

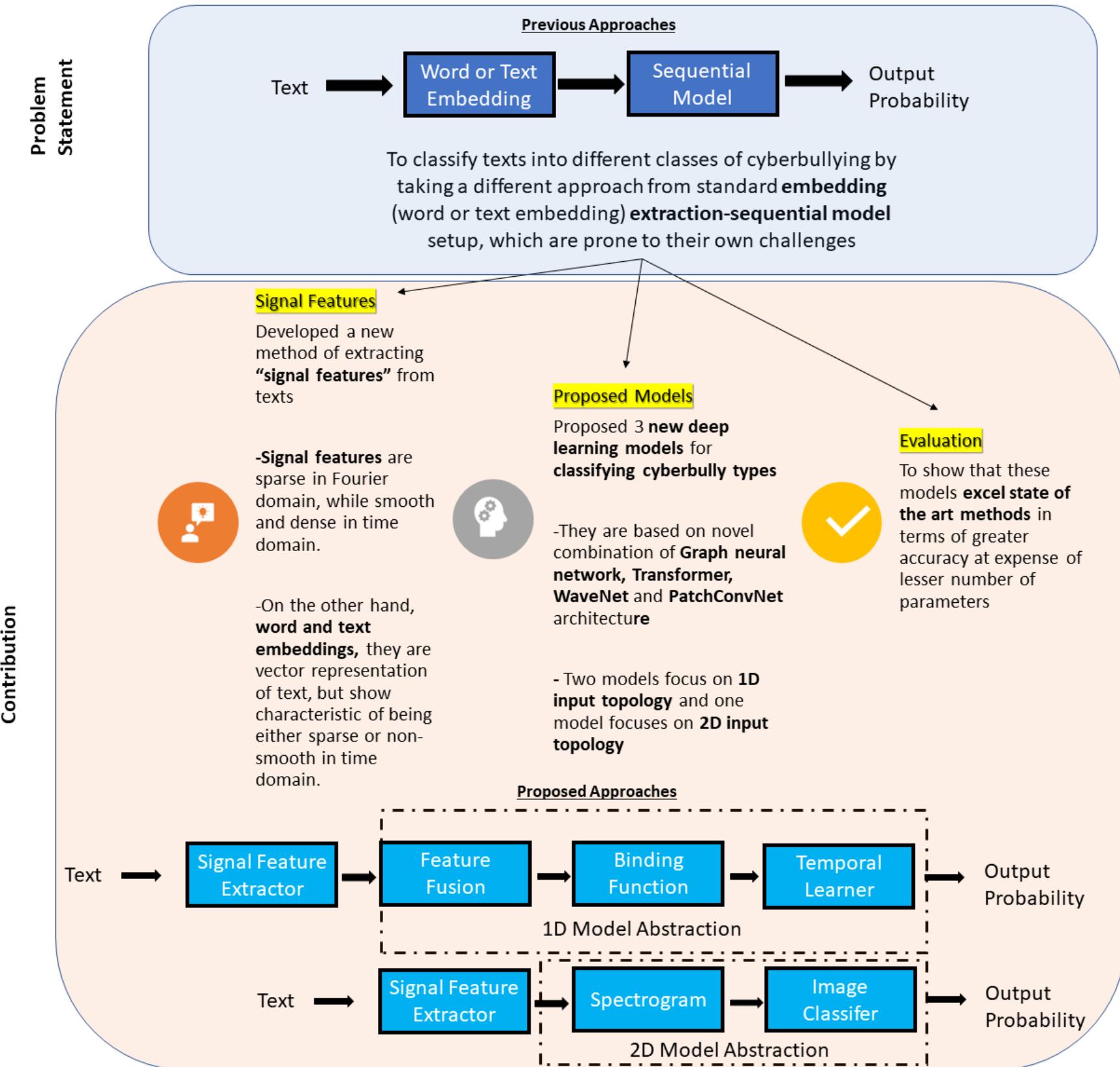
Supervised signal-features learning from cyber-text and novel deep learning techniques for fine-grained cyberbully classification

Defence Presentation

Junaid Iqbal Khan

**Department of Electronics & Information Engineering
Korea Aerospace University, South Korea**

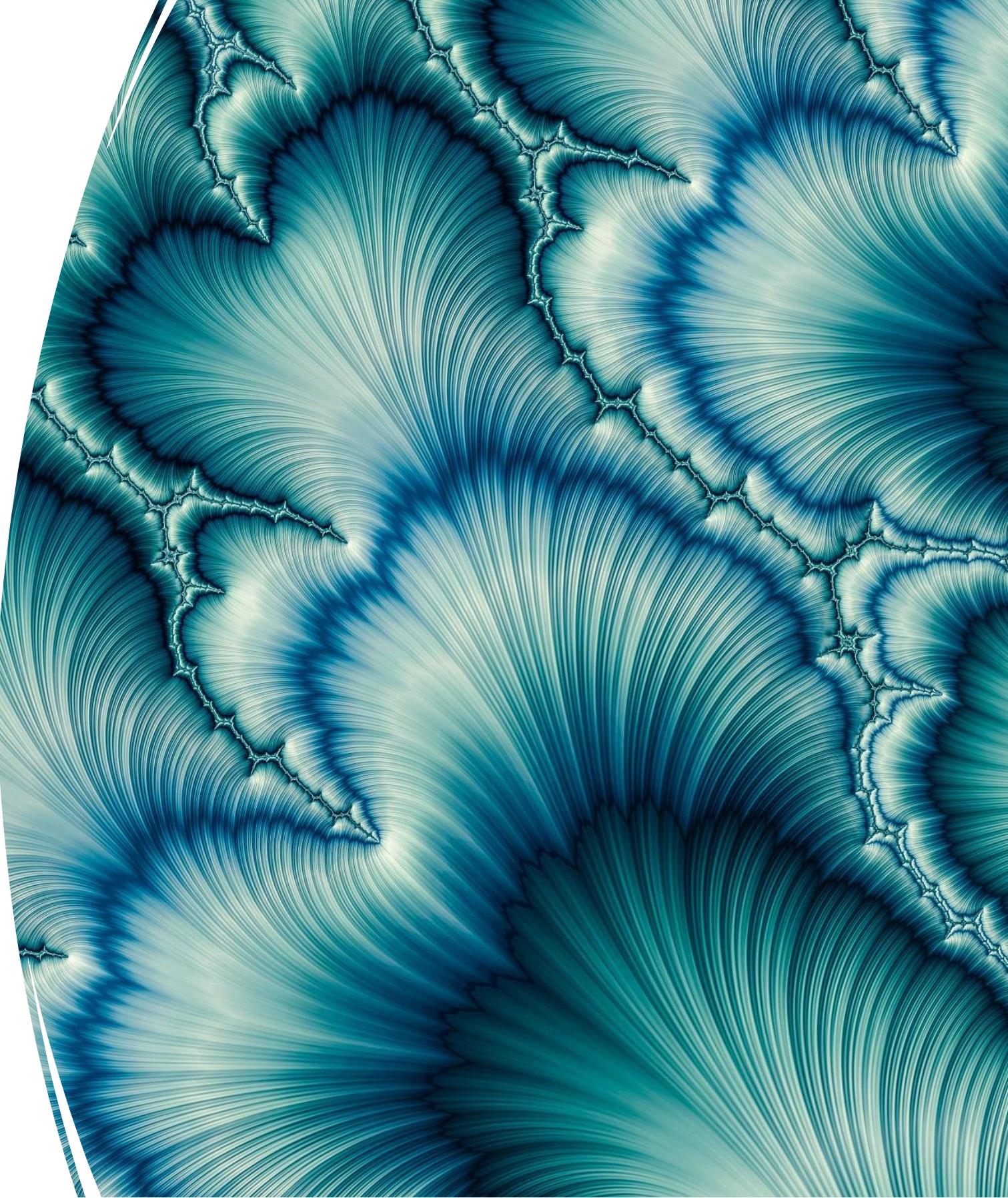
Key Original Concept at a Glance

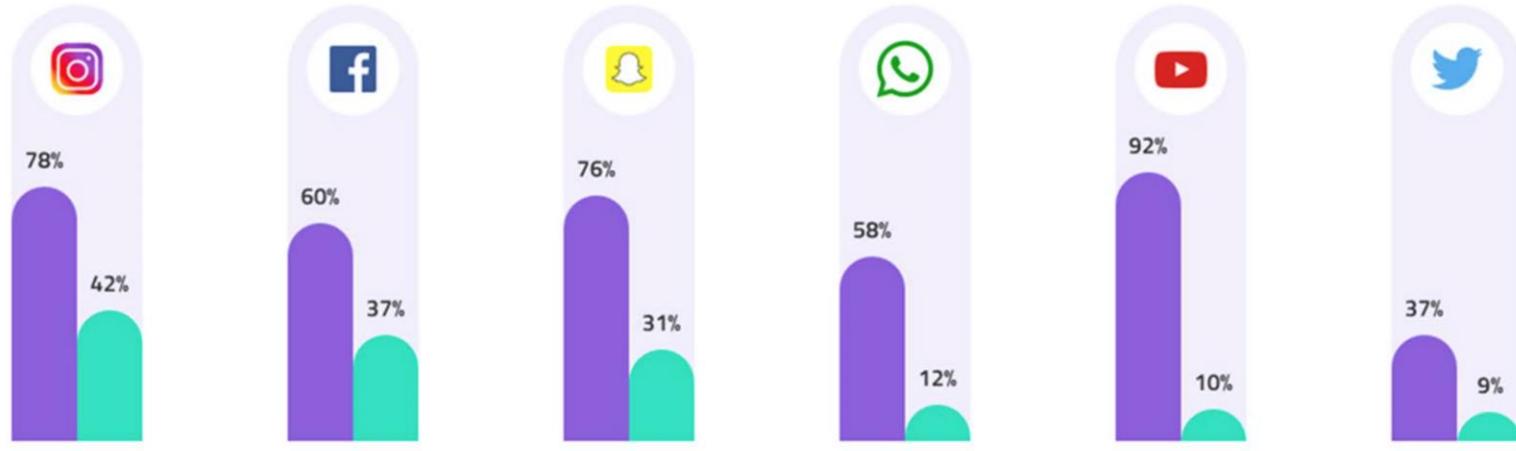


Contents



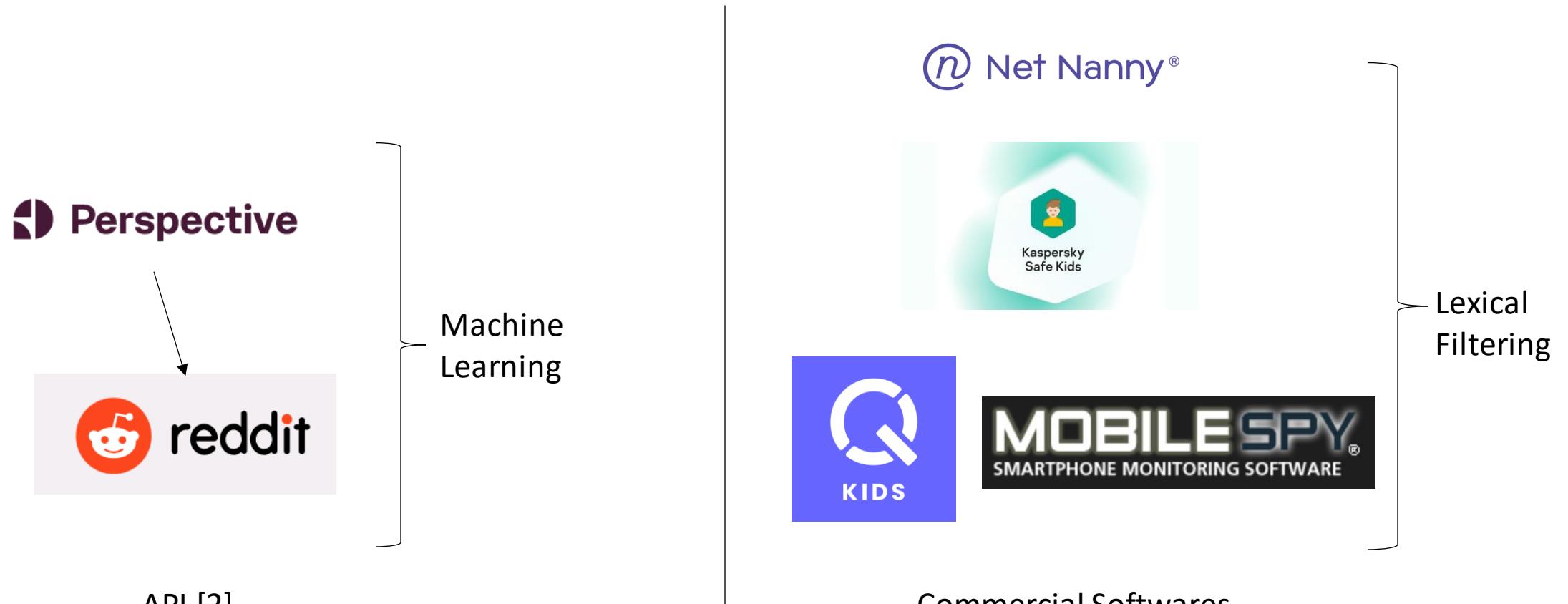
- **Introduction**





■ Percentage of all young people who use the platform ■ Percentage of young people who have experienced cyberbullying on the platform

Statistics on social media platforms where cyberbullying occurs [1]





Microsoft's Twitter AI Chatbot "Tay"

CNN

This AI chatbot is dominating social media with its frighteningly good essays

Possible issues. While ChatGPT successfully fielded a variety of questions submitted by CNN, some responses were noticeably off. In fact, Stack...



