

TRUE: Re-evaluating Factual Consistency Evaluation

Anonymous ACL submission

Abstract

001 Grounded text generation systems often gener- 002
003 ate text that contains factual inconsistencies, hindering their real-world applicability. 004
005 Automatically evaluating such inconsistencies may help to alleviate this limitation by ac- 006
007 celerating evaluation cycles, filtering inconsistent outputs and annotating large-scale training 008
009 data. While attracting increasing attention, such evaluation metrics are usually developed 010
011 and evaluated in silo for a single task or dataset. Moreover, previous meta-evaluation protocols 012
013 focused on system-level correlations with human annotations, which leave the example- 014
015 level accuracy of such metrics unclear. In this work, we introduce TRUE: a comprehensive 016
017 study of factual consistency metrics on a standardized collection of existing texts from diverse 018
019 tasks, manually annotated for factual consistency. Our standardization enables an example-level 020
021 meta-evaluation protocol that is more actionable and interpretable than previously reported 022
023 correlations, yielding clearer quality measures. Across diverse state-of-the-art metrics and 11 024
025 datasets we find that large-scale NLI and question generation-and-answering-based 026
027 approaches achieve strong and complementary results, and recommend them as a starting point 028
029 for future evaluations.¹

1 Introduction

030 A core issue in deploying text generation models for real-world applications is that they often 031
032 generate factually inconsistent text with respect to the input they are conditioned on, or even 033
034 completely “hallucinate” (Lee et al., 2018; Rohrbach et al., 2018; Maynez et al., 2020; Zhao et al., 2020) 035
036 as exemplified in Table 1.

037 To tackle such inconsistencies, one would like to detect them automatically by predicting whether 038
039 a generated text is factually consistent with respect to a grounding text (also referred to as the “knowledge”, or the “input”). Such capabilities attract 040
041

¹Our code will be made publicly available.

Summarization (Wang et al., 2020)

Input	Phyllis schlaflly, a leading figure in the us conservative movement, has died at her home in missouri, aged 92...
Summary	Us conservative activist phyllis schlaflly has died at the age of 87.

Fact Verification (Thorne et al., 2018)

Evidence	Ronald Bilius “Ron” Weasley is a character in J. K. Rowling’s Harry Potter fictional series.
Claim	Ron Weasley is a President.

Paraphrasing (Zhang et al., 2019)

Input	The tracks were produced by Tommy Lee , and feature Michael Beinhorn on drums.
Paraphrase	The tracks were produced by Michael Beinhorn and have Tommy Lee on drums.

Knowledge-Grounded Dialogue (Honovich et al., 2021)

Knowledge	The first flip trick called a kickflip, originally called a “magic flip,” was invented by professional skateboarder Rodney Mullen.
Response	I remember the first one was called magic flip. It was called a magic flip and was invented in the 60’s.

Table 1: Factual inconsistencies (in red) from various tasks which are part of the TRUE study. The corresponding parts in the input/grounding are in blue.

042 increasing attention (Zhou et al., 2021) as they enable both better evaluation and better generation 043
044 models via filtering training data (Gehrmann et al., 2021) or annotation of training data for controlled 045
046 generation (Rashkin et al., 2021b).

047 While automatically evaluating factual consistency is an active line of work, there is no single 048
049 agreed-upon meta-evaluation protocol for measuring the quality of such methods, and labeling 050
051 schemes vary in their granularity. Works are usually done in silo, introducing new datasets and 052
053 methods that target a specific task or domain, such as summarization (Falke et al., 2019; Kryscinski 054
055 et al., 2020; Wang et al., 2020; Scialom et al., 2021; Deutsch et al., 2021; Xie et al., 2021) or dialogue 056

(Dziri et al., 2021; Honovich et al., 2021; Nie et al., 2021; Qin et al., 2021). Comparing the robustness of such methods *across* tasks and datasets is therefore difficult, impeding progress on this subject.

This work presents a comprehensive study focusing on factuality evaluation, covering various metrics, tasks and datasets. To allow this, we consolidate 11 existing datasets annotated for factual consistency into a unified format, including pairs of a target text and a grounding source, with a binary annotation of whether the target text is factually consistent w.r.t its source. These datasets² cover summarization, knowledge-grounded dialogue, paraphrasing and fact verification. The proposed standardization enables us to properly compare factuality evaluation methods in a robust manner across these various tasks and domains.

Previous works on factuality assessment have mainly focused on measuring system-level correlations of the proposed metrics with human judgments (Pagnoni et al., 2021). Yet, these correlations are not useful for estimating the performance of a measured metric when making *binary* decisions, decoupled from specific system implementations. We aim to measure how well a method detects inconsistent texts (*recall*) and how often it falsely disregards consistent texts (*precision*), which can be easily computed using the aforementioned binary labeling scheme. Therefore, as a meta-evaluation protocol we report the Area Under the ROC Curve (ROC AUC) with respect to inconsistent example detection for each evaluation metric and dataset.

Our thorough evaluation of 12 metrics draws a clearer picture on the state of evaluating factuality. We show that Natural Language Inference (NLI) approaches, as well as Question Generation and Answering (QG-QA) approaches achieve significantly better³ results on a wide variety of tasks and datasets. We also show that NLI and QG-QA are complementary: combining the two yields even better results and hints that there is room for further improvement. Finally, we perform both quantitative and qualitative analysis of our results, finding that all approaches struggle with long inputs, labeling issues and personal statements – paving interesting avenues for future work.

To summarize, our contributions are as follows: (1) We argue that work on factuality evaluation

²We focus on English text-to-text tasks, and leave data-to-text (Parikh et al., 2020; Reiter and Thomson, 2020), multilingual and multimodal tasks to future work.

³We conduct significance testing, see section 4.

should be unified and generalized across tasks, and standardize 11 published datasets into a single labeling scheme to corroborate this. (2) We propose a meta-evaluation protocol that allows more actionable and interpretable quality measures than previously reported correlations. (3) We perform a meta-evaluation of 12 diverse metrics in this unified perspective, showing that large-scale NLI and QG-QA-based approaches achieve strong and complementary results across tasks. (4) We analyze our results both qualitatively and quantitatively, pointing at challenges like long inputs and personal statements to be addressed in future work.

2 Standardizing Factual Consistency

In this section we elaborate on our re-evaluation setup. We first formally define what factual consistency refers to in this work. We then detail the datasets we consider and how we standardize them. Finally, we discuss the meta-evaluation protocol we propose for measuring the performance of evaluation methods on the standardized datasets.

2.1 Definitions and Terminology

We define a text to be factually consistent w.r.t its grounding text if all the factual information it conveys is consistent with the factual information conveyed by the grounding text.⁴ While some previous works distinguished between inconsistent erroneous text to inconsistent correct text (Maynez et al., 2020), we take a strict approach, requiring the text to be faithful to its grounding text, regardless of the “correctness” w.r.t the “real world”. In other words, we consider only the information present in the input text, not external knowledge, to assess faithfulness. This enables a more well-defined task, since determining the truthfulness of a fact w.r.t a general “real world” is subjective and depends on the knowledge, values and beliefs of the subject (Heidegger, 2001). This definition follows similar strictness in Textual Entailment, Question Answering, Summarization and other tasks where comprehension is based on a given grounding text, irrespective of contradiction with other world knowledge. This is also in line with recent work on evaluating attribution in text generation (Rashkin et al., 2021a), where humans are required to judge whether a generated text is true according to a grounding text. We use the terms *consistent*,

⁴We exclude personal and social statements, such as opinions and chit chat from the scope of factual information.

Task	# Examples	Open Test	Cons.
Summarization			
- FRANK (Pagnoni et al., 2021)	671	+	33.2%
- SummEval (Fabbri et al., 2021)	1,600	-	81.6%
- MNBM (Maynez et al., 2020)	2,500	-	10.2%
- QAGS-CNN/DM (Wang et al., 2020)	235	-	48.1%
- QAGS-XSum (Wang et al., 2020)	239	-	48.5%
Dialogue			
- BEGIN (Dziri et al., 2021)	836	+	33.7%
- Q ² (Honovich et al., 2021)	1,088	-	57.7%
- DialFact (Gupta et al., 2021)	8,689	+	38.5%
Fact Verification			
- FEVER (Thorne et al., 2018)	18,209	-	35.1%
- VitaminC (Schuster et al., 2021)	63,054	+	49.9%
Paraphrasing			
- PAWS (Zhang et al., 2019)	8,000	+	44.2%

Table 2: Statistics for the datasets incorporated in TRUE. Cons. is the ratio of consistent examples.

grounded, faithful and factual interchangeably.

