | 3.40 _{±0.00} | 89.32 _{±0.86} | 69.45 _{±0.00} | | | | |
| | RepreGuard | 98.39 _{±0.31} | 99.36_{±0.07} | 98.53 _{±0.28} | 99.16_{±0.07} | 93.36 _{±0.27} | 79.88 _{±3.43} | 90.57_{±1.06} | 56.90_{±2.20} | 95.21_{±0.48} | 83.82_{±1.44} | | | | |
| Claude-instant | Roberta | 88.63 _{±1.34} | 25.80 _{±11.74} | 77.45 _{±4.73} | 23.70 _{±10.19} | 74.90 _{±1.76} | 14.58 _{±13.14} | 88.07 _{±3.89} | 36.56 _{±4.83} | 82.26 _{±2.98} | 25.84 _{±9.98} | | | | |
| | LRR | 86.33 _{±4.23} | 26.20 _{±0.00} | 87.46 _{±4.43} | 86.20 _{±0.00} | 81.19 _{±3.34} | 39.80 _{±0.00} | 63.74 _{±1.07} | 0.10 _{±0.00} | 79.68 _{±3.27} | 38.07 _{±0.00} | | | | |
| | DetectGPT | 53.95 _{±0.84} | 0.02 _{±0.06} | 57.90 _{±1.09} | 0.00 _{±0.00} | 54.33 _{±1.13} | 0.06 _{±0.07} | 55.97 _{±3.38} | 0.00 _{±0.00} | 55.54 _{±0.86} | 0.02 _{±0.03} | | | | |
| | Fast-Detect. | 81.27 _{±5.40} | 71.90 _{±0.00} | 82.35 _{±4.48} | 94.20 _{±0.00} | 80.85 _{±0.00} | 61.35 _{±0.35} | 0.40 _{±0.00} | 76.46 _{±3.73} | 62.08 _{±0.00} | | | | | |
| | Str-Detect. | 52.05 _{±0.00} | 0.01 _{±0.00} | 52.30 _{±0.00} | 0.01 _{±0.00} | 56.05 _{±0.00} | 0.01 _{±0.00} | 57.50 _{±0.00} | 0.01 _{±0.00} | 54.47 _{±0.40} | 0.01 _{±0.00} | | | | |
| | Binoculars | 92.36 _{±1.95} | 85.90 _{±0.00} | 93.18 _{±1.09} | 98.70_{±0.00} | 92.83_{±1.18} | 89.80_{±0.00} | 73.92 _{±1.97} | 3.40 _{±0.00} | 88.07 _{±1.00} | 69.45 _{±0.00} | | | | |
| | RepreGuard | 97.20 | | | | | | | | | | | | | |