## LLMs as Data Annotators: How Close Are We to Human Performance

**Anonymous ACL submission** 

#### Abstract

002

011

016

017

022

In NLP, fine-tuning LLMs is effective for various applications but requires high-quality annotated data. However, manual annotation of data is labor-intensive, time-consuming, and costly. Therefore, LLMs are increasingly used to automate the process, often employing in-context learning (ICL) in which some examples related to the task are given in the prompt for better performance. However, manually selecting context examples can lead to inefficiencies and suboptimal model performance. This paper presents comprehensive experiments comparing several LLMs, considering different embedding models, across various datasets for the Named Entity Recognition (NER) task. The evaluation encompasses models with approximately 7B and 70B parameters, including both proprietary and non-proprietary models. Furthermore, leveraging the success of Retrieval-Augmented Generation (RAG), it also considers a method that addresses the limitations of ICL by automatically retrieving contextual examples, thereby enhancing performance. The results highlight the importance of selecting the appropriate LLM and embedding model, understanding the trade-offs between LLM sizes and desired performance, and the necessity to direct research efforts towards more challenging datasets. The code, submitted as Supplementary Material, will be made publicly available after acceptance.

## 1 Introduction

Data annotation plays a crucial role in training machine learning (ML) models, especially in the era of Natural Language Processing (NLP). In NLP, data annotation typically involves annotating text data with relevant information, such as named entities, parts of speech, sentiment, intent, text classification, etc. The data annotation carries even more significance for fine-grained NLP tasks like token classification, where each token of a sentence has to be tagged with a gold label. In specialized domains such as Human Resource Management (HRM) or medical, organizations often possess large datasets that can be leveraged to enhance decision-making and operational efficiency through the use of LLM-based NLP approaches (Urlana et al., 2024). However, for these organizations to fully harness the power of LLMs through finetuning, they need high-quality annotated datasets. Traditional data annotation is a labor-intensive and costly process, especially when applied to large corpora. For example, in the case of HRM, annotating a dataset of 10,000 resumes for an information extraction task can be prohibitively time-consuming and requires significant human effort (Feng et al., 2021).

043

045

047

049

051

054

055

057

059

060

061

062

063

064

065

068

069

070

071

072

073

074

075

076

077

081

082

Nowadays, pre-trained LLMs (Devlin et al., 2019; Liu et al., 2019) can be cost-effectively fine-tuned on downstream tasks. These finetuned models are frequently used in scenarios where continuous LLM usage for inference is too expensive, such as when using API provided by propriety services (OpenAI, 2023; Team, 2024a), or when there is the need for tailored models to meet strict performance standards while maintaining the privacy of sensitive information, such as in specialized fields (Strohmeier, 2022; Karabacak and Margetis, 2023). With the advent of advanced LLMs such as GPT-4 (OpenAI, 2023), Qwen (Team, 2024b), and Llama (Touvron et al., 2023), researchers and practitioners are increasingly leveraging these models to enhance the data annotation process (Tan et al., 2024). Pre-trained on massive corpora, LLMs offer unprecedented capabilities for automating and streamlining annotation, improving scalability, and reducing costs (Wang et al., 2021).

Recent studies have demonstrated that LLMs (Wang et al., 2023; Naraki et al., 2024) can achieve performance comparable to human level in

084

- 0
- 0

0

100

117

118

119

121

122

123

124

125

126

127

128

129

data annotation for Named Entity Recognition Task (NER). However, most of these evaluations are conducted on widely used benchmark datasets such as CoNLL-2003 (Tjong Kim Sang and De Meulder, 2003) and WNUT-17 (Derczynski et al., 2017). For instance, between 2023 and 2025, the CoNLL-2003 dataset has been utilized in 191 studies, while WNUT-17 has been considered in 45. In contrast, more complex datasets like SKILLSPAN (Zhang et al., 2022a) and GUM (Zeldes, 2017) have been used significantly less frequently, in only 9 and 4 studies<sup>1</sup>, respectively. The results presented in this paper suggest that to gain a more comprehensive understanding of model performance regarding data annotation via LLMs in NER task, it is crucial to extend evaluations to more challenging datasets, which better reflect the complexities of real-world applications.

From a technical perspective, in the recent 101 literature, prompting (He et al., 2024a) and in-102 context learning (ICL) (Dong et al., 2024) are common approaches to leverage the LLMs for data annotation (Tan et al., 2024). ICL, which 105 is a technique where some solved examples of 106 the task are given within the prompt for better7 107 performance, is generally proven to be more 108 effective. However, selecting the right and relevant 109 examples to use as context for LLMs continues 110 to be a challenging task (Zhang et al., 2022b). 111 Manually choosing examples for each query 112 creates labor overhead, and more significantly, the 113 use of incorrect context examples may lead LLMs 114 to produce hallucinations (Yao et al., 2024) or 115 inaccurate outputs. 116

To address the above mentioned challenges, this paper presents the following contributions:

1. Comprehensive Evaluation of LLMs and Embeddings. It provides a comprehensive assessment of LLMs for data annotation in NER tasks, examining two distinct embedding models, as well as different techniques such as ICL and RAG, while utilizing datasets of varying complexity. It compares five models including proprietary models, such as gpt-40-mini, and open-source alternatives with approximately 7B and 70B parameters scale. 2. Trade-off Between LLM Sizes and Performances. The trade-off between LLM sizes and performance is demonstrated, which is further verified by the statistical tests. In fact, with the appropriate LLM and embedding models, there are no statistically significant differences in results between certain 7B and 70B models.

130

131

132

133

134

135

136

137

138

139

140

141

142

143

144

145

146

147

148

149

150

151

152

153

154

155

156

157

159

160

161

162

163

164

165

166

167

169

170

171

172

173

174

175

176

178

3. A RAG-Based Annotation Approach. To improve annotation quality and address the limitations of manual context selection in ICL, this paper considers a RAG based approach (Lewis et al., 2020). Instead of manually crafting in-context examples, the proposed method retrieves the most relevant samples based on similarity scores, enabling LLMs to generate more accurate annotations.

# 2 Related Work

In the recent past, there have been efforts by researchers to leverage the LLMs for data annotation (Tan et al., 2024). Wang et al. (2021) introduced the use of GPT-3 (Brown et al., 2020) for data annotation. The authors evaluated the quality of data generated by the GPT-3 against the human-labeled data. For each sentence to be annotated by the model, they construct a prompt consisting of several human-labeled examples along with the target sentence. They evaluate the performance in n-shot settings. Also, the authors report the performance of text classification and data generation tasks. Likewise, He et al. (2024a) leveraged the use of GPT-3.5 based models to annotate data. In comparison to the previous approach presented by Wang et al. (2021), the authors introduced the concept of chain-of-thought (CoT) (Wei et al., 2023) reasoning to annotate data. The authors simulate the human reasoning process to induce GPT-3.5 to motivate the annotated examples. They present the task description, specific examples, and the corresponding gold labels to GPT-3.5, and then ask the model to explain whether/why the given label is appropriate for that example. This enables the model to explain its choice of a specific label for the target sentence. Then, the authors construct the few-shot CoT prompts using the explanations generated by the model for data annotation.

To leverage the GPT model for the Named Entity Recognition (NER) task, Wang et al.

<sup>&</sup>lt;sup>1</sup>The statistics regarding the datasets usage is collected from https://paperswithcode.com/dataset/

250

251

252

253

254

255

256

257

258

259

260

261

262

263

264

265

266

267

268

269

270

271

272

273

274

275

276

277

229

179(2023) proposed a GPT-NER model. The main180contribution introduced by the authors is to181transform the NER into a text-generation task.182The authors used prompt engineering, where183prompts consist of three parts: (i) task description;184(ii) few-shot examples; and (iii) input sentence.185To choose few-shot context examples, they used186two different strategies: (i) random retrieval; and187(ii) k-NN based retrieval from training data.

In this work, the authors propose a retrievalbased approach for selecting context examples. 189 Specifically, for each training instance, the method 190 iterates through all tokens in a sentence to identify 191 the k-nearest neighbor (k-NN) tokens. The top 192 k retrieved tokens are then selected, and their 193 corresponding sentences are used as context. The 194 context examples are retrieved from the entire 195 training dataset. Furthermore, for sentences containing multiple entities, the algorithm runs 197 multiple times to ensure the extraction of all entities 198 within the sentence. 199

Following the work of Wang et al. (2021) and Wang et al. (2023), Naraki et al. (2024) also proposed a 201 LLMs based annotation for NER task. The authors 202 used the LLMs to clean noise and inconsistencies in the NER dataset, and then they merged the cleaned NER dataset with the original dataset to generate a more robust and diverse set of annotations. It is 206 worth mentioning that, in merging the annotations from LLM with human labels, preference is given to human-annotated examples compared to the LLM annotations. In addition, Bogdanov et al. 210 (2024) used the LLMs to create a general dataset 211 for NER tasks with a broad range of entity types. 212 The authors demonstrate a procedure that consists 213 of annotating raw data with an LLM to train a 214 task-specific foundation model for NER. Goel et al. 215 (2023) uses the same concept of data annotation 216 using LLMs, however, they do a case study on a 217 medical domain where they leverage the LLMs 218 for accelerating the annotation process along with 219 human input.

The research discussed above highlights the strong interest in using LLMs for dataset annotation, with most approaches relying on ICL. However, systematic evaluation on complex datasets remains limited, and selecting appropriate context examples for ICL is still a challenge. This study provides a comprehensive evaluation of LLMs for NER data annotation.

## 3 Methodology

## 3.1 Problem Definition

Given a dataset  $\mathcal{D} = \{S_i\}_{i=1}^n$ , where  $S_i$  represents the *i*-th sentence, with training, validation and test split given as  $\mathcal{D}_{train}$ ,  $\mathcal{D}_{valid}$  and  $\mathcal{D}_{test}$ . We divide  $\mathcal{D}_{train}$  into two disjoint subsets:  $\mathcal{X}$  (we call as sample space), from which we sample context examples, and  $\mathcal{T}$ , which will be annotated by the LLM. Formally, let  $\mathcal{X} \subset \mathcal{D}_{train}$  be a subset of size x, where x < n, and  $\mathcal{T} = \mathcal{D}_{train} \setminus \mathcal{X}$  be the remaining subset containing t sentences, where t = n - x. From  $\mathcal{X}$ , we select m examples, where m < x, to form the context set  $\mathcal{M}$ . The LLM uses all the m examples in  $\mathcal{M}$  as input context to annotate the t sentences in  $\mathcal{T}$ .

The NER task can be defined as the problem of learning an approximation function  $\tilde{f}_{\theta}$  that closely matches the real function  $f : S_{\mathcal{V}} \times \mathcal{V} \to \mathcal{C}$ , where  $S_{\mathcal{V}}$  represents the set of all the possible sentences composed only by words w in the vocabulary  $\mathcal{V}$ , and  $\mathcal{C}$  represents the set of possible entity categories. The real function f given: (*i*) a sentence  $S_i \in S_{\mathcal{V}}$ , and (*ii*) a word  $w \in \mathcal{V}$ , assigns w to its corresponding category  $c \in \mathcal{C}$ .

## 3.2 Data Annotation via LLMs

The methodology adopted in the proposed RAG approach is shown in Figure 1. This section discusses the steps followed in the proposed study. Section 3.2.1 explains the prompt template formation, while Section 3.2.2 presents the baseline approach, followed by ICL method in Section 3.2.3. Section 3.2.4 presents the proposed RAG technique, whereas the importance of structured outputs for NER task is discussed in Section 3.2.5.

