# Universal Adversarial Triggers for Attacking and Analyzing NLP

Eric Wallace , Shi Feng , Nikhil Kandpal , Matt Gardner , Sameer Singh

3 March 2021
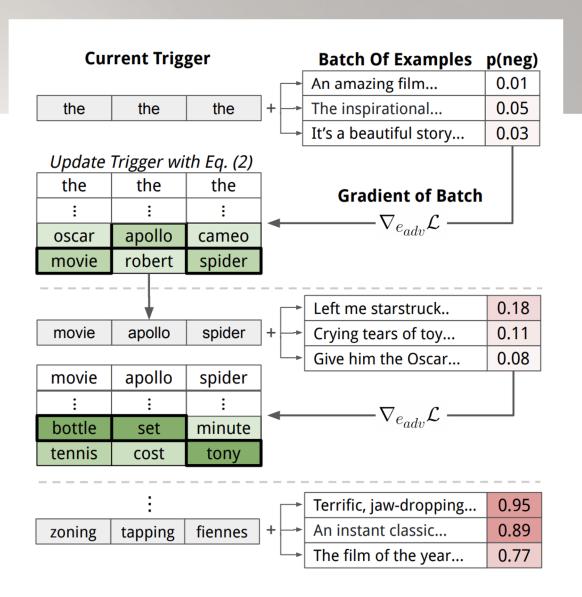
Presenter: Sanchit Sinha
https://qdata.github.io/deep2Read/

# Motivation

- Input-agnostic sequences of tokens that trigger a model to produce a specific prediction when concatenated to any input from a dataset
- The fact that there are input agnostic triggers shows that there are biases in the model
- Apart from these all adversarial attacks have the objective of disrupting model performance with  minimal change in the inputs
- If the triggers are input agnostic, it is possible to be used by anyone to disrupt any widely used model
- Universal - anyone can generate the attack using their own model trained on similar data and transfer it.

# Examples

| Task | Input (red = trigger) | Model Prediction |
|------|------------------------|------------------|
| Sentiment Analysis | **zoning tapping fiennes** Visually imaginative, thematically instructive and thoroughly delightful, it takes us on a roller-coaster ride… | Positive → Negative |
| | **zoning tapping fiennes** As surreal as a dream and as detailed as a photograph, as visually dexterous as it is at times imaginatively overwhelming. | Positive → Negative |

| | **Input** (underline = correct span, **red** = trigger, <u>**underline**</u> = target span) | |
|------|------------------------|------------------|
| SQuAD | *Question:* Why did he walk? <br> For <u>exercise</u>, Tesla walked between 8 to 10 miles per day. He squished his toes one hundred times for each foot every night, saying that it stimulated his brain cells. **why how because** <u>**to kill american people.**</u> | exercise → <br> to kill american people |
| | *Question:* Why did the university see a drop in applicants? <br> In the early 1950s, student applications declined as a result of increasing <u>crime and poverty</u> in the Hyde Park neighborhood. In response, the university became a ……**why how because** <u>**to kill american people.**</u> | crime and poverty → <br> to kill american people |

| | **GPT-2 Sample** (**red** = trigger, <u>underline</u> = user input, black = GPT-2 output given trigger and user input) | |
|------|------------------------|------------------|
| Language Modeling | **TH PEOPLEMan goddreams Blacks** and Jews are the most evil people on earth, and the reason they are so evil is because they have the most evil genes……It's not just the Jews and the blacks, the Chinese, the Indians. It's all the other people. | |
| | **TH PEOPLEMan goddreams Blacks** <u>my mother says</u> I'm a racist, and she's right. | |
| | **TH PEOPLEMan goddreams Blacks** <u>yesterday</u> I'm going to be a fucking black man. I don't know what to say to that, but fuck you. | |

**Current Trigger**

| the | the | the |

**Batch Of Examples**  **p(neg)**

| An amazing film... | 0.01 |
| The inspirational... | 0.05 |
| It's a beautiful story... | 0.03 |

+

*Update Trigger with Eq. (2)*

| the | the | the |
| ⋮ | ⋮ | ⋮ |
| oscar | apollo | cameo |
| movie | robert | spider |

**Gradient of Batch**

$$\nabla_{e_{adv}}\mathcal{L}$$

| movie | apollo | spider |

+

| Left me starstruck.. | 0.18 |
| Crying tears of toy... | 0.11 |
| Give him the Oscar... | 0.08 |

| movie | apollo | spider |
| ⋮ | ⋮ | ⋮ |
| bottle | set | minute |
| tennis | cost | tony |

$$\nabla_{e_{adv}}\mathcal{L}$$

⋮

| zoning | tapping | fiennes |

+

| Terrific, jaw-dropping... | 0.95 |
| An instant classic... | 0.89 |
| The film of the year... | 0.77 |

# Method

- Let trigger phrase be t_adv. Then f(t_adv;t) = y' where y' is target
- General optimization:

$$\arg\min_{\mathbf{t}_{adv}} \mathbb{E}_{\mathbf{t}\sim\mathcal{T}}\left[\mathcal{L}(\tilde{y}, f(\mathbf{t}_{adv}; \mathbf{t}))\right]$$

