

## Motivation and Introduction

The emergence of Large Language Models (LLMs) brings broad concern about the malicious usage of machine-generated text (MGT). Effective MGT detectors are urgently needed.

### Defects on existing detectors:

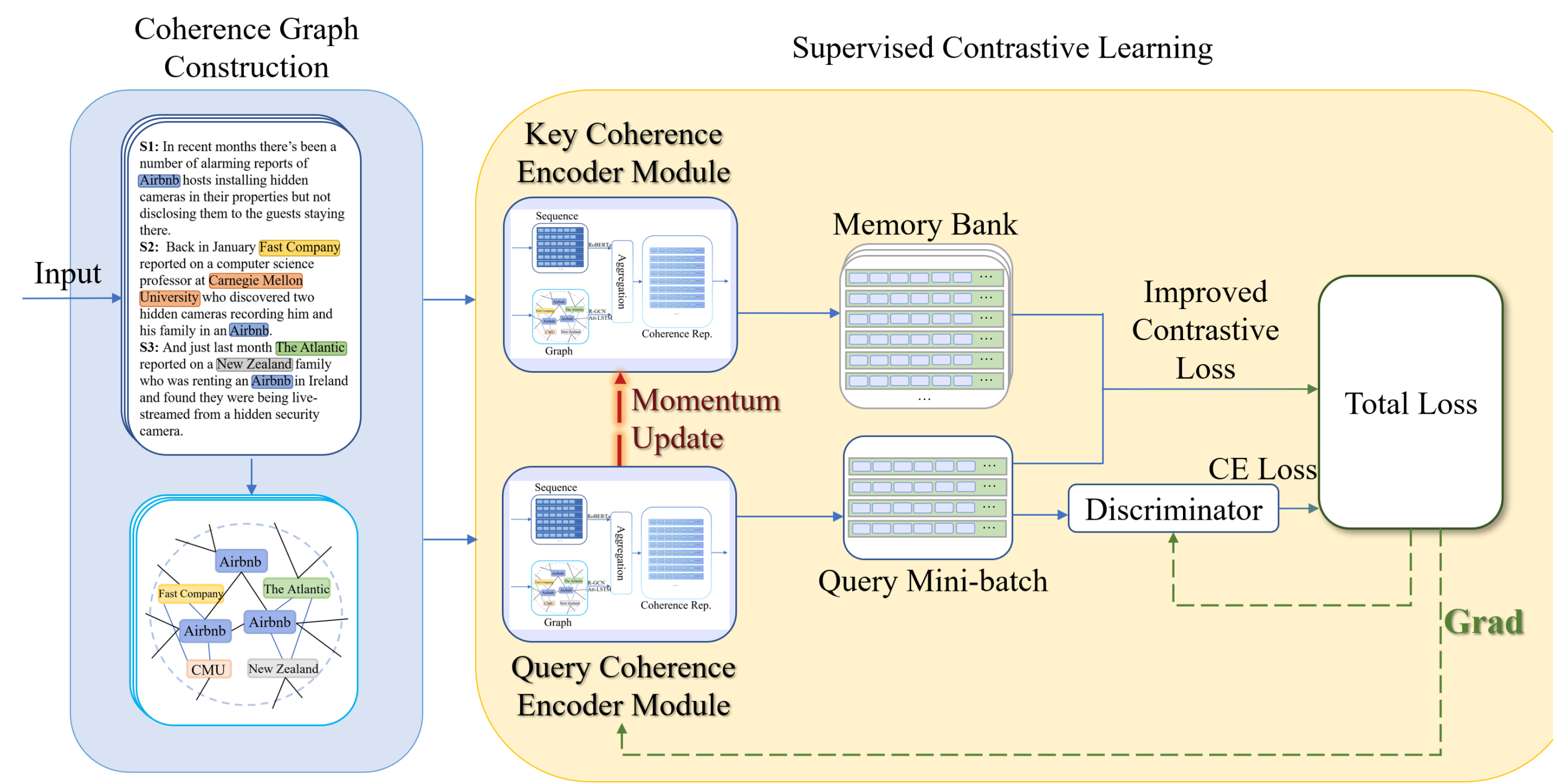
- Treat input documents as flat sequences of tokens while ignoring high-level linguistic representation of text structure
- Performance constrained by the amount of available annotated data

