

# BART for Span Detection and Classification

Cees Roele

## Objective

Simultaneous span detection and classification:

- Detect which spans in a text contain usage of a propaganda technique
- Identify which out of 20 manipulative techniques are used in a span

## Introduction

Given only the textual content of a meme, detect spans containing manipulative techniques and identify which of 20 techniques are used in each span. This multilabel sequence tagging task is subtask 2 of SemEval-2021 task 6[1].

Spans can overlap and can range over multiple sentences, as shown in figure 1 below. The ellipsis at the end stands for continuation of the second sentence.

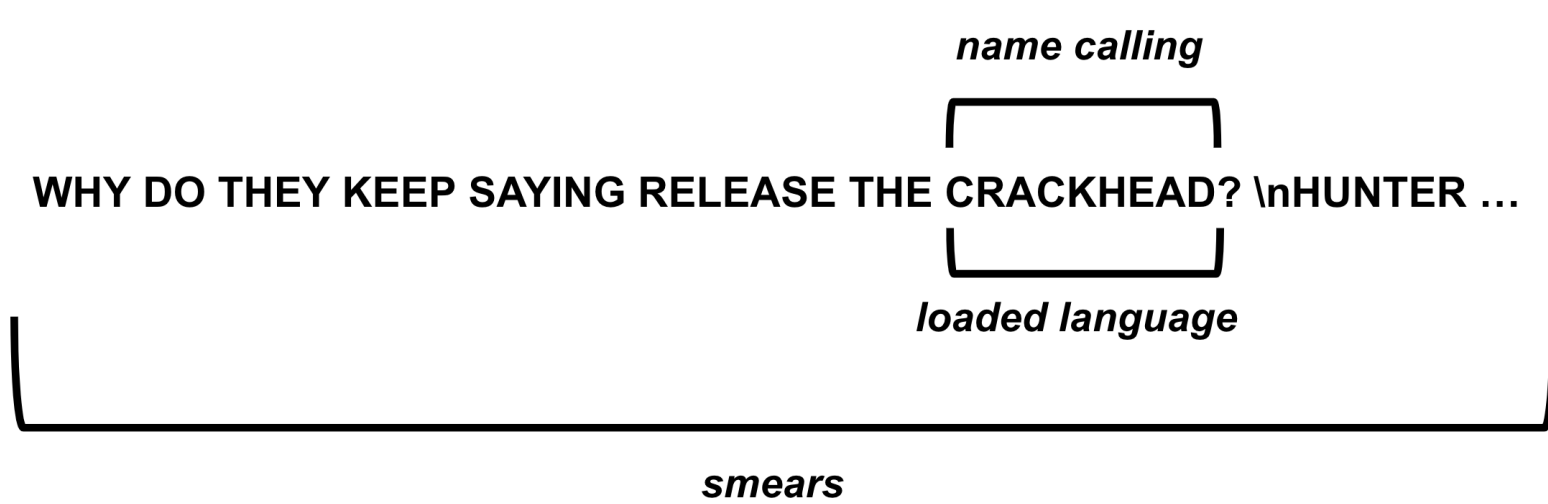


Figure 1: Overlapping spans and spans extending over multiple sentences

The solution chosen here and participating in task 6 as WVOQ is to *use a BART EncoderDecoder model to re-generate the input sequence with added XML-like tags to identify both spans and classes*.

## Sequence Generation with BART

BART is a denoising autoencoder built with a sequence-to-sequence model. [2]. BART uses a standard Transformer-based neural machine translation architecture to couple a bidirectional encoder with a left-to-right decoder. Pre-training BART was done by first corrupting text with an arbitrary noising function and then training a sequence-to-sequence model to reconstruct the original text.

Configuration parameters can tweak repetition, maximum output length, length penalty, and search algorithm.

## Generated Outcome

We generate the original sentence with XML-like markup to indicate spans and labels. Below the fragments from the example in figure 1 in the Introduction are labeled as *smears*, *loaded-language*, and *name-calling*.

```
<SMEARS>
WHY DO THEY KEEP SAYING RELEASE THE
<LOADED-LANGUAGE>
<NAME-CALLING>
CRACKHEAD
</NAME-CALLING>
</LOADED-LANGUAGE>
? \n HUNTER ...
</SMEARS>
```

This marked outcome defines both identification of spans and their classification.

## Conclusion

The described system offers a proof-of-concept for a novel approach to sequence tagging based on generating a version of a message with markup for labels. Drawback: *systemic errors* cannot be resolved through configuration of the standard system.

## Future Research

Resolve systemic errors:

- Extend decoder's generation algorithm to support matching of begin and end tags
- Adapt loss function to give priority to output tokens that are also in the input - with the exception of tags

## References

- [1] Dimitar Dimitrov, Bishr Bin Ali, Shaden Shaar, Firoj Alam, Fabrizio Silvestri, Hamed Firooz, Preslav Nakov, and Giovanni Da San Martino. Task 6 at SemEval-2021: Detection of persuasion techniques in texts and images. In *Proceedings of the 15th International Workshop on Semantic Evaluation, SemEval '21*, Bangkok, Thailand, 2021.
- [2] Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer Levy, Veselin Stoyanov, and Luke Zettlemoyer. BART: Denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 7871–7880, Online, July 2020. Association for Computational Linguistics.

## Important Results

- XML-like tag-generation for spans and classes works as a proof-of-concept
- The resulting score suffers from systemic errors which cannot be significantly decreased using standard configuration parameters for sequence generation

## BART Pre-Training

BART pre-training methodology includes multiple operations which are based on spans, rather than on single tokens like BERT.[2]

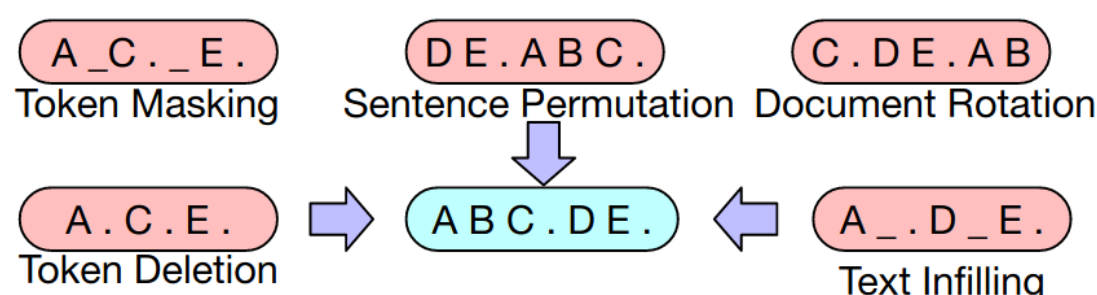


Figure 2: BART pre-training operations

Note:

- Span-based pre-training *might* help for span-detection.
- Pre-training doesn't cover matching of begin and end XML-like tags

## Results

The discussed WVOQ system's F1 score of 0.268 on subtask 2 on the test set scores about in the middle between the baseline and the highest ranking score.

Rank	Team	F1 score	Precision	Recall
1	Volta	.482	.501	.464
5	WVOQ	.268	.243	.299
	baseline	.010	.034	.006

Table 1: Subtask 2 scores on the test set

We observe the following systemic errors:

- Beginning and end tags in the generated text regularly don't match.
- Generated text contains changed words and even added words, which leads to faulty identifications of spans.

## Contact Information

- LinkedIn: <https://www.linkedin.com/in/ceesroele/>
- Email: cees.roele@gmail.com