



UNIVERSITY
of VIRGINIA

TextAttack:

Generalizing Adversarial Examples to Natural Language Processing

@ Science Academy Machine Learning Summer School
2021/06/24

Dr. Yanjun Qi

<http://www.cs.virginia.edu/yanjun/>

Table of Contents



Background on Natural Language Processing (NLP)



Background on Adversarial Examples (AE) in Vision



Background on Adversarial Examples (AE) in NLP



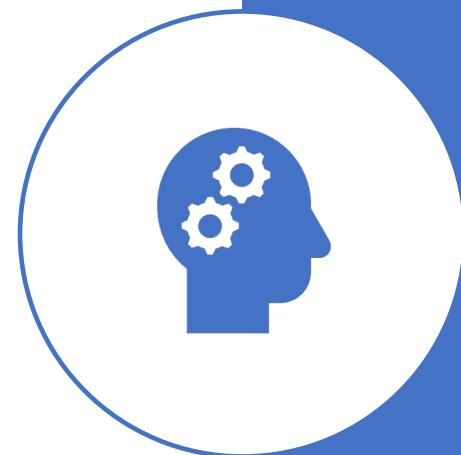
Textattack to rescue

Quality of the generated text

No standardized library

No clear insights what strategies work

No clear benefits whey studying AE NLP



Background:

Natural Language Processing and
Recent advances by Deep Learning

What is Natural language processing (NLP)

Wiki: is a field of computer science, artificial intelligence, and computational linguistics concerned with the interactions between computers and human (**natural**) languages.

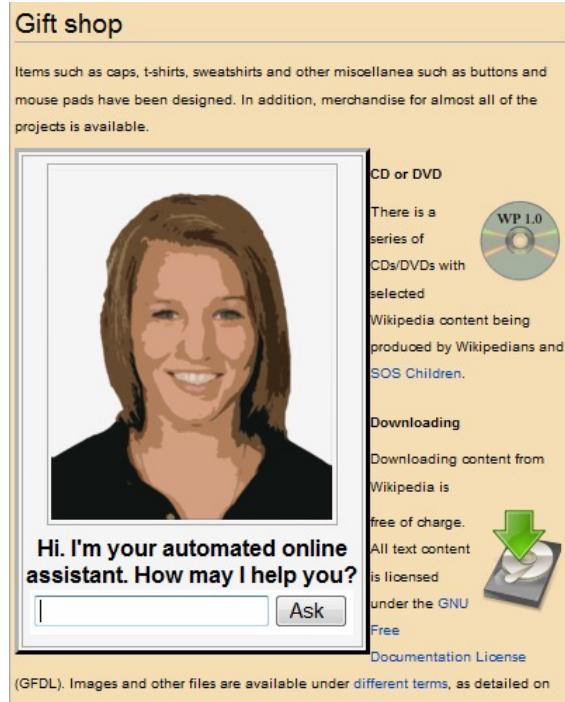


- Identify the **structure** and **meaning** of **words**, **sentences**, **texts** and **conversations**
- **Deep** understanding of **broad** language
- NLP is all around us

Machine translation

A screenshot of a Google search results page. The search query "buenas noches" is entered in the search bar. Below the search bar, the "All" tab is selected, along with other options like Images, Shopping, Apps, Videos, More, and Search tools. A message indicates "About 20,800,000 results (0.54 seconds)". The main result is a machine translation card. It shows "buenas noches" in Spanish on the left and "Goodnight" in English on the right. There are dropdown menus for both languages, microphone icons for voice input, and a speaker icon for audio output. Below the translation, there is a button labeled "3 more translations". At the bottom of the card, there is a link "Open in Google Translate".

Dialog Systems



Natural language instruction



Sentiment/Opinion Analysis

socialmention*

[Blogs](#) [Microblogs](#) [Bookmarks](#) [Images](#) [Video](#) [All](#)

starbucks

Search

[Advanced Search](#)
[Preferences](#)

20%
strength

2:1
sentiment

25%
passion

12%
reach

21 minutes avg. per mention

last mention 31 minutes ago

18 unique authors

0 retweets

Sentiment

positive		12
neutral		12
negative		5

Mentions about starbucks

Sort By:



Results:



Last Day

Results 1 - 15 of 29 mentions.

●  [pain.... all my homies know is pain](#)

submitted by /u/parkerblake204 to r/starbucks [link] [comments]

https://www.reddit.com/r/starbucks/comments/o6y66y/pain_all_my_homies_know_is_pain/
31 minutes ago - by /u/parkerblake204 on reddit

●  [What jobs provided tuition reimbursement ?](#)

I know Starbucks is one company that does this but I'm curious to find other jobs that also help you pay for college. submitted by /u/frozenfreddy7443 to r/college [...]

https://www.reddit.com/r/college/comments/o6y106/what_jobs_provided_tuition_reimbursement/
43 minutes ago - by /u/frozenfreddy7443 on reddit

●  [what is the proper build for a vanilla sweet cream cold brew?](#)

i used to do cold brew to the third line, ice, leave a little bit of room the the top, then sweet cream but one of the shifts told me that this was incorrect and the...

https://www.reddit.com/r/starbucks/comments/o6xr1h/what_is_the_proper_build_for_a_vanilla_sweet/
1 hour ago - by /u/chippyluvr on reddit

Question answering



[IBM 'Watson' computer wins at 'Jeopardy'](#)

Text Classification



BIDNESS ETC

Did you mean to attach files?
You wrote "is attached" in your message, but there are no files attached. Send anyway?

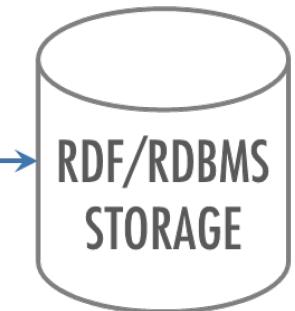
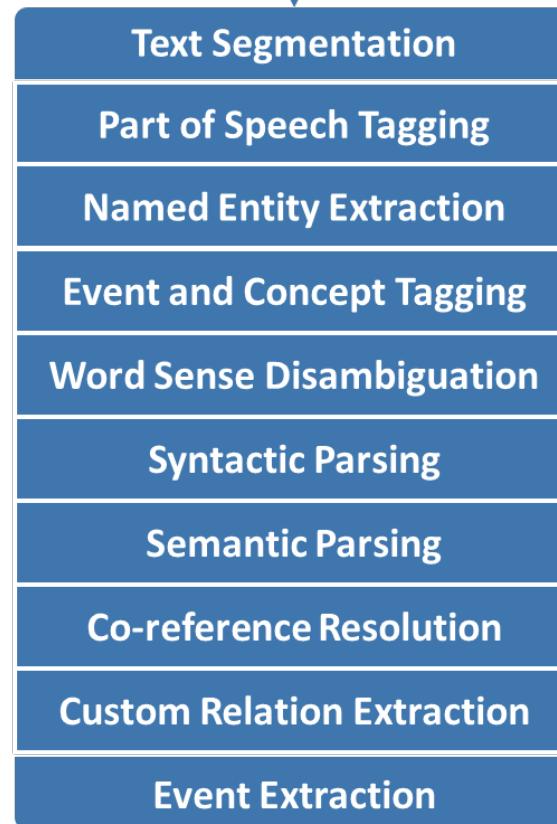
OK Cancel

1–21 of 21 < > Settings

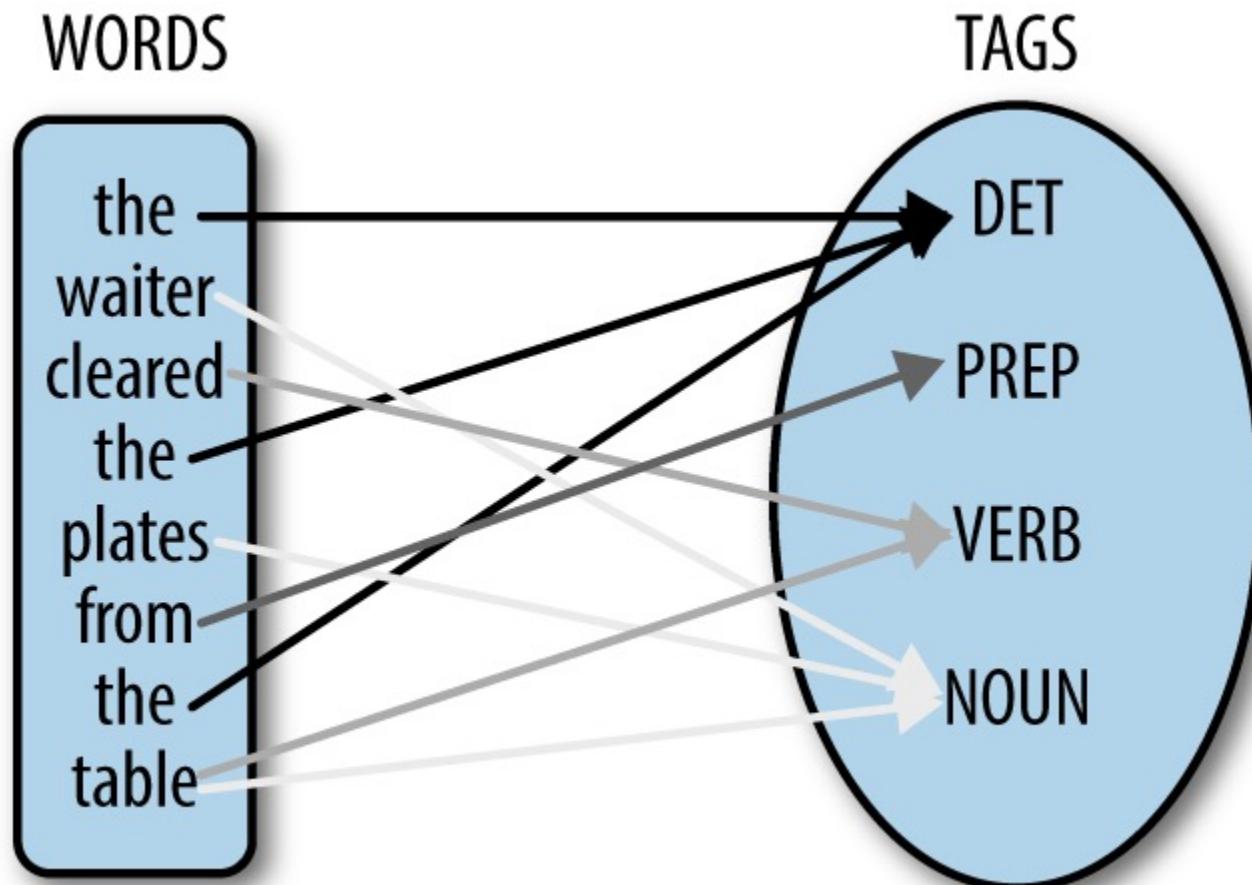
Primary	Social Google+ 1 new	Promotions Google Offers, Zagat 2 new	Updates Google Play 1 new
James, me (2)	Hiking Hiking trip on Saturday - Yay - so glad you can join. We should leave from I	3:14 pm	
Hannah Cho	Thank you - Keri - so good that you and Steve were able to come over. Thank you :	3:05 pm	
Jay Birdsong	School Upcoming school conference dates Hello everyone. A few people have	www.wired.com	



Classic NLP Pipeline
Includes a set of
Components for
Understanding Text



Part of speech tagging



Q: [Chris] = [Mr. Robin] ?

Christopher Robin is alive and well. He is the same person that you read about in the book, **Winnie the Pooh**. As a boy, **Chris** lived in a pretty home called **Cotchfield Farm**. When **Chris** was three years old, **his father** wrote a poem about **him**. The poem was printed in a magazine for others to read. **Mr. Robin** then wrote a book

Information Extraction

- Unstructured text to database entries

New York Times Co. named Russell T. Lewis, 45, president and general manager of its flagship New York Times newspaper, responsible for all business-side activities. He was executive vice president and deputy general manager. He succeeds Lance R. Primis, who in September was named president and chief operating officer of the parent.

Person	Company	Post	State
Russell T. Lewis	New York Times newspaper	president and general manager	start
Russell T. Lewis	New York Times newspaper	executive vice president	end
Lance R. Primis	New York Times Co.	president and CEO	start

Recent deep learning advances on NLP

Pre-2012

2013-16

2014-16

2016-NOW

Supervised predictors for each component

- BOW
- LSI
- Topic LDA
- ...

Word2Vec

- GloVe
- FastText)

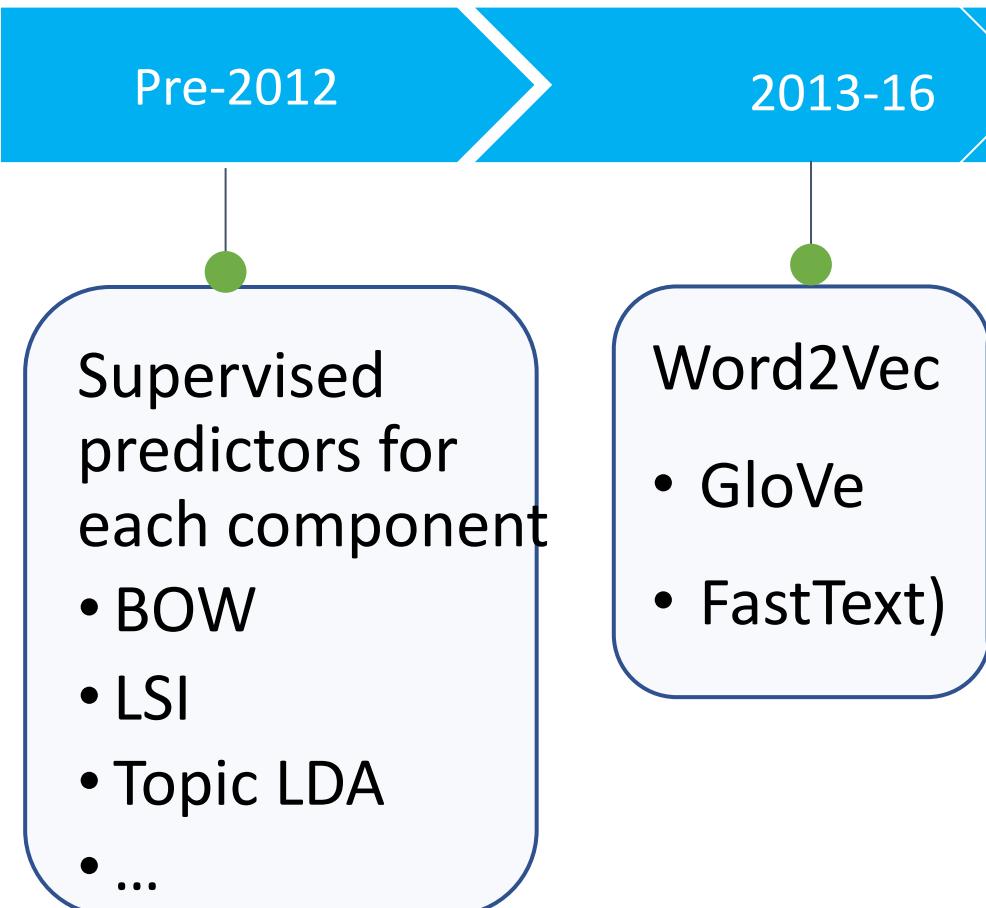
Recurrent

- LSTM
- Seq2Seq

- Attention
- Transformer (self-attention, attention only)

- BERT
- XLNet
- GPT-2 /3
- T5 ...

