Benchmarking Answer Verification Methods for Question Answering-Based Summarization Evaluation Metrics

Anonymous ACL submission

Abstract

answering-based Question summarization evaluation metrics must automatically determine whether the QA model's prediction is correct or not, a task known as answer verification. In this work, we benchmark the lexical answer verification methods which have been used by current QA-based metrics as well as two more sophisticated text comparison methods, BERTScore and LERC. We find that LERC out-performs the other methods in some settings while remaining statistically indistinguishable from lexical overlap in others. However, our experiments reveal that improved verification performance does not necessarily translate to overall QA-based metric quality: In some scenarios, using a worse verification method - or using none at all - has comparable performance to using the best verification method, a result that we attribute to properties of the datasets.

1 Introduction

001

006

016

017

034

040

A recent trend in summarization metrics is evaluating the quality of a summary via question answering (QA; Eyal et al., 2019; Scialom et al., 2019, 2021; Durmus et al., 2020; Wang et al., 2020; Deutsch et al., 2021a). These metrics compare the semantic content of two texts (e.g., the reference and candidate summaries) by generating questions from one and answering those questions against the other. The amount of common semantic content is proportional to the number of questions which are answered correctly.

A critical step of QA-based evaluation metrics is to verify whether the QA model's prediction is correct, a task known as answer verification (see Fig. 1). This helps to both suppress noisy output from the QA model as well as identify inconsistent information across the texts.

Answer verification is typically done by comparing the prediction to the expected answer by the exact match or token F_1 string comparison methods



Figure 1: In the answer verification task, the evaluation metrics score how likely two phrases (one the ground-truth answer and one the QA model's prediction) from different contexts have the same meaning.

(Rajpurkar et al., 2016). However, more sophisticated text comparison methods have been proposed in recent years, and it is unknown whether they provide a benefit in this particular scenario. 042

043

044

046

047

051

057

059

060

061

062

063

064

065

In this work, we benchmark various answer verification strategies for QA-based summarization evaluation metrics. Our goal is to understand whether methods that are more advanced than lexical overlap are better able to classify phrases as having the same or different meaning as well as whether any such improvements result in the overall QA-based metric being better at replicating human judgments of summary quality.

We analyze four answer verification methods, exact match, token F_1 , BERTScore (Zhang et al., 2020), and LERC, (Chen et al., 2020) in combination with two QA-based metrics, QAEval (Deutsch et al., 2021a) and FEQA (Durmus et al., 2020).

Based on a set of human annotations across two datasets, we find that LERC, in general, performs the best at the actual task of answer verification, although in some settings it is statistically indistinguishable from token F_1 (§4.1). However, our results also show that any such improvement in verification performance does not always translate 072

067

2

078

084 086

096

099

102 103

101

104 105

106

107

111

112

114

115

108 109 110

3 113

We define the answer verification task as the follow-

Definitions & Methods

the scope of this work.

to a better QA-based evaluation metric ($\S4.2$). which the QA pair was generated, a prediction, and We believe these results can be explained by the target text the prediction comes from, score properties of the QA metrics and the datasets.

When the QA model performance is high or the

verification task is in some sense easy to do, it may

not be necessary to have a sophisticated verifica-

tion method or even use one at all. Despite this, our

recommendation is still to do answer verification

with LERC as it can only improve performance,

although token F₁ may suffice in some situations.

The majority of summarization evaluation metrics

can be viewed as estimating how similar in mean-

ing two pieces of text are. For instance, ROUGE

(Lin, 2004) does this by calculating the number of overlapping n-grams between the two texts.

Instead of directly comparing the entire texts,

QA-based metrics identify specific phrases within

the texts which should be compared, as follows.

First, a set of questions is automatically generated from one text. Then, those questions are automat-

ically answered against a second text to obtain a

set of predicted answers. The final score is pro-

portional to the number of correct predictions, but

determining whether those predictions are correct -

the task of answer verification – is done by compar-

ing the text of the prediction to the expected answer.

