PrivacyGLUE: A Benchmark Dataset for General Language Understanding in Privacy Policies

Anonymous ACL submission

Abstract

001 Benchmarks for general language understanding have been rapidly developing in recent years of NLP research, with well-known ex-004 amples such as GLUE and SuperGLUE. While benchmarks have been proposed in the legal language domain, virtually no such benchmarks exist for privacy policies despite their increasing importance in modern digital life. This could be explained by privacy policies falling under the legal language domain, but we find evidence to the contrary that motivates 011 a separate benchmark for privacy policies. Con-012 sequently, we propose PrivacyGLUE as the first comprehensive benchmark of relevant and highquality privacy tasks for measuring general language understanding in the privacy language domain. Furthermore, we release performances 017 from the BERT, RoBERTa, Legal-BERT, Legal-**RoBERTa and PrivBERT transformer language** models and perform model-pair agreement analysis to detect PrivacyGLUE task examples where models benefited from domain specialization. Our findings show PrivBERT outperforms other models by an average of 2 - 3%over all PrivacyGLUE tasks, shedding light on the importance of in-domain pretraining for privacy policies. We believe PrivacyGLUE can accelerate NLP research and improve general language understanding for humans and AI algorithms in the privacy language domain.

1 Introduction

034

Data privacy is evolving into a critical aspect of modern life with the United Nations (UN) describing it as a *human right in the digital age* (Gstrein and Beaulieu, 2022). Despite its importance, several studies have demonstrated high barriers to the understanding of privacy policies (Obar and Oeldorf-Hirsch, 2020) and estimate that an average person would require ~200 hours annually to read through all privacy policies encountered in their daily life (McDonald and Cranor, 2008). To



Figure 1: UMAP visualization of BERT embeddings from Wikipedia, European Legislation (EURLEX) and company privacy policy documents with a total of 2.5M tokens per corpus

address this, studies such as Wilson et al. (2016) recommend training Artificial Intelligence (AI) algorithms on appropriate benchmark datasets to assist humans in understanding privacy policies.

In recent years, benchmarks have been gaining popularity in Machine Learning and Natural Language Processing (NLP) communities because of their ability to holistically evaluate model performance over a variety of representative tasks. GLUE (Wang et al., 2018) and SuperGLUE (Wang et al., 2019) are examples of popular NLP benchmarks which measure the natural language understanding capabilities of SOTA models. NLP benchmarks are also developing rapidly in language domains, with LexGLUE (Chalkidis et al., 2022) being an example of a recent benchmark hosting several difficult tasks in the legal language domain. Interestingly, we do not find similar NLP benchmarks in the privacy language domain for privacy policies. While

Task	Source	Task Type	Train/Dev/Test Instances	# Classes
OPP-115	Wilson et al. (2016)	Multi-label sequence classification	2,185/550/697	12
PI-Extract	Bui et al. (2021)	Multi-task token classification	2,579/456/1,029	3/3/3/3†
Policy-Detection	Amos et al. (2021)	Binary sequence classification	773/137/391	2
PolicyIE-A	Ahmad et al. (2021)	Multi-class sequence classification	4,109/100/1,041	5
PolicyIE-B	Ahmad et al. (2021)	Multi-task token classification	4,109/100/1,041	29/9 [†]
PolicyQA	Ahmad et al. (2020)	Reading comprehension	17,056/3,809/4,152	_
PrivacyQA	Ravichander et al. (2019)	Binary sequence classification	157,420/27,780/62,150	2

Table 1: Summary statistics of PrivacyGLUE benchmark tasks; † PI-Extract and PolicyIE-B consist of four and two subtasks respectively and the number of BIO token classes per subtask are separated by a forward slash character

this could be explained by privacy policies falling under the legal language domain due to their formal and jargon-heavy nature, we claim that privacy policies fall under a distinct language domain and cannot be subsumed under any other specialized NLP benchmark such as LexGLUE.

To investigate this claim, we gather documents from Wikipedia (Wikimedia Foundation, 2022), European Legislation (EURLEX; Chalkidis et al. 2019) and company privacy policies (Mazzola et al., 2022), with each corpus truncated to 2.5M tokens. Next, we feed these documents into BERT and gather contextualized embeddings, which are then projected to 2-dimensional space using UMAP (McInnes et al., 2018). In Figure 1, we observe that the three domain corpora cluster independently, providing evidence that privacy policies lie in a distinct language domain from both legal and wikipedia documents. With this motivation, we propose PrivacyGLUE as the first comprehensive benchmark for measuring general language understanding in the privacy language domain. Our main contributions are threefold:

- 1. Composition of seven high-quality and relevant PrivacyGLUE tasks, specifically OPP-115, PI-Extract, Policy-Detection, PolicyIE-A, PolicyIE-B, PolicyQA and PrivacyQA.
- 2. Benchmark performances of five transformer language models on all aforementioned tasks, specifically BERT, RoBERTa, Legal-BERT, Legal-RoBERTa and PrivBERT.
- 3. Model agreement analysis to detect PrivacyGLUE task examples where models benefited from domain specialization.

We release PrivacyGLUE as a fully configurable benchmark suite for straight-forward reproducibility and production of new results in our public GitHub repository¹. Our findings show that PrivBERT, the only model pretrained on privacy policies, outperforms other models by an average of 2 – 3% over all PrivacyGLUE tasks, shedding light on the importance of in-domain pretraining for privacy policies. Our model-pair agreement analysis explores specific examples where PrivBERT's privacy-domain pretraining provided both competitive advantage and disadvantage. By benchmarking holistic model performances, we believe PrivacyGLUE can accelerate NLP research into the privacy language domain and ultimately improve general language understanding of privacy policies for both humans and AI algorithms.

097

099

100

101

102

103

104

105

106

107

108

109

110

111

112

113

114

115

116

117

118

119

120

121

122

123

124

125

126

127

128

129

130

131

2 **Related work**

NLP benchmarks have been gaining popularity in recent years because of their ability to holistically evaluate model performance over a variety of representative tasks. GLUE (Wang et al., 2018) and SuperGLUE (Wang et al., 2019) are examples of benchmarks that evaluate SOTA models on a range of natural language understanding tasks. The GEM benchmark (Gehrmann et al., 2021) looks beyond text classification and measures performance in Natural Language Generation tasks such as summarization and data-to-text conversion. The XTREME (Hu et al., 2020) and XTREME-R (Ruder et al., 2021) benchmarks specialize in measuring crosslingual transfer learning on 40-50 typologically diverse languages and corresponding tasks. Popular NLP benchmarks often host public leaderboards with SOTA scores on supported tasks, thereby encouraging the community to apply new approaches for surpassing top scores.

