



CS779 Project

# Adversarial Techniques In NLP

Dept. of CSE, IIT Kanpur

**Under the guidance of**

Dr. Ashutosh Modi  
Assistant Professor  
Dept. of CSE

**Group 3:**

Abhishek Krishna- 20111002  
Deeksha Arora- 20111017  
Preeti Singh- 20111044  
Sambrant Maurya- 20111054  
Shruti Sharma- 20111061

# Adversarial attacks on text

<b>Original Input</b>	Connoisseurs of Chinese film will be pleased to discover that Tian's meticulous talent has not withered during his enforced hiatus.	Prediction: <b><u>Positive (77%)</u></b>
<b>Adversarial example [Visually similar]</b>	<b><u>Aonnoisseurs</u></b> of Chinese film will be pleased to discover that Tian's meticulous talent has not withered during his enforced hiatus.	Prediction: <b><u>Negative (52%)</u></b>
<b>Adversarial example [Semantically similar]</b>	Connoisseurs of Chinese <b><u>footage</u></b> will be pleased to discover that Tian's meticulous talent has not withered during his enforced hiatus.	Prediction: <b><u>Negative (54%)</u></b>

# Adversarial attacks on text

<b>Original Input</b>	Connoisseurs of Chinese film will be pleased to discover that Tian's meticulous talent has not withered during his enforced hiatus.	Prediction: <b><u>Positive (77%)</u></b>
<b>Adversarial example [Visually similar]</b>	<b><u>Aonnoisseurs</u></b> of Chinese film will be pleased to discover that Tian's meticulous talent has not withered during his enforced hiatus.	Prediction: <b><u>Negative (52%)</u></b>
<b>Adversarial example [Semantically similar]</b>	Connoisseurs of Chinese <b><u>footage</u></b> will be pleased to discover that Tian's meticulous talent has not withered during his enforced hiatus.	Prediction: <b><u>Negative (54%)</u></b>

- Neural Networks are susceptible to adversarial inputs.

**Textual adversaries:** Try to “fool” the Neural Network by introducing perturbations in the input text.

# Motivation



<b>Real Comment:</b> admitting im not going to read this (...)
Adv-comment: <i>hes a conservative from a few months ago</i>
<b>Prediction Change:</b> Real News → Fake News

Image 1: An adversarial comment causes misclassification of a fake news detector

# Motivation



NEWS

## Man arrested after ‘good morning’ post mistranslated by Facebook as ‘attack them’

Israeli police arrested a Palestinian man after his “good morning” post was translated by Facebook as “attack them.”

Image 1: An adversarial comment causes misclassification of a fake news detector

Image 2: Facebook NMT mistakes input word for another which differs by a single character in Arabic, and wrecks havoc

Image 1: <https://arxiv.org/pdf/2009.01048.pdf>,

Image 2: [Article from theguardian.com](https://www.theguardian.com/technology/2018/jun/05/facebook-mistranslation-arabic)

# Motivation



NEWS

## Man arrested after 'good morning' post mistranslated by Facebook as 'attack them'

Israeli police arrested a Palestinian man after his “good morning” post was translated by Facebook as “attack them.”

<b>Real Comment:</b> admitting im not going to read this (...)
<b>Adv-comment:</b> <i>hes a conservative from a few months ago</i>
<b>Prediction Change:</b> <b>Real News</b> → <b>Fake News</b>

Image 1: An adversarial comment causes misclassification of a fake news detector

Image 2: Facebook NMT mistakes input word for another which differs by a single character in Arabic, and wrecks havoc

Defence of Neural Networks against adversarial attacks is crucial.

Image 1: <https://arxiv.org/pdf/2009.01048.pdf>,

Image 2: [Article from theguardian.com](https://www.theguardian.com/technology/2018/jun/05/woody-allen-me-too-movement)



# Existing Research

- A recent and very active research area



# Existing Research

- A recent and very active research area
- Some latest attack methods:
  - DeepWordBug, Gao et.al (2018)[1]
  - TextBugger, Li et.al. (2018)[2]
  - \*Textfooler, Jin et.al. (2019)[3]
  - CAM-RWR (Pruthi et.al.,2019)[4]
  - \*BERT-Attack, Li et.al (2020)[5]
  - Bae, Garg et.al (2020)[6]
  - \*Adv-OLM, Malik et.al (2021)[7]

\*Performance benchmarks





# Designing Textual attacks

- A text attack is built on two main components:
  - **Search methods:** Finding which characters/words/sentences to perturb.
  - **Transformations:** Replacing the chosen characters/words/sentences with an aim to cause misprediction of the target model.