Therefore, instead of directly comparing the entire

contents of the two texts, QA-based metrics instead

reduce the scope of the problem to only comparing

Current QA-based metrics perform the answer

verification step by lexical comparison, either ex-

act match or token F1. Such metrics include QA-

Eval (Deutsch et al., 2021a), FEQA (Durmus et al.,

2020), and more (Eyal et al., 2019; Wang et al.,

2020; Scialom et al., 2019, 2021). However, any

such function which calculates the similarity of

arbitrary text can be used instead. This includes

embedding-based methods such as BERTScore

(Zhang et al., 2020) or metrics which have been

trained specifically to do this task, such as LERC

(Chen et al., 2020). Evaluating how these methods

perform as answer verification methods for QA-

based metrics compared to the lexical baselines is

ing: Given a question, answer, the source text from

specific pairs of phrases.

Related Work & Background

how similar the meanings of the answer and prediction are (see Fig. 1 for an example).¹ Answer verification is used by QA-based metrics to suppress noisy outputs from the QA model as well as identify when the OA prediction is correct with respect to the target text but incorrect with respect to the expected answer (e.g., unfaithful information). 116

117

118

119

120

121

122

123

124

125

126

127

128

129

130

131

132

133

134

135

136

137

138

139

140

141

142

143

144

145

146

147

148

149

150

151

152

153

154

155

156

157

We analyze four different answer verification methods.

Exact Match The exact match (EM) method compares the two phrases to see if they are identical (after light normalization). EM assigns a score of 1 if the phrases are identical and 0 otherwise.

Token F¹ The token F¹ comparison calculates an F_1 score based on the number of unigrams the two phrases have in common. This is equivalent to the F₁ variant of ROUGE-1.

BERTScore BERTScore (Zhang et al., 2020) compares two pieces of text by aligning the texts' tokens according to which pairs have the highest BERT embedding cosine similarity. We adapt BERTScore to answer verification by encoding the answer and prediction using their respective contexts, then calculating the BERTScore only between the two phrase encodings. Since the output of BERTScore is often in a narrow range of values, we rescale the scores by defining 0 and 1 as the 2.5th and 97.5th percentiles of the BERTScores calculated over the whole dataset.

LERC Chen et al. (2020) proposed LERC, a learned metric for scoring how similar the expected and predicted answers to a question are conditioned on the question and the target text the prediction comes from. All of the inputs are jointly encoded using a BERT-based classifier, which was finetuned on human annotations of meaning similarity. Because it was designed for scoring reading comprehension predictions, it does not use the source text. We rescale the output from LERC to be in the range [0, 1].

¹This is slightly different from the task defined by Chen et al. (2020) which does not include the source text because no such text exists in the standard definition of the reading comprehension task. However, we include it because the source text can be used to create a representation for the answer which may be better than using the question alone.

4 Experiments

158

160

161

162

164

165

168

169

170

171

172

173

174

175

176

177

178

179

181

182

184

187

188

189

190

193

194

195

196

197

198

199

201

202

206

The answer verification methods are evaluated independently (§4.1) as well as in combination with two QA-based metrics (§4.2), QAEval (Deutsch et al., 2021a) and FEQA (Durmus et al., 2020). QAEval measures the content quality of a summary (does the summary contain "summary-worthy" information) by using a reference summary as the source text and candidate summary as the target text. In contrast, FEQA estimates the faithfulness of the summary (does the summary contain information consistent with the input) by using the candidate summary as the source text and the input document as the target text.

