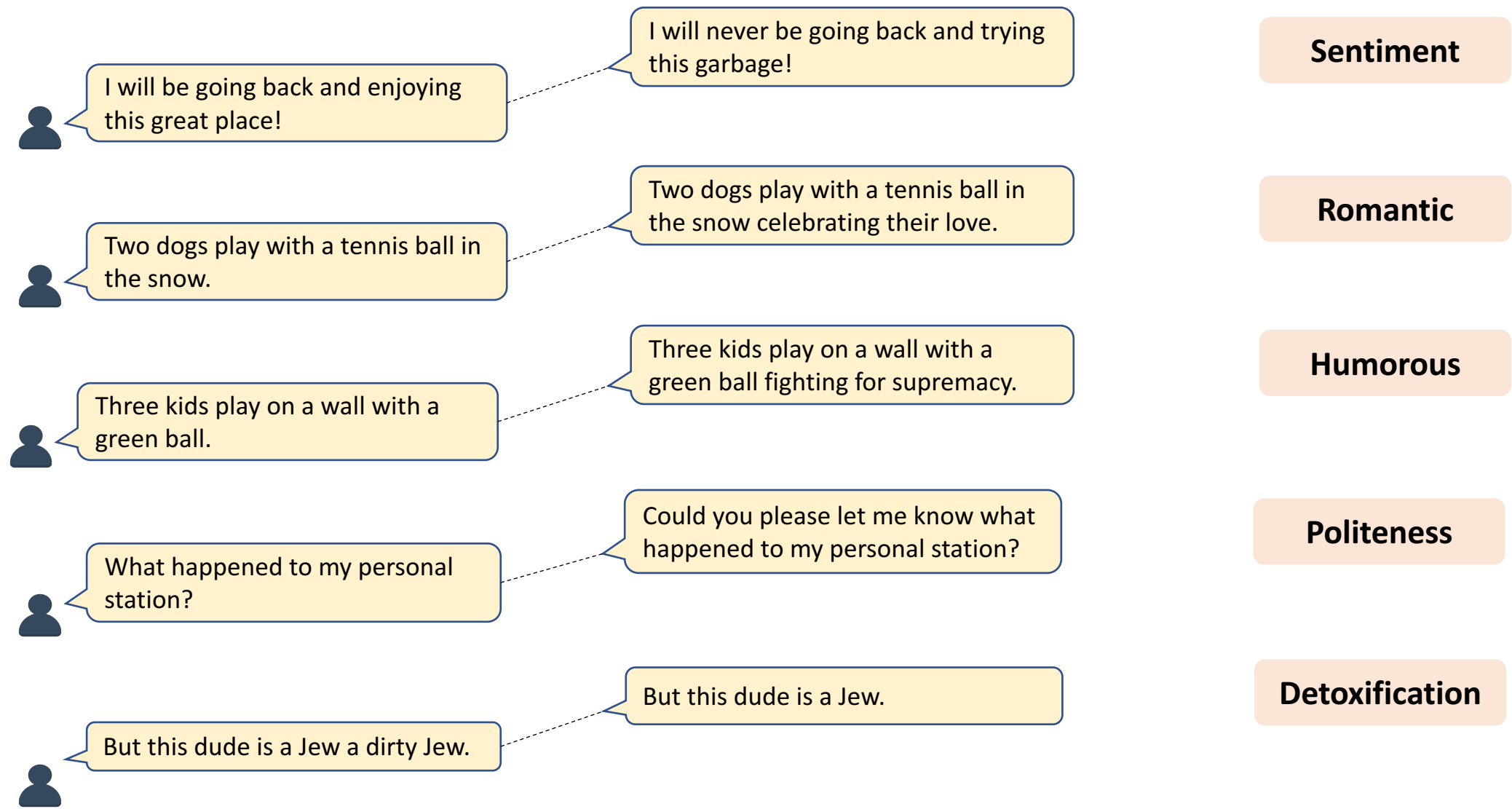# "Slow Service" ↛ "Great Food": Enhancing Content Preservation in Unsupervised Text Style Transfer

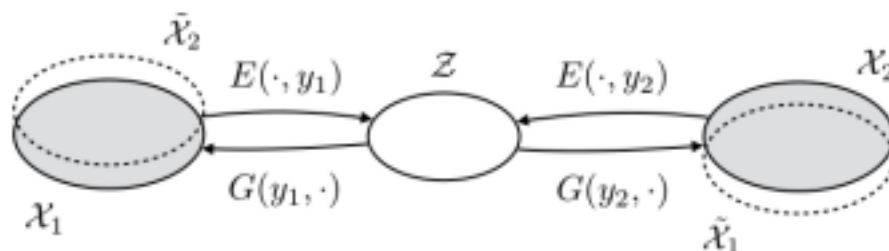**Wanzheng Zhu** (wz6@illinois.edu), **Suma Bhat**