# Designing Textual attacks

- A text attack is built on two main components:
  - **Search methods:** Finding which characters/words/sentences to perturb.
  - **Transformations:** Replacing the chosen characters/words/sentences with an aim to cause misprediction of the target model.
- Desirable properties of a good attack method:
  - Accuracy under attack: *low*
  - Attack success rate: *high*
  - Perturbation percentage: *low*
  - Avg no of queries: *low*



# Designing Textual attacks

- A text attack is built on two main components:
  - **Search methods:** Finding which characters/words/sentences to perturb.
  - **Transformations:** Replacing the chosen characters/words/sentences with an aim to cause misprediction of the target model.
- Desirable properties of a good attack method:
  - Accuracy under attack: *low*
  - Attack success rate: *high*
  - Perturbation percentage: *low*
  - Avg no of queries: *low*
- Most of the existing attacks don't score high on all 4 properties
  - Task at hand: Design an attack that does!

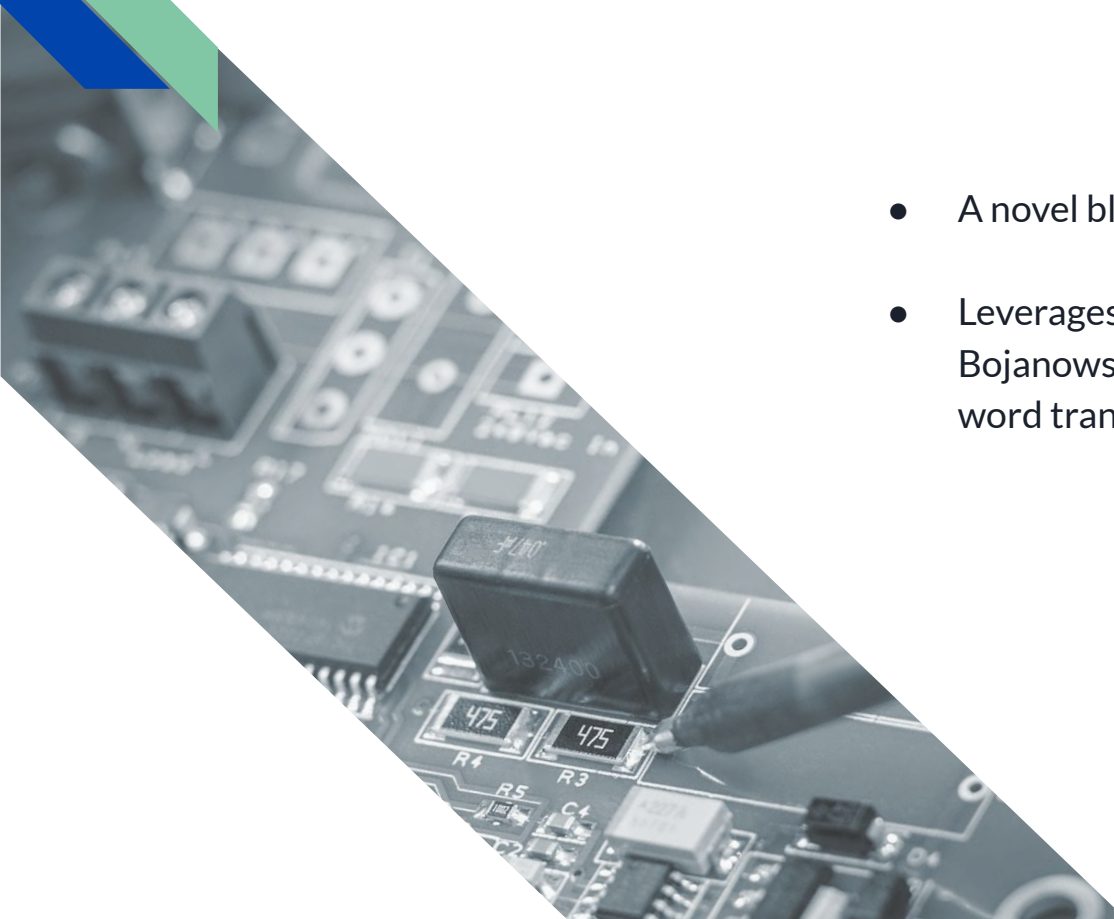


# Experiments

- Tried transformations such as Misspelling Oblivious Word (MOE), GloVe, Word2Vec, FastText and finally concluded to use FastText.
- Experimenting on search-methods available in TextAttack framework:
  - **Greedy Search:** Rapid, but poor performance
  - **Beam Search:** Descent performance, but slow for larger text inputs
  - **Alzantot genetic algorithm[12]:** Poor performance, slow
  - **Particle Swarm Optimization[11]:** Best performance, but the slowest
  - **GreedyWordSwapWIR[9]:** Good performance and fast
    - Provides a good tradeoff between performance and speed
- Compared our attack's performance against existing recipes.



# Introducing “FastAttack”

- A novel blackbox attack
  - Leverages “FastText” embeddings by Bojanowski et.al[8], (Facebook Research) for word transformations
- 



# FastText Embeddings

- FastText: Open source, free, light-weighted library, developed by Facebook Research.
- Allows to learn text representation and text classifiers.