OpenAI

New and Improved Content Moderation Tooling

ChatGPT



Examples

"Explain quantum computing in simple terms" →

"Got any creative ideas for a 10 year old's birthday?" →

"How do I make an HTTP request in Javascript?" →



Capabilities

Remembers what user said earlier in the conversation

Allows user to provide follow-up corrections

Trained to decline inappropriate requests



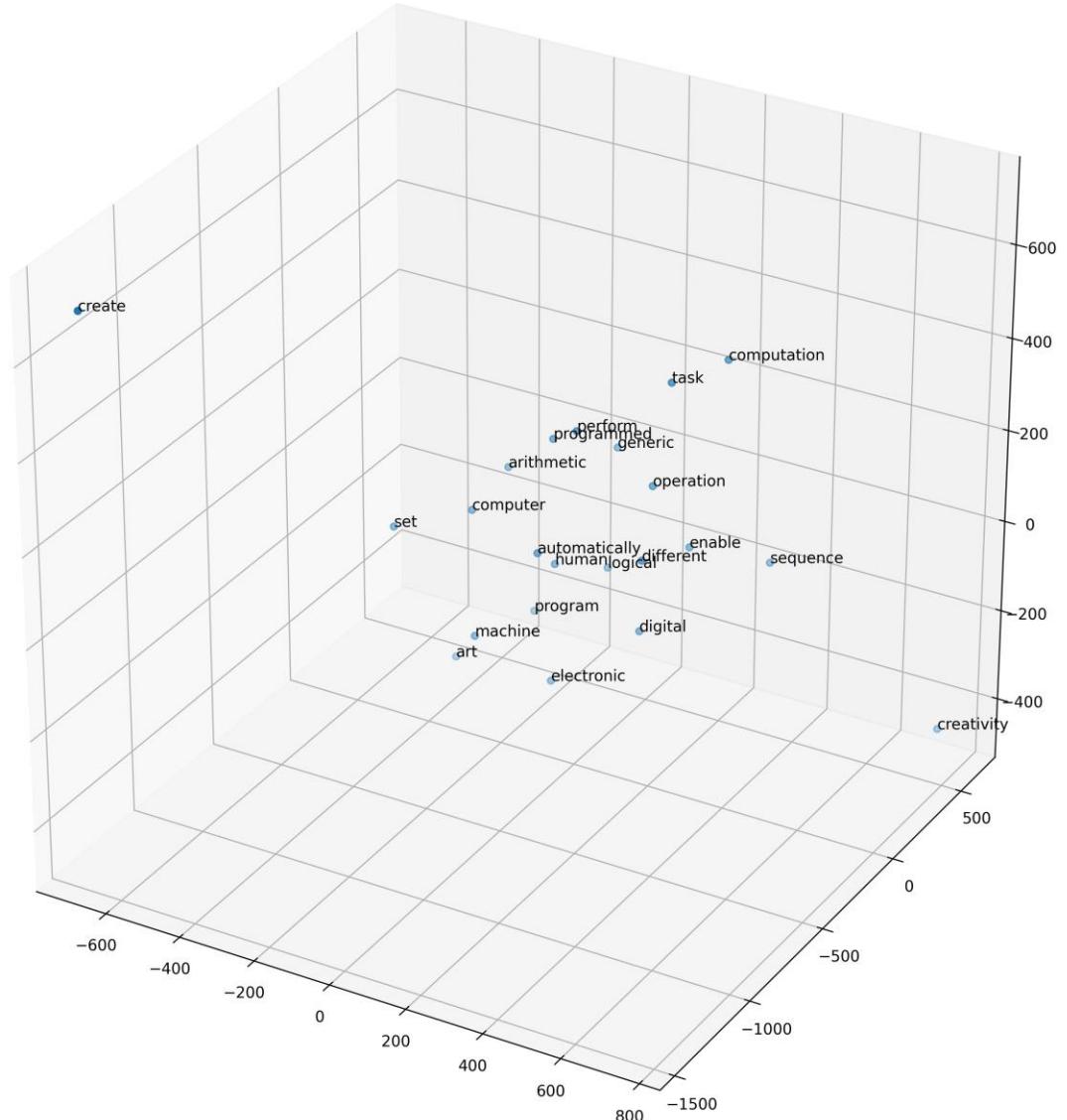
Limitations

May occasionally generate incorrect information

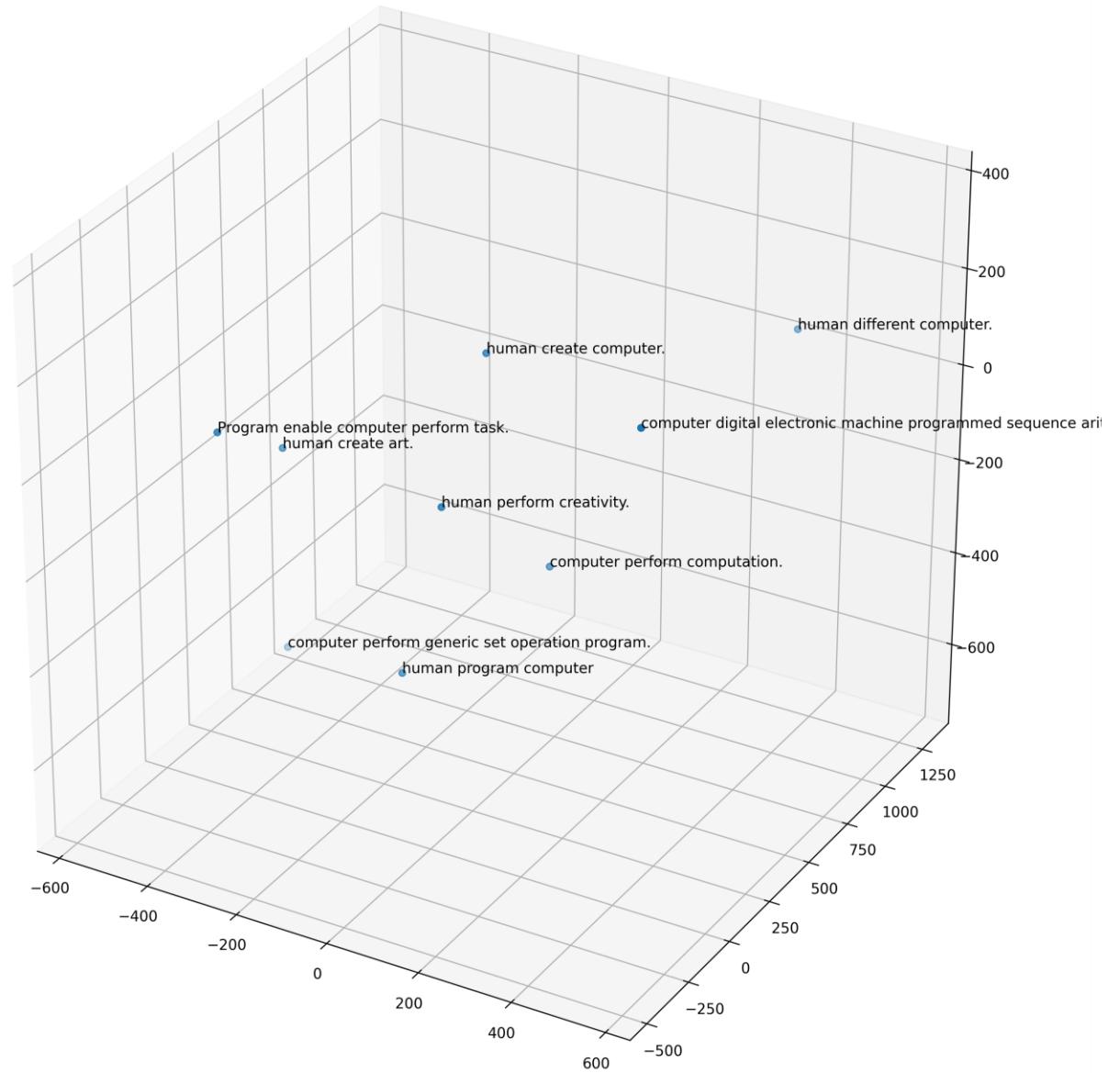
May occasionally produce harmful instructions or biased content

