# Interpretability & Security in NLP models
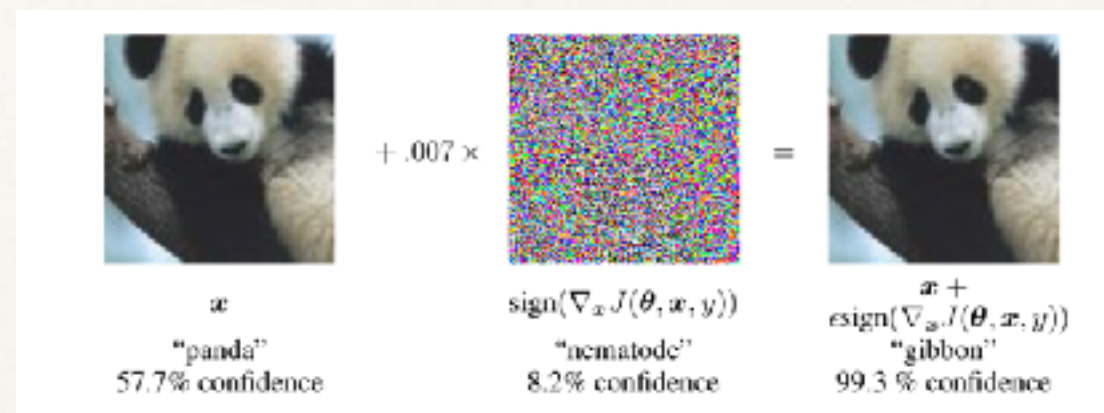
Dominika Basaj, 08/10/2018

## Agenda

- Introduction

- Related work

- Our research

- Future works

# Adversarial examples

In general:
- Inputs to machine learning models that an attacker has intentionally designed to cause the model to make a mistake
- First observed in 2014 by Szegedy et al. (2014)