# FastText Embeddings

- FastText: Open source, free, light-weighted library, developed by Facebook Research
- Allows to learn text representation and text classifiers.
- Additional ability to obtain word vectors for out-of-vocabulary words.
- Two models for computing word representations\*:
  - CBOW(Continuous Bag of Words)
  - Skip-gram

\*<https://fasttext.cc/docs/en/crawl-vectors.html>



# FastText Embeddings

- FastText: Open source, free, light-weighted library, developed by Facebook Research
- Allows to learn text representation and text classifiers.
- Additional ability to obtain word vectors for out-of-vocabulary words.
- Two models for computing word representations:
  - CBOW(Continuous Bag of Words)
  - Skip-gram
- CBOW: predict target word according to its context.
  - Context represented as BOW containing fixed size window around target.





# FastText Embeddings

- FastText: Open source, free, light-weighted library, developed by Facebook Research
- Allows to learn text representation and text classifiers.
- Additional ability to obtain word vectors for out-of-vocabulary words.
- Two models for computing word representations:
  - CBOW(Continuous Bag of Words)
  - Skip-gram
- CBOW: predict target word according to its context.
  - Context represented as BOW containing fixed size window around target.
- Skip-gram: learns to predict target word thanks to a near-by word.
  - In practice, skip-gram models works better with subword information than cbow.



# Embeddings used

- FastText, Word2Vec and GLoVE
  - Trained on the **text8** corpus\*, consisting of first 100,000,000 bytes of plain text from Wikipedia



# Embeddings used

- FastText, Word2Vec and GLoVE
  - Trained on the **text8** corpus\*, consisting of first 100,000,000 bytes of plain text from Wikipedia
- **Hyperparameters:**

Embedding	Model type	Vector size	Window size	Learning rate	Epochs	Threads	Training time (s)	Total words
FastText	Skipgram	100	5	0.025	5	3	932	17005207
Word2Vec	Skipgram	100	5	0.025	5	3	670	17005207
Glove	-	100	5	0.05	5	3	628	17005207



# Target models and Datasets

- **HuggingFace\* models with datasets:**
  - **Albert-Base-V2**
    - AG News, dataset ag\_news, split test
    - IMDB, dataset imdb, split test
    - Movie Reviews, dataset rotten\_tomatoes, split test
    - Yelp Polarity, dataset yelp\_polarity, split test
  - **Bert-Base-uncased**
    - AG News, dataset ag\_news, split test
    - IMDB, dataset imdb, split test
    - Movie Reviews, dataset rotten\_tomatoes, split test
    - Yelp Polarity, dataset yelp\_polarity, split test