2.2 Standardization Process

We include 11 datasets that contain human annotations w.r.t factual consistency in diverse tasks (Table 2). Other than the importance of covering a wide variety of error types, this also alleviates issues of rating quality which may vary across datasets (Denton et al., 2021).

To allow a unified evaluation framework we convert all annotations to binary labels that correspond to whether the entire target text is factual w.r.t the given grounding text or not. We note that a fine-grained annotation scheme, i.e., a typology of errors, was proposed for factual consistency (Pagnoni et al., 2021). While useful, most existing datasets do not include such labels. Moreover, while Machine Translation (MT) evaluation also showed value in fine-grained annotations (Freitag et al., 2021), it was proposed after years of improving MT to the level where coarse-grained annotation is insufficient. We argue that current grounded generation models are still at early stages w.r.t factual consistency, and that binary labeling is more beneficial now as it enables easier standardization across tasks and domains, with the goal of bringing researchers to collaborate on a shared methodology. Binary annotation also corresponds to practical applications where filtering out unfaithful predictions is desired, and is in-line with the recommendations for human evaluation of attribution in text generation by Rashkin et al. (2021a).

We next detail the 11 datasets included in TRUE.

2.2.1 Abstractive Summarization

FRANK Pagnoni et al. (2021) proposed a typology of factual errors, grounded in frame semantics (Fillmore, 1976; Palmer et al., 2005) and linguistic discourse theory (Brown and Yule, 1983). Based on this typology, they collected annotations for model-generated summaries on the

CNN/DailyMail (CNN/DM; Hermann et al., 2015) and XSum (Narayan et al., 2018) datasets, resulting in 2250 annotated system outputs. Each summary sentence was annotated by three annotators. We take the majority vote for each sentence to get a sentence-level label and consider a summary as consistent if all sentences are consistent.

SummEval SummEval (Fabbri et al., 2020) is a comprehensive study of evaluation metrics for text summarization. The authors collected human judgments for 16 model outputs on 100 articles taken from the CNN/DM dataset, using both extractive and abstractive models. Annotators were asked to rate summaries on a Likert scale from 1 to 5, over 4 dimensions: *consistency*, *coherence*, *fluency* and *relevance*. Each summary was scored by 5 crowd-workers and 3 expert annotators. We label summaries as consistent only if all the expert annotators gave a *consistency* score of 5.

MNBM Maynez et al. (2020) annotated system outputs for the XSum dataset (Narayan et al., 2018). They sampled 500 articles and annotated summaries generated by four different systems, as well as the gold summaries. Annotators were asked to assess whether the summary includes hallucinations. Judgments from three different annotators were collected for each document-summary pair. To convert to a binary-label format, we use the binary consistency decision of whether a summary contains no hallucinations, and assign a label by taking the majority vote of the three annotators.

QAGS Wang et al. (2020) collected judgments of factual consistency on generated summaries for CNN/DM and XSum. Annotators were presented with the summaries one sentence at a time, along with the article, and determined whether each sentence is factually consistent w.r.t the article. Each sentence was annotated by 3 annotators, using the majority vote as the final score. To convert to binary-label format, we consider a summary consistent only if all its sentences are consistent.

2.2.2 Dialogue Generation

BEGIN (Dziri et al., 2021) is a dataset for evaluating groundedness in knowledge-grounded dialogue systems, in which system outputs should be consistent with a grounding knowledge provided to the dialogue agent. BEGIN frames the task as NLI (Bowman et al., 2015), adopting the *entailment* and *contradiction* labels, and splitting the neutral

label into three sub-categories: *hallucination*, *off-topic* responses and *generic* responses. Dialogue responses were generated by fine-tuning two systems on the Wizard of Wikipedia (WOW) dataset (Dinan et al., 2019), in which responses should be grounded in a span of text from Wikipedia. The generated responses were split into sentences, and each sentence was annotated separately. To convert to a binary-label format, we treat entailed sentences as consistent and all others as inconsistent.

Q² Honovich et al. (2021) annotated 1,088 generated dialogue responses for binary factual consistency w.r.t the knowledge paragraph provided to the dialogue model, for two dialogue models trained on WOW. Responses were annotated using binary labels by 3 of the paper authors, one annotator per response. We use Q²'s labels without changes.

DialFact Gupta et al. (2021) introduced the task of fact-verification in dialogue and constructed a dataset of conversational claims paired with pieces of evidence from Wikipedia. They define three tasks: (1) detecting whether a response contains verifiable content (2) retrieving relevant evidence and (3) predicting whether a response is *supported* by the evidence, *refuted* by the evidence or if there is *not enough information* to determine. We use the verifiable (i.e., factual, rather than personal) responses annotated for the third task, treating *supported* annotations as consistent and the rest as inconsistent. In cases where several evidence were marked as required for verification, we concatenate all evidence sentences to be the grounding text.

2.2.3 Fact Verification

FEVER Thorne et al. (2018) introduced FEVER (Fact Extraction and VERification), a dataset for fact verification against textual sources. FEVER was constructed by extracting information from Wikipedia, generating claims from it using annotators, then classifying whether each claim is *supported* or *refuted* by Wikipedia. Claims can also be labeled with *NotEnoughInfo*, meaning that there is not enough information in Wikipedia to either verify or refute the claim. Given a claim, the task defined by FEVER is to first extract evidence, then to determine whether it supports or refutes the claim. In a slightly different framing, the latter stage in FEVER is to determine whether the claim is factually consistent or not w.r.t the evidence, which is aligned with what we measure in TRUE. We use the development set of the NLI version of FEVER

(Nie et al., 2019, 2020), treating *supported* claims as consistent and the rest as inconsistent.

VitaminC Schuster et al. (2021) derived a large-scale fact verification dataset from factual revisions to Wikipedia pages. Each example includes an evidence text from Wikipedia and a fact, with an annotation of whether the fact is supported, refuted or neutral w.r.t the evidence. The authors collected factual revisions to Wikipedia articles (pairs of “before” and “after” sentences), and asked annotators to write two facts for each pair: one that is *supported* by the first sentence and *refuted* by the second, and vice versa. When no explicit contradiction was present, the annotators wrote facts that are *neutral* w.r.t the evidence. Additional examples were created by revising examples from FEVER. We treat examples that include *supported* facts as consistent, and *refuted* or *neutral* facts as inconsistent.

2.2.4 Paraphrase Detection

PAWS Zhang et al. (2019) constructed a dataset for paraphrase identification with 108,463 paraphrase and non-paraphrase pairs with high lexical overlap, generated by controlled word swapping and back-translation, followed by judgments from human raters. Source sentences were drawn from Wikipedia and the Quora Question Pairs (QQP) corpus. We only use the examples with Wikipedia source sentences and view the binary paraphrase labels as consistency labels. We note that the definition of paraphrase is not equivalent to the definition of factual consistency, as a subset of a source text is not a paraphrase but may still be factually consistent with the source. However, PAWS was constructed such that non-paraphrases usually have contradicting meanings and is therefore relevant.

2.3 Meta-Evaluation

Previous work on evaluating factuality focused on measuring correlation with human judgements (Pagnoni et al., 2021). However, such numbers are not very informative when one is interested in evaluating the absolute performance of inconsistency detection methods that perform a *binary* decision w.r.t each input.

To conduct a more fine-grained evaluation at the single example level, we report the Receiver Operating Characteristic Area Under the Curve (ROC AUC) w.r.t binary detection of inconsistent examples.⁵ The ROC curve is created by plotting the

⁵This is equivalent to AUC w.r.t consistency detection.

true positive rate (TPR, a.k.a. the recall) against the *false positive rate* (FPR, a.k.a. the fallout) at different possible thresholds for each tested metric. Measuring ROC AUC evaluates the different metrics without setting a specific decision threshold.

For datasets with existing development/test split, we also tune a threshold for the binary consistency/inconsistency decision on the development set and report the test set accuracy using this threshold. We tune the thresholds by optimizing the geometric mean of TPR and 1-FPR: $\sqrt{\text{TPR} * (1 - \text{FPR})}$.

