

Fact-based Content Weighting for Evaluating Abstractive Summarisation

Xinnuo Xu[†], Ondřej Dušek[‡], Jingyi Li[†], Verena Rieser[†] and Ioannis Konstas[†]

[†]The Interaction Lab, MACS, Heriot-Watt University, Edinburgh, UK

[‡]Charles University, Faculty of Mathematics and Physics, Prague, Czechia

xx6, j1125, v.t.rieser, i.konstas@hw.ac.uk

odusek@ufal.mff.cuni.cz

Abstract

Abstractive summarisation is notoriously hard to evaluate since standard word-overlap-based metrics are biased towards specific words in the human reference. We introduce a new evaluation metric which abstracts away from the word-level and instead is based on fact-level content weighting, i.e. relating the facts of the document to the facts of the summary. We follow the assumption that a good summary will reflect all relevant facts, i.e. the ones present in the ground truth (human-generated reference summary). We confirm this hypothesis by showing that our weightings are highly correlated to human perception and compare favourably to the recent manual highlight-based metric of Hardy et al. (2019).

1 Introduction

Text summarisation compresses long textual documents into short summaries while retaining the most important information from the source. In contrast to extractive summarisation, which directly copies the most relevant fragments, abstractive summarization retains the most important facts and expresses them via paraphrasing, aggregating and even inferring new facts. Recent advances in neural decoders led to a number of single-document summarisation systems that exhibit some level of abstraction in their outputs, usually in the simplest form of paraphrasing (See et al. (2017); Narayan et al. (2018); Liu and Lapata (2019), *inter alia*).

Evaluating abstractive summarisation remains an open challenge (Schluter, 2017; Kryściński et al., 2019): First, decoders are amenable to pathogeniessuch as hallucination and/or omission of important information, which are hard to capture using existing evaluation metrics (Cao et al., 2018; Rohrbach et al., 2018; Dušek et al., 2020). Second, most datasets used for abstractive summarisation only contain a single reference summary,

e.g. (Narayan et al., 2018; Völske et al., 2017), which most existing automatic metrics evaluate against, e.g. ROUGE using exact n-gram overlap (Lin, 2004), and thus tend to downvote paraphrases.

We propose a new evaluation metric based on content weighting, where we abstract away from the particular surface form of the target summary, but represent it as facts using Semantic Role Labelling (SRL). In this way, we aim to better capture the semantic correctness of a summary, i.e. be more sensitive to hallucinations and omissions.¹

In particular, we weight the facts present in the source document according to the facts selected by a human-written summary. This *alignment* is conducted using contextual, rather than token-level, embeddings, e.g., BERT (Devlin et al., 2019). For evaluation, we measure whether an automatically generated summary is able to capture the same facts as the target. We also show that the computed weights correlate well with human perception. Our code is available at https://github.com/XinnuoXu/CorrFA_for_Summarizaion.

2 Related Work

The problem of reference bias has been addressed in several ways. First, metrics based on token-level or wider context embedding similarities which aim to better capture paraphrases but remain largely word-oriented, e.g. (Sun and Nenkova, 2019; Zhang et al., 2019; Zhao et al., 2019; Clark et al., 2019). Goodrich et al. (2019) come close to our approach by using entity and relation extraction, but their approach is limited to texts that lend themselves to be represented by RDF triples.

An alternative is manual evaluation against the source document. This entails selecting content either using domain experts, e.g., the PYRAMID method (Nenkova and Passonneau, 2004), factoids

¹Note that we do not make any claims about fluency, which we assume is less of a problem for neural text generation.

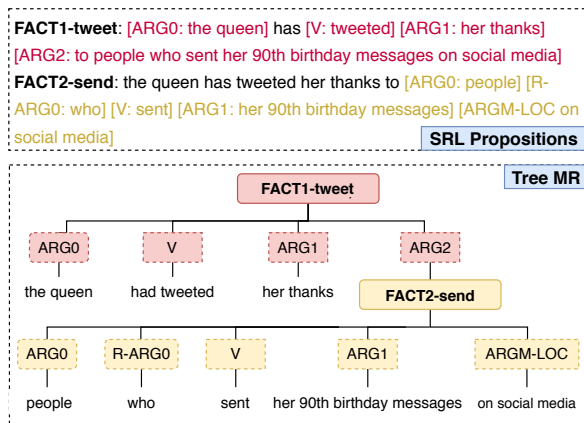


Figure 1: List of SRL propositions and corresponding tree MR with two facts for the sentence “The queen has tweeted her thanks to people who sent her 90th birthday messages on social media”.

(Teufel and van Halteren, 2004), or via crowdsourcing (Shapira et al., 2019; Hardy et al., 2019). However, evaluation based on a small human-labelled test set is noisy, time consuming, and costly. Xenou et al. (2019) propose a referenceless metric, which only checks properties of the summary, not its relation to the original document. Sun and Nenkova (2019) compare average token and sentence ELMo embeddings against the document and claim good (system-level) correlations.

Another option to avoid reference bias is question-based evaluation, either elicited manually (Clarke and Lapata, 2010; Narayan et al., 2018) or automatically (Scialom et al., 2019). However, it requires reference summaries as base for generating questions, thus only checking the summary contents indirectly.