### Our contributions:

- We model the text coherence with entity consistency and sentence interaction while statistically proving its distinctiveness in MGT detection, and we further introduce the linguistic feature at the input stage
- We introduce contrastive learning and improved contrastive loss into the MGT detector to alleviate data dependence
- We surprisingly find that MGTs originated from up-to-date language models could be easier to detect than those from previous models in our experiments



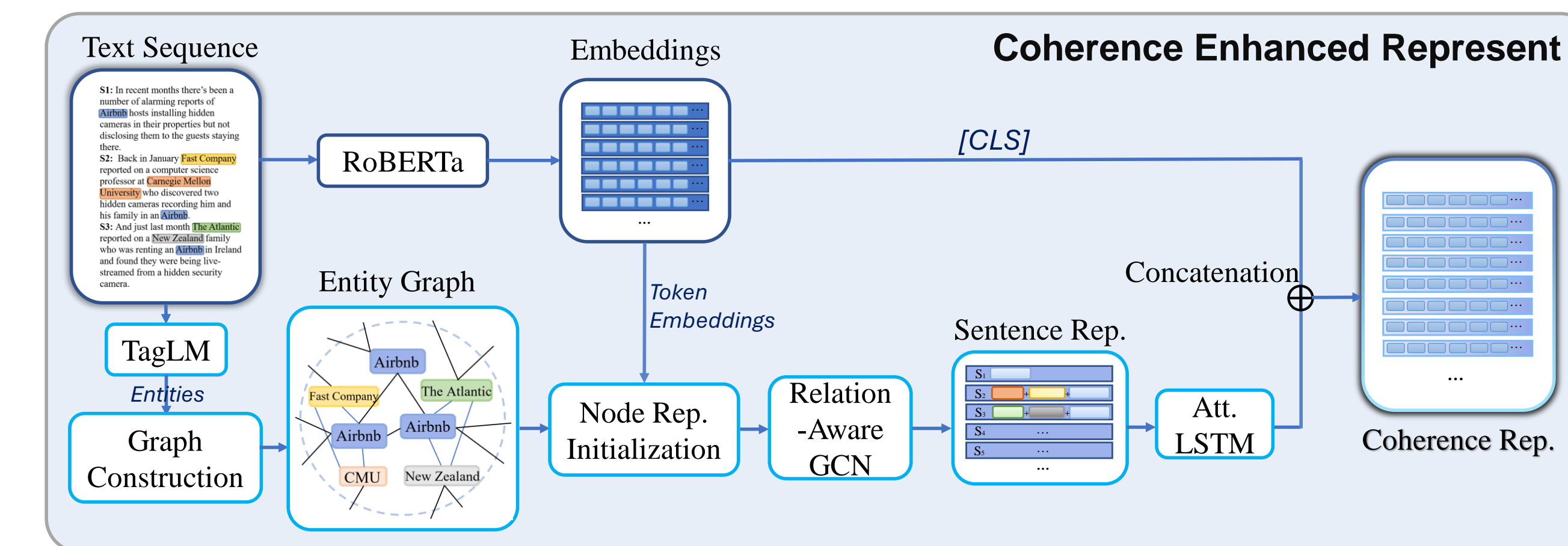
## CoCo Methodology



CoCo consists of two complements: Coherence Graph Construction and Contrastive Learning.

### Encoder Design: Coherence Graph Construction

We propose an innovative coherence encoder module (CEM), which is utilized to integrate coherence information into a semantic representation of text by propagating and aggregating information from different granularity via graph, to encode a coherence enhanced representation.



### Toward Low Resource Scenario: Contrastive Learning

To instance compactness and class separability in low-resource settings, we utilize MOCO<sup>[2]</sup> as backbone and come up with an improved contrastive loss (ICL) for dynamically adjusting the weight of negative pair similarity according to the hardness of negative samples.

$$\mathcal{L}_{ICL} = \sum_{j=1}^M \mathbf{1}_{y_i=y_j} \log \frac{S_{ij}}{\sum_{p \in \mathcal{P}(i)} S_{ip} + \sum_{n \in \mathcal{N}(i)} r f_{in} S_{in}}, \quad (1)$$

$$S_{ij} = \exp(D_q^i D_q^j / \tau), r f_{ij} = \beta \frac{D_q^i D_k^n}{\text{avg}(D_q^i D_k^1 : D_k^{|\mathcal{N}(i)|})},$$

where  $\mathcal{P}(i)$  is the positive set in which data has the same label with  $q_i$  and  $\mathcal{N}(i)$  is the negative set in which data has a different label from  $q_i$ .  $D_k$  is the key module representations and  $D_q$  is the query module representations.

