# **TPEval: A Novel Truth-Preserving Evaluation Method for Probing** LLMs Professional Factual Knowledge Mastery

**Anonymous ACL submission** 

#### Abstract

001 Using large language models (LLMs) to solve problems in professional fields (e.g., medicine) is emerging as a research hotspot, requiring LLMs to master sufficient domainspecific factual knowledge. Recently, several LLMs achieved notable performance on multiple professional-field evaluation benchmarks. 007 However, current benchmarks generally leverage common and fixed question formulations, 010 allowing LLMs to provide correct answers 011 based on surface-level patterns in questions without mastering the underlying knowledge. 012 In this paper, we focus on this problem. We propose a general truth-preserving evaluation framework (**TPEval**) to precisely probe LLMs' mastery of factual knowledge in professional fields through distinct representations of the 017 same knowledge. Specifically, for each piece of 019 knowledge, we convert its original expression into multiple truth-preserving statements with logical transformations, presenting the knowledge in diverse ways. By leveraging these statements, the proposed framework can more precisely estimate LLMs' mastery of the specified knowledge. Given the wealth of factual knowledge in medicine, we validate the effectiveness 027 of our framework in the medical domain. We curate 6,000+ clinical facts and generate eight statements for each fact using the proposed method, evaluating the mastery of LLMs. Experimental results indicate a notable decline in LLMs' performance as the number of statements per fact increases, suggesting insufficient knowledge mastery of LLMs. Our method can serve as an effective solution for probing LLMs' knowledge mastery in professional fields.

# 1 Introduction

042

Recent years have witnessed the rapid advancement of large language models (LLMs), which have achieved considerable performance in various downstream applications (Brown et al., 2020; Ouyang et al., 2022; Touvron et al., 2023; OpenAI,



Figure 1: An example of GPT-3.5-turbo verifying two contradictory medical factual statements.

043

045

047

049

051

054

057

060

061

063

064

065

067

068

069

071

2023; Madani et al., 2023; Boiko et al., 2023) and exhibited potential in several professional fields (e.g., medicine, finance, law). Solving problems in professional fields typically requires a comprehensive and in-depth mastery of the extensive factual knowledge within the specific domain. Recently, several studies (Petroni et al., 2019; Singh et al., 2023; Nori et al., 2023b; Wu et al., 2023) have reported that some LLMs (e.g., GPT-3.5-turbo) are capable of encoding domain-specific knowledge and largely surpass previous state-of-the-art models across various benchmark datasets within professional fields. Despite the considerable performance on existing evaluation benchmarks, it has also been observed that these LLMs are not practically applicable in real-world scenarios (Thirunavukarasu et al., 2023; Wornow et al., 2023; Li et al., 2023b), resulting in a gap between the evaluation and application. This paper aims to narrow this gap by more accurately evaluating current LLMs' proficiency in mastering domain-specific knowledge.

Several evaluation benchmark datasets have been proposed to evaluate LLMs' knowledge mastery in professional fields. Most of existing benchmark datasets (Hendrycks et al., 2020; Chen et al., 2021, 2022; Jin et al., 2021; Pal et al., 2022; Ben Abacha et al., 2017) leverage QA questions to evaluate LLMs' ability to answer questions using domain knowledge, while others also utilize tradi-



Figure 2: Framework of the proposed truth-preserving evaluation approach (**TPEval**) that probes LLMs' mastery of factual knowledge using multiple truth-preserving statements describing the knowledge in diverse ways.

tional NLP tasks (e.g., NER, sentiment analysis) (Zhang et al., 2022; Maia et al., 2018; Alvarado et al., 2015). However, these benchmark datasets typically evaluate LLMs' mastery of each fact with a fixed question formulation, which LLMs may have seen in pre-training. As a result, LLMs may directly derive answers based on surface-level patterns in questions without mastering the underlying factual knowledge. This issue is also known as data contamination (Sainz et al., 2023; Zhou et al., 2023). Figure 1 illustrates an example: GPT-3.5-turbo successfully verifies a medical factual statement but fails to check its negated version, suggesting the insufficiency of evaluating LLMs with uniform question formulations.

072

083

087

100

101

102

103

104

105

106

107

If an LLM masters a fact, it should understand various expressions of that fact. Motivated by this, we propose in this paper a novel truth-preserving evaluation approach (TPEval) to precisely probe LLMs' mastery of factual knowledge in professional fields. Figure 2 presents the framework of our proposed method. Specifically, for each piece of knowledge, we transform it into a seed logical expression and deduce a series of expressions based on this seed expression. The truth-preserving nature of deductive reasoning guarantees the correctness of the generated expressions. Finally, all the expressions are transformed back into statements, where LLMs are asked to determine the truthfulness of these statements. Compared to existing benchmarks, the proposed method evaluates the same knowledge through diverse expressions, thereby reducing the influence of LLMs in memorizing superficial patterns and leading to more accurate evaluation outcomes.

mains and is **adaptable** across diverse professional fields. Given the wealth of factual knowledge in the medical domain, we validate the effectiveness of our proposed method within this domain. Specifically, we select >6,000 pieces of clinical factual knowledge and generate eight statements for each of them employing the proposed method. Utilizing the generated statements, we evaluate a total of 14 LLMs, some of which (e.g., Gemini-pro) have achieved outstanding performance on existing medical benchmarks. Experimental results demonstrate that, though several LLMs perform well when evaluated with only one statement for each piece of knowledge, their performance sharply declines with the increasing number of statements. The results indicate that current LLMs have not mastered medical knowledge to the extent reflected by existing benchmarks. Moreover, we find that current LLMs generally perform worse when dealing with negative statements, suggesting they only have a surface-level mastery of medical knowledge. Our contributions are summarized as follows:

108

110

111

112

113

114

115

116

117

118

119

120

121

122

123

124

125

126

127

128

130

131

132

133

134

136

137

138

- We introduce a novel truth-preserving evaluation framework (**TPEval**) to evaluate LLMs' mastery of factual knowledge within professional domains. By evaluating LLMs with a series of truth-preserving statements, our method mitigates the impact caused by memorizing shallow cues and data contamination.
- Applying the proposed framework, we take the medical domain as an example and evaluate the mastery of LLMs on over 6,000 pieces of clinical medical knowledge.
- Furthermore, we compare LLMs' performance on different groups of statements along 142
- Our proposed method transcends specific do-

143three dimensions (knowledge type, statement144polarity, and expression form), shedding light145on developing domain-specific LLMs.

