

---

# How not to Lie with a Benchmark: Rearranging NLP Leaderboards

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 Comparison with a human is an essential requirement for a benchmark for it to  
2 be a reliable measurement of model capabilities. Nevertheless, the methods for  
3 model comparison could have a fundamental flaw - the arithmetic mean of separate  
4 metrics is used for all tasks of different complexity, different size of test and training  
5 sets.

6 In this paper, we examine popular NLP benchmarks' overall scoring methods and  
7 rearrange the models by geometric and harmonic mean (appropriate for averaging  
8 rates) according to their reported results. We analyze several popular benchmarks  
9 including GLUE, SuperGLUE, XGLUE, and XTREME. The analysis shows that  
10 e.g. human level on SuperGLUE is still not reached, and there is still room for  
11 improvement for the current models.

## 12 1 Introduction

13 The benchmarking approach has a rich history throughout computer science and is now the leading  
14 method in machine learning progress validation. In the field of Natural Language Processing (NLP),  
15 there exist at least 898 benchmarks<sup>1</sup>, the most prominent being GLUE, SuperGLUE, XGLUE, etc.,  
16 created within a single paradigm.

17 The increase in the number of publications and developments in the field of machine learning has led  
18 to the need for methodological development of standards for describing models and all stages of the  
19 experiment, including the collection and processing of preliminary data for training, reproducibility  
20 of results, testing conditions, and most importantly, the creation of common measurable criteria for  
21 evaluating intelligent systems, both natural and artificial (Chollet, 2019). Modern NLP benchmarks  
22 are substantively inherit to Turing test, i.e. test the model abilities with various intellectual tasks  
23 expressed with texts, and methodologically inherit the benchmark approach for measuring the  
24 computing performance, like SPEC<sup>2</sup>.

25 The delicate question of a general assessment of the model results on all the tasks is often solved  
26 by a method that is unforgivably simple for such a responsible task - the arithmetic mean for all  
27 the tasks. This does not take into account the scatter of results on different tasks, the different size  
28 of the task test sets (e.g. in SuperGLUE they differ a hundred times, compare 146 test samples in  
29 Winograd Schema and 10'000 test samples in ReCoRd(Wang et al., 2019a)), different susceptibility  
30 to leaks (Elangovan et al., 2021), including year of creation (Recognizing Textual Entailment data  
31 was collected in 2005 (Dagan et al., 2005), while BoolQ or CommitmentBank data was collected in  
32 2019(Clark et al., 2019), (De Marneffe et al., 2019)).

---

<sup>1</sup>according to <https://paperswithcode.com/area/natural-language-processing>

<sup>2</sup><https://www.spec.org/benchmarks.html>

33 In this article we present an analysis of the NLP benchmarks' results, using not the arithmetic mean,  
34 but the other metrics: geometric mean and harmonic mean. As F1 (harmonic mean) is frequently  
35 used to normalize Precision and Recall as they are fractions, optimizing the classifier threshold to  
36 maximize F1 leads to a more balanced balance between metrics than the arithmetic mean because it  
37 penalizes systems more for the smaller values(Sasaki et al., 2007). The geometric mean, as noted  
38 in the (Fleming and Wallace, 1986), is the preferred metric to the arithmetic mean in computing  
39 performance benchmarks when it comes to normalized values and percentages. The results change the  
40 usual idea of the models' order on leaderboards: humans still occupy the first place in the intellectual  
41 task solving, and the best results (1.5-2% worse than humans) belong to DeBerta(He et al., 2020),  
42 T5+Meena (Raffel et al., 2020) and McAlbert+DKM models<sup>3</sup>. Thus, the contribution of this paper is  
43 two-fold: 1) we present the reviewed approach model evaluation on multiple tasks 2) we re-arrange  
44 the currently existing leaderboards of most popular benchmarks.

45 This work is organized as follows: section 2 presents previous work on the topic, it is followed by  
46 section 3 with a description of the general methodology of cross-checking the results, including the  
47 scores for each benchmark. Next, an analysis of the results and discussion are presented in section 4,  
48 as well as a conclusion in section 5.

