# ENTITY-LEVEL FACTUAL CONSISTENCY OF ABSTRACTIVE TEXT SUMMARIZATION

## GROUP NO : 05

CS/2017/041 – SIRIWARDHANA B.D.S.A

CS/2017/020 – MUTHUKUMARA M.G.L.

# Entity-level Factual Consistency of Abstractive Text Summarization

Ensuring factual consistency of the generated summary with respect to the source document is a major challenge in abstractive text summarization. Recent researches reveal about 30% of the summaries generated by neural network sequence-to-sequence models have factual inconsistency. Factual inconsistency can happen at either the entity level or the relation level. This research paper focuses on entity-level factual consistency issue which is known as entity hallucination problem. An entity hallucination problem is named entities are that never contained in the source document that may be visible in the model generated summary. Although there is a widely used factual consistency metric which is known as the ROUGE score is inadequate to quantify factual consistency effectively. Three metrics are suggested to quantify factual consistency at the entity level in this research paper. The factual quality of summaries is analysed that is produced by the state-of-the-art BART model on three news data sets such as Newsroom, CNN/DailyMail and XSUM in this research. BART is a denoising autoencoder developed with a sequence-to-sequence model. Comparing to encoder-only pre-training such as in BERT or decoder-only pre-training such as in GPT-2, BART is an encoder-decoder transformer-based neural translation model that has been included five pre-taring techniques such as token masking, sentence permutation, document rotation, token deletion, and text infilling

Then certain strategies such as data filtering, multi-task learning and joint sequence generation are used to enhance performance on these metrics. The three metrics are proposed in this research are Precision-source($prec_s$), Entity-based data filtering and Precision-target and recall-target. Also, these metrics rely on the Spacy tool to perform Named-Entity Recognition (NER).

Precision-source quantifies the degree of hallucination with respect to the source document.

$prec_s$ = N (h ∩ s)/N(h)

N (h ∩ s) is the number of entities in the target summary that can detect a match in the source document.

N(h) is named entities in the generated summary.

Low $prec_s$ means high hallucination.

## Entity-based Data Filtering

If any of the entities in the generated summary which are not appear in the source document, The specific sentence that contains the entity is discarded. If there is only one sentence in the ground truth summary, it is required to be discarded and the document-summary pair is removed from the data set. Filtered data set that does not contain hallucination of entities ($prec_s$ = 1) is ensured from this technique. Here is the example below to the respective data set.

| | Newsroom | | | CNNDM | | | XSUM | | |
|---|---|---|---|---|---|---|---|---|---|
| | train | val | test | train | val | test | train | val | test |
| original | 922,500 (1.58) | 100,968 (1.60) | 100,933 (1.59) | 287,112 (3.90) | 13,368 (4.13) | 11,490 (3.92) | 203,540 (1.0) | 11,301 (1.0) | 11,299 (1.0) |
| after filtering | 855,975 (1.62) | 93,678 (1.64) | 93,486 (1.64) | 286,791 (3.77) | 13,350 (3.99) | 11,483 (3.77) | 135,155 (1.0) | 7,639 (1.0) | 7,574 (1.0) |

## Precision-target ($prec_t$) and Recall-target ($recall_t$)

Precision target ($prec_t$) and Recall-target ($recall_t$) is used to obtain a complete picture of entity-level accuracy of the generated summary since the precision source($prec_s$) metric do not capture the entity level accuracy of the generated summary. precision source quantifies the degree of entity hallucination concerning the source document.

$prec_t$ = N(h ∩ t)/N(h)

$recall_t$ = N(h∩t)/N(t)

N(t) is the number of named entities in the target summary

N(h ∩ t) is the number of named entities in the generated summary that can find a match in the ground truth summary.
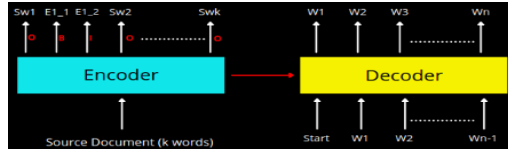
The F1 score can be computed as $F1_t = 2 \cdot prec_t \cdot recall_t /( prec_t + recall_t)$ to have a single quantifiable number.

## Multi-Task Learning

Also, furthermore an enhancement to the entity based data filtering is explored which is a way of classifying summary-worthy named entities in the source document (encoder side). A summary worthy named entity from the source that is also present in the ground truth summary is a salient entity. To achieve this, they label every token in the original document with a BIO scheme (B-Begin, I-Inside, O-outside) This is a pretty standard technique for labelling single/multi-word entities in a text segment.