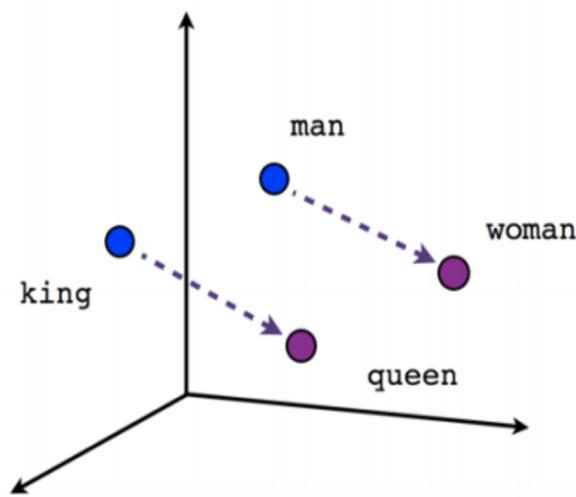
Recent deep learning advances on natural language



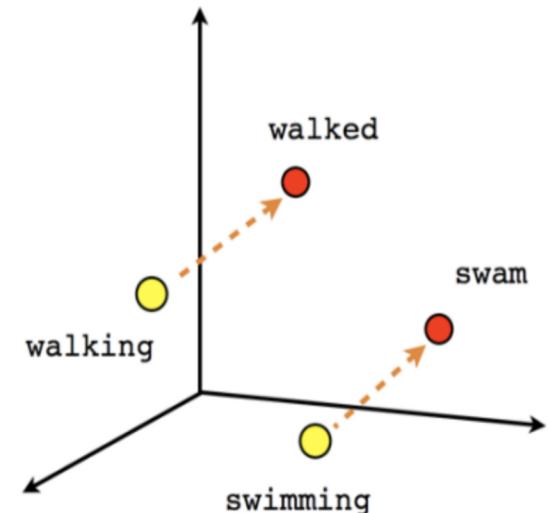
How to Represent A Word in DNN

- Basic approach – “one hot vector”

- Binary vector
- Length = | vocab |
- 1 in the position of the word id, the rest are 0
- However, does not represent word meaning
- Extremely high dimensional (there are over 200K words in the English language)
- Extremely sparse



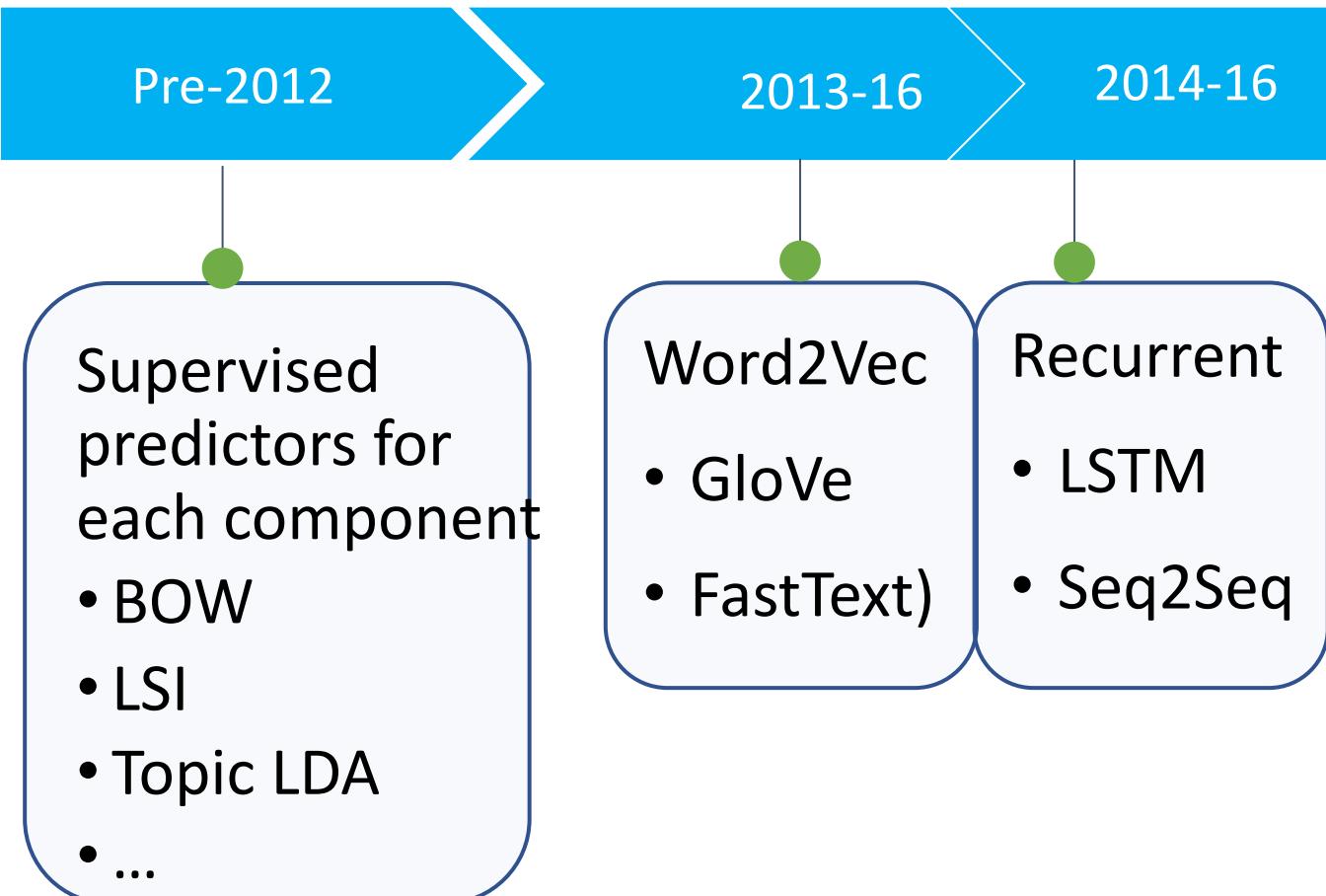
Male-Female



Verb tense

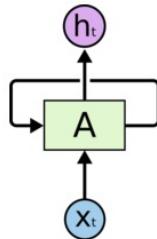
- Solution:
Distributional Word Embedding Vectors

Recent deep learning advances on natural language

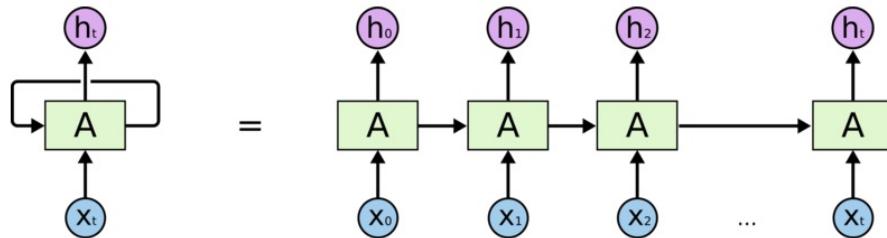


Recurrent Neural Networks

- Allow us to operate over sequences of vectors (**with variable length**)
- Allow Sequences in the input, as the output, or in the most general case both



Recurrent Neural Networks have loops.

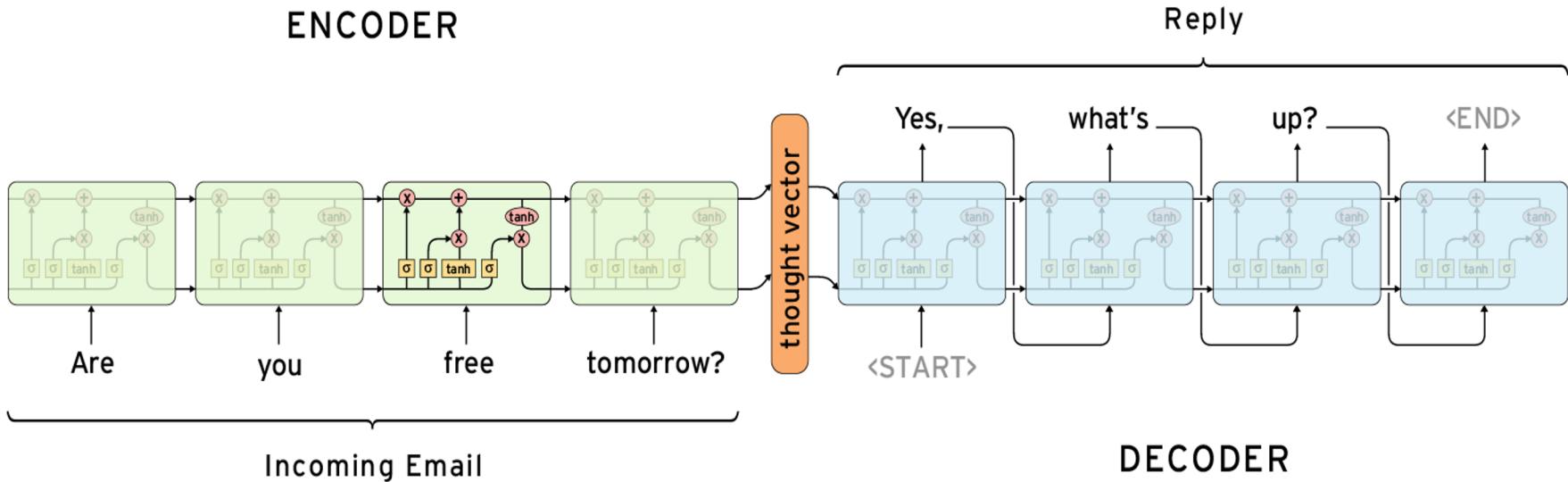
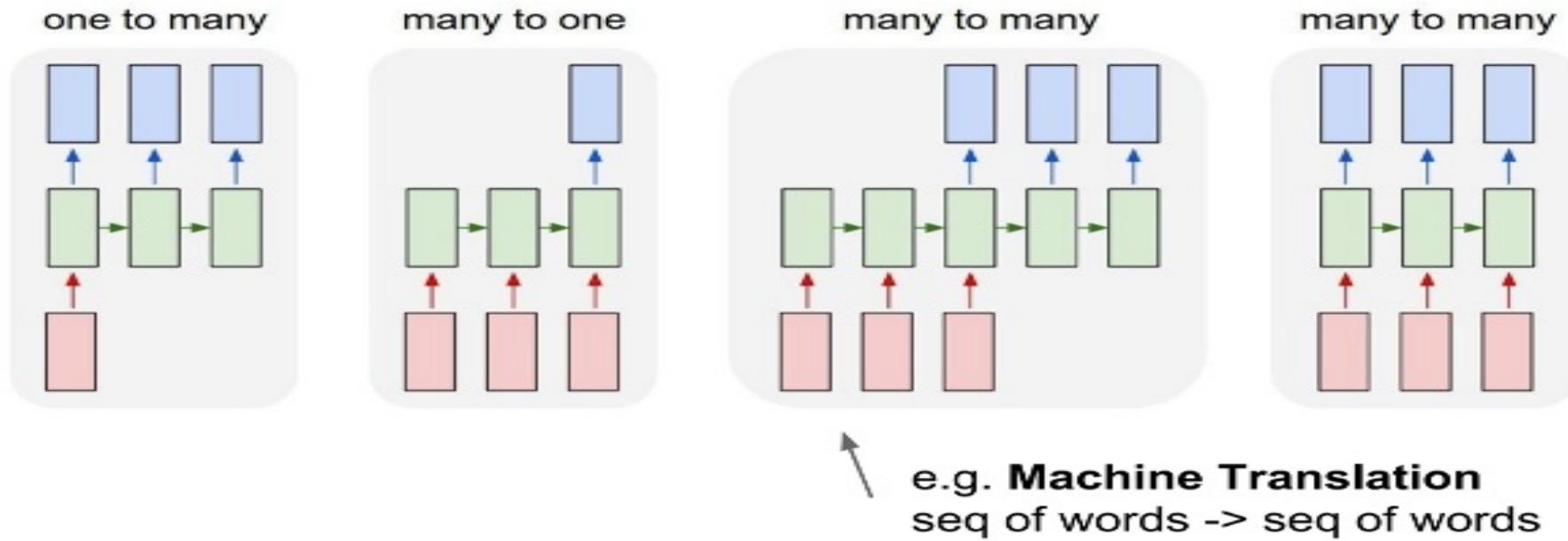


An unrolled recurrent neural network.

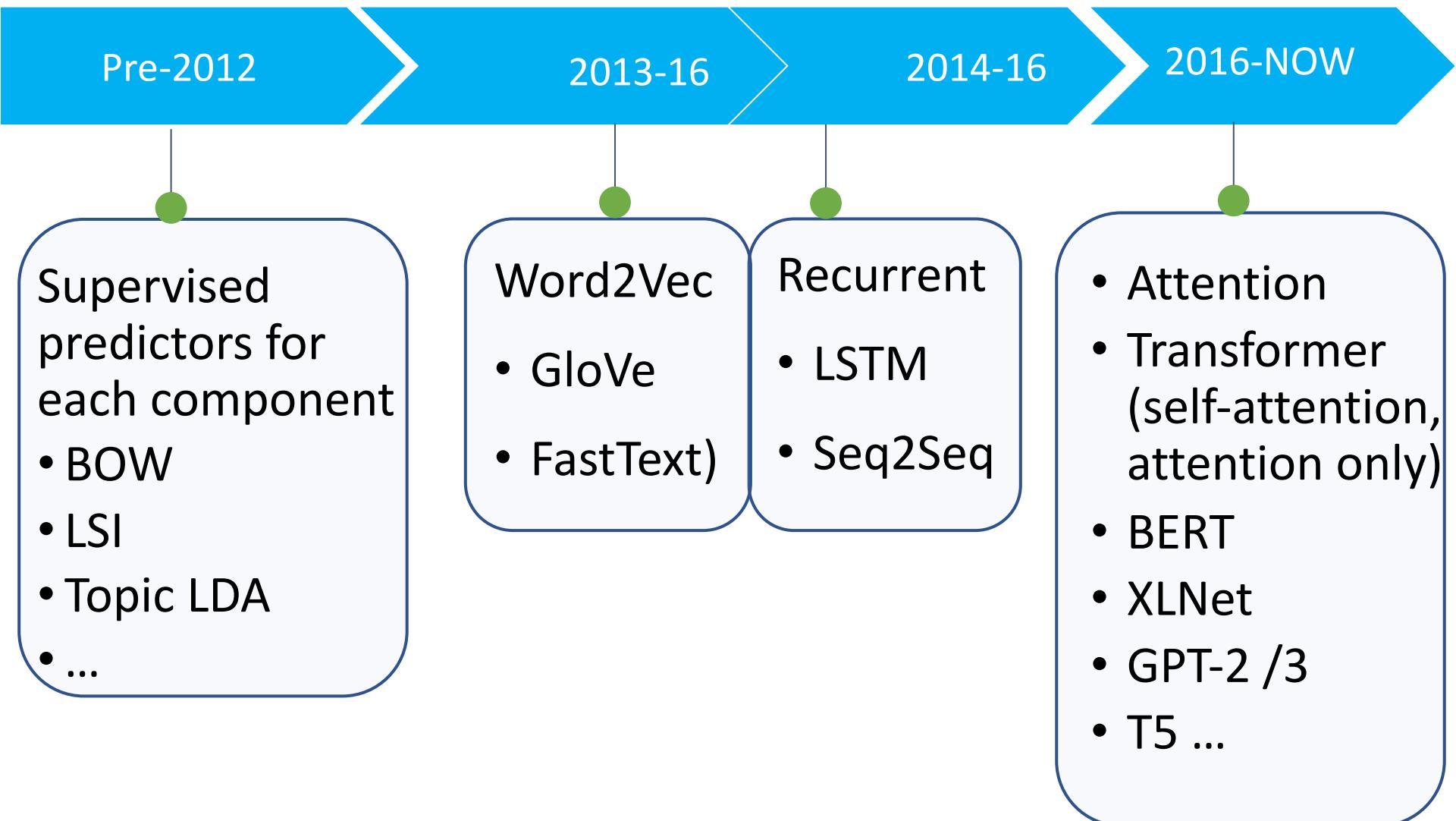
Recurrent Neural Networks are networks with loops in them, allowing information to persist.

Image Credits from Christopher Olah

Recurrent Neural Networks (RNNs) can handle



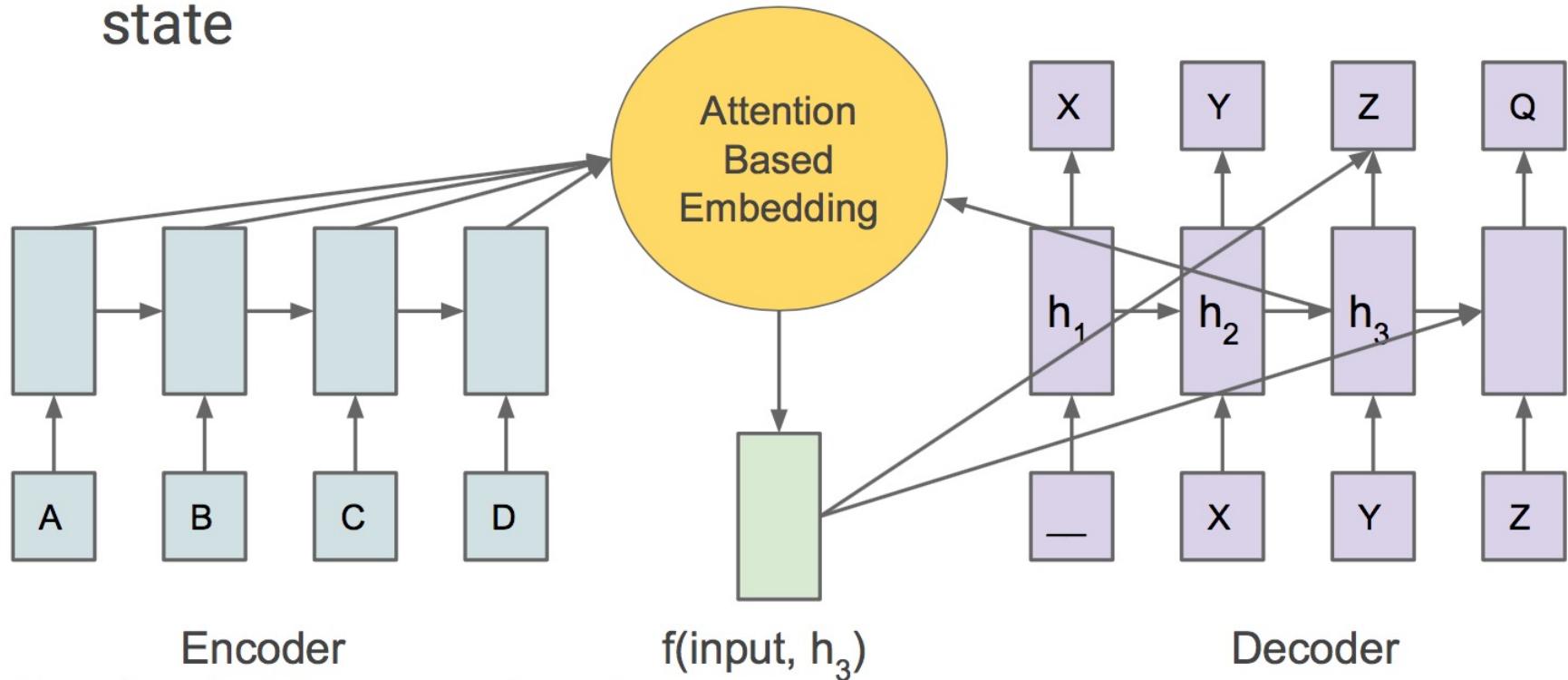
Recent deep learning advances on natural language