## 49 2 Previous Work

50 Evaluation and comparison of NLP models beget a rich history, rising with the Turing test (Turing,  
51 2009). The question-answering approach then evolved into the SQuAD task (Rajpurkar et al., 2016),  
52 comparing systems and annotator results by their ability to find answers to informative questions.  
53 The next step in the development and assessment of intelligent systems belongs to the benchmark  
54 methodology, which aims to bring the solution of the Natural Language Understanding problem  
55 closer - General Language Understanding Evaluation(Wang et al., 2019b).

56 The General Language Understanding Evaluation (GLUE) methodology includes:

- 57 1. a benchmark from N (11 in original GLUE) intellectual tasks of understanding a natural  
58 language, with a fixed division into training, validation and test data;
- 59 2. a set of diagnostic data designed exclusively for testing and analyzing the results of trained  
60 systems in relation to a wide range of categories found at various levels of natural language  
61 (morphological, lexical, syntactic, semantic);
- 62 3. averaged human performance evaluation on the tasks;
- 63 4. publicly available rating system and codebase to quickly reproduce results from publicly  
64 available systems and self-evaluate models.

65 GLUE has developed, by no means, a prolific method to model evaluation, and has already been  
66 reproduced several times in new language material: in Chinese (Xu et al., 2020), Korean (Park  
67 et al., 2021), Russian (Shavrina et al., 2020), Polish (Rybak et al., 2020), and French (Le et al.,  
68 2019) languages, and also jumpstarted two multilingual projects: XGLUE (Liang et al., 2020) and  
69 XTREME (Hu et al., 2020).

70 As stated in (Wang et al., 2019a), *"Lacking a fair criterion with which to weight the contributions of*  
71 *each task to the overall score, we opt for the simple approach of weighing each task equally, and for*  
72 *tasks with multiple metrics, first averaging those metrics to get a task score."* All the GLUE-based  
73 benchmarks follow this methodology.

74 However, apart from the GLUE format, other benchmarks have provided several alternatives to  
75 evaluate the overall model contribution.

76 KILT, a Benchmark for Knowledge Intensive Language Tasks (Petroni et al., 2020), avoids calculating  
77 the overall result, and also do not compare the result with the human level, but only provides metrics  
78 for individual tasks.

79 DecaNLP (McCann et al., 2018) makes a rating using not the average, but the sum of points for all  
80 tasks. This approach allows balancing the contributions of different tasks to the overall metric.

---

<sup>3</sup><https://www.iflytek.com/news/2118>

### 81 3 Method

82 We arrange the NLP benchmark results using the publicly available model scores for all the tasks to  
83 calculate new overall scores.

84 Other available options from Pythagorean means - harmonic mean and geometric mean - can also  
85 be considered: we have centered our research around 2 simple statistics that are widely used for  
86 averaging fractions (Iun Chou, 1969) or normalized values (Fleming and Wallace, 1986) among the  
87 possible alternatives. We did not consider other measures of central tendency, like median and mode,  
88 as the averaged samples more often consist of about 10 measurements, and on them such metrics can  
89 give the same results for competing systems.

- 90 • The arithmetic mean (AM) described in eq. 1 is calculated as the sum of the task scores (Xs)  
91 divided by the total number of tasks, referred to as N.
- 92 • The geometric mean (GM) described in eq. 2 is calculated as the N-th root of the product of  
93 all task scores (with the above conditions), where N is the number of values.
- 94 • The harmonic mean (HM) described in eq. 3 is calculated as the number of values N divided  
95 by the sum of the reciprocal of the values

$$AM = \frac{1}{n} \sum_{i=1}^n x_i = \frac{x_1 + x_2 + \dots + x_n}{n} \quad (1)$$

96

$$GM = \left( \prod_{i=1}^n x_i \right)^{\frac{1}{n}} = \sqrt[n]{x_1 x_2 \dots x_n} \quad (2)$$

$$HM = \frac{N}{\frac{1}{x_1} + \frac{1}{x_2} + \dots + \frac{1}{x_N}} \quad (3)$$

97 The sections below present the results of the leaderboard re-weighting as of May 2021.

#### 98 3.1 Reevaluating the Benchmarks

99 The GLUE (*11 tasks for English*), SuperGLUE (*10 tasks for English*), XGLUE (*11 tasks for 19*  
100 *languages*), and XTREME (*4 tasks for 40 languages*) provide different scoring metrics for each task,  
101 including Accuracy, F1, Matthew’s correlation coefficient, Exact Match, while the overall score is  
102 calculated by their simple average. In cases like these, the geometric mean is appropriate when the  
103 data contains values with different units of measure (Iun Chou, 1969).