Limited knowledge of world and events after 2021

OpenAI's ChatGPT

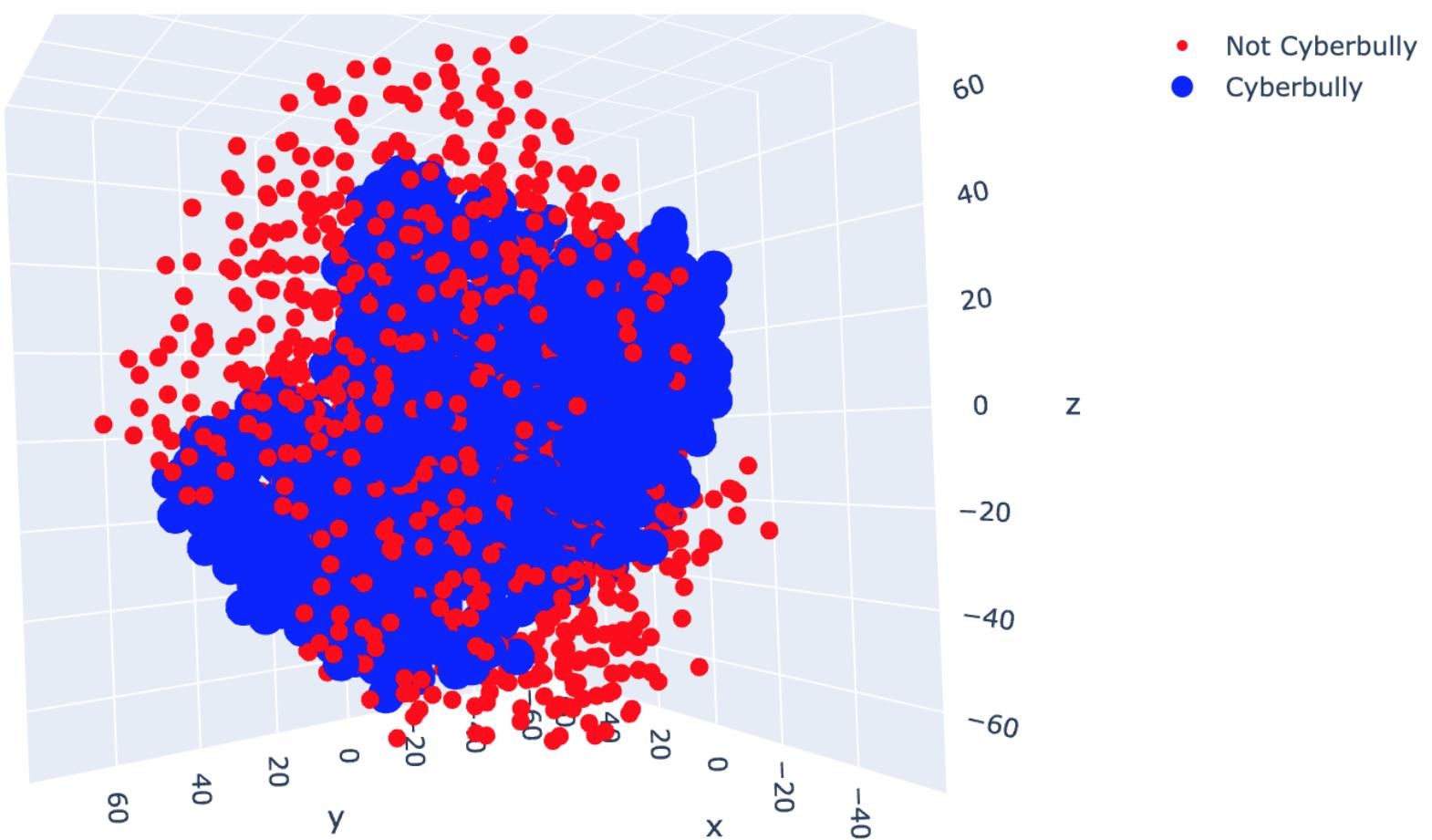


Word Embeddings

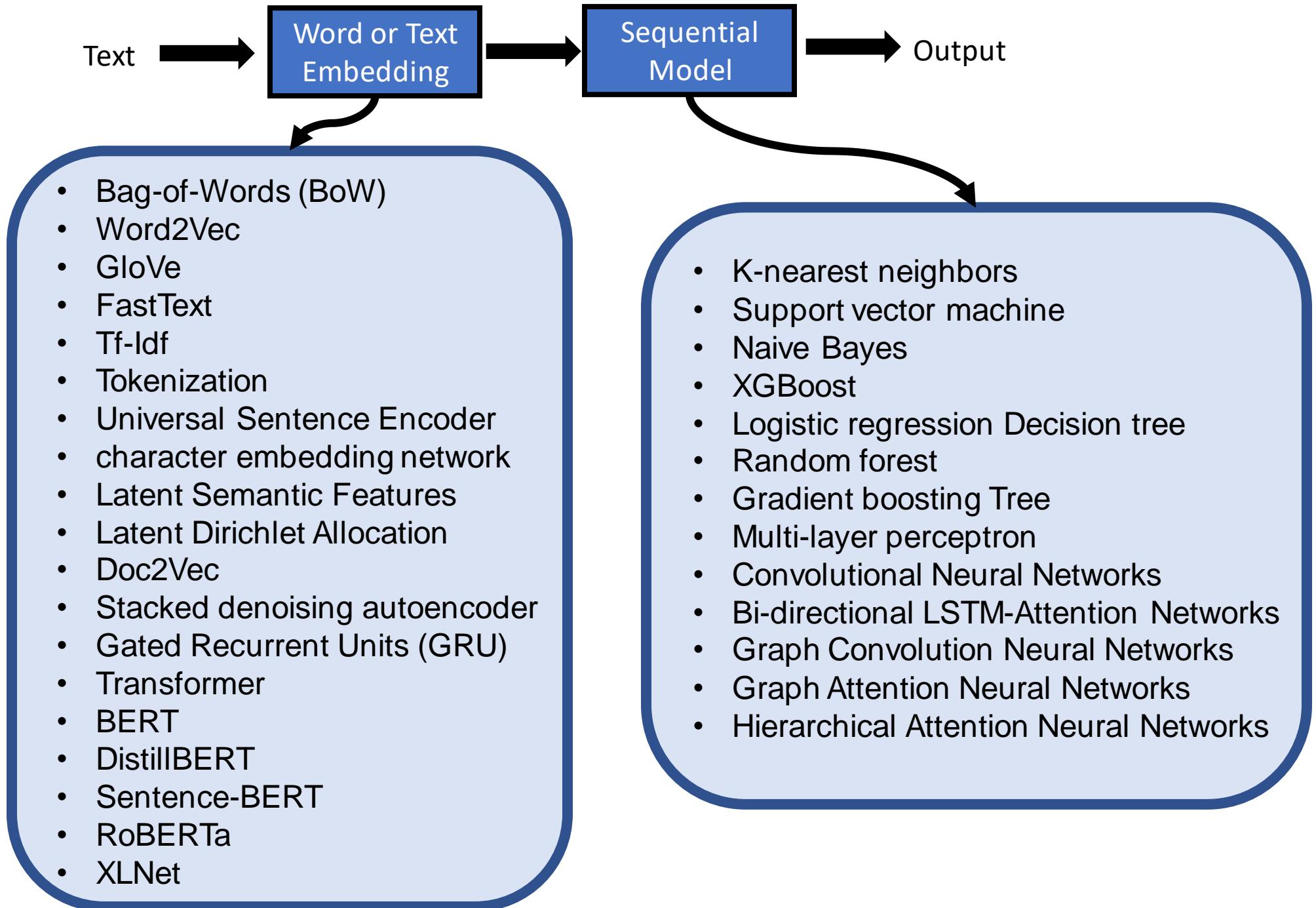


Text Embeddings

Comparison of Word and Textual Embeddings over same data

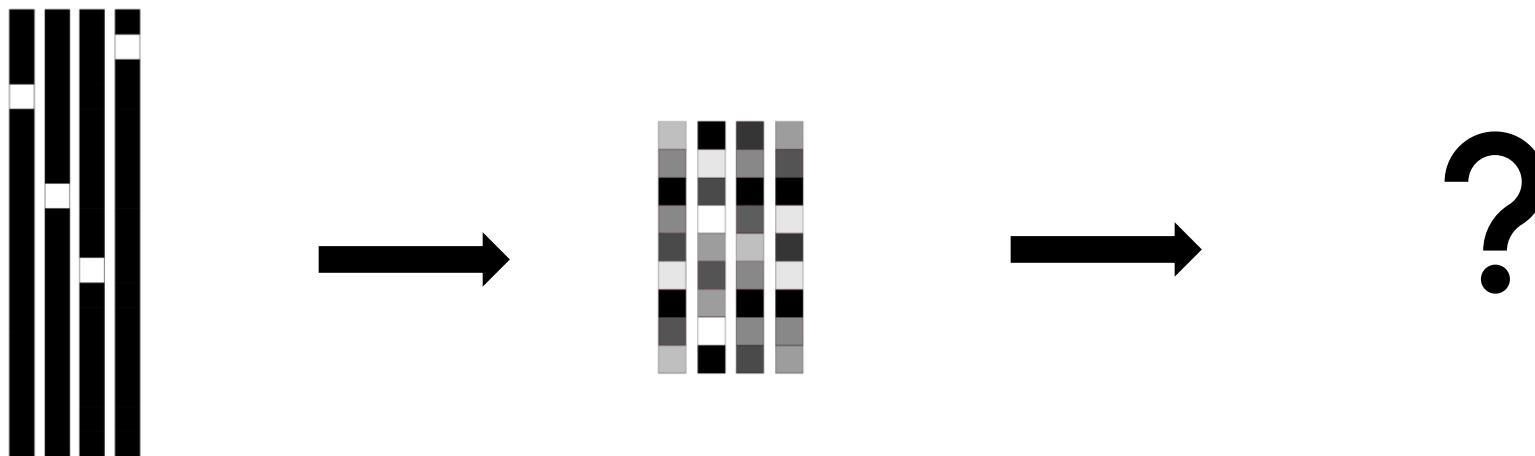


Comparison of Textual Embeddings for Cyberbully and Non-cyberbully related texts



Raw Pipeline of Current Data driven NLP (including Cyberbully Detection)

Research Motivation



One hot coded

- high dimensional
- sparse
- hard coded

Word2Vec

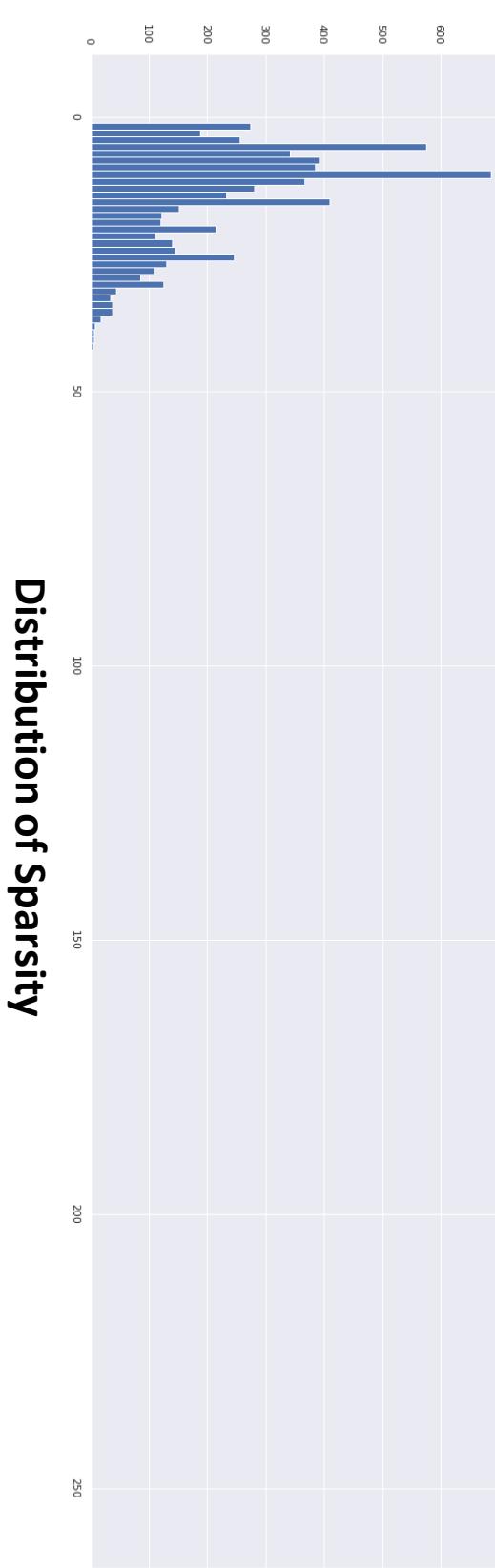
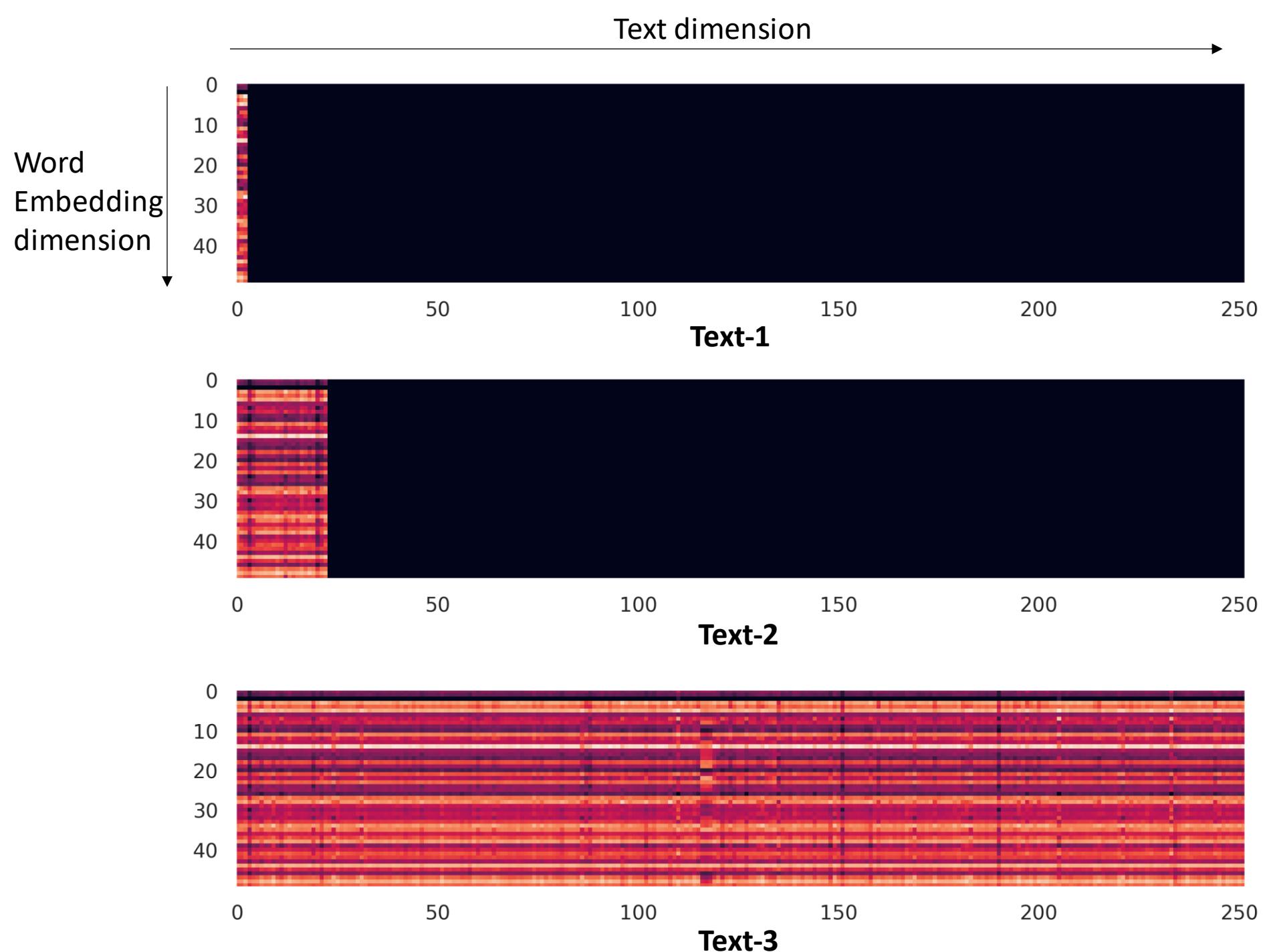
- interpretable in Euclidean space
- low dimensional
- dense (not dense when considered for textual embedding)
- learned from data