$x$
"panda"
57.7% confidence

$+ .007 \times$

$\text{sign}(\nabla_x J(\theta, x, y))$
"nematode"
8.2% confidence

$=$

$x + \epsilon\,\text{sign}(\nabla_x J(\theta, x, y))$
"gibbon"
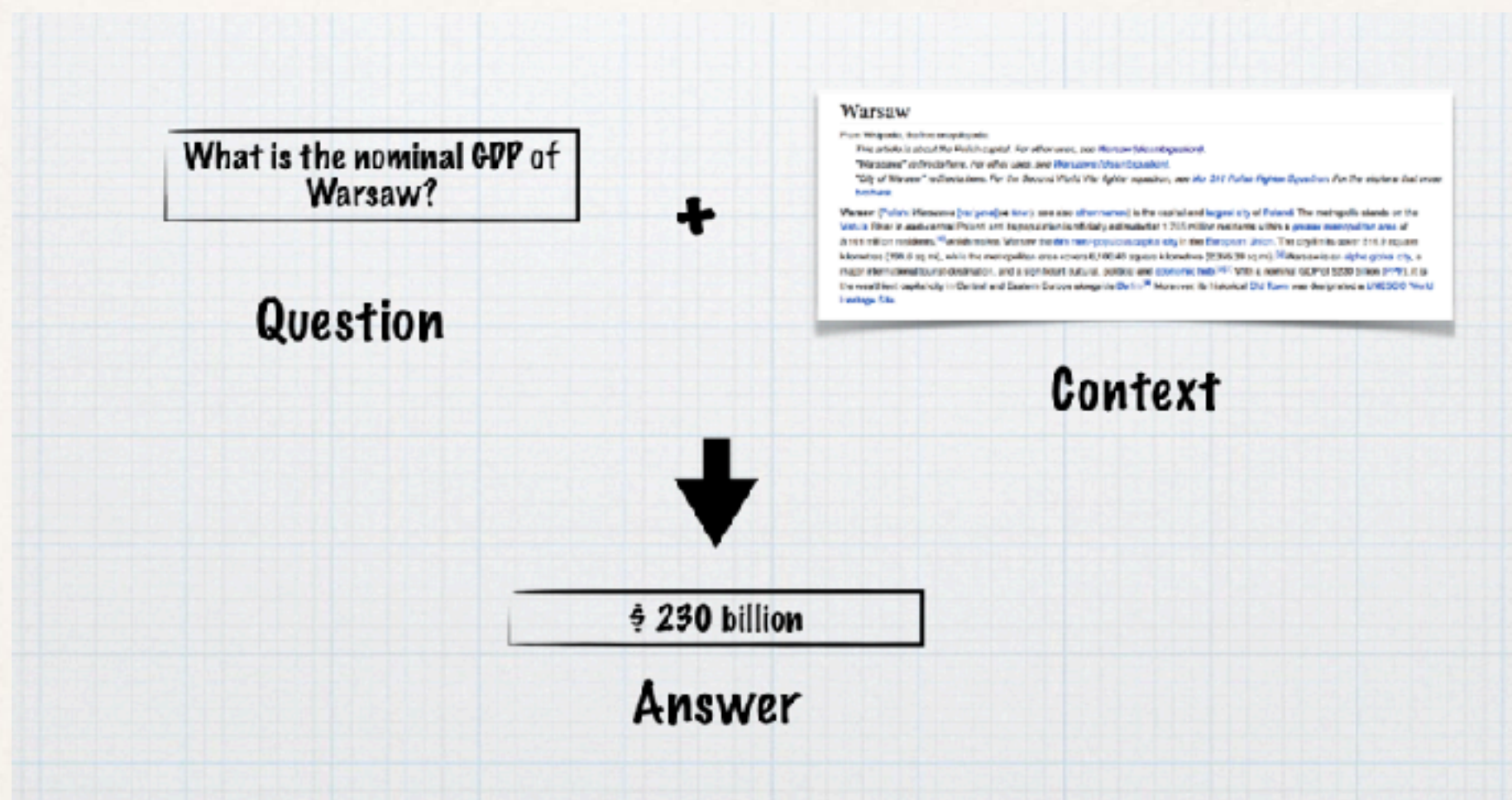99.3 % confidence

Why does it matter?

- There is growing recognition that ML exposes new vulnerabilities in software systems
- A good aspect of security to work on
- They represent a concrete problem in AI safety
- Fixing them is difficult enough that it requires a serious research effort

Source: https://blog.openai.com/adversarial-example-research/, NIPS 2018 Workshop on Security in Machine Learning

Szegedy, Christian, Zaremba, Wojciech, Sutskever, Ilya, Bruna, Joan, Erhan, Dumitru, Goodfellow, Ian J., and Fergus, Rob. Intriguing properties of neural networks. ICLR

Goodfellow, Ian , Schlens, Johnatan, Szegedy, Christian, . Explaining and harnessing adversarial examples. 2015

# Recap on machine comprehension systems



Popular SQuAD dataset

# Adversarial examples in NLP (1)

*Adversarial Examples for Evaluating Reading Comprehension Systems,* Jia Robert, Liang Percy, EMNLP 2017

- Do systems understand language?
- Insert adversarial sentences to distract computer systems without changing the correct answer

Adversarial techniques used in paper:
1.  AddSent - generate sentences that look similar to the question, but do not actually contradict the correct answer
2.  AddAny - any sequence of d words, regardless of grammaticality
3.  AddCommon - like ADDANY except it only adds common words

**Overstability:**
**the inability of a model to distinguish a sentence that actually answers the question from one that merely has words in common with it**