104 The harmonic mean of the task results as a better overall metric has the same grounding as introduction  
105 of the F-1 measure over precision and recall (Sasaki et al., 2007): the harmonic mean is more intuitive  
106 than the arithmetic mean when computing a mean of ratios. Given the set of metrics with a large  
107 scatter, the harmonic mean will be less than the arithmetic mean, penalizing the system more for the  
108 errors made.

109 As stated in (Dittmann and Maug, 2008), *"error measures are inherently subjective as they are*  
110 *determined by the loss function of the researcher or analyst who needs to choose a valuation*  
111 *procedure. Therefore, our analysis cannot establish which error measure should be used. Instead,*  
112 *our objective is to highlight the effects of the choice of error measure, so that researchers and*  
113 *analysts alike can draw their own conclusions about the error measure and, eventually, about the*  
114 *valuation methods they wish to use."* Nevertheless, the arithmetic, geometric and harmonic mean  
115 are all differently subjected to outliers in a data sample. The arithmetic mean always results in a  
116 higher values than the geometric mean or the harmonic mean, and the harmonic mean always results  
117 in a lower estimate than the geometric mean (Xia et al., 1999), thus, the harmonic and geometric  
118 mean tend more strongly toward the least values, and tend to mitigate the impact of large outliers and  
119 aggravate the impact of small ones.

120 The harmonic mean is the appropriate mean if the data is comprised of rates, while the geometric  
121 mean is used as an unbiased estimation when working with normalized ratios, for example, in  
122 finance (Dittmann and Maug, 2008) or computing benchmarks (Fleming and Wallace, 1986).

123 However, their applicability to a better summarization of the model performance to a single number  
 124 has been widely discussed, see (Smith, 1988), discussing performance computing:

- 125 • the **harmonic mean** is considered the appropriate metric to summarize benchmark results  
 126 expressed as rates,
- 127 • while **geometric mean** is applicable in case of the use of performance numbers that are  
 128 normalized with respect to one of the results being compared (see 4),
- 129 • and **arithmetic mean** should not be used as a summarizing metric with rates, making it the  
 130 worst choice for results accumulation.

131 When measuring the geometric and harmonic mean, the following assumptions were used:

- 132 • as the geometric mean does not accept zero values (also negative ones), we filled the cases  
 133 of metric lacking (for example, Sentence Retrieval in XTREME) with 0.00001 values;
- 134 • tasks with more than one metric measured (e.g. Accuracy and F1), are subjected to the aver-  
 135 aging operation when measuring the total score, as in the standard GLUE methodology: the  
 136 arithmetic mean of all metrics is taken for each task. Another modification of measurement  
 137 is potentially possible: for tasks with several metrics, take all of them at once and add them  
 138 as separate independent results to the averaging. Then tasks with separate high metrics will  
 139 have less weight on the total;
- 140 • among the best benchmark results, there were no negative values (MCC metric), but in  
 141 theory they could be, and that would prevent the calculation of the harmonic mean.
- 142 • the GLUE diagnostic dataset does not have a human evaluation score on the leaderboard  
 143 and therefore, is considered 0 (0.00001), though in the SuperGLUE benchmark the same  
 144 dataset has obtained its evaluation.

### 145 3.2 GLUE

146 GLUE benchmark (Wang et al., 2019b) combines 11 tasks in various text classification and question  
 147 answering.

148 **Overall score:** average of all the task results. If task has 2 main metrics, these metrics are averaged,  
 149 then added to the overall average.

150 **Human evaluation:** collected on reported human performance numbers from original datasets, not  
 151 exceeding 200 examples (heavily criticised in (Nangia and Bowman, 2019)). The human baseline  
 152 performance on the diagnostic set was provided by the project authors with the help of six NLP  
 153 researchers annotating 50 randomly selected sentence pairs.