3 Evaluation Metrics

We compare various standard as well as state-of-the-art approaches that measure factual consistency. This comparison should draw a clear picture of current research on this subject and directions for future work. For example, we expect that robust metrics should perform well across tasks and datasets. We next describe the different metrics tested as part of this study. We note that for all reference-based metrics, the grounding text serves as the reference. For metrics where the scores are not in the [0,1] range, we normalize scores to be in that range.

3.1 N-Gram Based Metrics

Standard N-Gram matching metrics such as BLEU (Papineni et al., 2002) ROUGE (Lin, 2004) and token-level F1 were shown to have weak correlation with factual consistency (Maynez et al., 2020; Honovich et al., 2021), with no exception on TRUE. For completeness, we report their performance in Table 9 in the appendix.

3.2 Model-Based Metrics

BERTScore (Zhang et al., 2020) aggregates similarity scores between the BERT contextual embedding of tokens in candidate and reference sentences. We report results for the BERTScore-precision variant as it showed better results in preliminary experiments. We use BERTScore version 0.3.11. with the DeBERTa-xl-MNLI model (He et al., 2021; Nangia et al., 2017), which is the recommended model as of the time of writing this paper.⁶

BLEURT (Sellam et al., 2020a,b) is a learned metric based on BERT (Devlin et al., 2019) for evaluating text generation. BLEURT includes additional pretraining on synthetic data followed by

⁶https://github.com/Tiiiger/bert_score

fine-tuning on human judgements to train a model that scores system outputs. We use the recommended BLEURT-20 checkpoint (Pu et al., 2021).⁷

FactCC (Kryscinski et al., 2020) is a BERT-based metric trained to verify factual consistency of summaries. Training data was synthetically generated by applying rule-based transformations to generate consistent and inconsistent summaries.

BARTScore (Yuan et al., 2021) evaluates text using probabilities from force-decoding with a BART model (Lewis et al., 2020). We use the version fine-tuned on the ParaBank2 dataset (Hu et al., 2019).

3.3 Natural Language Inference Metrics

ANLI The task of Textual Entailment (Dagan et al., 2006) or Natural Language Inference (NLI; Bowman et al., 2015) is to determine, given two sentences, a *hypothesis* and a *premise*, whether the *hypothesis* is entailed by the *premise*, contradicts it, or is neutral w.r.t it. The resemblance of NLI to factual consistency evaluation has led to utilizing NLI models for measuring factual consistency (Thorne et al., 2018; Maynez et al., 2020; Dziri et al., 2021). We trained an NLI model by fine-tuning T5-11B (Raffel et al., 2020) on the Adversarial NLI (ANLI; Nie et al., 2020) dataset. As suggested by Maynez et al. (2020), we compute the entailment probability with the grounding text as the premise and the generated text as the hypothesis and use it as the example-level factual consistency score.⁸

SUMMAC (Summary Consistency; Laban et al., 2021) is focused on evaluating factual consistency in summarization. They use NLI for detecting inconsistencies by splitting the document and summary into sentences and performing NLI on all document/summary sentence pairs, where the premise is a document sentence and the hypothesis is a summary sentence. They aggregate the NLI scores for all pairs by either taking the maximum score per summary sentence and averaging (SC_{ZS}) or by training a convolutional neural network to aggregate the scores (SC_{Conv}). We use the publicly available implementation⁹ and report results for SC_{ZS} as it performed better in our experiments.

⁷<https://github.com/google-research/bleurt/blob/master/checkpoints.md>

⁸More implementation details on the NLI model are available in Section B in the appendix.

⁹<https://github.com/tingofurro/summac>

	Ensemble	Q^2	ANLI	SC _{ZS}	F1	BLEURT	QuestEval	FactCC	BART _{score}	BERT _{score}
FRANK	91.2	87.8	89.4	89.1	76.1	82.8	84.0	76.4	86.1	84.3
SummEval	82.9	78.8	80.5	81.7	61.4	66.7	70.1	75.9	73.5	77.2
MNBM	76.6	68.7	77.9**	71.3	46.2	64.5	65.3	59.4	60.9	62.8
QAGS-C	87.7	83.5	82.1	80.9	63.8	71.6	64.2	76.4	80.9	69.1
QAGS-X	84.8	70.9	83.8	78.1	51.1	57.2	56.3	64.9	53.8	49.5
BEGIN	86.2	79.7	82.6	82.0	86.4	86.4	84.1	64.4	86.3	87.9
Q^2	82.8	80.9*	72.7	77.4	65.9	72.4	72.2	63.7	64.9	70.0
DialFact	90.4	86.1**	77.7	84.1	72.3	73.1	77.3	55.3	65.6	64.2
PAWS	91.2	89.7**	86.4	88.2	51.1	68.3	69.2	64.0	77.5	77.5
FEVER	94.7	88.4	93.2**	93.2	51.8	59.5	72.6	61.9	64.1	63.3
VitaminC	96.1	81.4	88.3**	97.9	61.4	61.8	66.5	56.3	63.2	62.5
Avg. w/o VitC, FEVER	86.0	80.7	81.5	81.4	63.8	71.4	71.4	66.7	72.2	71.4

Table 3: ROC AUC results for the different metrics on the TRUE development set. We exclude VitaminC and FEVER from the average calculation as SC_{ZS} was trained on VitaminC that includes examples from FEVER. The highest score in each row (excluding the Ensemble) is in bold and the aforementioned SC results are in strikethrough. Statistically significant results are indicated using * and ** for $p < 0.05$ and $p < 0.01$ respectively.

3.4 QG-QA Based Metrics

Durmus et al. (2020) and Wang et al. (2020) proposed to use Question Generation (QG) and Question Answering (QA) models to automatically evaluate factual consistency in abstractive summarization, showing promising results. Honovich et al. (2021) employed a similar approach for evaluating knowledge-grounded dialogue generation.

The steps of the QG-QA approach are as follows: (1) Questions are automatically generated for spans in the generated text, such that the answer to a question is its respective input span. (2) The generated questions are answered using a QA model on the grounding text, resulting in an answer span or a “no-answer” output. (3) For each question, the two answer spans from the grounding and the generated text are compared to get a score. (4) The scores for all questions are aggregated into a final score.

Q^2 (Honovich et al., 2021) is a QG-QA method that employs an NLI model to compare the two answers for each question, where the grounding text answer is the premise and the generated text answer is the hypothesis. We report results for a re-implementation of Q^2 using T5-11B as the backbone for the QG, QA and NLI models. While Honovich et al. (2021) validate each generated question by answering it using a QA model and comparing to the original extracted answer candidate using exact match, we relax this and instead use F1 token-overlap with a predefined threshold.¹⁰

QuestEval (Scialom et al., 2021) is a QG-QA method that measures both factual consistency and relevance (by reversing the roles of the generated and grounding texts). The authors trained a model that weights each generated question according to

¹⁰More implementation details are available in Section B in the appendix.

the relevance of its answer to appear in the generated text. Their results showed high correlation with human judgments in comparison to prior work on the SummEval benchmark (Fabbri et al., 2021). We use the publicly available version.¹¹

4 Results

We report the ROC AUC¹² of various metrics on the standardized datasets in Table 3. The ROC curves can be found in Figure 2 in the appendix. As all metrics operate in a “zero-shot” manner on all datasets (except for SUMMAC on VitaminC and FEVER) and no threshold tuning is required, we report results on the development sets.¹³ SC_{ZS} was trained on VitaminC which includes examples from FEVER, so we exclude those datasets from the average AUC calculation for a more fair comparison.

The results show that the NLI-based models (ANLI, SC_{ZS}) outperformed the other approaches on 6 datasets, with average AUC of 81.5 and 81.4 for ANLI and SC_{ZS}, respectively. Q^2 outperform the other approaches on 4 datasets, with an average AUC of 80.7. The next best method, BARTScore, had lower average AUC of 72.2. All other approaches scored 72 or lower on average across all datasets (excluding FEVER and VitaminC).

One outlier is BEGIN, which is the only dataset where simple metrics like F1 token overlap achieved scores higher than 80. We measured the average overlap between the grounding and target texts per dataset, and found that BEGIN exhibits a high difference between grounded and ungrounded texts in comparison to other datasets (Table 8 in appendix A), which explains this.