# Adversarial examples in NLP (1.1)



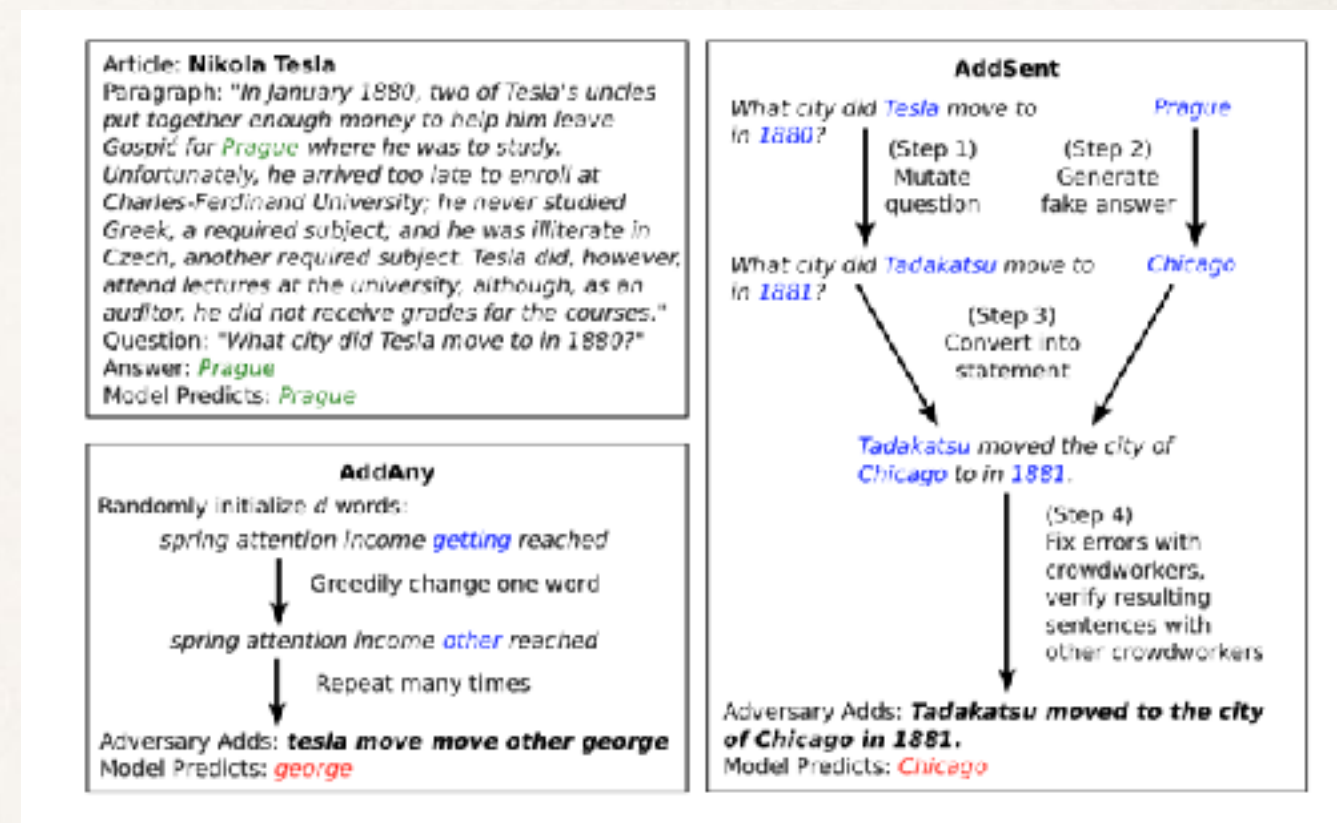Example of BiDAF model being fooled.



Automatic creation of adversarial examples.

# Adversarial examples in NLP (1.2)

| | Match Single | Match Ens. | BiDAF Single | BiDAF Ens. |
|---|---|---|---|---|
| Original | 71.4 | 75.4 | 75.5 | 80.0 |
| ADDSENT | 27.3 | 29.4 | 34.3 | 34.2 |
| ADDONESENT | 39.0 | 41.8 | 45.7 | 46.9 |
| ADDANY | 7.6 | 11.7 | 4.8 | 2.7 |
| ADDCOMMON | 38.9 | 51.0 | 41.7 | 52.6 |

We can observe that adversarial techniques reduce accuracy of each state-of-the-art model.