Mythical Embedding

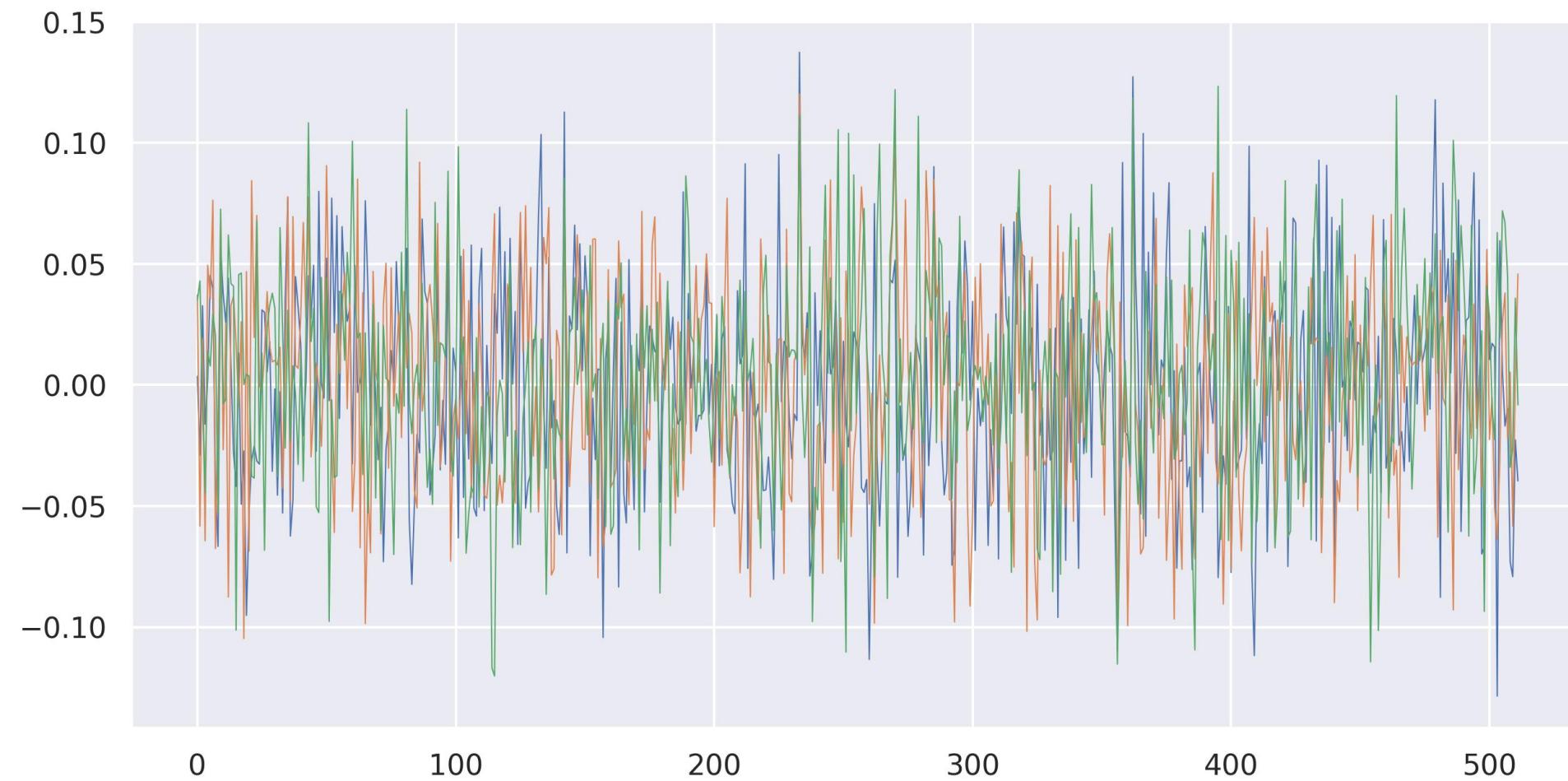
- low dimensional
- interpretable over real domain
- dense in real domain with respect to text
- smooth in real domain with respect to text
- sparse and definite in frequency domain
- learned from data
- invariant to size of text

- Limitations of State of Art Models



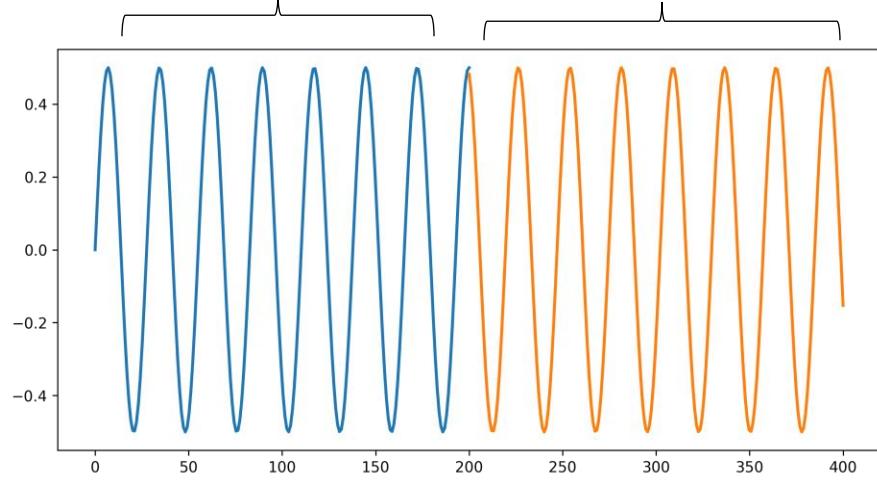


Fallacy of Word Embeddings: Sparsity in sequence format

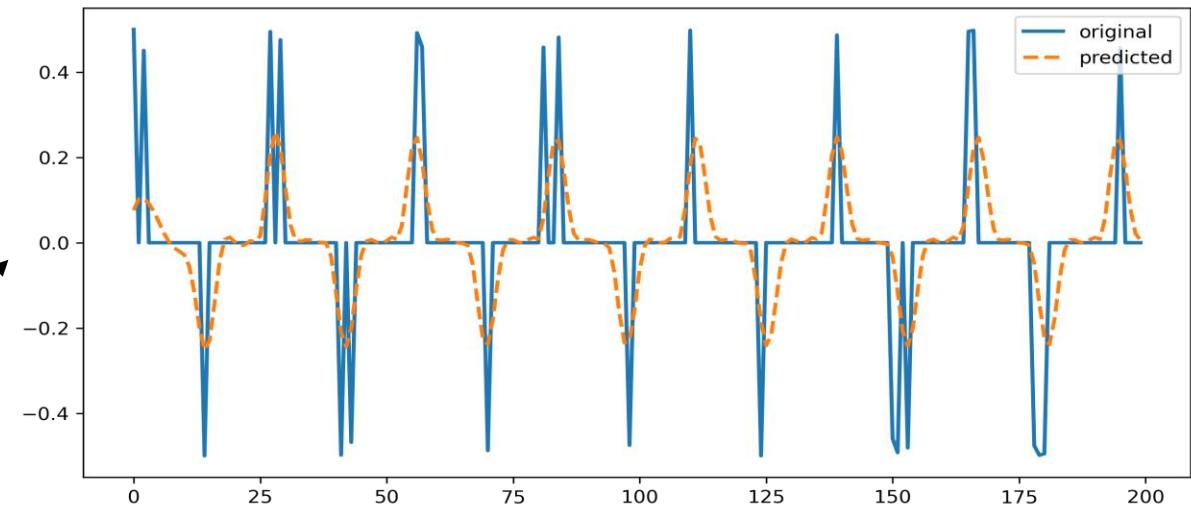
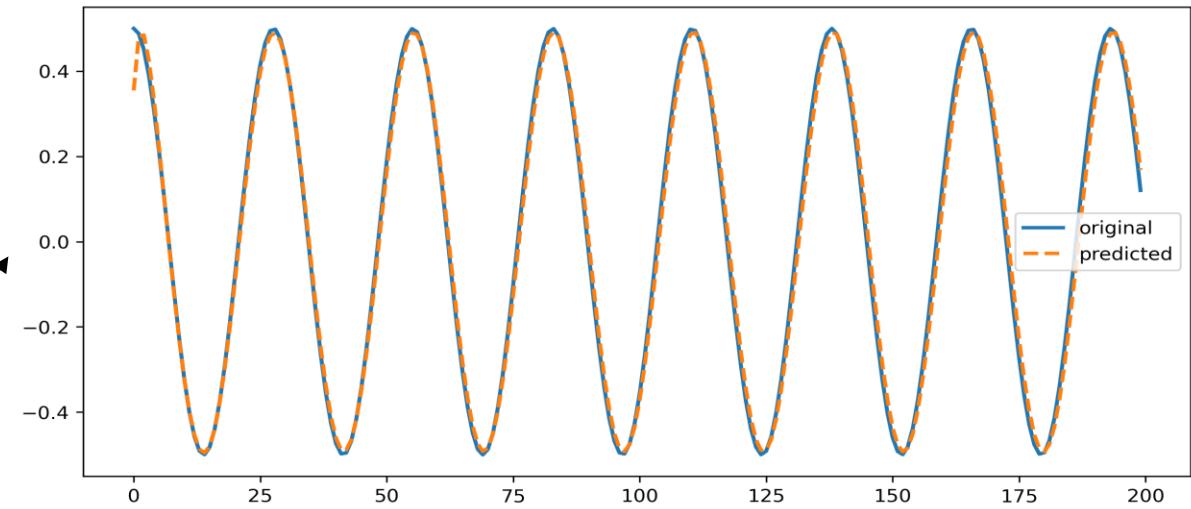
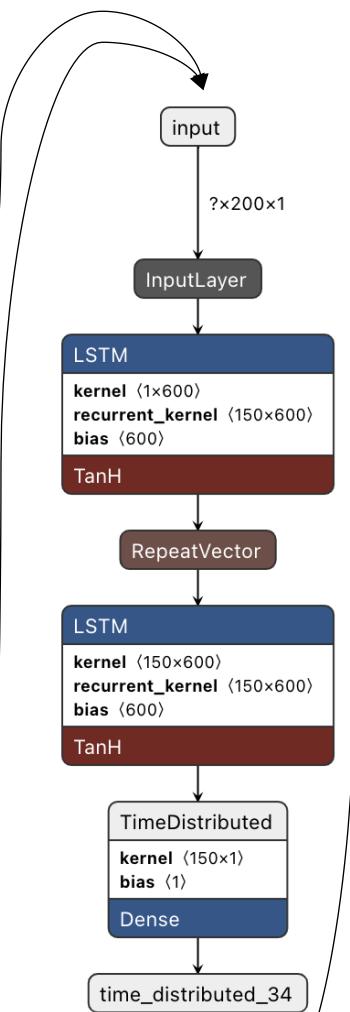
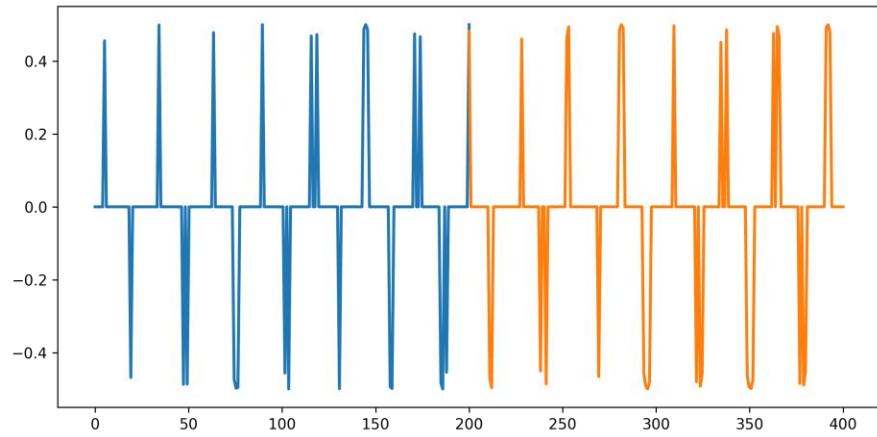