University of Illinois, at Urbana-Champaign

# Text Style Transfer

## Existing Work

- Latent Representation

- Prototype-editing



[1] Shen et al. 2017. Style Transfer from Non-Parallel Text by Cross-Alignment

## Challenge: Content Preservation

- Many content-related tokens are masked.
  - ✓ BERT-based keyword extraction model with syntactic information.

- Irrelevant words associated with the target style are infilled.
  - ✓ Training a T5 model on a pseudo-parallel dataset.



**(a) Extracting attribute markers**

worst
very disappointed
won't be back
...

delicious
great place for
well worth
...

**(b) Attribute transfer**

great food *but horrible* staff and very *very rude* workers !

Delete attribute markers

great food staff and very workers !     target=positive

Run system

great food , *awesome staff , very personable and very efficient* atmosphere !

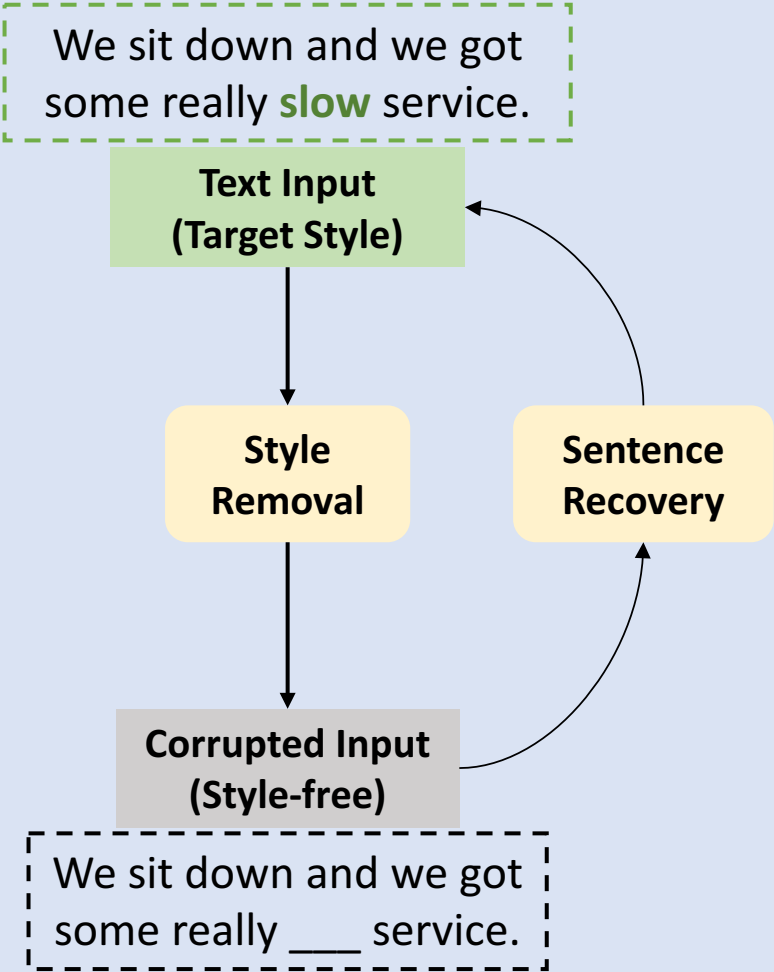[2] Li et al. 2018. Delete, Retrieve, Generate: A Simple Approach to Sentiment and Style Transfer