¹¹<https://github.com/ThomasScialom/QuestEval>

¹²Multipled by 100 for better readability.

¹³AUC and accuracy for the test sets are provided in Tables 10 and 11 in the appendix.

We follow Laban et al. (2021) and perform significance testing through bootstrap resampling (Efron, 1982), comparing the best method to the second-best method on each dataset. We perform interval comparison at $p = 0.05$ and $p = 0.01$ and find significantly best results on 6 datasets, 3 from Q^2 and 3 from ANLI.

Given that no single method outperformed the rest on all datasets, we hypothesize that the NLI and QG-QA based metrics are complementary. We test this by averaging the Q^2 , ANLI and SCz_s scores per example¹⁴ (Ensemble in Table 3). Indeed, averaging the three methods yields better results on most datasets and on average, with an increase of 4.5 in ROC AUC from the best single-metric result.

Our results show that a single metric can do well across all tasks and datasets, with all 3 best metrics scoring higher than 70 on all 11 datasets. This corroborates our hypothesis that evaluating factual consistency can be unified, and we hope such unified perspective will be adopted in future work to accelerate progress on the subject.

5 Analysis

Input Length. As QA and NLI models may struggle with long inputs (Kočíský et al., 2018; Pang et al., 2021; Yin et al., 2021; Shaham et al., 2022), metrics based on them may fail when handling long text. To study the effect of input length on the metrics performance, we unify all datasets¹⁵ and split examples into 6 bins according to the grounding length.¹⁶ We focus on the grounding as the target texts are usually short (see Table 6 in Appendix A). We measure AUC of the best 3 metrics according to their overall score for each length bin, sampling 1,000 examples per bin.

The results are shown in Figure 1. We find that there is a consistent degradation for texts longer than 200 tokens for all metrics, including SCz_s which is designed to better handle long text. We find it surprising that the ANLI-based model and Q^2 still do relatively well on the longest bin as they are required to perform end-to-end QA and NLI on texts with more than 500 tokens.

¹⁴Pairwise ensembles are reported in the appendix, Table 9.

¹⁵Excluding VitaminC as it is much larger than other datasets and might therefore distort results. Statistics regarding the grounding and target text lengths per dataset is in Appendix A.

¹⁶We measure length in tokens (before subword splitting) as different metrics use different subword tokenizations.

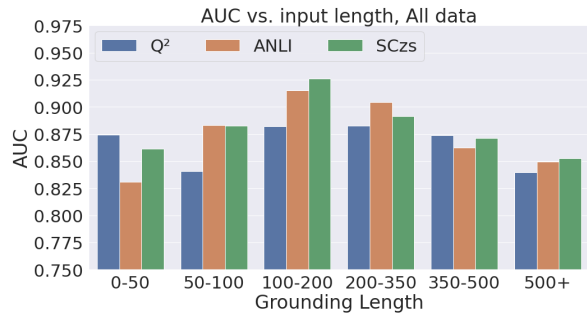


Figure 1: ROC AUC when splitting TRUE’s data according to the grounding length.

Model Size. Model-based metrics are expected to benefit from increasing model size. To quantify this we study the effect of using smaller models for the ANLI, BLEURT and BERTScore metrics. We compare the average ROC AUC of larger and smaller model variants for each metric. We find an advantage of 4.7, 3.7 and 1.3 average ROC AUC for the larger ANLI, BLEURT and BERTScore variants respectively, showing that larger models are important for evaluating factuality. The full results are in Table 7 in the appendix.

Qualitative Analysis. We conduct manual error analysis to point at weaknesses of the different metrics and present challenges posed by the task. We analyze 80 examples that were misclassified by all three best metrics, as well as 100 examples that were correctly classified by one or two of the three.

Out of the analyzed examples, many seem to have a wrong label. This is especially true for cases in which all best metrics failed, with annotation errors in 35/80 cases. For the cases where one or two metrics failed, we found annotation errors in 27/100 cases. To verify that the high annotation error rate is indeed a result of inspecting the “hardest” examples, we uniformly sample additional 100 examples, finding that only 10 had annotation errors. This is in line with the findings of Freitag et al. (2021), who showed that in some cases, metrics may be better than non-expert annotators.

Despite showing impressive results, the best-performing metrics fail to detect subtle inconsistencies, as presented in Table 4. This was the case for 21/180 analyzed examples. Metrics that aggregate scores across parts of a target text, such as Q^2 or SCz_s , might assign a high score for texts in which all but a small part is consistent. End-to-end NLI should predict “contradiction” even when only a small part of the text contradicts the grounding, but it may fail to do so. Applying a strict approach

Grounding	Generated Text	Explanation
The word "philately" is the English version of the French word "philatélie", coined by Georges Herpin in 1864.	The word philately is actually a french word coined by george herpin.	The word philately is an English word based on a French word, but not French. All best metrics misclassified this.
French police have interviewed presidential candidate francois fillon and his wife penelope over claims she was paid for fake work. They provided information that would help find the "truth", mr fillon said...	French presidential candidate francois fillon has said he and his wife penelope have been questioned by police over claims she worked illegally.	Most details are correct and the hallucination is subtle. In the case of Q^2 , most of the generated questions have the same answer based on the grounding and the generated text, therefore the overall score was high.
Stamp collecting is generally accepted as one of the areas that make up the wider subject of philately, which is the study of stamps.	I've never heard of stamps, but I do know that the word "philately" refers to the study of stamps.	The personal statement "I've never heard of stamps" is not factual and should not be evaluated.
Evidence suggests that cognitive behavioral therapy and a gradual increase in activity suited to individual capacity can be beneficial in some cases.	It has been suggested that cognitive behavioral therapy and gradual increase in exercise could help in some cases so I'm going to try that for now.	Similar to the previous examples - SummaC and ANLI falsely marked the text as inconsistent, probably due to the personal statement.

Table 4: Examples for the error analysis. The first two rows show cases of challenging inconsistencies, while the last two show dialogue responses containing non-factual personal statements.

in the aggregation step, like taking the minimum instead of the average, could potentially remedy this – with the price of having more false-negatives. Other errors are caused by domain-specific challenges, such as handling personal statements in dialogues. As shown in Table 4, such statements may be falsely classified as ungrounded. This was the case for 10/62 analyzed dialogue responses. A possible way to alleviate this would be to automatically exclude non-factual parts from the evaluation.

6 Related Work

Adding to the related work mentioned throughout the paper, works on unified evaluation of text generation across tasks include GEM (Gehrmann et al., 2021), where the focus is on evaluating system outputs and not the factuality evaluation methods as in TRUE. BEAMetrics (Scialom and Hill, 2021) proposes meta-evaluation protocols across tasks, but does not focus on factuality. When discussing factuality (“correctness”) they measure correlations, which are not sufficient as mentioned in Section 2.3. Other works on meta-evaluation of factuality across datasets include GO-FIGURE (Gabriel et al., 2021) FRANK (Pagnoni et al., 2021) and SummaC (Laban et al., 2021), however they all focus solely on summarization. To the best of our knowledge, our work is the first to generalize the discussion on evaluating factuality across tasks and datasets outside of summarization, and the first to show that large-scale QG-QA and NLI are highly complementary – setting stronger baselines for future work than previously published.

7 Discussion and Future Work

We discuss the main takeaways of the TRUE study, pointing at actionable insights for future work.

First, as QG-QA and NLI-based methods show better performance than other approaches, especially when combined together, we recommend model developers to use those methods for evaluation when factuality is a priority. As for metric developers, we recommend using those methods as baselines when proposing new metrics.

We also suggest reporting ROC AUC rather than correlations, as it is more interpretable and actionable. Our proposed binary annotation scheme allows to easily test new metrics across tasks and datasets, which would be useful for future work.

Finally, we encourage data curators to use the binary annotation scheme, which is inline with the recommendations of Rashkin et al. (2021a). Having said that, we do not rule out more detailed labeling schemes – but rather ask to provide a protocol for converting such labels into the more general binary format. We hope that future work will also address the challenges of long input text and personal statements in dialogue evaluation, which we point out in our analysis.

8 Conclusions

We presented TRUE, a meta-evaluation study for factual consistency. We standardized various datasets from diverse tasks into a unified labeling scheme to perform a thorough analysis of automatic evaluation methods, showing that NLI and QG-QA based approaches perform well across multiple tasks and datasets. We further show these methods are highly complementary – hinting at additional headroom for improvement while pointing on current limitations. We hope our results and methodology will encourage a more unified perspective in future work to foster progress towards more factual NLP applications.