BOI scheme can be represented as mathematically below,

$$\mathcal{L}_{\text{BIO}}^i(\theta(\text{enc}), x^i, z^i) = -\sum_{t=1}^{ts(i)} \log p_{\theta(\text{enc})}(z_t^i | x^i)$$

**Multi-task Learning**

In the past, it has been shown that if a decoder was not able to model an entity's representation, it would lead to hallucinations. The idea behind this loss is that this loss would force the encoder to model its representation so that it captures relevant information about summary-worthy entities. When the decoder gets this entity enriched representation, it can be built more accurate generations with lesser hallucinations.
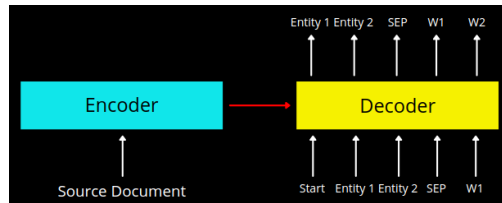
Apart from the BIO loss, they also use typical MLE loss for the training sequence generation model which would propagate loss from the decoder end.

$$\mathcal{L}_{\mathrm{MLE}}^{i}(\theta, x^i, y^i) = -\sum_{t-1}^{tt(i)} \log p_\theta(y_t^i | x^i, y_{<t}^i)$$

Here, theta, x, y, i are model parameters, input tokens, output tokens and ith token respectively. They minimize the joint loss L(i) = L(i)_MLE + αL(i)_ BIO, where a is a hyperparameter. They choose α between 0.1 to 0.5 based on the validation set.

## Join Salient Entity and Summary Generation (JAENS)

The researchers have developed a new way of generating summaries, where they train a sequence model to generate the sequence of summary-worthy named entities, followed by a special token, and then the summary tokens instead of just generating the summary. The idea behind this approach is that while generating summary tokens, the decoder can attend to summary-worthy entities via the self-attention mechanism. This allows it to jointly learn to identify the descriptive entities while learning to generate summaries.


**JAENS Architecture**

## Results

The table below reveals proposed exploration is great for hallucination alleviation. It demonstrates comparing models trained using original data with entity-based data filtering, with an additional classification task and with JAENS.