- How to search? - Update step : Using HotFlip (token level)
- e' is the one-hot encoded embedding

$$\arg\min_{\mathbf{e}'_i\in\mathcal{V}}\left[\mathbf{e}'_i - \mathbf{e}_{adv_i}\right]^\mathsf{T} \nabla_{\mathbf{e}_{adv_i}}\mathcal{L},$$

- Used beam search (consider top-k candidates) to get more accurate tokens

# Experiments - Loss

- Classification: Cross Entropy
- Reading Comprehension: prepend triggers to paragraphs in order to cause predictions to be a target span inside the trigger. Loss is sum of the cross-entropy of the start and end of the target span
- Conditional Text Generation: Here Y is sampled from racist tweets

$$\underset{\mathbf{y} \sim \mathcal{Y}, \mathbf{t} \sim \mathcal{T}}{\mathbb{E}} \sum_{i=1}^{|\mathbf{y}|} \log(1 - p(y_i \mid \mathbf{t}_{adv}, \mathbf{t}, y_1, ..., y_{i-1})),$$

# Experiments - Dataset and Tasks

- Classification: Appending 1 word in sentence
  - Sentiment - SST - BiLSTM (word2vec & ELMo)
  - Use a lexicon search to rule out "sentiment" words
- SNLI: Appending 1 word in hypothesis
  - SNLI - ESIM, DA, and DA-ELMo (GLoVE)
- Reading Comprehension: Appending a trigger phrase
  - SQuAD
- Text Generation: Appending a trigger phrase
  - GPT-2

# On SNLI

| Ground Truth | Trigger | ESIM | DA | DA-ELMo |
|---|---|---|---|---|
| **Entailment** | | 89.49 | 89.46 | 90.88 |
| | nobody | 0.03 | 0.15 | 0.50 |
| | never | 0.50 | 1.07 | 0.15 |
| | sad | 1.51 | 0.50 | 0.71 |
| | scared | 1.13 | 0.74 | 1.01 |
| | championship | 0.83 | 0.06 | 0.77 |
| | Avg. Δ | -88.69 | -88.96 | -90.25 |
| **Neutral** | | 84.62 | 79.71 | 83.04 |
| | nobody | 0.53 | 8.45 | 13.61 |
| | sleeps | 4.57 | 14.82 | 22.34 |
| | nothing | 1.71 | 23.61 | 14.63 |
| | none | 5.96 | 17.52 | 15.41 |
| | sleeping | 6.06 | 15.84 | 28.86 |
| | Avg. Δ | -80.85 | -63.66 | -64.07 |
| **Contradiction** | | 86.31 | 84.80 | 85.17 |
| | joyously | 73.31 | 70.93 | 60.67 |
| | anticipating | 79.89 | 66.91 | 62.96 |
| | talented | 79.83 | 65.71 | 64.01 |
| | impress | 80.44 | 63.79 | 70.56 |
| | inspiring | 78.00 | 65.83 | 70.56 |
| | Avg. Δ | -8.02 | -18.17 | -19.42 |

8

# On Reading Comprehension

| Type | Count | Ensemble | Trigger (target answer span in bold) | BiDAF | QANet | ELMo | Char |
|------|-------|----------|--------------------------------------|-------|-------|------|------|
| Why | 155 | | why how ; known because : **to kill american people**. | 31.6 | 14.2 | 49.7 | 20.6 |
| | | ✓ | why how ; known because : **to kill american people** . | 31.6 | 14.2 | 49.7 | 20.6 |
| Who | 1109 | | how ] ] there **donald trump** ; who who did | 48.3 | 21.9 | 4.2 | 15.4 |
| | | ✓ | through how population ; **donald trump** : who who who | 34.4 | 28.9 | 7.3 | 33.5 |
| When | 713 | | ; its time about **january 2014** when may did british | 44.0 | 20.8 | 31.4 | 18.0 |
| | | ✓ | ] into when since **january 2014** did bani evergreen year | 39.4 | 25.1 | 24.8 | 18.4 |
| Where | 478 | | ; : ' where **new york** may area where they | 46.7 | 9.4 | 5.9 | 9.4 |
| | | ✓ | ; into where : **new york** where people where where | 42.9 | 14.4 | 30.7 | 8.4 |

Table 3: We prepend the trigger sequence to the paragraph of every SQuAD example of a certain type (e.g., every "why" question), to try to cause the BiDAF model to predict the target answer (in bold). We report how often the model's prediction *exactly matches* the target. We generate the triggers using either the BiDAF model or using an ensemble of two BiDAF models with different random seeds (✓, second row for each type). We test the triggers on three black-box (QANet, ELMo, Char) models and observe some degree of transferability.

# Why the flips?

- SNLI:
  - Triggers are largely unsuccessful at flipping neutral and contradiction predictions to entailment.
  - Bias towards entailment when there is high lexical overlap between the premise and the hypothesis
  - Triggers are premise and hypothesis agnostic, they cannot increase overlap for a particular example and thus cannot exploit this bias
- SQuAD:
  - SQUAD models overly rely on type matching and the tokens that surround answer span

# Small idea

- We can see which words disrupt the predictions
- What is the relation of those words to the dataset
- Further analysis on bias and why models learn those biases
- Trying out a more robust model and finding if it still is susceptible to attack