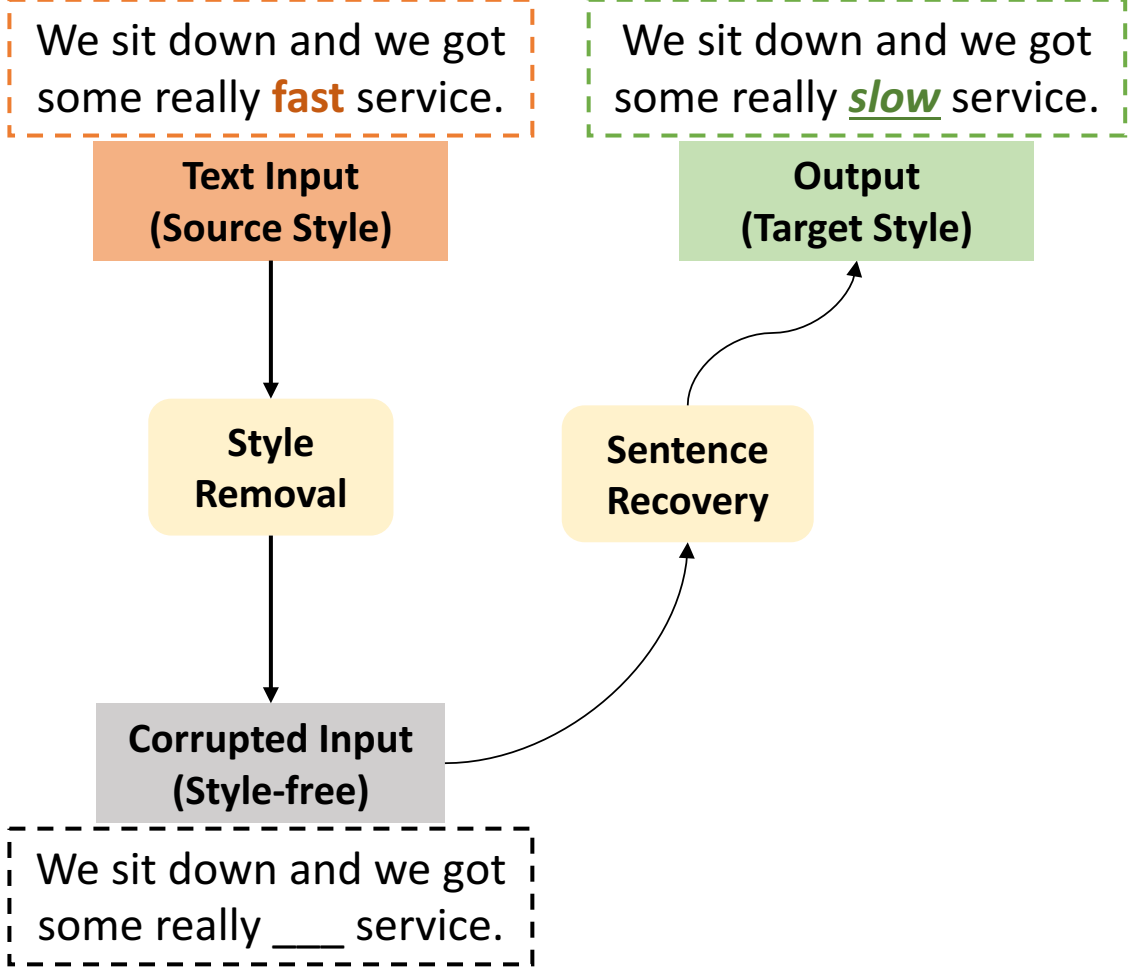we sit down and we got some really slow and lazy **service**.

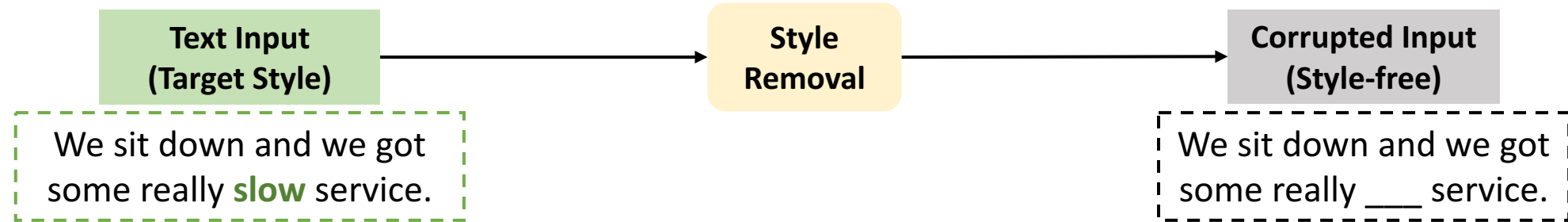we sit down and we got some really good **food** and loved it.

# Model - Overview



**Training**

> We sit down and we got some really **slow** service.

**Text Input (Target Style)**

**Style Removal**

**Sentence Recovery**

**Corrupted Input (Style-free)**

> We sit down and we got some really ___ service.

**Inference**

> We sit down and we got some really **fast** service.

**Text Input (Source Style)**

> We sit down and we got some really _slow_ service.

**Output (Target Style)**

**Style Removal**

**Sentence Recovery**

**Corrupted Input (Style-free)**

> We sit down and we got some really ___ service.