| | training data | Rouge1 | Rouge2 | RougeL | macro $prec_s$ | micro $prec_s$ | macro $prec_t$ | micro $prec_t$ | macro $recall_t$ | micro $recall_t$ | macro $F1_t$ | micro $F1_t$ |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Newsroom | original | $47.7_{\pm0.2}$ | $35.0_{\pm0.3}$ | $44.1_{\pm0.2}$ | $97.2_{\pm0.1}$ | $97.0_{\pm0.1}$ | $65.4_{\pm0.3}$ | $62.9_{\pm0.4}$ | $70.8_{\pm0.3}$ | $68.5_{\pm0.2}$ | $68.0_{\pm0.2}$ | $65.6_{\pm0.3}$ |
| | + filtering | $47.7_{\pm0.1}$ | $35.1_{\pm0.1}$ | $44.1_{\pm0.1}$ | $98.1_{\pm0.1}$ | $98.0_{\pm0.0}$ | $66.5_{\pm0.1}$ | $63.8_{\pm0.1}$ | $70.2_{\pm0.2}$ | $67.7_{\pm0.3}$ | $68.3_{\pm0.1}$ | $65.7_{\pm0.1}$ |
| | + classification | $47.7_{\pm0.2}$ | $35.1_{\pm0.1}$ | $44.2_{\pm0.2}$ | $98.1_{\pm0.1}$ | $98.0_{\pm0.0}$ | $67.2_{\pm0.4}$ | $64.2_{\pm0.4}$ | $70.3_{\pm0.2}$ | $67.8_{\pm0.4}$ | $68.7_{\pm0.3}$ | $65.9_{\pm0.4}$ |
| | JAENS | $46.6_{\pm0.5}$ | $34.3_{\pm0.3}$ | $43.2_{\pm0.3}$ | $\mathbf{98.3}_{\pm0.1}$ | $\mathbf{98.3}_{\pm0.1}$ | $\mathbf{69.5}_{\pm1.6}$ | $\mathbf{67.3}_{\pm1.2}$ | $68.9_{\pm1.5}$ | $66.8_{\pm1.6}$ | $\mathbf{69.2}_{\pm0.1}$ | $\mathbf{67.0}_{\pm0.2}$ |
| CNNDM | original | $43.7_{\pm0.1}$ | $21.1_{\pm0.1}$ | $40.6_{\pm0.1}$ | $99.5_{\pm0.1}$ | $99.4_{\pm0.1}$ | $66.0_{\pm0.4}$ | $66.5_{\pm0.4}$ | $74.7_{\pm0.7}$ | $75.4_{\pm0.6}$ | $70.0_{\pm0.2}$ | $70.7_{\pm0.3}$ |
| | + filtering | $43.4_{\pm0.2}$ | $20.8_{\pm0.1}$ | $40.3_{\pm0.2}$ | $\mathbf{99.9}_{\pm0.0}$ | $\mathbf{99.9}_{\pm0.0}$ | $66.2_{\pm0.4}$ | $66.6_{\pm0.3}$ | $74.1_{\pm0.6}$ | $74.9_{\pm0.6}$ | $69.9_{\pm0.2}$ | $70.5_{\pm0.2}$ |
| | + classification | $43.5_{\pm0.2}$ | $20.8_{\pm0.2}$ | $40.4_{\pm0.2}$ | $\mathbf{99.9}_{\pm0.0}$ | $\mathbf{99.9}_{\pm0.0}$ | $67.0_{\pm0.6}$ | $67.5_{\pm0.5}$ | $74.7_{\pm0.2}$ | $75.5_{\pm0.1}$ | $70.6_{\pm0.3}$ | $71.3_{\pm0.3}$ |
| | JAENS | $42.4_{\pm0.6}$ | $20.2_{\pm0.2}$ | $39.5_{\pm0.5}$ | $\mathbf{99.9}_{\pm0.0}$ | $\mathbf{99.9}_{\pm0.0}$ | $\mathbf{67.9}_{\pm0.7}$ | $\mathbf{68.4}_{\pm0.6}$ | $\mathbf{75.1}_{\pm0.7}$ | $\mathbf{76.4}_{\pm0.7}$ | $\mathbf{71.3}_{\pm0.2}$ | $\mathbf{72.2}_{\pm0.2}$ |
| XSUM | original | $\mathbf{45.6}_{\pm0.1}$ | $\mathbf{22.5}_{\pm0.1}$ | $\mathbf{37.2}_{\pm0.1}$ | $93.9_{\pm0.1}$ | $93.6_{\pm0.2}$ | $74.1_{\pm0.2}$ | $73.3_{\pm0.2}$ | $80.1_{\pm0.1}$ | $80.3_{\pm0.3}$ | $77.0_{\pm0.1}$ | $76.6_{\pm0.2}$ |
| | + filtering | $45.4_{\pm0.1}$ | $22.2_{\pm0.1}$ | $36.9_{\pm0.1}$ | $98.2_{\pm0.0}$ | $98.2_{\pm0.1}$ | $77.9_{\pm0.2}$ | $77.3_{\pm0.2}$ | $79.4_{\pm0.2}$ | $79.6_{\pm0.2}$ | $78.6_{\pm0.1}$ | $78.4_{\pm0.2}$ |
| | + classification | $45.3_{\pm0.1}$ | $22.1_{\pm0.0}$ | $36.9_{\pm0.1}$ | $98.3_{\pm0.1}$ | $98.2_{\pm0.1}$ | $78.6_{\pm0.3}$ | $\mathbf{78.0}_{\pm0.3}$ | $79.5_{\pm0.3}$ | $79.8_{\pm0.4}$ | $\mathbf{79.1}_{\pm0.1}$ | $\mathbf{78.9}_{\pm0.1}$ |
| | JAENS | $43.4_{\pm0.7}$ | $21.0_{\pm0.3}$ | $35.5_{\pm0.4}$ | $\mathbf{99.0}_{\pm0.1}$ | $\mathbf{99.0}_{\pm0.1}$ | $77.6_{\pm0.9}$ | $77.1_{\pm0.6}$ | $79.5_{\pm0.6}$ | $80.0_{\pm0.5}$ | $78.5_{\pm0.2}$ | $78.5_{\pm0.1}$ |

data filtering facilitates higher $prec_s$ scores, and it reveals entity hallucination can be reduced by this simple technique. Data filtering mainly enhances other entity-level metrics: $prec_t$, $recall_t$ and $F1_t$. Including the classification task (multi-task) or JAENS to data filtering is a reason to further improvement of the performance on $prec_t$ and $recall_t$. The scores are shown in percentages, averages over 5 runs with standard deviations.