References

- 649 Samuel R. Bowman, Gabor Angeli, Christopher Potts,
650 and Christopher D. Manning. 2015. [A large anno-](#)
651 [tated corpus for learning natural language inference.](#)
652 In *Proceedings of the 2015 Conference on Empiri-*
653 *cal Methods in Natural Language Processing*, pages
654 632–642, Lisbon, Portugal. Association for Compu-
655 tational Linguistics.
- 656 Gillian Brown and George Yule. 1983. *Discourse Anal-*
657 *ysis*. Cambridge University Press.
- 658 Ido Dagan, Oren Glickman, and Bernardo Magnini.
659 2006. The pascal recognising textual entailment
660 challenge. In *Machine Learning Challenges. Eval-*
661 *uating Predictive Uncertainty, Visual Object Classi-*
662 *fication, and Recognising Tectual Entailment*, pages
663 177–190, Berlin, Heidelberg. Springer Berlin Hei-
664 delberg.
- 665 Emily Denton, Mark Díaz, Ian Kivlichan, Vinodku-
666 mar Prabhakaran, and Rachel Rosen. 2021. Whose
667 ground truth? accounting for individual and collec-
668 tive identities underlying dataset annotation. *arXiv*
669 *preprint arXiv:2112.04554*.
- 670 Daniel Deutsch, Tania Bedrax-Weiss, and Dan Roth.
671 2021. Towards question-answering as an automatic
672 metric for evaluating the content quality of a sum-
673 mary. *Transactions of the Association for Computa-*
674 *tional Linguistics*, 9:774–789.
- 675 Jacob Devlin, Ming-Wei Chang, Kenton Lee, and
676 Kristina Toutanova. 2019. [BERT: Pre-training of](#)
677 [deep bidirectional transformers for language under-](#)
678 [standing.](#) In *Proceedings of the 2019 Conference*
679 *of the North American Chapter of the Association*
680 *for Computational Linguistics: Human Language*
681 *Technologies, Volume 1 (Long and Short Papers)*,
682 pages 4171–4186, Minneapolis, Minnesota. Associ-
683 ation for Computational Linguistics.
- 684 Emily Dinan, Stephen Roller, Kurt Shuster, Angela
685 Fan, Michael Auli, and Jason Weston. 2019. Wizard
686 of Wikipedia: Knowledge-powered conversational
687 agents. In *Proceedings of the International Confer-*
688 *ence on Learning Representations (ICLR)*.
- 689 Esin Durmus, He He, and Mona Diab. 2020. [FEQA: A](#)
690 [question answering evaluation framework for faith-](#)
691 [fulness assessment in abstractive summarization.](#) In
692 *Proceedings of the 58th Annual Meeting of the Asso-*
693 *ciation for Computational Linguistics*, pages 5055–
694 5070, Online. Association for Computational Lin-
695 guistics.
- 696 Nouha Dziri, Hannah Rashkin, Tal Linzen, and David
697 Reitter. 2021. [Evaluating groundedness in dialogue](#)
698 [systems: The begin benchmark.](#)
- 699 Bradley Efron. 1982. *The jackknife, the bootstrap and*
700 *other resampling plans*. SIAM.
- 701 Alexander R Fabbri, Wojciech Kryściński, Bryan
702 McCann, Caiming Xiong, Richard Socher,
and Dragomir Radev. 2020. Summeval: Re-
evaluating summarization evaluation. *arXiv*
preprint arXiv:2007.12626.
- Alexander R Fabbri, Wojciech Kryściński, Bryan
McCann, Caiming Xiong, Richard Socher, and
Dragomir Radev. 2021. Summeval: Re-evaluating
summarization evaluation. *Transactions of the Asso-*
ciation for Computational Linguistics, 9:391–409.
- Tobias Falke, Leonardo F. R. Ribeiro, Prasetya Ajie
Utama, Ido Dagan, and Iryna Gurevych. 2019.
[Ranking generated summaries by correctness: An in-](#)
[teresting but challenging application for natural lan-](#)
[guage inference.](#) In *Proceedings of the 57th Annual*
Meeting of the Association for Computational Lin-
guistics, pages 2214–2220, Florence, Italy. Associa-
tion for Computational Linguistics.
- C. Fillmore. 1976. Frame semantics and the nature of
language *. *Annals of the New York Academy of Sci-*
ences, 280.
- Markus Freitag, George Foster, David Grangier, Viresh
Ratnakar, Qijun Tan, and Wolfgang Macherey. 2021.
Experts, errors, and context: A large-scale study of
human evaluation for machine translation. *arXiv*
preprint arXiv:2104.14478.
- Saadia Gabriel, Asli Celikyilmaz, Rahul Jha, Yejin
Choi, and Jianfeng Gao. 2021. [GO FIGURE: A](#)
[meta evaluation of factuality in summarization.](#) In
Findings of the Association for Computational Lin-
guistics: ACL-IJCNLP 2021, pages 478–487, On-
line. Association for Computational Linguistics.
- Sebastian Gehrmann, Tosin Adewumi, Karmanya
Aggarwal, Pawan Sasanka Ammanamanchi,
Anuluwapo Aremu, Antoine Bosselut, Khy-
athi Raghavi Chandu, Miruna-Adriana Clinciu,
Dipanjan Das, Kaustubh Dhole, Wanyu Du,
Esin Durmus, Ondřej Dušek, Chris Chinenye
Emezue, Varun Gangal, Cristina Garbacea, Tat-
sunori Hashimoto, Yufang Hou, Yacine Jernite,
Harsh Jhamtani, Yangfeng Ji, Shailza Jolly, Mi-
hir Kale, Dhruv Kumar, Faisal Ladhak, Aman
Madaan, Mounica Maddela, Khyati Mahajan,
Saad Mahamood, Bodhisattwa Prasad Majumder,
Pedro Henrique Martins, Angelina McMillan-
Major, Simon Mille, Emiel van Miltenburg, Moin
Nadeem, Shashi Narayan, Vitaly Nikolaev, Andre
Niyongabo Rubungo, Salomey Osei, Ankur Parikh,
Laura Perez-Beltrachini, Niranjan Ramesh Rao,
Vikas Raunak, Juan Diego Rodriguez, Sashank
Santhanam, João Sedoc, Thibault Sellam, Samira
Shaikh, Anastasia Shimorina, Marco Antonio
Sobrevilla Cabezudo, Hendrik Strobelt, Nishant
Subramani, Wei Xu, Diyi Yang, Akhila Yerukola,
and Jiawei Zhou. 2021. [The GEM benchmark: Nat-](#)
[ural language generation, its evaluation and metrics.](#)
In *Proceedings of the 1st Workshop on Natural*
Language Generation, Evaluation, and Metrics
(GEM 2021), pages 96–120, Online. Association for
Computational Linguistics.