## 3.2.1 Prompt Formation

In NLP, crafting an effective prompt for LLMs is a crucial task, as an ill-formed prompt could lead to poor performance. Different LLMs, whether opensource or proprietary, tend to respond differently to variations in prompt (Errica et al., 2024). This work adopts a similar approach to prompt design presented in (He et al., 2024b; Wang et al., 2023), i.e. structuring our prompts around three key components, also visible in Figure 1: (*i*) *Task Description.* This component clearly defines the task the LLM is expected to perform; (*ii*) *Context.* This component provides task-related examples that help the LLM to better understand the problem, while also clarifying the expected input/output



Figure 1: Workflow of the proposed approach.  $\mathcal{D}_{train}$  denotes the training data,  $\mathcal{X}$  denotes the few human annotated examples, whereas  $\mathcal{T}$  denotes the training instances to be annotated by LLM. For each entry  $\mathcal{T}_i \in \mathcal{T}$ , we extract  $\mathcal{M}$  context examples from a vector store using a retriever module. Then, given an input sentence, the final prompt to LLM consists of the task description, the context examples in  $\mathcal{M}$ , and input sentence.

format; and *(iii) Input*. This final component presents the LLM with the specific examples to be annotated. The prompt structures adopted in the experiments are outlined in Appendix G, while several prompt examples are reported in Appendix H.

## 3.2.2 Zero-shot Data Annotation

279

290

291

295

298

304

307

In the zero-shot setting (refers to the baseline), the LLM receives only task descriptions and entity categories from the dataset. The task description explains the task, whereas entity categories provide information about the classes that the LLM has to use for annotation. Providing entity categories in the prompt allows the LLM to produce consistent output annotation as in the training set. For instance, in the CoNLL-2003 dataset, person and organization categories are labelled as PER and ORG respectively. Thus, the prompt to the LLM includes PER and ORG to annotate entities in the person and organization categories, respectively. However, in zero-shot data annotation, the lack of context examples hinders the model's understanding, often leading to suboptimal performance. Nonetheless, this setting allows to evaluate the general knowledge of LLM on a task.

#### 3.2.3 In-Context Learning

In ICL, the prompts given to LLMs are enhanced by including not only a task description and entity categories but also contextual examples. These examples aid the models in better understanding the task at hand. As detailed in Section 3.1,  $\mathcal{D}_{train}$ is is split into  $\mathcal{X}$  and  $\mathcal{T}$ . From  $\mathcal{X}$ , the selection of  $\mathcal{M}$  can be approached in two ways: either through manual cherry-picking or by random sampling. However, manually selecting  $\mathcal{M}$  can be both timeconsuming and subjective, which contradicts the rationale of the proposed study. Therefore, we opt to randomly sample  $\mathcal{M}$  from  $\mathcal{X}$ , although it does not guarantee whether the selected context examples  $\mathcal{M}$  are semantically close to the input text  $\mathcal{T}_i$ , which is a limitation of this approach. 308

309

310

311

312

313

314

315

316

317

318

319

320

#### 3.2.4 Retrieval-Based Approach

To overcome the limitations of the previously 321 mentioned approaches, this paper introduces a 322 retrieval-based method for automatically selecting 323 relevant context examples. As outlined in 324 Section 3.1, the proposed RAG-based approach 325 first generates embedding representations for all examples in  $\mathcal{X}$ , which are then stored in a vector 327 database (Douze et al., 2024) for subsequent retrieval, as illustrated in Figure 1. Subsequently, 329 for each sentence  $\mathcal{T}_i \in \mathcal{T}$ , its embedding representation is generated, and the most similar 331  $\mathcal{M}$  examples are retrieved from  $\mathcal{X}$  stored in the 332 vector database.  $\mathcal{M}$  is then used as context for the 333 LLM to provide the most relevant examples for annotating the input text  $T_i$ . 335

341

342

347

353

361

363

371

374

375

377

381

## **3.2.5** Structured Output from LLMs

For a label-sensitive task like NER, getting a structured output from a LLM is a crucial step. In the NER task, as defined in Section 3.1, each token in a sentence is tagged with a corresponding label. Hence, preserving the tokenlabel correspondence in the output is necessary for the LLMs. The most recent LLMs are based on a decoder architecture that, while being suitable for sequence-to-sequence tasks, encounters challenges when tackling the NER task due to the potential misalignment between tokens and labels (UI Haq et al., 2024). In fact, recent studies on NER (Li et al., 2024; Liu et al., 2024; Wang et al., 2023) have shown that the decoder architecture presents structural inconsistencies in the output. Recently, OpenAI (OpenAI, 2023) released a feature for the latest GPT-4 based models which guarantees to follow the structured output format<sup>2</sup>. To solve the token-label misalignment problem, in this study, we leverage the latest feature of StructuredOutput released by OpenAI. However, it is important to note that despite the inclusion of such features in the latest LLMs, including Qwen (Team, 2024b) and Llama (Touvron et al., 2023) based models, they still exhibit inconsistencies in their output, unlike the gpt-4o-mini-2024-07-18.

## 4 Experimental Setup

### 4.1 Datasets

In this study, to evaluate the performance of the proposed methodology and assess the capabilities of LLMs, four datasets are considered, with their statistics summarized in Table 1 of Appendix A. Each dataset presents unique challenges for LLMs in performing NER tasks, allowing this study to comprehensively analyze the ability of LLMs to handle diverse entity types, from well-structured entities to complex, ambiguous, and domainspecific annotations.

**CoNLL-2003** The CoNLL-2003 (Tjong Kim Sang and De Meulder, 2003) dataset consists of four general entity types. Entities in this dataset typically follow structured patterns, making them relatively easier for LLMs to identify and classify.

WNUT-17 The WNUT-17 (Derczynski et al., 2017) dataset contains six categories of rare entities. This dataset is particularly challenging due to its

noisy text, sparse entity occurrences, and limited labeled examples per category. Improving recall on this dataset remains a significant challenge for LLMs. 383

384

385

388

389

390

391

393

394

395

397

398

399

400

401

402

403

404

405

406

407

408

409

410

411

412

413

414

415

416

417

418

419

420

421

422

423

424

425

426

427

428

**GUM** The GUM (Zeldes, 2017) dataset is a richly annotated corpus designed for multiple NLP tasks, including NER. It captures linguistic phenomena across various domains and genres, making it a valuable resource for evaluating model performance. The dataset includes eleven distinct named entity types. Compared to CoNLL-2003 and WNUT-17, GUM presents a higher level of complexity by incorporating a diverse set of entity types spanning multiple domains.

**SKILLSPAN** The SKILLSPAN (Zhang et al., 2022a) dataset is composed of a single entity type. Unlike traditional entities, soft skills do not follow a fixed syntactic or semantic structure, making them inherently ambiguous. These entities can range from single tokens to multitoken expressions, increasing the complexity of annotation and information extraction tasks for LLMs.

#### 4.2 Approaches Under Study

In the empirical assessment of the datasets annotated by LLMs, the zero-shot data annotation approach is chosen as the baseline since it provides no context about the task to the LLM. This zeroshot setting allows the evaluation of the LLM's general knowledge of the task. Moreover, ICL and RAG-based approaches, detailed in Section 3.2.3 and Section 3.2.4 respectively, are considered. For both, experiments are conducted with three different numbers of context examples: (*i*) 25, (*ii*) 50, and (*iii*) 75. Experiments are conducted on a 30% sample of the training set  $\mathcal{D}_{\text{train}}$ , while the ablation study in Appendix D examines the effects of 10% and 20% sample sizes.

This paper considers five different LLMs<sup>3</sup>: (*i*) gpt-4o-mini-2024-07-18, (*ii*) Qwen2.5-72B-Instruct, (*iii*) Llama3.5-70B-Instruct, (*iv*) Qwen2.5-7B-Instruct, and (*v*) Llama3.1-8B-Instruct, and two embeddings models: (*i*) the text-embedding-3-large model<sup>4</sup>, and (*ii*) the sentence transformer all-MiniLM-L6- $v^2$  model (Reimers and Gurevych, 2019).

<sup>&</sup>lt;sup>2</sup>https://openai.com/index/

 $introducing\-structured\-outputs\-in\-the\-api$ 

<sup>&</sup>lt;sup>3</sup>The models are referred to by their base names, such as Qwen2.5-72B for Qwen2.5-72B-Instruct, and so on.

<sup>&</sup>lt;sup>4</sup>https://platform.openai.com/docs/guides/embeddings

Throughout the remainder 429 of the paper, text-embedding-3-large will be referred 430 to as OpenAI, and sentence transformer 431 all-MiniLM-L6-v2 will be referred to as ST. 432 Implementation details of results are reported 433 Appendix **B**. 434

## 4.3 NER Evaluation Process

435

436

437

438

439

440

441

442

443

444

445

446

447

448

449

450

451

452

453

454

455

456

457

458

459

460

461

462

463

464

To assess the quality of annotations generated by LLMs, the RoBERTa model (Liu et al., 2019) is fine-tuned on LLM-annotated datasets, leveraging its proven effectiveness in NER tasks (Zhou et al., 2022; Zhang et al., 2022a). Initially, an LLM is employed to automatically annotate sentences in  $\mathcal{T} \subset \mathcal{D}_{train}$ , using strategies from Section 3.2. This process generates annotations for  $\mathcal{T}$ , resulting in a new training set,  $\hat{\mathcal{T}}$ , with  $|\hat{\mathcal{T}}| = |\mathcal{T}|$ . This annotated set is then used to fine-tune the RoBERTa model (Liu et al., 2019). Model selection is performed on the validation set,  $\mathcal{D}_{valid}$ , and the final evaluation results are based on the test set,  $\mathcal{D}_{test}$ . To ensure robustness and mitigate the impact of random initialization, we average the results across five different seed values. The  $F_1$  score is used to assess the performances of the models.

## 5 Results and Analysis

This section presents the quantitative results of this study, as well as its analysis. Qualitative results are reported in Appendix F, while Appendix E reports the statistical tests to support the findings.

## 5.1 Quantitative Results

Figure 2 presents the overall results of the experiments, while the corresponding detailed outcomes are reported in Appendix C. Specifically, the heatmaps present the  $F_1$  scores obtained on the test set for different datasets, comparing several models and methods used in the proposed study.

The CoNLL-2003 dataset, which contains named 465 entities like persons, organizations, and locations, 466 is relatively well-structured, making it easier for 467 LLMs to generate high-quality annotations. The 468 gpt-4o-mini model with OpenAI embeddings 469 emerges as the top performer (also shows statistical 470 significance over other models as detailed in 471 472 Appendix E), achieving an  $F_1$  score of 89.72 with 75 context examples, which is just 2.7% below 473 human-level annotation. Among the  $\sim 70B$ 474 models, Qwen2.5-72B with OpenAI embeddings 475 performs comparably to gpt-4o-mini with an 476

 $F_1$  score of 89.34, while Llama3.5-70B with 477 ST embeddings lags slightly behind with an  $F_1$ 478 score of 87.33. At the  $\sim 7B$  scale, Qwen2.5-7B 479 with ST embeddings significantly outperforms 480 its counterparts, achieving an  $F_1$  score of 87.94, 481 while Llama3.1-8B with OpenAI embeddings 482 scores 84.91. This suggests that smaller models 483 can still perform competitively when paired with 484 appropriate embedding methods. Interestingly, 485 the heatmap reveals that context size plays 486 a crucial role-gpt-4o-mini and Qwen2.5-70B 487 benefit significantly from larger context sizes of 488 75 examples, while Llama3.5-70B performs best 489 at a slightly lower context size. This suggests that 490 different models have varying levels of context 491 saturation, where additional examples may not 492 always improve performance linearly. 493

494

495

496

497

498

499

500

501

502

503

504

505

506

507

508

509

510

511

512

513

514

515

516

517

518

519

520

521

522

523

524

525

527

The WNUT-17 dataset, which focuses on lowfrequency and emerging entities, presents a significant challenge due to limited training samples for each entity. However, Qwen2.5-70B with OpenAI embeddings achieves the highest  $F_1$  score of 53.72, slightly outperforming gpt-40-mini, which attains an  $F_1$  score of 53.43. The Llama3.5-70B model exhibits inconsistent performance, scoring 51.18 with ICL at 75 context examples, suggesting that it struggles to generalize well for rare entity detection. At the  $\sim 7B$ scale, Qwen2.5-7B with ST embeddings achieves an  $F_1$  score of 49.48, significantly outperforming Llama3.1-8B, which scores 44.42. This highlights that ST embeddings provide a crucial advantage for smaller models. Compared to human-level annotation, which achieves an  $F_1$  score of 54.93, the best-performing LLM reduces the gap to just 1.21%, which is the smallest performance gap between human and LLM annotation across all datasets used in the experiments. This suggests that RAG-based annotation is highly effective in adapting to rare entity recognition, particularly when combined with larger models and strong embeddings.

