# Clickbait Spoiler Generation

Rajasvi Vinayak Sharma | A59012988 | rvsharma@ucsd.edu

## Abstract

- Clickbait spoiling aims at generating short texts that satisfy the curiosity induced by a clickbait post. This project is a derivative of Clickbait Challenge 2023 at SemEval 2023. This work aims to solve 2 subtasks from challenge i.e. spoiler classification (phrase, passage,multi-line) and generation.
- I propose a 2 stage approach where firstly a classification model is used to identify spoiler type then we use appropriate spoiler generating model.
- Main point being that we use separate models for different spoiler types i.e. one for phrase and another for passage type after classification stage.
- Primarily use Question Answering and Passage Retrieval/Ranking models for spoiler generation..
- This work only covers passage and phrase spoiler generation yet as multi line is still in process. Finally for QA models, I fine-tuned different Hugging face models like roberta, deberta, bert etc.

## Introduction

| Clickbait tweet | Spoiler |
|---|---|
| **Lifehacker** @lifehacker<br>How to keep your workout clothes from stinking: lifehac.kr/57YOuEZ | "washing [them]" |
| **New York Post** @nypost<br>Just how safe are NYC's water fountains? nyp.st/2yHSGnr | "The Post independently tested eight water fountains in New York City's most frequented parks, and found that all met or exceeded the state's guidelines for water quality." |
| **CNBC** @CNBC<br>A Harvard nutritionist and brain expert says she avoids these 5 foods that "weaken memory and focus." (via @CNBCMakeIt) cnb.cx/2TG6zeX | "1. Added sugar" [...]<br>"2. Fried foods" [...]<br>"3. High-glycemic-load carbohydrates" [...]<br>"4. Alcohol" [...]<br>"5. Nitrates" [...] |

- Clickbaits often don't contain enough relevant information and meant to have attractive titles to draw user's attention. They often are used as source to show advertisement or convey something obvious.
- Generally, clickbait spoilers can be of 3 types: phrase, passage, and multi-line step based.
- As seen from example, all 3 spoiler types have different text structure, length etc. which demands that we have different methods for generating them.
- Based on research paper provided by organizers of challenge, I worked on replicating solution for phrase and passage type spoiler generation till now.
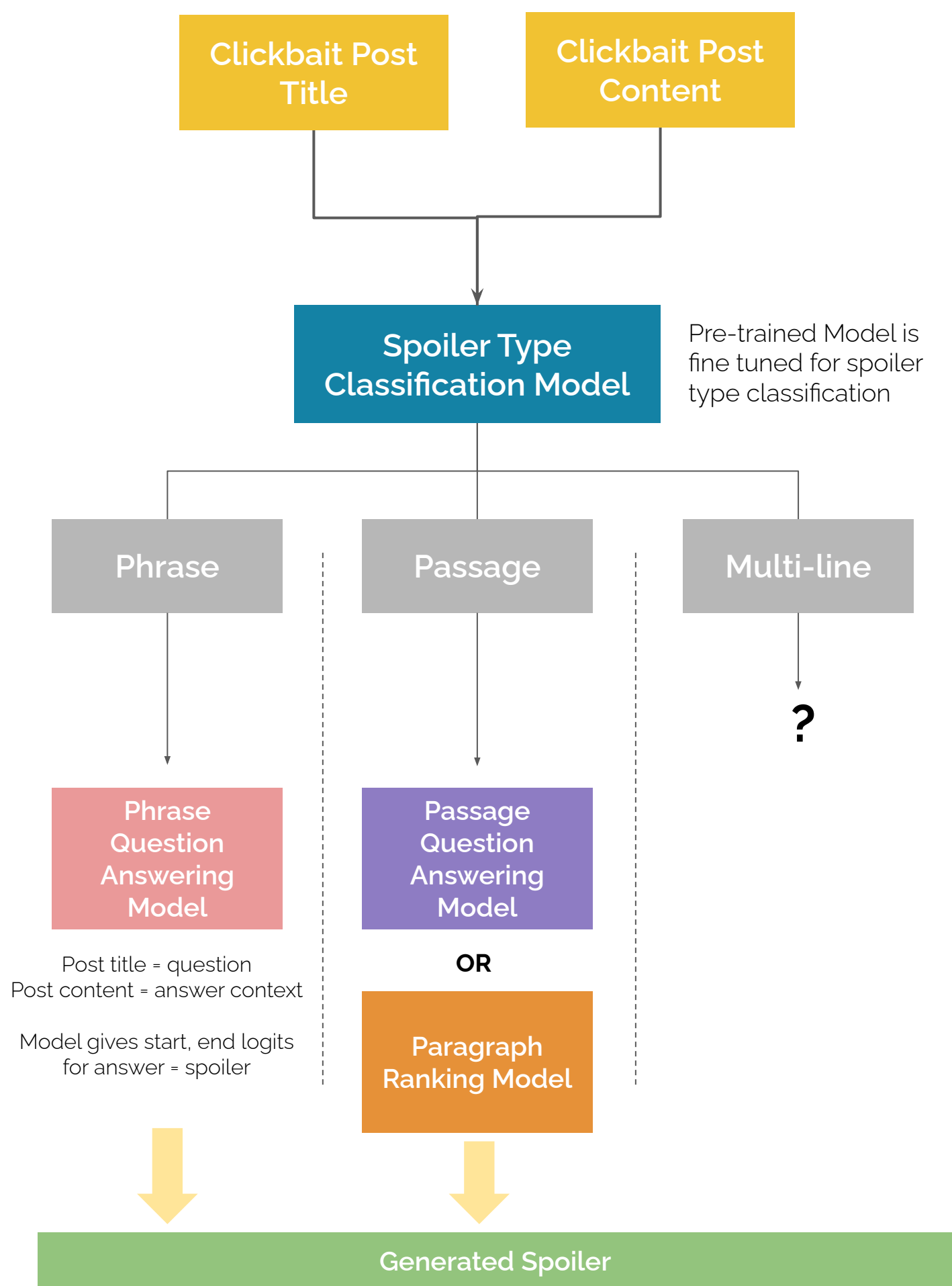
## Related Work