3 Content Weighting

3.1 Fact Representation

We represent facts in a sentence by adapting SRL (Palmer et al., 2005), which roughly captures “who did what to whom” in terms of predicates and their arguments. Given a list of parsed propositions for a sentence,² each predicate-argument structure is considered as one separate *fact*, where the predicate stands for the event and its arguments are mapped to actors, recipients, time, place, etc (see Fig. 1). Following a simple observation that arguments can function as separate predicates themselves, we construct a hierarchical tree structure for the whole sentence. We create the tree meaning representa-

²We use the SRL implementation of He et al. (2018) found in <https://allennlp.org> with 86.49 test F1 on the Ontonotes 5.0 dataset.

tion (MR) from the list of facts by choosing the fact with the largest coverage as the root and recursively build sub-trees by replacing arguments with their corresponding sub-facts (ARG2 in FACT1 is replaced by FACT2 in Fig. 1).³

3.2 Automatic Content Weighting

We compute argument and fact weights by measuring the similarity of facts/arguments in the original document and the target summary based on their BERT word embeddings (for content words only) and their distance in the tree MR. We denote tokens of a document D and its summary S as $\mathbf{t}^D = \{t_1^D, t_2^D, \dots, t_n^D\}$ and $\mathbf{t}^S = \{t_1^S, t_2^S, \dots, t_m^S\}$. To get their corresponding contextual embeddings e_k^D and e_k^S , we concatenate the two texts,⁴ feed them into a pre-trained BERT model (Devlin et al., 2019) and take the contextualized embedding output from its last Transformer layer.

Argument-based weighting: We first represent the summary and the document as two sequences of leaf arguments⁵ $\{A_1^D, A_2^D, \dots, A_N^D\}$ and $\{A_1^S, A_2^S, \dots, A_M^S\}$ respectively, and weight the i -th leaf argument in the document as:

$$w_i^a = \text{avg}_{j=1 \dots M} \text{cosdist}(E_i^D, E_j^S) \quad (1)$$

i.e. the average embedding cosine distance to all arguments in the summary. Argument embeddings E_i^D and E_j^S are average embeddings of content-word tokens belonging to the arguments.⁶

$$E_i^* = \text{avg}_{k \in A_i^*, k \notin \text{stops}} e_k^* \quad (2)$$

$*$ $\in \{D, S\}$, “stops” denotes a list of stopwords.

Fact-based weighting: We can represent the summary and the document as two sequences of facts $\{F_1^D, F_2^D, \dots, F_{N'}^D\}$ and $\{F_1^S, F_2^S, \dots, F_{M'}^S\}$, and weight the i -th fact in the document by its average distance to facts in the summary:

$$w_i^f = \text{avg}_{j=1 \dots M'} d_{ij}^f \quad (3)$$

³We avoid using sentence-level MRs such as AMR (Banasescu et al., 2013), since current state-of-the-art performance of parsers is far behind compared to the simpler SRL task.

⁴By concatenating, the information in each text can be embedded in each other through self-attention. This is useful since the summary sometimes contains additional and/or common-sense knowledge not captured in the document.

⁵For example, in Fig. 1, ARG0, V, ARG1 in FACT1, and all the arguments in FACT2 are leaf arguments in the sentence, whereas ARG2 in FACT1 is not.

⁶For example, in Fig. 1, “her” and “thanks” are two tokens directly attached to the argument ARG1 of FACT1. Thus, the embedding for ARG1 of FACT1 is the average embedding of these two tokens.

The fact-level distance d_{ij}^f is defined on top of argument weighting:

$$d_{ij}^f = \text{avg}_{A_i^D \in F_i^D, A_k^S \in F_j^S} \beta_{il} \beta_{jk} [\text{cosdist}(E_i^D, E_k^S)]_{>\gamma} \quad (4)$$

It is computed as the average cosine distance over embeddings of all leaf arguments in the subtrees of fact F_i^D in the document and fact F_j^S in the summary, which is (1) filtered by a threshold γ to discard argument pairs with weak semantic relation⁷ and (2) weighted by MR tree distances of arguments to facts: $\beta_{il} = \frac{1}{\sqrt{\text{treedist}(F_i, A_l)}}$.⁸

4 Content-weighting-based Metrics

We now use these weights to introduce two metrics: **Corr-F** (fact-level) and **Corr-A** (argument-level). Let \mathbf{w}_{gold}^f and \mathbf{w}_{cand}^f denote the fact-level content weights calculated using the procedure from Section 3 based on human-reference and system-generated summaries, respectively. Similarly, \mathbf{w}_{gold}^a and \mathbf{w}_{cand}^a denote the argument-level weights. Corr-F is then the Pearson Correlation Coefficient (PCC) between \mathbf{w}_{gold}^f and \mathbf{w}_{cand}^f . Corr-A is PCC between \mathbf{w}_{gold}^a and \mathbf{w}_{cand}^a . In other words, Corr-F and Corr-A indicate whether the generated summary focuses on the informative main points in the document (i.e. the same points as the reference summary), on two different levels of granularity.

5 Metrics Evaluation

We validate our Corr-F and Corr-A metrics by collecting human judgements. In the following, we (1) collect content highlights from human judges using the Amazon Mechanical Turk platform⁹ and calculate manual content weighting based on them, (2) calculate correlations of the manual content weights with our automatic content weights, (3) compare our metrics against existing reference-based ROUGE (Lin, 2004) and BERTScore (Zhang et al., 2019), as well as the referenceless manual HROUGE score (Hardy et al., 2019).¹⁰

We use the extreme summarisation dataset (XSum; Narayan et al., 2018), which consists of

⁷In this work, we set the threshold to 0.6.

⁸E.g., in Fig. 1, $\text{treedist}(\text{FACT1}, \text{“ARG1: her thanks”}) = 1$, $\text{treedist}(\text{FACT1}, \text{“ARG0: people”}) = 2$, $\text{treedist}(\text{FACT2}, \text{“ARG0: people”}) = 1$.

⁹Using the interface from <https://github.com/sheffieldnlp/highres>.

¹⁰Note that Corr-F/A are calculated with content weighting with respect to the reference. Therefore, strictly speaking, Corr-F/A are different to all existing metrics but still share some properties with them. We show the correlation between Corr-F/A and existing metrics in terms of relative system ranking, rather than a head-to-head metrics comparison.

BBC articles and accompanying single-sentence summaries, i.e. sub-headlines of the original articles, professionally written by the authors of the articles. Due to the abstractive nature of the summaries, factoid content selection on phrase level is required beyond sentence-level extraction or token-level matching, making this dataset a popular test bed for abstractive summarisation.