761	Prakhar Gupta, Chien-Sheng Wu, Wenhao Liu, and Caiming Xiong. 2021. Dialfact: A benchmark for fact-checking in dialogue. <i>arXiv preprint arXiv:2110.08222</i> .	816
762		817
763		818
764		819
		820
765	Pengcheng He, Xiaodong Liu, Jianfeng Gao, and Weizhu Chen. 2021. Deberta: Decoding-enhanced bert with disentangled attention . In <i>International Conference on Learning Representations</i> .	821
766		822
767		823
768		824
769	Martin Heidegger. 2001. On the essence of truth. <i>The Nature of Truth: Classic and Contemporary Perspectives</i> , 1:295–316.	825
770		826
771		827
772	Karl Moritz Hermann, Tomas Kocisky, Edward Grefenstette, Lasse Espeholt, Will Kay, Mustafa Suleyman, and Phil Blunsom. 2015. Teaching machines to read and comprehend . In <i>Advances in Neural Information Processing Systems</i> , volume 28. Curran Associates, Inc.	828
773		829
774		830
775		831
776		832
777		833
778	Or Honovich, Leshem Choshen, Roei Aharoni, Ella Neeman, Idan Szpektor, and Omri Abend. 2021. q²: Evaluating factual consistency in knowledge-grounded dialogues via question generation and question answering . In <i>Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing</i> , pages 7856–7870, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.	834
779		835
780		836
781		837
782		838
783		839
784		840
785		841
786		842
787	J. Edward Hu, Abhinav Singh, Nils Holzenberger, Matt Post, and Benjamin Van Durme. 2019. Large-scale, diverse, paraphrastic bitexts via sampling and clustering . In <i>Proceedings of the 23rd Conference on Computational Natural Language Learning (CoNLL)</i> , pages 44–54, Hong Kong, China. Association for Computational Linguistics.	843
788		844
789		845
790		846
791		847
792		848
793		849
794		850
795		851
796	Tomáš Kočiský, Jonathan Schwarz, Phil Blunsom, Chris Dyer, Karl Moritz Hermann, Gábor Melis, and Edward Grefenstette. 2018. The NarrativeQA reading comprehension challenge . <i>Transactions of the Association for Computational Linguistics</i> , 6:317–328.	852
797		853
798		854
799		855
800		856
801	Wojciech Kryscinski, Bryan McCann, Caiming Xiong, and Richard Socher. 2020. Evaluating the factual consistency of abstractive text summarization . In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)</i> , pages 9332–9346, Online. Association for Computational Linguistics.	857
802		858
803		859
804		860
805		861
806		862
807	Philippe Laban, Tobias Schnabel, Paul N Bennett, and Marti A Hearst. 2021. Summac: Re-visiting nli-based models for inconsistency detection in summarization . <i>arXiv preprint arXiv:2111.09525</i> .	863
808		864
809		865
810		866
811	Zhenzhong Lan, Mingda Chen, Sebastian Goodman, Kevin Gimpel, Piyush Sharma, and Radu Soricut. 2019. Albert: A lite bert for self-supervised learning of language representations . <i>arXiv preprint arXiv:1909.11942</i> .	867
812		868
813		869
814		870
815		871
		872
		873
		874
		875
		876
		877
		878
		879
		880
		881
		882
		883
		884
		885
		886
		887
		888
		889
		890
		891
		892
		893
		894
		895
		896
		897
		898
		899
		900
		901
		902
		903
		904
		905
		906
		907
		908
		909
		910
		911
		912
		913
		914
		915
		916
		917
		918
		919
		920
		921
		922
		923
		924
		925
		926
		927
		928
		929
		930
		931
		932
		933
		934
		935
		936
		937
		938
		939
		940
		941
		942
		943
		944
		945
		946
		947
		948
		949
		950
		951
		952
		953
		954
		955
		956
		957
		958
		959
		960
		961
		962
		963
		964
		965
		966
		967
		968
		969
		970
		971
		972
		973
		974
		975
		976
		977
		978
		979
		980
		981
		982
		983
		984
		985
		986
		987
		988
		989
		990
		991
		992
		993
		994
		995
		996
		997
		998
		999
		1000

874			
875			
876			
877			
878			
879	Artidoro Pagnoni, Vidhisha Balachandran, and Yulia		
880	Tsvetkov. 2021. Understanding factuality in abstrac-		
881	tive summarization with FRANK: A benchmark for		
882	factuality metrics . In <i>Proceedings of the 2021 Con-</i>		
883	<i>ference of the North American Chapter of the Asso-</i>		
884	<i>ciation for Computational Linguistics: Human Lan-</i>		
885	<i>guage Technologies</i> , pages 4812–4829, Online. As-		
886	sociation for Computational Linguistics.		
887	Martha Palmer, Daniel Gildea, and Paul Kingsbury.		
888	2005. The Proposition Bank: An Annotated Cor-		
889	pus of Semantic Roles . <i>Computational Linguistics</i> ,		
890	31(1):71–106.		
891	Richard Yuanzhe Pang, Alicia Parrish, Nitish Joshi,		
892	Nikita Nangia, Jason Phang, Angelica Chen,		
893	Vishakh Padmakumar, Johnny Ma, Jana Thompson,		
894	He He, et al. 2021. Quality: Question answer-		
895	ing with long input texts, yes! <i>arXiv preprint</i>		
896	<i>arXiv:2112.08608</i> .		
897	Kishore Papineni, Salim Roukos, Todd Ward, and Wei-		
898	Jing Zhu. 2002. Bleu: a method for automatic eval-		
899	uation of machine translation . In <i>Proceedings of</i>		
900	<i>the 40th Annual Meeting of the Association for Com-</i>		
901	<i>putational Linguistics</i> , pages 311–318, Philadelphia,		
902	Pennsylvania, USA. Association for Computational		
903	Linguistics.		
904	Ankur Parikh, Xuezhi Wang, Sebastian Gehrmann,		
905	Manaaf Faruqui, Bhuwan Dhingra, Diyi Yang, and		
906	Dipanjan Das. 2020. ToTTo: A controlled table-to-		
907	text generation dataset . In <i>Proceedings of the 2020</i>		
908	<i>Conference on Empirical Methods in Natural Lan-</i>		
909	<i>guage Processing (EMNLP)</i> , pages 1173–1186, On-		
910	line. Association for Computational Linguistics.		
911	Amy Pu, Hyung Won Chung, Ankur Parikh, Sebastian		
912	Gehrmann, and Thibault Sellam. 2021. Learning		
913	compact metrics for MT . In <i>Proceedings of the 2021</i>		
914	<i>Conference on Empirical Methods in Natural Lan-</i>		
915	<i>guage Processing</i> , pages 751–762, Online and Punta		
916	Cana, Dominican Republic. Association for Compu-		
917	tational Linguistics.		
918	Libo Qin, Tianbao Xie, Shijue Huang, Qiguang Chen,		
919	Xiao Xu, and Wanxiang Che. 2021. Don't be contra-		
920	dicted with anything! CI-ToD: Towards benchmark-		
921	ing consistency for task-oriented dialogue system .		
922	In <i>Proceedings of the 2021 Conference on Empiri-</i>		
923	<i>cal Methods in Natural Language Processing</i> , pages		
924	2357–2367, Online and Punta Cana, Dominican Re-		
925	public. Association for Computational Linguistics.		
926	Colin Raffel, Noam Shazeer, Adam Roberts, Kather-		
927	ine Lee, Sharan Narang, Michael Matena, Yanqi		
928	Zhou, Wei Li, and Peter J. Liu. 2020. Exploring		
929	the limits of transfer learning with a unified text-to-		
930	text transformer . <i>Journal of Machine Learning Re-</i>		
931	<i>search</i> , 21(140):1–67.		
	Hannah Rashkin, Vitaly Nikolaev, Matthew Lamm,	932	
	Michael Collins, Dipanjan Das, Slav Petrov, Gau-	933	
	rav Singh Tomar, Iulia Turc, and David Re-	934	
	itter. 2021a. Measuring attribution in natu-	935	
	ral language generation models . <i>arXiv preprint</i>	936	
	<i>arXiv:2112.12870</i> .	937	
	Hannah Rashkin, David Reitter, Gaurav Singh Tomar,	938	
	and Dipanjan Das. 2021b. Increasing faithfulness	939	
	in knowledge-grounded dialogue with controllable	940	
	features . In <i>Proceedings of the 59th Annual Meet-</i>	941	
	<i>ing of the Association for Computational Linguistics</i>	942	
	<i>and the 11th International Joint Conference on Natu-</i>	943	
	<i>ral Language Processing (Volume 1: Long Papers)</i> ,	944	
	pages 704–718, Online. Association for Computa-	945	
	tional Linguistics.	946	
	Ehud Reiter and Craig Thomson. 2020. Shared task on	947	
	evaluating accuracy . In <i>Proceedings of the 13th In-</i>	948	
	<i>ternational Conference on Natural Language Gener-</i>	949	
	<i>ation</i> , pages 227–231, Dublin, Ireland. Association	950	
	for Computational Linguistics.	951	
	Anna Rohrbach, Lisa Anne Hendricks, Kaylee Burns,	952	
	Trevor Darrell, and Kate Saenko. 2018. Object hal-	953	
	lucination in image captioning . In <i>Proceedings of</i>	954	
	<i>the 2018 Conference on Empirical Methods in Natu-</i>	955	
	<i>ral Language Processing</i> , pages 4035–4045.	956	
	Tal Schuster, Adam Fisch, and Regina Barzilay. 2021.	957	
	Get your vitamin C! robust fact verification with con-	958	
	trastive evidence . In <i>Proceedings of the 2021 Con-</i>	959	
	<i>ference of the North American Chapter of the Asso-</i>	960	
	<i>ciation for Computational Linguistics: Human Lan-</i>	961	
	<i>guage Technologies</i> , pages 624–643, Online. Asso-	962	
	ciation for Computational Linguistics.	963	
	Thomas Scialom, Paul-Alexis Dray, Sylvain Lamprier,	964	
	Benjamin Piwowarski, Jacopo Staiano, Alex Wang,	965	
	and Patrick Gallinari. 2021. QuestEval: Summa-	966	
	rization asks for fact-based evaluation . In <i>Procee-</i>	967	
	<i>dings of the 2021 Conference on Empirical Methods</i>	968	
	<i>in Natural Language Processing</i> , pages 6594–6604,	969	
	Online and Punta Cana, Dominican Republic. Asso-	970	
	ciation for Computational Linguistics.	971	
	Thomas Scialom and Felix Hill. 2021. Beametrics: A	972	
	benchmark for language generation evaluation eval-	973	
	uation . <i>arXiv preprint arXiv:2110.09147</i> .	974	
	Thibault Sellam, Dipanjan Das, and Ankur Parikh.	975	
	2020a. BLEURT: Learning robust metrics for text	976	
	generation . In <i>Proceedings of the 58th Annual Meet-</i>	977	
	<i>ing of the Association for Computational Linguistics</i> ,	978	
	pages 7881–7892, Online. Association for Computa-	979	
	tional Linguistics.	980	
	Thibault Sellam, Amy Pu, Hyung Won Chung, Sebas-	981	
	tian Gehrmann, Qijun Tan, Markus Freitag, Dipan-	982	
	jan Das, and Ankur Parikh. 2020b. Learning to eval-	983	
	uate translation beyond English: BLEURT submis-	984	
	sions to the WMT metrics 2020 shared task . In <i>Pro-</i>	985	
	<i>ceedings of the Fifth Conference on Machine Trans-</i>	986	
	<i>lation</i> , pages 921–927, Online. Association for Com-	987	
	putational Linguistics.	988	