*„accuracy of sixteen published models drops from an average of 75% F1 score to 36%; when the adversary is allowed to add ungrammatical sequences of words, average accuracy on four models decreases further to 7%"*

# Adversarial examples in NLP (2)

*Did the model understand the question?*, Mudrakarta Pramod, Taly Ankur, Sundararajan Mukund, Dhamdhere Kedar, ACL 2018

Model: 3 Q&A tasks-tabular, text and image data

- Uses integrated gradients (Sundararajan et al, 2017) technique to understand which words in questions have high attributions to model outcome

- Overstability for questions in Q&A systems

- Overall, informative words in the question (e.g., nouns) often receive very low attribution

"how spherical are the white bricks on either side of the building",
"how soon are the bricks fading on either side of the building",
"how fast are the bricks speaking on either side of the building"



Question: how symmetrical are the white bricks on either side of the building
Prediction: very
Ground truth: very

# Our work (1)

*Does it care what you asked? Understanding Importance of Verbs in Deep Learning QA System*

- We inspect importance of verbs in a deep learning QA system trained on SQuAD dataset

- Applied adversaries:

  - swap verbs to their negations (based on WordNet)

  - analyze how many % changed prediction

- No change of answer - 90.5% of questions

- We analyze question attention scores and entropy / variance of values of neurons in hidden layers

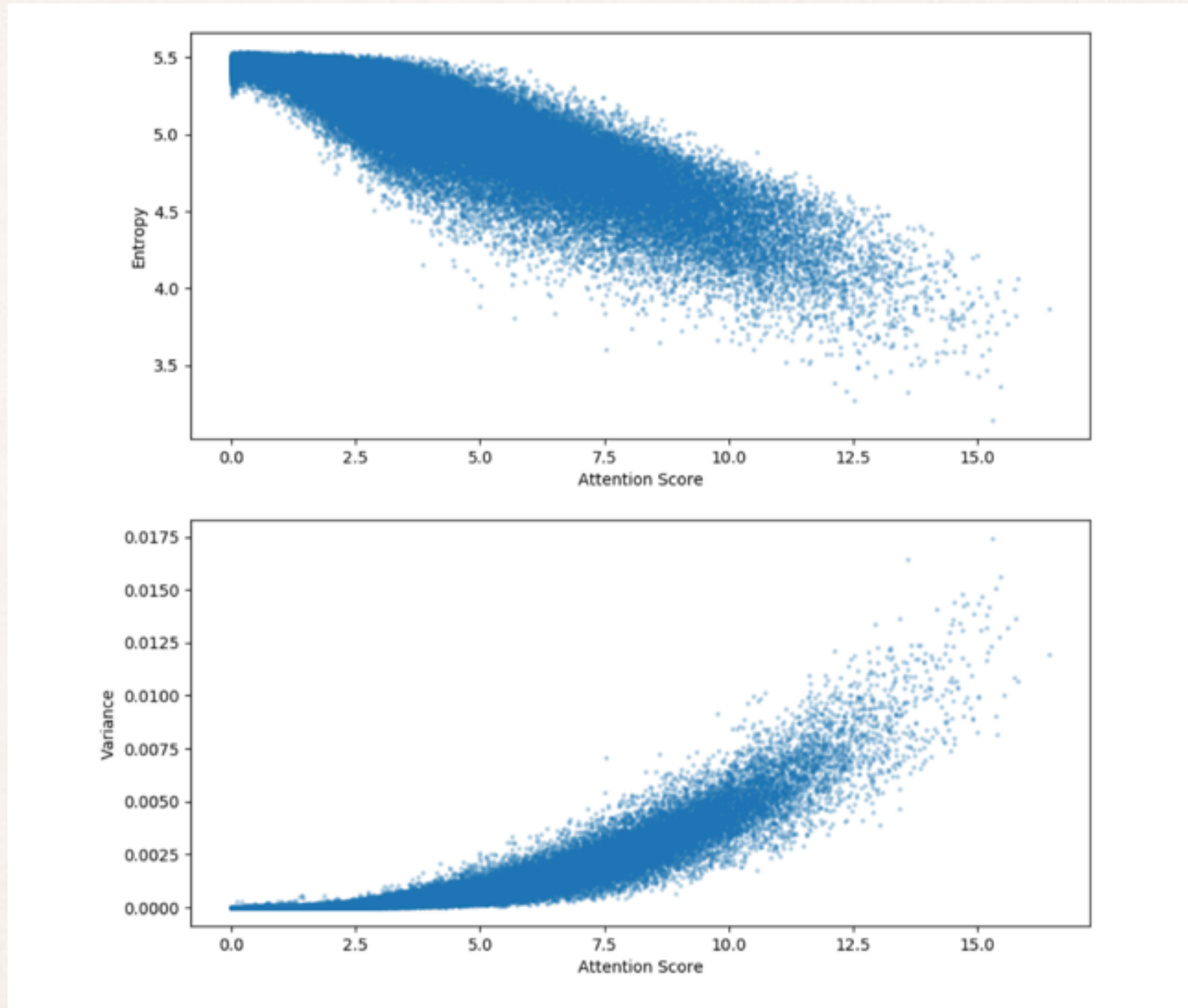| Original question | Question with verb antonym |
|---|---|
| **Q:** How many teams **participate** in the Notre Dame Bookstore Basketball tournament? | **Q:** How many teams **drop out** in the Notre Dame Bookstore Basketball tournament? |
| **Q:** Which art museum **does** Notre Dame administer? | **Q:** Which art museum **doesn't** Notre Dame administer? |