We use the outputs of three recent abstractive summarization systems as evaluation targets for our metrics: (i) the Pointer-Generator model (PTGEN; See et al., 2017); (ii) the Topic-aware Convolutional Sequence-to-Sequence model (TCONVS2S; Narayan et al., 2018) and (iii) the abstractive summarization model using pretrained BERT encoders (BERTSUMABS; Liu and Lapata, 2019).¹¹

5.1 Manual Annotation Collection

Manual Content Highlighting: By extending the framework of Hardy et al. (2019), we collect manual content highlights on fact and argument levels, where we present human judges with the source document and the gold summary, with one fact/argument typeset in bold. The judges are required to select phrases or sentences in the document that support the bolded fact/argument (see Figure 4-9 in Appendix B). In both cases, judges are allowed to select parts of the text with any granularity. We limit the number of allowed continuous chunks and the maximum number of words to encourage highlights of fact/argument level.¹² We employ 3 judges per document in both cases. We use the same 50 articles and gold summaries sampled from the XSum test set as Hardy et al. (2019).

Manual Content Weighting Calculation:

Argument Level: Given a document D and a summary S , we define the weight of each token t_k^D with respect to a summary argument A_j^S as:

$$w_{kj}^t = \frac{\text{NumH}(t_k^D, A_j^S)}{\text{NumA}(A_j^S)} \quad (5)$$

$\text{NumH}(t_k^D, A_j^S)$ denotes the number of times token t_k was selected and $\text{NumA}(A_j^S)$ is the total number of annotators who were shown A_j^S bolded. We use token weights to compute manual argument-level weights \mathbf{w}_{man}^a (parallel to Eq. 1):

$$w_{man,i}^a = \text{avg}_{j=1\dots M} \text{avg}_{t_k^D \in A_j^D} w_{kj}^t \quad (6)$$

¹¹For the first two, we use candidate summaries provided by the authors. For the third, we generated summaries by training a model with code and data offered by the authors.

¹²We allow 4 chunks of max. 50 words total for fact-level and 5 chunks of max. 20 words for argument-level annotation.

Granularity	PCC-W	PCC-S
Argument-level	0.3326	0.4762
Fact-level	0.3129	0.7291

Table 1: Correlation of automatic content weighting and selection with human highlights.

Fact Level: By adapting Eq. 5, we calculate a weight w_{ki}^t for each token in document D w.r.t. bolded fact F_i^S in the summary S . The weight w_{ij}^f between fact F_i^D in the document and F_j^S in its summary is calculated using Eq. 6. We use Eq. 3 to get the manual fact content weighting w_{man}^f .

5.2 Agreement with Manual Weighting

Correlation: We evaluate how automatic content weighting w_{gold}^a and w_{gold}^f correlates with manual content weighting w_{man}^a and w_{man}^f . Using the Pearson Correlation Coefficient directly over the content weights (PCC-W), we evaluate the correlation between content weights assigned by human judges and automatically calculated weights – PCC(w_{gold}^* , w_{man}^*). As a more extreme form of weighting, we compute the correlation between content “selected” (i.e. ignoring computed weights) by human judges and the automatic mechanism (PCC-S); we set the value to 1 if the weight is over 0, meaning the fact/argument is selected.

While content-weighting correlations are just moderate, content-selection correlations are strong, especially the fact-based (Table 1). In other words, the automatic method attends to facts human judges consider important, but weighs them differently.

System-level Agreement: We check system-level agreement on Corr-F and Corr-A metrics when using automatic vs. manual content weighting (Table 2): We compute fact/argument-level content weights w_{cand}^* for each system (cf. Section 4). We then calculate Corr-F and Corr-A of w_{cand}^* against both w_{man}^* (manual weighting) and w_{gold}^* (automatic weighting) on the 50 articles with human annotation introduced in Section 5.1.

The Corr-F metric shows the same system-level ordering for both manual and automatic content weighting. Furthermore, both manual and automatic content weighting agree that TCONVS2S and PTGEN achieve similar performance but are strongly outperformed by BERTSUMABS.

5.3 Comparison to existing metrics

Corr-F/A vs. referenceless metrics: HROUGE score (Hardy et al., 2019) is a content-weighting-based referenceless evaluation metric. Unlike our

Model	Corr-F	Corr-A
Manual content weighting – w_{cand}^* vs. w_{man}^*		
TCONVS2S	0.2274	0.2464
PTGEN	0.2180	0.2433
BERTSUMABS	0.2508	0.2662
Automatic content weighting – w_{cand}^* vs. w_{gold}^*		
TCONVS2S	0.6203	0.6280
PTGEN	0.5822	0.5727
BERTSUMABS	0.6714	0.6533

Table 2: System-level scores for manual and automatic content weighting on 50 human-annotated documents.

Model	Unigram		Bigram	
	Pre	Rec	Pre	Rec
TCONVS2S	7.64	5.37	3.16	2.08
PTGEN	7.62	6.42	3.25	2.61
BERTSUMABS	8.24	6.25	3.29	2.41

Table 3: HROUGE on 50 human-annotated documents.

approach, it operates on token level and is entirely based on manual annotation. The evaluation results in Table 3 show that Corr-F/A’s ranking is identical to HROUGE’s unigram and bigram precision, with Corr-F also assigning similar proportions.¹³

Corr-F/A vs. reference-based metrics:

ROUGE (Lin, 2004) and BERTScore (Zhang et al., 2019) are both reference-based metrics, which compute a similarity score for each token in the candidate sentence with each token in the reference sentence. However, instead of exact matches as used in ROUGE, BERTScore computes token similarity using contextual embeddings. Comparing to ROUGE and BERTScore on the full XSum test set (see Table 4) shows full agreement on system ordering for both metrics.

6 Discussion

6.1 Error Analysis

We now provide examples demonstrating the strength and weaknesses of Corr-F/A by analysing system outputs where BERTScore and Corr-F/A demonstrate different ordering.