061

062 063

- 089

094

¹Repository will be made public post-acceptance. Anonymous repository: https://anonymous.4open.science/ r/f4293357886f671347fa69fae3650543

While the aforementioned benchmarks focus on problem types such as natural language understanding and generation, other benchmarks focus on language domains. The LexGLUE benchmark (Chalkidis et al., 2022) is an example of a benchmark that evaluates models on tasks from the legal language domain. LexGLUE consists of seven English-language tasks that are representative of the legal language domain and chosen based on size and legal specialization. Chalkidis et al. (2022) benchmarked several models such as BERT (Devlin et al., 2019) and Legal-BERT (Chalkidis et al., 2020), where Legal-BERT has a similar architecture to BERT but was pretrained on diverse legal corpora. A key finding of LexGLUE was that Legal-BERT outperformed other models which were not pretrained on legal corpora. In other words, they found that an in-domain pretrained model outperformed models that were pretrained out-of-domain.

132

133

134

135

136

137

138

139

140

141

142

143

144

145

146

147

148

149

150

151

152

154

155

156

158

160

161

162

163

164

165

166

167

168

169

170

171

172

173

174

In the privacy language domain, we tend to find isolated datasets from specialized studies. Zimmeck et al. (2019), Wilson et al. (2016), Bui et al. (2021) and Ahmad et al. (2021) are examples of studies that introduce annotated corpora for privacy-practice sequence and token classification tasks, while Ravichander et al. (2019) and Ahmad et al. (2020) release annotated corpora for privacy-practice question answering. Amos et al. (2021) is another recent study that released an annotated corpus of privacy policies. As of writing, no comprehensive NLP benchmark exists for general language understanding in privacy policies, making PrivacyGLUE the first consolidated NLP benchmark in the privacy language domain.

3 Datasets and Tasks

The PrivacyGLUE benchmark consists of seven natural language understanding tasks originating from six datasets in the privacy language domain. Summary statistics, detailed label information and representative examples are shown in Table 1, Table 5 (Appendix A) and Table 6 (Appendix B) respectively.

175**OPP-115** Wilson et al. (2016) was the first study176to release a large annotated corpus of privacy poli-177cies. A total of 115 privacy policies were selected178based on their corresponding company's popularity179on Google Trends. The selected privacy policies180were annotated with 12 data privacy practices on181a paragraph-segment level by experts in the pri-

vacy domain. As noted by Mousavi Nejad et al. (2020), one limitation of Wilson et al. (2016) was the lack of publicly released training and test data splits which are essential for machine learning and benchmarking. To address this, Mousavi Nejad et al. (2020) released their own training, validation and test data splits for researchers to easily reproduce OPP-115 results. PrivacyGLUE utilizes the "Majority" variant of data splits released by Mousavi Nejad et al. (2020) to compose the OPP-115 task. Given an input paragraph segment of a privacy policy, the goal of OPP-115 is to predict one or more data practice categories.

182

183

184

185

186

187

188

189

190

191

192

193

194

195

196

197

198

199

201

202

203

204

205

206

207

209

210

211

212

213

214

215

216

217

218

219

220

221

222

223

224

225

226

227

228

229

230

PI-Extract Bui et al. (2021) focuses on enhanced data practice extraction and presentation to help users better understand privacy policies. As part of their study, they released the PI-Extract dataset consisting of 4.1K sentences (97K tokens) and 2.6K expert-annotated data practices from 30 privacy policies in the OPP-115 dataset. Expert annotations were performed on a token-level for all sentences of selected privacy policies. PI-Extract is broken into four subtasks, where spans of tokens are independently tagged using the BIO scheme commonly used in Named Entity Recognition (NER). Subtasks I, II, III and IV require the classification of token spans for data-related entities that are collected, not collected, not shared and shared respectively. In the interest of diversifying tasks in PrivacyGLUE, we composed PI-Extract as a multi-task token classification problem where all four PI-Extract subtasks are to be jointly learned.

Policy-Detection Amos et al. (2021) developed a crawler for automated collection and curation of privacy policies. An important aspect of their system is the automated classification of documents into privacy policies and non-privacy-policy documents encountered during web crawling. To train such a privacy policy classifier, Amos et al. (2021) performed expert annotations of commonly encountered documents during web crawls and classified them into the aforementioned categories. The Policy-Detection dataset was released with a total of 1.3K annotated documents and is utilized in PrivacyGLUE as a binary sequence classification task.

PolicyIE Inspired by Wilson et al. (2016) and Bui et al. (2021), Ahmad et al. (2021) created PolicyIE, an English corpus composed by 5.3K sentence-level and 11.8K token-level data practice

Model	Source	# Params	Vocab. Size	Pretraining corpora †
BERT	Devlin et al. (2019)	110M	30K	Wikipedia, BC (16 GB)
RoBERTa	Liu et al. (2019)	125M	50K	Wikipedia, BC, CC-News, OWT (160 GB)
Legal-BERT	Chalkidis et al. (2020)	110M	30K	Legislation, Court Cases, Contracts (12 GB)
Legal-RoBERTa [‡]	Geng et al. (2021)	125M	50K	Patents, Court Cases (5 GB)
PrivBERT [‡]	Srinath et al. (2021)	125M	50K	Privacy policies (17 GB)

Table 2: Summary of models used in the PrivacyGLUE benchmark; all models used are base-sized variants of BERT/RoBERTa architectures; † BC = BookCorpus, CC-News = CommonCrawl-News, OWT = OpenWebText; ‡ models were initialized with the pretrained RoBERTa model

annotations over 31 privacy policies from websites and mobile applications. PolicyIE was designed to be used for machine learning in NLP, to ultimately make data privacy concepts easier for users to understand. We split the PolicyIE corpus into two tasks, namely *PolicyIE-A* and *PolicyIE-B*. Given an input sentence, PolicyIE-A entails multi-class data practice classification while PolicyIE-B entails multi-task token classification over distinct subtasks I and II, which require the classification of token spans for entities that participate in privacy practices and their conditions/purposes respectively. The motivation for composing PolicyIE-B as a multi-task problem is similar to that of PI-Extract.

239

241

242

243

245

PolicyQA Ahmad et al. (2020) argue in favour 246 of short-span answers to user questions for long 247 privacy policies. They release PolicyQA, a dataset 248 249 of 25k reading comprehension examples curated from the OPP-115 corpus from Wilson et al. (2016). 250 Furthermore, they provide 714 human-written gues-251 tions optimized for a wide range of privacy policies. The final question-answer annotations follow the SQuAD-1.0 format (Rajpurkar et al., 2016), which improves the ease of adaptation into NLP pipelines. We utilize PolicyQA as PrivacyGLUE's reading comprehension task. 257

PrivacyQA Similar to Ahmad et al. (2020), Ravichander et al. (2019) argue in favour of annotated question-answering data for training NLP 260 models to answer user questions about privacy poli-261 cies. They correspondingly released PrivacyQA, 262 a corpus composed by 1.75K questions and more 263 than 3.5K expert annotated answers. Unlike PolicyQA, PrivacyQA proposes a binary sequence clas-265 sification task where a question-answer pair is classified as either relevant or irrelevant. Correspond-267 ingly, we treat PrivacyQA as a binary sequence classification task in PrivacyGLUE.