# 2 Related Work

146

LLMs in Professional Fields Recently, several 147 famous LLMs, such as GPT-4, are reported to 148 have achieved considerable performance on evalua-149 tion benchmarks across various professional fields. For example, in the medical domain, well-known 151 LLMs Gemini-pro, Flan-PaLM, and GPT-4 achieve 152 accuracies of 67.0, 67.6, and 90.2 on a medical 153 exam benchmark MedQA (Pal and Sankarasubbu, 154 2024; Singhal et al., 2023; Nori et al., 2023b), 155 largely surpassing previous SOTA performance 156 (Liévin et al., 2023). In the financial domain, GPT-157 4 achieves notable performance on two financial QA datasets FinQA (Chen et al., 2021) and ConvFinQA (Chen et al., 2022) by 78.0 and 76.5, re-160 spectively (Li et al., 2023a), outperforming SOTA 161 models by around ten percents. However, sev-162 eral studies (Thirunavukarasu et al., 2023; Wornow 164 et al., 2023; Li et al., 2023b) demonstrate that these LLMs are not yet applicable in real-world scenarios. 165 Therefore, we aim to study in this paper the causes 166 of the gap between LLMs' benchmark performance 167 and their insufficient practical effectiveness. 168

Evaluation Benchmarks for LLMs Various 169 evaluation benchmarks have been developed in 170 recent years to examine LLMs' mastery and application of domain-specific knowledge. Current 172 evaluation benchmarks can be categorized into two 173 types: (1) QA-based benchmarks that assess LLMs 174 with multiple-choice questions, such as MMLU 175 (Hendrycks et al., 2020), FinQA (Chen et al., 176 2021), ConvFinQA (Chen et al., 2022), MedQA 177 (Jin et al., 2021), MedMCQA (Pal et al., 2022), and 178 LiveQA (Ben Abacha et al., 2017); (2) benchmarks 179 that combine traditional NLP tasks with domainspecific corpora, such as CBLUE (Zhang et al., 181 2022), FiQA (Maia et al., 2018), and FIN (Alvarado et al., 2015). However, these benchmarks 183 test LLMs' knowledge mastery with common ques-185 tions, allowing LLMs to answer solely based on surface-level patterns. This paper tackles this issue by introducing a truth-preserving evaluation method, which generates multiple statements describing the same fact in various forms. 189

#### **3** Truth-preserving Evaluation Method

#### 3.1 Principle of the Method

In this section, we introduce the principle of our proposed truth-preserving evaluation approach (TPEval), which systematically probes LLMs' mastery of domain-specific factual knowledge based on truth-preserving statement generation. We denote the language-form expression of a piece of factual knowledge as S. The existing evaluation methods generally examine whether an LLM  $\mathcal{M}$ has mastered the fact as follows:

$$m'_{\rm S} = f_{\rm S}(\mathcal{M}) \tag{1}$$

190

191

192

193

194

195

196

198

199

200

201

202

203

206

207

208

210

211

212

213

214

215

216

217

218

219

220

221

222

223

224

225

226

227

229

230

231

232

233

234

where  $f_{\rm S}$  refers to the evaluation question generated based on S, and  $m'_{\rm S} \in \{0,1\}$  denotes the evaluation result:  $m'_{\rm S} = 1$  when  $\mathcal{M}$  answers the question correctly, otherwise  $m'_{\rm S} = 0$ . In contrast, the proposed TPEval method evaluate  $\mathcal{M}$  by generating K truth-preserving statements  $\{S_i\}_{i=1}^K$  based on the same piece of factual knowledge:

$$\mathbf{p} = g(\mathbf{S}) \tag{2}$$

$$[\mathbf{q}_1, \mathbf{q}_2, \cdots, \mathbf{q}_K] = \text{Deduce}(\mathbf{p}) \tag{3}$$

$$S_i = g^{-1}(q_i), 1 \le i \le K$$
 (4)

$$\mathbf{m}_{\mathrm{S}} = [f_{\mathrm{S}_{1}}(\mathcal{M}), f_{\mathrm{S}_{2}}(\mathcal{M}), \cdots, f_{\mathrm{S}_{K}}(\mathcal{M})] \quad (5)$$

where g denotes a mapping that projects the original statement S into the associated logical form p (seed expression), and  $g^{-1}$  is its inverse operation. Deduce refers to the deductive reasoning process, and  $\{q_i\}_{i=1}^{K}$  are truth-preserving logical expressions deduced from the seed expression p. Compared with the former method, the proposed method mitigates the impact of the model memorizing specific patterns in the original expression, where the properties of deductive reasoning guarantee the validity of generated questions.

# **3.2** Truth-preserving Evaluation Framework

Built on the truth-preserving principle, we design a novel evaluation framework to evaluate LLMs' mastery of factual knowledge in professional fields more precisely.

**Statement Generation** We primarily consider factual knowledge that can be expressed by a triplet: (A, R, B), where A and B are entities, and R is the relation between them. Such type of factual knowledge is usually expressed as "A is/has the [attribute/relation] of/with B" in the language form



Figure 3: An overview of the truth-preserving logical transformation module in the proposed TPEval framework.

(S). The logical form of this expression is formulated as p = R(A, B), where R(A, B) is a predicate that denotes "A has the relation R with B". For example, Drug\_of(A, B) presents a medical fact "A is the therapeutic medication of disease B."

235

237

241

242

243

245

246

247

254

257

258

261

262

263

269

270

271

We leverage deductive reasoning to transform the seed expression p into multiple truth-preserving expressions. Figure 3 depicts the whole transformation process. Specifically, we combine pwith each piece of commonsense knowledge  $c_i$ to obtain a conclusion  $q_i$  based on syllogism. The commonsense knowledge encompasses the principles and rules in the specific domain, serving as the foundation for solving problems in this domain. In our framework, we focus on a type of commonsense knowledge represented as  $c_i =$  $\forall X \forall Y(R(X,Y) \Rightarrow Q_i(X,Y))$ . To illustrate, consider a medical commonsense knowledge instance: "If a drug is a therapeutic drug of a disease, then the drug can be used to treat a person who suffers from the disease." Here, the phrases "a drug ... for a disease" and "the drug ... from the disease" can be denoted as R(X, Y) and  $Q_i(X, Y)$ , respectively. Combining the factual knowledge p with each  $c_i$ , we have:

$$\forall X \forall Y(R(\mathbf{X}, \mathbf{Y}) \Rightarrow Q_i(\mathbf{X}, \mathbf{Y})) \qquad (c_i)$$

$$\frac{R(A,B)}{\therefore Q_i(A,B)} \qquad (p)$$

It is worth noting that if  $\mathcal{M}$  incorrectly predicts the truth value of  $Q_i(X, Y)$ ), it indicates that the LLM either lacks the corresponding factual knowledge or commonsense knowledge. In our framework, we choose commonsense knowledge that is clear, common, and simple as much as possible to ensure that the generated expressions are easy to understand. We generate K/2 expressions through this process and also create another K/2 expressions based on the double negation rule:  $q_{i+K/2} = \neg(\neg q_i), 1 \le i \le K/2$ . This process



Figure 4: A comparison between the statement generation procedures based on positive and negative triplets.

yields a total of K truth-preserving expressions.