## Experiment and Analysis on Comparison, Ablation and Robustness

We conduct our main experiments on two public datasets and two self-constructed GPT-3.5 datasets<sup>[3]</sup>, against seven baselines and SOTA. Also, an ablation study and robustness test are implemented. More additional experiments are in the paper. Here are some **key findings**:

- CoCo surpasses the state-of-the-art methods in MGT detection in both settings
- GROVER Dataset is the hardest to detect while GPT-3.5 datasets are surprisingly easy
- Coherence graph and contrastive learning module both contribute orthogonally
- CoCo shows comparable robustness to perturbations to some extent

Dataset Size	GROVER				GPT-2			
	Limited Dataset (500 examples)		Full Dataset		Limited Dataset (500 examples)		Full Dataset	
Metric	ACC	F1	ACC	F1	ACC	F1	ACC	F1
GPT2	0.5747 ± 0.0217	0.4394 ± 0.0346	0.8274 ± 0.0091	0.8003 ± 0.0141	0.5380 ± 0.0067	0.4734 ± 0.0182	0.8913 ± 0.0066	0.8839 ± 0.0078
XLNet	0.5660 ± 0.0265	0.4707 ± 0.0402	0.8156 ± 0.0079	0.7493 ± 0.0073	0.6551 ± 0.0083	0.5715 ± 0.0095	0.9091 ± 0.0091	0.9027 ± 0.0111
RoBERTa	0.6621 ± 0.0133	0.5895 ± 0.0231	0.8772 ± 0.0029	0.8171 ± 0.0048	0.8223 ± 0.0088	0.7978 ± 0.0085	0.9402 ± 0.0039	0.9384 ± 0.0044
DualCL	0.5835 ± 0.0857	0.4628 ± 0.1076	0.7574 ± 0.0855	0.6388 ± 0.1300	0.6039 ± 0.1367	0.5435 ± 0.0903	0.8023 ± 0.1120	0.8046 ± 0.1530
CE+SCL	0.6870 ± 0.0142	0.5961 ± 0.0197	0.8782 ± 0.0044	0.8202 ± 0.0057	0.8355 ± 0.0046	0.8127 ± 0.0067	0.9408 ± 0.0006	0.9390 ± 0.0009
GLTR	0.3370	0.4935	0.6040	0.5182	0.7755	0.7639	0.7784	0.7691
DetectGPT	0.5910	0.4258	0.6142	0.5018	0.7941	0.6982	0.7939	0.7002
CoCo	<b>0.6993 ± 0.0119</b>	<b>0.6125 ± 0.0159</b>	<b>0.8826 ± 0.0018</b>	<b>0.8265 ± 0.0036</b>	<b>0.8530 ± 0.0019</b>	<b>0.8410 ± 0.0018</b>	<b>0.9457 ± 0.0004</b>	<b>0.9452 ± 0.0004</b>

Dataset Size	GPT-3.5 Unmixed		GPT-3.5 Mixed	
	Limited Dataset (500 examples)	Full Dataset	Limited Dataset (500 examples)	Full Dataset
Metric	ACC	F1	ACC	F1
GPT2	0.9023 ± 0.0095	0.8920 ± 0.0073	0.9917 ± 0.0056	0.9905 ± 0.0042
XLNet	0.9107 ± 0.0068	0.9037 ± 0.0064	0.9620 ± 0.0043	0.9634 ± 0.0068
RoBERTa	0.9670 ± 0.0084	0.9681 ± 0.0077	0.9928 ± 0.0035	0.9913 ± 0.0040
CE+SCL	0.9823 ± 0.0053	0.9703 ± 0.0070	0.9944 ± 0.0023	0.9943 ± 0.0031
GLTR	0.9255	0.9287	0.9350	0.9358
DetectGPT	0.9220	0.8744	0.9245	0.8991
CoCo	<b>0.9889 ± 0.0044</b>	<b>0.9791 ± 0.0062</b>	<b>0.9972 ± 0.0015</b>	<b>0.9957 ± 0.0020</b>