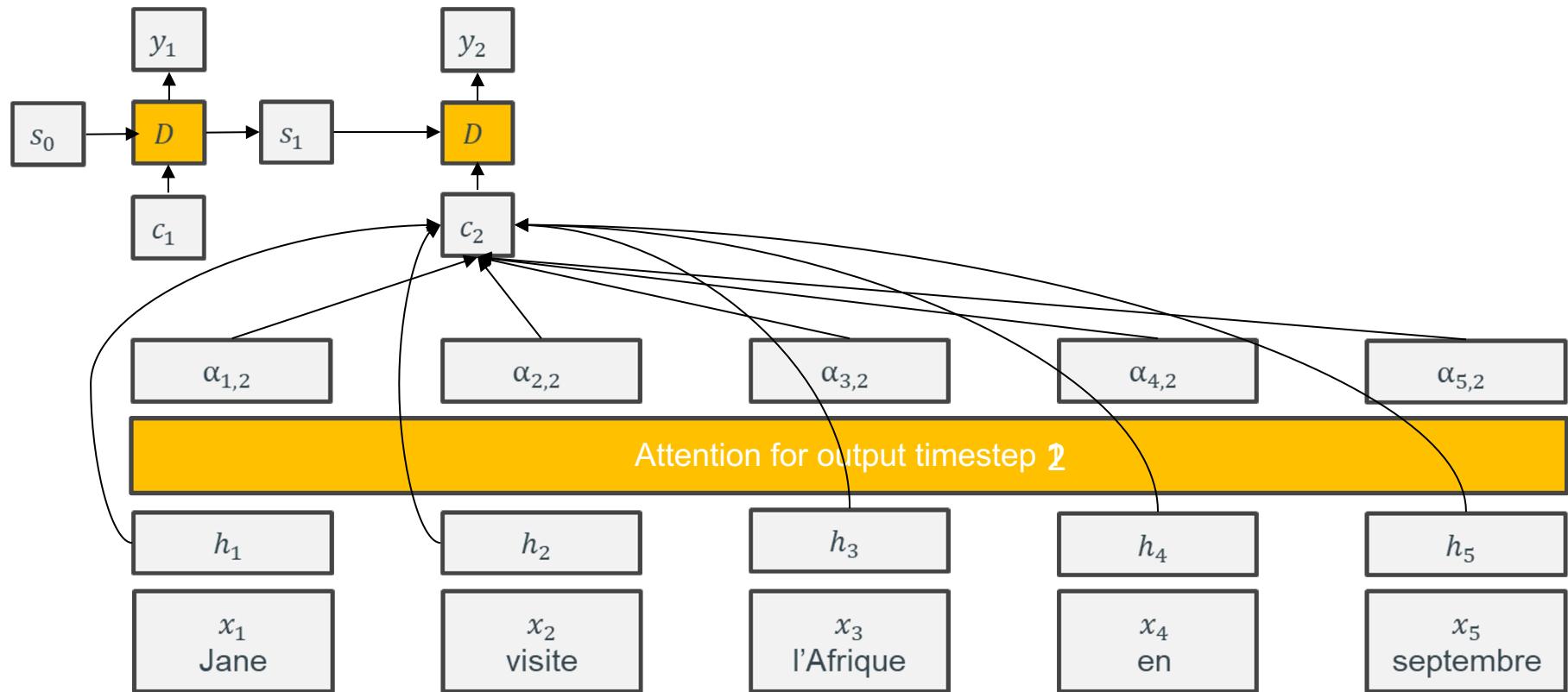
Attention Trick:

Seq2Seq with Attention

- Embedding used to predict output, and compute next hidden state



The attention module gives us a weight for each input.



Self-attention creates attention layers mapping from a sequence to itself.

The FBI is chasing a criminal on the run .

The FBI is chasing a criminal on the run .

The FBI is chasing a criminal on the run .

The FBI is chasing a criminal on the run .

The FBI is chasing a criminal on the run .

The FBI is chasing a criminal on the run .

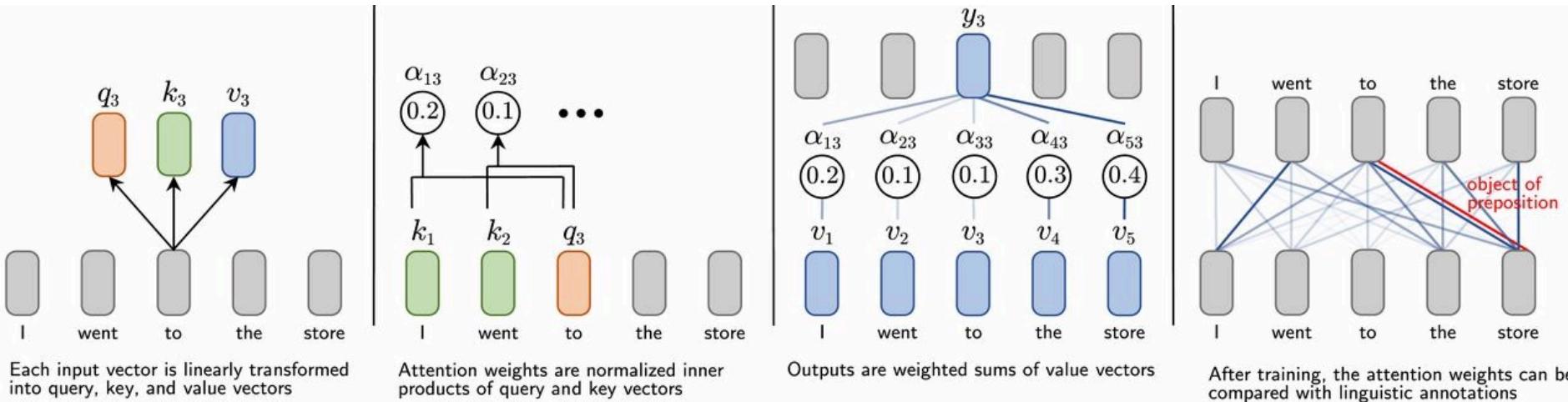
The FBI is chasing a criminal on the run .

The FBI is chasing a criminal on the run .

The FBI is chasing a criminal on the run .

The FBI is chasing a criminal on the run .

The FBI is chasing a criminal on the run .



Transformer: Exploiting Self Attentions

- Uses 3 kinds of attention
 - Encoder self-attention.
 - Decoder self-attention.
 - Encoder-decoder multi-head attention.

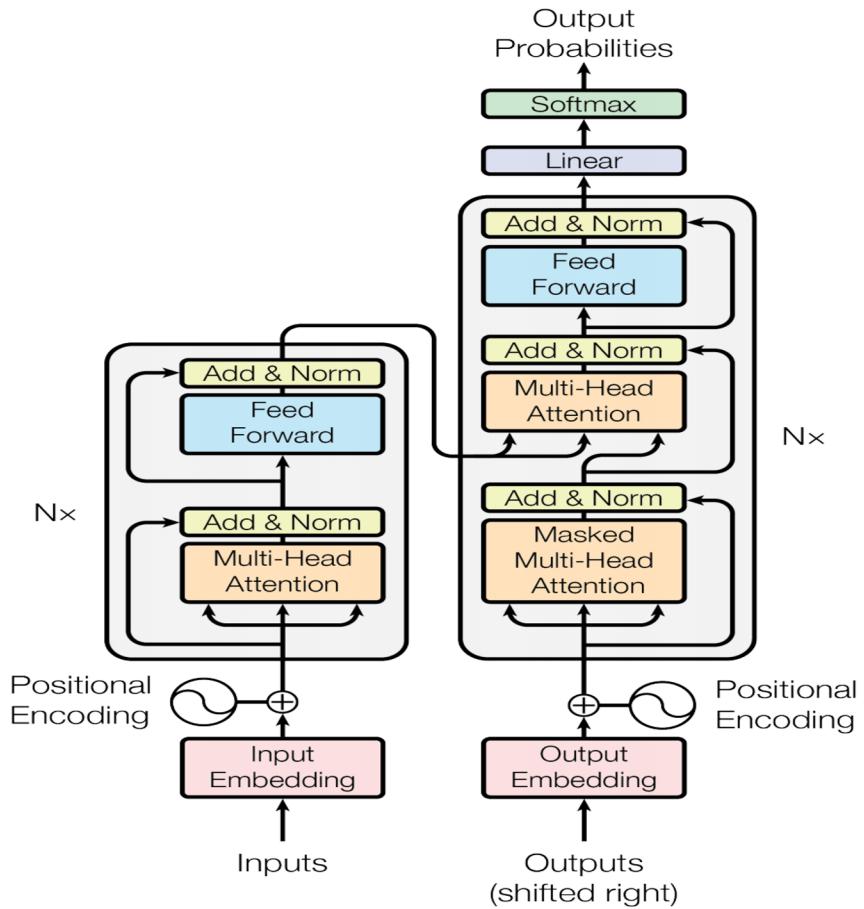


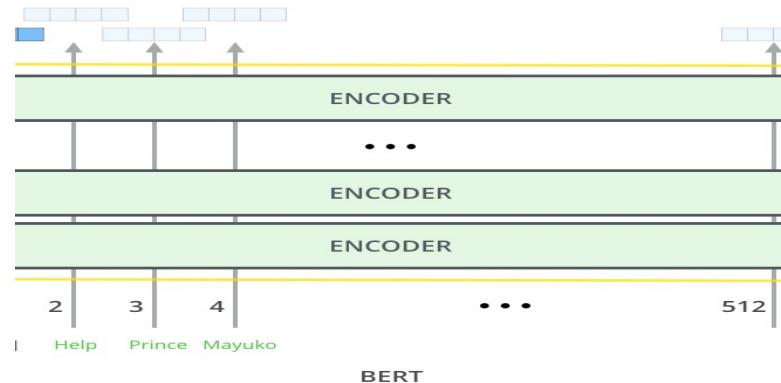
Figure 1: The Transformer - model architecture.



BERT: Bidirectional Encoder Representations from Transformers
Pre-trained transformer encoder for sentence embedding

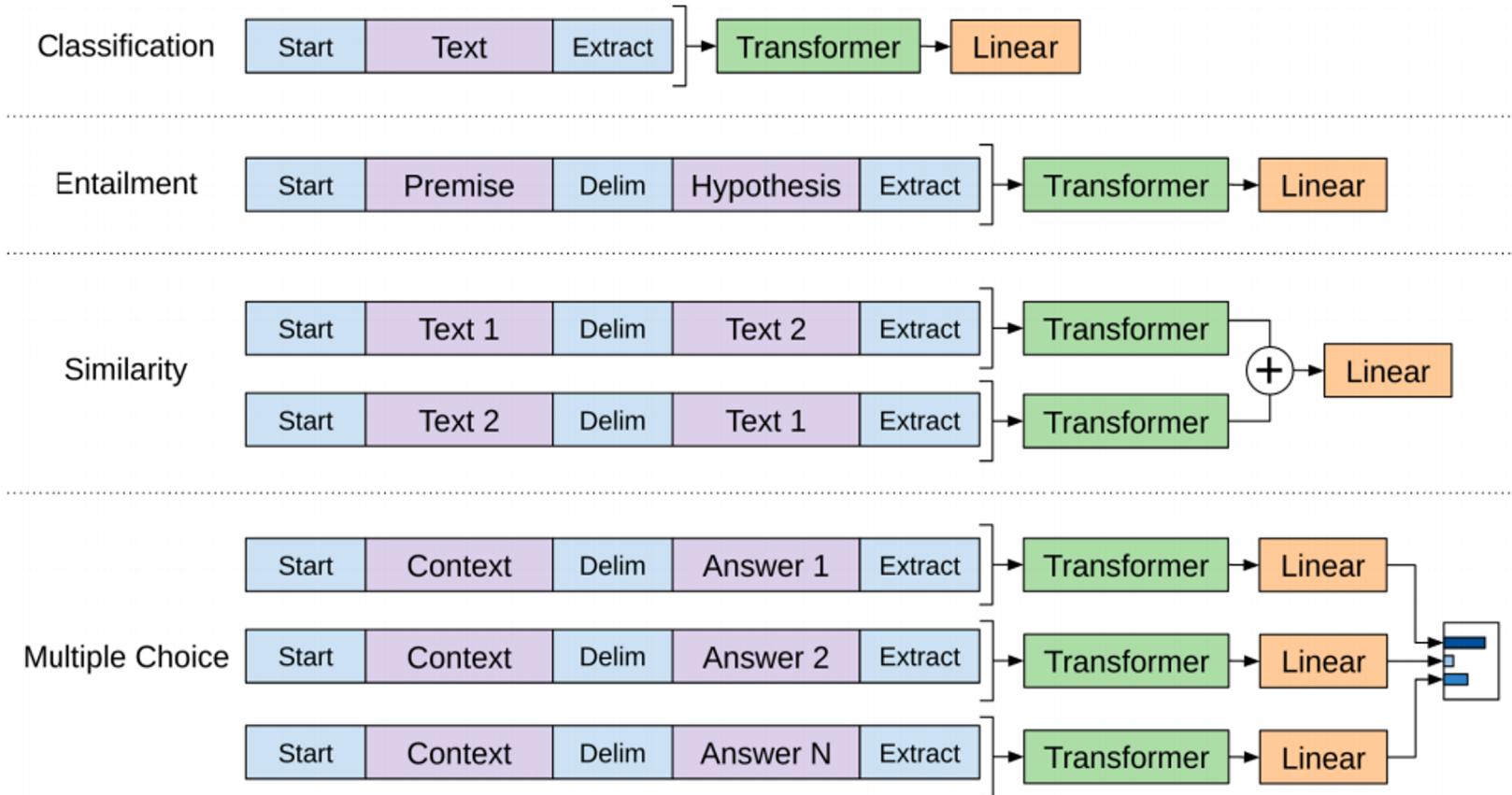


Notable pre-trained NLP models



BERT's architecture is just a transformer's encoder stack.

As with BERT, you can use the pretrained GPT models for any task.
Different tasks use the OpenAI transformer in different ways.



GPT: generative pre-training,

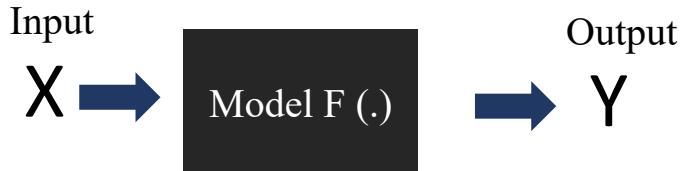
GPT's architecture is just a transformer's decoder stack.



Background: Adversarial Examples in Vision

Background: Machine Learning

- Machine Learning: to find **model $F(\cdot)$** that can **generalize** from observed data to **unseen** data



Model $F(\cdot)$ generalizes to Unseen X'

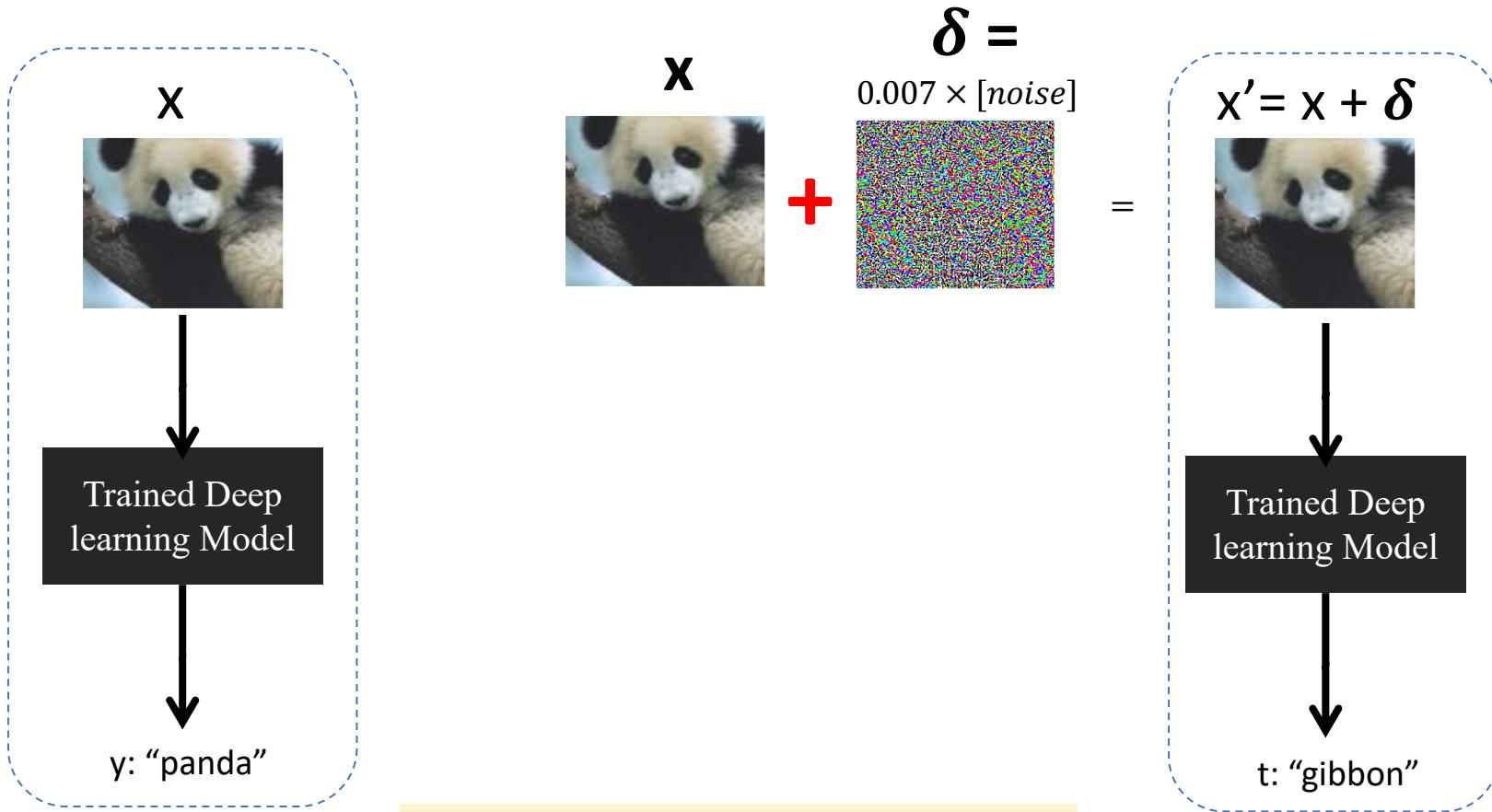
For instance:



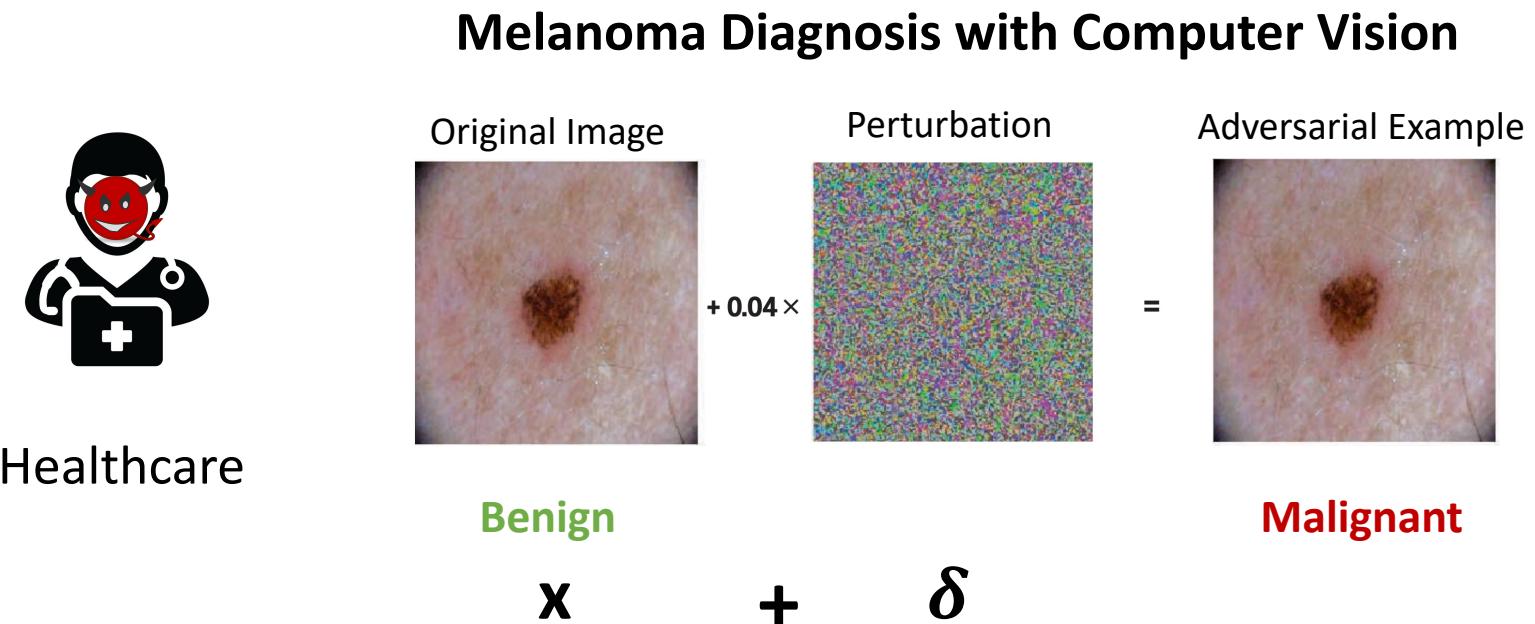
Trained Deep learning Model

“panda”

Background: Adversarial Examples

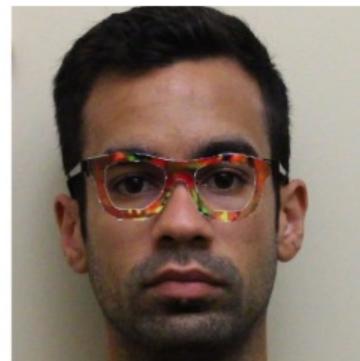


Deep Learning Classifiers are Easily Fooled



Samuel G Finlayson et al. "Adversarial attacks on medical machine learning", *Science*, 2019.

Classifiers Under Attack: Adversary Adapts



Accessorize to a Crime: Real and Stealthy Attacks on State-of-the-Art Face Recognition

Mahmood Sharif
Carnegie Mellon University
Pittsburgh, PA, USA
mahmoods@cmu.edu

Sruti Bhagavatula
Carnegie Mellon University
Pittsburgh, PA, USA
srutib@cmu.edu

Michael K. Reiter
University of North Carolina
Chapel Hill, NC, USA
reiter@cs.unc.edu

Luo Bauer
Carnegie Mellon University
Pittsburgh, PA, USA
lbauer@cmu.edu

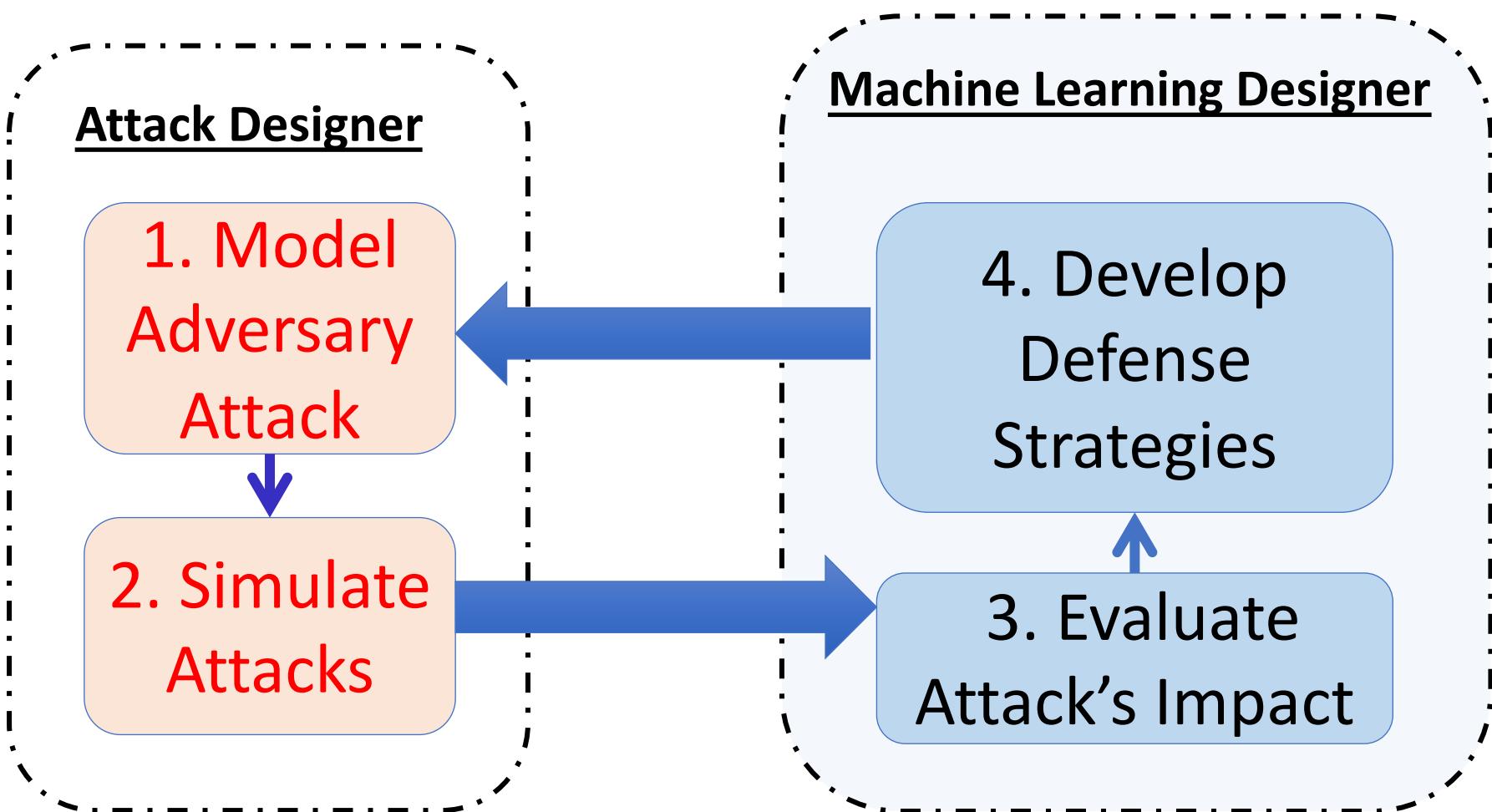
ACM CCS 2016

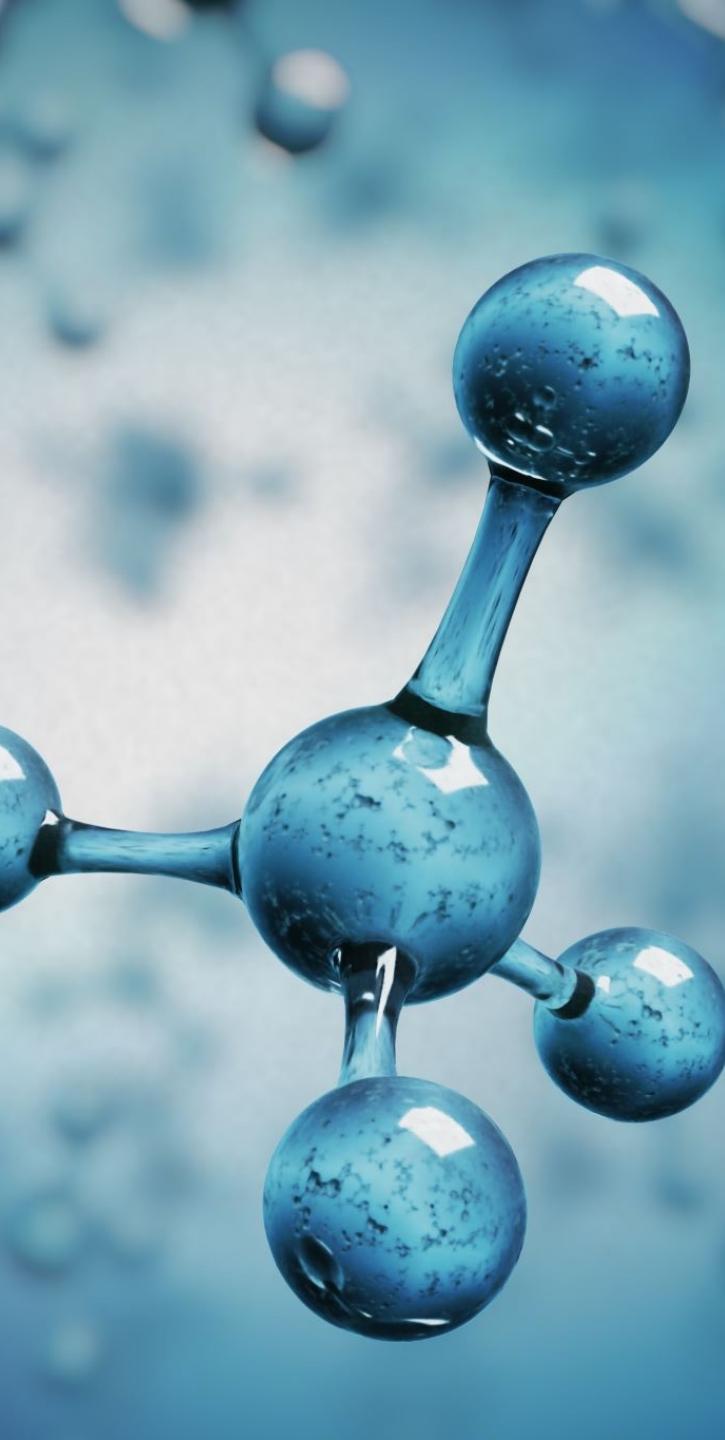
Actual images

Recognized faces

Mahmood Sharif et al. "Accessorize to a Crime: Real and Stealthy Attacks on State-of-the-Art Face Recognition", In CCS, 2016.

Goal of Adversarial Machine Learning

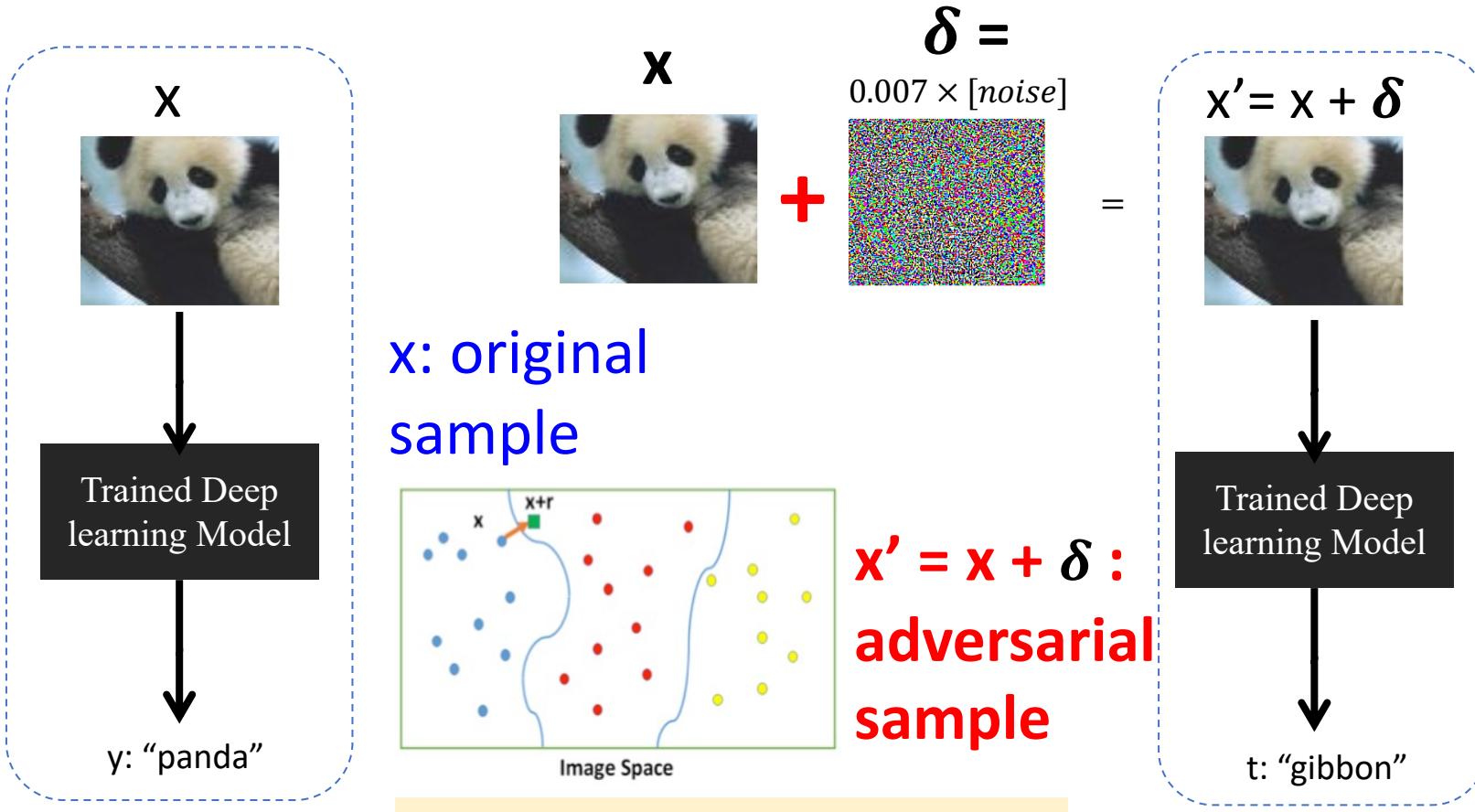




Terminology

- **Adversarial examples** or **adversarial perturbations**: changes to inputs that fool a trained model;
- We call “a program that repeatedly generates adversarial examples for a target model” as an **adversarial attack**
- We name a model’s resistance to adversarial examples as **adversarial robustness**

Background: Adversarial Examples



How to find δ ? Definition in Vision

Let F be our neural network represented as a function. To generate adversarial example for given image $x \in \mathbb{R}^m$ and its label $y \in \{1, \dots, K\}$, we want to find some perturbation $\delta \in \mathbb{R}^m$ such that

$$F(x + \delta) \neq y \quad \text{Misclassification term}$$

At the same time, we want $\|\delta\|_p \leq \epsilon$ where $\epsilon > 0$ and p is typically 2 or ∞ .

Distance term

Background: Different variations of Adversarial Examples

$$\begin{array}{c} \text{Original Example} \\ \text{“1”} \\ 100\% \text{ confidence} \end{array} + \begin{array}{c} \text{BIM} \\ \text{Perturbations} \\ \text{JSMA} \end{array} = \begin{array}{c} \text{“4”} \\ 100\% \\ \text{“2”} \\ 99.9\% \\ \text{“2”} \\ 83.8\% \end{array}$$

$\mathbf{x} + \boldsymbol{\delta}$

Adversarial Examples

36

How to search for δ ?

Adversarial Attack as Optimization Problem

Generating an adversarial example therefore becomes an optimization problem where

(3) How to define $L(\cdot)$?

$$\text{maximize } L(F, x + \delta, y)$$

(1) How to define δ ?

$$\text{subject to } \|\delta\|_p \leq \epsilon$$

(4) How to
optimize this?

(2) How to limit δ ?

Many different variations of formulations to search for δ ?

(3) How to define $L(.)$?