Strengths: (1) **Corr-F/A are more sensitive to content-level hallucination than BERTScore.**

Summaries with facts/arguments never mentioned in the original document get much lower Corr-F/A scores than summaries with content that appears in the document verbatim or as a paraphrase. Example 1 in Table 5 shows Corr-F/A penalizing the incorrect fact “to become the next president” generated by BERTSUMABS, while giving higher scores to TCONVS2S which paraphrased “abdicate” with

¹³We computed HROUGE for BERTSUMABS using <https://github.com/sheffielddnlp/highres>.

Model	CorrF/A		CorrF/A(L)		ROUGE			BERTScore		
	Corr-F	Corr-A	Corr-F	Corr-A	R1	R2	RL	P	R	F1
TCONVS2S	0.616	0.636	0.700	0.650	31.89	11.54	25.75	0.613	0.573	0.591
PTGEN	0.596	0.623	0.664	0.620	29.70	9.21	23.24	0.577	0.566	0.570
BERTSUMABS	0.655	0.683	0.715	0.670	38.53	16.09	30.80	0.628	0.616	0.621

Table 4: Summarisation models evaluated using Corr-F/A on full test set, with ROUGE and BERTScore scores. Note that Corr-F/A(L) is Corr-F/A calculated using a lower-performing SRL tool (He et al., 2017, see Section 6.2).

# Source	Summary	Corr-F	Corr-A	BS-F1
<i>Ground truth</i>	Japan’s emperor Akihito has expressed his desire to abdicate in the next few years, public broadcaster NHK reports.			
1 BERTSUMABS	Japan’s emperor Akihito is considering whether to become the next president of the country, reports say.	0.68	0.68	0.67
TCONVS2S	Japan’s emperor Akihito has announced that he will step down in the Japanese capital, Tokyo.	0.81	0.71	0.67
<i>Ground truth</i>	Dick Advocaat has resigned as Sunderland boss, with the team yet to win in the Premier League this season.			
2 BERTSUMABS	Sunderland manager Dick Advocaat has left the club by mutual consent after only eight games in charge.	0.60	0.66	0.65
PTGEN	Sunderland have appointed former boss Dick Advocaat as manager at the end of the season to sign a new deal.	0.26	0.34	0.65
<i>Ground truth</i>	A Chinese space capsule carrying three crew members has returned to Earth following a 13-day mission.			
3 BERTSUMABS	China has successfully landed its first ever space flight, in a move hailed as a “historic moment”.	0.56	0.67	0.53
TCONVS2S	China has successfully launched the first ever robotic mission to date for the first time in its history.	0.85	0.68	0.51
<i>Ground truth</i>	A council plans to employ its own staff to help young people with mental health problems.			
4 BERTSUMABS	A new academy to train people with mental health problems is to be set up in West Berkshire.	0.82	0.68	0.64
TCONVS2S	A new academy for children with mental health problems is being launched in West Berkshire.	0.73	0.56	0.67

Table 5: Examples of system outputs where Corr-F/A and BERTScore-F1 disagree on system ordering.

“step down”. (2) **Corr-F/A better identify phrases**, especially those containing extra content mentioned in the document but not in the ground-truth summary. Example 2 in Table 5 shows that Corr-F/A do not penalize BERTSUMABS for generating the argument “after only eight games in charge”, which is mentioned in the document.

Weaknesses: (1) **Corr-F is weaker in identifying token-level hallucination**,¹⁴ as in Example 3 in Table 5. Corr-F gives a higher score to TCONVS2S output with one hallucinated token “robotic”. However, Corr-A’s more fine-grained approach works slightly better in this case. (2) **Corr-F/A tend to under-score summaries containing content mentioned in the ground truth but only touched briefly in the document.** In Example 4 in Table 5, Corr-F/A score the output of TCONVS2S lower, even though it correctly captures “an academy for children with mental health”, which is mentioned only once in the document.

In sum, Corr-F/A is less dependent on the reference summary by also considering the source document, and thus has less of a reference bias than BERTScore. In addition, Corr-F/A helps to identify ungrounded facts, i.e. content-level hallucinations, which is important for identifying misinformation in automated news reporting.

6.2 Robustness of Corr-F/A

As noted in Section 3.1, Corr-F/A is based on publicly available SRL tools. To demonstrate the robustness of our metrics, we evaluate the same sys-

¹⁴Token-level hallucination means an incorrect token within an otherwise correct fact structure. Content-level hallucination happens when whole facts or arguments are hallucinated.

tem outputs with Corr-F/A calculated using a lower-performing SRL tool (He et al., 2017).¹⁵ The results are shown as Corr-F/A(L) in Table 4 and show full agreement with Corr-F/A in terms of system ordering. However, the better performing original SRL system widens the margin between systems.

7 Conclusions and Future Work

We present an automatic evaluation framework for abstractive summarisation, which is low-cost and robust, as it does not rely on expert annotators nor is susceptible to crowdsourcing noise. Using fact representations, we are able to capture semantically similar, but at the same time distant in surface form, content in the summary that aligns with arbitrarily far-apart parts of the input document, casting our metric to be directly interpretable. Our metric is more sensitive to perturbations of the facts in the target summary, which resemble common hallucination phenomena of neural decoders (see Figure 2-3 in Appendix A for examples). In the future, we intend to investigate different meaning representation formalisms, such as AMR (Banarescu et al., 2013) and Dynamic Syntax (Kempson et al., 2001) and extend to other datasets (e.g. multiple-reference summarization) and tasks (e.g. response generation in dialogue).

Acknowledgements

This research received funding from the EPSRC project MaDrIgAL (EP/N017536/1) and Charles University project PRIMUS/19/SCI/10. We would like to acknowledge the AWS Cloud Credits for Research programme.

¹⁵F1 on the Ontonotes 5.0 dataset is 81.6% ($\delta = -4.89$).