The experiments are run on two datasets, TAC'08 (Dang and Owczarzak, 2008) and Summ-Eval (Fabbri et al., 2021). These datasets have summaries generated by 58 and 16 models for 48 and 100 inputs, respectively, which are annotated with expert judgments. Both QAEval and FEQA are evaluated on SummEval because it contains annotations for both summary quality and faithfulness, whereas only QAEval is evaluated on TAC'08 since it does not have faithfulness judgments.²

4.1 Answer Verification Performance

First, we examine how well each answer verification method accurately scores manually labeled answer pairs from the summarization datasets. For each QA metric and dataset combination, we ran the metric on the summaries, then randomly sampled 200 QA predictions (making 600 total). Each prediction and expected answer were manually annotated by the authors for whether or not the two phrases share the same meaning.

Ideally, the answer verification methods should both successfully classify phrases based on their meaning as well as provide a score close to 1 for phrases with the same meaning and close to 0 with different meanings. These properties are quantified by the binary classification accuracy (assigning labels based on a threshold which maximizes this score) as well as the mean squared error (MSE) of the predicted scores, show in Table 1.

We find that LERC is the only method with the best (or tied for the best) performance across all three metric-dataset combinations. Despite LERC's significant improvement on the SummEval data with QAEval predictions, it is statistically indistinguishable from F_1 on the same dataset with

ion.

		QA	FE	FEQA		
Ans. Verif.	TAC'08		SummEval		SummEval	
	Acc	MSE	Acc	MSE	Acc	MSE
Majority Cls	51.5	.49	78.5	.22	56.5	.44
EM	64.5	.36	78.5	.46	76.0	.24
F_1	84.0	.19	79.5	.25	91.0	.10
BERTScore	<u>81.0</u>	.16	79.5	.20	82.5	.16
LERC	85.0	.13	88.0	.11	<u>88.5</u>	.09

Table 1: The binary accuracies and mean squared errors of the answer verification methods evaluated on three metric-dataset combinations with 200 manually labeled examples each. Underlined values are statistically indistinguishable from those in bold under a single-tailed pairwise permutation test with $\alpha = 0.05$.

FEQA predictions. We believe this can be explained by which texts are being compared for each metric. FEQA compares the generated summary to the input document. Recent summarization models are known to copy heavily from the input with little high-level abstraction or rephrasing, so comparing phrases with token F_1 is likely to be quite successful. In contrast, QAEval compares the reference and generated summaries. The reference summaries are written by humans, and thus more likely to contain information from the input document which is expressed differently. In such a scenario, the learned metric, LERC, shows strong improvements over F_1 .

In general, we find that when BERTScore and LERC do improve over F_1 , they do so by identifying paraphrases that have no tokens in common, which sometimes requires world knowledge. Examples of this are included in Appendix B.

4.2 Overall Metric Evaluation

Next, we investigate whether the differences in classification performance of the verification methods translate to downstream improvements in the overall quality of the QA-based metrics. To do so, we evaluate different variants of the metrics that use each answer verification method. For both QAEval and FEQA, the final score for the summary is the output of the answer verification method averaged over all of the QA pairs.³

QAEval For QAEval, we report the standard system- and summary-level correlations of the metrics' scores to human judgments in Table 2 (due to space constraints, we refer the reader to Deutsch

237

239

207

³QAEval can also predict a question is unanswerable. In such cases, the score of the prediction is 0.

et al. (2021b) for definitions of the correlations). We also compare against the standard BERTScore and ROUGE metrics as well as a QAEval variant which uses no answer verification by always marking the phrases as correct if the QA model predicts the question is answerable, denoted QAEval-IsAns.

240

241

242

244

245

247

249

254

256

257

260 261

267

269

270

273

274

278

281

283

287

288

In general, all of the answer verification methods work comparably well, although BERTScore and LERC do statistically improve over the lexical methods in some settings, but not by large margins. We believe the performance of QAEval-IsAns offers an explanation as follows.

Answer verification is not necessary if the QA model is perfect and the summaries are faithful (i.e., the QA prediction is always correct). For SummEval, Deutsch et al. (2021a) demonstrated that QAEval's QA performance was reasonable, and the summaries are very faithful with an average consistency score of 4.7 / 5 according to Fabbri et al. (2021). Therefore, it may be difficult to demonstrate an improvement with any answer verification method even if it is high quality since the need for answer verification is low. Indeed, we see QAEval-IsAns statistically ties the best methods.