The GUM dataset presents a unique challenge due to its diverse entity types, requiring models to generalize across various linguistic structures. Qwen2.5-70B with ST embeddings achieves the best  $F_1$  score of 55.11, significantly surpassing gpt-4o-mini, which attains an  $F_1$  score of 52.28, and Llama3.5-70B with OpenAI embeddings, which achieves an  $F_1$  score of 48.33. At the ~ 7B scale, Qwen2.5-7B with OpenAI embeddings



Figure 2: Heatmaps of the  $F_1$  scores across four datasets. The color scale represents performance, with red indicating higher scores reaching human-level, and blue indicating lower scores starting from the lowest performing model

achieves an  $F_1$  score of 44.48, outperforming L1ama3.1-8B, which scores 43.91. However, both models show a notable performance drop compared to their larger counterparts, suggesting that smaller models struggle with datasets with diverse entities. The 3.15% gap between the best-performing LLM and human-level annotation highlights that GUM remains a challenging dataset for LLMs. The heatmap further suggests that model performance fluctuates significantly depending on context size and embedding choice.

528

529

532

533

534

536

The SKILLSPAN dataset is the most difficult 539 540 among those evaluated, as it requires understanding nuanced skill mentions across various job contexts. 541 gpt-4o-mini with OpenAI embeddings performs 542 the best, achieving an  $F_1$  score of 34.06 with 543 75 context examples, but this is still far from 544 human-level annotation. At the  $\sim 70B$  scale, Qwen2.5-70B with ST embeddings achieves an 546  $F_1$  score of 32.35 with 50 context examples, outperforming Llama3.5-70B, which achieves an Among  $\sim 7B$  models,  $F_1$  score of 27.55. 550 Qwen2.5-7B with OpenAI embeddings achieves an  $F_1$  score of 29.67, significantly surpassing 551 Llama3.1-8B, which scores 22.88. This suggests that embedding choice plays a crucial role in skill extraction tasks. Notably, the gap between 554

human annotation and the best-performing LLM is much larger in this dataset compared to others, indicating that LLMs struggle with skillbased entity recognition. This could be due to the complexity of contextual skill interpretation, requiring deeper domain knowledge and better understanding capabilities. 555

556

557

559

560

561

563

564

565

566

567

568

570

571

572

573

574

575

576

577

578

579

## 5.2 Different Sample Space Choices

This section examines the impact of sample space choices, denoted as  $\mathcal{X}$  in Section 3.1, using the proposed RAG-based approach as overall it performs better than ICL. The experiments are conducted on the SKILLSPAN dataset with the gpt-4o-mini model and OpenAI embeddings. As shown in Figure 3, for smaller dataset splits, the RAG-based approach exhibits greater variability, similar to the behavior seen with ICL. This suggests that as the sample space for selecting context examples decreases, the performance of the RAGbased approach converges more closely with that of ICL. More detailed results are reported in Appendix D.

## 6 Discussion

**Performance of LLMs** The performance of different LLMs in our study reveals interesting



Figure 3:  $F_1$  scores for different context sizes (25, 50, and 75) and sample spaces (10% and 20%) for the RAG and ICL approach on the SKILLSPAN dataset, using the gpt-4o-mini model. The plot indicates that with a smaller sample size, the RAG approach performs comparably to ICL.

insights. Across all datasets, RAG-based approaches improve annotation quality, with gpt-4o-mini and OpenAI embeddings achieving the best results. In contrast, ICL struggles in datasets with sparse or ambiguous entities, particularly SKILLSPAN. While all models perform well on CoNLL-2003, performance declines as entity structures become more complex, such as in GUM and SKILLSPAN.

580

581

582

583

588

Embeddings The Effect of choice of embeddings for retrieval of context for LLMs plays a crucial role in annotation quality in 591 retrieval-based methods. OpenAI embeddings lead 592 to better  $F_1$  scores compared to smaller-scale ST 593 embeddings especially for gpt-4o-mini model. 594 This effect is particularly evident in WNUT-17 and 595 GUM, where entity distributions are more diverse, 596 and high-quality embeddings improve retrieval 597 effectiveness. In contrast, SKILLSPAN remains 598 challenging across all embedding strategies, suggesting that current embedding techniques 600 struggle with soft skill representation due to the abstract nature of the entities.

Effect of Model Size Larger models generally perform better, but retrieval quality is equally Qwen2.5-7B slightly outperforms critical. Llama3.1-8B 606 and performs comparably to Llama3.5-70B with proper embeddings, 607 indicating that architecture and training data impact annotation beyond parameter count. Statistical tests in Appendix E support this finding. 610

Effect of Dataset Complexity Breaking down
results per dataset, CoNLL-2003 shows minimal
variance across methods, as structured entities
are well-represented in training data. WNUT-17
benefits the most from retrieval-based methods,

as rare entities require additional context for accurate recognition. GUM's diverse entity types pose a challenge for ICL, but RAGbased methods significantly improve performance. Finally, SKILLSPAN remains the most difficult dataset, with lower performance across all methods, underscoring the limitations of LLMs and embeddings in capturing the semantics of soft skills. 616

617

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

## 7 Conclusions and Future Works

This study systematically evaluates the effectiveness of LLMs for data annotation across diverse datasets-CoNLL-2003, four WNUT-17, GUM, and SKILLSPAN of varying It compares RAG in different complexity. embedding strategies, ICL, and a baseline approach. The results demonstrate that RAG-based methods consistently outperform both ICL and the baseline across all datasets, significantly reducing the performance gap with human-level annotation.

A key finding is that dataset complexity plays a crucial role in model performance. For structured datasets like CoNLL-2003, LLMs perform exceptionally well, with models such as gpt-4o-mini and Qwen2.5-72B achieving results within 3% of human-level annotation. Conversely, performance deteriorates as dataset complexity increases. The SKILLSPAN dataset, which requires nuanced skill recognition, presents the greatest challenge, with LLMs struggling to capture implicit skill mentions.

Our analysis also highlights the importance of context size and embedding choice in retrieval-augmented annotation. We observe that larger models such as Qwen2.5-72B and gpt-4o-mini benefit from larger context sizes, while smaller models like Qwen2.5-7B can still perform competitively when paired with highquality sentence embeddings. However, models exhibit context saturation effects, where additional examples do not always lead to linear performance improvements.

Future works will focus on enhancing the performance of LLMs for complex datasets, particularly in specialized domains. In addition, future works will expand the study to more LLMs and to different NLP tasks.

765

766

767

768

713

714

## 663 Limitations

- In this study, we evaluate LLMs for data annotation tasks and introduce a RAG-based approach with different embedding models to enhance performance on NER datasets. However, our work has several limitations that highlight areas for future research.
- First, our experiments focus solely on NER
  tasks. While this provides a solid foundation
  for evaluation, extending the analysis to other
  NLP tasks, such as text classification or question
  answering, would offer a more comprehensive
  understanding of the proposed methodology's
  applicability and generalizability.
- Second, for the proof of concept, we employ a naïve RAG approach for context selection. Future work could explore more sophisticated retrieval techniques, such as adaptive retrieval strategies, re-ranking mechanisms, or hybrid approaches combining dense and sparse retrieval, to further optimize performance.
- Third, our study does not explicitly examine the
  biases introduced by LLMs in the data annotation
  process. Given the growing concerns about fairness
  and model biases, a deeper investigation into how
  LLMs influence annotation patterns, especially
  in diverse and underrepresented datasets, could
  provide valuable insights.

## References

- Sergei Bogdanov, Alexandre Constantin, Timothée
   Bernard, Benoit Crabbé, and Etienne Bernard. 2024.
   Nuner: Entity recognition encoder pre-training via Ilmannotated data. *Preprint*, arXiv:2402.15343.
- 696Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie697Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind698Neelakantan, Pranav Shyam, Girish Sastry, Amanda699Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen700Krueger, Tom Henighan, Rewon Child, Aditya701Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens702Winter, Christopher Hesse, Mark Chen, Eric Sigler,703Mateusz Litwin, Scott Gray, Benjamin Chess, Jack704Clark, Christopher Berner, Sam McCandlish, Alec705Radford, Ilya Sutskever, and Dario Amodei. 2020.706Language models are few-shot learners. Preprint,707arXiv:2005.14165.
- William Jay Conover. 1999. *Practical Nonparametric*Statistics, volume 350. John Wiley & Sons.
- Leon Derczynski, Eric Nichols, Marieke van Erp, and
  Nut Limsopatham. 2017. Results of the WNUT2017
  shared task on novel and emerging entity recognition.

In Proceedings of the 3rd Workshop on Noisy Usergenerated Text, pages 140–147, Copenhagen, Denmark. Association for Computational Linguistics.

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. Bert: Pre-training of deep bidirectional transformers for language understanding. *ArXiv*, abs/1810.04805.

Qingxiu Dong, Lei Li, Damai Dai, Ce Zheng, Jingyuan Ma, Rui Li, Heming Xia, Jingjing Xu, Zhiyong Wu, Baobao Chang, Xu Sun, Lei Li, and Zhifang Sui. 2024. A survey on in-context learning. *Preprint*, arXiv:2301.00234.

Matthijs Douze, Alexandr Guzhva, Chengqi Deng, Jeff Johnson, Gergely Szilvasy, Pierre-Emmanuel Mazaré, Maria Lomeli, Lucas Hosseini, and Hervé Jégou. 2024. The faiss library.

Federico Errica, Giuseppe Siracusano, Davide Sanvito, and Roberto Bifulco. 2024. What did i do wrong? quantifying llms' sensitivity and consistency to prompt engineering. *Preprint*, arXiv:2406.12334.

Steven Y. Feng, Varun Gangal, Jason Wei, Sarath Chandar, Soroush Vosoughi, Teruko Mitamura, and Eduard Hovy. 2021. A survey of data augmentation approaches for NLP. In *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*, pages 968–988, Online. Association for Computational Linguistics.

Akshay Goel, Almog Gueta, Omry Gilon, Chang Liu, Sofia Erell, Lan Huong Nguyen, Xiaohong Hao, Bolous Jaber, Shashir Reddy, Rupesh Kartha, Jean Steiner, Itay Laish, and Amir Feder. 2023. Llms accelerate annotation for medical information extraction. *Preprint*, arXiv:2312.02296.

Xingwei He, Zhenghao Lin, Yeyun Gong, A-Long Jin, Hang Zhang, Chen Lin, Jian Jiao, Siu Ming Yiu, Nan Duan, and Weizhu Chen. 2024a. Annollm: Making large language models to be better crowdsourced annotators. *Preprint*, arXiv:2303.16854.