- 989 Uri Shaham, Elad Segal, Maor Ivgi, Avia Efrat, Ori
990 Yoran, Adi Haviv, Ankit Gupta, Wenhan Xiong,
991 Mor Geva, Jonathan Berant, and Omer Levy. 2022.
992 [Scrolls: Standardized comparison over long lan-
993 guage sequences.](#)
- 994 James Thorne, Andreas Vlachos, Oana Cocarascu,
995 Christos Christodoulopoulos, and Arpit Mittal. 2018.
996 [The fact extraction and VERification \(FEVER\)
997 shared task.](#) In *Proceedings of the First Workshop on
998 Fact Extraction and VERification (FEVER)*, pages
999 1–9, Brussels, Belgium. Association for Computa-
1000 tional Linguistics.
- 1001 Alex Wang, Kyunghyun Cho, and Mike Lewis. 2020.
1002 [Asking and answering questions to evaluate the fac-
1003 tual consistency of summaries.](#) In *Proceedings of
1004 the 58th Annual Meeting of the Association for Com-
1005 putational Linguistics*, pages 5008–5020, Online.
1006 Association for Computational Linguistics.
- 1007 Yuexiang Xie, Fei Sun, Yang Deng, Yaliang Li, and
1008 Bolin Ding. 2021. [Factual consistency evaluation
1009 for text summarization via counterfactual estimation.](#)
1010 In *Findings of the Association for Computational
1011 Linguistics: EMNLP 2021*, pages 100–110, Punta
1012 Cana, Dominican Republic. Association for Computa-
1013 tional Linguistics.
- 1014 Wenpeng Yin, Dragomir Radev, and Caiming Xiong.
1015 2021. [DocNLI: A large-scale dataset for document-
1016 level natural language inference.](#) In *Findings of
1017 the Association for Computational Linguistics: ACL-
1018 IJCNLP 2021*, pages 4913–4922, Online. Associa-
1019 tion for Computational Linguistics.
- 1020 Weizhe Yuan, Graham Neubig, and Pengfei Liu. 2021.
1021 [BARTScore: Evaluating generated text as text gen-
1022 eration.](#) In *Advances in Neural Information Process-
1023 ing Systems*.
- 1024 Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q.
1025 Weinberger, and Yoav Artzi. 2020. [Bertscore: Eval-
1026 uating text generation with bert.](#) In *International
1027 Conference on Learning Representations*.
- 1028 Yuan Zhang, Jason Baldridge, and Luheng He. 2019.
1029 [PAWS: Paraphrase adversaries from word scram-
1030 bling.](#) In *Proceedings of the 2019 Conference of
1031 the North American Chapter of the Association for
1032 Computational Linguistics: Human Language Tech-
1033 nologies, Volume 1 (Long and Short Papers)*, pages
1034 1298–1308, Minneapolis, Minnesota. Association
1035 for Computational Linguistics.
- 1036 Zheng Zhao, Shay B Cohen, and Bonnie Webber. 2020.
1037 [Reducing quantity hallucinations in abstractive sum-
1038 marization.](#) *arXiv preprint arXiv:2009.13312*.
- 1039 Chunting Zhou, Graham Neubig, Jiatao Gu, Mona
1040 Diab, Francisco Guzmán, Luke Zettlemoyer, and
1041 Marjan Ghazvininejad. 2021. [Detecting halluci-
1042 nated content in conditional neural sequence gen-
1043 eration.](#) In *Findings of the Association for Computa-
1044 tional Linguistics: ACL-IJCNLP 2021*, pages 1393–
1045 1404, Online. Association for Computational Lin-
1046 guistics.

A Additional Data Statistics

Tables 5 and 6 presents statistics regarding the length of the grounding text and the generated text for TRUE’s datasets, respectively.

Dataset	Min len.	Max len.	Median len.	Avg len.
FRANK	102	1005	550	548
SummEval	100	540	367	359
MNBM	8	10315	287	383
QAGS-CNNDM	73	360	325	318
QAGS-XSUM	218	520	339	351
BEGIN	7	64	23	23
Q^2	6	71	21	23
DialFact	4	174	22	26
PAWS	5	37	21.0	21
FEVER	8	286	44	59
VitaminC	1	265	26	28

Table 5: Grounding length statistics for TRUE.

Dataset	Min len.	Max len.	Median len.	Avg len.
FRANK	2	126	40	41
SummEval	5	133	61	63
MNBM	2	52	19	19
QAGS-CNNDM	23	85	47	49
QAGS-XSUM	9	31	18	18
BEGIN	5	40	13	14
Q^2	7	44	15	16
DialFact	4	69	16	17
PAWS	5	37	21	21
FEVER	2	36	8	8
VitaminC	1	103	12	13

Table 6: Generated text length statistics for TRUE.

Model	Avg. ROC AUC
ANLI-T5-11B	81.5 (+4.7)
ANLI-T5-Large	76.8
BLEURT-20	71.4 (+3.7)
BLEURT-20-D6	67.7
BERTScore P - deberta-xl-mnli	71.4 (+1.3)
BERTScore P - roberta-large	70.1

Table 7: Ablation study comparing the average ROC AUC results for models with different sizes. “BERTScore P” stands for BERTScore Precision.

B Implementation Details

We train all models using the t5x library.¹⁷

QG-QA For our reimplementaion of Q^2 (Honovich et al., 2021) we use T5-11B as the pretrained model for QG, QA and NLI, while Honovich et al. (2021) used T5-Base, ALBERT (Lan et al., 2019), and RoBERTa (Liu et al., 2019) for the QG, QA and NLI models, respectively. We use a maximum length of 2048 tokens for the input. We set the F1 token overlap threshold to 0.54 by tuning it on a held-out dataset. We use beam search with a beam size of 4 to generate multiple questions, and use the first question that passes the validation threshold.

¹⁷<https://github.com/google-research/t5x>

NLI We fine-tune a T5-11B model on ANLI (Nie et al., 2020) for 25K steps with a learning rate of 10^{-4} and a batch size of 32. During inference we use a maximum input length of 2048 tokens.

C Ablation Study

Table 7 presents the results of an ablation study testing the effect of model size for different model-based metrics.

D ROC Curves

Figure 2 presents the ROC curves for the different datasets studied in TRUE, using the best-performing metrics.