Fallacy of Text Embeddings: Non-smoothness in sequence format

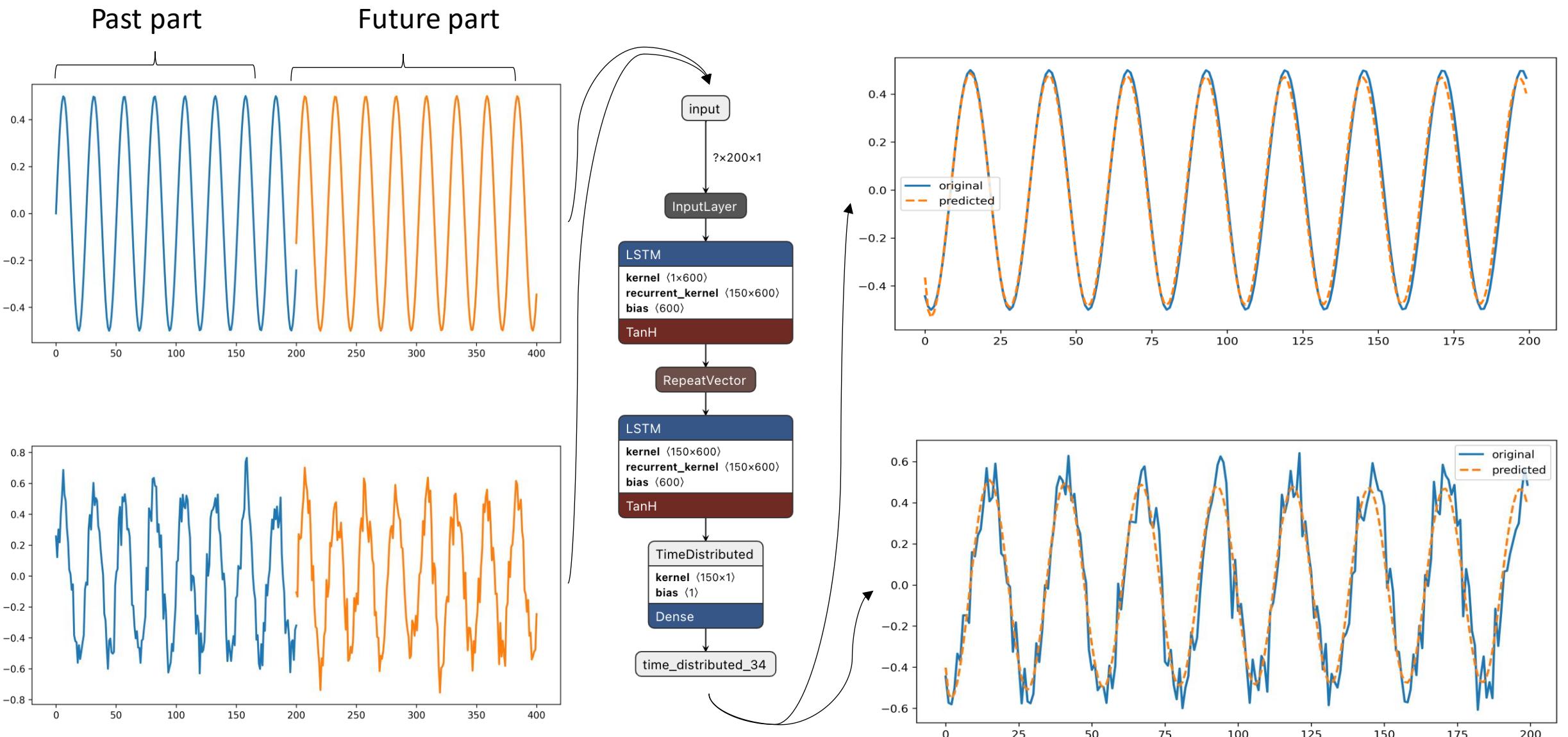
Past part



Future part

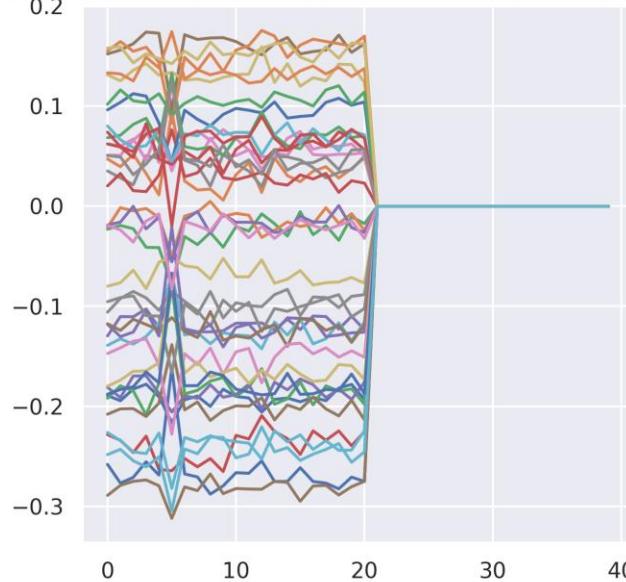


Proposal: Sequential models are not resilient at approximating sparse sequences

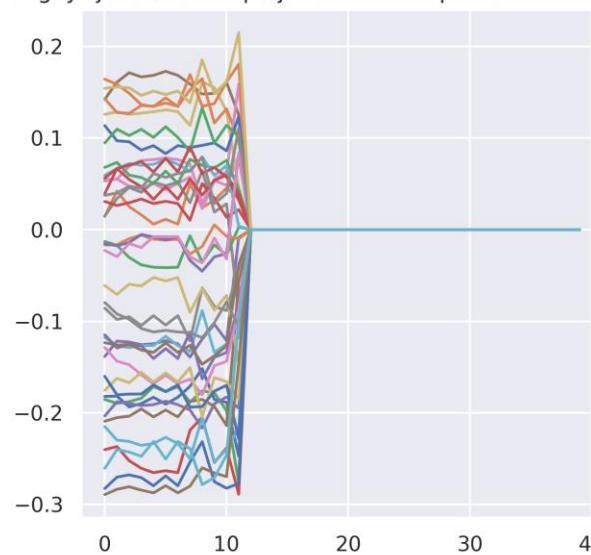


Proposal: Sequential Models are not resilient at approximating non-smooth sequences

Is class shaming justified while condemning rape threats?? Make a gay, trans joke, and these people get triggered. Majority is not expected tolerate the type of shaming that's insulting to minorities just because it's majority.

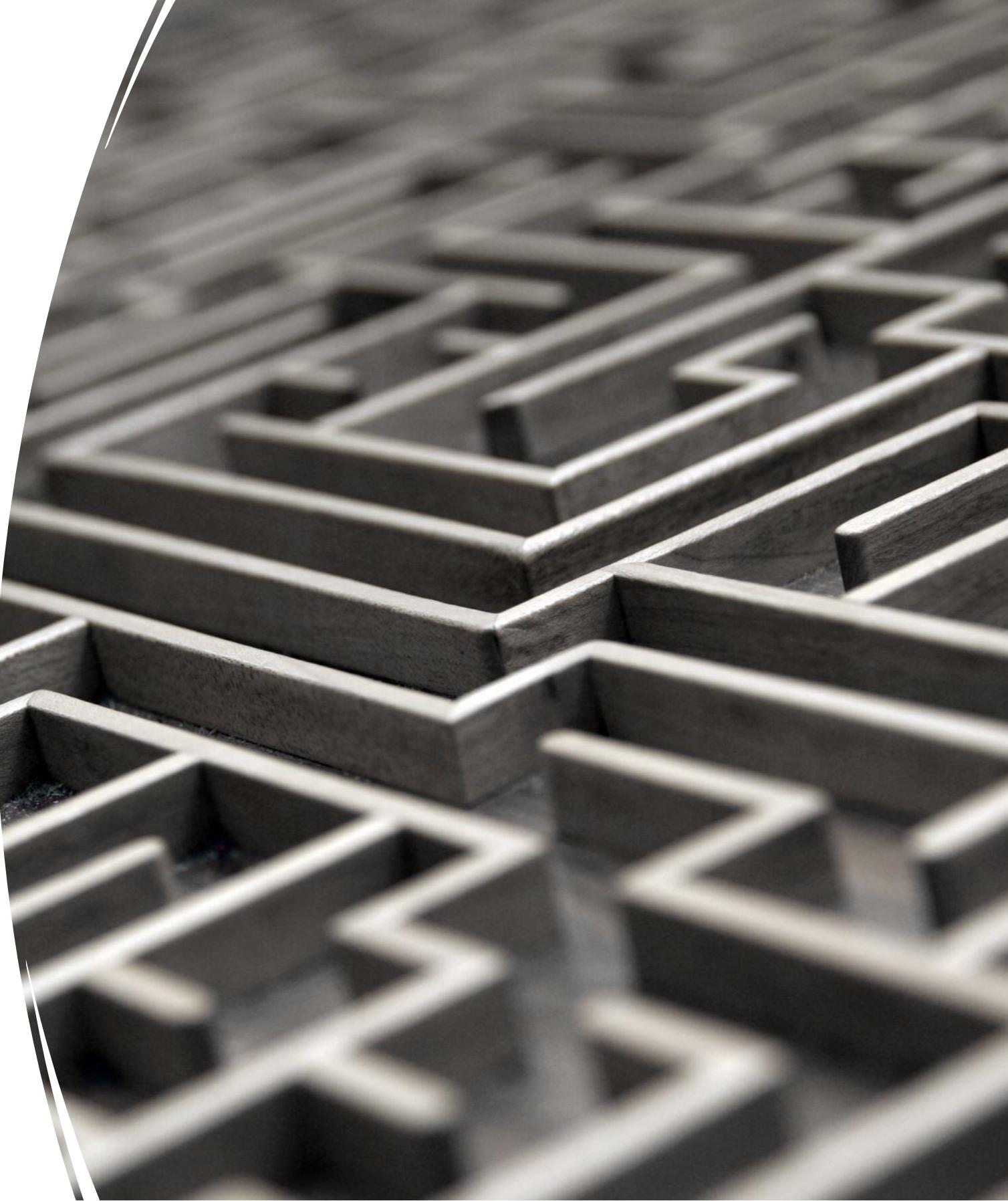


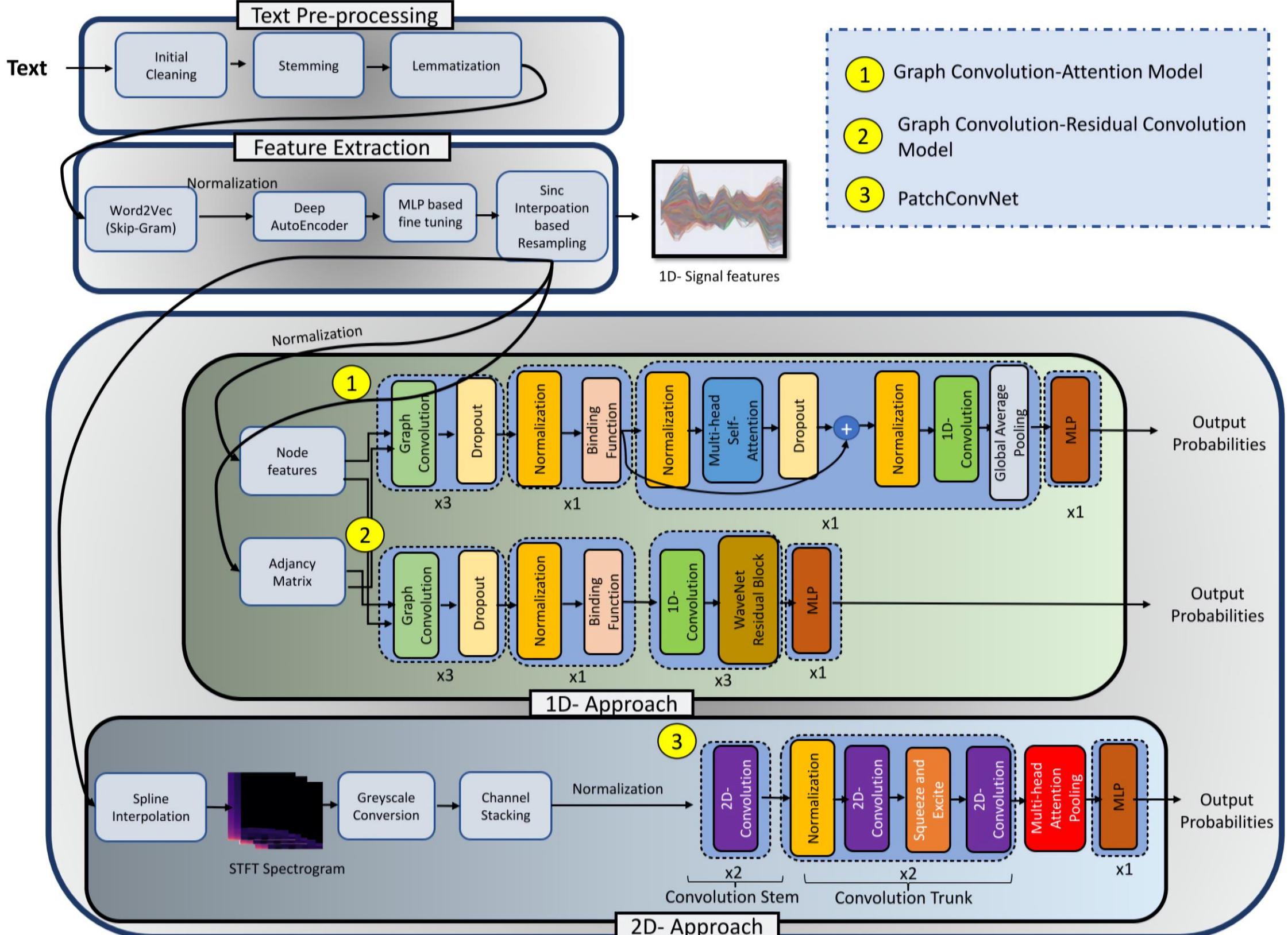
Just heard someone make a "gay" joke and a rape joke within a span of 2 minutes. Very cool... #overheardatUST



Proposal: Lack of uniformness in “window of information” for same context but different length sentences

- Proposed Approach





Functional Pipeline

- Evaluation



No.	Name	Nature	Classes
1	Fine Grained Cyberbullying Dataset	Single-label, Multi-class	6
2	Toxic Comment Classification Dataset	Multi-label, Multi-class	7

Datasets Considered

Snapshot of Dataset-1

	tweet_text	cyberbullying_type
1	In other words #katandandre, your food was crapilicious! #mkr	not_cyberbullying
2	Why is #aussietv so white? #MKR #theblock #ImACelebrityAU #today #sunrise #studio10 #Neighbours #WonderlandTen #etc	not_cyberbullying
3	@XochitlSuckks a classy whore? Or more red velvet cupcakes?	not_cyberbullying
4	@Jason_Gio meh. :P thanks for the heads up, but not too concerned about another angry dude on twitter.	not_cyberbullying
5	@RudhoeEnglish This is an ISIS account pretending to be a Kurdish account. Like Islam, it is all lies.	not_cyberbullying
6	@Raja5aab @Quickileaks Yes, the test of god is that good or bad or indifferent or weird or whatever, it all proves gods existence.	not_cyberbullying
7	Itu sekolah ya bukan tempat bully! Ga jauh kaya neraka	not_cyberbullying
8	Karma. I hope it bites Kat on the butt. She is just nasty. #mkr	not_cyberbullying
9	@stockputout everything but mostly my priest	not_cyberbullying
10		

Snapshot of Dataset-2

comment_text	toxic	severe_toxic	obscene	threat	insult	identity_hate
I Explanation	0	0	0	0	0	0
D'aww! He matches this background colour I'm seemingly stuck with. Thanks. (talk) 21:51, Januar 0	0	0	0	0	0	0
Hey man, I'm really not trying to edit war. It's just that this guy is constantly removing relevant inf 0	0	0	0	0	0	0
"	0	0	0	0	0	0
You, sir, are my hero. Any chance you remember what page that's on?	0	0	0	0	0	0
"	0	0	0	0	0	0
COCKSUCKER BEFORE YOU PISS AROUND ON MY WORK	1	1	1	0	1	0
Your vandalism to the Matt Shirvington article has been reverted. Please don't do it again, or you 0	0	0	0	0	0	0
Sorry if the word 'nonsense' was offensive to you. Anyway, I'm not intending to write anything in 0	0	0	0	0	0	0
I alignment on this subject and which are contrary to those of DuLithgow	0	0	0	0	0	0
"	0	0	0	0	0	0