References

- Laura Banarescu, Claire Bonial, Shu Cai, Madalina Georgescu, Kira Griffitt, Ulf Hermjakob, Kevin Knight, Philipp Koehn, Martha Palmer, and Nathan Schneider. 2013. [Abstract meaning representation for sembanking](#). In *Proceedings of the 7th Linguistic Annotation Workshop and Interoperability with Discourse*, pages 178–186.
- Ziqiang Cao, Furu Wei, Wenjie Li, and Sujian Li. 2018. [Faithful to the Original: Fact Aware Neural Abstractive Summarization](#). In *AAAI*, New Orleans, LA, USA. ArXiv: 1711.04434.
- Elizabeth Clark, Asli Celikyilmaz, and Noah A Smith. 2019. [Sentence Mover’s Similarity: Automatic Evaluation for Multi-Sentence Texts](#). In *ACL*, Florence, Italy.
- James Clarke and Mirella Lapata. 2010. [Discourse constraints for document compression](#). *Computational Linguistics*, 36(3):411–441.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. [BERT: Pre-training of deep bidirectional transformers for language understanding](#). In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4171–4186, Minneapolis, Minnesota. Association for Computational Linguistics.
- Ondřej Dušek, Jekaterina Novikova, and Verena Rieser. 2020. [Evaluating the state-of-the-art of end-to-end natural language generation: The E2E NLG Challenge](#). *Computer Speech & Language*, 59:123 – 156.
- Ben Goodrich, Vinay Rao, Mohammad Saleh, and Peter J. Liu. 2019. [Assessing The Factual Accuracy of Generated Text](#). In *KDD*, Anchorage, AK, USA. ArXiv: 1905.13322.
- Hardy, Shashi Narayan, and Andreas Vlachos. 2019. [HighRES: Highlight-based reference-less evaluation of summarization](#). In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 3381–3392, Florence, Italy.
- Luheng He, Kenton Lee, Omer Levy, and Luke Zettlemoyer. 2018. [Jointly predicting predicates and arguments in neural semantic role labeling](#). In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 364–369, Melbourne, Australia.
- Luheng He, Kenton Lee, Mike Lewis, and Luke Zettlemoyer. 2017. [Deep semantic role labeling: What works and what’s next](#). In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 473–483, Vancouver, Canada. Association for Computational Linguistics.
- Ruth M Kempson, Wilfried Meyer-Viol, and Dov M Gabbay. 2001. *Dynamic syntax: The flow of language understanding*. Blackwell Oxford.
- Wojciech Kryściński, Nitish Shirish Keskar, Bryan McCann, Caiming Xiong, and Richard Socher. 2019. [Neural text summarization: A critical evaluation](#). In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 540–551, Hong Kong, China.
- Chin-Yew Lin. 2004. [ROUGE: A package for automatic evaluation of summaries](#). In *Text Summarization Branches Out*, pages 74–81, Barcelona, Spain. Association for Computational Linguistics.
- Yang Liu and Mirella Lapata. 2019. [Text summarization with pretrained encoders](#). In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 3728–3738, Hong Kong, China.
- Shashi Narayan, Shay B Cohen, and Mirella Lapata. 2018. [Don’t give me the details, just the summary! topic-aware convolutional neural networks for extreme summarization](#). In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 1797–1807, Brussels, Belgium.
- Ani Nenkova and Rebecca Passonneau. 2004. [Evaluating content selection in summarization: The pyramid method](#). In *Proceedings of the Human Language Technology Conference of the North American Chapter of the Association for Computational Linguistics: HLT-NAACL 2004*, pages 145–152, Boston, Massachusetts, USA. Association for Computational Linguistics.
- Martha Palmer, Daniel Gildea, and Paul Kingsbury. 2005. [The proposition bank: An annotated corpus of semantic roles](#). *Computational Linguistics*, 31(1):71–106.
- Anna Rohrbach, Lisa Anne Hendricks, Kaylee Burns, Trevor Darrell, and Kate Saenko. 2018. [Object Hallucination in Image Captioning](#). In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 4035–4045, Brussels, Belgium.
- Natalie Schluter. 2017. [The limits of automatic summarisation according to ROUGE](#). In *Proceedings of the 15th Conference of the European Chapter of the Association for Computational Linguistics: Volume 2, Short Papers*, pages 41–45, Valencia, Spain. Association for Computational Linguistics.
- Thomas Scialom, Sylvain Lamprier, Benjamin Piwowarski, and Jacopo Staiano. 2019. [Answers](#)

- Unite! Unsupervised Metrics for Reinforced Summarization Models. In *2019 Conference on Empirical Methods in Natural Language Processing (EMNLP) and 9th International Joint Conference on Natural Language Processing (IJCNLP)*, Hong Kong. ArXiv: 1909.01610.
- Abigail See, Peter J. Liu, and Christopher D. Manning. 2017. *Get To The Point: Summarization with Pointer-Generator Networks*. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics*, pages 1073–1083, Vancouver, Canada.
- Ori Shapira, David Gabay, Yang Gao, Hadar Ronen, Ramakanth Pasunuru, Mohit Bansal, Yael Amsterdamer, and Ido Dagan. 2019. *Crowdsourcing Lightweight Pyramids for Manual Summary Evaluation*. In *NAACL*, Minneapolis, MN, USA. ArXiv: 1904.05929.
- Simeng Sun and Ani Nenkova. 2019. *The Feasibility of Embedding Based Automatic Evaluation for Single Document Summarization*. In *2019 Conference on Empirical Methods in Natural Language Processing (EMNLP) and 9th International Joint Conference on Natural Language Processing (IJCNLP)*, Hong Kong.
- Simone Teufel and Hans van Halteren. 2004. *Evaluating information content by factoid analysis: Human annotation and stability*. In *Proceedings of the 2004 Conference on Empirical Methods in Natural Language Processing*, pages 419–426, Barcelona, Spain. Association for Computational Linguistics.
- Michael Völske, Martin Potthast, Shahbaz Syed, and Benno Stein. 2017. *TL;DR: Mining Reddit to learn automatic summarization*. In *Proceedings of the Workshop on New Frontiers in Summarization*, pages 59–63, Copenhagen, Denmark. Association for Computational Linguistics.
- Stratos Xenouelas, Prodromos Malakasiotis, Marianna Apidianaki, and Ion Androutsopoulos. 2019. *SUMQE: a BERT-based Summary Quality Estimation Model*. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 6007–6013, Hong Kong, China.
- Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q. Weinberger, and Yoav Artzi. 2019. *BERTScore: Evaluating text generation with BERT*. *arXiv preprint arXiv:1904.09675*.
- Wei Zhao, Maxime Peyrard, Fei Liu, Yang Gao, Christian M. Meyer, and Steffen Eger. 2019. *MoverScore: Text Generation Evaluating with Contextualized Embeddings and Earth Mover Distance*. In *2019 Conference on Empirical Methods in Natural Language Processing (EMNLP) and 9th International Joint Conference on Natural Language Processing (IJCNLP)*, Hong Kong. ArXiv: 1909.02622.