On TAC'08, we expect it should be easier to show answer verification helps since Deutsch et al. (2021a) showed the QA performance is poor, suggesting answer verification could help to suppress noisy predictions. Indeed, we do see QAEval-IsAns is statistically out-performed by the verification methods. We suspect the improvements are larger at the system-level than the summarylevel because the system quality is estimated over a larger number of QA pairs than an individual summary's quality is. A larger number of questions reduces any noise introduced by the verification methods, resulting in a more accurate estimate of summary quality and a better metric.

FEQA We report the direct correlations between the FEQA metrics' scores and human judgments across all of the summaries in Table 3, including those for ROUGE, BERTScore, as well as FactCC (Kryscinski et al., 2020). FactCC is a learned model to predict the factual consistency between two texts that was trained on synthetically generated data.

Among the FEQA variants, F_1 is the best or indistinguishable from LERC. This result is expected given how similarly they perform at answer verification on this QA metric and dataset split. This is again likely due to the fact that the summarization models copy heavily from the input documents, so

Motric	TAC	C'08	Sum	SummEval		
WICHIC	Sys	Sum	Sys	Sum		
BERTScore	$.68^{\dagger}$	$.40^{\dagger}$.75†	.27†		
ROUGE-1	.60	.39†	.50	.20		
ROUGE-2	.67	$.39^{\dagger}$.43	.14		
QAEval-IsAns	.63	.37	$.70^{\dagger}$.26 [†]		
QAEval-EM	.74 [†]	.29	$.77^{\dagger}$.19		
QAEval-F1	.68	.36	$.77^{\dagger}$.22		
QAEval-BERTScore	$.68^{\dagger}$	<u>.38</u> †	$.77^{\dagger}$	$.26^{\dagger}$		
QAEval-LERC	$.68^{\dagger}$.39 [†]	.80 †	<u>.24</u> [†]		

Table 2: System- and summary-level Kendall's τ (results with Pearson and Spearman are included in Appendix A). Underlined QAEval values are statistically indistinguishable from the best in bold. Values marked with † are statistically indistinguishable from the best metric overall. Statistical testing done using the single-tailed PERM-BOTH permutation test (Deutsch et al., 2021b) with $\alpha = 0.05$.

Metric	r	ρ	au
ROUGE-1	.13	.13	.11
ROUGE-2	.25	.25	.19
BERTScore	.17	.17	.14
FactCC	$.34^{\dagger}$	$.36^{\dagger}$.29†
FactCCX	.29	.31	.24
FEQA-EM	.17	.14	.11
$FEQA-F_1$.20	.16	.13
FEQA-BERTScore	.15	.12	.10
FEQA-LERC	.18	<u>.15</u>	<u>.12</u>

Table 3: The Pearson r, Spearman ρ , and Kendall τ correlations on the SummEval dataset. Values in bold are the best FEQA variants with those underlined being statistically indistinguishable. \dagger marks the best results across all metrics.

the expected answers and QA model predictions are likely to be quite lexically similar. Overall, the FEQA correlations are still lower than those by FactCC by a large margin.

5 Conclusion

In this work, we benchmarked four different answer verification methods for QA-based summarization evaluation metrics. Although we were able to identify some methods perform better than others at verification, any such improvement does not necessarily translate a better overall metric quality. We hypothesize that several factors, including the quality of the QA model and properties of the datasets, likely explain this result. Despite this, our recommendation is that practitioners use LERC, although token F_1 may be sufficient in some scenarios.