Single label

Multi label

Precision: measure of the accuracy of the classifier when it predicts the positive class. High precision means that the model is good at identifying positive examples and minimizing false positives

Recall: measure of the classifier's ability to correctly identify the positive class. High recall means that the model is good at identifying all of the positive examples in the data and minimizing false negatives

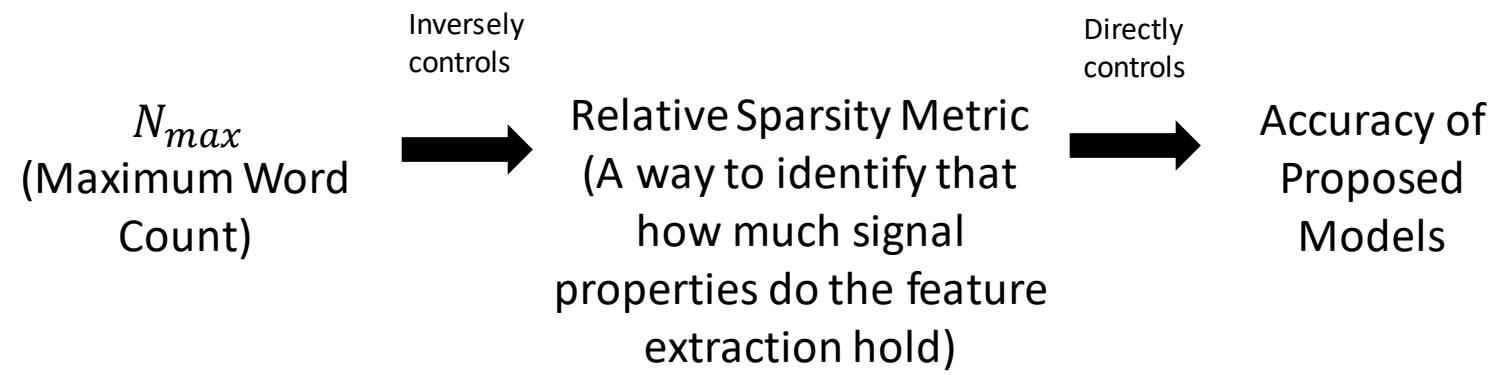
F_1 score: harmonic mean of precision and recall, and is a good metric to use when you want to balance precision and recall

Accuracy: ratio of predictions that are correct

Hamming Loss: A measure of the number of incorrect predictions made by a classifier. It is the fraction of the wrong labels to the total number of labels.

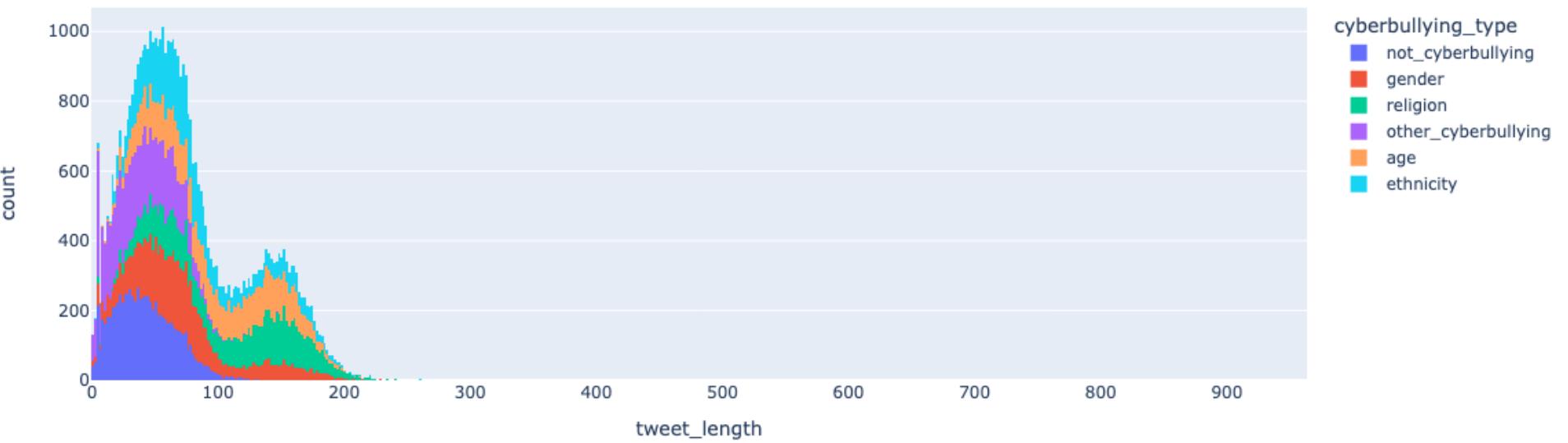
Total parameters: sum of trainable parameters (TP) and pretrained parameters (PP), which inversely accounts for inference time of deep learning models.

Performance Metrics

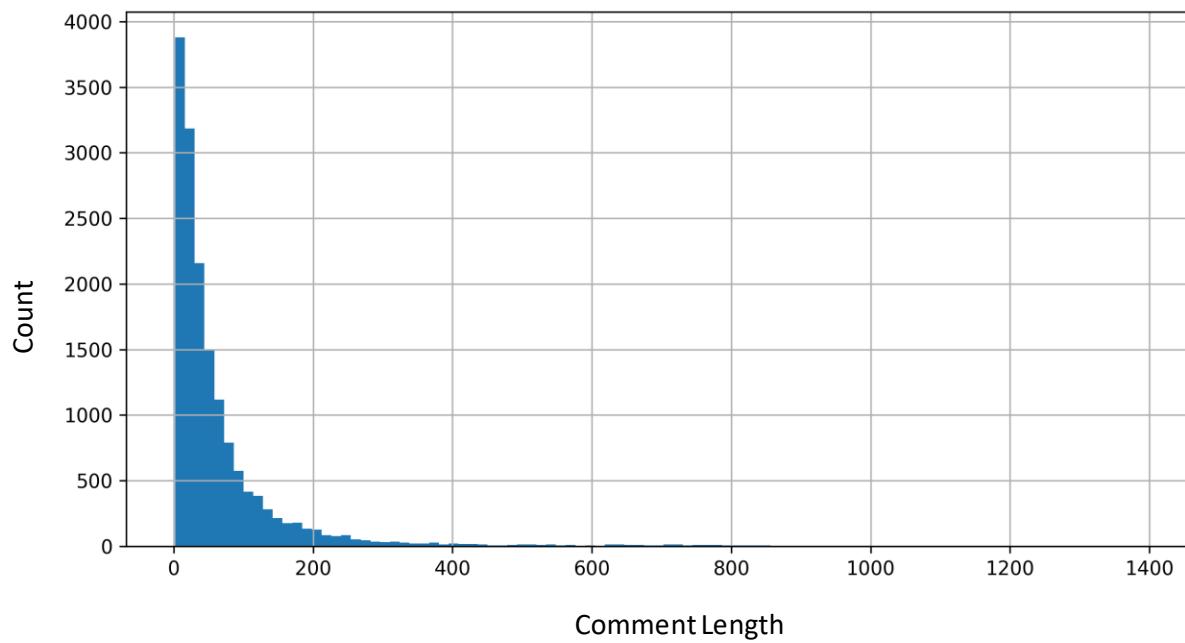


Main Hypothesis of Proposed Scheme

Dataset-1



Dataset-2



Word Count per Text Distribution

ID	Description	TP	PP	LR	Epochs
M_1	BERT-MLP	5.26×10^5	1.09×10^8	0.001	50
M_2	DistillBERT-MLP	5.26×10^5	6.68×10^7	0.001	50
M_3	SentenceBERT-GCN	1.35×10^4	2.25×10^7	0.001	50
M_4	SentenceBERT-GAN	8.31×10^6	2.25×10^7	0.001	50
M_5	GCN-Attention (1D)	1.38×10^5	0	0.001	50
M_6	GCN-Residual (1D)	2.31×10^6	0	0.001	50
M_7	PatchConvnet (2D)	1.38×10^6	0	0.001	50

Proposed models have 10 to 1000 times less total parameter count than benchmark models

Evaluation Models with trainable parameter (TP) and pretrained parameter (PP)

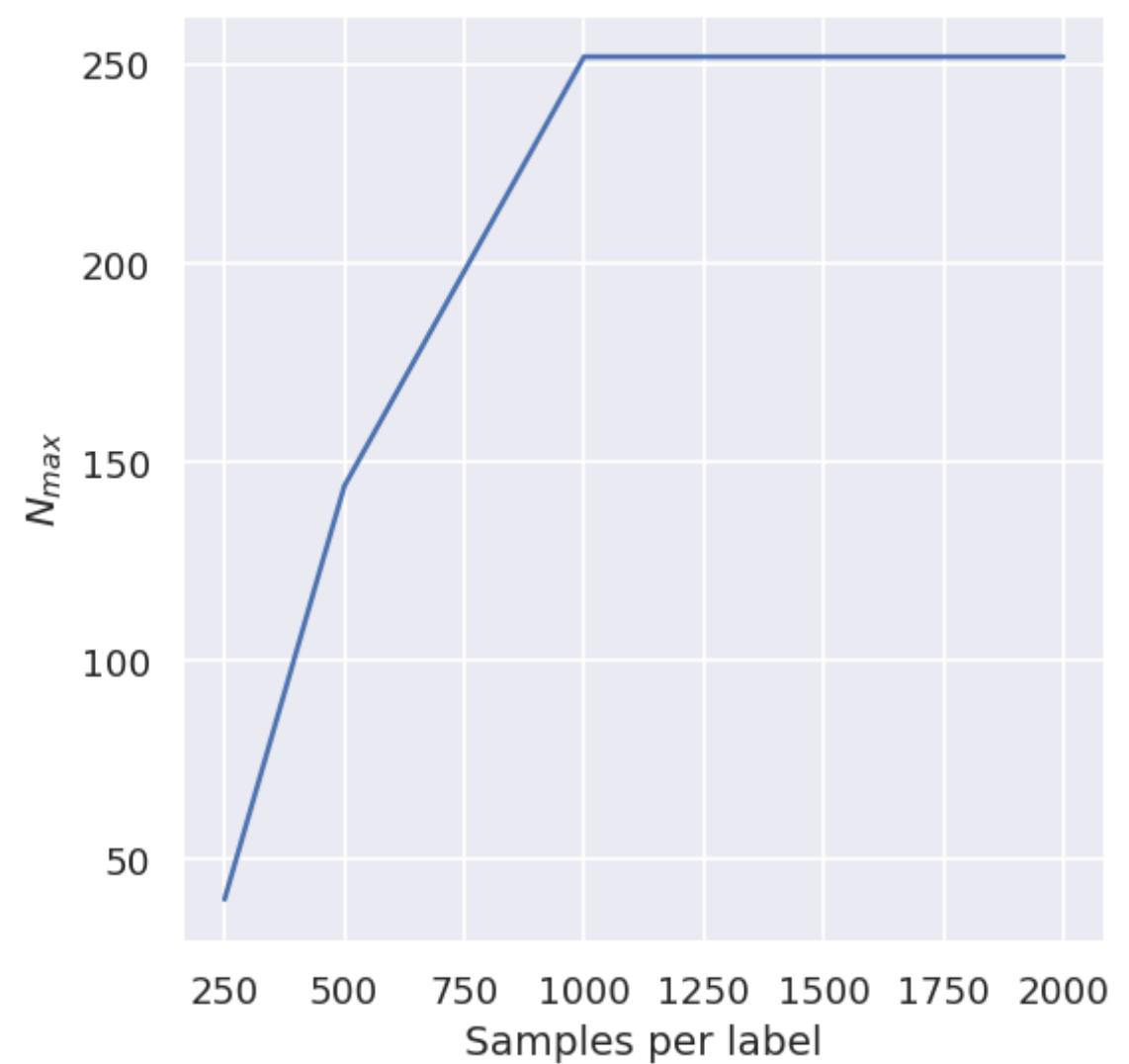
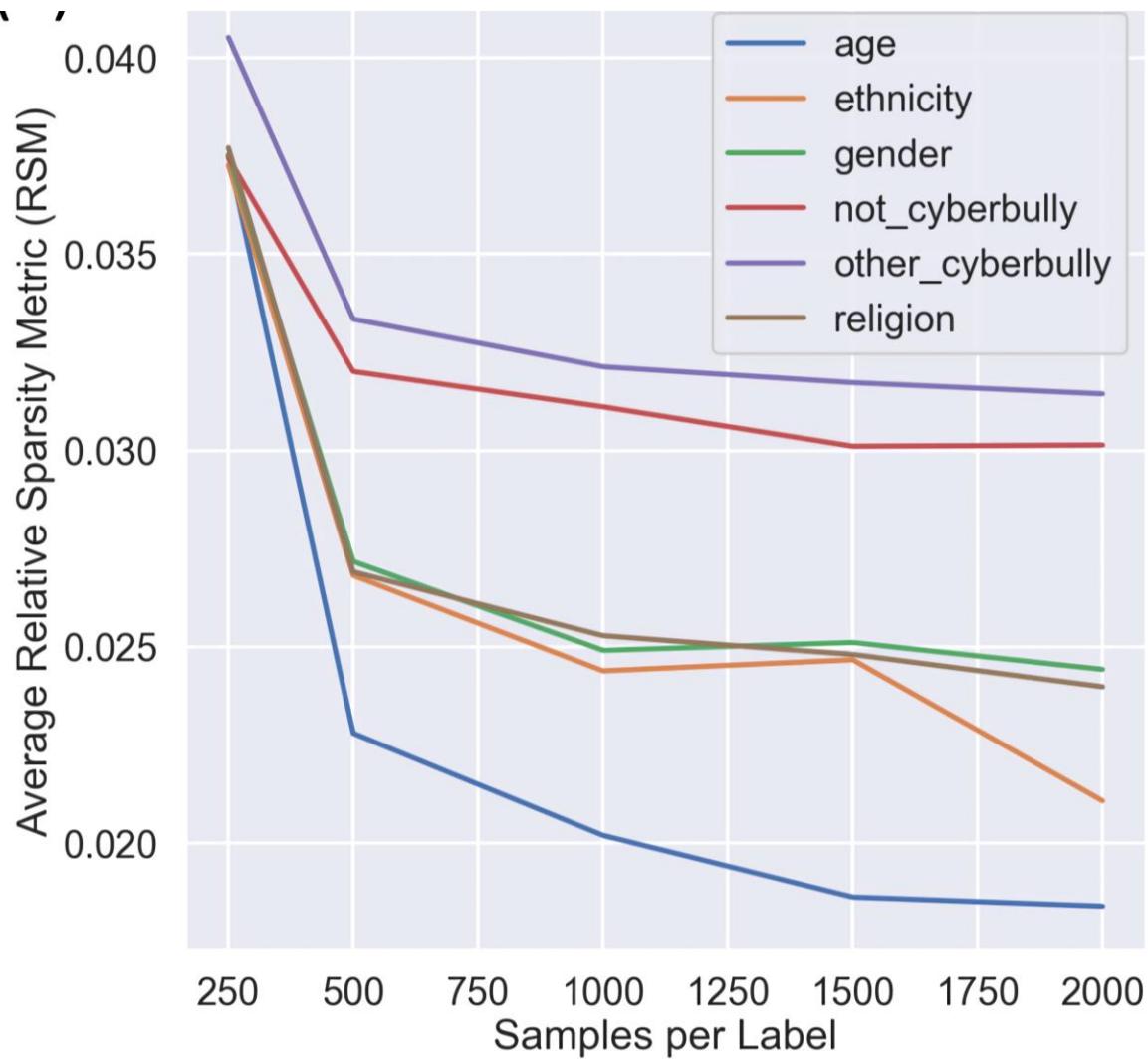
Metric	M_1	M_2	M_3	M_4	M_5	M_6	M_7
250 samples per label ($N_{max} = 40$)							
A	70.95	72.97	66.67	73.66	76.35	74.66	57.72
P	72.85	78.96	65.67	72.92	76.95	74.91	57.99
R	70.55	73.34	66.89	74.15	75.2	73.51	58.23
F	71.18	70.64	64.28	73.24	75.3	73.83	56.74
500 samples per label ($N_{max} = 144$)							
A	73.46	73.46	68.8	76.24	83.80	83.08	68.16
P	72.74	76.22	71.5	75.39	84.17	83.52	66.31
R	73.76	72.06	69.91	75.57	84.29	83.58	66.15
F	72.92	72.10	69.97	75.42	84.87	83.43	64.41
1000 samples per label ($N_{max} = 252$)							
A	71.97	73.41	64.21	75.82	85.41	84.22	82.01
P	75.04	73.30	65.16	76.00	85.14	85.11	86.22
R	72.16	72.47	65.74	75.91	85.15	84.24	81.71
F	71.05	72.28	64.91	75.78	84.46	84.25	82.01
1500 samples per label ($N_{max} = 252$)							
A	72.94	68.76	67.47	79.14	88.30	87.36	88.82
P	74.29	69.62	66.41	79.75	88.05	87.17	88.95
R	73.78	68.09	66.81	79.46	88.38	87.45	88.85
F	73.18	67.92	65.02	79.57	88.15	87.28	88.47
2000 samples per label ($N_{max} = 252$)							
A	72.63	73.49	70.19	77.82	87.42	86.37	87.96
P	72.84	74.5	71.69	77.71	87.46	86.18	88.08
R	73.05	74.42	70.39	77.75	87.39	86.38	87.79
F	72.60	74.56	69.55	77.64	87.52	86.18	87.73