Subsequently, each logical expression  $q_i$  is transformed into a statement  $S_i$  along with a label  $l_i \in {\mathbf{T}, \mathbf{F}}$ . For  $q_i$  where  $i \leq K/2$ , we transform it using the predicate  $Q_i$  and set  $l_i$  as the true value of p. For  $q_{i+K/2}$  (double negation), it is transformed using the negated predicate  $\neg Q_i$ , with  $l_i$  set to the opposite value of p. It's worth noting that  $S_{i+K/2}$  can be derived from  $S_i$  by incorporating negative words (e.g., "not"). These statements are applied to evaluate  $\mathcal{M}$ 's mastery of the same factual knowledge.

**Generation from Negative Triplets** The statement generation method above is designed to generate statements based on positive knowledge triplets. As a result, p always holds in this scenario, causing statements generated by positive predicates to always be true, while those generated by negative predicates are consistently false. Therefore, generating statements exclusively from positive triplets could introduce bias, as LLMs may predict outcomes solely based on the presence of negation cues. Moreover, recognizing the absence of certain relations between entities is crucial for LLMs (e.g., a drug cannot be used to treat a specific disease). Therefore, we also generate statements from nega-

299

323

327

330

335

336

337

339

tive triplets, where these statements share the same formats as those generated from the corresponding positive triplets but have opposite labels.

Figure 4 compares the generation procedure based on positive triplets with negative ones. Specifically, for each p = R(A, B), we sample an A' that satisfies  $\neg R(A', B)$ . Then, we treat  $\neg R$ 306 as a new predicate and employ the same method to generate truth-preserving statements. To ensure consistency in the formats of the generated statements with those from positive triplets, we only choose commonsense knowledge where the inverse 311 proposition  $(\forall X \forall Y (\neg R(\mathbf{X}, \mathbf{Y}) \Rightarrow \neg Q_i(\mathbf{X}, \mathbf{Y})))$ also holds. Consequently, for every  $q_i = Q_i(A, B)$ , 313 we have  $q'_i = \neg Q_i(A', B)$  holding true as well. 314 Thus,  $S'_i$  can be derived by replacing A with A' in  $S_i$ , and the corresponding label is exactly opposite to the original label:  $l'_i = \neg l_i$ . By generating 317 statements from both positive and negative triplets, 318 we can effectively prevent LLMs from verifying 319 statements solely based on surface-level patterns and guarantee the completeness of the proposed evaluation framework. 322

> **LLM Evaluation** The proposed truth-preserving evaluation principle does not restrict the types of evaluation questions. In our framework, we evaluate LLMs with statement verification questions, asking LLMs to determine whether the given statement  $S_i$  is true or false:

$$f_{\mathbf{S}_i}(\mathcal{M}) = \mathbb{1}(\mathcal{M}(\mathbf{S}_i) = l_i), 1 \le i \le K \quad (6)$$

Where  $\mathcal{M}(S_i) \in \{\mathbf{T}, \mathbf{F}\}$  denotes the LLM's prediction of  $S_i$ 's truth value,  $\mathbb{1}(\cdot)$  represents the characteristic function that equals 1 when the enclosed expression is true, and 0 otherwise.

We measure  $\mathcal{M}$ 's performance on a dataset that includes N pieces of knowledge by two modes: multi-statement **average** evaluation and multi-statement **joint** evaluation. These two modes calculate accuracies in different granularities:

$$a_{\text{avg}} = \frac{1}{N} \frac{1}{K} \sum_{i=1}^{N} \sum_{j=1}^{K} f_{\text{S}_{j}^{i}}(\mathcal{M})$$
(7)

$$a_{\text{joint}} = \frac{1}{N} \sum_{i=1}^{N} \prod_{j=1}^{K} f_{\mathbf{S}_{j}^{i}}(\mathcal{M})$$
(8)

Here,  $S_j^i$  denotes the  $j^{th}$  statement derived from the ith piece of knowledge. The former calculates the accuracy of  $\mathcal{M}$  across all statements. In contrast, the latter calculates the accuracy across knowledge,



Figure 5: The knowledge structure of the proposed disease-centric medical knowledge base DiseK.

considering a piece of knowledge being correctly answered if and only if **all the related statements** are judged correctly.

# 4 Experiments

# 4.1 Experiment Setup

**Dataset Generation** We select the medical domain to validate the effectiveness of the proposed method because it encompasses a wide range of factual knowledge. Moreover, several existing LLMs have been reported to achieve impressive performance on various medical benchmarks (Nori et al., 2023a; Singhal et al., 2023; Nori et al., 2023b). One of the primary goals in the medical domain is to diagnose and treat diseases. Motivated by this, we construct a large-scale disease-centric medical knowledge base, **DiseK**, covering 1,000 highfrequency diseases across four crucial knowledge aspects closely related to disease diagnosis and treatment. LLMs must master this fundamental medical knowledge to assist doctors in diagnosing and treating corresponding diseases (Wu et al., 2018; Liang et al., 2019). Figure 5 depicts the knowledge structure of the proposed knowledge base. Detailed statistics and annotation details of DiseK are provided in Appendix A.

DiseK contains more than 24k pieces of factual knowledge, each of them can be represented as a triplet (A, R, D), denoting "A is the R of disease D." Here, R corresponds to one of four knowledge aspects: symptoms, affected sites, therapeutic drugs, and surgical procedures, and A is an entity of R. To reduce computational cost in evaluation, we select a positive triplet (A, R, D) and a negative triplet (A',  $\neg R$ , D) for each pair of (R, D), resulting in 3,167 positive triplets and 3,167 negative triplets<sup>1</sup>. Subsequently, we generate a truth-

<sup>&</sup>lt;sup>1</sup>Some diseases may not have certain aspects of knowledge,

381	preserving evaluation dataset TPDiseK using the
382	proposed method, where each fact is evaluated by
383	K = 8 statements. We initially transform a piece
384	of knowledge into four truth-preserving statements
385	based on commonsense knowledge. Two of the
386	generated statements are derived through both triv-
387	ial reasoning $(S_1)$ and reverse reasoning $(S_2)$ , ex-
388	hibiting minimal deviation from the original repre-
389	sentation. The remaining two statements ( $S_3$ and
390	S <sub>4</sub> ) are generated by applying the factual knowl-
391	edge to specific medical cases, thereby evaluating
392	LLMs' capability of handling specific problems
393	with the knowledge acquired. Another four state-
394	ments with opposite labels $(S_5 \text{ to } S_8)$ are generated
395	by negating these four statements. The statement
396	templates are meticulously designed to ensure
397	they are easily understandable and faithfully
398	express the meaning of corresponding logical
399	expressions. To summarize, TPDiseK consists of
400	6,334 knowledge triplets (positive and negative),
401	each comprising 8 statements for evaluation. More
402	details of TPDiseK are provided in Appendix B.