\*<https://huggingface.co/textattack>



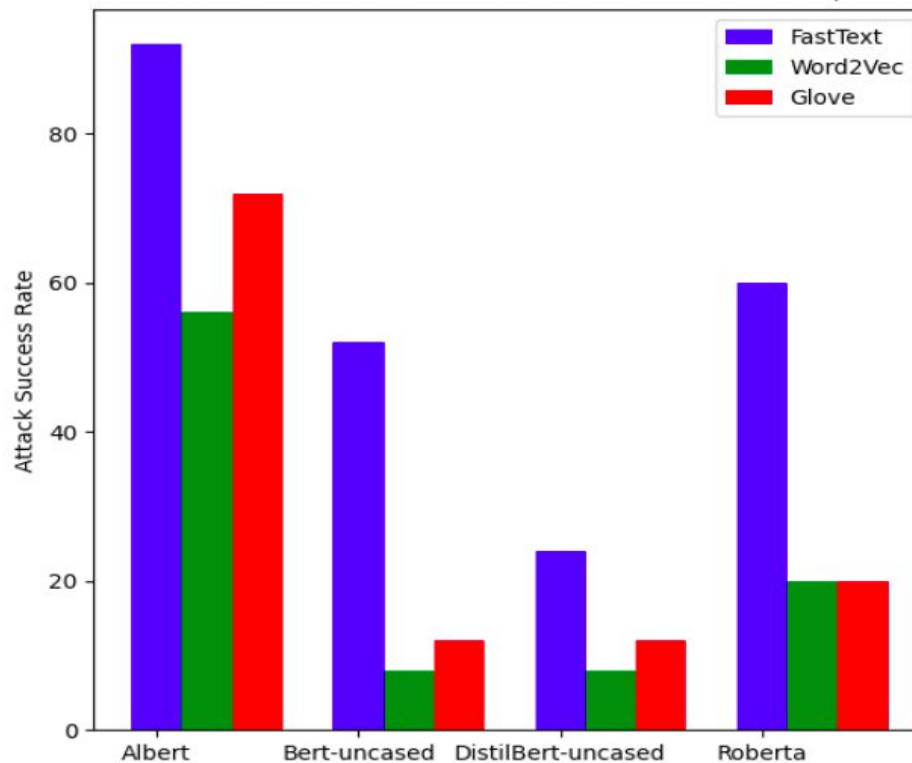
# Target models and Datasets

- **HuggingFace\* models with datasets:**
  - **Distilbert-Base-cased**
    - Quora Question Pairs, dataset glue, qqp, split validation
    - SST-2, dataset glue, sst2, split validation
  - **Distilbert-Base-uncased**
    - AG News, dataset ag\_news, split test
    - IMDB, dataset imdb, split test
  - **Roberta-Base**
    - AG News, dataset ag\_news, split test
    - IMDB, dataset imdb, split test
    - Movie Reviews, dataset rotten\_tomatoes, split test
    - SST-2, dataset glue, sst2, split validation

\*<https://huggingface.co/textattack>

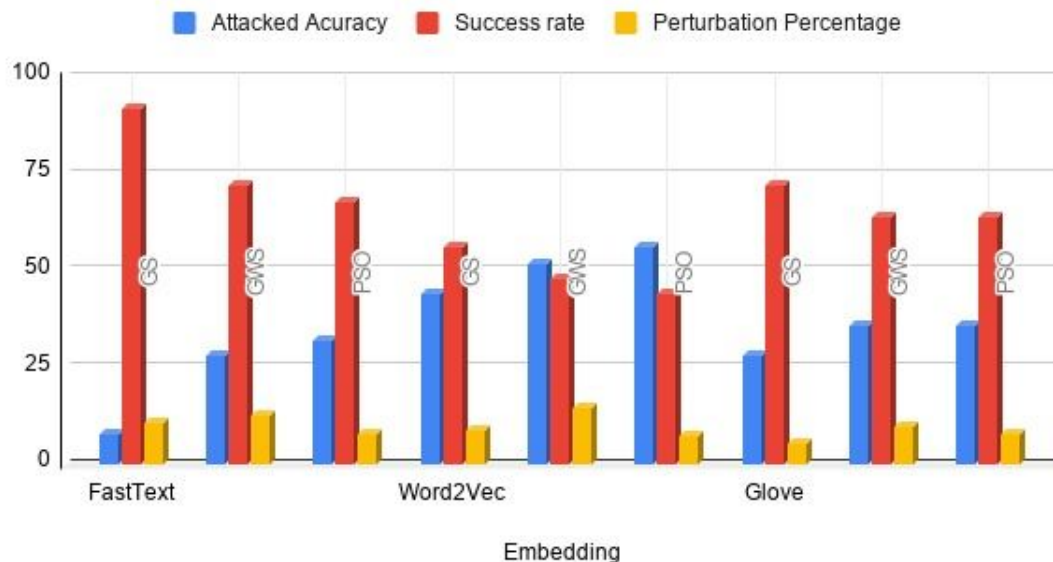
# FastText vs Word2Vec vs GloVe

FastText vs Word2Vec vs GloVe Attack Success Rate on AG News, Greedy Strategy



# FastText vs Word2Vec vs GloVe

Comparison of Fasttext, Word2Vec and Glove on yelp polarity



**Search methods:** GS: Greedy Search, GWS: GreedyWordSwapWIR[9], PSO: ParticleSwarmOptimization[11]



# FastAttack implementation

- **Search method:** Uses GreedyWordSwapWIR algorithm[9] with weighted saliency to find the words **W** to be perturbed.
  - A tradeoff between speed and accuracy





# FastAttack implementation

- **Search method:** Uses GreedyWordSwapWIR algorithm[9] with weighted saliency to find the words  $\mathbf{W}$  to be perturbed.
  - A tradeoff between speed and accuracy
- **Transformation technique:** Replace the chosen word  $\mathbf{w} \in \mathbf{W}$  with its closest neighbor  $\mathbf{w}'$  in the FastText embedding under the following constraints:
  - $\mathbf{w}'$  should not have  $\mathbf{w}$  as a substring.
  - $\mathbf{w}$  should not be a stop-word as defined in the NLTK corpora\*.
  - For the same sentence, the same word should not be modified twice.
  - The case of the replaced word should be preserved.

\*<https://www.nltk.org/book/ch02.html>



# FastAttack implementation

- **Search method:** Uses GreedyWordSwapWIR algorithm[9] with weighted saliency to find the words  $\mathbf{W}$  to be perturbed.
  - A tradeoff between speed and accuracy
- **Transformation technique:** Replace the chosen word  $\mathbf{w} \in \mathbf{W}$  with its closest neighbor  $\mathbf{w}'$  in the FastText embedding under the following constraints:
  - $\mathbf{w}'$  should not have  $\mathbf{w}$  as a substring.
  - $\mathbf{w}$  should not be a stop-word as defined in the NLTK corpora.
  - For the same sentence, the same word should not be modified twice.
  - The case of the replaced word should be preserved.
- **Goal function:** UntargetedClassification
- We use the TextAttack framework, Morris et.al.[10] (2020) to frame our attack.