154 **Rearranging the scores:** the results of geometric and harmonic mean rearrangement are presented  
 155 in Tab. 1. GLUE benchmark seem to be the most reordered of all the ratings considered: the best  
 156 result by geometric and harmonic means belongs to humans, DeBERTa and McAlberty+DKM got a 1  
 157 point demotion, and the other models got severely rearranged their places.

N	Name	AM	HM	GM	CoLA	SST-2	MRPC Mean	STS-B Mean	QQP Mean	MNLI m	MNLI mm	QNLI
16	Human	87,10	86,16	86,91	66,40	97,80	83,55	92,65	69,95	92,00	92,80	91,20
1	DeBERTa Mac	90,80	84,78	86,25	71,50	97,50	93,00	92,75	83,50	91,90	91,60	99,20
2	Albert +DKM	90,70	84,70	86,13	74,80	97,00	93,55	92,70	82,65	91,30	91,10	97,80
6	T5	90,30	84,48	85,92	71,60	97,50	91,60	92,95	82,85	92,20	91,90	96,90
4	PING-AN	90,60	84,26	85,83	73,50	97,20	93,00	92,70	83,55	91,60	91,30	97,50
5	ERNIE	90,40	84,27	85,75	74,40	97,50	92,45	92,80	83,05	91,40	91,00	96,60

Table 1: Top results of ranking GLUE benchmark with geometric mean. N – original model rank on the leaderboard. MNLI m and MNLI mm correspond to MultiNLI Matched MultiNLI Mismatched, other task abbreviations correspond to their GLUE leaderboard designations accordingly.

### 158 3.3 SuperGLUE

159 SuperGLUE (Wang et al., 2019a) is the sophisticated version of the GLUE benchmark, combining 10  
 160 tasks with a higher demand for higher intellectual abilities. Task data must be available under various  
 161 licenses that allow use and redistribution for research purposes.

162 **Overall score:** average of all the task results. If task has 2 main metrics, these metrics are averaged,  
 163 then added to the overall average.

164 **Human evaluation:** ready-made estimates for WiC, MultiRC, RTE, and ReCoRD datasets, the  
 165 other tasks being evaluated by the project creators with the help of crowdworker annotators through  
 166 Amazon’s Mechanical Turk.

167 **Rearranging the scores:** Re-weighting the results using the geometric mean and harmonic mean  
 168 again makes significant changes to the original ranking: the top-3 result (human) is ranked top-1, the  
 169 DeBerta and T5 models are shifted down 1 position, PAI ALbert and Nezha Plus models swap their  
 170 places, see Tab. 2.

N	Model	AM	HM	GM	BoolQ	CB Mean	COPA	MRC	RCD	RTE	WiC	WSC	AX-b	AX-g mean
3	Human	89,80	87,96	88,73	89,00	97,35	100,00	66,85	91,50	93,60	80,00	100,00	76,6	99,5
1	DeBERTa	90,30	86,89	87,60	90,40	96,65	98,40	75,95	94,30	93,20	77,50	95,90	66,7	93,55
2	T5+ Meena	90,20	86,42	87,10	91,30	96,70	97,40	75,65	93,85	92,70	77,90	95,90	66,5	89,35
4	T5	89,30	85,89	86,57	91,20	95,35	94,80	75,70	93,75	92,50	76,90	93,80	65,6	92,3
6	PAI ALbert	86,10	85,24	85,78	88,10	94,40	91,80	69,65	88,65	88,80	74,10	93,20	75,6	98,75
5	Nezha plus	86,70	81,30	82,29	87,80	95,20	93,60	69,85	89,85	89,10	74,60	93,20	58,00	80,75

Table 2: Top results of ranking SuperGLUE benchmark with geometric mean. N – original model rank on the leaderboard MRC stands for MultiRC averaged metric, RCD - ReCoRD averaged metric.

### 171 3.4 XTREME

172 The XTREME benchmark (Hu et al., 2020) covers 40 typologically diverse languages from 12  
 173 language families and includes 9 tasks that require analysis of different levels of syntax or semantics.