4 Experimental setup

The PrivacyGLUE benchmark was tested using the BERT, RoBERTa, Legal-BERT, Legal-RoBERTa and PrivBERT models which are summarized in Table 2. We describe the models used and task-specific approaches, and provide details on our benchmark configuration in Appendix E.

270

271

272

273

274

275

278

279

280

281

285

287

289

290

291

292

293

294

295

296

297

298

299

300

301

302

303

304

305

306

4.1 Models

BERT Proposed by Devlin et al. (2019), BERT is perhaps the most well-known transformer language model. BERT utilizes the WordPiece tokenizer (Wu et al., 2016) and is case-insensitive. It is pretrained with the Masked Language Model (MLM) and Next Sentence Prediction (NSP) tasks on the Wikipedia and BookCorpus corpora.

RoBERTa Liu et al. (2019) proposed RoBERTa as an improvement to BERT. RoBERTa uses dynamic token masking and eliminates the NSP task during pretraining. Furthermore, it uses a case sensitive byte-level Byte-Pair Encoding (Sennrich et al., 2016) tokenizer and is pretrained on larger corpora. Liu et al. (2019) reported improved results on various benchmarks using RoBERTa over BERT.

Legal-BERT Chalkidis et al. (2020) proposed Legal-BERT by pretraining BERT from scratch on legal corpora consisting of legislation, court cases and contracts. The sub-word vocabulary of Legal-BERT is learned from scratch using the Sentence-Piece (Kudo and Richardson, 2018) tokenizer to better support legal terminology. Legal-BERT was the best overall performing model in the LexGLUE benchmark as reported in Chalkidis et al. (2022).

Legal-RoBERTa Inspired by Legal-BERT, Geng et al. (2021) proposed Legal-RoBERTa by further pretraining RoBERTa on legal corpora, specifically patents and court cases. Legal-RoBERTa

3(

.

310

311

312 313

- 314 315
- 31

317

3

319

322

328

329

334

335

336

340

341

342

345

350

4.2 Task-specific approaches

corpus.

tuning legal domain tasks.

is pretrained on less legal data than Legal-BERT

while producing similar results on downstream fine-

PrivBERT Due to the scarcity of large corpora

in the privacy domain, Srinath et al. (2021) pro-

posed PrivaSeer, a novel corpus of 1M English

language website privacy policies crawled from

the web. They subsequently proposed PrivBERT

by further pretraining RoBERTa on the PrivaSeer

Given the aforementioned models and tasks, we now describe our task-specific fine-tuning and evaluation approaches. Given an input sequence $s = \{w_1, w_2, \dots, w_N\}$ consisting of N sequential sub-word tokens, we feed s into a transformer encoder and obtain a contextual representation $\{h_0, h_1, \ldots, h_N\}$ where $h_i \in \mathbb{R}^D$ and D is the output dimensionality of the transformer encoder. Here, h_0 refers to the contextual embedding for the starting token which is [CLS] for BERT-derived models and <s> for RoBERTa-derived models. For PolicyQA and PrivacyQA, the input sequence s is composed by concatenating the question and context/answer pairs respectively. The concatenated sequences are separated by a separator token, which is [SEP] for BERT-derived models and </s> for RoBERTa-derived models.

4.2.1 Sequence classification

The h_0 embedding is fed into a class-wise sigmoid classifier (1) and softmax classifier (2) for multilabel and binary/multi-class tasks respectively. The classifier has weights $W \in \mathbb{R}^{D \times C}$ and bias $b \in \mathbb{R}^{C}$ and is used to predict the probability vector $y \in \mathbb{R}^{C}$, where *C* refers to the number of output classes. We fine-tune models end-to-end by minimizing the binary cross-entropy loss and cross-entropy loss for multi-label and binary/multi-class tasks respectively.

$$y = \text{sigmoid}(W^{\top}h_0 + b) \tag{1}$$

$$y = \operatorname{softmax}(W^{\top}h_0 + b) \tag{2}$$

We report the macro and micro-average F₁ scores for all sequence classification tasks since the former ignores class imbalance while the latter takes it into account.

4.2.2 Multi-task token classification

Each $h_i \in \{h_1, h_2, ..., h_N\}$ token embedding is fed into *J* independent softmax classifiers with weights $W_j \in \mathbb{R}^{D \times C_j}$ and bias $b_j \in \mathbb{R}^{C_j}$ to predict the token probability vector $y_{ij} \in \mathbb{R}^{C_j}$, where C_j refers to the number of output BIO classes per subtask $j \in$ $\{1, 2, ..., J\}$. We fine-tune models end-to-end by minimizing the cross-entropy loss across all tokens and subtasks.

$$y_{ij} = \operatorname{softmax}(W_j^{\mathsf{T}} h_i + b_j) \tag{3}$$

355

356

357

360

361

362

363

364

365

366

367

368

369

370

371

372

373

374

375

376

377

378

379

380

381

382

383

384

385

386

387

390

391

392

393

395

396

397

399

We report the macro and micro-average F_1 scores for all multi-task token classification tasks by averaging the respective metrics for each sub-task. Furthermore, we ignore cases where B or I prefixes are mismatched as long as the main token class is correct.

4.2.3 Reading comprehension

Each $h_i \in \{h_1, h_2, ..., h_N\}$ token embedding is fed into two independent linear layers with weights $W_j \in \mathbb{R}^D$ and bias $b_j \in \mathbb{R}$ where $j \in \{1, 2\}$. These linear outputs are then concatenated per layer and a softmax function is applied to form a probability vector y_j across all tokens for answer-start and answer-end token probabilities respectively. We fine-tune models end-to-end by minimizing the cross-entropy loss on the gold answer-start and answer-end indices.

$$y_j = \operatorname{softmax}\left(\left[W_j \cdot h_1 + b_j \dots W_j \cdot h_N + b_j\right]\right) \quad (4)$$

Similar to SQuAD (Rajpurkar et al., 2016), we report the sample F_1 and exact match accuracy for our reading comprehension task. It is worth noting that Rajpurkar et al. (2016) refer to their reported F_1 score as a macro-average, whereas we refer to it as the sample-average as we believe this is a more accurate term.