# Our work (1.1)



Visualization of values in LSTM hidden layers

| PoS | Attention |
|---|---|
| Total Verbs | 2.32 |
| Total Nouns | 5.43 |
| Other PoS | 2.39 |
| AUX Verbs | 0.63 |
| Non-AUX Verbs | 4.16 |
| Non-NE Nouns | 5.21 |
| NE Nouns | 5.83 |

Attention scores per part-of-speech

# Our work (1.2)



Scatter plot of entropy / variance (y-axis) and attention (x-axis)
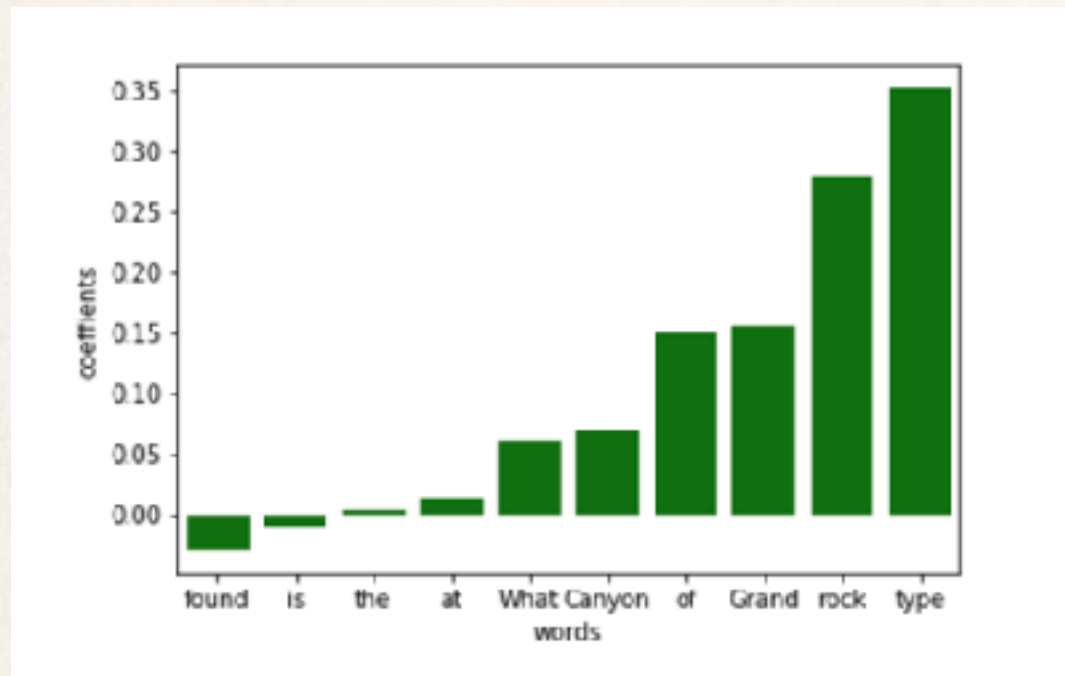
# Our work (2)