# Performance of FastAttack

# An Illustration

Sci/tech (98%) --> World (73%) (**FastAttack**)

E-mail scam targets police chief Wiltshire Police warns about "phishing" after its fraud squad chief was targeted.

E-mail scam targets police chief Wiltshire Police warns about "furnishing" after its bribery squad chief was targeted.

Sci/tech (98%) --> World (54%) (**TextFooler**)

E-mail scam targets police chief Wiltshire Police warns about "phishing" after its fraud squad chief was targeted.

E-mail scam targets police chief Wiltshire Police warns about "phishing" after its hoax battalion leiter was targeted.

# An Illustration

Sci/tech (98%) --> World (73%) (**FastAttack**)

E-mail scam targets police chief Wiltshire Police warns about "phishing" after its fraud squad chief was targeted.

E-mail scam targets police chief Wiltshire Police warns about "furnishing" after its bribery squad chief was targeted.

Sci/tech (98%) --> World (54%) (**TextFooler**)

E-mail scam targets police chief Wiltshire Police warns about "phishing" after its fraud squad chief was targeted.

E-mail scam targets police chief Wiltshire Police warns about "phishing" after its hoax battalion leiter was targeted.

Sci/tech (98%) --> World (78%) (**DeepWordBug**)

E-mail scam targets police chief Wiltshire Police warns about "phishing" after its fraud squad chief was targeted.

E-mail scam targets police chief Wiltshire Police warns about "phishin" after its fruad squad chief was targeted.

Sci/tech (98%) --> World (82%) (**BERT-Attack**)

E-mail scam targets police chief Wiltshire Police warns about "phishing" after its fraud squad chief was targeted.

E-mail scam targets police chief Wiltshire Police warns about "phitism" after its fraudulent squad chief was targeted.

# FastAttack vs CAM-RWR (Pruthi et.al.,2019)

Model	Dataset	Attack method	Successful attacks (Out of 25)	Original accuracy	Attacked accuracy	Success rate	Perturbed %	Avg num queries	Attack time (in seconds)
ALBERT-base-v2	YELP-polarity	fastattack	<b>19</b>	100%	<b>24%</b>	<b>76%</b>	9.47%	<b>241.72</b>	<b>139.27</b>
		cam-rwr	2		92%	8%	<b>4.28%</b>	1666.0	924.62
	IMDB	fastattack	<b>17</b>	89.29%	<b>28.57%</b>	<b>68.0%</b>	2.3%	<b>420.24</b>	<b>271.88</b>
		cam-rwr	8		60.71%	32.0%	<b>0.53%</b>	3843.64	2493.46
BERT-base-uncased	YELP-polarity	fastattack	<b>18</b>	100%	<b>28.0%</b>	<b>72.0%</b>	11.07%	<b>241.48</b>	<b>128.77</b>
		cam-rwr	1		96.0%	4.0%	<b>6.67%</b>	1665.96	848.89
	IMDB	fastattack	<b>20</b>	92.59%	<b>18.52%</b>	<b>80.0%</b>	5.61%	<b>428.4</b>	<b>251.87</b>
		cam-rwr	10		55.56%	40.0%	<b>0.48%</b>	3868.6	2192.64

# FastAttack vs Faster Genetic Algorithm (Jia et.al.,2019)

Model	Dataset	Attack method	Successful attacks (Out of 25)	Original accuracy	Attacked accuracy	Success rate	Perturbed %	Avg num queries	Attack time (in seconds)
ALBERT-base-v2	YELP-polarity	fastattack	<b>19</b>	100%	<b>24%</b>	<b>76%</b>	<b>11.43%</b>	<b>242.76</b>	<b>114.78</b>
		FGA_Jia	17		32%	68%	11.6%	4762.84	5833.25
	IMDB	fastattack	<b>16</b>	96.15%	<b>34.62%</b>	<b>64.0%</b>	<b>17.56%</b>	<b>63.52</b>	<b>28.86</b>
		FGA_Jia	10		57.69%	40.0%	18.3%	3358.08	3900.52
BERT-base-uncased	YELP-polarity	fastattack	<b>21</b>	100%	<b>16%</b>	<b>84%</b>	12.51%	<b>238.84</b>	<b>100.20</b>
		FGA_Jia	18		28%	72%	<b>10.12%</b>	5123.12	4985.12
	IMDB	fastattack	<b>16</b>	92.59%	<b>33.33%</b>	<b>64.0%</b>	21.46%	<b>65.16</b>	<b>32.84</b>
		FGA_Jia	9		59.26%	36.0%	<b>15.85%</b>	3292.84	3728.22