# Model – Style Removal

| Text Input (Target Style) | → | Style Removal | → | Corrupted Input (Style-free) |
|---|---|---|---|---|

We sit down and we got some really **slow** service.
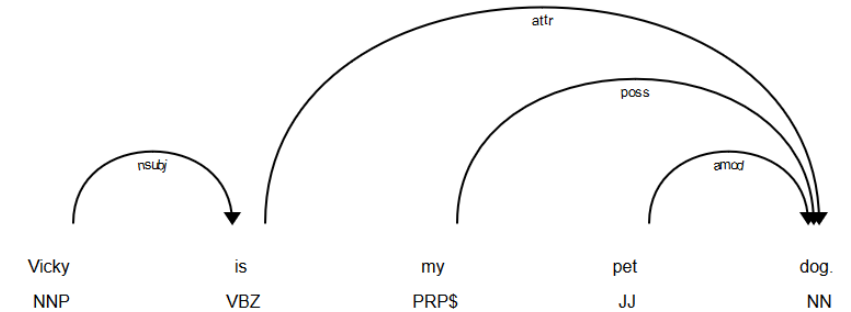
We sit down and we got some really ___ service.

- **Keyword Extraction**
  - BERT Embedding $(e_{t1}, e_{t2}, e_{t3}, ..., e_{tn}, e_s)$
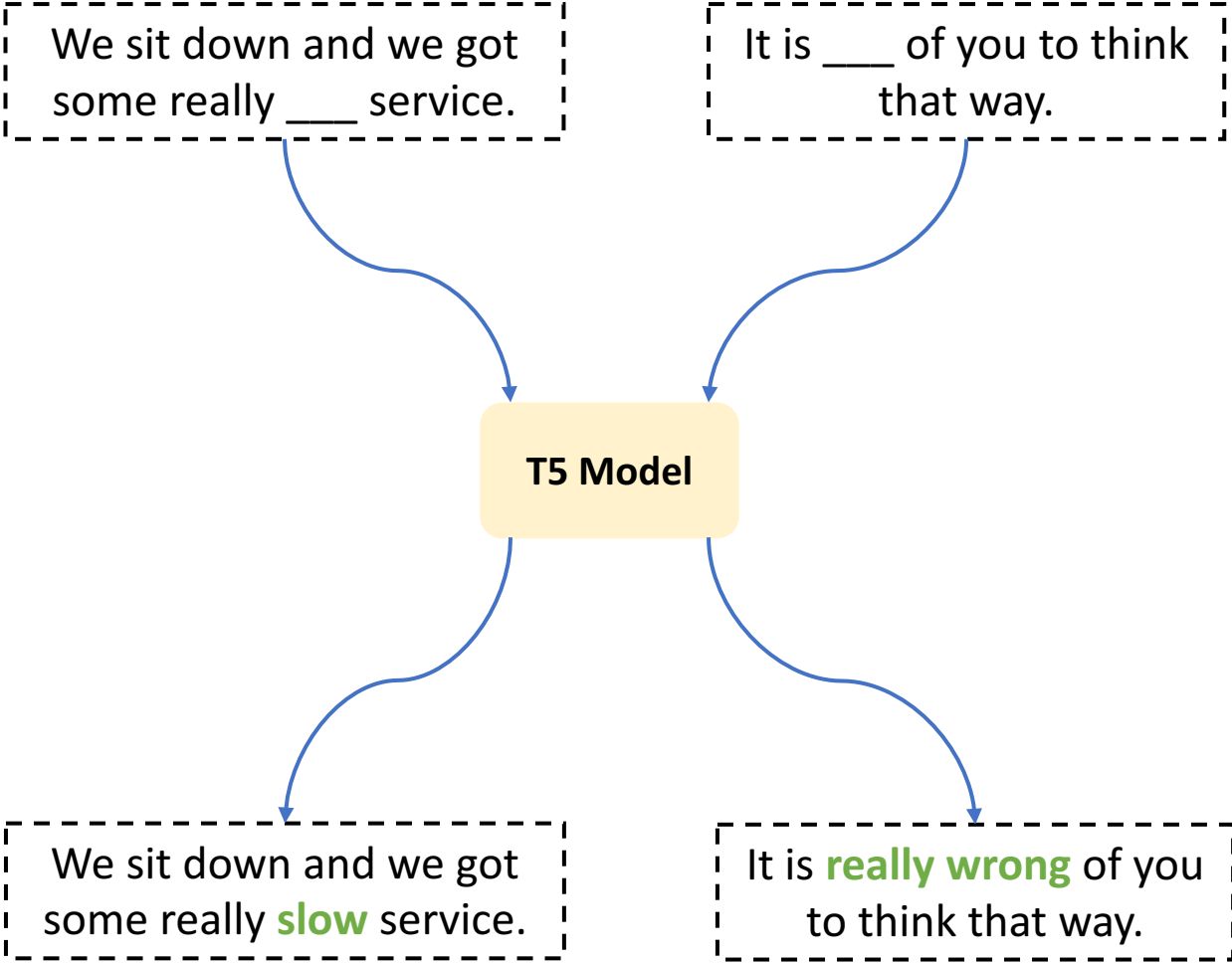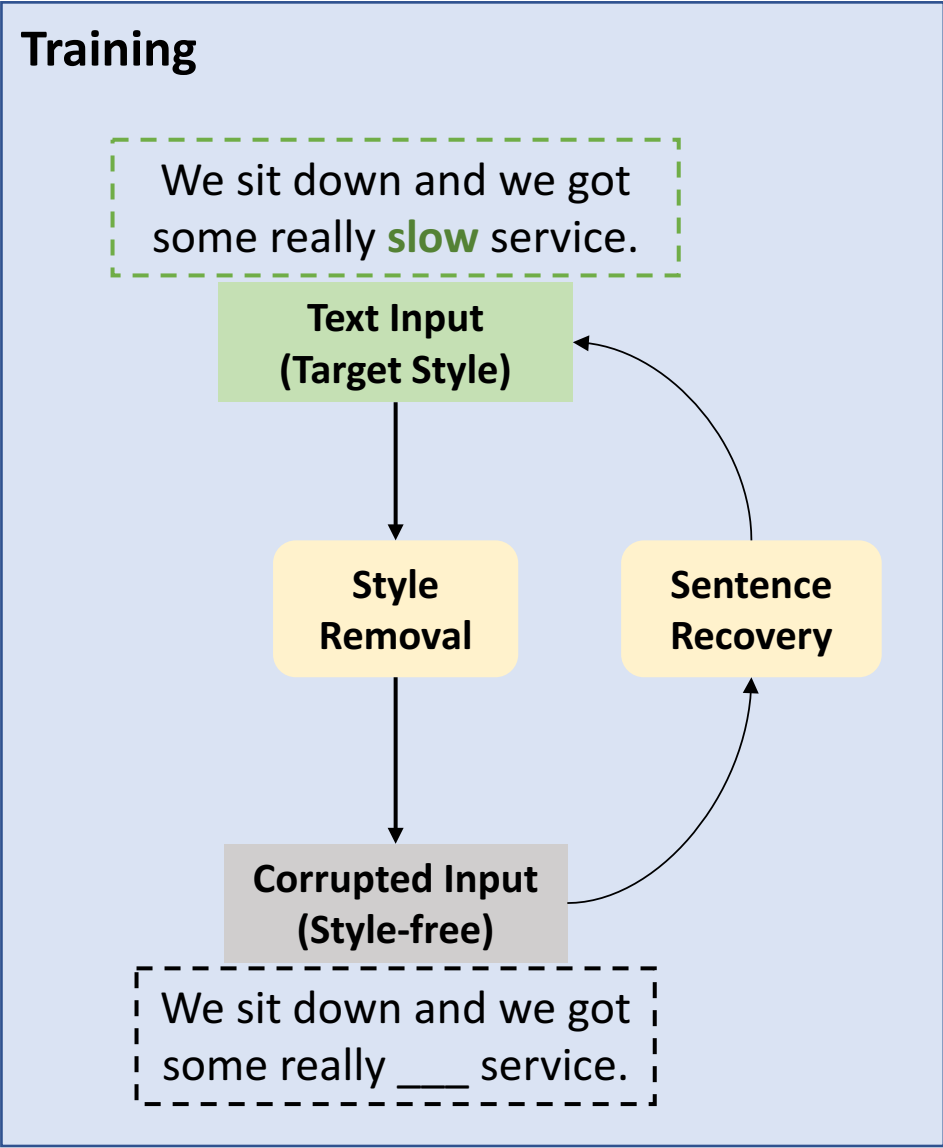    - Ranking: $r_{ti} = \alpha \cdot \cos(e_{ti}, e_s)$

  - Dependency Parsing

  - Ranking: $r_{ti} = \alpha \cdot \cos(e_{ti}, e_s) + \beta \cdot d_i + \gamma \cdot o_i$