Xingwei He, Zhenghao Lin, Yeyun Gong, A-Long Jin, Hang Zhang, Chen Lin, Jian Jiao, Siu Ming Yiu, Nan Duan, and Weizhu Chen. 2024b. AnnoLLM: Making large language models to be better crowdsourced annotators. In Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 6: Industry Track), pages 165–190, Mexico City, Mexico. Association for Computational Linguistics.

Hugging Face. 2023. Transformers APIs. https: //huggingface.co/docs/transformers/index. Accessed: 2023-01-21.

Benoit Jacob, Skirmantas Kligys, Bo Chen, Menglong Zhu, Matthew Tang, Andrew Howard, Hartwig Adam, and Dmitry Kalenichenko. 2018. Quantization and training of neural networks for efficient integerarithmetic-only inference. In 2018 IEEE/CVF

823

- *Conference on Computer Vision and PatternRecognition*, pages 2704–2713.
- 771 Mert Karabacak and Konstantinos Margetis. 2023.
  772 Embracing large language models for medical applications: Opportunities and challenges. *Cureus*, 15(5):e39305.
- Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio
  Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich
  Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel,
  Sebastian Riedel, and Douwe Kiela. 2020. Retrievalaugmented generation for knowledge-intensive nlp
  tasks. In Advances in Neural Information Processing
  Systems, volume 33, pages 9459–9474. Curran
  Associates, Inc.
- Yinghao Li, Rampi Ramprasad, and Chao Zhang. 2024.
  A simple but effective approach to improve structured language model output for information extraction. *Preprint*, arXiv:2402.13364.
- Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du,
  Mandar Joshi, Danqi Chen, Omer Levy, Mike
  Lewis, Luke Zettlemoyer, and Veselin Stoyanov.
  2019. Roberta: A robustly optimized bert pretraining
  approach. ArXiv, abs/1907.11692.
- Yu Liu, Duantengchuan Li, Kaili Wang, Zhuoran Xiong,
  Fobo Shi, Jian Wang, Bing Li, and Bo Hang. 2024.
  Are llms good at structured outputs? a benchmark
  for evaluating structured output capabilities in llms. *Information Processing & Management*, 61(5):103809.
- Yuji Naraki, Ryosuke Yamaki, Yoshikazu Ikeda, Takafumi Horie, and Hiroki Naganuma. 2024.
  Augmenting ner datasets with llms: Towards automated and refined annotation. *Preprint*, arXiv:2404.01334.
- OpenAI. 2023. Gpt-4 technical report. *Preprint*, arXiv:2303.08774.

803

804

- D. Pereira, Anabela Afonso, and Fátima Medeiros. 2015. Overview of friedman's test and post-hoc analysis. *Communications in Statistics - Simulation and Computation*, 44:2636–2653.
- Nils Reimers and Iryna Gurevych. 2019. Sentencebert: Sentence embeddings using siamese bert-networks. *Preprint*, arXiv:1908.10084.
- 810 Stefan Strohmeier. 2022. Handbook of Research on
  811 Artificial Intelligence in Human Resource Management.
  812 Edward Elgar Publishing.
- Zhen Tan, Dawei Li, Alimohammad Beigi, Song Wang,
  Ruocheng Guo, Amrita Bhattacharjee, Bohan Jiang,
  Mansooreh Karami, Jundong Li, Lu Cheng, and Huan
  Liu. 2024. Large language models for data annotation:
  A survey. ArXiv, abs/2402.13446.
- 818 Gemini Team. 2024a. Gemini: A family of highly capable multimodal models. *Preprint*, arXiv:2312.11805.
- 821Qwen Team. 2024b. Qwen2.5: A party of foundation822models.

Maksim Terpilowski. 2019. scikit-posthocs: Pairwise multiple comparison tests in python. *The Journal of Open Source Software*, 4(36):1169.

Erik F. Tjong Kim Sang and Fien De Meulder. 2003. Introduction to the CoNLL-2003 shared task: Languageindependent named entity recognition. In *Proceedings* of the Seventh Conference on Natural Language Learning at HLT-NAACL 2003, pages 142–147.

Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurelien Rodriguez, Armand Joulin, Edouard Grave, and Guillaume Lample. 2023. Llama: Open and efficient foundation language models. *Preprint*, arXiv:2302.13971.

Muhammad Uzair Ul Haq, Paolo Frazzetto, Alessandro Sperduti, and Giovanni Da San Martino. 2024. Improving soft skill extraction via data augmentation and embedding manipulation. In *Proceedings of the 39th ACM/SIGAPP Symposium on Applied Computing*, SAC '24, page 987–996, New York, NY, USA. Association for Computing Machinery.

Ashok Urlana, Charaka Vinayak Kumar, Ajeet Kumar Singh, Bala Mallikarjunarao Garlapati, Srinivasa Rao Chalamala, and Rahul Mishra. 2024. Llms with industrial lens: Deciphering the challenges and prospects – a survey. *Preprint*, arXiv:2402.14558.

Shuhe Wang, Xiaofei Sun, Xiaoya Li, Rongbin Ouyang, Fei Wu, Tianwei Zhang, Jiwei Li, and Guoyin Wang. 2023. Gpt-ner: Named entity recognition via large language models. *Preprint*, arXiv:2304.10428.

Shuohang Wang, Yang Liu, Yichong Xu, Chenguang Zhu, and Michael Zeng. 2021. Want to reduce labeling cost? GPT-3 can help. In *Findings of the Association for Computational Linguistics: EMNLP 2021*, pages 4195–4205, Punta Cana, Dominican Republic. Association for Computational Linguistics.

Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed Chi, Quoc Le, and Denny Zhou. 2023. Chain-of-thought prompting elicits reasoning in large language models. *Preprint*, arXiv:2201.11903.

Jia-Yu Yao, Kun-Peng Ning, Zhen-Hui Liu, Mu-Nan Ning, Yu-Yang Liu, and Li Yuan. 2024. Llm lies: Hallucinations are not bugs, but features as adversarial examples. *Preprint*, arXiv:2310.01469.

Amir Zeldes. 2017. The GUM corpus: Creating multilayer resources in the classroom. *Language Resources and Evaluation*, 51(3):581–612.

Mike Zhang, Kristian Nørgaard Jensen, Sif Dam Sonniks, and Barbara Plank. 2022a. Skillspan: Hard and soft skill extraction from english job postings. In North American Chapter of the Association for Computational Linguistics.

Zhuosheng Zhang, Aston Zhang, Mu Li, and Alex Smola. 2022b. Automatic chain of thought prompting in large language models. *Preprint*, arXiv:2210.03493.

880	Ran Zhou, Xin Li, Ruidan He, Lidong Bing, Erik
881	Cambria, Luo Si, and Chunyan Miao. 2022. MELM:
882	Data augmentation with masked entity language
883	modeling for low-resource NER. In Proceedings
884	of the 60th Annual Meeting of the Association for
885	Computational Linguistics (Volume 1: Long Papers),
886	pages 2251-2262, Dublin, Ireland. Association for
887	Computational Linguistics.

## A Datasets Statistics

Table 1: Statistics of the datasets considered in this study. The average entity length refers to the average number of tokens for each entity.

Dataset		Sentences			Tokens	Avg. Entity Length		
2	Train	Validation	Test	Train	Validation	Test	· · · · · · · · · · · · · · · · · · ·	
CoNLL-2003	14041	3250	3453	203621	51362	46435	1.60	
WNUT-2017	3394	1008	1287	62730	15734	23394	1.73	
GUM	1435	615	805	29392	12688	17437	3.15	
SKILLSPAN	3074	1396	1522	92621	39923	42541	4.72	

Table 1 highlights the complexity of entity mentions across different datasets, as reflected in their average entity length. CoNLL-2003 and WNUT-2017 contain relatively short entities, with average lengths of 1.60 and 1.73 tokens, respectively, indicating that most entities are single-token mentions. In contrast, GUM exhibits greater complexity, with an average entity length of 3.15 tokens, suggesting the presence of multi-token entities. SKILLSPAN is the most complex dataset, with an average entity length of 4.72 tokens, implying more intricate entity structures that require advanced modeling techniques for accurate recognition.

- 6 Moreover, we discuss below the entity information for each dataset.
  - **CoNLL-2003** The CoNLL-2003 (Tjong Kim Sang and De Meulder, 2003) dataset consists of general entity types: (*i*) PERSON ; (*ii*) ORGANIZATION; (*iii*) LOCATION; and (*iv*) MISCELLANEOUS. Entities in this dataset typically follow structured patterns, making them relatively easier for LLMs to identify and classify.
- **WNUT-17** The WNUT-17 (Derczynski et al., 2017) dataset contains six categories of rare entities: (*i*) PERSON; (*ii*) CORPORATION; (*iii*) LOCATION; (*iv*) CREATIVE\_WORK; (*v*) GROUP; and (*vi*) PRODUCT. This dataset is particularly challenging due to its noisy text, sparse entity occurrences, and limited labeled examples per category.
- GUM The GUM (Zeldes, 2017) dataset is a richly annotated corpus designed for multiple NLP tasks,
  including NER. The dataset includes eleven distinct named entity types: (*i*) ABSTRACT; (*ii*) ANIMAL;
  (*iii*) EVENT; (*iv*) OBJECT; (*v*) ORGANIZATION; (*vi*) PERSON; (*vii*) PLACE; (*viii*) PLANT;
  (*ix*) QUANTITY; (*x*) SUBSTANCE; and (*xi*) TIME.
- SKILLSPAN The SKILLSPAN (Zhang et al., 2022a) dataset is composed of a single entity type,
   SOFTSKILLS, extracted from job descriptions. Unlike traditional entities, soft skills do not follow a fixed
   syntactic or semantic structure, making them inherently ambiguous.

## **B** Implementation Details

To perform experiments for data annotation with gpt-4o-mini, the model is accessed via the API service 913 provided by OpenAI. To ensure reproducible results, the temperature is set to 0 and a seed value of 42914 is used. Furthermore, the system fingerprint fp\_1bb46167f9 is reported as noted during API access. 915 For data annotation generation using Qwen (Team, 2024b) and Llama (Touvron et al., 2023) based 916 models, the HuggingFace (Hugging Face, 2023) implementation is utilized. The instructed fine-tuned 917 variants of the open-source models are employed in the proposed study. The models are used only for inference, with 4-bit quantization (Jacob et al., 2018). The experiments with billion scale models are 920 conducted on an A100 GPU with a seed value of 42. All experiments to fine-tune NER task are performed with the RoBERTa model, available via HuggingFace (Hugging Face, 2023), are conducted in a python 921 environment, on an RTX A5000 GPU. The experiments are performed using the following five seed 922 values: [23112, 13215, 6465, 42, 5634]. Moreover, the statistical significance tests are performed with the 923 help of scikit-posthocs (Terpilowski, 2019) library available in python. 924

889

900

901

902

903

904

# C Complete Results

Table 2: The  $F_1$ , precision and recall along with standard deviation are reported on the test set. The values are averaged over five different random initializations. #Ex. represents the number of context examples used. Baseline refers to the use of LLM with no context examples.