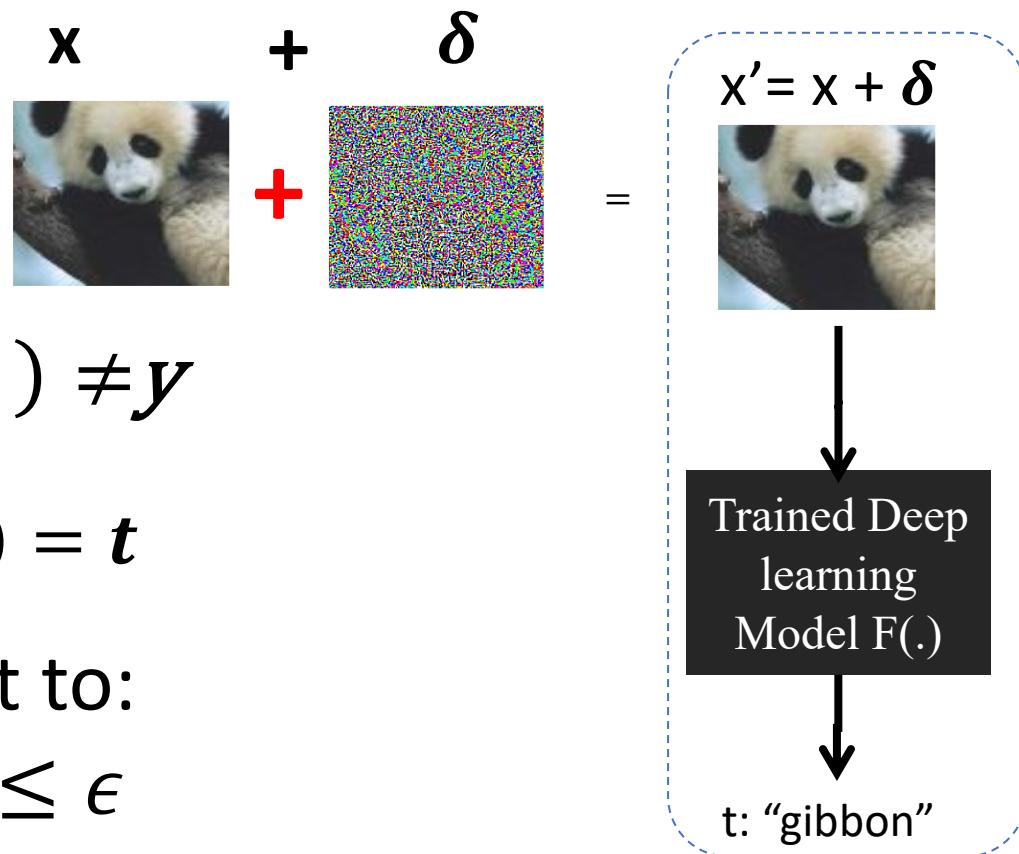
untargeted $F(x + \delta) \neq y$

targeted $F(x + \delta) = t$

Subject to:

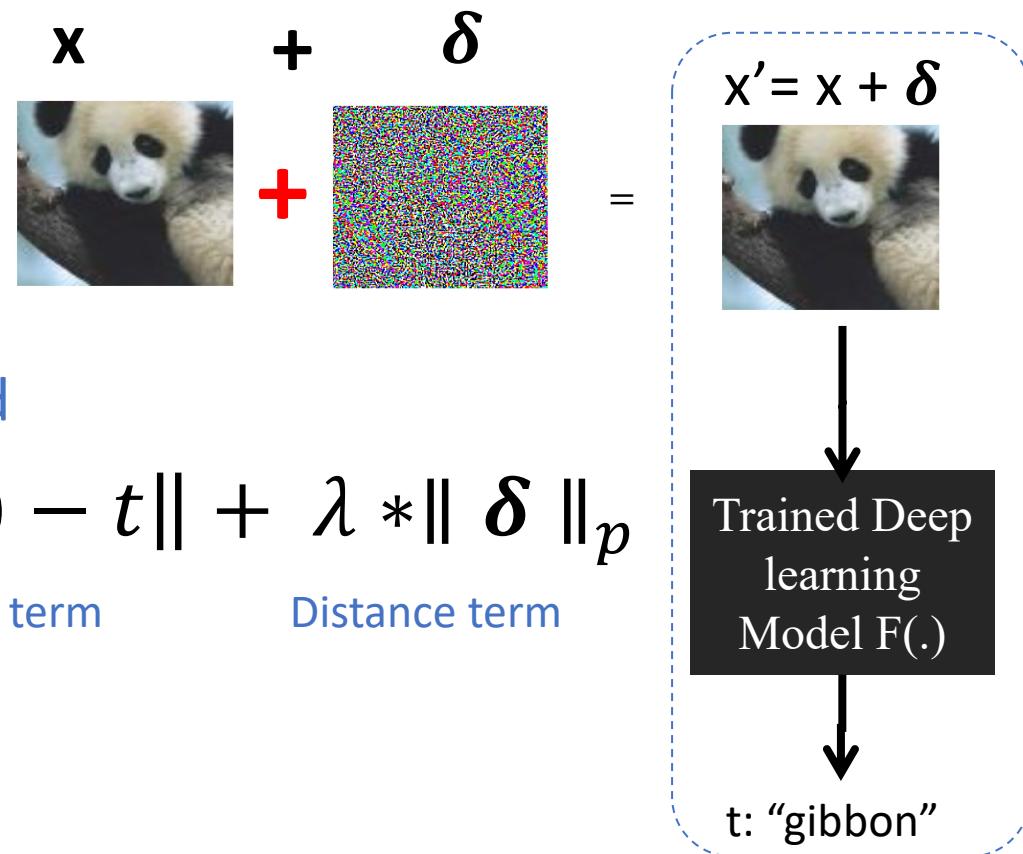
$$\|\delta\|_p \leq \epsilon$$

$$\epsilon > 0$$



Many different variations of formulations to search for x' from x

(3) How to define $L(\cdot)$?



Popular Attacks in Vision: FGSM

Goodfellow et al. (2014) proposed a fast approximation algorithm called Fast Gradient Sign Method (FGSM).

Perturbation δ was found by:

(4) How to optimize to get best δ ?
$$\delta = \varepsilon * \text{sign}(\nabla L(F, x, y))$$

Popular Attacks in Vision: PGD

Projected gradient descent (PGD) attack, proposed by Madry et al. (2017), can be viewed as an iterative version of FGSM.

Adversarial example is found by:

(4) How to optimize to get best δ ?

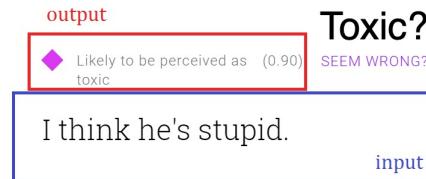
$$x^{t+1} = \text{Clip}_\epsilon(x^t + \alpha * \text{sign}(\nabla L(F, x^t, y)))$$

Background:

Adversarial Examples in NLP



Toxicity Identification



Sentiment Classification

"I love this movie.
I've seen it many times
and it's still awesome."



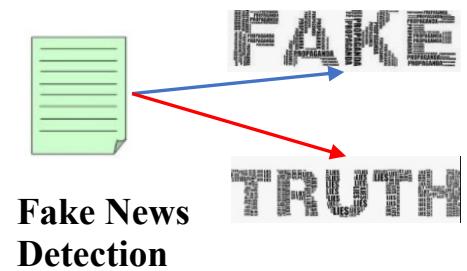
"This movie is bad.
I don't like it at all.
It's terrible."

Authorship Detection

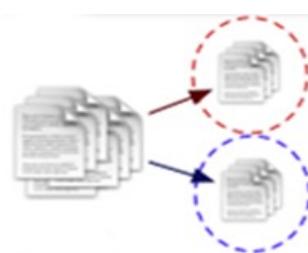


Rowling?

NLP Computer System needs Trustworthiness and Robustness



Spam Detection



Electronic Medical Records

What about Adversarial Examples in NLP?

Naturally, we are interested if we can borrow the previous vision formulation to NLP.

But, we face some difficulties.

Images are continuous while text is discrete. This leads to significant difference in how adversarial examples are generated.

What about AE in NLP?

Adversarial Attack as Optimization Problem

- Images are continuous while text is discrete. This leads to significant difference in how adversarial examples are generated.

Generating an adversarial example therefore becomes an optimization problem where

(3) How to define $L(\cdot)$?

$$\text{maximize } L(F, x + \delta, y)$$

(1) How to define δ ?

$$\text{subject to } \|\delta\|_p \leq \epsilon$$

(4) How to
optimize this?

(2) How to limit δ ?

What about AE in NLP?

- Images are continuous while text is discrete. This leads to significant difference in how adversarial examples are generated.

Generating an adversarial example therefore becomes an optimization problem where

$$\begin{aligned} & \text{maximize} && L(F, \boxed{x + \delta}, y) \\ & \text{subject to} && \|\delta\|_p \leq \epsilon \end{aligned}$$

(1) How do you define δ for text?

What about AE in NLP?

- Images are continuous while text is discrete. This leads to significant difference in how adversarial examples are generated.

Generating an adversarial example therefore becomes an optimization problem where

$$\begin{aligned} & \text{maximize} && L(F, x + \delta, y) \\ & \text{subject to} && \|\delta\|_p \leq \epsilon \end{aligned}$$

(2) How do you limit how much to change the text?

What about AE in NLP?

Adversarial Attack as Optimization Problem

- Images are continuous while text is discrete. This leads to significant difference in how adversarial examples are generated.

Generating an adversarial example therefore becomes an optimization problem where

$$\begin{aligned} \text{(3) How to define } L(\cdot)? & \quad \text{maximize} \quad L(F, x + \delta, y) \\ & \quad \text{subject to} \quad \|\delta\|_p \leq \epsilon \end{aligned}$$

What about AE in NLP?

Adversarial Attack as Optimization Problem

- Images are continuous while text is discrete. This leads to significant difference in how adversarial examples are generated.

Generating an adversarial example therefore becomes an optimization problem where

$$\begin{aligned} & \text{maximize} && L(F, x + \delta, y) \\ & \text{subject to} && \|\delta\|_p \leq \epsilon \end{aligned}$$

(4) How to
optimize this?

(1) How do you define δ for NLP?

=> Four main types of perturbations

1. **Character substitution:** add, remove, or modify characters until the prediction changes.
2. **Word insertion or removal:** add or remove words until prediction changes.
3. **Paraphrase:** train a model to paraphrase sentences; iteratively paraphrase it until prediction score changes.
4. **Synonym substitution:** swap out words in the input for a direct substitution until prediction changes.

Most successful technique (so far)

$$T(x): \cancel{x + \delta}$$

Transformation term

(1) How do you define δ for NLP?

- Idea 1: examples that are *almost* visually indistinguishable to humans (**mispellings**)

Input, x :

“True Grit” was the best **movie**
I’ve seen **since** I was a small boy.

Prediction: **Positive ✓**

Perturbation, x_{adv} :

“True Grit” was the best **moive**
I’ve seen **snice** I was a small boy.

Prediction: **Negative X**

Useful, but easy to defend against:

- Pass inputs into a **spell-checker** before feeding them into the model
- Or, train an RNN to correct inputs

(1) How do you define δ for NLP?

- **Idea 4:** examples that are indistinguishable *in meaning* to the original input (semantics-preserving changes)

Input, x :

“True Grit” was the best movie
I’ve seen since I was a **small boy**.

Prediction: **Positive ✓**

Perturbation, x_{adv} :

“True Grit” was the best movie
I’ve seen since I was a **wee lad**.

Prediction: **Negative X**



Most successful technique (so far)

AE NLP
literature
is messy
(chaotic)

1. Many generate examples are bad

2. No standard library

3. No clear benchmarking insights (which strategy?)

4. No clear benefits

Our Solution:
TextAttack to Rescue

1. Many generate
examples are bad

AE NLP
literature
is messy
(chaotic)

(2) How do you limit how much to change from the seed text?

- Images are continuous while text is discrete. This leads to significant difference in how adversarial examples are generated.

Generating an adversarial example therefore becomes an optimization problem where

$$\begin{aligned} & \text{maximize} && L(F, x + \delta, y) \\ & \text{subject to} && \|\delta\|_p \leq \epsilon \end{aligned}$$

(2) How do you limit how much to change the text?

Input, x :

“True Grit” was the best movie I’ve seen since I was a **small boy**.

Prediction: **Positive ✓**

Perturbation, x_{adv} :

“True Grit” was the best movie I’ve seen since I was a **wee lad**.

Prediction: **Negative X**

Bad examples of adversarial perturbations in NLP

Perturbation, x_{adv} :

different **semantics** than original input



“True Grit” was the **worst** movie I’ve seen since I was a small boy.

violates **grammar** (unlike the original input)



“True Grit” was the best movie I’ve seen since I **were boy small**.

this is just **suspicious** – nobody talks like that!



“True Grit” was the best movie I’ve seen since I was a **minuscule youngster**.

Constraints to ensure “valid” examples

→ How to measure δ for NLP?

- **Idea 1:** what is the cosine similarity between the sentence embeddings of \mathbf{x} and \mathbf{x}_{adv} ?
 - (we can obtain sentence embeddings from the Universal Sentence Encoder, for example)
- **Idea 2:** Use a grammar checker to sure that we didn't introduce any grammatical errors in \mathbf{x}_{adv} .

Evaluation of grammar



LanguageTool
proofreading software

- We evaluated syntax with **LanguageTool**, an open-source grammar checker
- Detected more errors in x_{adv} than x in 35% to 70% of samples (depending on datasets)

	(Jin et al., 2019)			(Alzantot et al., 2018)	
	IMDB	Yelp	MR	IMDB	Yelp
% with errors	61.8%	71.6%	35.5%	42.7%	43.8%

Let $T(x)$ be perturbation and $C_i(x)$ be a constraint,

$$C_1(T(x)) \wedge C_2(T(x)) \wedge \dots \wedge C_m(T(x))\}$$

We propose a taxonomy of Constraints to control

Constraint	Human Evaluation Method	Automatic Evaluation Proxy
Semantic	Ask humans whether meaning is preserved	Universal Sentence Encoder
Grammatical	Ask humans to find grammatical errors in shuffled mix	LanguageTool
Character Overlap	Ask humans how similar x, x_{adv} look	Edit distance, BLEU, METEOR
Non-suspicious	Ask humans to identify suspicious in shuffled mix	GROVER

Our Analysis paper: Reevaluating Adversarial Examples in Natural Language

- 2020 [EMNLP Findings](#)

Our evaluation reveals two concerns:

1. Literature's comparisons between past attacks are problematic. What is really necessary is comparison ***with the same constraints***
2. Even once constraints are standardized, researchers chose too lax thresholds! → We ***actually asked humans*** (via Amazon Turk) to provide guidance on the best threshold for each constraint.

Human Study Standardized Constraints Enables Better/ Truthful Comparisons

Constraints Search Method	TFADJUSTED		TEXTFOOLER	
	TEXTFOOLER	GENETICATTACK	TEXTFOOLER	GENETICATTACK
Semantic Preservation	4.06	4.11	-	-
Grammatical Error %	0	0	-	-
Non-suspicion Score	58.8	56.9	-	-
Attack Success %	10.6	12.0	91.1	95.0
Perturbed Word %	11.1	11.0	18.9	17.2
Num Queries	27.1	4431.6	77.0	3225.7

Table 7: Comparison of the search methods from GENETICATTACK and TEXTFOOLER with two sets of constraints (TEXTFOOLER and TFADJUSTED). Attacks were run on 1000 samples against BERT fine-tuned on the MR dataset. GENETICATTACK’s genetic algorithm is more successful than TEXTFOOLER’s greedy strategy, albeit much less efficient.

Our Analysis paper: Reevaluating Adversarial Examples in Natural Language
• 2020 [EMNLP Findings](#)

AE NLP
literature
is messy
(chaotic)

2. No standard library

Problems with Current NLP Attack Ecosystem

Many attacks, but Each
implemented and
benchmarked in separate
codebases (if released at all)

- Hard to trust literature comparisons because implementation differences can affect results
- hard to benchmark

Problems with Current NLP Attack Ecosystem

Many attacks, but Each implemented and benchmarked in [separate codebases](#) (if released at all)

- Hard to trust literature comparisons because implementation differences can affect results
- hard to benchmark

Challenging to develop new attacks re-using existing components

- Lots of overlap between attacks (e.g. synonym substitution techniques), but little standardization or re-usability

Problems with Current NLP Attack Ecosystem

Many attacks, but Each implemented and benchmarked in separate codebases (if released at all)

- Hard to trust literature comparisons because implementation differences can affect results
- hard to benchmark

Challenging to develop new attacks re-using existing components

- Lots of overlap between attacks (e.g. synonym substitution techniques), but little standardization or re-usability

Difficult to utilize attacks and attack components for improving models

- Attack implementations are almost never model-agnostic
- Adversarial training code is usually unreleased or non-existent

Generating an adversarial example therefore becomes an optimization problem where

(3) How to define $L(\cdot)$?

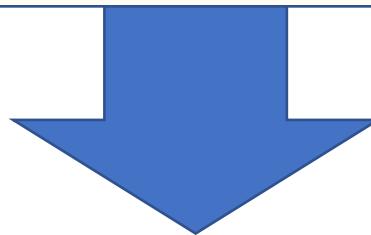
$$\text{maximize } L(F, x + \delta, y)$$

(4) How to
optimize this?

(1) How to define δ ?

$$\text{subject to } \|\delta\|_p \leq \epsilon$$

(2) How to limit δ ?



(4) How to
optimize this?

(3) Goal Function term

$$\text{maxime: } L(F, T(x), y)$$

$T(x)$

(1) Transformation term

Subject To:

$$C_1(T(x)) \wedge C_2(T(x)) \wedge \cdots \wedge C_m(T(x))$$

(2) Constraints' term

Standardize Generating NLP adversarial examples

Four Components Framework:

1. Transformation: mechanism for generating potential adversarial examples
2. Constraints: linguistic requirements for valid adversarial examples
3. Goal Function: defines end-goal for adversarial attack
4. Search Algorithm: method for finding sequence of transformations that produce valid adversarial examples defined by goal function and constraints

(4) How to
optimize this?