A Fact-level Content Weighting Examples

Fig. 2 and 3 show examples for documents weighted using Corr-F/Corr-A with respect to different summaries.

In Fig. 2, the left column shows one document weighted by the reference summary and two system-generated summaries from BERTSUMABS and TCONVS2S respectively (summaries are shown in the right column). As we can see, there are 4 relatively important facts in the document weighted by the reference summary. BERTSUMABS and TCONVS2S capture 3 and 2 out of 4, respectively. Other than the important facts highlighted by the reference summary, TCONVS2S also assigns high weights to other facts; that leads to the hallucinated generation and lower Corr-F/Corr-A scores. On the other hand, BERTSUMABS’s summary weighs facts in the document in a similar way to the reference summary, which lead to a strongly related summary and high Corr-F and Corr-A scores.

In Fig. 3, there are 5 relatively important facts in the document weighted by the reference summary. BERTSUMABS and TCONVS2S capture 4 and 3 out of 5, respectively. Both systems miss the fact “Pope Francis, who has taken a more liberal stance on homosexuality”. However, the weight of this fact given by BERTSUMABS’s output is higher than with TCONVS2S’s. The Corr-F and Corr-A are lower for TCONVS2S due to misweighting of informative facts in the document.

B Annotation Interface

We provide the following illustrations of the human annotation interface:

- Annotation interface for manual content weighting examples, including the instructions, for fact-level (Fig. 4 and 5) and argument-level (Fig. 6 and 7) annotation,
- Examples of human annotation results for fact (Fig. 9) and argument (Fig. 8) level.

Please refer to the individual figure captions for detailed descriptions.

turia pitt received burns to 65 % of her body and was told she would never compete again after the 2011 ultra - marathon in western australia after 200 operations , she completed the 226 km @ 140 mile @ hawaii event on sunday she has been hailed on social media as an " amazing role model " find out how to get into triathlon in our special guide ms pitt completed a 3.8 km swim , 180 km bike ride and 42 km run in a time of 14:37:30 the mining engineer , author and motivational speaker completed the event using custom - made gear and brake levers on her bike to accommodate the severe injuries to her hands she also wore special race clothing to deal with the extreme heat and humidity because of my burns , i ca n't regulate my own body temperature so i 've had to make some adjustments , or use standard tri gear in different ways , " she revealed last month i 'm going to need things like cooling sleeves and white suits so i do n't overheat after her encounter with the bushfire ms pitt spent 864 days in hospital and underwent many operations to treat her injuries she made her comeback in the ironman australia triathlon in may. competing in ironman has ultimately showed me that i literally can do anything i put my mind to , " she said on her blog ms pitt was widely praised on social media as an inspiration such a great example of what can be done through hard work belief in one 's self and good a support system , " one person wrote on her facebook page you continue to amaze and inspire , " said another

Reference:

An australian runner who suffered like threatening burns when she was trapped by a bushfire during a race has completed the hawaii ironman, seen as the world's toughest triathlon

turia pitt received burns to 65 % of her body and was told she would never compete again after the 2011 ultra - marathon in western australia after 200 operations , she completed the 226 km @ 140 mile @ hawaii event on sunday she has been hailed on social media as an " amazing role model " find out how to get into triathlon in our special guide ms pitt completed a 3.8 km swim , 180 km bike ride and 42 km run in a time of 14:37:30 the mining engineer , author and motivational speaker completed the event using custom - made gear and brake levers on her bike to accommodate the severe injuries to her hands she also wore special race clothing to deal with the extreme heat and humidity because of my burns , i ca n't regulate my own body temperature so i 've had to make some adjustments , or use standard tri gear in different ways , " she revealed last month i 'm going to need things like cooling sleeves and white suits so i do n't overheat after her encounter with the bushfire ms pitt spent 864 days in hospital and underwent many operations to treat her injuries she made her comeback in the ironman australia triathlon in may. competing in ironman has ultimately showed me that i literally can do anything i put my mind to , " she said on her blog ms pitt was widely praised on social media as an inspiration such a great example of what can be done through hard work belief in one 's self and good a support system , " one person wrote on her facebook page you continue to amaze and inspire , " said another

BertSumAbs:

An australian runner who suffered severe burns in a bushfire in hawaii has completed an ironman triathlon

Corr-F: 0.96 Corr-A: 0.88

turia pitt received burns to 65 % of her body and was told she would never compete again after the 2011 ultra - marathon in western australia after 200 operations , she completed the 226 km @ 140 mile @ hawaii event on sunday she has been hailed on social media as an " amazing role model " find out how to get into triathlon in our special guide ms pitt completed a 3.8 km swim , 180 km bike ride and 42 km run in a time of 14:37:30 the mining engineer , author and motivational speaker completed the event using custom - made gear and brake levers on her bike to accommodate the severe injuries to her hands she also wore special race clothing to deal with the extreme heat and humidity because of my burns , i ca n't regulate my own body temperature so i 've had to make some adjustments , or use standard tri gear in different ways , " she revealed last month i 'm going to need things like cooling sleeves and white suits so i do n't overheat after her encounter with the bushfire ms pitt spent 864 days in hospital and underwent many operations to treat her injuries she made her comeback in the ironman australia triathlon in may. competing in ironman has ultimately showed me that i literally can do anything i put my mind to , " she said on her blog ms pitt was widely praised on social media as an inspiration such a great example of what can be done through hard work belief in one 's self and good a support system , " one person wrote on her facebook page you continue to amaze and inspire , " said another