403 **Evaluation Setting** We primarily assess LLMs with the **five-shot in-context learning** strategy 404 (Brown et al., 2020), where five demonstrative 405 question-answer pairs are presented before the test 406 question, guiding LLMs to produce answers consis-407 tent with the provided examples. We also examined 408 the zero-shot performance of LLMs and found that 409 the trend is similar to that observed in the five-shot 410 setting. Therefore, we provide the zero-shot results 411 for complement in Appendix D. We report the accu-412 racies measured by the average and joint evaluation 413 modes introduced in Sec 3.2. Note that the joint 414 accuracy can be regarded as the proportion of 415 factual knowledge truly mastered by LLMs. We 416 present more details in Appendix C. 417

Evaluated Models In our study, we evaluate a total of 14 well-known general and medical-419 domain-specific LLMs on the proposed TPDiseK dataset: (1) general LLMs: ChatGLM (6B) (Du et al., 2022), Bloomz-mt (7.1B) (Muennighoff et al., 2023), Llama2 (7B,70B) (Touvron et al., 2023), Vicuna (7B,13B) (Zheng et al., 2023), GPT-3.5-turbo (Ouyang et al., 2022), ERNIE-Bot-turbo (Sun et al., 2021) and Gemini-pro (Team et al., 2023); (2) medical-domain-specific LLMs: Pulse (7B) (Zhang et al., 2023), ClinicalCamel (70B) (Toma et al., 2023), Meditron (7B,70B) (Chen et al.,

418

420

421

422

423

424

425

426

427

428

429

Models	<b>DiseK</b> (K=1)	<b>TPDiseK</b> ( <i>K</i> =8)
Random	50.0	50.0
ChatGLM-6B	52.6	48.7
Llama2-7B	52.8	51.7
Pulse-7B	54.9	51.4
Bloomz-mt-7.1B	52.3	48.8
Vicuna-7B	59.0	52.0
Meditron-7B	50.0	50.0
Vicuna-13B	61.2	54.0
Llama2-70B	65.9	56.5
ClinicalCamel-70B	74.4	64.4
Meditron-70B	70.3	57.0
Med42-70B	71.7	64.1
ERNIE-Bot-turbo	69.8	56.7
GPT-3.5-turbo	73.5	60.5
Gemini-pro	78.0	73.1

Table 1: Comparison of LLMs' average accuracy on the original medical evaluation dataset (DiseK) with that on the truth-preserving dataset (TPDiseK) generated by the proposed framework TPEval. K: the number of evaluated statements per knowledge triplet.

2023) and Med42 (70B) (Christophe et al., 2023). We have not assess LLMs that are either too expensive (e.g., GPT-4 (OpenAI, 2023)) or not publicly available (e.g., MedPaLM (Singhal et al., 2023)).

430

431

432

433

434

435

436

437

438

439

440

441

442

443

444

445

446

447

448

449

450

451

452

453

454

455

456

# 4.2 Results

#### 4.2.1 Overall Performance

We initially compare the performance of LLMs on the original dataset DiseK with that on the truth-preserving dataset TPDiseK. LLMs' performance on DiseK is measured by their performance on the first type of statements  $(S_1)$ , which expresses the knowledge in a trivial way (A is the R of B). The experimental results are provided in Table 1 and measured by average accuracy. We observe that LLMs  $\leq$ 13B generally perform poorly on both datasets, while several 70B LLMs and commercial LLMs (Gemini-pro, ERNIE-Botturbo, GPT-3.5-turbo) achieve notable performance on the trivial representation of DiseK. However, their performance significantly declines when using the proposed truth-preserving evaluation method. LLMs like GPT-3.5-turbo, Meditron-70B, and ClinicalCamel-70B exhibit a performance decline of over 10% on TPDiseK compared to their performance on DiseK. Even the best-performing Gemini-pro demonstrates a performance gap of 4.9% between DiseK and TPDiseK.

such as therapeutic medication or surgeries.



Figure 6: Performance of LLMs evaluated by increasing number of statements per triplet. Left: joint accuracy (proportion of triplets that all the related statements are correctly predicted); Right: average accuracy. Dotted lines: LLMs  $\leq$ 13B; Solid lines: LLMs >13B.

486

487

Subsequently, we study the performance of LLMs by gradually increasing the number of truthpreserving statements per knowledge triplet from 1 (DiseK) to 8. Statements are gradually added to the evaluation in the order of  $S_1$  to  $S_8$  defined in Sec 4.1. The experimental results are presented in Figure 6, showing accuracies for both joint evaluation (proportion of triplets where all related statements are correctly predicted) and average evaluation. The experimental results indicate a significant decrease in the joint accuracy of LLMs as the number of truth-preserving statements increases, declining much faster than their average accuracy. Surprisingly, some famous LLMs like GPT-3.5-turbo, ERNIE-Bot-turbo, and Llama2-70B even achieve comparable joint accuracies with LLMs <13B (dotted lines) when evaluated by all the statements. The results suggest that while these LLMs may memorize more surface-level patterns of knowledge expressions than smaller models, they do not truly master a broader range of knowledge. Gemini-pro and Med42-70B significantly outperform other LLMs regarding joint accuracy, suggesting a relatively comprehensive mastery of medical knowledge than others.

4.2.2 Performance across Statement Types

To further investigate current LLMs' performance on handling different types of statements, we categorize statements based on three criteria: (1) the type of knowledge triplets (positive/negative); (2) statement polarity (positive/negative); (3) the expression forms of statements determined by predicates  $Q_i, 1 \le i \le 4$ . We only present the top-5 performing LLMs here for convenience and provide the results of other LLMs in Appendix D. The performance is measured by average accuracy.

488

489

490

491

492

493

494

495

496

497

498

499

500

501

502

503

504

505

507

509

510

511

512

513

514

515

516

517

**Knowledge Type** Figure 7a presents LLMs' performance on statements generated from different types of knowledge triplets. We observe that LLMs exhibit varying degrees of proficiency in different types of factual knowledge (positive/negative). Gemini-pro and Med42-70B perform significantly better on negative triplets, indicating that they are more accurate in determining the absence of a relationship between two medical entities. It is worth noting that Med42 performs slightly worse than random guessing on positive triplets, suggesting a tendency to indicate no given relationship between two medical entities. In contrast, the other three LLMs achieve more balanced performance between different triplet types.

**Statement Polarity** Figure 7b examines how LLMs handle statements of different polarities: positive ( $S_1$  to  $S_4$ ) and negative ( $S_5$  to  $S_8$ ). The results show that **current LLMs generally perform significantly worse on negative statements**. Some LLMs' performance (Meditron-70B, GPT-3.5-turbo) even approaches or is inferior to random guessing. Gemini-pro significantly outperforms others on negative statements, achieving the most balanced performance on the two polarities.







(a) Performance on statements from different types of knowledge triplets.

(b) Performance on statements with different polarities.

(c) Performance on statements with varied expression forms (denoted by  $Q_i$ ).