- Clickbait dataset was proposed in organizers paper which also mentioned the use of QA and Passage Retrieval for spoiler generation.
  https://webis.de/downloads/publications/papers/hagen_2022a.pdf
- Prior to this there was clickbait detection challenge in 2017.

## Dataset

- The dataset contains the clickbait posts and manually cleaned versions of the linked documents, and extracted spoilers for each clickbait post.
- Contains 3,200 posts for training and 800 posts for validation. Limited data per spoiler type.

## Methodology



## Results

**Clickbait Post Title:**
Watch as a romantic waterfall wedding proposal goes horribly wrong

**Generated Spoiler:**
the ring was gone

**Actual Spoiler:**
in a flash, the ring was gone, falling into the cold waters below

---

**Clickbait Post Title:**
Will Xbox One be getting Playstation 4.5 style upgrades?

**Generated Spoiler:**
he doesn't believe in the piece by piece upgrading of systems, and you won't be seeing that with the Xbox One

**Actual Spoiler:**
Xbox head Phil Spencer told Gamespot he doesn't believe in the piece by piece upgrading of systems, and you won't be seeing that with the Xbox One.

---

**Clickbait Post Title:**
Trouble in paradise? Rose Leslie doesn't want to work with Game Of Thrones beau Kit Harington again

**Generated Spoiler:**
Leslie wants to achieve her own success away from the gossip-mongers who feed off her high-profile relationship, which means working on individual projects and carving out a career on her own accord

**Actual Spoiler:**
wants to achieve her own success

---

**Clickbait Post Title:**
This is how much coffee Americans drinks every day

**Generated Spoiler:**
2.1 coffee drinks per day

**Actual Spoiler:**
2.1 coffee drinks

## Web Demo



Please contact owner for latest link

## Evaluation

Accuracy is primarily used for evaluation of spoiler type classification task.

| Model | Accuracy |
|---|---|
| bert-base-uncased | 0.66 |
| roberta-base | 0.75 |

BLEU-4, Exact Match and F1 score are primarily used for evaluation of phrase spoiler generation task.

| Spoiler Type | Model | BLEU-4 | F-1 | Exact match |
|---|---|---|---|---|
| Phrase | Roberta | 36.1 | 70.1 | 59.0 |
| Phrase | minilm | 36.3 | 68.4 | 56.7 |
| Passage | Roberta | 37.1 | 45.8 | 12.4 |
| Passage | minilm | 31.3 | 41.7 | 14.3 |

## Conclusion & Future work

- Main aim of this project is to generate spoiler given clickbait post title and article content.
- As a result, we implemented a 2 stage model approach for generating spoiler for phrase and passage based spoilers.
- Model crossed baseline for spoiler type classification but still not at par with spoiler generation.
- Initial experiments with paragraph retrieval model for passage type spoiler gave better results than equivalent baseline mentioned in original paper.
- Future work primarily includes devising method for generating multi-line spoiler.
- Currently, our approach is close to baseline but isn't better than the paper by absolute margin in some areas.

## References

[1] Matthias Hagen, Maik Fröbe, Artur Jurk, and Martin Potthast. Clickbait Spoiling via Question Answering and Passage Retrieval. In Smaranda Muresan, Preslav Nakov, and Aline Villavicencio, editors, 60th Annual Meeting of the Association for Computational Linguistics (ACL 2022), pages 7025-7036, May 2022. Association for Computational Linguistics.

[2] Martin Potthast, Tim Gollub, Matthias Hagen, and Benno Stein. The Clickbait Challenge 2017: Towards a Regression Model for Clickbait Strength. CoRR, abs/1812.10847, December 2018.

[3] Lidor Ivan, Shira Dvir Gvirsman, Mario Haim, and Martin Potthast. Don't Take the Bait: Users' Engagement with Clickbait and Its Effect on Editorial Considerations. In 71st Annual International Communication Association Conference (ICA 2021), May 2021.