TConvS2S:

An australian runner become the first person to win a race for the first time in almost 30 years

Corr-F: 0.67 Corr-A: 0.73

Figure 2: A document (left) weighted with respect to a reference summary and two system outputs (right), with Corr-F/Corr-A scores. The colour represents the sum of argument- and fact-level weights for each token (Eqs. 3 and 4). The darker the colour, the more important the fact is.

the french government proposed senior diplomat laurent stefanini for the post in january but the vatican is yet to respond to approve the choice the vatican usually responds within six weeks to approve such a new ambassador the nomination of mr stefanini was seen as a litmus test for pope francis , who has taken a more liberal stance on homosexuality a french government spokesman said there had been negotiations with the vatican over the appointment france has chosen its ambassador to the vatican . this choice was stefanini and that remains the french proposal , " said spokesman stephane le foll observers say most vatican appointments are confirmed within six weeks and that this long silence should be read as a rejection the vatican traditionally makes no statement if it intends to decline a nomination mr stefanini served in the holy see as a deputy ambassador in the french embassy from 2001 to 2005 he was described by the country 's foreign ministry as " one of our best diplomats " france legalised same - sex marriage in 2013 , despite opposition from the catholic church pope francis is regarded as more tolerant of homosexuality than previous popes . " who am i to judge ? " , he said in 2013

Reference:

France has said it will not back down over its nomination of an openly gay ambassador to the Vatican.

the french government proposed senior diplomat laurent stefanini for the post in january but the vatican is yet to respond to approve the choice the vatican usually responds within six weeks to approve such a new ambassador the nomination of mr stefanini was seen as a litmus test for pope francis , who has taken a more liberal stance on homosexuality a french government spokesman said there had been negotiations with the vatican over the appointment france has chosen its ambassador to the vatican . this choice was stefanini and that remains the french proposal , " said spokesman stephane le foll observers say most vatican appointments are confirmed within six weeks and that this long silence should be read as a rejection the vatican traditionally makes no statement if it intends to decline a nomination mr stefanini served in the holy see as a deputy ambassador in the french embassy from 2001 to 2005 he was described by the country 's foreign ministry as " one of our best diplomats " france legalised same - sex marriage in 2013 , despite opposition from the catholic church pope francis is regarded as more tolerant of homosexuality than previous popes . " who am i to judge ? " , he said in 2013

BertSumAbs:

France has said it is considering whether to appoint a French ambassador to the Vatican as a replacement for the right-wing politician.

Corr-F: 0.73 Corr-A: 0.59

the french government proposed senior diplomat laurent stefanini for the post in january but the vatican is yet to respond to approve the choice the vatican usually responds within six weeks to approve such a new ambassador the nomination of mr stefanini was seen as a litmus test for pope francis , who has taken a more liberal stance on homosexuality a french government spokesman said there had been negotiations with the vatican over the appointment france has chosen its ambassador to the vatican . this choice was stefanini and that remains the french proposal , " said spokesman stephane le foll observers say most vatican appointments are confirmed within six weeks and that this long silence should be read as a rejection the vatican traditionally makes no statement if it intends to decline a nomination mr stefanini served in the holy see as a deputy ambassador in the french embassy from 2001 to 2005 he was described by the country 's foreign ministry as " one of our best diplomats " france legalised same - sex marriage in 2013 despite opposition from the catholic church pope francis is regarded as more tolerant of homosexuality than previous popes . " who am i to judge ? " , he said in 2013

TConvS2S:

The Vatican has announced the appointment of a new ambassador to the Vatican.

Corr-F: 0.69 Corr-A: 0.39

Figure 3: Another document (left) weighted with respect to a reference summary and two system outputs (right), with Corr-F/Corr-A scores (see Fig. 2 for details).

Main Points Highlight

This task requires **three** steps:

Step 1: Highlighting Informative Phrases in a Document that Support the Bolded Part in a Given Summary the Best

In this step, you will be given two components a **summary with one bolded main point** and a **document**.

You should:

- 1) read the summary and identify the bolded main point
- 2) read the document
- 3) select sentences/phrases that support the bolded main point the best

An example of highlights:

Summary:
researchers have identified **a gene that may put people at greater risk of strokes and heart attacks**

(The main point you should identify is: a gene that may put people at greater risk of strokes and heart attacks)

Document:
...
Writing in PLOS ONE they say the **gene fault may encourage the formation of blood clots - the ultimate cause of most heart attacks and strokes**
...
Around one in 10 people in the Caucasian population carries **this variation of the gene , named PIA2**
...
They found **individuals with PIA2 were more likely to have a stroke**
...
the scientists show **PIA2 is also linked to an increased risk of heart attacks** in people under 45
...

You can highlight **4 sentences/phrases** at most. The maximum combined length of all highlights is **50 words**.

Figure 4: The instruction for fact-level human highlight annotation.

Please don't refresh the page.

Instructions

Your task is to highlight informative phrases in the document that support the bolded part in the given summary the best.

The maximum combined length of all highlighted phrases is **50 words**.

To highlight, use your mouse to select phrases from the document, and click on the pen icon.

To delete a group of highlights, right click on it and confirm .

Summary:
the queen has tweeted her thanks to people who sent her 90th birthday messages on social media

" I am most grateful for the many digital messages of goodwill I have received and would like to thank you all for your kindness , " she wrote.

The monarch , whose milestone birthday was marked with numerous events , signed off the rare message " Elizabeth R " .

The Queen sent her first ever tweet in 2014 when she opened a new exhibition at the Science Museum in London.