# FastAttack vs TextFooler (Jin et.al.,2019)

Model	Dataset	Attack method	Successful attacks (Out of 25)	Original accuracy	Attacked accuracy	Success rate	Perturbed %	Avg num queries	Attack time (in seconds)
ALBERT-base-v2	YELP-polarity	fastattack	19	100%	24.0%	76.0%	11.43%	<b>242.76</b>	<b>210.15</b>
		textfooler	<b>24</b>		<b>4.0%</b>	<b>96%</b>	<b>9.11%</b>	480.92	316.80
	IMDB	fastattack	22	89.29%	10.71%	88.0%	<b>5.2%</b>	<b>434.92</b>	<b>426.02</b>
		textfooler	<b>25</b>		0.0%	100%	8.35%	674.2	473.95
BERT-base-uncased	YELP-polarity	fastattack	21	100%	16.0%	84.0%	12.51%	<b>238.84</b>	<b>137.26</b>
		textfooler	<b>25</b>		<b>0.0%</b>	<b>100%</b>	<b>10.89%</b>	479.08	237.07
	IMDB	fastattack	21	92.59%	14.81%	84.0%	<b>5.29%</b>	<b>419.72</b>	<b>342.94</b>
		textfooler	<b>25</b>		<b>0.0%</b>	<b>100.0%</b>	7.99%	663.64	400.18



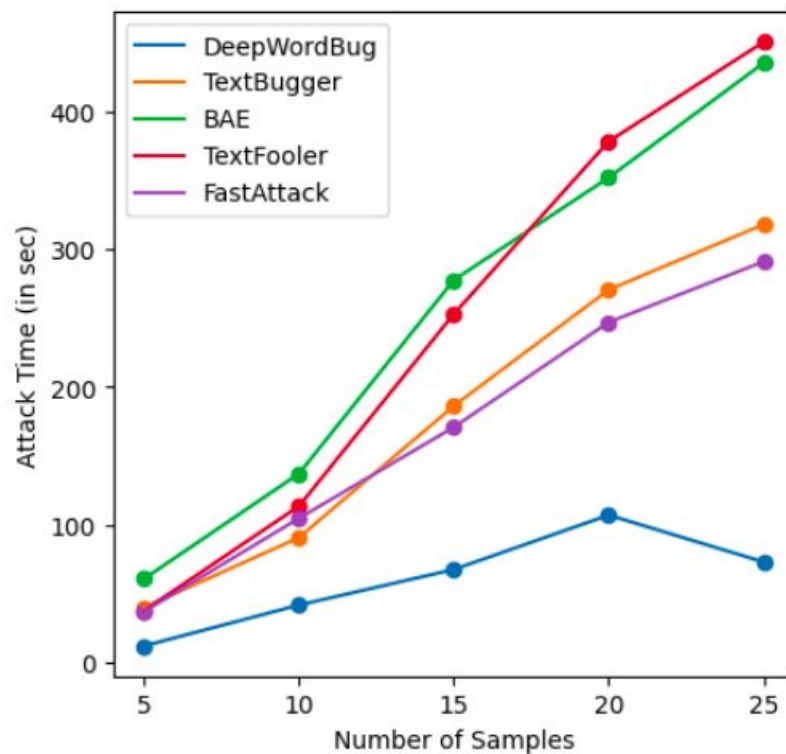
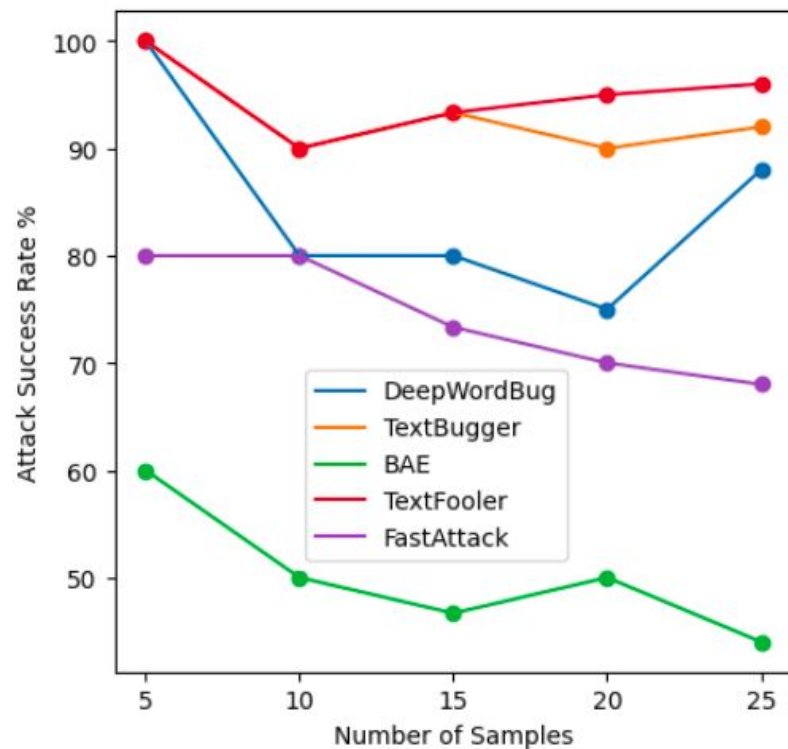
FastAttack vs BAE (Garg et.al.,2020)