299

300

302

303

304

305

306

291

References

307

308

309

310

311

312

313

314

315

319

321

322

323

324

326

327

328

329

337

339

340

341

342

344

345

347

351 352

353

354

357

- Anthony Chen, Gabriel Stanovsky, Sameer Singh, and Matt Gardner. 2020. MOCHA: A Dataset for Training and Evaluating Generative Reading Comprehension Metrics. In Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP), pages 6521–6532, Online. Association for Computational Linguistics.
- Hoa Trang Dang and Karolina Owczarzak. 2008. Overview of the TAC 2008 Update Summarization Task. In *Proc. of the Text Analysis Conference* (*TAC*).
- Daniel Deutsch, Tania Bedrax-Weiss, and Dan Roth. 2021a. Towards Question-Answering as an Automatic Metric for Evaluating the Content Quality of a Summary. *Transactions of the Association for Computational Linguistics*, 9:774–789.
- Daniel Deutsch, Rotem Dror, and Dan Roth. 2021b. A Statistical Analysis of Summarization Evaluation Metrics using Resampling Methods. *Transactions* of the Association for Computational Linguistics, 9.
- Esin Durmus, He He, and Mona Diab. 2020. FEQA:
 A Question Answering Evaluation Framework for Faithfulness Assessment in Abstractive Summarization. In Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics, ACL 2020, Online, July 5-10, 2020, pages 5055–5070. Association for Computational Linguistics.
- Matan Eyal, Tal Baumel, and Michael Elhadad. 2019. Question Answering as an Automatic Evaluation Metric for News Article Summarization. In Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL-HLT 2019, Minneapolis, MN, USA, June 2-7, 2019, Volume 1 (Long and Short Papers), pages 3938–3948. Association for Computational Linguistics.
- A. R. Fabbri, Wojciech Kryscinski, Bryan McCann, R. Socher, and Dragomir Radev. 2021. SummEval: Re-evaluating Summarization Evaluation. *Transactions of the Association for Computational Linguistics*, 9:391–409.
- Wojciech Kryscinski, Bryan McCann, Caiming Xiong, and Richard Socher. 2020. Evaluating the Factual Consistency of Abstractive Text Summarization. In Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP), pages 9332–9346, Online. Association for Computational Linguistics.
- 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.
- Pranav Rajpurkar, Jian Zhang, Konstantin Lopyrev, and Percy Liang. 2016. SQuAD: 100,000+ Questions

for Machine Comprehension of Text. In Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing, EMNLP 2016, Austin, Texas, USA, November 1-4, 2016, pages 2383–2392. The Association for Computational Linguistics. 363

364

366

367

369

370

372

373

374

375

376

377

378

379

381

382

383

384

387

388

389

390

391

392

393

394

- Thomas Scialom, Paul-Alexis Dray, Patrick Gallinari, Sylvain Lamprier, Benjamin Piwowarski, Jacopo Staiano, and Alex Wang. 2021. Questeval: Summarization Asks for Fact-Based Evaluation. *arXiv preprint arXiv:2103.12693*.
- Thomas Scialom, Sylvain Lamprier, Benjamin Piwowarski, and Jacopo Staiano. 2019. Answers Unite! Unsupervised Metrics for Reinforced Summarization Models. 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 3246–3256, Hong Kong, China. Association for Computational Linguistics.
- Alex Wang, Kyunghyun Cho, and Mike Lewis. 2020. Asking and Answering Questions to Evaluate the Factual Consistency of Summaries. In Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics, ACL 2020, Online, July 5-10, 2020, pages 5008–5020. Association for Computational Linguistics.
- Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q. Weinberger, and Yoav Artzi. 2020. BERTScore: Evaluating Text Generation with BERT. In 8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020. OpenReview.net.



Figure 2: The distributions of score values for three metrics on the SummEval dataset for ground-truth answer and QA model prediction pairs from QAEval with the same (blue) and different (orange) meanings.

A Additional Results

Fig. 2 contains the distributions of score values for token F_1 , BERTScore, and LERC on the Summ-Eval dataset grouped by phrases that have and do no have the same meaning. LERC most confidently separates the positive and negative examples. F_1 performs similarly, except it fails in a large number of cases when the two phrases have no tokens in common. BERTScore tends to mix the scores of the positive and negative classes, although they are separated on average.