#Fv	Method		CoNLL20	003		WNUT-17 GUM			GUM			SKILLSPAN		
πιΔλ.	Methou	Р	R	$\mathbf{F_1}$	P	R	$\mathbf{F_1}$	P	R	$\mathbf{F_1}$	P	R	$\mathbf{F_1}$	
	Human	$91.09_{\pm 0.49}$	$93.17_{\pm 0.17}$	$92.12_{\pm 0.33}$	$65.21_{\pm 2.32}$	$47.48_{\pm 1.83}$	$54.93_{\pm 1.67}$	$55.07_{\pm 0.31}$	$61.86_{\pm 0.44}$	$58.26_{\pm 0.19}$	$54.30_{\pm 1.60}$	$55.38_{\pm 1.75}$	$54.79_{\pm 0.26}$	
							gpt-4o-mi	ni-2024-07-1	8					
	Baseline	$64.65_{\pm 0.85}$	$80.37_{\pm 0.50}$	$71.66_{\pm 0.41}$	$47.35_{\pm 2.46}$	$55.18_{\pm 2.84}$	$50.88_{\pm 1.14}$	$20.32_{\pm 5.26}$	$13.93_{\pm 2.76}$	$16.42_{\pm 3.34}$	$11.09_{\pm 0.97}$	$17.83_{\pm 2.02}$	$13.59_{\pm 0.52}$	
25	ICL RAG w/ST RAG w/OpenAI	$\begin{array}{c} 76.48 _{\pm 0.43} \\ 84.48 _{\pm 1.04} \\ 87.35 _{\pm 0.65} \end{array}$	$\begin{array}{c} 82.06_{\pm 0.35} \\ 88.99_{\pm 0.65} \\ 90.71_{\pm 0.34} \end{array}$	$\begin{array}{c} 79.17 _{\pm 0.25} \\ 86.68 _{\pm 0.85} \\ \textbf{89.00} _{\pm 0.29} \end{array}$	$\begin{array}{c} 53.18 \pm 3.22 \\ 51.42 \pm 2.63 \\ 52.26 \pm 2.24 \end{array}$	$\begin{array}{c} 52.24_{\pm 2.73} \\ 50.98_{\pm 1.52} \\ 49.75_{\pm 1.51} \end{array}$	$\begin{array}{c} \textbf{52.58}_{\pm 0.78} \\ 51.14_{\pm 1.01} \\ 50.93_{\pm 0.93} \end{array}$	$\begin{array}{c} 44.06 {\scriptstyle \pm 0.69} \\ 46.09 {\scriptstyle \pm 0.66} \\ 47.04 {\scriptstyle \pm 0.23} \end{array}$	$\begin{array}{c} 52.04_{\pm 1.57} \\ 54.38_{\pm 1.07} \\ 57.56_{\pm 1.44} \end{array}$	$\begin{array}{c} 47.71 _{\pm 0.79} \\ 49.89 _{\pm 0.70} \\ 51.77 _{\pm 0.66} \end{array}$	$\begin{array}{c} 21.23 \pm 1.46 \\ 20.29 \pm 0.78 \\ 21.26 \pm 1.69 \end{array}$	$\begin{array}{c} 45.26_{\pm 1.72} \\ 49.47_{\pm 1.80} \\ 56.74_{\pm 1.37} \end{array}$	$\begin{array}{c} 28.86_{\pm 1.24} \\ 28.77_{\pm 0.93} \\ \textbf{30.91}_{\pm 1.94} \end{array}$	
50	ICL RAG w/ST RAG w/OpenAI	$\begin{array}{c} 79.77_{\pm 0.34} \\ 86.73_{\pm 1.03} \\ 87.43_{\pm 0.48} \end{array}$	$\begin{array}{c} 82.64_{\pm 0.49} \\ 89.29_{\pm 0.84} \\ 91.39_{\pm 0.16} \end{array}$	$\begin{array}{c} 81.18_{\pm 0.29} \\ 87.99_{\pm 0.90} \\ \textbf{89.36}_{\pm 0.27} \end{array}$	$\begin{array}{c} 55.75_{\pm 2.80} \\ 53.74_{\pm 3.02} \\ 56.53_{\pm 2.35} \end{array}$	$\begin{array}{c} 49.53_{\pm 3.07} \\ 48.74_{\pm 4.44} \\ 50.29_{\pm 2.64} \end{array}$	$\begin{array}{c} 52.33_{\pm 0.97}\\ 50.90_{\pm 1.34}\\ \textbf{53.14}_{\pm 0.75}\end{array}$	$\begin{array}{c} 45.12_{\pm 0.82} \\ 46.46_{\pm 1.34} \\ 47.32_{\pm 0.92} \end{array}$	$\begin{array}{c} 54.35_{\pm 2.02} \\ 55.46_{\pm 1.21} \\ 58.44_{\pm 1.21} \end{array}$	$\begin{array}{c} 49.28 \pm 0.88 \\ 50.56 \pm 1.29 \\ \textbf{52.28} \pm 0.65 \end{array}$	$\begin{array}{c} 20.56 {\scriptstyle \pm 0.89} \\ 22.22 {\scriptstyle \pm 1.47} \\ 23.88 {\scriptstyle \pm 1.09} \end{array}$	$\begin{array}{c} 47.42_{\pm 2.01} \\ 52.60_{\pm 1.41} \\ 54.28_{\pm 2.26} \end{array}$	$\begin{array}{c} 28.66_{\pm 0.85}\\ 31.20_{\pm 1.32}\\ \textbf{33.13}_{\pm \textbf{0.77}}\end{array}$	
75	ICL RAG w/ST RAG w/OpenAI	$78.74_{\pm 1.02}\\86.91_{\pm 0.31}\\88.07_{\pm 0.35}$	$\frac{83.17_{\pm 0.55}}{89.25_{\pm 0.44}}$ 91.44 <sub>±0.28</sub>	$\begin{array}{c} 80.89_{\pm 0.66} \\ 88.06_{\pm 0.26} \\ \textbf{89.72}_{\pm 0.25} \end{array}$	$51.90_{\pm 4.29} \\ 53.80_{\pm 1.75} \\ 55.72_{\pm 4.22}$	$\begin{array}{c} 52.85 \\ \pm 1.95 \\ 51.79 \\ \pm 1.88 \\ 51.71 \\ \pm 3.34 \end{array}$	$\begin{array}{c} 52.24_{\pm 1.76} \\ 52.73_{\pm 0.80} \\ \textbf{53.43}_{\pm 0.54} \end{array}$	$\begin{array}{r} 44.40 \scriptstyle{\pm 0.63} \\ 47.22 \scriptstyle{\pm 0.98} \\ 47.04 \scriptstyle{\pm 1.29} \end{array}$	$\begin{array}{c} 53.89_{\pm 1.79} \\ 55.57_{\pm 0.43} \\ 58.19_{\pm 1.18} \end{array}$	$\begin{array}{c} 48.67_{\pm 0.69} \\ 51.05_{\pm 0.60} \\ \textbf{52.02}_{\pm 1.15} \end{array}$	$\begin{array}{r} 20.84 \pm 1.59 \\ 21.39 \pm 0.87 \\ 24.66 \pm 1.34 \end{array}$	$\begin{array}{r} 52.06_{\pm 1.01} \\ 52.85_{\pm 1.10} \\ 55.39_{\pm 3.19} \end{array}$	$\begin{array}{c} 29.73 \pm 1.58 \\ 30.43 \pm 0.73 \\ \textbf{34.06} \pm 0.88 \end{array}$	