174 **Overall score:** 2-step averaging: 1) calculating average for each task on all languages 2) calculating  
 175 average on all tasks.

176 **Human evaluation:** 2-step averaging:

177 1. step 1: ready-made estimates from the original datasets taken and extrapolated to all  
 178 unestimated languages; besides, for some datasets there were no original estimates provided  
 179 (POS) and an empirical estimate of 97% was taken based on (Manning, 2011); no estimates  
 180 for NER and sentence retrieval tasks;

181 2. step 2: all the task results averaged together.

182 **Rearranging the scores:** the results of applying the geometric mean and harmonic mean did not  
 183 change the current ranking of the models - the quality spread between them is high enough for the  
 184 metrics averaging them to retain the current order.

N	Model	AM	HM	GM	Sentence-pair Classification	Structured Prediction	Question Answering	Sentence Retrieval
1	Human	93,30	93,13	93,21	95,10	97,00	87,80	0.00001
2	VECO	81,10	81,27	81,70	88,60	75,40	72,40	92,10
3	ERNIE-M	80,90	81,11	81,52	87,90	75,60	72,30	91,90
4	T-ULRv2	80,70	80,91	81,25	88,80	75,40	72,90	89,30
5	Anonymous3	79,90	80,12	80,50	88,20	74,60	71,70	89,00
6	Polyglot	77,80	78,02	78,56	87,80	72,90	67,40	88,30

Table 3: Top results of ranking XTREME benchmark with geometric mean. N – original model rank on the leaderboard; the averaged task scores are shown by the column markings.

### 185 3.5 XGLUE

186 The XGLUE benchmark (Liang et al., 2020) consists of 11 problems in 19 languages and evaluates  
 187 the performance of multilingual pre-trained systems in terms of their ability to cross-language  
 188 understanding and natural language generation.

189 **Overall score:** 2-step averaging: 1) calculating average for each task on all languages 2) calculating  
 190 average on all tasks.

191 **Human evaluation:** not provided.

192 **Rearranging the scores:** Since the human level is not measured in the benchmark, we can only  
 193 compare the 2 present models with each other. Tab. 4 shows the results - the difference in the quality  
 194 of the models is large enough to preserve their ranking on all averaging metrics.

N	Model	AM	HM	GM	NER	POS	NC	MLQA	XNLI	PAWS-X	QADSM	WPR	QAM
1	FILTER	80,10	79,61	79,86	82,60	81,60	83,50	76,20	83,90	93,80	71,40	74,70	73,40
2	Unicoder Baseline	76,10	75,45	75,80	79,70	79,60	83,50	66,00	75,30	90,10	68,40	73,90	68,90

  

N	Model	AM	HM	GM	QG	NTG
1	Unicoder Baseline	10,70	10,65	9,10	10,60	10,70
2	MP-Tune	8,70	8,70	7,10	8,10	9,40

Table 4: Top results of ranking XGLUE benchmark with geometric mean. N – original model rank on the leaderboard; the first 3 rows correspond to NLU tasks, the last 3 rows - to the NLG tasks.

## 195 4 Results and Discussion

196 The results show that the ranking of results within a single leaderboard can fluctuate significantly. So,  
 197 in GLUE, the first place in terms of the harmonic and geometric mean belongs to the result occupying  
 198 the 16th line in the arithmetic mean. In SuperGLUE, the permutation is not so striking - the third  
 199 result is on the 1st place. On the XTREME and XGLUE benchmarks system ranking is preserved.

200 Since all three averaging metrics considered are subject to different biases, we present the statistical  
 201 measurements of the top-3 SuperGLUE results in Tab. 5. Human results have the highest total points  
 202 for all tasks (as in the DecaNLP methodology), while the standard deviation and variance are greater  
 203 than top-2 and top-3 models.

Model	AM	GM	HM	Sum	Var	Std
Human	89,8	<b>88,73</b>	<b>87,96</b>	<b>894,40</b>	130,31	11,42
DeBerta	<b>90,3</b>	87,60	86,89	882,55	117,46	10,84
T5 + Meena	90,2	87,10	86,42	877,25	<b>112,39</b>	<b>10,60</b>

Table 5: Measuring the statistics of the top-3 SuperGLUE results. Sum is a sum of all the task scores, Var and Std are variance and standard deviation on the task scores respectively. Notable results are highlighted in bold.