(3) Goal Function term

$T(x)$

maxime: $L(F, T(x), y)$

(1) Transformation term

Subject To:

$C_1(T(x)) \wedge C_2(T(x)) \wedge \dots \wedge C_m(T(x))\}$

(2) Constraints' term

Our Paper: TextAttack: A Framework for Adversarial Attacks, Data
Augmentation, and Adversarial Training in NLP

• 2020 [EMNLP Demo](#)

$T(x)$:

Transformation term

all of these are TextAttack **transformations**
(textattack.transformations)

Transformation: Word Substitution centered

- **Thesaurus:** Look up the word in a thesaurus
- **Hybrid:** Search for nearest neighbors in the *counter-fitted* embedding space (*Mrkšić et al, 2016*)
- **Embeddings:** Search for nearest-neighbors in the embedding space

$T(x)$:

Transformation term

all of these are TextAttack **transformations**
(textattack.transformations)

Transformation
by Perturbing
with synonyms

Lexical knowledge base

- WordNet ([Miller, 1995](#))
- HowNet ([Dong et al., 2010](#))

Word embedding space

- Counter-fitted
- GloVe

Masked language model

- BERT ([Devlin et al., 2018](#))
- RoBERTa ([Liu et al., 2019](#))

$$C_1(T(x)) \wedge C_2(T(x)) \wedge \dots \wedge C_m(T(x))\}$$

Constraints' term

all of these are TextAttack **constraints**
(textattack.constraints)

Constraints

Examples:

1. **Word Embedding Similarity**: When replacing x_i with x'_i via counter-fitted word embeddings, we require that embedding of x_i and x'_i satisfy minimum cosine similarity.
2. **Part-of-speech Consistency**: To preserve fluency, we require that the two words being swapped have the same part-of-speech.
3. **Sentence Encoding Similarity**: Using sentence encoders trained for semantic textual similarity, we compare the sentence encodings of original text x and perturbed text x' .

Our Analysis paper: Reevaluating Adversarial Examples in Natural Language
• 2020 [EMNLP Findings](#)

$$L(F, T(x), y)$$

Goal Function term

TextAttack **goal functions**
(textattack.goal_functions)

Goal Function:

- A way to know whether an example successfully fools the model.

Untargeted $F(x + \delta) \neq y$

Targeted $F(x + \delta) = t$

Many more special scoring for e.g. Seq2Seq outputs

Search Algorithm to find the best $T(\mathbf{x})$

- A way to **search** the space of transformations for a valid, successful adversarial example.
- Details in next section

(4) How to
optimize this?

(3) Goal Function term

$T(\mathbf{x})$

maxime: $L(F, T(\mathbf{x}), y)$

(1) Transformation term

Subject To:

$C_1(T(\mathbf{x})) \wedge C_2(T(\mathbf{x})) \wedge \cdots \wedge C_m(T(\mathbf{x}))$

(2) Constraints' term

The TextAttack Framework

NLP attacks can be constructed from four components:

- 1. transformation** (`textattack.transformations.Transformation`)
- 2. constraint(s)** (`list(textattack.constraints.Constraint)`)
- 3. search method** (`textattack.search_methods.SearchMethod`)
- 4. goal function** (`textattack.goal_functions.GoalFunction`)

(4) How to
optimize this?

(3) Goal Function term

$T(x)$

maxime: $L(F, T(x), y)$

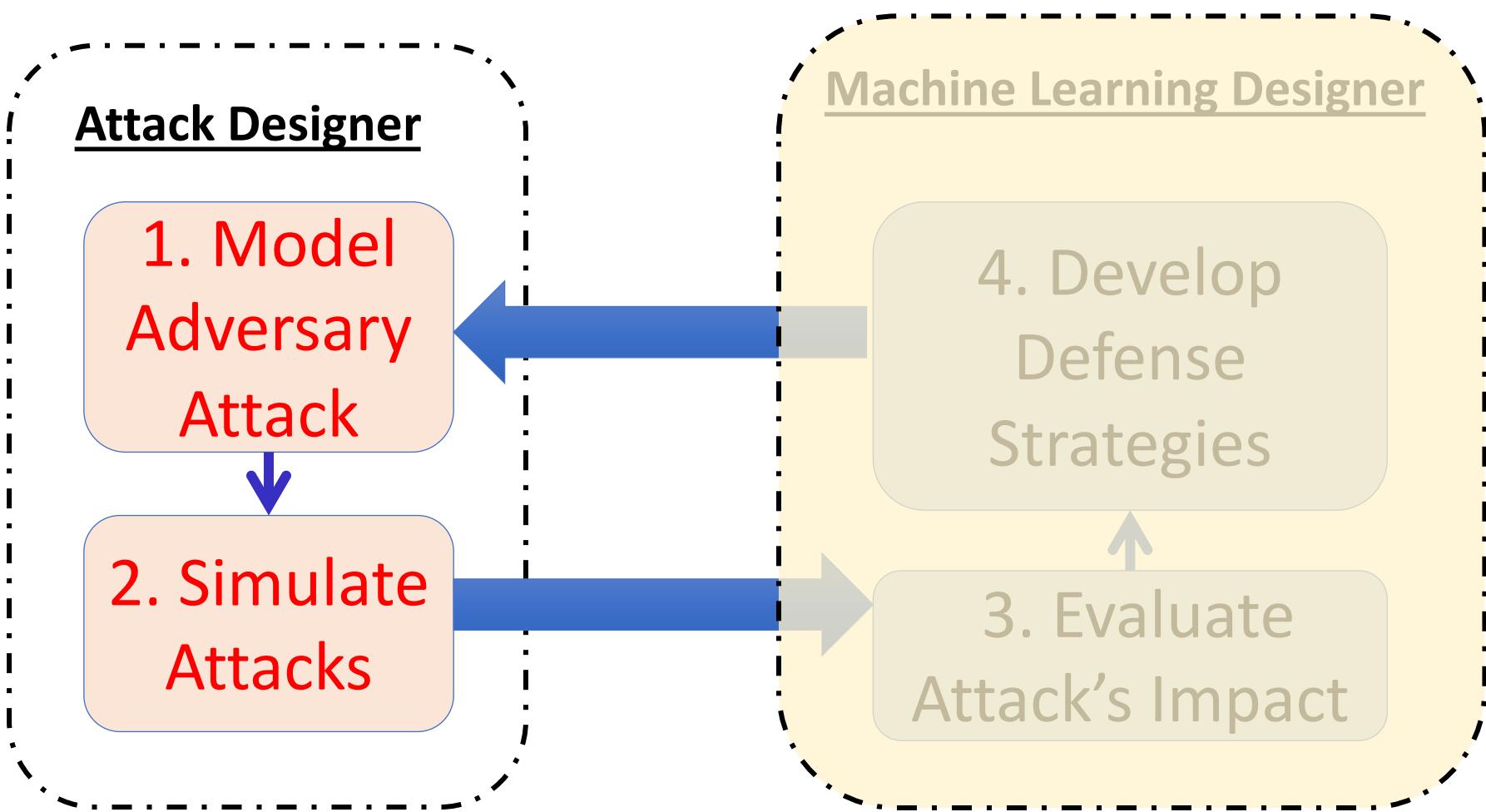
(1) Transformation term

Subject To:

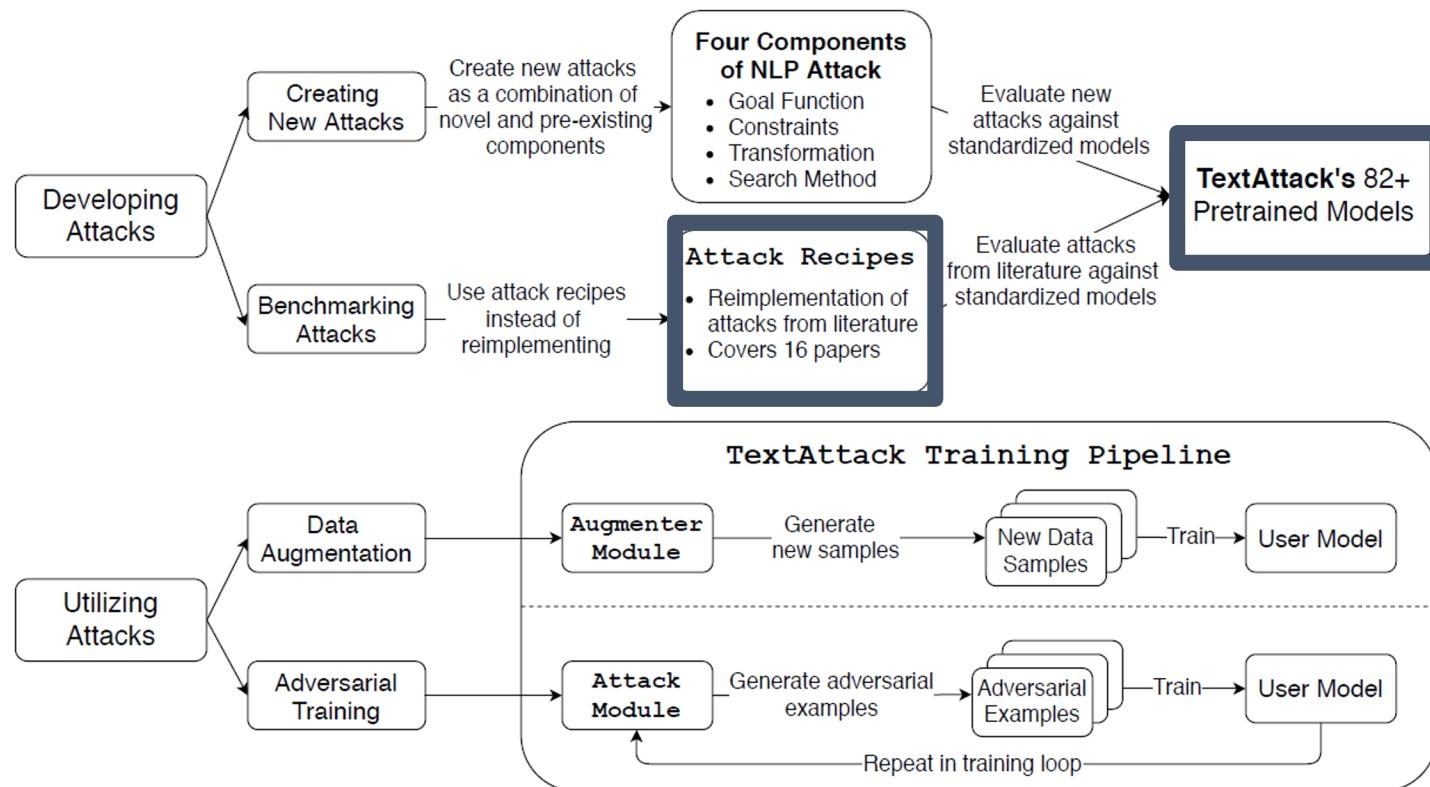
$C_1(T(x)) \wedge C_2(T(x)) \wedge \cdots \wedge C_m(T(x))\}$

(2) Constraints' term

Goal of Adversarial Machine Learning



TextAttack's Features



Is BERT Really Robust? A Strong Baseline for Natural Language Attack on Text Classification and Entailment

Di Jin,^{1*} Zhijing Jin,^{2*} Joey Tianyi Zhou,³ Peter Szolovits¹

¹Computer Science & Artificial Intelligence Laboratory, MIT

²University of Hong Kong

³A*STAR, Singapore

jindi15@mit.edu, zhijing.jin@connect.hku.hk, zhouty@ihpc.a-star.edu.sg, psz@mit.edu

Abstract

Machine learning algorithms are often vulnerable to adversarial examples that have imperceptible alterations from the original counterparts but can fool the state-of-the-art models. It is helpful to evaluate or even improve the robustness of these models by exposing the maliciously crafted adversarial examples. In this paper, we present **TEXTFOOLER**, a simple but strong baseline to generate adversarial text. By applying it to two fundamental natural language tasks, text classification and textual entailment, we successfully attacked three target models, including the powerful pre-trained BERT, and the widely used convolutional and recurrent neural networks. We demonstrate three advantages of this framework: (1) effective—it outperforms previous attacks by success rate and perturbation rate, (2) utility-preserving—it preserves semantic content, grammaticality, and correct types classified by humans, and (3) efficient—it generates adversarial text with computational complexity linear to the text length.¹

Classification Task: Is this a *positive* or *negative* review?

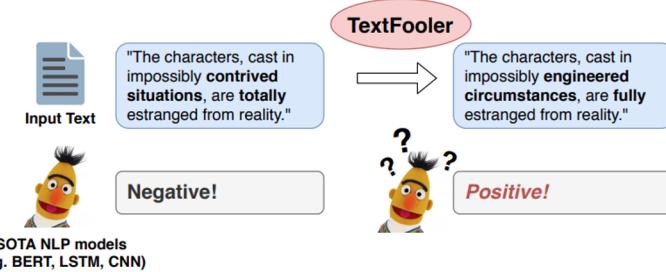


Figure 1: Our model TextFooler slightly change the input text but completely altered the prediction result.

al. 2013; Carlini and Wagner 2018), it is still challenging to deal with text data due to its discrete nature. Formally, besides the ability to fool the target models, outputs of a natural

Algorithm 1 Adversarial Attack by TEXTFOOLER

Input: Sentence example $X = \{w_1, w_2, \dots, w_n\}$, the corresponding ground truth label Y , target model F , sentence similarity function $\text{Sim}(\cdot)$, sentence similarity threshold ϵ , word embeddings Emb over the vocabulary Vocab .

Output: Adversarial example X_{adv}

- 1: Initialization: $X_{\text{adv}} \leftarrow X$
- 2: **for** each word w_i in X **do**
- 3: Compute the importance score I_{w_i} via Eq. (2)
- 4: **end for**
- 5:
- 6: Create a set W of all words $w_i \in X$ sorted by the descending order of their importance score I_{w_i} .
- 7: Filter out the stop words in W .
- 8: **for** each word w_j in W **do**
- 9: Initiate the set of candidates CANDIDATES by extracting the top N synonyms using $\text{CosSim}(\text{Emb}_{w_j}, \text{Emb}_{\text{word}})$ for each word in Vocab .
- 10: $\text{CANDIDATES} \leftarrow \text{POSFilter}(\text{CANDIDATES})$
- 11: $\text{FINCANDIDATES} \leftarrow \{\}$
- 12: **for** c_k in CANDIDATES **do**
- 13: $X' \leftarrow \text{Replace } w_j \text{ with } c_k \text{ in } X_{\text{adv}}$
- 14: **if** $\text{Sim}(X', X_{\text{adv}}) > \epsilon$ **then**
- 15: Add c_k to the set FINCANDIDATES
- 16: $Y_k \leftarrow F(X')$
- 17: $P_k \leftarrow F_{Y_k}(X')$
- 18: **end if**
- 19: **end for**
- 20: **if** there exists c_k whose prediction result $Y_k \neq Y$ **then**
- 21: In FINCANDIDATES , only keep the candidates c_k whose prediction result $Y_k \neq Y$
- 22: $c^* \leftarrow \underset{c \in \text{FINCANDIDATES}}{\text{argmax}} \text{Sim}(X, X'_{w_j \rightarrow c})$
- 23: $X_{\text{adv}} \leftarrow \text{Replace } w_j \text{ with } c^* \text{ in } X_{\text{adv}}$
- 24: **return** X_{adv}
- 25: **else if** $P_{Y_k}(X_{\text{adv}}) > \min_{c_k \in \text{FINCANDIDATES}} P_k$ **then**
- 26: $c^* \leftarrow \underset{c_k \in \text{FINCANDIDATES}}{\text{argmin}} P_k$
- 27: $X_{\text{adv}} \leftarrow \text{Replace } w_j \text{ with } c^* \text{ in } X_{\text{adv}}$
- 28: **end if**
- 29: **end for**
- 30: **return** None

Four Components in Action

TextFooler method proposed by Jin et al. (2019)

Search Algorithm: Greedy with Word Importance Ranking

Transformation: Counter-fitted embedding word swap

Constraint #3: Cosine similarity of sentence embeddings

Goal Function: Untargeted attack for classification

Algorithm 1 Adversarial Attack by TEXTFOOLER

Input: Sentence example $X = \{w_1, w_2, \dots, w_n\}$, the corresponding ground truth label Y , target model F , sentence similarity function $\text{Sim}(\cdot)$, sentence similarity threshold ϵ , word embeddings Emb over the vocabulary Vocab .