#E	Mathad		CoNLL20	)03		WNUT-1	17	GUM SK		SKILLSPAN			
#EX.	Methou	Р	R	$\mathbf{F}_1$	P	R	$\mathbf{F}_1$	P	R	$\mathbf{F_1}$	P	R	$\mathbf{F}_1$
	Human	$91.09_{\pm 0.49}$	$93.17_{\pm 0.17}$	$92.12_{\pm 0.33}$	$65.21_{\pm 2.32}$	$47.48_{\pm 1.83}$	$54.93_{\pm 1.67}$	$55.07_{\pm 0.31}$	$61.86_{\pm 0.44}$	$58.26_{\pm 0.19}$	$54.30_{\pm 1.60}$	$55.38_{\pm 1.75}$	$54.79_{\pm 0.26}$
							Qwen2.5	-72B-Instruc	t				
	Baseline	$26.97_{\pm 0.28}$	$60.80_{\pm 1.25}$	$37.36_{\pm 0.30}$	$16.40_{\pm 1.30}$	$41.19_{\pm 1.23}$	$23.43_{\pm 1.42}$	$6.32_{\pm 0.21}$	$27.46 \pm 0.97$	$10.28 \pm 0.35$	$4.89{\scriptstyle \pm 0.41}$	$13.79_{\pm 2.15}$	$7.21_{\pm 0.69}$
25	ICL RAG w/ST RAG w/OpenAI	$\begin{array}{c} 74.57_{\pm 0.79} \\ 81.87_{\pm 0.72} \\ 84.81_{\pm 1.16} \end{array}$	$\begin{array}{c} 83.57_{\pm 0.82} \\ 89.90_{\pm 0.46} \\ 91.68_{\pm 0.64} \end{array}$	$\begin{array}{c} 78.81 _{\pm 0.30} \\ 85.69 _{\pm 0.56} \\ \textbf{88.11} _{\pm 0.82} \end{array}$	$\begin{array}{r} 45.58_{\pm 2.66} \\ 46.55_{\pm 2.70} \\ 48.33_{\pm 2.82} \end{array}$	$\begin{array}{c} 59.47_{\pm 3.09} \\ 45.68_{\pm 1.43} \\ 49.88_{\pm 1.89} \end{array}$	$\begin{array}{c} \textbf{51.49}_{\pm 0.92} \\ 46.06_{\pm 1.33} \\ 49.05_{\pm 1.86} \end{array}$	$\begin{array}{c} 41.69_{\pm 1.10} \\ 47.79_{\pm 1.05} \\ 47.16_{\pm 0.46} \end{array}$	$\begin{array}{c} 55.80 \scriptstyle{\pm 1.00} \\ 60.15 \scriptstyle{\pm 0.99} \\ 59.83 \scriptstyle{\pm 0.64} \end{array}$	$\begin{array}{c} 47.73 \pm 1.06 \\ \textbf{53.26} \pm 0.89 \\ 52.74 \pm 0.16 \end{array}$	$\begin{array}{c} 17.06_{\pm 2.18} \\ 18.25_{\pm 2.05} \\ 18.23_{\pm 1.90} \end{array}$	$\begin{array}{c} 31.17_{\pm 1.63} \\ 47.93_{\pm 2.08} \\ 50.07_{\pm 4.49} \end{array}$	$\begin{array}{c} 22.01_{\pm 2.13} \\ 26.40_{\pm 2.37} \\ \textbf{26.63}_{\pm 1.90} \end{array}$
50	ICL RAG w/ST RAG w/OpenAI	$\begin{array}{c} 77.48 \scriptstyle{\pm 0.51} \\ 84.30 \scriptstyle{\pm 0.98} \\ 85.96 \scriptstyle{\pm 1.44} \end{array}$	$\begin{array}{c} 83.34_{\pm 0.53}\\ 91.49_{\pm 0.85}\\ 92.32_{\pm 0.30}\end{array}$	$\begin{array}{c} 80.30_{\pm 0.43} \\ 87.74_{\pm 0.77} \\ \textbf{89.02}_{\pm 0.66} \end{array}$	$\begin{array}{c} 45.04 \pm 1.78 \\ 45.60 \pm 2.82 \\ 48.66 \pm 2.91 \end{array}$	$\begin{array}{c} 59.31_{\pm 1.71} \\ 56.63_{\pm 1.52} \\ 57.09_{\pm 2.42} \end{array}$	$\begin{array}{c} 51.17_{\pm 1.16} \\ 50.45_{\pm 1.14} \\ \textbf{52.46}_{\pm 1.23} \end{array}$	$\begin{array}{r} 44.30_{\pm 1.09} \\ 48.83_{\pm 1.45} \\ 47.33_{\pm 0.81} \end{array}$	$\begin{array}{c} 57.69_{\pm 1.35} \\ 60.55_{\pm 1.05} \\ 61.23_{\pm 0.44} \end{array}$	$\begin{array}{c} 50.12_{\pm 1.14} \\ \textbf{54.06}_{\pm 1.25} \\ 53.38_{\pm 0.63} \end{array}$	$\begin{array}{c} 17.51 {\scriptstyle \pm 0.85} \\ 21.32 {\scriptstyle \pm 1.82} \\ 23.44 {\scriptstyle \pm 2.19} \end{array}$	$\begin{array}{c} 33.82 \pm 1.01 \\ 55.79 \pm 4.56 \\ 52.32 \pm 2.82 \end{array}$	$\begin{array}{c} 23.06_{\pm 0.86}\\ 30.84_{\pm 2.52}\\ \textbf{32.35}_{\pm 2.54}\end{array}$
75	ICL RAG w/ST RAG w/OpenAI	$\begin{array}{c} 77.50_{\pm 0.68} \\ 87.46_{\pm 0.39} \\ 86.77_{\pm 0.54} \end{array}$	$\begin{array}{c} 83.60_{\pm 0.74} \\ 91.95_{\pm 0.29} \\ 92.05_{\pm 0.72} \end{array}$	$\begin{array}{c} 80.43_{\pm 0.67} \\ 89.65_{\pm 0.31} \\ \textbf{89.34}_{\pm 0.61} \end{array}$	$\begin{array}{c} 52.14_{\pm 2.27} \\ 48.36_{\pm 3.25} \\ 48.56_{\pm 2.08} \end{array}$	$\begin{array}{c} 55.62_{\pm 3.11} \\ 55.51_{\pm 1.88} \\ 60.22_{\pm 1.52} \end{array}$	$\begin{array}{c} 53.72_{\pm 0.80} \\ 51.58_{\pm 1.04} \\ 53.72_{\pm 0.71} \end{array}$	$\begin{array}{c} 47.60_{\pm 0.77} \\ 50.29_{\pm 0.27} \\ 47.24_{\pm 1.27} \end{array}$	$\begin{array}{c} 57.08_{\pm 1.22} \\ 60.97_{\pm 0.51} \\ 60.34_{\pm 0.57} \end{array}$	$\begin{array}{c} 51.91_{\pm 0.80}\\ \textbf{55.11}_{\pm 0.17}\\ 52.98_{\pm 0.76}\end{array}$	$\begin{array}{c} 20.81_{\pm 1.15} \\ 20.99_{\pm 2.12} \\ 19.95_{\pm 0.74} \end{array}$	$\begin{array}{c} 48.26_{\pm 2.80} \\ 49.99_{\pm 1.23} \\ 50.74_{\pm 1.47} \end{array}$	$\begin{array}{c} 29.05_{\pm 1.26} \\ \textbf{29.52}_{\pm 2.10} \\ 28.62_{\pm 0.77} \end{array}$
							Llama3.5	-70B-Instruc	t				
	Baseline	$23.56_{\pm 0.10}$	$63.25_{\pm 0.17}$	$34.33_{\pm 0.15}$	$16.35_{\pm 0.74}$	$54.65_{\pm 0.42}$	$25.16_{\pm 0.84}$	$6.44_{\pm 0.08}$	$27.79_{\pm 0.35}$	$10.46_{\pm 0.13}$	$3.51_{\pm 0.08}$	$24.30_{\pm 0.64}$	$6.14_{\pm 0.12}$
25	ICL RAG w/ST RAG w/OpenAI	$\begin{array}{c} 73.59_{\pm 0.78} \\ 83.15_{\pm 1.42} \\ 68.32_{\pm 3.99} \end{array}$	$\begin{array}{c} 78.73_{\pm 1.03} \\ 86.37_{\pm 0.90} \\ 87.50_{\pm 1.82} \end{array}$	$\begin{array}{c} 76.06_{\pm 0.41} \\ 84.72_{\pm 0.54} \\ 76.65_{\pm 2.19} \end{array}$	$\begin{array}{r} 48.77_{\pm 2.20} \\ 36.68_{\pm 1.32} \\ 43.52_{\pm 4.33} \end{array}$	$\begin{array}{c} 47.66_{\pm 5.18} \\ 49.10_{\pm 3.77} \\ 44.71_{\pm 3.86} \end{array}$	$\begin{array}{c} \textbf{48.00}_{\pm 2.14} \\ \textbf{41.89}_{\pm 0.99} \\ \textbf{43.82}_{\pm 1.29} \end{array}$	$\begin{array}{c} 18.26_{\pm 2.80} \\ 43.09_{\pm 1.10} \\ 42.46_{\pm 1.75} \end{array}$	$\begin{array}{c} 41.83_{\pm 0.98} \\ 50.88_{\pm 2.31} \\ 48.87_{\pm 4.60} \end{array}$	$\begin{array}{c} 25.34_{\pm 2.68} \\ \textbf{46.63}_{\pm 0.89} \\ 45.29_{\pm 1.50} \end{array}$	$\begin{array}{c} 17.04_{\pm 0.52} \\ 19.62_{\pm 1.44} \\ 19.59_{\pm 1.52} \end{array}$	$\begin{array}{c} 45.86_{\pm 2.86} \\ 46.47_{\pm 1.76} \\ 42.16_{\pm 1.49} \end{array}$	$\begin{array}{c} 24.84_{\pm 0.95} \\ \textbf{27.55}_{\pm 1.21} \\ 26.73_{\pm 1.65} \end{array}$
50	ICL RAG w/ST RAG w/OpenAI	$\begin{array}{c} 76.13_{\pm 1.12} \\ 83.87_{\pm 0.69} \\ 68.36_{\pm 1.53} \end{array}$	$\begin{array}{c} 76.79_{\pm 1.24} \\ 88.57_{\pm 0.88} \\ 89.08_{\pm 0.75} \end{array}$	$\begin{array}{c} 76.44_{\pm 0.30} \\ \textbf{86.15}_{\pm 0.28} \\ 77.35_{\pm 0.97} \end{array}$	$\begin{array}{c} 50.24_{\pm 2.81} \\ 42.92_{\pm 2.03} \\ 44.14_{\pm 1.97} \end{array}$	$\begin{array}{c} 48.90_{\pm 2.24} \\ 48.79_{\pm 2.99} \\ 51.28_{\pm 2.94} \end{array}$	$\begin{array}{c} \textbf{49.48}_{\pm1.08} \\ \textbf{45.57}_{\pm0.76} \\ \textbf{47.36}_{\pm0.64} \end{array}$	$\begin{array}{c} 35.67_{\pm 1.83} \\ 43.76_{\pm 1.50} \\ 43.70_{\pm 2.43} \end{array}$	$\begin{array}{c} 48.79_{\pm 3.19} \\ 50.24_{\pm 2.00} \\ 49.70_{\pm 1.83} \end{array}$	$\begin{array}{c} 41.12_{\pm 1.07} \\ \textbf{46.73}_{\pm 0.49} \\ 46.45_{\pm 1.40} \end{array}$	$\begin{array}{c} 16.09_{\pm 0.97} \\ 17.69_{\pm 0.66} \\ 18.37_{\pm 2.42} \end{array}$	$\begin{array}{c} 44.15_{\pm 4.16} \\ 46.11_{\pm 4.63} \\ 44.41_{\pm 4.44} \end{array}$	$\begin{array}{c} 23.51_{\pm 0.71} \\ 25.50_{\pm 0.54} \\ \textbf{25.77}_{\pm 1.75} \end{array}$
75	ICL RAG w/ST RAG w/OpenAI	$74.94_{\pm 1.03} \\ 85.70_{\pm 0.60} \\ 76.99_{\pm 1.57}$	$\begin{array}{c} 75.15_{\pm 1.03} \\ 89.03_{\pm 0.55} \\ 87.46_{\pm 1.39} \end{array}$	$\begin{array}{c} 75.04_{\pm 0.70} \\ \textbf{87.33}_{\pm 0.23} \\ 81.87_{\pm 0.67} \end{array}$	$50.78_{\pm 1.74}$ $47.41_{\pm 3.89}$ $49.43_{\pm 4.27}$	$51.69_{\pm 2.43} \\ 51.36_{\pm 1.97} \\ 48.16_{\pm 5.99}$	$\begin{array}{c} \textbf{51.18}_{\pm 1.05} \\ 49.18_{\pm 1.61} \\ 48.39_{\pm 2.17} \end{array}$	$\begin{array}{r} \hline 39.62_{\pm 1.64} \\ 45.84_{\pm 1.19} \\ 44.46_{\pm 0.61} \end{array}$	$\begin{array}{r} 47.88_{\pm 3.05} \\ 50.36_{\pm 1.12} \\ 52.96_{\pm 1.65} \end{array}$	$\begin{array}{c} 43.30_{\pm1.39} \\ 47.98_{\pm0.68} \\ \textbf{48.33}_{\pm0.64} \end{array}$	$\begin{array}{r} 17.55_{\pm 1.05} \\ 18.87_{\pm 1.35} \\ 9.51_{\pm 1.65} \end{array}$	$51.80_{\pm 1.68} \\ 51.17_{\pm 2.06} \\ 47.74_{\pm 3.51}$	$\begin{array}{c} 26.19_{\pm 1.14} \\ 27.52_{\pm 1.18} \\ 15.83_{\pm 2.47} \end{array}$

Table 3: The  $F_1$ , precision and recall along with standard deviation are reported on the test set. The values are averaged over five different random initializations. #Ex. represents the number of context examples used. Baseline refers to the use of LLM with no context examples.