554

555

556

557

558

559

560

561

562

563

564

565

566

567

568

569

570

571

572

573

574

575

576

577

578

579

581

582

583

584

585

586

587

589

Figure 7: LLMs performance (average accuracy) grouped by different statement types according to three criteria. Only the top 5 performing LLMs are selected for convenience.

**Expression Form** We further investigate LLMs' performance on various expression forms in Figure 7c, where each  $Q_i$  denotes a specific expression form. Generally, LLMs perform notably worse on the last expression form than other forms. It may be because the expression form combines factual knowledge with specific cases and involves a contrapositive transform from the third expression. Therefore, this expression form occurs **less frequently** in existing medical corpora, making it less likely to be predicted by surface-level cues. Gemini-pro is the only LLM that achieves over 60% accuracy on every single expression form, suggesting its comprehensive mastery of medical factual knowledge compared to other LLMs.

#### 5 Discussion

518

519

520

522

524

527

528

530

532

533

**Effectiveness of Truth-preserving Evaluation** 534 The experimental results reveal that the LLMs' mastery of medical knowledge, as assessed by the proposed truth-preserving method, is notably 537 lower than that evaluated by the traditional methods. 538 Moreover, their performance declines sharply as the number of statements per knowledge triplet in-540 creases. These findings suggest that current LLMs 541 generate responses by memorizing surface-level 542 patterns of knowledge expressions without gen-543 uinely mastering the underlying knowledge. Furthermore, the examination across various statement types indicates that these LLMs struggle to manipulate factual knowledge through basic transforma-547 tions (e.g., negation), instead relying on specific 549 forms of knowledge presentation. Some LLMs also prefer to memorize specific types of knowledge. All these findings demonstrate that the proposed truth-preserving evaluation method can serve as an effective solution to evaluate LLMs' mastery of 553

factual knowledge in professional fields.

Insights into Developing Domain-specific LLMs The experimental results demonstrate that current LLMs lack an in-depth mastery of medical factual knowledge. Therefore, training LLMs to genuinely acquire domain knowledge is crucial rather than merely memorizing surface-level patterns in corresponding expressions. Expanding the training data by expressing factual knowledge in diverse ways may potentially enhance LLMs' knowledge mastery in professional fields.

#### 6 Conclusion

Understanding LLMs' mastery of domain knowledge is crucial for their application in real-world scenarios. In this paper, we introduce a novel truthpreserving evaluation method (TPEval) to systematically evaluate LLMs' proficiency in factual knowledge of professional fields. The proposed method evaluates the same knowledge using multiple statements that present it in diverse ways, leading to a more precise estimation of LLMs' knowledge mastery. We investigate the proposed method in the medical domain based on >6,000 knowledge triplets from a medical knowledge base. The results reveal a significant drop in the performance of existing LLMs when assessed using the proposed TPEval method compared to traditional evaluation methods, suggesting that they lack an indepth mastery of medical factual knowledge. Our method can serve as an effective solution for evaluating the knowledge mastery of LLMs in professional fields, shedding light on developing domainspecific LLMs. In the future, we aim to refine this method further by integrating it with more question formats, such as question answering, and applying it to more professional domains.

# Limitations

590

592

593

596

601

606

610

611

614

616

617

619

625

626

627

628

629

631

634

637

641

One of the major limitation of our study is that we only verify the effectiveness of the proposed method in the medical field due to space constraints. However, the principle of our method is decoupled with specific professional fields, and the experimental results in this paper is already sufficient to demonstrate the effectiveness of our method. In future, we will utilize the proposed method to build more truth-preserving datasets in other professional fields to promote related researches.

> Moreover, while our study evaluated several well-known general and medical-domain-specific LLMs, some other notable models like GPT-4 and MedPaLM were excluded. This was due to either their high costs (it would require \$800 to evaluate GPT-4 on TPDiseK) or their unavailability for public access (e.g., MedPaLM). We will keep evaluating other LLMs in future if feasible.

# References

- Julio Cesar Salinas Alvarado, Karin Verspoor, and Timothy Baldwin. 2015. Domain adaption of named entity recognition to support credit risk assessment. In *Proceedings of the Australasian Language Technology Association Workshop 2015*, pages 84–90.
- Asma Ben Abacha, Eugene Agichtein, Yuval Pinter, and Dina Demner-Fushman. 2017. Overview of the medical question answering task at trec 2017 liveqa. In *TREC 2017*.
- Daniil A Boiko, Robert MacKnight, Ben Kline, and Gabe Gomes. 2023. Autonomous chemical research with large language models. *Nature*, 624(7992):570– 578.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. 2020. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901.
- Zeming Chen, Alejandro Hernández Cano, Angelika Romanou, Antoine Bonnet, Kyle Matoba, Francesco Salvi, Matteo Pagliardini, Simin Fan, Andreas Köpf, Amirkeivan Mohtashami, et al. 2023. Meditron-70b: Scaling medical pretraining for large language models. *arXiv preprint arXiv:2311.16079*.
- Zhiyu Chen, Wenhu Chen, Charese Smiley, Sameena Shah, Iana Borova, Dylan Langdon, Reema Moussa, Matt Beane, Ting-Hao Huang, Bryan Routledge, and William Yang Wang. 2021. FinQA: A dataset of numerical reasoning over financial data. In *Proceedings* of the 2021 Conference on Empirical Methods in Natural Language Processing, pages 3697–3711, Online

and Punta Cana, Dominican Republic. Association for Computational Linguistics.