Output: Adversarial example X_{adv}

```

1: Initialization:  $X_{\text{adv}} \leftarrow X$ 
2: for each word  $w_i$  in  $X$  do
3:   Compute the importance score  $I_{w_i}$  via Eq. (2)
4: end for
5:
6: Create a set  $W$  of all words  $w_i \in X$  sorted by the descending order of their importance score  $I_{w_i}$ .
7: Filter out the stop words in  $W$ .
8: for each word  $w_i$  in  $W$  do
9:   Initiate the set of candidates CANDIDATES by extracting the top  $N$  synonyms using  $\text{CosSim}(\text{Emb}_{w_j}, \text{Emb}_{\text{word}})$  for each word in  $\text{Vocab}$ .
10:  CANDIDATES  $\leftarrow \text{POSFilter}(\text{CANDIDATES})$ 
11:  FINCANDIDATES  $\leftarrow \{\}$ 
12:  for  $c_k$  in CANDIDATES do
13:     $X' \leftarrow \text{Replace } w_i \text{ with } c_k \text{ in } X_{\text{adv}}$ 
14:    if  $\text{Sim}(X', X_{\text{adv}}) > \epsilon$  then
15:      Add  $c_k$  to the set FINCANDIDATES
16:       $Y_k \leftarrow F(X')$ 
17:       $P_k \leftarrow F_{Y_k}(X')$ 
18:    end if
19:  end for
20:  if there exists  $c_k$  whose prediction result  $Y_k \neq Y$  then
21:    In FINCANDIDATES, only keep the candidates  $c_k$  whose prediction result  $Y_k \neq Y$ 
22:     $c^* \leftarrow \underset{c \in \text{FINCANDIDATES}}{\text{argmax}} \text{Sim}(X, X'_{w_j \rightarrow c})$ 
23:     $X_{\text{adv}} \leftarrow \text{Replace } w_j \text{ with } c^* \text{ in } X_{\text{adv}}$ 
24:  return  $X_{\text{adv}}$ 
25: else if  $P_{Y_k}(X_{\text{adv}}) > \min_{c_k \in \text{FINCANDIDATES}} P_k$  then
26:    $c^* \leftarrow \underset{c_k \in \text{FINCANDIDATES}}{\text{argmin}} P_k$ 
27:    $X_{\text{adv}} \leftarrow \text{Replace } w_j \text{ with } c^* \text{ in } X_{\text{adv}}$ 
28: end if
29: end for

```

Constraint #1: Cosine similarity of word embeddings

Constraint #2: Consistent part-of-speech

Four Components Standardized 18 Attacks:

	Alzantot et al. (2018)	Jin et al. (2019)
Goal Function	UntargetedClassification	UntargetedClassification
Search Method	GeneticAlgorithmWordSwap	GreedyWordSwapWordImportanceRanking
Transformation	WordSwapEmbedding(embedding='cf')	WordSwapEmbedding(embedding='cf')
Constraints	<ul style="list-style-type: none">WordsPerturbedPercentage(max_perc=20)WordEmbeddingDistance(max_mse=0.5)GoogleLanguageModel(n_per_index=4)	<ul style="list-style-type: none">WordEmbeddingDistance(min_cos_sim=0.5)PartOfSpeech(verb_noun_swap=True)UniversalSentenceEncoder(metric='angular', thresh=0.904458599)

Installing TextAttack

 Github PyTest passing pypi package 0.2.12

pip install textattack

<https://github.com/QData/TextAttack>

QData / TextAttack

Watch ▾

28

Star

1.5k

Fork

183

Code

Issues 39

Pull requests 9

Actions

Projects 6

Wiki

...

master ▾

Go to file

Add file ▾

Code ▾

About



qianjun make attack-recipes more easier to find in doc ... ✓ 6 days ago 2,055



.github Update run-pytest.yml

4 months ago



docs make attack-recipes more easier to find in doc

6 days ago



examples isort format of attack_camembert

8 months ago



tests [TEST] Update expected test output

14 days ago



textattack [DOC, CODE] fix documentation and minor b...

15 days ago



.gitignore fix issues to pass tests

20 days ago



readthedocs.yml fix readthedocs module load

7 months ago

TextAttack  is a Python framework for adversarial attacks, data augmentation, and model training in NLP
<https://textattack.readthedocs.io/en/latest/>

[textattack.readthedocs.i...](https://textattack.readthedocs.io/en/latest/)

nlp security

machine-learning

natural-language-processing

data-augmentation

We have also
shared 82
Pretrained
Models

Integration
with
HuggingFace's
[Model Hub](#)
and [nlp](#) library

- Can attack any model on the model hub on any dataset from nlp

TextAttack has
82 pretrained
models on its
[Model Hub page](#)

- Models: BERT,
DistilBERT, ALBERT,
BART, RoBERTa, XLNet
- Trained on all [GLUE](#) tasks

TextAttack



An open-source, model-agnostic library for **attacking NLP models and standardizing evaluations**

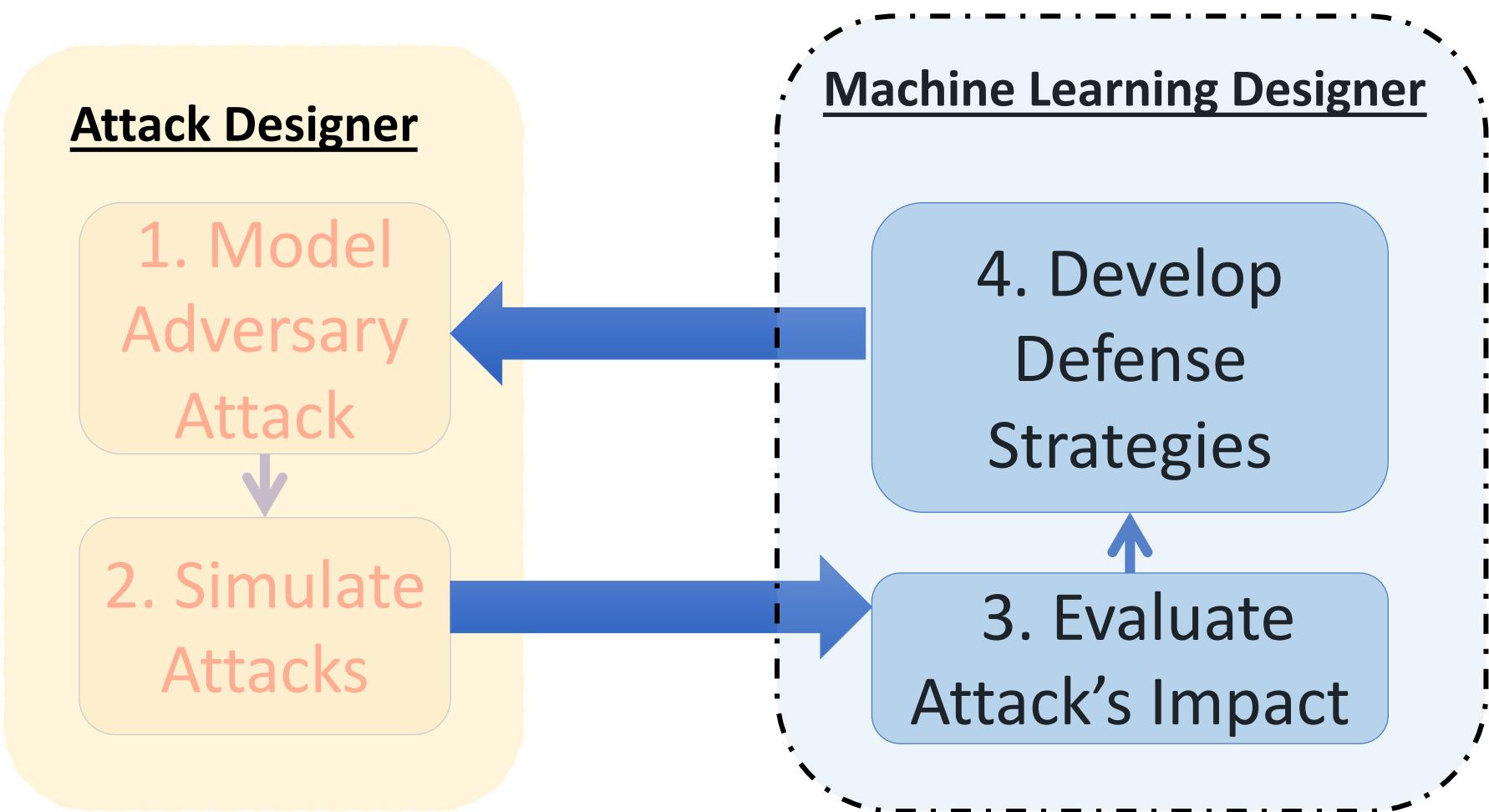
Contains 18 popular white-box and black-box attacks

82 Pre-trained models on popular datasets across multiple Types of NLP tasks

AE NLP
literature
is messy
(chaotic)

3. No clear benchmarking
insights

Goal of Adversarial Machine Learning



Search Method

Typically, one word replacement is not enough to change the model's prediction. Instead, a set of word replacements is necessary.

(4) How to
optimize this?

(3) Goal Function term

$T(x)$

maxime: $L(F, T(x), y)$

(1) Transformation term

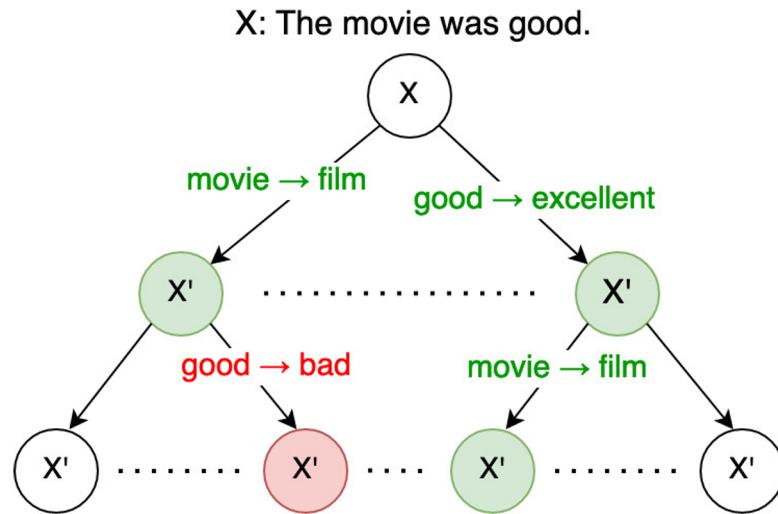
Subject To:

$C_1(T(x)) \wedge C_2(T(x)) \wedge \cdots \wedge C_m(T(x))$

(2) Constraints' term

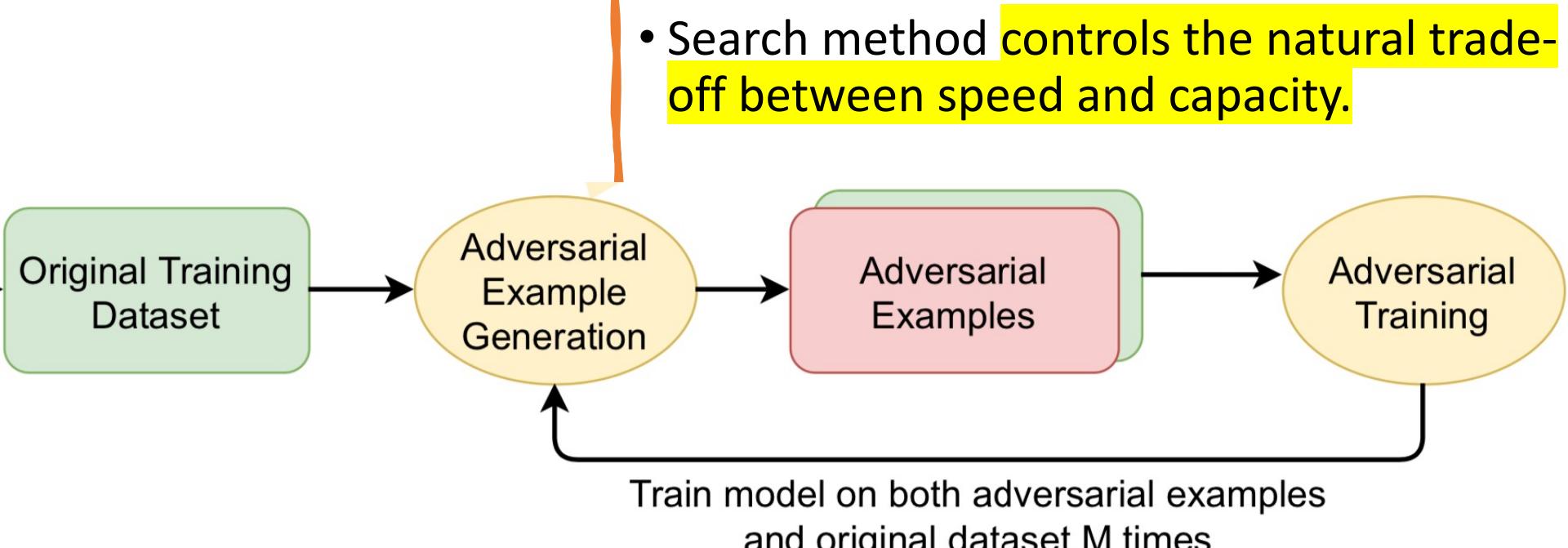
Search Method

Typically, one word replacement is not enough to change the model's prediction. Instead, a set of word replacements is necessary.



Motivation

- Adversarial training uses both clean examples and adversarial examples to train robust models.
- Two criteria we need to consider when constructing attacks for adversarial training are:
 1. Speed
 2. Capability
- Search method controls the natural trade-off between speed and capacity.

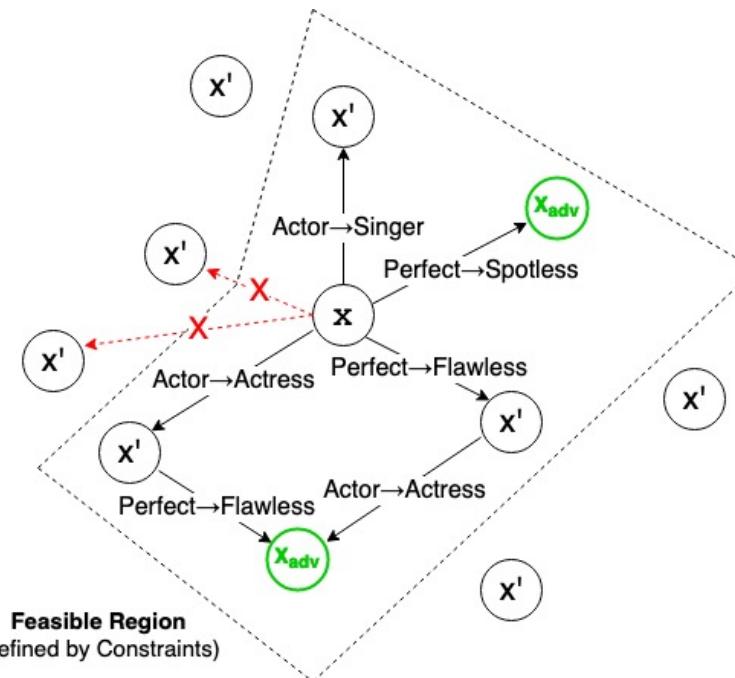


Search Algorithm

Search Algorithm: A way to **search** the space of transformations for a valid, successful adversarial example.

Why a search algorithm?

- We need to find set of transformations that successfully produce x_{adv}
- Combinatorial search problem with heuristic $score(x)$ provided by goal function



Our Analysis paper: Searching for a Search Method: Benchmarking
Search Algorithms for Generating NLP Adversarial Examples

•2020 [EMNLP BlackBoxNLP](#)

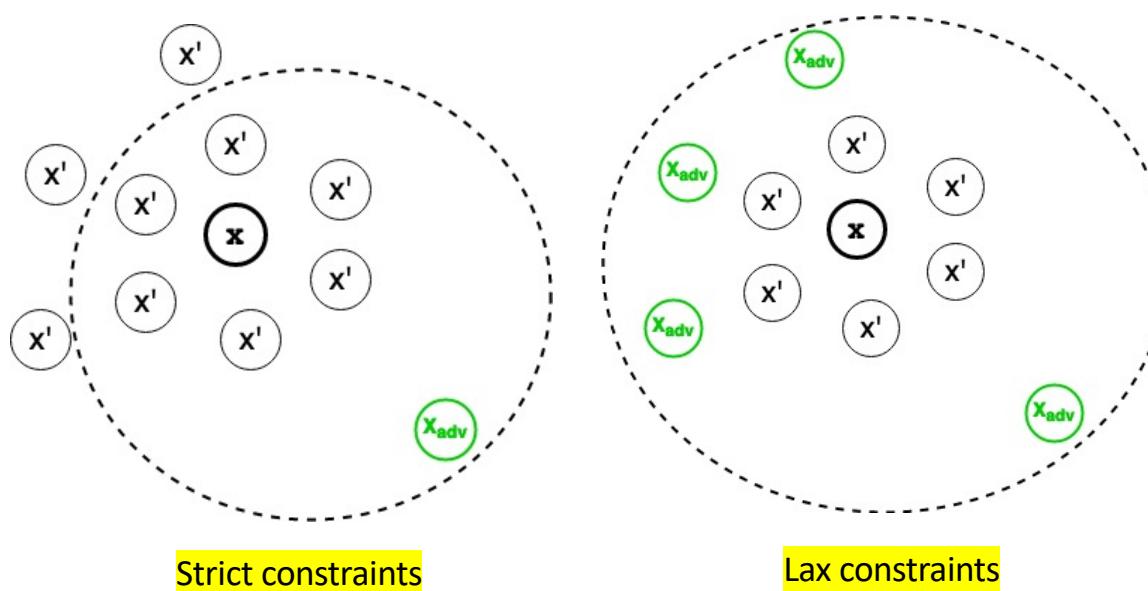
Search Space

Search Algorithm: A way to **search** the space of transformations for a valid, successful adversarial example.

Search space defined by transformation and constraints

Let $T(x)$ be our transformation and $C_i(x)$ be a constraint,

$$S(x) = \{T(x) | C_1(T(x)) \wedge C_2(T(x)) \wedge \dots \wedge C_m(T(x))\}$$



Our Analysis paper: Searching for a Search Method: Benchmarking
Search Algorithms for Generating NLP Adversarial Examples

•2020 [EMNLP BlackBoxNLP](#)

Search Algorithms from Literature

Beam Search (Ebrahimi et al., 2017)

Greedy with Word Importance Ranking

- UNK (Gao et al., 2018)
- DEL (Jin et al., 2019)
- PWWS (Ren et al., 2019)

Genetic Algorithm (Alzantot et al., 2018),

Particle Swarm Optimization (Zang et al., 2020)

Problems in Current Literature

Inconsistent search
space for comparisons

Lack of comprehensive
performance
benchmark for search
algorithm

Lack of comprehensive
speed benchmark for
search algorithm

Our Analysis paper: Searching for a Search Method: Benchmarking
Search Algorithms for Generating NLP Adversarial Examples

•2020 [EMNLP BlackBoxNLP](#)

Benchmarking Insights



Optimal search for absolute performance is beam search with beam width of 8.



When within a small query budget,
greedy with word importance
ranking is most effective

For two constraint settings across three datasets,
the relative differences between the attack success
rates of greedy with word importance ranking and
the success rates of beam search are less than 20%.