#E	Mathad		CoNLL20	003		WNUT-1	7	GUM SKI		SKILLSPA	KILLSPAN		
#EX.	Methou	P	R	$\mathbf{F}_1$	P	R	$\mathbf{F}_1$	P	R	$\mathbf{F_1}$	P	R	$\mathbf{F_1}$
	Human	$91.09_{\pm 0.49}$	$93.17_{\pm 0.17}$	$92.12_{\pm 0.33}$	$65.21_{\pm 2.32}$	$47.48_{\pm 1.83}$	$54.93_{\pm 1.67}$	$55.07_{\pm 0.31}$	$61.86_{\pm 0.44}$	$58.26_{\pm 0.19}$	$54.30_{\pm 1.60}$	$55.38_{\pm 1.75}$	$54.79_{\pm 0.26}$
							Qwen2.5	-7B-Instruct					
	Baseline	$21.79_{\pm 1.28}$	$62.11_{\pm 0.44}$	$32.24_{\pm 1.44}$	$20.95_{\pm 2.83}$	$44.08_{\pm 3.68}$	$28.36_{\pm 3.17}$	$3.27_{\pm 0.22}$	$14.10_{\pm 1.03}$	$5.31_{\pm 0.37}$	$5.41_{\pm 1.05}$	$35.29_{\pm 2.73}$	$9.35_{\pm 1.58}$
25	ICL RAG w/ST RAG w/OpenAI	$\begin{array}{c} 70.22 {\scriptstyle \pm 1.45} \\ 83.81 {\scriptstyle \pm 0.67} \\ 84.05 {\scriptstyle \pm 1.15} \end{array}$	$\begin{array}{c} 75.96_{\pm 1.49} \\ 89.68_{\pm 0.57} \\ 90.85_{\pm 0.31} \end{array}$	$\begin{array}{c} 72.95 {\scriptstyle \pm 0.30} \\ 86.64 {\scriptstyle \pm 0.45} \\ \textbf{87.32} {\scriptstyle \pm 0.65} \end{array}$	$\begin{array}{r} 47.79_{\pm 3.40} \\ 37.82_{\pm 3.45} \\ 50.22_{\pm 3.43} \end{array}$	$\begin{array}{c} 47.02_{\pm 2.78} \\ 49.68_{\pm 3.12} \\ 41.75_{\pm 4.88} \end{array}$	$\begin{array}{c} \textbf{47.29}_{\pm 1.63} \\ \textbf{42.81}_{\pm 2.14} \\ \textbf{45.30}_{\pm 1.73} \end{array}$	$\begin{array}{c} 28.31_{\pm 1.01} \\ 35.90_{\pm 1.86} \\ 34.63_{\pm 1.09} \end{array}$	$\begin{array}{c} 44.01_{\pm 0.97} \\ 50.09_{\pm 2.04} \\ 49.97_{\pm 1.00} \end{array}$	$\begin{array}{c} 34.43_{\pm 0.57} \\ \textbf{41.80}_{\pm 1.65} \\ 40.89_{\pm 0.52} \end{array}$	$\begin{array}{c} 14.12 {\scriptstyle \pm 0.88} \\ 15.99 {\scriptstyle \pm 0.38} \\ 20.45 {\scriptstyle \pm 1.35} \end{array}$	$\begin{array}{c} 54.89_{\pm 1.48} \\ 55.06_{\pm 0.93} \\ 54.60_{\pm 4.83} \end{array}$	$\begin{array}{c} 22.44_{\pm 1.08} \\ 24.77_{\pm 0.42} \\ \textbf{29.67}_{\pm 1.29} \end{array}$
50	ICL RAG w/ST RAG w/OpenAI	$\begin{array}{c} 72.55_{\pm 1.01} \\ 85.78_{\pm 0.69} \\ 80.90_{\pm 1.79} \end{array}$	$\begin{array}{c} 78.54_{\pm 0.32} \\ 90.21_{\pm 0.38} \\ 91.55_{\pm 0.36} \end{array}$	$\begin{array}{c} 75.42 {\scriptstyle \pm 0.57} \\ \textbf{87.94} {\scriptstyle \pm 0.43} \\ 85.89 {\scriptstyle \pm 1.13} \end{array}$	$\begin{array}{r} 47.95 {\scriptstyle \pm 3.18} \\ 52.14 {\scriptstyle \pm 4.79} \\ 41.97 {\scriptstyle \pm 2.87} \end{array}$	$\begin{array}{c} 49.36_{\pm 3.94} \\ 44.00_{\pm 3.95} \\ 48.62_{\pm 5.64} \end{array}$	$\begin{array}{c} \textbf{48.54}_{\pm 2.51} \\ \textbf{47.41}_{\pm 0.49} \\ \textbf{44.75}_{\pm 0.98} \end{array}$	$\begin{array}{c} 33.51 {\scriptstyle \pm 0.76} \\ 39.38 {\scriptstyle \pm 1.31} \\ 34.63 {\scriptstyle \pm 1.32} \end{array}$	$\begin{array}{c} 43.59_{\pm 1.18} \\ 49.85_{\pm 1.69} \\ 50.61_{\pm 1.67} \end{array}$	$\begin{array}{c} 37.88_{\pm 0.56} \\ \textbf{43.97}_{\pm 0.44} \\ 41.11_{\pm 1.40} \end{array}$	$\begin{array}{c} 15.13 \scriptstyle{\pm 0.90} \\ 17.36 \scriptstyle{\pm 1.08} \\ 18.12 \scriptstyle{\pm 0.94} \end{array}$	$\begin{array}{c} 52.64_{\pm 2.98} \\ 51.06_{\pm 2.65} \\ 56.68_{\pm 3.93} \end{array}$	$\begin{array}{c} 23.47_{\pm 1.01} \\ 25.87_{\pm 1.02} \\ \textbf{27.41}_{\pm 0.73} \end{array}$
75	ICL RAG w/ST RAG w/OpenAI	$\begin{array}{c} 81.36_{\pm 1.19} \\ 86.54_{\pm 1.93} \\ 81.67_{\pm 1.51} \end{array}$	$\begin{array}{c} 75.72_{\pm 1.00} \\ 88.40_{\pm 1.15} \\ 90.96_{\pm 0.30} \end{array}$	$\begin{array}{c} 78.43_{\pm 0.63} \\ \textbf{87.44}_{\pm 0.87} \\ 86.06_{\pm 0.88} \end{array}$	$\begin{array}{r} 47.90_{\pm 5.76} \\ 52.39_{\pm 4.73} \\ 48.73_{\pm 1.31} \end{array}$	$\begin{array}{c} 47.78_{\pm 3.57} \\ 47.17_{\pm 1.99} \\ 47.85_{\pm 2.49} \end{array}$	$\begin{array}{c} 47.51_{\pm 1.97} \\ \textbf{49.48}_{\pm 1.30} \\ 48.25_{\pm 1.30} \end{array}$	$\begin{array}{r} 34.23_{\pm 2.14} \\ 40.14_{\pm 1.39} \\ 39.56_{\pm 1.25} \end{array}$	$\begin{array}{c} 46.40_{\pm 1.10} \\ 48.15_{\pm 1.09} \\ 50.87_{\pm 1.25} \end{array}$	$\begin{array}{c} 39.39_{\pm1.77} \\ 43.76_{\pm0.73} \\ \textbf{44.48}_{\pm0.55} \end{array}$	$\begin{array}{c} 12.98_{\pm 1.17} \\ 18.34_{\pm 0.46} \\ 14.07_{\pm 0.80} \end{array}$	$\begin{array}{c} 51.23_{\pm 6.20} \\ 46.32_{\pm 3.08} \\ 61.07_{\pm 0.86} \end{array}$	$\begin{array}{c} 20.68_{\pm1.81} \\ \textbf{26.25}_{\pm0.76} \\ 22.86_{\pm1.06} \end{array}$
							Llama-3.1	1-8B-Instruct					
	Baseline	$22.98_{\pm 0.67}$	$74.87_{\pm 0.48}$	$35.17_{\pm 0.83}$	$11.06_{\pm 2.70}$	$36.38_{\pm 10.40}$	$16.88_{\pm 4.16}$	$6.98_{\pm0.03}$	$28.22_{\pm 0.18}$	$11.19_{\pm 0.05}$	$3.03_{\pm 0.21}$	$20.37_{\pm 5.76}$	$5.22_{\pm 0.32}$
25	ICL RAG w/ST RAG w/OpenAI	$\begin{array}{c} 63.86 _{\pm 0.95} \\ 78.44 _{\pm 1.18} \\ 69.03 _{\pm 1.02} \end{array}$	$\begin{array}{c} 75.71_{\pm 1.61} \\ 86.16_{\pm 0.88} \\ 86.41_{\pm 2.27} \end{array}$	$\begin{array}{c} 69.26_{\pm 0.69} \\ \textbf{82.11}_{\pm 0.86} \\ 76.73_{\pm 1.02} \end{array}$	$\begin{array}{c} 35.94_{\pm 3.54} \\ 36.82_{\pm 4.64} \\ 32.83_{\pm 3.20} \end{array}$	$\begin{array}{c} 51.58_{\pm 2.70} \\ 43.86_{\pm 7.22} \\ 48.82_{\pm 7.41} \end{array}$	$\begin{array}{c} \textbf{42.23}_{\pm 2.38} \\ 39.38_{\pm 2.26} \\ 38.89_{\pm 2.78} \end{array}$	$\begin{array}{c} 33.95_{\pm 1.97} \\ 39.94_{\pm 2.23} \\ 40.77_{\pm 1.82} \end{array}$	$\begin{array}{c} 41.74_{\pm 3.02} \\ 46.48_{\pm 0.96} \\ 49.07_{\pm 2.80} \end{array}$	$\begin{array}{c} 37.39_{\pm1.85} \\ 42.92_{\pm1.20} \\ \textbf{44.45}_{\pm0.54} \end{array}$	$\begin{array}{c} 12.40_{\pm 0.85} \\ 14.95_{\pm 1.88} \\ 12.16_{\pm 0.97} \end{array}$	$\begin{array}{c} 33.63_{\pm 6.65} \\ 42.25_{\pm 6.45} \\ 41.55_{\pm 3.20} \end{array}$	$\begin{array}{c} 17.93_{\pm 0.66} \\ \textbf{21.87}_{\pm 1.51} \\ 18.79_{\pm 1.31} \end{array}$
50	ICL RAG w/ST RAG w/OpenAI	$\begin{array}{c} 67.78_{\pm 1.48} \\ 79.29_{\pm 3.86} \\ 69.98_{\pm 1.50} \end{array}$	$\begin{array}{c} 76.79_{\pm 0.69} \\ 86.85_{\pm 2.00} \\ 87.05_{\pm 2.15} \end{array}$	$\begin{array}{c} 72.01_{\pm 1.13} \\ \textbf{82.82}_{\pm 1.41} \\ 77.56_{\pm 0.74} \end{array}$	$\begin{array}{c} 40.49_{\pm 1.76} \\ 40.04_{\pm 4.26} \\ 39.75_{\pm 1.82} \end{array}$	$\begin{array}{c} 48.82_{\pm 2.88} \\ 48.60_{\pm 4.26} \\ 49.53_{\pm 2.63} \end{array}$	$\begin{array}{c} \textbf{44.20}_{\pm 1.03} \\ \textbf{43.59}_{\pm 1.04} \\ \textbf{44.03}_{\pm 0.76} \end{array}$	$\begin{array}{c} 36.43_{\pm 1.51} \\ 39.73_{\pm 1.81} \\ 40.89_{\pm 1.62} \end{array}$	$\begin{array}{c} 42.53_{\pm 2.12} \\ 46.54_{\pm 2.12} \\ 46.96_{\pm 2.27} \end{array}$	$\begin{array}{c} 39.22_{\pm 1.32} \\ 42.81_{\pm 0.99} \\ \textbf{43.66}_{\pm 0.67} \end{array}$	$\begin{array}{c} 12.94_{\pm 1.13} \\ 15.13_{\pm 0.96} \\ 12.09_{\pm 1.18} \end{array}$	$\begin{array}{c} 35.45_{\pm 2.12} \\ 47.13_{\pm 1.48} \\ 42.63_{\pm 3.13} \end{array}$	$\begin{array}{c} 18.90_{\pm 1.01} \\ \textbf{22.88}_{\pm 0.99} \\ 18.76_{\pm 1.19} \end{array}$
75	ICL RAG w/ST RAG w/OpenAI	$\begin{array}{c} 71.44_{\pm 2.02} \\ 82.36_{\pm 2.15} \\ 74.58_{\pm 2.51} \end{array}$	$\begin{array}{c} 77.86_{\pm 2.41} \\ 87.69_{\pm 1.70} \\ 85.21_{\pm 0.99} \end{array}$	$\begin{array}{c} 74.47_{\pm 1.00} \\ \textbf{84.91}_{\pm 0.88} \\ 79.51_{\pm 1.02} \end{array}$	$\begin{array}{c} 39.13_{\pm 1.16} \\ 39.92_{\pm 2.09} \\ 41.85_{\pm 2.52} \end{array}$	$\begin{array}{c} 48.70_{\pm 2.37} \\ 50.34_{\pm 3.12} \\ 47.17_{\pm 6.92} \end{array}$	$\begin{array}{c} 43.35_{\pm 0.84} \\ 44.42_{\pm 0.59} \\ 43.96_{\pm 2.54} \end{array}$	$\begin{array}{c} 34.33_{\pm 1.27} \\ 41.14_{\pm 1.47} \\ 41.99_{\pm 1.01} \end{array}$	$\begin{array}{c} 41.30_{\pm 2.37} \\ 43.71_{\pm 3.36} \\ 46.13_{\pm 2.68} \end{array}$	$\begin{array}{c} 37.43_{\pm 0.53} \\ 42.30_{\pm 1.57} \\ \textbf{43.91}_{\pm 0.93} \end{array}$	$\begin{array}{c} 13.37_{\pm 1.12} \\ 12.77_{\pm 0.10} \\ 10.42_{\pm 0.96} \end{array}$	$\begin{array}{c} 37.83_{\pm 4.68} \\ 45.34_{\pm 0.80} \\ 49.98_{\pm 3.64} \end{array}$	$\begin{array}{c} 19.67_{\pm 1.20} \\ \textbf{19.92}_{\pm 0.08} \\ 17.22_{\pm 1.29} \end{array}$

Table 4: The  $F_1$ , precision and recall along with standard deviation are reported on the test set. The values are averaged over five different random initializations. #Ex. represents the number of context examples used. Baseline refers to the use of LLM with no context examples.