5 Results

After running the PrivacyGLUE benchmark with 10 random seeds, we collect results on the testsets of all tasks. Figure 2 shows the respective results in a graphical form while Table 7 in Appendix C shows the numerical results in a tabular form. In terms of absolute metrics, we observe that PrivBERT outperforms other models for all PrivacyGLUE tasks. We apply the Mann-Whitney Utest (Mann and Whitney, 1947) over random seed metric distributions and find that PrivBERT significantly outperforms other models on six out of seven PrivacyGLUE tasks with p <= 0.05, where



 \blacksquare Macro F₁ \blacksquare Micro F₁ \blacksquare Sample F₁ \blacksquare Exact Match

Figure 2: Test-set results of the PrivacyGLUE benchmark where points indicate mean performance and error bars indicate standard deviation over 10 random seeds; *** implies $p \le 0.001$, ** implies 0.001 , * implies <math>0.01 given an alternative hypothesis that PrivBERT has a greater performance metric than all other models in a task using the Mann-Whitney U-test

Policy-Detection was the task where the significance threshold was not met. We utilize the Mann-Whitney U-test because it does not require a normal distribution for test-set metrics, an assumption which has not been extensively validated for deep neural networks (Dror et al., 2019).

In Figure 2, we observe large differences between the two representative metrics for OPP-115, Policy-Detection, PolicyIE-A, PrivacyQA and PolicyQA. For the first four of the aforementioned tasks, this is because of data imbalance resulting in the micro-average F_1 being significantly higher since it can be skewed by the metric of the majority class. For PolicyQA, this occurs because the EM metric requires exact matches and is therefore much stricter than the sample F₁ metric. Furthermore, we observe an exceptionally large standard deviation on PI-Extract metrics compared to other tasks. This can be attributed to data imbalance between the four subtasks of PI-Extract, with the NOT_COLLECT and NOT_SHARE subtasks having less than 100 total examples each.

We apply the arithmetic, geometric and harmonic means to aggregated metric means and standard deviations as shown in Table 3. With this, we observe the following general ranking of models

Model	A-Mean		G-Mean		H-Mean	
	μ	σ	μ	σ	μ	σ
BERT	67.5	1.1	64.6	0.9	61.1	0.6
RoBERTa	69.0	1.2	66.4	0.7	63.2	0.3
Legal-BERT	67.9	1.1	64.9	0.8	61.2	0.4
Legal-RoBERTa	68.5	1.3	65.7	0.8	62.3	0.4
PrivBERT	70.8	1.2	68.3	0.8	65.2	0.5

Table 3: Macro-aggregation of means (μ) and standard deviations (σ) per model using the arithmetic mean (A-Mean), geometric mean (G-Mean) and harmonic mean (H-Mean)

from best to worst: PrivBERT, RoBERTa, Legal-RoBERTa, Legal-BERT and BERT. Interestingly, models derived from RoBERTa generally outperformed models derived from BERT. Using the arithmetic mean for simplicity, we observe that PrivBERT outperforms all other models by 2 - 3%.

426

427

428

429

430

431

432

433

434

435

436

437

6 Discussion

With the PrivacyGLUE benchmark results, we revisit our privacy vs. legal language domain claim from Section 1 and discuss our model-pair agreement analysis for detecting PrivacyGLUE task examples where models benefited from domain spe-

423

424

425

400



Figure 3: Model-pair agreement analysis of PrivBERT against other models over all PrivacyGLUE tasks; bars represent proportions of examples per model-pair and task which fell into categories P and O; all models on the x-axis are compared against PrivBERT

Category P	Category O
ID: 1978	ID: 33237
Question: Who can see my information?	Question: Could the wordscapes app contain mal-
Answer: We do not sell or rent your personal infor-	ware?
mation to third parties for their marketing purposes	Answer: We encrypt the transmission of all informa-
without your explicit consent.	tion using secure socket layer technology (SSL).
Label: Relevant	Label: Relevant

Table 4: Test-set examples from PrivacyQA that fall under categories P and O for PrivBERT vs. BERT

cialization.

438

439

440

441

442

443

444

445

446

447

448

449

450

451

452

6.1 Privacy vs. legal language domain

We initially provided evidence from Figure 1 suggesting that the privacy language domain is distinct from the legal language domain. We believe that our PrivacyGLUE results further support this initial claim. If the privacy language domain was subsumed under the legal language domain, we could have observed Legal-RoBERTa and Legal-BERT performing competitively with PrivBERT. Instead, we observed that the legal models underperformed compared to both PrivBERT and RoBERTa, further indicating that the privacy language domain is distinct and requires its own NLP benchmark.

6.2 Model-pair agreement analysis

453 PrivBERT, the top performing model, differentiates454 itself from other models by its in-domain pretrain-

ing on the PrivaSeer corpus (Srinath et al., 2021). Therefore, we can infer that PrivBERT incorporated *knowledge* of privacy policies through its pretraining and became specialized for fine-tuning tasks in the privacy language domain. We investigate this specialization using model-pair agreement analysis to detect examples where PrivBERT had a competitive advantage over other models. Consequently, we detect examples where PrivBERT was disadvantaged due to its in-domain pretraining.

We compare $10 \times 10 = 100$ random seed combinations for all test-set pairs between PrivBERT and other models. Each prediction-pair can be classified into one of four mutually exclusive categories (B, P, O and N) shown below. Categories B and N represent examples that are either not challenging or too challenging for both PrivBERT and the other model respectively. Categories P and O are more interesting for us since they indicate examples where

455

- 474 PrivBERT had a competitive advantage and disad475 vantage over the other model respectively. There476 fore, we focus on categories P and O in our analysis.
 477 We classify examples over all random seed combi478 nations and take the majority occurrence for each
 479 category within its distribution.
- 480 **Category B:** Both PrivBERT and the other model 481 were correct, i.e. (PrivBERT, Other Model)
- 482 Category P: PrivBERT was correct and the other
 483 model was wrong, i.e. (PrivBERT, ¬ Other Model)
- 484 Category O: Other model was correct and
 485 PrivBERT was wrong, i.e. (¬ PrivBERT, Other
 486 Model)

488

489

490

491

492

493

494

495

496

497

498

499

501

505

506

507

509

510

511

512

513

514

515

516

517

518

519

Category N: Neither PrivBERT nor the other model was correct, i.e. $(\neg \text{PrivBERT}, \neg \text{Other Model})$

Figure 3 shows a relative distribution of majority categories across model-pairs and PrivacyGLUE tasks. We observe that category P is always greater than category O, which correlates with PrivBERT outperforming all other models. We also observe that category P is often the greatest when compared against BERT, implying that PrivBERT has the most competitive advantage over BERT. Surprisingly, we also observe category O is often the greatest when compared against BERT, implying that BERT has the highest absolute advantage over PrivBERT. This is an insightful observation since we would have expected BERT to have the least competitive advantage given its lowest overall PrivacyGLUE performance.