Dataset-1 evaluation (scores)

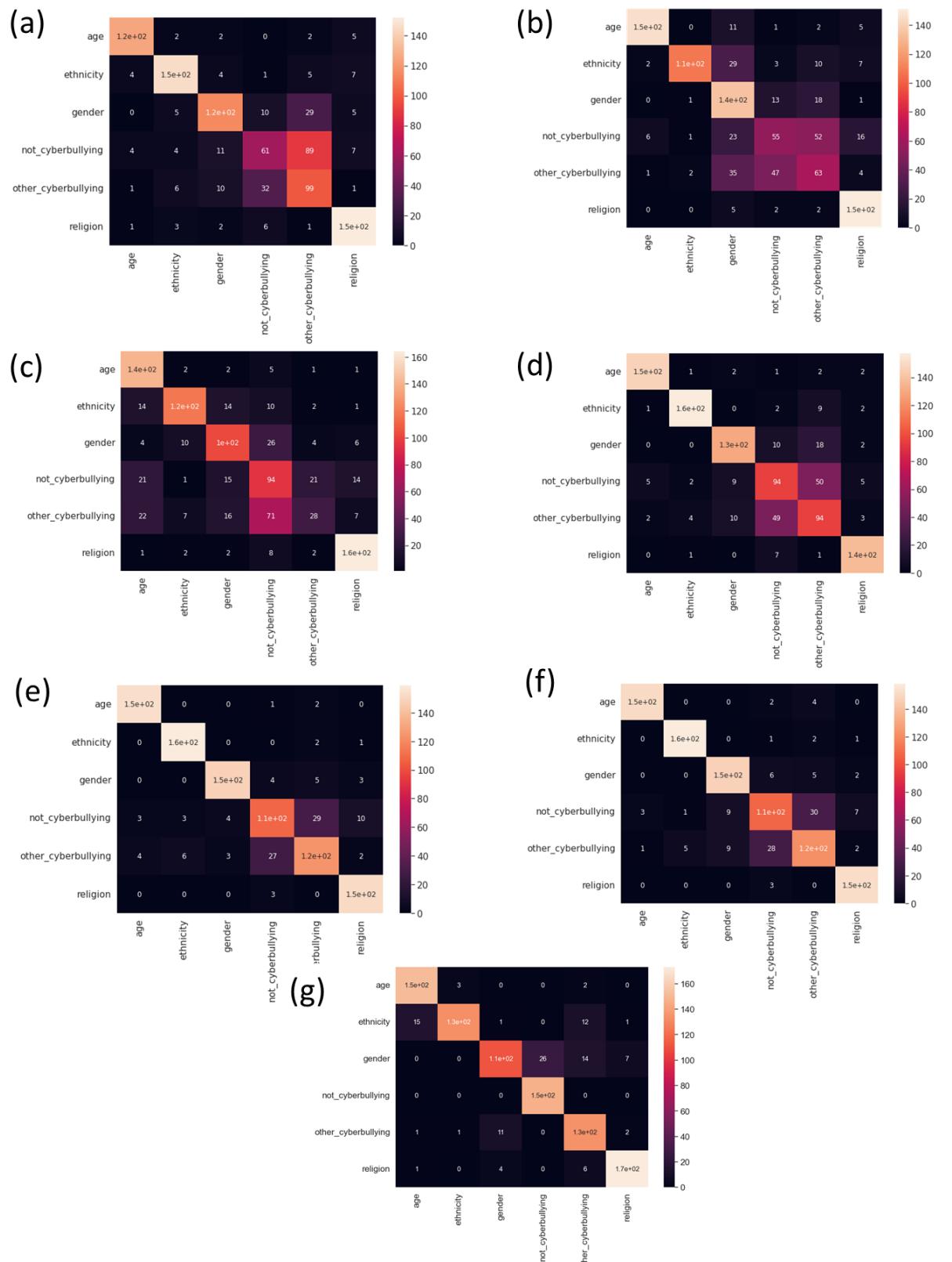
Models	H	P	R	F	A
	2000 Samples ($N_{max} = 797$)				
M_1	0.0675	0.8218	0.8472	0.8282	0.8078
M_2	0.0534	0.853	0.8587	0.8533	0.8177
M_3	0.0802	0.7988	0.8054	0.7981	0.7832
M_4	0.0679	0.818	0.8806	0.8333	0.7851
M_5	0.0953	0.8004	0.8004	0.8004	0.8005
M_6	0.0891	0.8122	0.8322	0.8128	0.8208
M_7	0.0879	0.8226	0.8226	0.8226	0.8226
4000 Samples ($N_{max} = 834$)					
M_1	0.0431	0.8781	0.9034	0.8841	0.8659
M_2	0.0398	0.8950	0.9102	0.9003	0.8686
M_3	0.0798	0.8071	0.8049	0.7989	0.7908
M_4	0.0494	0.8668	0.9056	0.8764	0.8455
M_5	0.0337	0.9262	0.9305	0.9244	0.8726
M_6	0.0360	0.9341	0.9320	0.9262	0.8512
M_7	0.0459	0.8940	0.9101	0.8976	0.8659

Dataset-2 evaluation (percentages)

A: Accuracy
 P: Precision
 R: Recall
 F: F_1 score
 H: Hamming Loss



Relationship between N_{max} and Relative Sparsity Metric



Classification Confusion Matrices (Dataset-1)

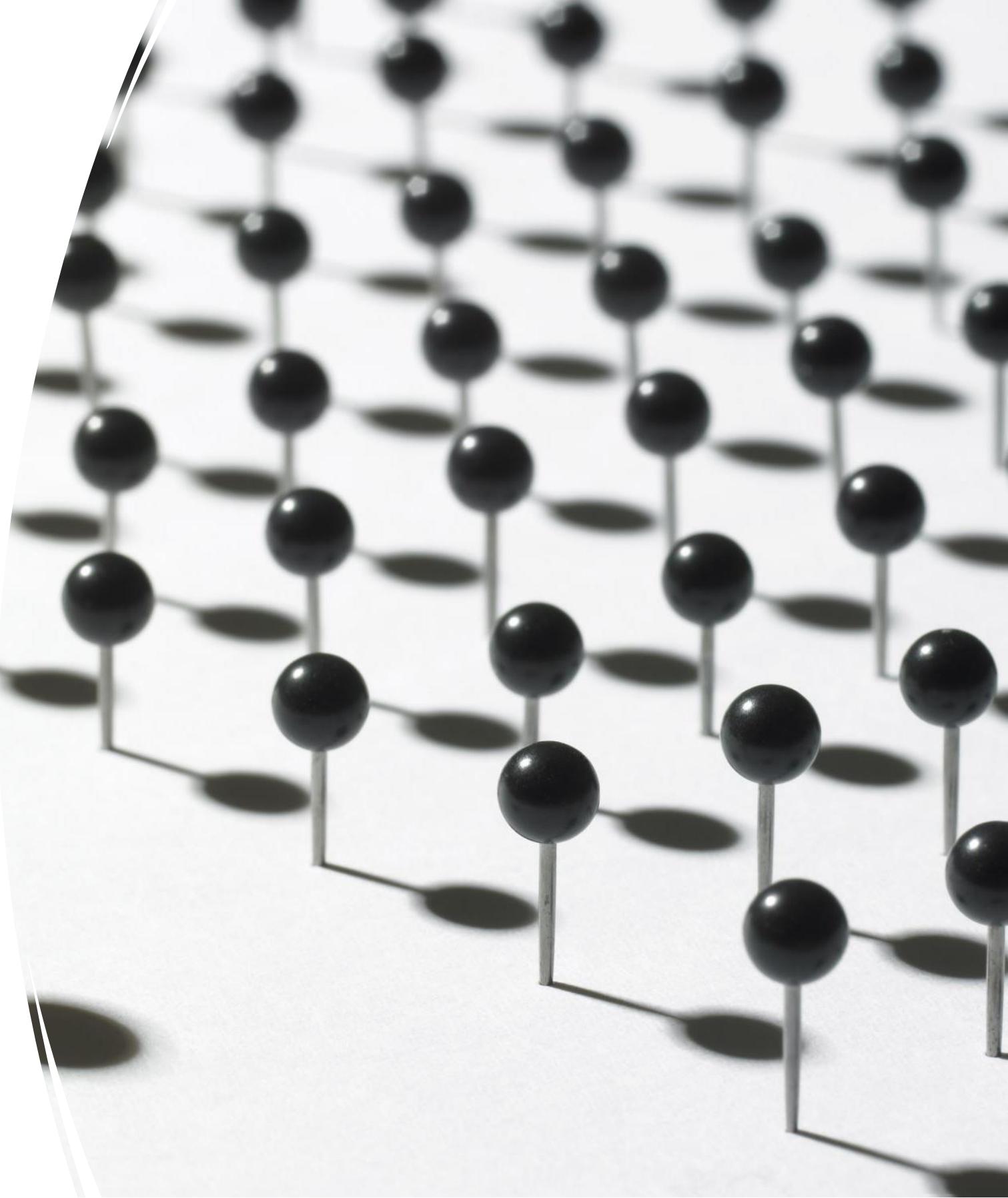
- (a) BERT-MLP
- (b) DistillBERT-MLP
- (c) SBERT-GCN
- (d) SBERT-GAN
- (e) GCN-Attention *
- (f) GCN-Residual *
- (g) PatchConvNet *