Model	ACC		F1	
	RoBERTa	CoCo	RoBERTa	CoCo
CoCo (Plain)	0.6635	0.5901	0.6993	0.6125
CoCo (Sentence Nodes)	0.6635	0.5901	0.6993	0.6125
CoCo (Coherence)	0.6635	0.5901	0.6993	0.6125
CoCo (Coherence+LSTM)	0.6635	0.5901	0.6993	0.6125
CoCo (Coherence+LSTM+SCL)	0.6635	0.5901	0.6993	0.6125
CoCo	0.7843	0.6684	0.7843	0.6684

### Preliminary Explore on the Detectable Feature in GPT-3.5 Dataset

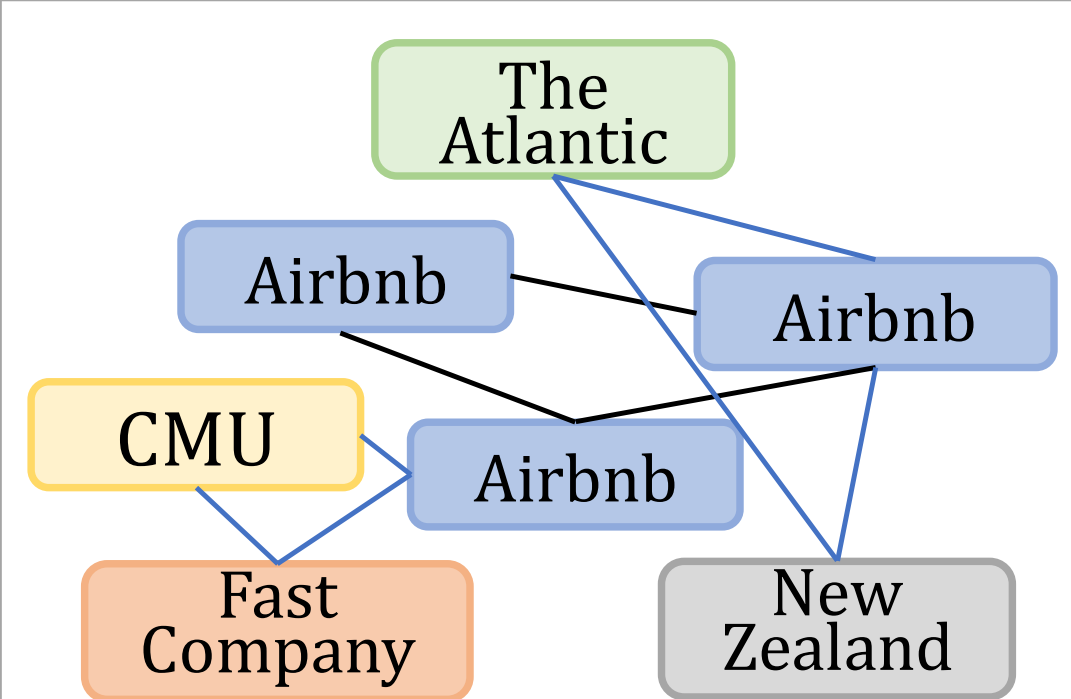
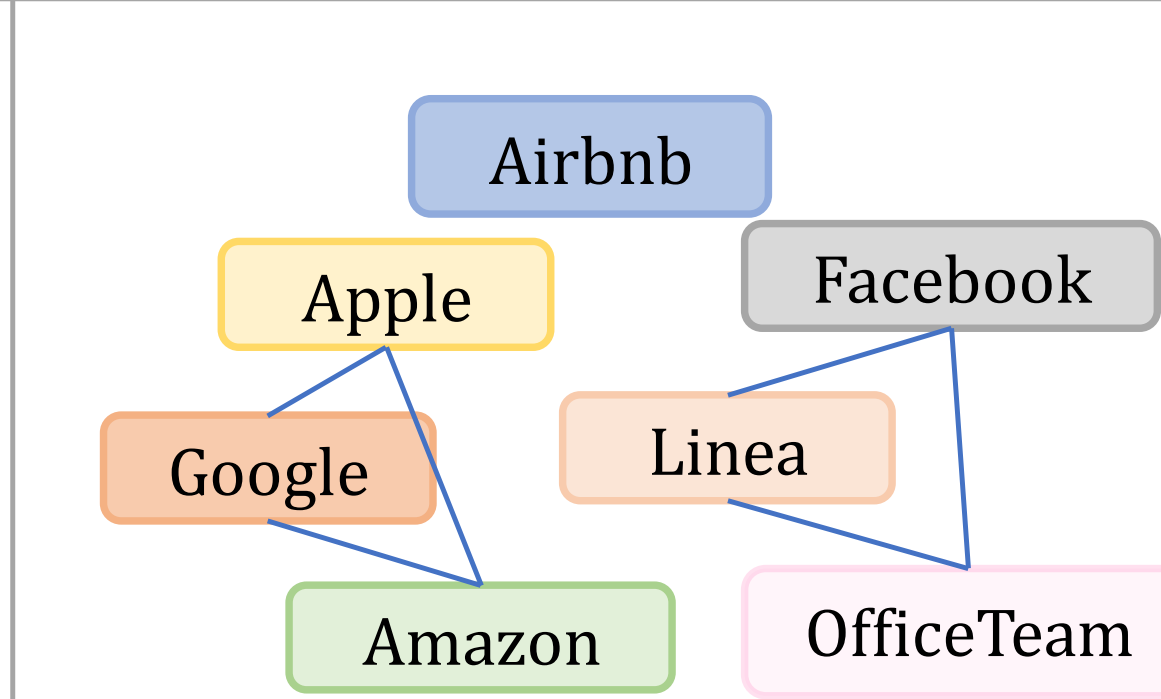
We probe the statistical interpretation behind the GPT-3.5 dataset and try to answer the question: *Why the MGTs by GPT-3.5 are relatively easy to detect?* We count the N-gram coverage of the supporters in Transformers-Interpret and the token coverage from the Statistic Cue.