Britain 's longest-serving monarch celebrated her 90th birthday on 21 April , and a host of events were held over three months , from April to June.

The Queen has two birthdays - her real birthday on 21 April , and her official birthday held on a Saturday in June - a tradition going back 250 years. It was introduced to try to ensure better weather for the monarch 's official celebrations.

Summary

Words left
50 words.

Chunks left
4 chunks.

Highlighted Phrases:

Figure 5: The human annotation interface for fact level. Human judges are required to highlight content in the document that is supporting the fact printed in bold "The Queen has tweeted her thanks" (FACT1 of the summary in Figure 1 in the paper).

Noun Phrases Highlight

This task requires **three** steps:

Step 1: Highlighting Informative Phrases in a Document that Support the Bolded Part in a Given Summary the Best

In this step, you will be given two components a **summary with one bolded noun phrase** and a **document**.

You should:

- 1) read the summary and identify the bolded noun phrase
- 2) read the document
- 3) select **Noun Phrases** that support the bolded noun phrase the best

Examples of highlights:

Example 1:

Summary:

researchers have identified a **gene** that may put people at greater risk of strokes and heart attacks
(The noun phrase you should identify is: a gene)

Document:

...Writing in PLOS ONE they say the **gene fault** may encourage the formation of blood clots - the ultimate cause of most heart attacks and strokes ...
Around one in 10 people in the Caucasian population carries **this variation of the gene, named P1A2** ...
They found individuals with **P1A2** were more likely to have a stroke ...
the scientists show **P1A2** is also linked to an increased risk of heart attacks in people under 45 ...

You can highlight **5 noun phrases** at most. The maximum combined length of all highlights is **20 words**. Please note that for some cases, the bolded noun phrase is **paraphrased** in the document. You should also highlight the paraphrased phrases.

Figure 6: The instruction for argument-level human highlight annotation.

Please don't refresh the page.

Instructions

Your task is to highlight informative phrases in the document that support the bolded part in the given summary the best.

The maximum combined length of all highlighted phrases is **20 words**.

To highlight, use your mouse to select phrases from the document, and click on the pen icon.

To delete a group of highlights, right click on it and confirm.

Summary

Words left
20 words.

Chunks left
5 chunks.

Highlighted Phrases:

Summary:
the queen has tweeted her thanks to people who sent her 90th birthday messages on **social media**

" I am most grateful for the many digital messages of goodwill I have received and would like to thank you all for your kindness , " she wrote.

The monarch , whose milestone birthday was marked with numerous events , signed off the rare message " Elizabeth R " .

The Queen sent her first ever tweet in 2014 when she opened a new exhibition at the Science Museum in London.

Britain 's longest-serving monarch celebrated her 90th birthday on 21 April , and a host of events were held over three months , from April to June.

The Queen has two birthdays - her real birthday on 21 April , and her official birthday held on a Saturday in June - a tradition going back 250 years. It was introduced to try to ensure better weather for the monarch 's official celebrations.

Her official birthday this year was 11 June and the annual Trooping the Colour was held on Horse Guards Parade , followed by an RAF flypast which the Royal Family watched from the balcony of Buckingham Palace.

The following day the Queen hosted the Patron 's Lunch , a street party for some 10,000 people along The Mall which recognised her patronage of more than 600 organisations in the UK and around the Commonwealth.

Queen Elizabeth II at 90 Find out more about Queen Elizabeth II on BBCiWonder

[Click to submit](#)

Summary

Words left
20 words.

Chunks left
5 chunks.

Highlighted Phrases:

Figure 7: The human annotation interface for argument level. Human judges are required to highlight content in the document that is supporting the phrase printed in bold "on social media" (argument ARGM-LOC of FACT2 of the summary in Figure 1 in the paper).

Article

" i am most grateful for the many digital messages of goodwill i have received and would like to thank you all for your kindness , " she wrote . the monarch , whose milestone birthday was marked with numerous events , signed off the rare message " elizabeth r " . the queen sent her first ever tweet in 2014 when she opened a new exhibition at the science museum in london . britain 's longest-serving monarch celebrated her 90th birthday on 21 april , and a host of events were held over three months , from april to june . the queen has two birthdays - her real birthday on 21 april , and her official birthday held on a saturday in june - a tradition going back 250 years . it was introduced to try to ensure better weather for the monarch 's official celebrations . her official birthday this year was 11 june and the annual trooping the colour was held on horse guards parade , followed by an raf flypast which the royal family watched from the balcony of buckingham palace . the following day the queen hosted the patron 's lunch , a street party for some 10,000 people along the mall which recognised her patronage of more than 600 organisations in the uk and around the commonwealth . queen elizabeth ii at 90 find out more about queen elizabeth ii on bbc iWonder

Figure 8: Human highlight annotation for the argument ARG1 of FACT1 “her thanks” of the summary in Figure 1 in the paper.

Article

" i am most grateful for the many digital messages of goodwill i have received and would like to thank you all for your kindness , " she wrote . the monarch , whose milestone birthday was marked with numerous events , signed off the rare message " elizabeth r " . the queen sent her first ever tweet in 2014 when she opened a new exhibition at the science museum in london . britain 's longest-serving monarch celebrated her 90th birthday on 21 april , and a host of events were held over three months , from april to june . the queen has two birthdays - her real birthday on 21 april , and her official birthday held on a saturday in june - a tradition going back 250 years . it was introduced to try to ensure better weather for the monarch 's official celebrations . her official birthday this year was 11 june and the annual trooping the colour was held on horse guards parade , followed by an raf flypast which the royal family watched from the balcony of buckingham palace . the following day the queen hosted the patron 's lunch , a street party for some 10,000 people along the mall which recognised her patronage of more than 600 organisations in the uk and around the commonwealth . queen elizabeth ii at 90 find out more about queen elizabeth ii on bbc iWonder

Figure 9: Human highlight annotation for the FACT1 “The Queen has tweeted her thanks” of the summary in Figure 1 in the paper.