- Conclusion



- I illuminated on **some fallacies** of the **state of the art models**
- Leveraging this, I proposed a novel “**signal feature**” extraction scheme
- I further developed **three new deep learning models** for cyberbully classification
- I **evaluated** my proposed approaches over different datasets, as well as different benchmark models
- I showed that proposed approaches have greater accuracy scores with advantage of lesser parameter count

- **References**



1. <https://www.ditchthelabel.org/>
2. <https://perspectiveapi.com/>
3. Neetu Rani, Prasenjit Das, and Amit Kumar Bhardwaj. Rumor, misinformation among web: a contemporary review of rumor detection techniques during different web waves. *Concurrency and Computation: Practice and Experience*, 34(1):e6479, 2022.
4. Meaghan C McHugh, Sandra L Saperstein, and Robert S Gold. Omg u# cyberbully! an exploration of public discourse about cyberbullying on twitter. *Health Education & Behavior*, 46(1): 97-105, 2019.
5. Cynthia Van Hee, Gilles Jacobs, Chris Emmery, Bart Desmet, Els Lefever, Ben Verhoeven, Guy De Pauw, Walter Daelemans, and Veronique Hoste. Automatic detection of cyberbullying in social media text. *PloS one*, 13(10):e0203794, 2018.
6. Spiros V Georgakopoulos, Sotiris K Tasoulis, Aristidis G Vrahatis, and Vassilis P Plagianakos. Convolutional neural networks for toxic comment classification. In *Proceedings of the 10th hellenic conference on artificial intelligence*, pages 1-6, 2018.
7. Semiu Salawu, Yulan He, and Joanna Lumsden. Approaches to automated detection of cyberbullying: A survey. *IEEE Transactions on Affective Computing*, 11(1):3-24, 2017.
8. SV Drishya, S Saranya, JI Sheeba, and S Pradeep Devaneyan. Cyberbully image and text detection using convolutional neural networks. *CiIT Int. J. Fuzzy Syst*, 11(2):25-30, 2019.
9. JI Sheeba and S Pradeep Devaneyan. Impulsive intermodal cyber bullying recognition from public nets. *International Journal of Advanced Research in Computer Science*, 9(3), 2018.
10. Jason Wang, Kaiqun Fu, and Chang-Tien Lu. Sosnet: A graph convolutional network approach to fine-grained cyberbullying detection. In *2020 IEEE International Conference on Big Data (Big Data)*, pages 1699-1708. IEEE, 2020.]
11. Defu Cao, Yujing Wang, Juanyong Duan, Ce Zhang, Xia Zhu, Congrui Huang, Yunhai Tong, Bixiong Xu, Jing Bai, Jie Tong, et al. Spectral temporal graph neural network for multivariate time-series forecasting. *Advances in neural information processing systems*, 33:17766-17778, 2020.
12. Lei Yang and Hongdong Zhao. Sound classification based on multihead attention and support vector machine. *Mathematical Problems in Engineering*, 2021, 2021.
13. Sandeep Kumar Pandey, Hanumant Singh Shekhawat, and SRM Prasanna. Emotion recognition from raw speech using wavenet. In *TENCON 2019-2019 IEEE Region 10 Conference (TENCON)*, pages 1292-1297. IEEE, 2019.
14. Rui Zhao, Anna Zhou, and Kezhi Mao. Automatic detection of cyberbullying on social networks based on bullying features. In *Proceedings of the 17th international conference on distributed computing and networking*, pages 1-6, 2016.
15. Hugo Touvron, Matthieu Cord, Alaaeldin El-Nouby, Piotr Bojanowski, Armand Joulin, Gabriel Synnaeve, and Herve Jegou. Augmenting convolutional networks with attention-based aggregation. *arXiv preprint arXiv:2112.13692*, 2021.
16. Salim Alami and Omar Elbeqqali. Cybercrime profiling: Text mining techniques to detect and predict criminal activities in microblog posts. In *2015 10th International Conference on Intelligent Systems: Theories and Applications (SITA)*, pages 1-5. IEEE, 2015.
17. April Kontostathis, Kelly Reynolds, Andy Garron, and Lynne Edwards. Detecting cyberbullying: query terms and techniques. In *Proceedings of the 5th annual acm web science conference*, pages 195-204, 2013.
18. Pete Burnap and Matthew L Williams. Cyber hate speech on twitter: An application of machine classification and statistical modeling for policy and decision making. *Policy & internet*, 7(2): 223-242, 2015.

19. Kun Wang, Yanpeng Cui, Jianwei Hu, Yu Zhang, Wei Zhao, and Luming Feng. Cyberbullying detection, based on the fasttext and word similarity schemes. *ACM Transactions on Asian and Low-Resource Language Information Processing (TALLIP)*, 20(1):1-15, 2020.
20. Steven Zimmerman, Udo Kruschwitz, and Chris Fox. Improving hate speech detection with deep learning ensembles. In *Proceedings of the eleventh international conference on language resources and evaluation (LREC 2018)*, 2018.
21. Gretel Liz De la Pena Sarracen, Reynaldo Gil Pons, Carlos Enrique Muniz Cuza, and Paolo Rosso. Hate speech detection using attention-based lstm. *EVALITA evaluation of NLP and speech tools for Italian*, 12:235, 2018.
22. Nijia Lu, Guohua Wu, Zhen Zhang, Yitao Zheng, Yizhi Ren, and Kim-Kwang Raymond Choo. Cyberbullying detection in social media text based on character-level convolutional neural network with shortcuts. *Concurrency and Computation: Practice and Experience*, 32(23):e5627, 2020.
23. Raunak Joshi and Abhishek Gupta. Performance comparison of simple transformer and res-cnn-bilstm for cyberbullying classification. *arXiv preprint arXiv:2206.02206*, 2022.
24. Jianwei Zhang, Taiga Otomo, Lin Li, and Shinsuke Nakajima. Cyberbullying detection on twitter using multiple textual features. In *2019 IEEE 10th International Conference on Awareness Science and Technology (iCAST)*, pages 1-6. IEEE, 2019.
25. Rui Zhao and Kezhi Mao. Cyberbullying detection based on semantic-enhanced marginalized denoising auto-encoder. *IEEE Transactions on Affective Computing*, 8(3):328-339, 2016.
26. Lu Cheng, Ruocheng Guo, Yasin N Silva, Deborah Hall, and Huan Liu. Modeling temporal patterns of cyberbullying detection with hierarchical attention networks. *ACM/IMS Transactions on Data Science*, 2(2):1-23, 2021.
27. Suyu Ge, Lu Cheng, and Huan Liu. Improving cyberbullying detection with user interaction. In *Proceedings of the Web Conference 2021*, pages 496-506, 2021.
28. Sudhanshu Mishra, Shivangi Prasad, and Shubhanshu Mishra. Multilingual joint fine-tuning of transformer models for identifying trolling, aggression and cyberbullying at trac 2020. In *Proceedings of the Second Workshop on Trolling, Aggression and Cyberbullying*, pages 120-125, 2020.
29. Tasnim Ahmed, Shahriar Ivan, Mohsinul Kabir, Hasan Mahmud, and Kamrul Hasan. Performance analysis of transformer-based architectures and their ensembles to detect trait-based cyberbullying. *Social Network Analysis and Mining*, 12(1):1-17, 2022.
30. Bandeh Ali Talpur and Declan O'Sullivan. Multi-class imbalance in text classification: A feature engineering approach to detect cyberbullying in twitter. In *Informatics*, volume 7, page 52. MDPI, 2020.
31. Nabi Rezvani, Amin Beheshti, and Alireza Tabebordbar. Linking textual and contextual features for intelligent cyberbullying detection in social media. In *Proceedings of the 18th International Conference on Advances in Mobile Computing & Multimedia*, pages 3-10, 2020.
32. Joan Plepi and Lucie Flek. Perceived and intended sarcasm detection with graph attention networks. *arXiv preprint arXiv:2110.04001*, 2021.
33. Devin Soni and Vivek Singh. Time reveals all wounds: Modeling temporal dynamics of cyberbullying sessions. In *12th International AAAI Conference on Web and Social Media, ICWSM 2018*, pages 684-687. AAAI press, 2018.
34. Vimala Balakrishnan, Shahzaib Khan, and Hamid R Arabnia. Improving cyberbullying detection using twitter users' psychological features and machine learning. *Computers & Security*, 90:101710, 2020.

35. Krishanu Maity, Abhishek Kumar, and Sriparna Saha. A multi-task multi-modal framework for sentiment and emotion aided cyberbully detection. *IEEE Internet Computing*, 2022.
36. Hugo Rosa, David Matos, Ricardo Ribeiro, Luisa Coheur, and Joao P Carvalho. A “deeper” look at detecting cyberbullying in social networks. In *2018 international joint conference on neural networks (IJCNN)*, pages 1-8. IEEE, 2018.
37. Lihao Ge and Teng-Sheng Moh. Improving text classification with word embedding. In *2017 IEEE International Conference on Big Data (Big Data)*, pages 1796-1805. IEEE, 2017.
38. Lingfei Wu, Ian EH Yen, Kun Xu, Fangli Xu, Avinash Balakrishnan, Pin-Yu Chen, Pradeep Ravikumar, and Michael J Witbrock. Word mover’s embedding: From word2vec to document embedding. *arXiv preprint arXiv:1811.01713*, 2018.
39. Joseph Lilleberg, Yun Zhu, and Yanqing Zhang. Support vector machines and word2vec for text classification with semantic features. In *2015 IEEE 14th International Conference on Cognitive Informatics & Cognitive Computing (ICCI*CC)*, pages 136-140. IEEE, 2015.
40. Thomas N Kipf and Max Welling. Semi-supervised classification with graph convolutional networks. *arXiv preprint arXiv:1609.02907*, 2016.
41. Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30, 2017.
42. Aaron van den Oord, Sander Dieleman, Heiga Zen, Karen Simonyan, Oriol Vinyals, Alex Graves, Nal Kalchbrenner, Andrew Senior, and Koray Kavukcuoglu. Wavenet: A generative model for raw audio. *arXiv preprint arXiv:1609.03499*, 2016.
43. Jianfeng Zhao, Xia Mao, and Lijiang Chen. Speech emotion recognition using deep 1d & 2d cnn lstm networks. *Biomedical signal processing and control*, 47:312-323, 2019.
44. Ankit Thakkar, Dhara Mungra, Anjali Agrawal, and Kinjal Chaudhari. Improving the performance of sentiment analysis using enhanced preprocessing technique and artificial neural network. *IEEE transactions on affective computing*, 2022.
45. Nltk wordnet. https://www.nltk.org/_modules/nltk/stem/wordnet.html.
46. Jesus Selva. Convolution-based trigonometric interpolation of band-limited signals. *IEEE transactions on signal processing*, 56(11):5465-5477, 2008.
47. Junaid Iqbal Khan and Usman Zabit. On two fourier transform-based methods for estimation of displacement and parameters of self-mixing interferometry over major optical feedback regimes. *IEEE Sensors Journal*, 21(9):10610-10617, 2021.
48. Min Lin, Qiang Chen, and Shuicheng Yan. Network in network. *arXiv preprint arXiv:1312.4400*, 2013.
49. Jingshan Huang, Binqiang Chen, Bin Yao, and Wangpeng He. Ecg arrhythmia classification using stft-based spectrogram and convolutional neural network. *IEEE access*, 7:92871-92880, 2019.
50. fine-grained balanced cyberbullying dataset. https://ieee-dataport.org/open-access/_fine-grained-balanced-cyberbullying-dataset, .
51. toxic comment classification challenge. <https://www.kaggle.com/c/jigsaw-toxic-comment-classification-challenge> .
52. Sentencetransformer. <https://huggingface.co/sentence-transformers/all-MiniLM-L6-v2>, .

53. Scott M Lundberg and Su-In Lee. A unified approach to interpreting model predictions. *Advances in neural information processing systems*, 30, 2017.
54. John Muradeli. ssqueezepy. *GitHub*. Note: <https://github.com/OverLordGoldDragon/ssqueezepy/>, 2020. doi: 10.5281/zenodo.5080508.
55. Zhou Wang, Eero P Simoncelli, and Alan C Bovik. Multiscale structural similarity for image quality assessment. In *The Thirty-Seventh Asilomar Conference on Signals, Systems & Computers, 2003*, volume 2, pages 1398-1402. ieee, 2003.
56. Jinyin Chen, Yi-tao Yang, Ke-ke Hu, Hai-bin Zheng, and Zhen Wang. Dad-mcnn: Ddos attack detection via multi-channel cnn. In *Proceedings of the 2019 11th International Conference on Machine Learning and Computing*, pages 484-488, 2019.
57. Nitin Kumar Chauhan and Krishna Singh. Impact of variation in number of channels in cnn classification model for cervical cancer detection. In *2021 9th International Conference on Reliability, Infocom Technologies and Optimization (Trends and Future Directions)(ICRITO)*, pages 1-6. IEEE, 2021.