204 To further explore the rating results of the GLUE and SuperGLUE benchmarks, we have conducted a  
 205 series of experiments with normalizing the model performance with human scores, fully transferring  
 206 the standard methodology for computing performance, in which the geometric mean is the approved  
 207 metric. The results are presented in Appendix 1. As can be concluded from the table with normalized  
 208 values, in the case of SuperGLUE, the results obtained by the new ranking method are confirmed. In  
 209 the case of GLUE, the geometric mean shows that the normalized ratio of models to the human level  
 210 asserts a high level of artificial solutions over the human level.

211 The following topics remain debatable and need special attention of the community:

- 212 1. Different metrics for obtaining the average value (arithmetic, geometric, harmonic) have  
 213 different restrictions on the accepted values (for example, not every one can take negative or  
 214 zero values). At the same time, metrics that take zero and negative values are actively used  
 215 in measuring various skills - MCC metric on SuperGLUE diagnostics can be negative, other  
 216 metrics can be equal to or greater than zero, and they are averaged altogether. Potentially,  
 217 the issue of a fair metric will raise the problem of revising the use of some individual metrics  
 218 for evaluating tasks. The differences in the metric scale for different tasks can pose problems  
 219 for the computation of the total score and some metrics can be scaled or normalized. We  
 220 may consider rescaling MCC score so that it is between 0 and 1.
- 221 2. Correct averaging of the overall score for multilingual benchmarks creates additional prob-  
 222 lems while performing the averaging operation in 2 stages: for all languages and all tasks.

223 As a result, the score, consisting of one number, becomes less and less informative and more  
224 prone to outliers.

225 3. Nevertheless, the competitive side of benchmarks is the driving force behind the progress  
226 in the field of machine learning, and besides all problems, it is still not worth giving up  
227 leaderboards with a single metric.

228 4. In addition to the problem of the main averaging metric, we left outside of the scope the  
229 problem that was also discovered within the framework of this study: human benchmark  
230 scores on various tasks were obtained in a very different way, and always on a smaller  
231 sample than the full test set. For a fair comparison of humans and machines, the test results  
232 should be normalized by the same number of test items, and it is worth revising the human  
233 evaluation and re-performing it on all test items using more annotators.

## 234 5 Conclusion

235 In this paper we present an alternative methods to arrange the popular NLP benchmark results,  
236 elaborating on several task evaluation. We analyze popular benchmark averaging methods and  
237 provide new insight into model comparison. Namely, we obtain the following results:

- 238 • for popular benchmarks GLUE and SuperGLUE we can conclude that their overall score is  
239 subject to bias due to outliers; the alternative arrangemend methods end with significantly  
240 different ordering of the results;
- 241 • rebuilding leaderboards using other metrics (geometric or harmonic mean) allows one to  
242 conclude that human result is the first in the rankings;
- 243 • in XGLUE leaderboard human result is obtained by extrapolation from one language  
244 to others, while in practice the level of problem-solving by native speakers of different  
245 languages varies;
- 246 • the last finding could be extended to other multilingual benchmarks also.

247 An unbiased view of the overall score in the benchmarks is a necessity for a community to target  
248 language model development on the complex improvement of their quality, not the partial results in  
249 narrow tasks. The expansion towards multilingual and multimodal models makes this issue more and  
250 more urgent and we hope our help to foster research in this direction.