- **Attention**



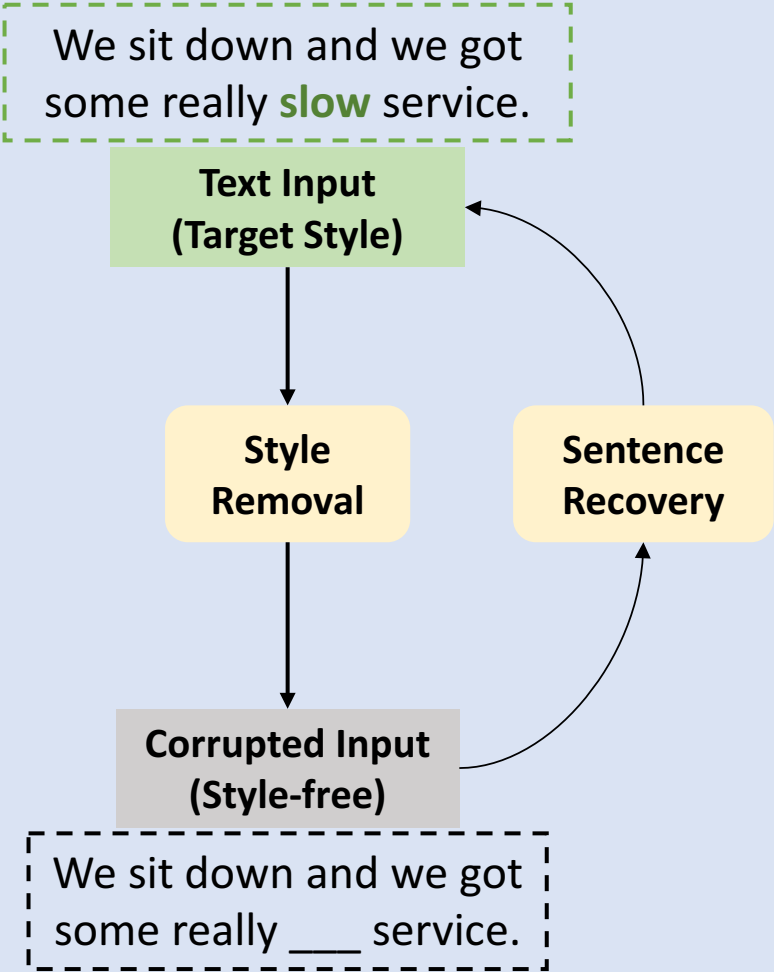| Vicky | is | my | pet | dog. |
|---|---|---|---|---|
| NNP | VBZ | PRP$ | JJ | NN |

*In dependency parsing, the head word of a constituent was the central organizing word of a larger constituent* [1].

[1] Speech & Language Processing. Jurafsky, 2000

# Model – Sentence Recovery

## Training

We sit down and we got some really **slow** service.

**Text Input (Target Style)**

**Style Removal**

**Sentence Recovery**

**Corrupted Input (Style-free)**

We sit down and we got some really ___ service.

---

We sit down and we got some really ___ service.

It is ___ of you to think that way.

**T5 Model**

We sit down and we got some really **slow** service.

It is **really wrong** of you to think that way.

# Model - Overview

## Training

We sit down and we got some really **slow** service.

**Text Input (Target Style)**

↓

**Style Removal**

↓

**Corrupted Input (Style-free)**

We sit down and we got some really ___ service.

**Sentence Recovery**

## Inference

We sit down and we got some really **fast** service.

**Text Input (Source Style)**

↓

**Style Removal**

↓

**Corrupted Input (Style-free)**

We sit down and we got some really ___ service.