642

643

644

645

646

647

648

649

650

651

652

653

654

655

656

657

658

659

660

661

663

664

665

666

667

668

669

670

671

672

673

674

675

676

677

678

679

680

681

682

683

684

685

686

687

688

689

690

691

692

693

694

695

- Zhiyu Chen, Shiyang Li, Charese Smiley, Zhiqiang Ma, Sameena Shah, and William Yang Wang. 2022. ConvFinQA: Exploring the chain of numerical reasoning in conversational finance question answering. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 6279– 6292, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.
- Clément Christophe, Avani Gupta, Nasir Hayat, Praveen Kanithi, Ahmed Al-Mahrooqi, Prateek Munjal, Marco Pimentel, Tathagata Raha, Ronnie Rajan, and Shadab Khan. 2023. Med42 - a clinical large language model.
- Zhengxiao Du, Yujie Qian, Xiao Liu, Ming Ding, Jiezhong Qiu, Zhilin Yang, and Jie Tang. 2022. GLM: General language model pretraining with autoregressive blank infilling. In Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 320–335, Dublin, Ireland. Association for Computational Linguistics.
- Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. 2020. Measuring massive multitask language understanding. arXiv preprint arXiv:2009.03300.
- Di Jin, Eileen Pan, Nassim Oufattole, Wei-Hung Weng, Hanyi Fang, and Peter Szolovits. 2021. What disease does this patient have? a large-scale open domain question answering dataset from medical exams. *Applied Sciences*, 11(14):6421.
- Xianzhi Li, Samuel Chan, Xiaodan Zhu, Yulong Pei, Zhiqiang Ma, Xiaomo Liu, and Sameena Shah. 2023a. Are ChatGPT and GPT-4 general-purpose solvers for financial text analytics? a study on several typical tasks. In Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing: Industry Track, pages 408–422, Singapore. Association for Computational Linguistics.
- Yinheng Li, Shaofei Wang, Han Ding, and Hang Chen. 2023b. Large language models in finance: A survey. In *Proceedings of the Fourth ACM International Conference on AI in Finance*, pages 374–382.
- Huiying Liang, Brian Y Tsui, Hao Ni, Carolina CS Valentim, Sally L Baxter, Guangjian Liu, Wenjia Cai, Daniel S Kermany, Xin Sun, Jiancong Chen, et al. 2019. Evaluation and accurate diagnoses of pediatric diseases using artificial intelligence. *Nature medicine*, 25(3):433–438.
- Valentin Liévin, Andreas Geert Motzfeldt, Ida Riis Jensen, and Ole Winther. 2023. Variational opendomain question answering. In *International Conference on Machine Learning*, pages 20950–20977. PMLR.

811

755

Ali Madani, Ben Krause, Eric R Greene, Subu Subramanian, Benjamin P Mohr, James M Holton, Jose Luis Olmos Jr, Caiming Xiong, Zachary Z Sun, Richard Socher, et al. 2023. Large language models generate functional protein sequences across diverse families. *Nature Biotechnology*, pages 1–8.

703

704

705

709

710

711

712

714

715

716

717

718

719

720

721

722

723

724

725

729

731

732

733

737

740

741

742

743

744

745

746

747

749

750

751

754

- Macedo Maia, Siegfried Handschuh, André Freitas, Brian Davis, Ross McDermott, Manel Zarrouk, and Alexandra Balahur. 2018. Www'18 open challenge: financial opinion mining and question answering. In *Companion proceedings of the the web conference* 2018, pages 1941–1942.
- Niklas Muennighoff, Thomas Wang, Lintang Sutawika, Adam Roberts, Stella Biderman, Teven Le Scao, M Saiful Bari, Sheng Shen, Zheng Xin Yong, Hailey Schoelkopf, Xiangru Tang, Dragomir Radev, Alham Fikri Aji, Khalid Almubarak, Samuel Albanie, Zaid Alyafeai, Albert Webson, Edward Raff, and Colin Raffel. 2023. Crosslingual generalization through multitask finetuning. In Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 15991–16111, Toronto, Canada. Association for Computational Linguistics.
- Harsha Nori, Nicholas King, Scott Mayer McKinney, Dean Carignan, and Eric Horvitz. 2023a. Capabilities of gpt-4 on medical challenge problems. *arXiv preprint arXiv:2303.13375*.
- Harsha Nori, Yin Tat Lee, Sheng Zhang, Dean Carignan, Richard Edgar, Nicolo Fusi, Nicholas King, Jonathan Larson, Yuanzhi Li, Weishung Liu, et al. 2023b. Can generalist foundation models outcompete special-purpose tuning? case study in medicine. arXiv preprint arXiv:2311.16452.
- OpenAI. 2023. Gpt-4 technical report.
  - Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. 2022. Training language models to follow instructions with human feedback. *Advances in Neural Information Processing Systems*, 35:27730–27744.
  - Ankit Pal and Malaikannan Sankarasubbu. 2024. Gemini goes to med school: Exploring the capabilities of multimodal large language models on medical challenge problems & hallucinations.
  - Ankit Pal, Logesh Kumar Umapathi, and Malaikannan Sankarasubbu. 2022. Medmcqa: A large-scale multisubject multi-choice dataset for medical domain question answering. In *Proceedings of the Conference on Health, Inference, and Learning,* volume 174 of *Proceedings of Machine Learning Research,* pages 248–260. PMLR.
  - Fabio Petroni, Tim Rocktäschel, Sebastian Riedel, Patrick Lewis, Anton Bakhtin, Yuxiang Wu, and Alexander Miller. 2019. Language models as knowledge bases? 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 2463–2473, Hong Kong, China. Association for Computational Linguistics.

- Oscar Sainz, Jon Campos, Iker García-Ferrero, Julen Etxaniz, Oier Lopez de Lacalle, and Eneko Agirre. 2023. NLP evaluation in trouble: On the need to measure LLM data contamination for each benchmark. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 10776–10787, Singapore. Association for Computational Linguistics.
- Gagandeep Singh, Yue Pan, Jesus Andres-Ferrer, Miguel Del-Agua, Frank Diehl, Joel Pinto, and Paul Vozila. 2023. Large scale sequence-to-sequence models for clinical note generation from patient-doctor conversations. In *Proceedings of the 5th Clinical Natural Language Processing Workshop*, pages 138– 143, Toronto, Canada. Association for Computational Linguistics.
- Karan Singhal, Shekoofeh Azizi, Tao Tu, S Sara Mahdavi, Jason Wei, Hyung Won Chung, Nathan Scales, Ajay Tanwani, Heather Cole-Lewis, Stephen Pfohl, et al. 2023. Large language models encode clinical knowledge. *Nature*, pages 1–9.
- Yu Sun, Shuohuan Wang, Shikun Feng, Siyu Ding, Chao Pang, Junyuan Shang, Jiaxiang Liu, Xuyi Chen, Yanbin Zhao, Yuxiang Lu, et al. 2021. Ernie 3.0: Large-scale knowledge enhanced pre-training for language understanding and generation. *arXiv preprint arXiv:2107.02137*.
- Gemini Team, Rohan Anil, Sebastian Borgeaud, Yonghui Wu, Jean-Baptiste Alayrac, Jiahui Yu, Radu Soricut, Johan Schalkwyk, Andrew M Dai, Anja Hauth, et al. 2023. Gemini: a family of highly capable multimodal models. *arXiv preprint arXiv:2312.11805*.
- Arun James Thirunavukarasu, Darren Shu Jeng Ting, Kabilan Elangovan, Laura Gutierrez, Ting Fang Tan, and Daniel Shu Wei Ting. 2023. Large language models in medicine. *Nature medicine*, 29(8):1930– 1940.
- Augustin Toma, Patrick R Lawler, Jimmy Ba, Rahul G Krishnan, Barry B Rubin, and Bo Wang. 2023. Clinical camel: An open-source expert-level medical language model with dialogue-based knowledge encoding. *arXiv preprint arXiv:2305.12031*.
- Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. 2023. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*.
- Michael Wornow, Yizhe Xu, Rahul Thapa, Birju Patel, Ethan Steinberg, Scott Fleming, Michael A Pfeffer, Jason Fries, and Nigam H Shah. 2023. The shaky foundations of large language models and foundation models for electronic health records. *npj Digital Medicine*, 6(1):135.