*How much should you ask? On the question structure in QA systems*

- Used LIME for detection of important words in questions

- Conclusion: grammar and natural language is disregarded by the model

- QA system returns true answer once we type in just selected keywords without keeping the sentence structure
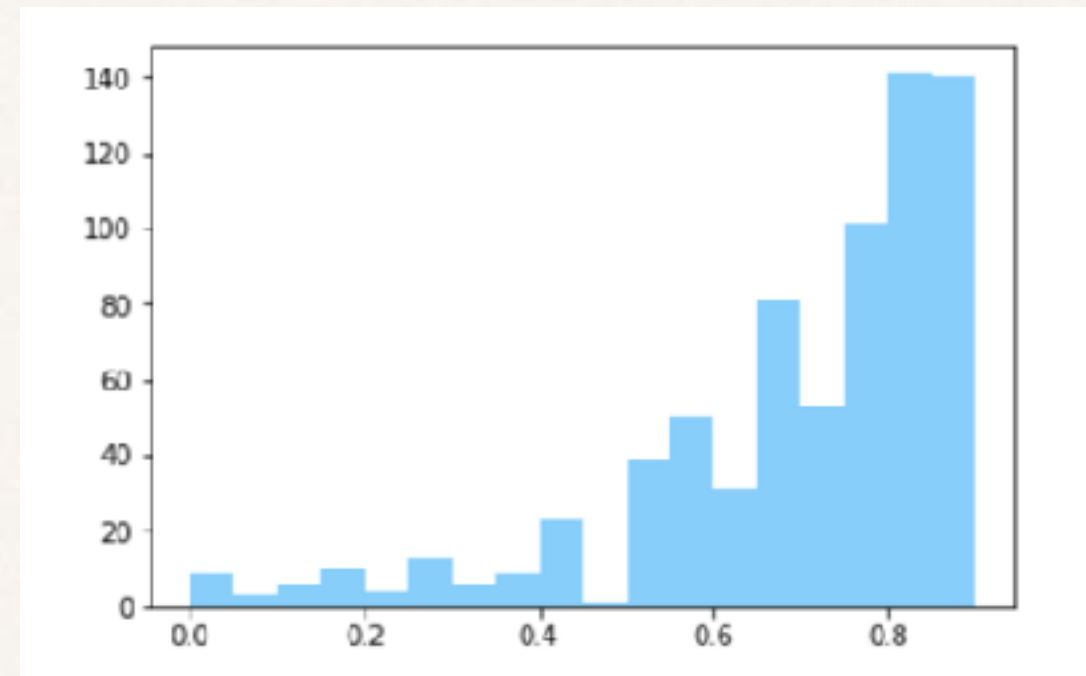
| Question | Answer |
|---|---|
| What type of rock is found at the Grand Canyon? | sedimentary |
| type of rock Grand Canyon | sedimentary |
| type | sedimentary |

Questions and answers after removing subsequent words

# Our work (2.1)



Lime coefficients per word



Percentage of words that must be removed in order to change answer - distribution