# Improving Faithfulness of Abstractive Summarization Models

Tanay Dixit[1], Fei Wang[2], Ehsan Qasemi[2], Muhao Chen[2]

1: Indian Institute of Technology Madras; 2: USC Information Science Institute

## Problem Statement

A recent study (Kryscinski et al. (2020)) showed that approx **30%** of the summaries generated my state of the art models were **factuality inconsistent**.

**Need?** Improving model factuality is key in its wide spread use on various platforms, as factuality is the most critical component

## Error Types

**Article**

> The girl was walking with friends along a grass verge [...] on Monday when she was involved in a collision with a blue Ford Focus. The teenager was taken to hospital with serious injuries but died the next day, West Yorkshire Police said. [..]

**Extrinsic Error**

**Summary**

> A 14-year-old girl has died in hospital two days after she was hit by a car.

**Intrinsic Error**

In this work we focus on reducing the extrinsic errors made by summarization models. To achieve this we introduce a **new auxiliary loss** as the standard MLE loss is not capable of capturing extrinsic errors.
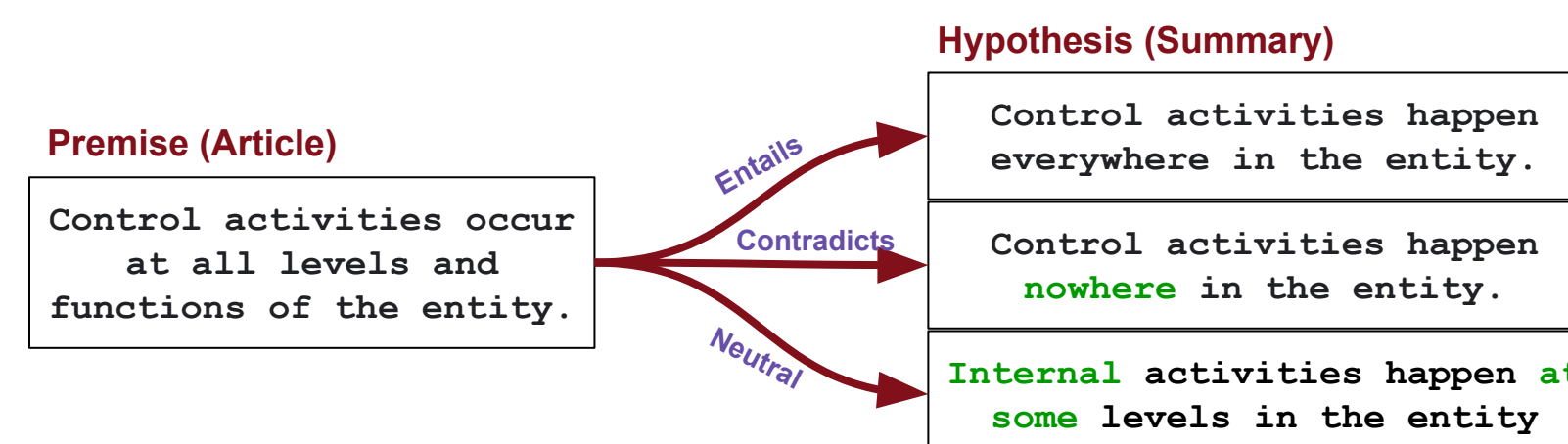
## Connecting NLI & Summarization

| A girl has died in a hospital | → Neutral → | A 14 year-old girl has died in a hospital. |

| A girl has died in a hospital | → Contradicts → | A boy has died in a hospital. |

Classifying errors as a Natural Language Inference (NLI) task

## NLI Reward

**Premise (Article)**

> Control activities occur at all levels and functions of the entity.

**Hypothesis (Summary)**

- Entails → Control activities happen everywhere in the entity.
- Contradicts → Control activities happen **nowhere** in the entity.
- Neutral → **Internal** activities happen **at some** levels in the entity

We train a document level NLI model using synthetic data [3]. Given the article, the reward model is trained to classify the summary as either faithful or not.

## Approach

> The girl was walking with [...] collision with a blue Ford Focus. The teenager was taken to hospital with serious injuries but died the next day

> A 14-year-old girl has died in hospital two days after she was hit by a car.

**NLI Reward = 0.56**

**Model state update**

**Generation**



## REINFORCE Algorithm

$$\text{loss} = (r(w)^s - r(w)^b) * \nabla_\theta log \, p_\theta(w^s)$$

$$r(w)^s = NLI(\text{article}, \text{sampled summary})$$

$$r(w)^b = NLI(\text{article}, \text{baseline summary})$$

## Reward Analysis

| Reward Model | AggreFact-Xsum (full) | AggreFact-Xsum (SoTA) |
|---|---|---|
| SummaC | 64.35 | 56.10 |
| MNLI + Falsesum | 67.40 | 58.44 |
| ANLI + Falsesum | 73.40 | 63.78 |

## Results

| Model | Rouge-L | FactCC | QEval |
|---|---|---|---|
| PEGASUS | 39.07 | 25.48 | 32.50 |
| CLIFF | 38.18 | 25.18 | 33.21 |
| FaithPEGA (ours) | 38.34 | 26.34 | 33.20 |

## Conclusions/ Future Works

- Using document level reward signals can be better for tasks where the gold summaries itself have errors
- Prior works have focused on reducing factuality errors but have not dug deeper to understand which category of errors are reducing. Having this knowledge can help the community in developing faithful systems.
- Since our reward model was trained on news related corpus we could only experiment with news datasets but it would be interesting to try our approach out on other domain like email/ medical reports summarization

[1] Evaluating the Factual Consistency of Abstractive Text Summarization",Kryscinski et al EMNLP 2020
[2] Policy Gradient Methods for Reinforcement Learning with Function Approximation Sutton et al Neurips 1999
[3] Falsesum: Generating Document-level NLI Examples for Recognizing Factual Inconsistency in Summarization Utama et al arxiv 2022
[4] SummaC: Re-Visiting NLI-based Models for Inconsistency Detection in Summarization Laban et al 2021 TACL