To investigate PrivBERT's competitive advantage and disadvantage against BERT, we extract several examples from categories P and O in the PrivacyQA task for brevity. Two interesting examples are listed in Table 4 and additional examples can be found in Table 8 in Appendix D. From Table 4, we speculate that PrivBERT specializes in example 1978 because it contains several privacyspecific terms such as "third parties" and "explicit consent". On the other hand, we speculate that BERT specializes in example 33237 since it contains more generic information regarding encryption and SSL, which also happens to be a topic in BERT's Wikipedia pretraining corpus as seen in Figure 1 and Table 2.

> Looking at further examples in Table 8, we can also observe that all sampled category P exam

ples have the Relevant label while many sampled category O examples have the Irrelevant label. On further analysis of the PrivacyQA testset, we find that 71% of category P examples have the Relevant label and 61% of category O samples have the Irrelevant label. We can infer that PrivBERT specializes in the minority Relevant label while BERT specializes in the majority Irrelevant label as the former label could require more privacy knowledge than the latter. 522

523

524

525

526

527

528

529

530

531

532

533

534

535

536

537

538

539

540

541

542

543

544

545

546

547

548

549

550

551

552

553

554

555

556

557

558

559

560

561

562

563

564

565

566

567

568

569

7 Conclusions and further work

In this paper, we describe the importance of data privacy in modern digital life and observe the lack of a NLP benchmark in the privacy language domain despite its distinctness. To address this, we propose PrivacyGLUE as the first comprehensive benchmark for measuring general language understanding in the privacy language domain. We release benchmark performances from the BERT, RoBERTa, Legal-BERT, Legal-RoBERTa and PrivBERT transformer language models. Our findings show that PrivBERT outperforms other models by an average of 2 - 3% over all PrivacyGLUE tasks, shedding light on the importance of in-domain pretraining for privacy policies. We apply model-pair agreement analysis to detect PrivacyGLUE examples where PrivBERT's pretraining provides competitive advantage and disadvantage. By benchmarking holistic model performances, we believe PrivacyGLUE can accelerate NLP research into the privacy language domain and ultimately improve general language understanding of privacy policies for both humans and AI algorithms.

Looking forward, we envision several ways to further our study. Firstly, we intend to apply deeplearning explainability techniques such as Integrated Gradients (Sundararajan et al., 2017) on examples from Table 4, to explore PrivBERT's and BERT's token-level attention attributions for categories P and O. Additionally, we intend to benchmark large prompt-based transformer language models such as T5 (Raffel et al., 2020) and T0 (Sanh et al., 2022), as they incorporate large amounts of knowledge from the various sequenceto-sequence tasks that they were trained on. Finally, we plan to continue maintaining our PrivacyGLUE GitHub repository and host new model results from the community.

570 Limitations

To the best of our knowledge, our study has two 571 main limitations. While we provide performances 572 from transformer language models, our study does 573 not provide human expert performances on Priva-574 cyGLUE. This would have been a valuable contribution to judge how competitive language models 576 are against human expertise. However, this limitation can be challenging to address due to the 578 difficulty in finding experts and high costs for their services. Additionally, our study only focuses on English language privacy tasks and omits multilingual scenarios. Multilingual tasks would have been very interesting and relevant to explore, but also involve significant complexity since privacy experts for non-English languages may be harder 585 to find. 586

Ethics Statement

587

591

594

597

601

607

611

612

Original work attribution

All datasets used to compose PrivacyGLUE are publicly available and originate from previous studies. We cite these studies in our paper and include references for them in our GitHub repository. Furthermore, we clearly illustrate how these datasets were used to form the PrivacyGLUE benchmark.

Social impact

PrivacyGLUE could be used to produce fine-tuned transformer language models, which could then be utilized in downstream applications to help users understand privacy policies and/or answer questions regarding them. We believe this could have a positive social impact as it would empower users to better understand lengthy and complex privacy policies. That being said, application developers should perform appropriate risk analyses when using fine-tuned transformer language models. Important points to consider include the varying performance ranges on PrivacyGLUE tasks and known examples of implicit bias, such as gender and racial bias, that transformer language models incorporate through their large-scale pretraining (Bender et al., 2021).

Software licensing

We release source code for PrivacyGLUE under
version 3 of the GNU General Public License (GPL3.0). We chose GPL-3.0 as it is a strong copyleft
license that protects user freedoms such as the freedom to use, modify and distribute software.

References

Wasi Ahmad, Jianfeng Chi, Tu Le, Thomas Norton, Yuan Tian, and Kai-Wei Chang. 2021. Intent classification and slot filling for privacy policies. In Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers), pages 4402–4417, Online. Association for Computational Linguistics. 618

619

620

621

622

623

624

625

626

627

628

629

630

631

632

633

634

635

636

637

638

639

640

641

642

643

644

645

646

647

648

649

650

651

652

653

654

655

656

657

658

659

660

661

662

663

664

665

666

667

668

669

670

671

672

673

- Wasi Ahmad, Jianfeng Chi, Yuan Tian, and Kai-Wei Chang. 2020. PolicyQA: A reading comprehension dataset for privacy policies. In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 743–749, Online. Association for Computational Linguistics.
- Ryan Amos, Gunes Acar, Elena Lucherini, Mihir Kshirsagar, Arvind Narayanan, and Jonathan Mayer. 2021. Privacy Policies over Time: Curation and Analysis of a Million-Document Dataset. In *Proceedings of The Web Conference 2021*, WWW '21, page 22. Association for Computing Machinery.
- Emily M. Bender, Timnit Gebru, Angelina McMillan-Major, and Shmargaret Shmitchell. 2021. On the dangers of stochastic parrots: Can language models be too big? In *Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency*, FAccT '21, page 610–623, New York, NY, USA. Association for Computing Machinery.
- Lukas Biewald. 2020. Experiment tracking with weights and biases. Software available from wandb.com.
- Duc Bui, Kang G. Shin, Jong-Min Choi, and Junbum Shin. 2021. Automated extraction and presentation of data practices in privacy policies. *Proceedings on Privacy Enhancing Technologies*, 2021(2):88–110.
- Ilias Chalkidis, Manos Fergadiotis, Prodromos Malakasiotis, Nikolaos Aletras, and Ion Androutsopoulos. 2020. LEGAL-BERT: The muppets straight out of law school. In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 2898– 2904, Online. Association for Computational Linguistics.
- Ilias Chalkidis, Manos Fergadiotis, Prodromos Malakasiotis, and Ion Androutsopoulos. 2019. Large-scale multi-label text classification on EU legislation. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 6314– 6322, Florence, Italy. Association for Computational Linguistics.
- Ilias Chalkidis, Abhik Jana, Dirk Hartung, Michael Bommarito, Ion Androutsopoulos, Daniel Katz, and Nikolaos Aletras. 2022. LexGLUE: A benchmark dataset for legal language understanding in English. In Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 4310–4330, Dublin, Ireland. Association for Computational Linguistics.