Model	Dataset	Attack method	Successful attacks (Out of 25)	Original accuracy	Attacked accuracy	Success rate	Perturbed %	Avg num queries	Attack time (in seconds)
ALBERT-base-v2	YELP-polarity	fastattack	17	100%	32%	68%	8.19%	242.28	146.80
		BAE_garg	11		56%	44%	4.72%	287.24	438.28
	IMDB	fastattack	22	89.29%	10.71%	88.0%	5.2%	434.92	584.61
		BAE_garg	16		32.14%	64.0%	2.73%	512.32	1051.52
BERT-base-uncased	YELP-polarity	fastattack	21	100%	16%	84%	12.81%	238.84	252.68
		BAE_garg	11		56%	44%	6.44%	314.88	764.20
	IMDB	fastattack	21	92.59%	14.81%	84.0%	5.29%	419.72	268.29
		BAE_garg	12		48.15%	48.0%	2.97%	561.52	1153.58

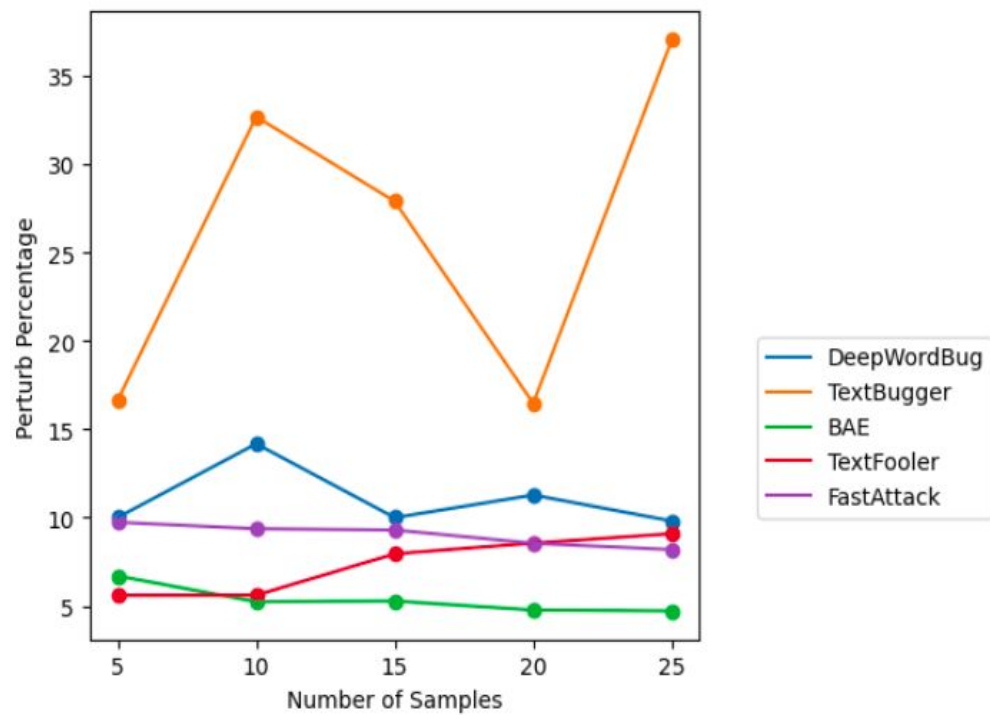
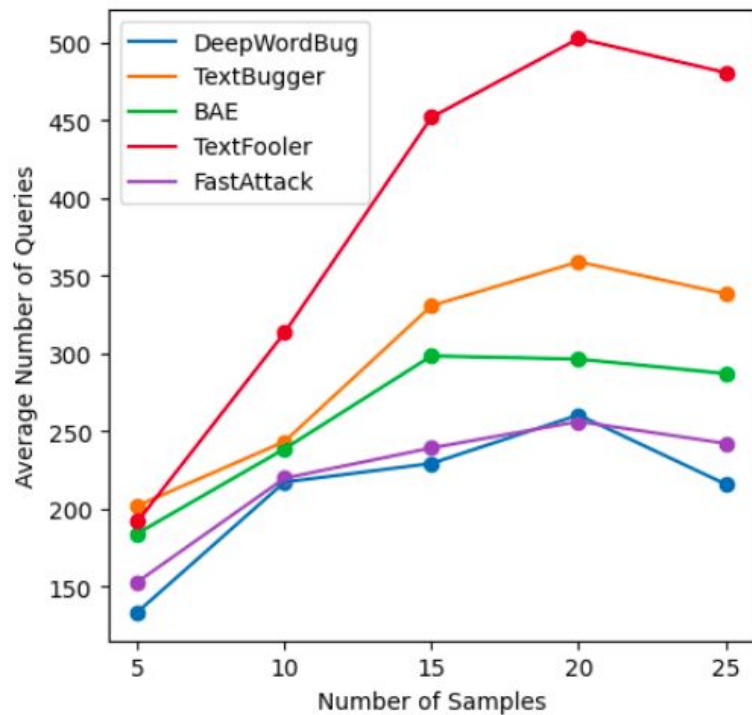
FastAttack vs BERT-Attack (Li et.al.,2020)