N-gram Coverage	MGT	HWT	Token	Productivity	Coverage
γ <sub>1</sub>	0.6659	0.6377	according	0.6923	0.3126
γ <sub>2</sub>	0.4250	0.3630	where	0.6842	0.1998
γ <sub>3</sub>	0.2883	0.2076	they	0.6316	0.3837
γ <sub>4</sub>	0.2019	0.1372			
γ <sub>5</sub>	0.1425	0.0935			

- More consecutive spans of tokens act as an indicator for MGT than HWT
- No existing vulnerability in our dataset since trade-off between productivity and coverage
- The Easy-to-detect nature of GPT-3.5 texts might originate from language patterns

### References

- [1] Grosz B J, Sidner C L. Attention, intentions, and the structure of discourse[J]. Computational linguistics, 1986.
- [2] He K, Fan H, Wu Y, et al. Momentum contrast for unsupervised visual representation learning[C]//Proceedings of the IEEE/CVF conference on computer vision and pattern recognition. 2020: 9729-9738.
- [3] CoCo GPT-3.5 Machine-Generated Text Datasets, [https://huggingface.co/datasets/ZachW/MGTDetect\\_CoCo](https://huggingface.co/datasets/ZachW/MGTDetect_CoCo)

	Human-Written Text (HWT)	Machine-Generated Text (MGT)
Sequence	<p><b>S1:</b> In recent months, there's been a number of alarming reports of Airbnb hosts installing hidden cameras in their properties but not disclosing them to the guests staying there.</p> <p><b>S2:</b> Back in January, Fast Company reported on a computer science professor at Carnegie Mellon University who discovered two hidden cameras recording him and his family in an Airbnb.</p> <p><b>S3:</b> And just last month, The Atlantic reported on a New Zealand family who was renting an Airbnb in Ireland and found they were being live-streamed from a hidden security camera.</p>	<p><b>S1:</b> Anyone who finds a video of someone on Airbnb will probably fall under the new category of hidden cameras, which can be found only in a large part of every Airbnb listing, and you're never alone.</p> <p><b>S2:</b> Apple, Google, and Amazon combined to find the most hidden camera listings in December 2018.</p> <p><b>S3:</b> The electronics giant's Facebook, the mapping app, and the mobile messaging company Linea formed an OfficeTeam unit that can find the video even if someone's not using them, and can track real-time activity.</p>
Graph		

### Coherence Modeling based on Centering Theory<sup>[1]</sup>

"Coherence of texts could be modeled by sentence interaction around center entities." We build a coherence graph, treat entities as nodes and co-occurrence relationship of entities as edges.