788

690

675

676

677

679

697

701

704 705 706

707

708

710 711

712

713

715 716

717

719

720

721

724

725

726 727 728

729

731

732 733 734

A simple and language independent subword tokenizer and detokenizer for neural text processing. In

4411-4421. PMLR.

abs/2109.06862.

phy & Technology, 35(1):1-38.

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and

Kristina Toutanova. 2019. BERT: Pre-training of

deep bidirectional transformers for language under-

standing. In Proceedings of the 2019 Conference of

the North American Chapter of the Association for

Computational Linguistics: Human Language Tech-

nologies, Volume 1 (Long and Short Papers), pages

4171-4186, Minneapolis, Minnesota. Association for

Rotem Dror, Segev Shlomov, and Roi Reichart. 2019. Deep dominance - how to properly compare deep

neural models. In Proceedings of the 57th Annual

Meeting of the Association for Computational Lin-

guistics, pages 2773-2785, Florence, Italy. Associa-

Aggarwal, Pawan Sasanka Ammanamanchi,

Anuoluwapo Aremu, Antoine Bosselut, Khyathi Raghavi Chandu, Miruna-Adriana Clinciu,

Dipanjan Das, Kaustubh Dhole, Wanyu Du, Esin

Durmus, Ondřej Dušek, Chris Chinenye Emezue,

Varun Gangal, Cristina Garbacea, Tatsunori

Hashimoto, Yufang Hou, Yacine Jernite, Harsh Jham-

tani, Yangfeng Ji, Shailza Jolly, Mihir Kale, Dhruv

Kumar, Faisal Ladhak, Aman Madaan, Mounica

Maddela, Khyati Mahajan, Saad Mahamood, Bod-

hisattwa 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 Shimo-

rina, Marco Antonio Sobrevilla Cabezudo, Hendrik

Strobelt, Nishant Subramani, Wei Xu, Diyi Yang,

GEM benchmark: Natural 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.

transformer models may not always help. CoRR,

Saibo Geng, Rémi Lebret, and Karl Aberer. 2021. Legal

Oskar J Gstrein and Anne Beaulieu. 2022. How to pro-

Junjie Hu, Sebastian Ruder, Aditya Siddhant, Gra-

ham Neubig, Orhan Firat, and Melvin Johnson.

2020. XTREME: A massively multilingual multi-

task benchmark for evaluating cross-lingual gener-

alisation. In Proceedings of the 37th International Conference on Machine Learning, volume 119 of

Proceedings of Machine Learning Research, pages

Taku Kudo and John Richardson. 2018. SentencePiece:

tect privacy in a datafied society? a presentation of

multiple legal and conceptual approaches. Philoso-

Akhila Yerukola, and Jiawei Zhou. 2021.

Sebastian Gehrmann, Tosin Adewumi, Karmanya

Computational Linguistics.

tion for Computational Linguistics.

Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing: System Demonstrations, pages 66-71, Brussels, Belgium. Association for Computational Linguistics.

- Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Dangi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019. Roberta: A robustly optimized bert pretraining approach. arXiv preprint arXiv:1907.11692.
- Ilva Loshchilov and Frank Hutter. 2017. Fixing weight decay regularization in adam. CoRR. abs/1711.05101.

Henry B Mann and Donald R Whitney. 1947. On a test

Luca Mazzola, Andreas Waldis, Atreya Shankar, Dia-

mantis Argyris, Alexander Denzler, and Michiel

Van Roey. 2022. Privacy and customer's education:

Nlp for information resources suggestions and expert

finder systems. In HCI for Cybersecurity, Privacy

and Trust, pages 62-77, Cham. Springer International

Aleecia M McDonald and Lorrie Faith Cranor. 2008.

The cost of reading privacy policies. Isjlp, 4:543.

L. McInnes, J. Healy, and J. Melville. 2018. UMAP:

Najmeh Mousavi Nejad, Pablo Jabat, Rostislav

Nedelchev, Simon Scerri, and Damien Graux. 2020.

Establishing a strong baseline for privacy policy clas-

sification. In IFIP International Conference on ICT

Systems Security and Privacy Protection, pages 370-

Jonathan A Obar and Anne Oeldorf-Hirsch. 2020. The

biggest lie on the internet: Ignoring the privacy poli-

cies and terms of service policies of social network-

ing services. Information, Communication & Society,

Adam Paszke, Sam Gross, Francisco Massa, Adam

Lerer, James Bradbury, Gregory Chanan, Trevor

Killeen, Zeming Lin, Natalia Gimelshein, Luca

Antiga, Alban Desmaison, Andreas Kopf, Edward

Yang, Zachary DeVito, Martin Raison, Alykhan Te-

jani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang,

Junjie Bai, and Soumith Chintala. 2019. Pytorch: An

imperative style, high-performance deep learning li-

brary. In Advances in Neural Information Processing

Colin Raffel, Noam Shazeer, Adam Roberts, Kather-

ine Lee, Sharan Narang, Michael Matena, Yanqi

Zhou, Wei Li, and Peter J. Liu. 2020. Exploring the

limits of transfer learning with a unified text-to-text

transformer. Journal of Machine Learning Research,

Systems, volume 32. Curran Associates, Inc.

Dimension Reduction. ArXiv e-prints.

Uniform Manifold Approximation and Projection for

cal statistics, pages 50-60.

Publishing.

383. Springer.

23(1):128-147.

21(140):1-67.

10

The

of whether one of two random variables is stochasti-

cally larger than the other. The annals of mathemati-

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, pages 2383–2392, Austin, Texas. Association for Computational Linguistics.