Ji Wu, Xien Liu, Xiao Zhang, Zhiyang He, and Ping Lv. 2018. Master clinical medical knowledge at certificated-doctor-level with deep learning model. *Nature communications*, 9(1):4352.

812

813

814 815

816 817

818

819

820

822

823 824

825

826

827

830

831

833

838

839

840 841

842

843

844

- Shijie Wu, Ozan Irsoy, Steven Lu, Vadim Dabravolski, Mark Dredze, Sebastian Gehrmann, Prabhanjan Kambadur, David Rosenberg, and Gideon Mann. 2023. Bloomberggpt: A large language model for finance. arXiv preprint arXiv:2303.17564.
- Ningyu Zhang, Mosha Chen, Zhen Bi, Xiaozhuan Liang, Lei Li, Xin Shang, Kangping Yin, Chuanqi Tan, Jian Xu, Fei Huang, Luo Si, Yuan Ni, Guotong Xie, Zhifang Sui, Baobao Chang, Hui Zong, Zheng Yuan, Linfeng Li, Jun Yan, Hongying Zan, Kunli Zhang, Buzhou Tang, and Qingcai Chen. 2022. CBLUE: A Chinese biomedical language understanding evaluation benchmark. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 7888–7915, Dublin, Ireland. Association for Computational Linguistics.
  - Xiaofan Zhang, Kui Xue, and Shaoting Zhang. 2023. Pulse: Pretrained and unified language service engine.
  - Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric Xing, et al. 2023. Judging llm-as-a-judge with mt-bench and chatbot arena. *arXiv preprint arXiv:2306.05685*.
  - Kun Zhou, Yutao Zhu, Zhipeng Chen, Wentong Chen, Wayne Xin Zhao, Xu Chen, Yankai Lin, Ji-Rong Wen, and Jiawei Han. 2023. Don't make your llm an evaluation benchmark cheater. *arXiv preprint arXiv:2311.01964*.

# A Details of DiseK

Knowledge Aspects	#Uniq	Avg.
#Symptoms	4,163	12.37
#Affected Sites	387	0.75
#Therapeutic Drugs	1,782	7.89
<b>#Surgical Procedures</b>	1,932	3.40

Table 2: Statistics of the proposed knowledge base **DiseK**. #Uniq: the number of unique entities that appear in the knowledge bases. Avg: the average number of entities associated to each disease. Note that some diseases may affect multiple organs and do not have specific affected sites.

We propose in this paper a large-scale diseasecentric medical knowledge base **DiseK**, which involves a total of 1,000 high-frequency dieases along with 4 knowledge aspects that are essential to the diagnosis and treatment of diseases. The four knowledge aspects in DiseK are listed below:

- **Symptoms**: Physical or mental feature that indicates the presence of the disease.
- Affected sites: Specific parts of the body that are impacted or harmed by the disease.
- **Therapeutic Drugs**: Pharmaceutical substances prescribed to manage, alleviate, or cure the symptoms and effects of the disease.
- **Surgical Procedures**: Medical procedures that treat the disease, involving the cutting, repairing, or removal of body parts.

The diseases in DiseK are selected by filtering out top 1,000 most frequently occurring diseases among approximately four million medical records gathered from over 100 hospitals across five cities. After that, we ask 20 medical experts to annotate the related knowledge of these diseases from the four aspects introduced above. To ensure the quality of annotation, we first leverage a retrieval module to automatically retrieve the related knowledge from medical books, literature, and the Internet. Then the experts are asked to recheck the retrieved content, and revise the incorrect parts. We find that the utilized retrieve-and-check annotation framework can alleviate the burden of annotators while ensuring consistency in the annotations.

The proposed medical knowledge base is annotation by 20 medical experts over a period of approximately 3 months. The medical experts are employees in our company and have obtained medical practitioner licenses. The data collection process is approved by an ethics review board, and we have obtained approval from the person in charge to use the annotated data for research.

# **B** Details of TPDiseK

Triplet	Туре	Commonsense Expression
Pos	Triv.	$R(\mathbf{X}, \mathbf{Y}) \to R(\mathbf{X}, \mathbf{Y})$
	Rev.	$R(\mathbf{X}, \mathbf{Y}) \to R^{-1}(\mathbf{Y}, \mathbf{X})$
	Spec.	$R(\mathbf{X}, \mathbf{Y}) \to (Has(\mathbf{Y}) \to P(\mathbf{Y}))$
	Contr.	$R(\mathbf{X}, \mathbf{Y}) \to (\neg P(\mathbf{Y}) \to \neg Has(\mathbf{Y}))$
Neg	Triv.	$\neg R(\mathbf{X}, \mathbf{Y}) \rightarrow \neg R(\mathbf{X}, \mathbf{Y})$
	Rev.	$\neg R(\mathbf{X}, \mathbf{Y}) \rightarrow \neg R^{-1}(\mathbf{Y}, \mathbf{X})$
	Spec.	$\neg R(\mathbf{X}, \mathbf{Y}) \rightarrow \neg (Has(\mathbf{Y}) \rightarrow P(\mathbf{Y}))$
	Contr.	$\neg R(\mathbf{X}, \mathbf{Y}) \rightarrow \neg (\neg P(\mathbf{Y}) \rightarrow \neg Has(\mathbf{Y}))$

Table 3: The expressions of commonsense knowledge leveraged in our experiments. All the expressions hold for all X and all Y. Pos: positive triplets; Neg: negative triplets.  $R^{-1}$ : The R of the disease Y include X. Has(Y): A patient **only** suffers from the disease Y.

DiseK contains 24,413 disease-related knowledge triplets in total. Assessing LLMs using all of these triplets would lead to a large computational cost. To balance the scale of evaluation and the computation efficiency, we select a single knowledge triplet (A, R, D) and a negative knowledge triplet  $(A', \neg R, D)$  for each pair of (R, D), resulting in a total of 6,334 knowledge triplets. Each positive triplet can be directly expressed by the statement "A is a R of the disease D".

Following the principle of truth-preserving evaluation, we choose four pieces of commonsense knowledge for deductive reasoning: (1) **Trivial** reasoning: the generated statement is exactly the same as the seed statement; (2) **Reverse** reasoning: we reverse the original expression, resulting in "The Rof D include A"; (3) **Specialization**: we combine the triplet with a specific case, such as "If a patient only suffers from disease Y, he/she (has the symptom of P(Y)/has lesions in P(Y)/P(Y) can be used to treat the patient)."; (4) **Contrapositive**: we conduct contrapositive transformation based on the third to derive this piece of commonsense. We also design similar deductive rules for negative triplets.

All of the commonsense knowledge used in our framework are listed in Table 3. Based on these pieces of commonsense knowledge and double negation, we generate a total of 8 statements for each triplet. Finally, the generated TPDiseK contains a total of 50,672 statements. We list all the utilized statement templates in Table 4.