Dataset	Pos ROUGE_L	Neg ROUGE_L	ROUGE_L diff	Pos F1	Neg F1	F1 diff
FRANK	0.105	0.060	0.045	0.165	0.103	0.062
SummEval	0.181	0.141	0.041	0.282	0.244	0.038
MNBM	0.044	0.047	0.003	0.079	0.084	0.006
QAGS-CNNNDM	0.215	0.170	0.045	0.281	0.249	0.031
QAGS-XSUM	0.051	0.050	0.002	0.082	0.080	0.002
BEGIN	0.465	0.159	0.306	0.553	0.207	0.346
Q ²	0.228	0.169	0.059	0.368	0.264	0.104
DialFact	0.302	0.200	0.102	0.394	0.249	0.144
PAWS	0.832	0.734	0.098	0.938	0.934	0.003
FEVER	0.174	0.179	0.005	0.276	0.258	0.018
VitaminC	0.314	0.270	0.044	0.362	0.290	0.072

Table 8: Average overlap between the generated text and the grounding, measured using ROUGE-L and simple F1 token-overlap, taking the grounding to be the reference text. The “Pos” columns contain the statistics for the grounded text, while the “Neg” columns contain the statistics for the ungrounded text.

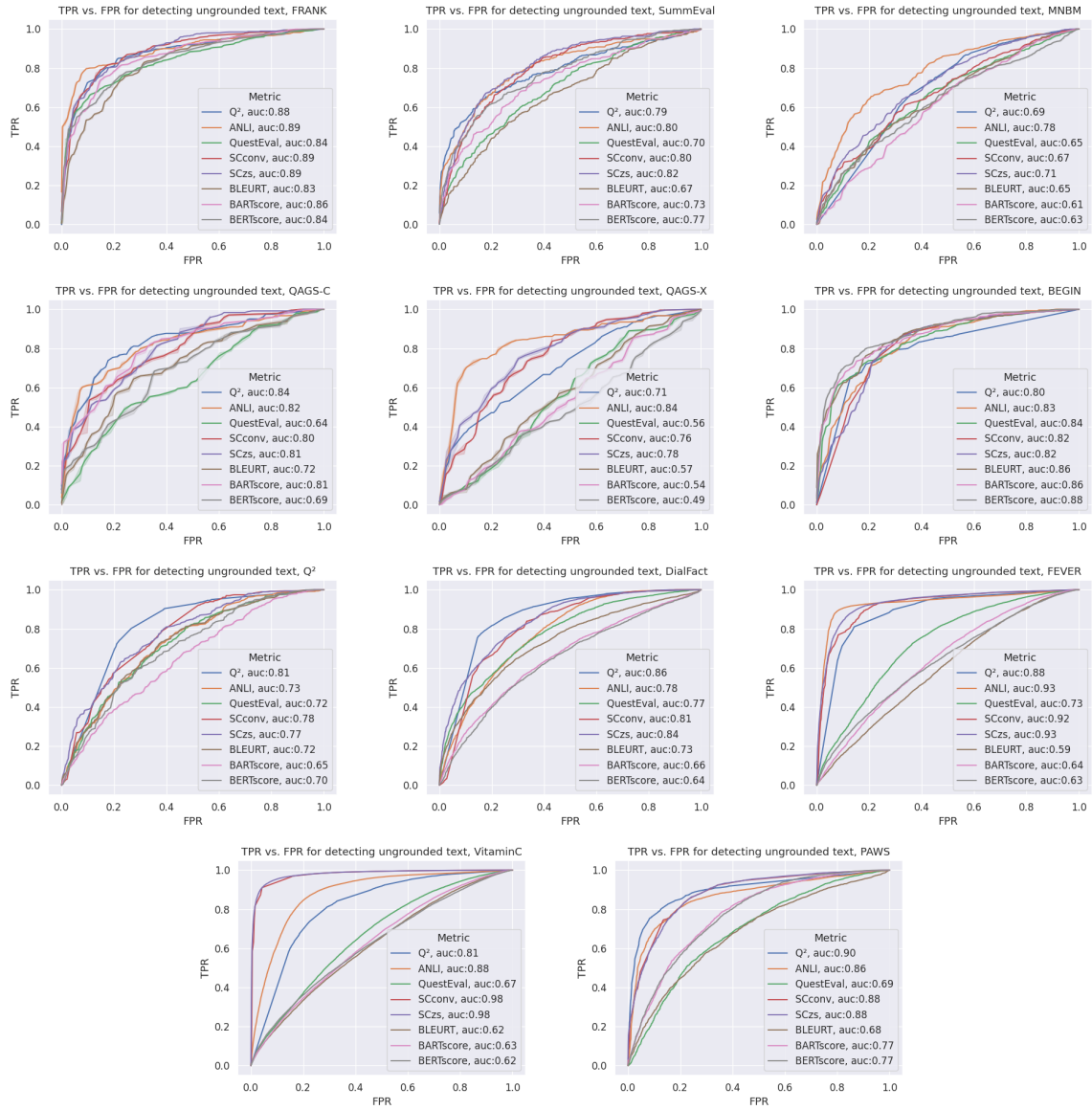


Figure 2: ROC curves for the best performing methods.

	ANLI+Q ²	ANLI+SC _{ZS}	Q ² +SC _{ZS}	SC _{Conv}	ROUGE-L	BLEU4
FRANK	89.6	91.1	90.4	88.9	80.1	78.0
SummEval	80.7	83.0	82.0	79.8	68.8	60.2
MNBM	75.6	77.1	74.6	67.2	47.5	49.3
QAGS-C	86.0	84.7	86.4	79.6	67.1	63.9
QAGS-X	81.8	85.1	79.3	76.1	52.9	48.6
BEGIN	85.7	82.1	85.7	81.6	86.4	84.6
Q²	83.0	76.9	83.9	77.5	66.8	64.3
DialFact	89.4	84.5	90.2	81.2	71.2	72.5
PAWS	90.5	89.7	91.4	88.2	82.2	77.3
FEVER	94.0	94.6	93.9	86.7	49.9	51.1
VitaminC	90.3	96.4	96.5	97.5	59.9	59.6
Avg. w/o VitC, FEVER	84.7	83.8	84.9	80.0	69.2	66.5

Table 9: ROC AUC results for metrics that were not reported in Table 3.

	Ensemble	Q ²	ANLI	SC _{ZS}	BLEURT	QuestEval	FactCC	BART _{score}	BERT _{score}
FRANK	90.8	87.8	89.2	88.6	83.2	86.4	73.9	88.3	86.0
BEGIN	85.9	78.0	82.8	84.2	82.2	81.4	65.0	83.7	86.0
DialFact	88.6	85.0	75.9	82.1	72.2	76.3	55.1	65.5	64.3
PAWS	92.4	90.1	87.3	89.7	67.1	70.1	65.1	77.3	76.4
VitaminC	96.7	83.4	89.6	98.4	63.0	67.8	56.8	64.1	63.5
Avg. w/o VitC	89.4	85.2	83.8	86.2	76.2	78.5	64.8	78.7	78.2

Table 10: ROC AUC results for the different metrics on the TRUE test set. We exclude VitaminC from the average calculation as SC_{ZS} was trained on VitaminC. The highest score in each row (excluding the Ensemble) is in bold and the aforementioned SC results are in strikethrough.

	Ensemble	Q ²	ANLI	SC _{ZS}	BLEURT	QuestEval	FactCC	BART _{score}	BERT _{score}
FRANK	83.0	81.5	82.0	79.0	76.6	73.0	72.1	80.7	75.6
BEGIN	76.8	74.1	76.8	78.9	74.3	73.4	62.09	74.8	78.1
DialFact	80.9	78.1	68.4	74.2	67.1	69.0	52.5	58.6	60.2
PAWS	84.8	84.1	82.1	82.3	62.9	64.8	60.7	70.9	69.8
VitaminC	92.1	77.5	83.9	94.2	59.0	63.3	55.5	59.8	58.0
Avg. w/o VitC	81.4	79.4	77.3	78.6	70.2	70.0	62.1	71.3	70.9

Table 11: Accuracy results for the different metrics on the TRUE test set. Thresholds were tuned on the corresponding development sets. We exclude VitaminC from the average calculation as SC_{ZS} was trained on VitaminC. The highest score in each row (excluding the Ensemble) is in bold and the aforementioned SC results are in strikethrough.