790

791

793

796

799

806

807

811

812

814

815

816

817

818

827

832

836

841

844

845

846

- Abhilasha Ravichander, Alan W Black, Shomir Wilson, Thomas Norton, and Norman Sadeh. 2019. Question answering for privacy policies: Combining computational and legal perspectives. 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 4949–4959, Hong Kong, China. Association for Computational Linguistics.
- Sebastian Ruder, Noah Constant, Jan Botha, Aditya Siddhant, Orhan Firat, Jinlan Fu, Pengfei Liu, Junjie Hu, Dan Garrette, Graham Neubig, and Melvin Johnson. 2021. XTREME-R: Towards more challenging and nuanced multilingual evaluation. In Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing, pages 10215–10245, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Victor Sanh, Albert Webson, Colin Raffel, Stephen Bach, Lintang Sutawika, Zaid Alyafeai, Antoine Chaffin, Arnaud Stiegler, Arun Raja, Manan Dey, M Saiful Bari, Canwen Xu, Urmish Thakker, Shanya Sharma Sharma, Eliza Szczechla, Taewoon Kim, Gunjan Chhablani, Nihal Nayak, Debajyoti Datta, Jonathan Chang, Mike Tian-Jian Jiang, Han Wang, Matteo Manica, Sheng Shen, Zheng Xin Yong, Harshit Pandey, Rachel Bawden, Thomas Wang, Trishala Neeraj, Jos Rozen, Abheesht Sharma, Andrea Santilli, Thibault Fevry, Jason Alan Fries, Ryan Teehan, Teven Le Scao, Stella Biderman, Leo Gao, Thomas Wolf, and Alexander M Rush. 2022. Multitask prompted training enables zero-shot task generalization. In International Conference on Learning Representations.
 - Rico Sennrich, Barry Haddow, and Alexandra Birch. 2016. Neural machine translation of rare words with subword units. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1715–1725, Berlin, Germany. Association for Computational Linguistics.
- Mukund Srinath, Shomir Wilson, and C Lee Giles. 2021. Privacy at scale: Introducing the PrivaSeer corpus of web privacy policies. In Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers), pages 6829–6839, Online. Association for Computational Linguistics.
- Mukund Sundararajan, Ankur Taly, and Qiqi Yan. 2017. Axiomatic attribution for deep networks. *CoRR*, abs/1703.01365.

Alex Wang, Yada Pruksachatkun, Nikita Nangia, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel Bowman. 2019. Superglue: A stickier benchmark for general-purpose language understanding systems. In *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc. 848

849

850

851

852

853

854

855

856

857

858

859

860

861

862

863

864

865

866

867

868

869

870

871

872

873

874

875

876

877

878

879

881

882

883

884

885

886

887

888

889

890

891

892

893

894

895

896

897

Alex Wang, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel Bowman. 2018. GLUE: A multi-task benchmark and analysis platform for natural language understanding. In Proceedings of the 2018 EMNLP Workshop BlackboxNLP: Analyzing and Interpreting Neural Networks for NLP, pages 353–355, Brussels, Belgium. Association for Computational Linguistics.

Wikimedia Foundation. 2022. Wikimedia downloads.

- Shomir Wilson, Florian Schaub, Aswarth Abhilash Dara, Frederick Liu, Sushain Cherivirala, Pedro Giovanni Leon, Mads Schaarup Andersen, Sebastian Zimmeck, Kanthashree Mysore Sathyendra, N. Cameron Russell, Thomas B. Norton, Eduard Hovy, Joel Reidenberg, and Norman Sadeh. 2016. The creation and analysis of a website privacy policy corpus. In Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 1330–1340, Berlin, Germany. Association for Computational Linguistics.
- Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick von Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gugger, Mariama Drame, Quentin Lhoest, and Alexander M. Rush. 2020. Transformers: State-of-the-art natural language processing. In Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations, pages 38–45, Online. Association for Computational Linguistics.
- Yonghui Wu, Mike Schuster, Zhifeng Chen, Quoc V Le, Mohammad Norouzi, Wolfgang Macherey, Maxim Krikun, Yuan Cao, Qin Gao, Klaus Macherey, et al. 2016. Google's neural machine translation system: Bridging the gap between human and machine translation. *arXiv preprint arXiv:1609.08144*.
- Sebastian Zimmeck, Peter Story, Daniel Smullen, Abhilasha Ravichander, Ziqi Wang, Joel R Reidenberg, N Cameron Russell, and Norman Sadeh. 2019. Maps: Scaling privacy compliance analysis to a million apps. *Proc. Priv. Enhancing Tech.*, 2019:66.

A Detailed label information

arty Collection/Use,
neric, Policy
ation, Third Party
Jser Choice/Control
on,
tion
<pre>{B,I}-data-collector,</pre>
-retained,
arer, {B,I}-data-shared,
-against, {B,I}-action, O
,I}-method,
Jser Choice/Control on, tion {B,I}-data-collector, -retained, arer, {B,I}-data-shared -against, {B,I}-action, ,I}-method,

Table 5: Breakdown of labels for each PrivacyGLUE task; PolicyQA is omitted from this table since it is a reading comprehension task and does not have explicit labels like other tasks

B PrivacyGLUE task examples

Task	Input	Target
OPP-115	Revision Date: March 24th 2015	Introductory/Generic, Policy Change
PI-Extract	We may collect and share your IP ad- dress but not your email address with our business partners .	Subtask-I: 0 0 0 0 0 B-COLLECT I-COLLECT I-COLLECT 0 0 0 0 0 0 0 0 0 0 0 0 Subtask-II: 0 0 0 0 0 0 0 0 0 0 0 0 B-NOT_COLLECT I-NOT_COLLECT I-NOT_COLLECT 0 0 0 0 0 Subtask-II: 0 0 0 0 0 0 0 0 0 0 0 0 B-NOT_SHARE I-NOT_SHARE I-NOT_SHARE 0 0 0 0 0 Subtask-IV: 0 0 0 0 0 B-SHARE I-SHARE I-SHARE 0 0 0 0 0 0 0 0 0
Policy-Detection	Log in through another service: * Facebook * Google	Not Policy
PolicyIE-A	To backup and restore your Pocket AC camera log	data-collection-usage
PolicyIE-B	Access to your personal information is restricted .	Subtask-I: 0 0 B-data-provider B-data-protected I-data-protected 0 B-action 0 Subtask-II: B-method 0 0 0 0 0 0 0
PolicyQA	Question: How do they secure my data? Context: Users can visit our site anonymously	Answer: Users can visit our site anonymously
PrivacyQA	Question: What information will you collect about my usage? Answer: Location information	Relevant