In Table 4, we report the system- and summarylevel correlations on TAC'08 and SummEval with Pearson's r and Spearman's ρ correlation coefficients in addition to the Kendall's τ which was presented in the main body of the paper. The other coefficients lead to a similar conclusion to that which we made with Kendall's τ : All answer verification methods perform comparably well, and when BERTScore or LERC does improve over a lexical baseline, it is not by a large margin. Further, using no verification method (QAEval-IsAns) largely performs equally well as QAEval variants which do use a verification step on the SummEval dataset, but not on TAC'08.

B Example BERTScore/LERC Improvements

Table 5 contains example expected answer and QA model prediction pairs for which BERTScore and LERC improve over exact match and token F_1 . We see that the improvements come from better identifying when the phrases are paraphrases of each other, which sometimes involves world knowledge.

417

418

419

420

421

422

423

424

425

426

427

428

396

400

401

402

403 404

405

406

407

408

409

TAC'08						SummEval						
Metric	System-Level			Summary-Level			System-Level			Summary-Level		
	r	ρ	au	r	ρ	au	r	ρ	τ	r	ρ	au
BERTScore	.83	$.85^{\dagger}$	$.68^{\dagger}$	$.50^{\dagger}$	$.50^{\dagger}$	$.40^{\dagger}$.84†	.91†	.75†	.37†	.35†	.27†
ROUGE-1	.79	.80	.60	.49†	$.48^{\dagger}$.39†	.61	.62	.50	.28	.26	.20
ROUGE-2	.83	$.87^{+}$.67	$.48^{\dagger}$	$.48^{\dagger}$.39†	.64	.60	.43	.23	.19	.14
ROUGE-L	.74	.77	.57	.46	.45	.36	.61	.48	.32	.21	.18	.14
ROUGE-SU4	.80	.83	.63	.49†	$.48^{\dagger}$.39†	.62	.56	.38	.23	.19	.15
QAEval-IsAns	.87	.82	.63	$.48^{\dagger}$.47	.37	.76	<u>.86</u> †	<u>.70</u> †	.33†	<u>.32</u> †	$.26^{\dagger}$
QAEval-EM	.92 †	.89 †	.74 †	.35	.35	.29	<u>.80</u> †	<u>.91</u> †	<u>.77</u> †	.23	.23	.19
QAEval-F1	<u>.90</u> †	<u>.86</u> †	.68	.46	.45	.36	<u>.82</u> †	<u>.91</u> †	<u>.77</u> †	.30	.29	.22
QAEval-BERTScore	<u>.90</u> †	<u>.85</u> †	<u>.68</u> †	<u>.49</u> †	<u>.48</u> †	<u>.38</u> †	.84 †	<u>.89</u> †	<u>.77</u> †	.36 †	.34 †	.26†
QAEval-LERC	<u>.89</u> †	$.85^{\dagger}$	$.68^{\dagger}$.50 [†]	.49 [†]	.39 [†]	$.81^{\dagger}$.93 †	.80 †	.33†	<u>.31</u> [†]	<u>.24</u> †

Table 4: System- and summary-level correlations using Pearson's r, Spearman's ρ , and Kendall's τ .

Answer	Prediction	BSc	LERC
EU	European Union	0.73	0.84
a smaller leftist guerilla group	National Liberation Army	0.48	0.10
six-time Olympic gold medalist	Usain Bolt	0.34	0.35
Luis Enrique's side	Barcelona	0.40	0.18
emergency responders	paramedics	0.20	0.67
the child	toddler	0.38	0.45

Table 5: Examples where BERTScore and LERC improve over F_1 (all examples have an F_1 score of 0). Successfully classing these phrases requires paraphrasing (e.g., "the child" and "toddler") and, in some cases, world knowledge (e.g., Usain Bolt had won six gold medals when the article was written).