## 251 References

- 252 François Chollet. 2019. On the measure of intelligence. *arXiv preprint arXiv:1911.01547*.
- 253 Christopher Clark, Kenton Lee, Ming-Wei Chang, Tom Kwiatkowski, Michael Collins, and Kristina  
254 Toutanova. 2019. Boolq: Exploring the surprising difficulty of natural yes/no questions. *arXiv  
255 preprint arXiv:1905.10044*.
- 256 Ido Dagan, Oren Glickman, and Bernardo Magnini. 2005. The pascal recognising textual entailment  
257 challenge. In *Machine Learning Challenges Workshop*, pages 177–190. Springer.
- 258 Marie-Catherine De Marneffe, Mandy Simons, and Judith Tonhauser. 2019. The commitmentbank:  
259 Investigating projection in naturally occurring discourse. In *proceedings of Sinn und Bedeutung*,  
260 volume 23, pages 107–124.
- 261 Ingolf Dittmann and Ernst Maug. 2008. Biases and error measures: How to compare valuation  
262 methods.
- 263 Aparna Elangovan, Jiayuan He, and Karin Verspoor. 2021. Memorization vs. generalization: quanti-  
264 fying data leakage in nlp performance evaluation. *arXiv preprint arXiv:2102.01818*.
- 265 Philip J. Fleming and John J. Wallace. 1986. How not to lie with statistics: The correct way to  
266 summarize benchmark results. *Commun. ACM*, 29(3):218–221.
- 267 Pengcheng He, Xiaodong Liu, Jianfeng Gao, and Weizhu Chen. 2020. Deberta: Decoding-enhanced  
268 bert with disentangled attention. *arXiv preprint arXiv:2006.03654*.

- 269 Junjie Hu, Sebastian Ruder, Aditya Siddhant, Graham Neubig, Orhan Firat, and Melvin Johnson. 2020.  
270 Xtreme: A massively multilingual multi-task benchmark for evaluating cross-lingual generalisation.  
271 In *International Conference on Machine Learning*, pages 4411–4421. PMLR.
- 272 Hang Le, Loïc Vial, Jibril Frej, Vincent Segonne, Maximin Coavoux, Benjamin Lecouteux, Alexandre  
273 Allauzen, Benoît Crabbé, Laurent Besacier, and Didier Schwab. 2019. Flaubert: Unsupervised  
274 language model pre-training for french. *arXiv preprint arXiv:1912.05372*.
- 275 Yaobo Liang, Nan Duan, Yeyun Gong, Ning Wu, Fenfei Guo, Weizhen Qi, Ming Gong, Linjun  
276 Shou, Daxin Jiang, Guihong Cao, et al. 2020. Xglue: A new benchmark dataset for cross-lingual  
277 pre-training, understanding and generation. *arXiv preprint arXiv:2004.01401*.
- 278 Ya lun Chou. 1969. *Statistical analysis: With business and economic applications*.
- 279 Christopher D Manning. 2011. Part-of-speech tagging from 97% to 100%: is it time for some  
280 linguistics? In *International conference on intelligent text processing and computational linguistics*,  
281 pages 171–189. Springer.
- 282 Bryan McCann, Nitish Shirish Keskar, Caiming Xiong, and Richard Socher. 2018. The natural  
283 language decathlon: Multitask learning as question answering. *arXiv preprint arXiv:1806.08730*.
- 284 Nikita Nangia and Samuel R Bowman. 2019. Human vs. muppet: A conservative estimate of human  
285 performance on the glue benchmark. *arXiv preprint arXiv:1905.10425*.
- 286 Sungjoon Park, Jihyung Moon, Sungdong Kim, Won Ik Cho, Jiyeon Han, Jangwon Park, Chisung  
287 Song, Junseong Kim, Yongsook Song, Taehwan Oh, et al. 2021. Klue: Korean language under-  
288 standing evaluation. *arXiv preprint arXiv:2105.09680*.
- 289 Fabio Petroni, Aleksandra Piktus, Angela Fan, Patrick Lewis, Majid Yazdani, Nicola De Cao, James  
290 Thorne, Yacine Jernite, Vladimir Karpukhin, Jean Maillard, et al. 2020. Kilt: a benchmark for  
291 knowledge intensive language tasks. *arXiv preprint arXiv:2009.02252*.
- 292 Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi  
293 Zhou, Wei Li, and Peter J. Liu. 2020. Exploring the limits of transfer learning with a unified  
294 text-to-text transformer.
- 295 Pranav Rajpurkar, Jian Zhang, Konstantin Lopyrev, and Percy Liang. 2016. Squad: 100,000+  
296 questions for machine comprehension of text.
- 297 Piotr Rybak, Robert Mroczkowski, Janusz Tracz, and Ireneusz Gawlik. 2020. Klej: comprehensive  
298 benchmark for polish language understanding. *arXiv preprint arXiv:2005.00630*.
- 299 Yutaka Sasaki et al. 2007. The truth of the f-measure. 2007.
- 300 Tatiana Shavrina, Alena Fenogenova, Anton Emelyanov, Denis Shevelev, Ekaterina Artemova,  
301 Valentin Malykh, Vladislav Mikhailov, Maria Tikhonova, Andrey Chertok, and Andrey Evlampiev.  
302 2020. Russiansuperglue: A russian language understanding evaluation benchmark. *arXiv preprint*  
303 *arXiv:2010.15925*.
- 304 James E. Smith. 1988. Characterizing computer performance with a single number. *Communications*  
305 *of the ACM*, 31(10):1202–1206.
- 306 Alan M Turing. 2009. Computing machinery and intelligence. In *Parsing the turing test*, pages  
307 23–65. Springer.
- 308 Alex Wang, Yada Pruksachatkun, Nikita Nangia, Amanpreet Singh, Julian Michael, Felix Hill, Omer  
309 Levy, and Samuel R Bowman. 2019a. Superglue: A stickier benchmark for general-purpose  
310 language understanding systems. *arXiv preprint arXiv:1905.00537*.
- 311 Alex Wang, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel R Bowman. 2019b.  
312 Glue: A multi-task benchmark and analysis platform for natural language understanding. In *7th*  
313 *International Conference on Learning Representations, ICLR 2019*.
- 314 Da-Feng Xia, Sen-Lin Xu, and Feng Qi. 1999. A proof of the arithmetic mean-geometric mean-  
315 harmonic mean inequalities. *RGMA research report collection*, 2(1).