Table 6: Representative examples of each PrivacyGLUE benchmark task

899

Task	Metric [†]	BERT	RoBERTa	Legal-BERT	Legal-RoBERTa	PrivBERT
OPP-115	m-F ₁	$78.4_{\pm 0.6}$	$79.5_{\pm 1.1}$	$79.6_{\pm 1.0}$	$79.1_{\pm 0.7}$	82.1 $_{\pm 0.5}$
	μ -F ₁	$84.0_{\pm 0.5}$	$85.4_{\pm 0.5}$	$84.3_{\pm 0.7}$	$84.7_{\pm 0.3}$	$87.2_{\pm 0.4}$
PI_Extract	m-F ₁	$60.0_{\pm 2.7}$	$62.4_{\pm 4.4}$	$59.5_{\pm 3.0}$	$60.5_{\pm 3.9}$	66.4 _{±3.4}
FI-Extract	μ -F ₁	$60.0_{\pm 2.7}$	$62.4_{\pm 4.4}$	$59.5_{\pm 3.0}$	$60.5_{\pm 3.9}$	66.4 $_{\pm 3.4}$
Policy-Detection	m-F ₁	$85.3_{\pm 1.8}$	$86.9_{\pm 1.3}$	$86.6_{\pm 1.0}$	$86.4_{\pm 2.0}$	$87.3_{\pm 1.1}$
	μ -F ₁	$92.1_{\pm 1.2}$	$92.7_{\pm 0.8}$	$92.7_{\pm 0.5}$	$92.4_{\pm 1.3}$	$92.9_{\pm 0.8}$
PolicyIE A	m-F ₁	$72.9_{\pm 1.7}$	$73.2_{\pm 1.6}$	$73.2_{\pm 1.5}$	$73.5_{\pm 1.5}$	75 .3 $_{\pm 2.2}$
T ONCYTE-A	μ -F ₁	$84.7_{\pm 1.0}$	$84.8_{\pm 0.6}$	$84.7_{\pm 0.5}$	$84.8_{\pm 0.3}$	86 . $2_{\pm 1.0}$
PolicyIE-B	m-F ₁	$50.3_{\pm 0.7}$	$52.8_{\pm 0.6}$	$51.5_{\pm0.7}$	$53.5_{\pm 0.5}$	55.4 $_{\pm 0.7}$
T ONCYTE-D	μ -F ₁	$50.3_{\pm 0.5}$	$54.5_{\pm 0.7}$	$52.2_{\pm 1.0}$	$53.6_{\pm 0.9}$	55.7 _{±1.3}
PolicyQA	s-F ₁	$55.7_{\pm 0.5}$	$57.4_{\pm 0.4}$	$55.3_{\pm 0.7}$	$56.3_{\pm 0.6}$	59 .3 $_{\pm 0.5}$
	EM	$28.0_{\pm 0.9}$	$30.0_{\pm 0.5}$	$27.5_{\pm 0.6}$	$28.6_{\pm 0.9}$	$\textbf{31.4}_{\pm 0.6}$
PrivacyOA	m-F ₁	$53.6_{\pm 0.8}$	$54.4_{\pm 0.3}$	$53.6_{\pm 0.8}$	$54.4_{\pm 0.5}$	55.3 _{±0.6}
ThracyQA	μ -F ₁	$90.0_{\pm 0.1}$	$90.2_{\pm 0.0}$	$90.0_{\pm 0.1}$	$90.2_{\pm 0.1}$	90.2 $_{\pm 0.1}$

C PrivacyGLUE benchmark results

Table 7: Test-set results of the PrivacyGLUE benchmark; \dagger m-F₁ refers to macro-average F₁, μ -F₁ refers to the micro-average F₁, s refers to sample-average F₁, EM refers to the exact match accuracy, metrics are reported as percentages with the following format: mean_{±standard deviation}

D Additional PrivacyQA examples from categories P and O

Category P	Category O
ID: 9227 Question: Will the app use my data for marketing purposes? Answer: We will never share with or sell the infor- mation gained through the use of Apple HealthKit, such as age, weight and heart rate data, to advertisers or other agencies without your authorization. Label: Relevant	ID: 8749 Question: Will my fitness coach share my informa- tion with others? Answer: Develop new services. Label: Irrelevant
ID: 10858Question: What information will this app have access to of mine?Answer: Information you make available to us when you open a Keep account, as set out above;Label: Relevant	 ID: 47271 Question: Who will have access to my medical information? Answer: 23andMe may share summary statistics, which do not identify any particular individual or contain individual-level information, with our qualified research collaborators. Label: Irrelevant
 ID: 18704 Question: Does it share my personal information with others? Answer: We may also disclose Non-Identifiable Information: Label: Relevant 	ID: 54904 Question: What data do you keep and for how long? Answer: We may keep activity data on a non- identifiable basis to improve our services. Label: Irrelevant
 ID: 45935 Question: Will my test results be shared with any third party entities? Answer: 23andMe may share summary statistics, which do not identify any particular individual or contain individual-level information, with our qualified research collaborators. Label: Relevant 	ID: 57239 Question: Do you sell any of our data? Answer: (c) Advertising partners: to enable the lim- ited advertisements on our service, we may share a unique advertising identifier that is not attributable to you, with our third party advertising partners, and advertising service providers, along with certain tech- nical data about you (your language preference, coun- try, city, and device data), based on our legitimate interest. Label: Relevant
 ID: 50467 Question: Can I delete my personally identifying information? Answer: (Account Deletion), we allow our customers to delete their accounts at any time. Label: Relevant 	 ID: 59334 Question: Does the app protect my account details from being accessed by other people? Answer: Note that chats with bots and Public Accounts, and communities are not end-to-end encrypted, but we do encrypt such messages when sent to the Viber servers and when sent from the Viber servers to the third party (the Public Account owner and/or additional third party tool (eg CRM solution) integrated by such owner). Label: Irrelevant

Table 8: Additional test-set examples from PrivacyQA that fall under categories P and O for PrivBERT vs. BERT; note that these examples are not paired and can therefore be compared in any order between categories

903 904

905

906

907

908

909

910 911

912

913

914 915

916

917

926

928

930

931

933

934

935

936

E Benchmark configuration

We run PrivacyGLUE benchmark tasks with the following configuration:

- We train all models for 20 epochs with a batch size of 16. We utilize a linear learning rate scheduler with a warmup ratio of 0.1 and peak learning rate of 3 × 10⁻⁵. We utilize AdamW (Loshchilov and Hutter, 2017) as our optimizer. Finally, we monitor respective metrics on the validation datasets and utilize early stopping if the validation metric does not improve for 5 epochs.
 - We use Python v3.8.13, CUDA v11.7, PyTorch v1.12.1 (Paszke et al., 2019) and Transformers v4.19.4 (Wolf et al., 2020) as our core software dependencies.
- We following use the 918 HuggingFace model 919 tags: bert-base-uncased, roberta-base, nlpaueb/legal-bert-base-uncased, 921 saibo/legal-roberta-base, 922 mukund/privbert for BERT, RoBERTa, 923 Legal-BERT, Legal-RoBERTa and PrivBERT 924 respectively. 925
 - We use 10 random seeds for each benchmark run, i.e. {0, 1, 2, 3, 4, 5, 6, 7, 8, 9}. This provides a distribution of results that can be used for statistical significance testing.
 - We run the PrivacyGLUE benchmark on a Lambda workstation with 4 × NVIDIA RTX A4000 (16 GB VRAM) GPUs for ~180 hours.
 - We use Weights and Biases v0.13.3 (Biewald, 2020) to monitor model metrics during training and for intermediate report generation.