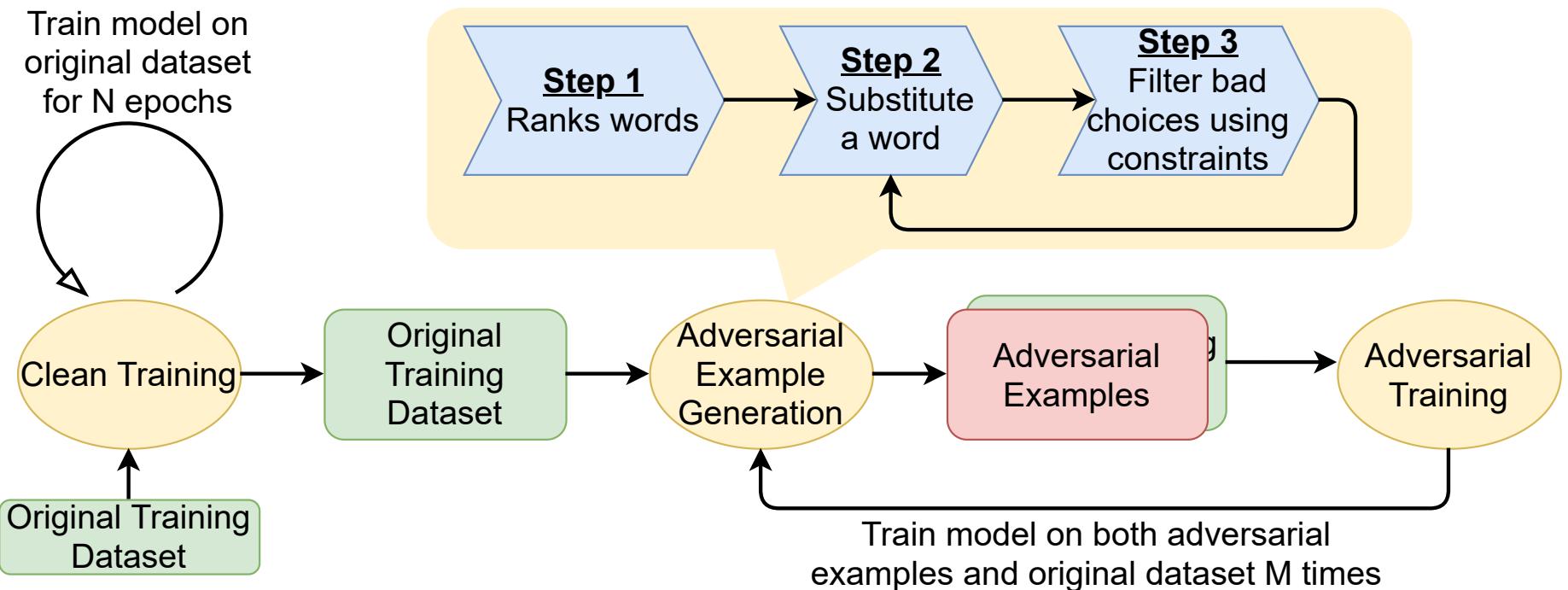


If only aiming for attack success,
Search algorithms matter less than
transformations and constraints.

Although changing the search methods did not
change attack success rate by more than 20%,
changing the constraints changed attack success
rate by over 60%.

AE NLP
literature
is messy
(chaotic)

4. No clear benefits



Adversarial Training in recent NLP Literature

Collection of recent works on **adversarial training** do not study
whether it defends against **adversarial attacks** proposed in literature!

Adversarial Training in recent NLP Literature

Collection of recent works on **adversarial training** do not study whether it defends against **adversarial attacks** proposed in literature!

Adversarial Training

- *FreeLB* (Zhu et al., 2019)
- *SMART* (Jiang et al., 2019)
- *ALUM* (Liu et al., 2020)

Adversarial Training in recent NLP Literature

Collection of recent works on **adversarial training** do not study whether it defends against **adversarial attacks** proposed in literature!

Adversarial Training

- *FreeLB* ([Zhu et al., 2019](#))
- *SMART* ([Jiang et al., 2019](#))
- *ALUM* ([Liu et al., 2020](#))

Adversarial Attacks

- [Alzantot et al. \(2018\)](#)
- *TextFooler* ([Jin et al., 2019](#))
- *PWWS* ([Ren et al., 2019](#))
- [Zang et al. \(2020\)](#)
- *BAE* ([Garg and Ramakrishnan, 2020](#))
- *BERT-Attack* ([Li et al., 2020](#))
- *CLARE* ([Li et al., 2021](#))

Adversarial Training in NLP

Collection of recent works on **adversarial training** do not study whether it defends against **adversarial attacks** proposed in literature!

Adversarial Training

- *FreeLB* ([Zhu et al., 2019](#))
- *SMART* ([Jiang et al., 2019](#))
- *ALUM* ([Liu et al., 2020](#))

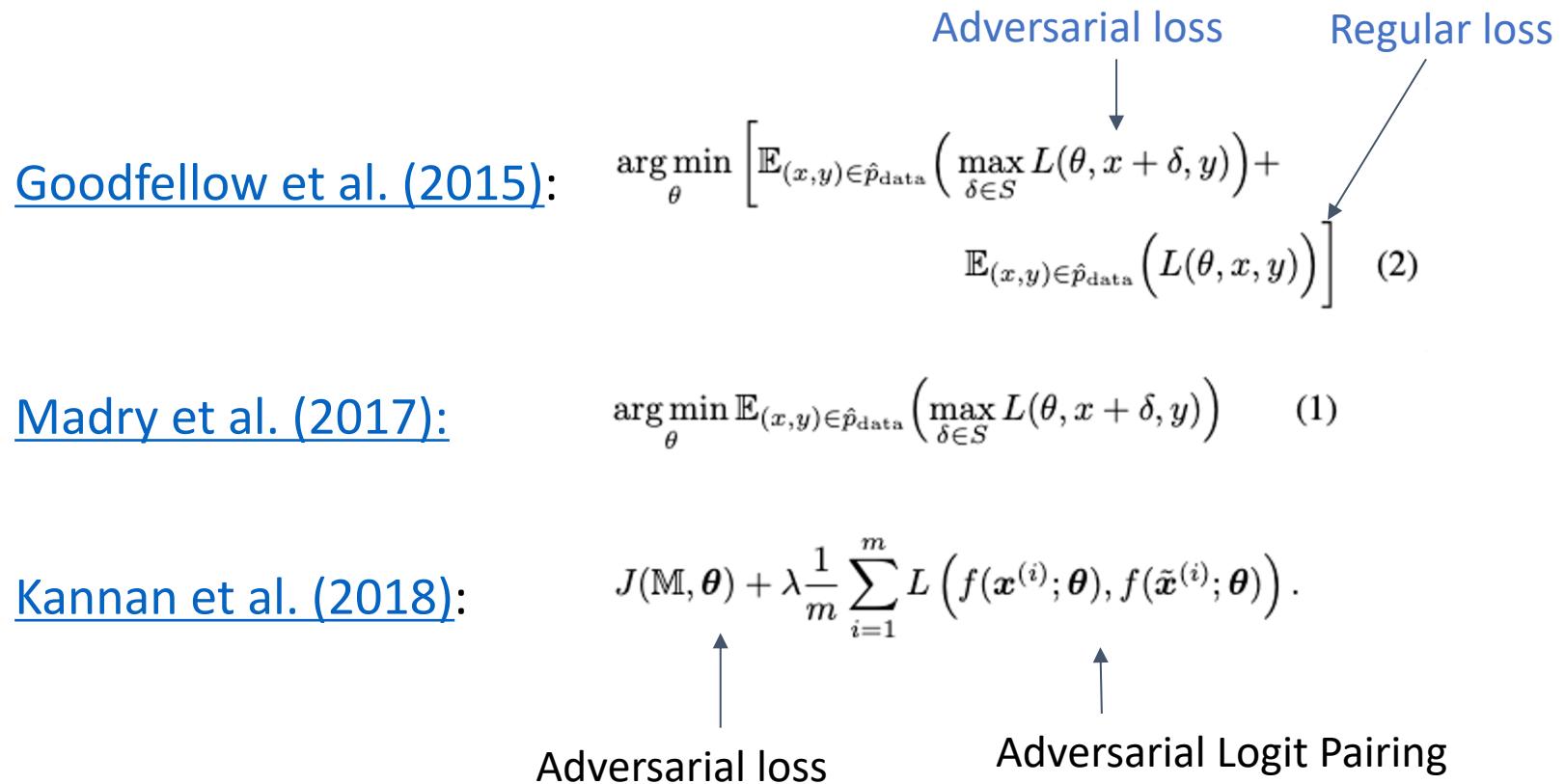


Adversarial Attacks

- [Alzantot et al. \(2018\)](#)
- [TextFooler](#) ([Jin et al., 2019](#))
- *PWWS* ([Ren et al., 2019](#))
- [Zang et al. \(2020\)](#)
- *BAE* ([Garg and Ramakrishnan, 2020](#))
- *BERT-Attack* ([Li et al., 2020](#))
- *CLARE* ([Li et al., 2021](#))

Recent NLP Adversarial Training (above) add perturbations in the embedding, instead of in the input space; They don't evaluate robustness well.

Adversarial Training for Robustness



We propose faster attacks that suit for vanilla adversarial training.

- Many engineering tricks to make vanilla adversarial training feasible in NLP
- We observe that
 - Adversarial training can help improve adversarial robustness against attacks that were *not used to train the model*.
 - Adversarial training can provide a regularization effect and improve the model's *standard accuracy* and *cross-domain generalization*.
 - Adversarial training can improve the model's *interpretability*.

TextAttack Rescues Messy AE NLP literature

1. Many generated examples are bad
2. No standard library
3. No clear benchmarking insights
4. No clear benefits

Who is TextAttack for?

- **researchers** who want to implement new NLP attacks or compare them in a standardized framework
- **any machine learning practitioner** who want to understand their limitations of NLP models and/or use adversarial training to make their models better
- **anyone training an NLP model** who wants to apply data augmentation to increase test-set accuracy by 1-2%

<http://trustworthymachinelearning.org>

1. To Fool / Evade Learned Models

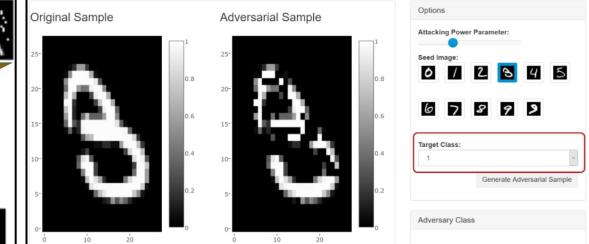
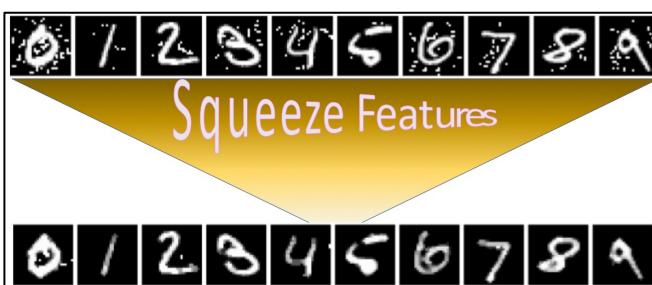
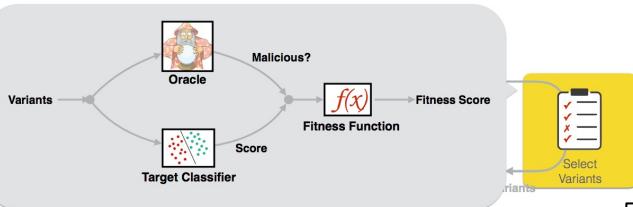
2. To Detect fooling/ Evasion

3. To Defend Against Evasion

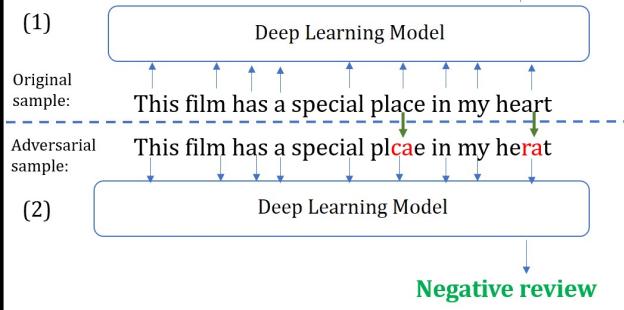
4. To Visualize and Benchmarking

5. To Understand Theoretically

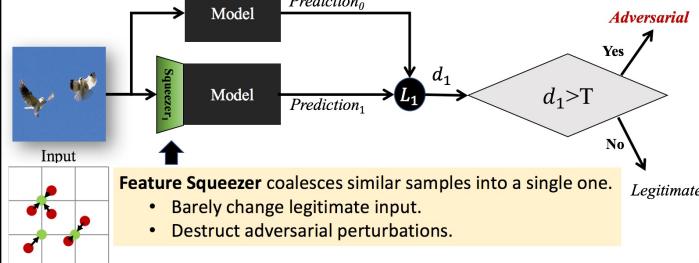
Automated Evasion Approach Based on Genetic Programming



Positive review



Detection Framework



(X, d'_1)

f_1 Machine-learning classifier

g_1



(X_1, d_1)

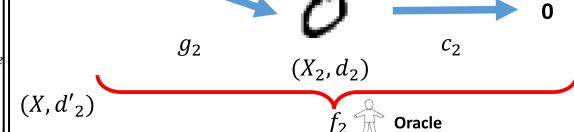
$c_1 \rightarrow 0$

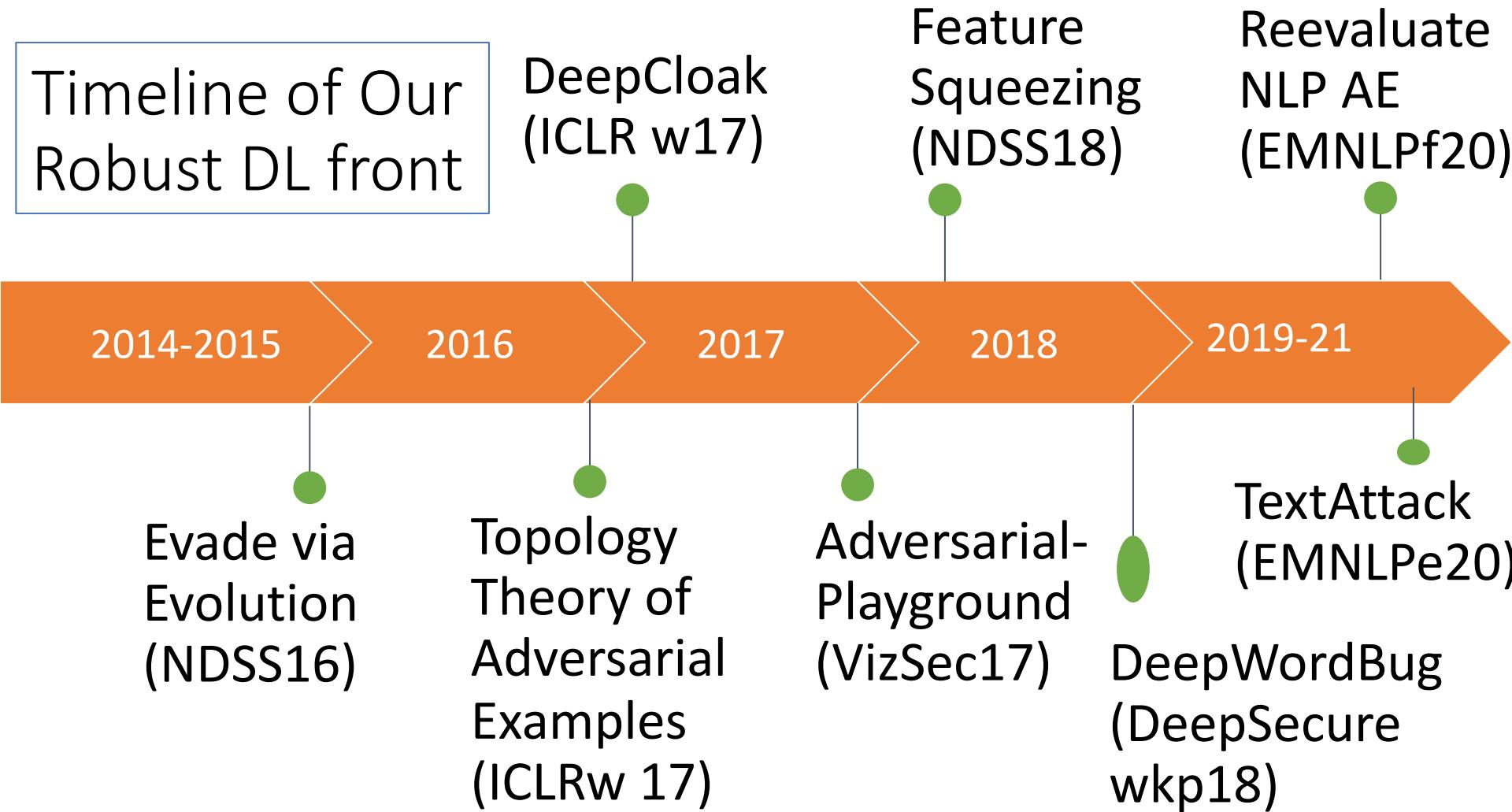
g_2



(X_2, d_2)

$c_2 \rightarrow 0$





Acknowledgements

My Students on this project:



UVA Computer Science Dept. Security Research Group: Prof. David Evans

UVA Computer Science Dept. Software Safety Group: Prof. Matthew B Dwyer

UVA Computer Science Dept. NLP Research Group: Prof. Yangfeng Ji

UVA Computer Science Dept. Software Engineering Group: Prof. Mary Lou Soffa



6/28/21

Yanjun Qi / UVA CS



106



Thank you

What can I do with TextAttack?

- run **standardized attack recipes** on models & datasets (yours or ours)
- visualize attack results using the command line, Visdom, W&B, etc.
- or, use the infrastructure of TextAttack to develop and benchmark **your own NLP attacks**
- or, use the components of TextAttack for **adversarial training**
- or, use the components from TextAttack for **data augmentation**