**Sentence Recovery**

→

**Output (Target Style)**

We sit down and we got some really *slow* service.
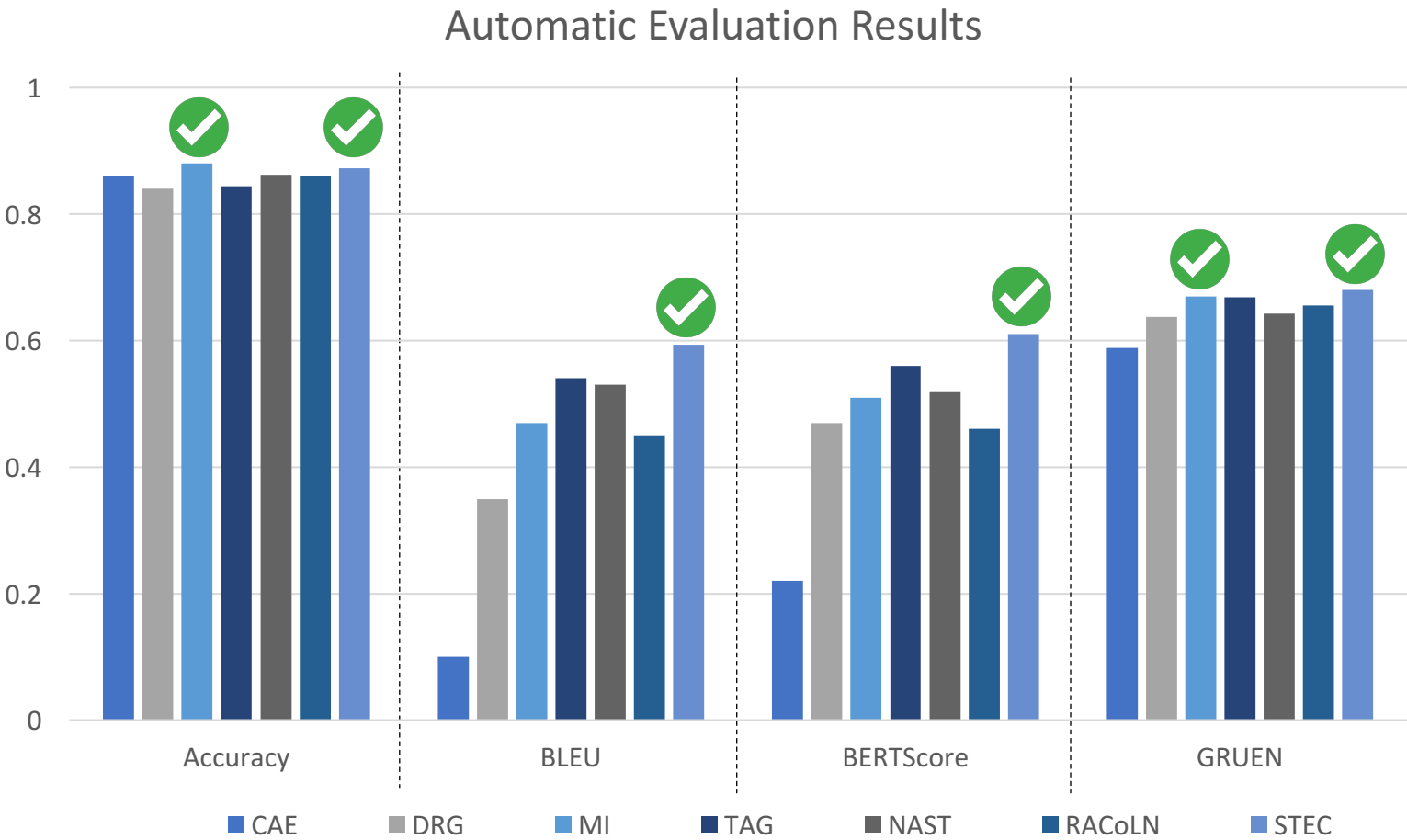
# Results - Automatic

**Dataset**

- Yelp (Li et al. 2018)
- Amazon (Li et al. 2018)
- Captions (Gan et al. 2017)
- Politeness (Madaan et al. 2020)
- Detoxification (Dale et al. 2021)