316 Liang Xu, Hai Hu, Xuanwei Zhang, Lu Li, Chenjie Cao, Yudong Li, Yechen Xu, Kai Sun, Dian  
317 Yu, Cong Yu, et al. 2020. Clue: A chinese language understanding evaluation benchmark. *arXiv*  
318 *preprint arXiv:2004.05986*.

- 319 1. For all authors...
- 320 (a) Do the main claims made in the abstract and introduction accurately reflect the paper's  
321 contributions and scope? [Yes]
- 322 (b) Did you describe the limitations of your work? [Yes]
- 323 (c) Did you discuss any potential negative societal impacts of your work? [Yes]
- 324 (d) Have you read the ethics review guidelines and ensured that your paper conforms to  
325 them? [Yes]
- 326 2. If you are including theoretical results...
- 327 (a) Did you state the full set of assumptions of all theoretical results? [N/A]
- 328 (b) Did you include complete proofs of all theoretical results? [N/A]
- 329 3. If you ran experiments...
- 330 (a) Did you include the code, data, and instructions needed to reproduce the main experi-  
331 mental results (either in the supplemental material or as a URL)? [Yes]
- 332 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they  
333 were chosen)? [Yes]
- 334 (c) Did you report error bars (e.g., with respect to the random seed after running experi-  
335 ments multiple times)? [N/A]
- 336 (d) Did you include the total amount of compute and the type of resources used (e.g., type  
337 of GPUs, internal cluster, or cloud provider)? [N/A]
- 338 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 339 (a) If your work uses existing assets, did you cite the creators? [Yes]
- 340 (b) Did you mention the license of the assets? [Yes]
- 341 (c) Did you include any new assets either in the supplemental material or as a URL? [Yes]
- 342 (d) Did you discuss whether and how consent was obtained from people whose data you're  
343 using/curating? [N/A]
- 344 (e) Did you discuss whether the data you are using/curating contains personally identifiable  
345 information or offensive content? [N/A]
- 346 5. If you used crowdsourcing or conducted research with human subjects...
- 347 (a) Did you include the full text of instructions given to participants and screenshots, if  
348 applicable? [N/A]
- 349 (b) Did you describe any potential participant risks, with links to Institutional Review  
350 Board (IRB) approvals, if applicable? [N/A]
- 351 (c) Did you include the estimated hourly wage paid to participants and the total amount  
352 spent on participant compensation? [N/A]

353 —  
354  
355