Model	Dataset	Attack method	Successful attacks (Out of 15)	Original accuracy	Attacked accuracy	Success rate	Perturbed %	Avg num queries	Attack time (in seconds)
ALBERT-base-v2	YELP-polarity	fastattack	14	78.95%	5.26%	93.33%	14.74%	<b>34.6</b>	<b>12.19</b>
		bert-attack	<b>15</b>		0.0%	<b>100%</b>	<b>12.52%</b>	110.8	123.67
	IMDB	fastattack	15	93.75%	0.0%	100%	4.45%	<b>352.0</b>	<b>253.93</b>
		bert-attack	15		0.0%	100%	<b>1.75%</b>	347.27	404.62
BERT-base-uncased	YELP-polarity	fastattack	14	78.95%	5.26%	93.33%	14.74%	<b>34.6</b>	<b>12.02</b>
		bert-attack	<b>15</b>		<b>0.0%</b>	<b>100%</b>	<b>12.52%</b>	110.8	121.36
	IMDB	fastattack	14	93.75%	6.25%	93.33%	5.34%	<b>365.33</b>	<b>238.02</b>
		bert-attack	<b>15</b>		<b>0.0%</b>	<b>100.0%</b>	<b>1.75%</b>	347.27	420.75

## More comparisons



# More comparisons



# Summarising Results

- FastText Embeddings are high-performance embeddings, outperforming Word2Vec and GloVE embeddings in most cases
- FastAttack is a simple and highly potent attack
  - No complex transformers, no language models
  - Still gives some of the latest attacks a run for their money
  - Is indeed “fast”
  - Replacement is done with valid words only
    - Replacements are always available, even for Out-of-Vocabulary words

## Member contributions

Roll No.	Member Name	Contribution
20111002	Abhishek Krishna	Trained Embeddings, Comparisons among Embeddings, FastAttack implementation, Comparisons with previous attack methods, Presentation
20111017	Deeksha Arora	Trained Embeddings, Comparisons among word Embeddings, Comparisons with previous attack methods, Report
20111044	Preeti Singh	Comparisons among Embeddings on DistilBERT-base, Report
20111054	Sambhrant Maurya	Trained Embeddings, Comparisons among embeddings, FastAttack implementation, Comparisons with previous attack methods, Presentation
20111061	Shruti Sharma	Trained Misspelling Oblivious Word Embeddings (MOE), Comparisons among Embeddings on RoBERTa-base, Report

# References

- [1] Ji Gao, Jack Lanchantin, Mary Lou Soffa, Yanjun Qi, “Black-box Generation of Adversarial Text Sequences to Evade Deep Learning Classifiers”, (2018)
- [2] J. Li, S. Ji, T. Du, B. Li, and T. Wang, “Textbugger: Generating adversarial text against real-world applications,” Proceedings 2019 Network and Distributed System Security Symposium, 2019.
- [3] D. Jin, Z. Jin, J. T. Zhou, and P. Szolovits, “Is bert really robust? a strong baseline for natural language attack on text classification and entailment,” 2020.
- [4] D. Pruthi, B. Dhingra, and Z. C. Lipton, “Combating adversarial misspellings with robust word recognition,” 2019.
- [5] L. Li, R. Ma, Q. Guo, X. Xue, and X. Qiu, “Bert-attack: Adversarial attack against bert using bert,” 2020.
- [6] S. Garg and G. Ramakrishnan, “Bae: Bert-based adversarial examples for text classification,” 2020.
- [7] V. Malik, A. Bhat, and A. Modi, “Adv-olm: Generating textual adversaries via olm,” 2021.
- [8] P. Bojanowski, E. Grave, A. Joulin, and T. Mikolov, “Enriching word vectors with subword information,” 2017.
- [9] S. Ren, Y. Deng, K. He, and W. Che, “Generating natural language adversarial examples through probability weighted word saliency,” 2019
- [10] J. Morris, E. Lifland, J. Y. Yoo, J. Grigsby, D. Jin, and Y. Qi, “TextAttack: A framework for adversarial attacks, data augmentation, and adversarial training in NLP,” 2020
- [11] Zang et.al. (2020), Word-level Textual Adversarial Attacking as Combinatorial Optimization