885

887

888

889

890

891

892

893

894

895

896

897

898

899

900

901

902

903

904

905

906

907

908

909

910

911

912

913

914

915

916

917

877

Knowledge Aspect Statement ID State		Statement Template
	$S_1$	A is a common symptom of B.
	$S_2$	The common symptoms of B include A.
	C	If a patient only suffers from B, then he/she is likely to have
	$\mathfrak{z}_3$	symptoms of A.
	C	If a patient does not have the symptoms of A, then it is unlikely
C	$5_4$	that he/she only suffer from B.
Symptom	$S_5$	A is not a common symptom of B.
	$S_6$	The common symptoms of B do not include A.
	C	If a patient only suffers from B, then he/she is unlikely to have the
	$\mathfrak{S}_7$	symptoms of A.
	C	If a patient have the symptoms of A, then it is unlikely that he/she
	$\mathfrak{28}$	only suffer from B.
	$S_1$	A' is the affected site of B.
	$S_2$	The affected sites of B include A.
	C	If a patient only suffers from B, then he/she may have lesions in
	$\mathfrak{d}_3$	А.
	C	If a patient does not have lesions in A, then it is unlikely that
Affected Site	$\mathfrak{2}_4$	he/she only suffers from B.
Affected Site	$S_5$	A is not the affected site of B.
	$S_6$	The affected sites of B do not include A.
	<b>C</b> _	If a patient only suffers from B, then he/she is unlikely to have
	57	lessions in A.
	S.	If a patient have lesions in A, then it is unlikely that he/she only
	58	suffers from B.
	$\mathrm{S}_1$	A is a common therapeutic drug for B.
	$\mathrm{S}_2$	The common therapeutic drugs used to treat B include A.
	S	If a patient only suffers from B, then A can be used to treat his/her
	53	condition.
	S.	If A cannot be used to treat a patient's condition, then it is unlikely
Therapeutic Drug	64	that he/she only suffers from B.
Therapeutic Drug	$S_5$	A is not a common therapeutic drug for B.
	$\mathrm{S}_6$	The common therapeutic drugs used to treat B do not include A.
	$S_7$	If a patient only suffers from B, then it is unlikely that A can be
		used to treat his/her condition.
	$S_8$	If A can be used to treat a patient's condition, then it is unlikely
		that he/she only suffers from B.
	$S_1$	A is a common surgical procedure for B.
	$\mathrm{S}_2$	The common surgical procedures used to treat B include A.
Surgical Procedure	$\mathbf{S}_{2}$	If a patient only suffers from B, then A can be used to treat his/her
	~ 5	condition.
	$\mathrm{S}_4$	If A cannot be used to treat a patient's condition, then it is unlikely
		that he/she only suffers from B.
	$\mathbf{S}_{5}$	A is not a common surgical procedure for B.
	$\mathrm{S}_6$	The common surgical procedures used to treat B do not include A.
	$S_7$	If a patient only suffers from B, then it is unlikely that A can be
		used to treat his/her condition.
	$S_8$	If A can be used to treat a patient's condition, then it is unlikely
		that he/she only suffers from B.

Table 4: The statement templates we use in our evaluation framework. A: an entity from the specified knowledge aspect that has/does not have the relation with the disease. B: the name of the given disease.

Categories	Keywords
True	True, Entailed, Correct, Yes
False	False, Contradicted, Wrong, No

Table 5: The keywords we utilize to extract answers from LLMs' responses.

# C More Details of Evaluation Setting

918

941

943

945

In our implementation, we form the statement ver-919 920 ification question based on the template: "[Statement], is the statement above true or false? Please 921 922 answer True or False." For the five-shot setting, we randomly choose another five diseases, and follow the similar method applied in constructing 924 TPDiseK to form demonstrative examples. It is 925 worth noting that we always leverage example state-926 ments in the same format of the test statement to 927 achieve the best performance. For the zero-shot set-928 ting, we directly examine LLMs with the generated verification question. Complex prompting strategies such as chain-of-thought are not applied in 931 our study, as the evaluation statements are crafted 932 to be straightforward and easily understandable, allowing for verification without the need for complex logical reasoning. In the inference process, 935 we use greedy search for most of LLMs. However, some commercial LLMs (e.g, GPT-3.5-turbo) do 937 not support greedy search, and we use their default 939 generation setting to make a relative fair comparison across LLMs. 940

> We recognize the answer from models' response based on keyword recognition since the words/phrases used to express True and False are limited. We listed all of the keywords we applied to recognize answers in Table 5.

**D** Complementary Experiments

# D.1 Performance of all LLMs across Statement Types

We provide the detailed performance of all LLMs across different types of triplets, polarities, and ex-950 pression forms in Figure 8, 9, and 10, respectively. 951 The experimental results support the conclusions 952 made in our paper: the evaluated LLMs generally 954 perform worse on the negative and contrapositive types statements. Smaller LLMs perform signifi-955 cantly worse than larger LLMs on positive triplets, positive statements, and the first three expression forms. 958

#### **D.2** Zero-shot Performance of LLMs

As introduced in our paper, we also study the 960 zero-shot performance of LLMs on the proposed 961 TPDiseK dataset. The experimental results are dis-962 played in Figure 11, 12, 13, and 14, respectively. 963 The experimental results show that LLMs gener-964 ally achieve lower performance under the zero-shot 965 setting (< 10%), and LLMs' performance under 966 the zero-shot setting declines faster than that under 967 the five-shot setting. Some 70B LLMs (LLama2-968 70B and Meditron-70B) even achieve performance 969 close to random guessing, indicating that they have 970 a poor medical knowledge manipulation perfor-971 mance under the zero-shot setting. Nevertheless, 972 the results under the zero-shot setting exhibit the 973 similar trend compared to the five-shot setting, 974 demonstrating the correctness of our conclusions 975 based on the five-shot setting. 976



Figure 8: Average accuracy on statements generated from different types of knowledge triplets.



Figure 9: Average accuracy on statements with varied expression polarities. positive statements:  $S_1$  to  $S_4$ ; negative statements:  $S_5$  to  $S_8$ .



Figure 10: Average accuracy on statements with varied expression forms. Each expression form is determined by a specific predicate  $Q_i$ .



Figure 11: **Zero-shot** performance of LLMs evaluated by increasing number of truth-preserving statements. Left: joint accuracy; Right: average accuracy. Dotted lines: LLMs  $\leq$ 13B; Solid lines: LLMs >13B.



Figure 12: Zero-shot average accuracy on statements generated from different types of knowledge triplets.



Figure 13: Zero-shot average accuracy on statements with varied expression polarities. positive statements:  $S_1$  to  $S_4$ ; negative statements:  $S_5$  to  $S_8$ .



Figure 14: **Zero-shot** average accuracy on statements with varied expression forms. Each expression form is determined by a specific predicate  $Q_i$ .