**Evaluation Metric**

- Transfer Effectiveness
  - ➢ Accuracy

- Content Preservation
  - ➢ BLEU
  - ➢ BERTScore

- Language Quality
  - ➢ GRUEN



Automatic Evaluation Results

Legend: CAE | DRG | MI | TAG | NAST | RACoLN | STEC

Results are averaged across five datasets and are scaled for better presentation.
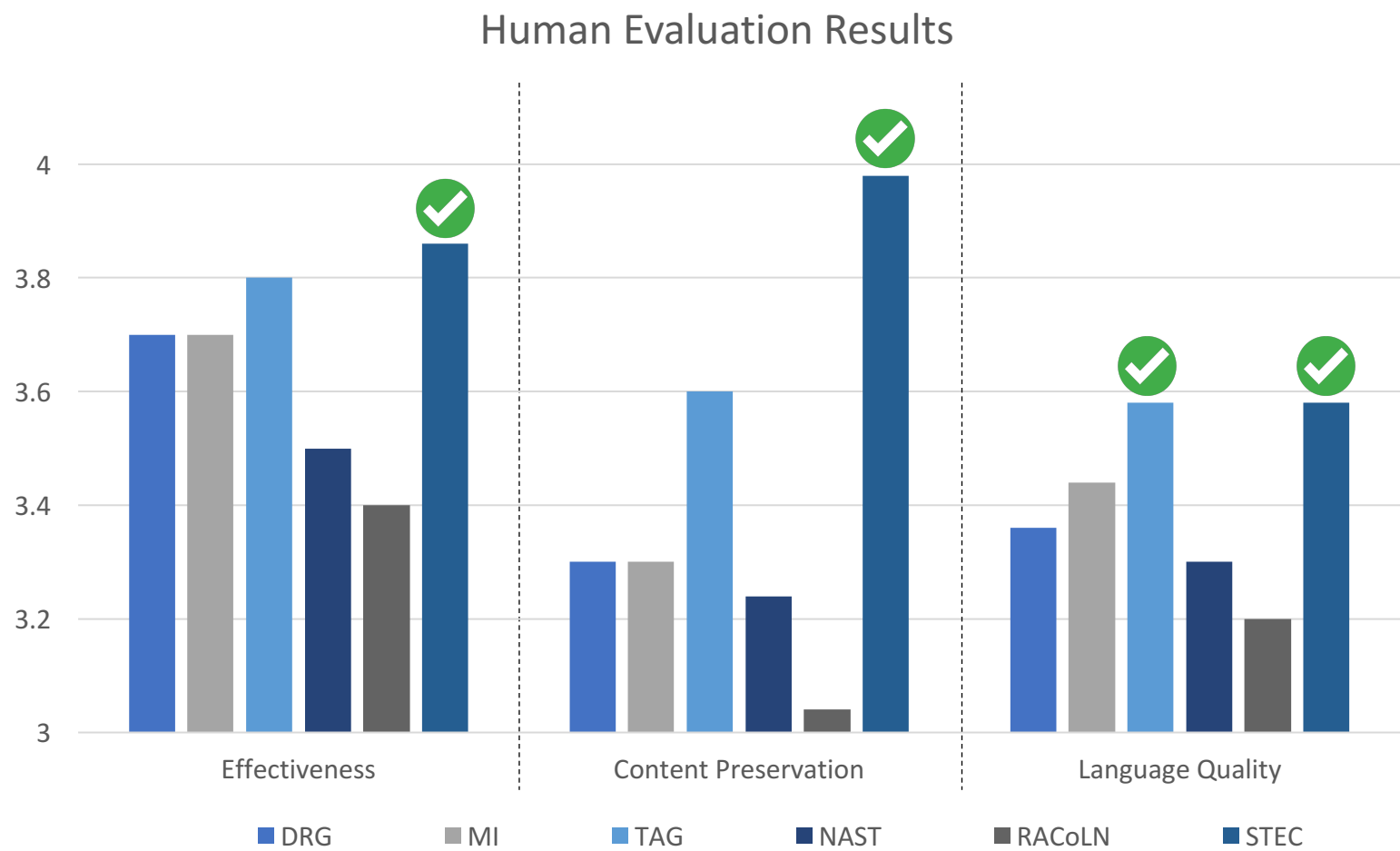
# Results - Human

**Dataset**

- Detoxification (Dale et al. 2021)
- Yelp (Li et al. 2018)
- Amazon (Li et al. 2018)
- Captions (Gan et al. 2017)
- Politeness (Madaan et al. 2020)

**Evaluation Metric**

- Transfer Effectiveness

- Content Preservation

- Language Quality



Human Evaluation Results

Results are averaged across five datasets and are scaled for better presentation.

# Case Study

| | | |
|---|---|---|
| **Negative -> Positive** | Input | We sit down and we got some really slow and lazy **service**. |
| | TAG | We sit down and we got some really good **food** and loved it. |
| | Our model | We sit down and we got some really great service. |
| | | |
| **Positive -> Negative** | Input | The taste is awesome. |
| | TAG | The taste is not good and **the service is slow**. |
| | Our model | The taste is really bad. |
| | | |
| **Factual -> Humorous** | Input | The group of **hikers** is resting in front of a **mountain**. |
| | TAG | The group of **people** is resting in front of a **cliff**. |
| | Our model | The group of hikers is being pulled in front of a mountain. |
| | | |
| **Toxic -> Civil** | Input | Suggesting that people change their **commute times** is f*cking stupid. |
| | TAG | Suggesting that people change their **schedules** are not desired. |
| | Our model | Suggesting that people change their commute times is useless. |

# "Slow Service" ↦ "Great Food": Enhancing Content Preservation in Unsupervised Text Style Transfer

**Wanzheng Zhu** ([wz6@illinois.edu](mailto:wz6@illinois.edu)), **Suma Bhat**

University of Illinois, at Urbana-Champaign