#### D **Further Results on Different Sample Space Choices**

Tables 5 examine the influence of sample space  $\mathcal{X}$  and context size  $\mathcal{M}$  on entity recognition performance using the best-performing model, gpt-4o-mini, on the SKILLSPAN dataset. Increasing the context size from 25 to 75 generally improves the  $F_1$  score, though gains diminish beyond 50 examples. RAG consistently outperforms ICL in recall and  $F_1$  score, demonstrating its effectiveness in leveraging external knowledge, while ICL achieves higher precision but lower recall, suggesting a more conservative prediction approach. At a 10% sample space, ICL delivers competitive results, but as it increases to 20%, RAG maintains a clear advantage, achieving the highest  $F_1$  score of 32.39% at a context size of 75. Notably, for smaller dataset splits, RAG exhibits greater variability, similar to ICL, suggesting that when fewer examples are available, their performances converge. These findings underscore the importance of context size and external knowledge availability in optimizing RAG-based methods.

Sample Space	Context Size	Precision	Recall	F1 Score
		RAG		
	25	$21.83 \pm 1.22$	$56.94 \pm 1.17$	$31.53 \pm 1.18$
10%	50	$22.44_{\pm 1.32}$	$56.46_{\pm 2.46}$	$32.07_{\pm 1.13}$
	75	$22.82 \pm 0.58$	$55.82_{\pm 1.40}$	$32.34_{\pm 0.41}$
	25	$20.26 \pm 1.55$	$54.46_{\pm 3.71}$	$29.45_{\pm 1.29}$
20%	50	$21.00 \pm 0.81$	$57.40_{\pm 1.57}$	$30.74_{\pm 0.86}$
	75	$22.69 \pm 0.46$	$56.31_{\pm 2.06}$	$32.39_{\pm 0.60}$
		ICL		
	25	$22.57_{\pm 1.49}$	$48.72_{\pm 3.95}$	$30.74_{\pm 0.87}$
10%	50	$23.62_{\pm 0.85}$	$50.73_{\pm 1.33}$	$32.21_{\pm 0.67}$
	75	$23.12_{\pm 1.16}$	$51.09_{\pm 4.19}$	$31.76_{\pm 0.89}$
	25	$19.35_{\pm 1.57}$	$45.83_{\pm 4.74}$	$27.05_{\pm 0.66}$
20%	50	$22.02_{\pm 1.56}$	$51.17_{\pm 1.15}$	$30.76_{\pm 1.39}$
	75	$22.89_{\pm 1.11}$	$49.78_{\pm 2.89}$	$31.32_{\pm 0.99}$

Table 5: Study comparing RAG and ICL methods at different size of sample spaces (10% and 20%) and context sizes (25, 50, and 75). Experiments were conducted on the SKILLSPAN dataset using the gpt-4o-mini-2024-07-18 model. The results are presented with standard deviations, showing how performance metrics vary across sampling choices and conte

#### **Statistical Significance Test** E

This study evaluated various large language models across multiple datasets, considering different embeddings and examples as context. While some models clearly outperformed others in the results, the differences in predictions might not be statistically significant for certain models. Therefore, to determine the statistical significance of our findings, we conducted a non-parametric test. This test helps us assess whether there are significant differences among the models and, if so, identify which models differ statistically from each other.

The Friedman test (Pereira et al., 2015) is a non-parametric statistical test used to detect differences in performance across multiple related samples — in this case, different models evaluated over multiple datasets. It ranks the performance scores among datasets and assesses whether the rank distributions differ significantly among models. Let N be the number of datasets, K the number of models, and  $R_i$  be the sum of ranks for each model j. The Friedman test statistic  $chi_F^2$ , which follows a chi-square distribution, is calculated as follows:

$$\chi_F^2 = \frac{12N}{k(k+1)} \sum_{j=1}^k R_j^2 - 3N(k+1).$$
<sup>(1)</sup>

If the test statistic exceeds the critical value for a significance level  $\alpha = 0.01$ , we reject the null hypothesis, indicating that there are significant differences in performance among the models. If significant differences are found, the post-hoc Conover (Conover, 1999) test is performed to discover pair-wise statistical

937

926

927

928

929

930

931

932

933

934

935

936

938 939 940

941

942

943

944

945

946

947

948

949

950

951

952



Figure 4: Critical Difference diagram of average score ranks. The models connected with horizontal line shows no statistical difference. The models with lower ranks shows superior performance than those of higher ranks.

differences among models while adjusting for multiple comparisons. This test evaluates whether specificmodels differ significantly in performance.

Given that the Friedman test produces a test statistic of 114.42 with a *p*-value of  $7.71^{-18}$ , we reject the null 956 hypothesis, suggesting that at least one model shows a statistically significant difference in performance. 957 Consequently, we conducted the post-hoc Conover test. Figure 4 presents the statistical significance of the model rankings, with significant pairwise differences highlighted accordingly. The x-axis indicates the average rank of each model, where lower ranks closer to the left signify better performance. Each colored node corresponds to a particular model, labeled with its respective rank, while the black horizontal bars connecting multiple nodes highlight groups of models that do not show statistically significant differences at the specified confidence level. The top-performing combination is gpt4omini-OpenAI, 963 with an average rank of 1.9, indicating it consistently outperformed other approaches. Other strong 964 performers include Qwen2.5-72B-OpenAI (3), gpt-4o-mini-ST (3.8), and Qwen2.5-72B-ST (4.3). These 965 models have lower rankings and are clustered towards the left. In contrast, Llama3.1-8B-ICL (14), Llama3.1-8B-OpenAI (13), and Qwen2.5-7B-ICL (11) have the highest ranks, suggesting they performed the worst in comparison. These models do not overlap with the higher-ranked ones, highlighting their statistically inferior performance. Interestingly, Llama3.1-8B-ST shows no statistical differences when 969 compared to Llama3.5-70B, whether using ICL or RAG with OpenAI embedding. Similarly, Qwen2.5-7B, 970 when utilizing RAG with either OpenAI or ST embeddings, exhibits no statistical differences compared to Llama3.5-70B using ST embeddings and Qwen2.5-72B using ICL. These tests highlight a crucial aspect: a trade-off when addressing the NER task. Indeed, larger models, such as those with 70B parameters, may 973 not necessarily offer better performance than smaller models like Llama3.1-8B-ST or Owen2.5-7B. This 974 suggests that the additional computational resources required for bigger models might not always justify 975 their use, especially if smaller models can achieve statistically similar results.

## F Qualitative Analysis

977

This study broadly explores the efficacy of LLMs for data annotation tasks. Four different datasets of varying complexity are chosen. From Table 2, it is observed that the performance of LLMs decreases as dataset complexity increases. The performance of LLMs on the SKILLSPAN dataset is significantly lower than human annotation, suggesting that even the latest available LLMs struggle to annotate data when the task is complex. For instance, soft skills lack clear or distinct definitions, making the task more challenging. Similarly, the GUM dataset also poses challenges for LLMs due to its entity diversity. On the other hand, in the case of the WNUT-17 and CoNLL-2003 datasets, which consist of simpler entities (more details are reported in Section 4.1), annotations are easier to extract for an LLM given its prior knowledge. Furthermore, the quality of context in LLMs plays a major role, particularly in data annotation tasks, as indicated by Tables 2, 3, and 4, where the RAG-based approach significantly outperforms its counterpart. Moreover, for simpler datasets, the RAG-based approach achieves performance comparable to human annotation.

990To gain better insights into the performance of the proposed RAG-based approach, Table 6 presents the<br/>qualitative results for the SKILLSPAN dataset annotated by gpt-4o-mini. In this dataset, data annotation<br/>performance remains far below human-level, suggesting that the LLM struggles to extract sufficient

Table 6: Qualitative analysis of soft skills annotations on dataset samples using gpt-4o-mini-2024-07-18. The output of the best-performing model is reported. The highlighted texts in the first column are gold labels, while those in the other columns are the corresponding LLM-generated annotations.

№	Human	Baseline	ICL	RAG		
1.	Very good understanding	Very good understanding	Very good <mark>understanding</mark>	Very good <mark>understanding</mark>		
	of test automation	of test automation	of test automation	<mark>of test automation</mark>		
	frameworks.	frameworks.	frameworks.	frameworks.		
2.	Must have excellent verbal	Must have excellent verbal	Must have excellent verbal	Must have excellent verbal		
	and written skills being	and written skills being	and written skills being	and written skills being		
	able to communicate	able to communicate	able to communicate	able to communicate		
	effectively on both a	effectively on both a	effectively on both a	effectively on both a		
	technical and business	technical and business	technical and business	technical and business		
	level Ability to work under	level Ability to work under	level Ability to work under	level Ability to work under		
	pressure to resolve issues	pressure to resolve issues	pressure to resolve issues	pressure to resolve issues		
	affecting the production	affecting the production	affecting the production	affecting the production		
	services.	services.	services.	services.		
3.	Must have excellent work	Must have excellent work	Must have excellent work	Must have excellent work		
	ethic and be detail oriented	ethic and be detail oriented	ethic and be detail oriented	ethic and be detail oriented		
	and be able to work	and be able to work	and be able to work	and be able to work		
	independently.	independently.	independently.	independently.		
4.	Technical Skills Core Java.	Technical Skills Core Java.	Technical Skills <mark>Core</mark> Java.	Technical Skills <mark>Core</mark> Java.		
5.	You will work with	You will work with	You will work with	You will work with		
	the business to define	the business to define	the business to define	the business to define		
	requirements and have	requirements and have	requirements and have	requirements and have		
	excellent communication	excellent communication	excellent communication	excellent communication		
	skills to interpret these into	skills to interpret these into	skills to interpret these into	skills to interpret these into		
	consolidated development	consolidated development	consolidated development	consolidated development		
	scopes.	scopes.	scopes.	scopes.		

information from the context examples when the task is difficult. From Tables 2, 3, and 4, it is observed that LLM-generated annotations improve recall, whereas precision is compromised. Table 6 shows that in examples 1 and 4, the LLM incorrectly annotates soft skills that are not identified by human annotators, whereas in examples 2 and 3, the annotations are nearly identical to human annotations. In Example 5, the RAG-based approach performs comparably to human annotation, while both the baseline and ICL fail to do so.

## 

## **G** Prompt

This section presents the prompts used to generate the response of LLMs. These prompts are carefully1000synthesized to encompass all the components required to get structured output for both: (i) baseline, and1001(ii) in-context learning models.1002

#### **Baseline Prompt Structure**

#### **Task Description**

You are an advanced Named-Entity Recognition (NER) system.

Your task is to analyze the given sentence or passage, identify, extract, and classify specific named entities according to the following predefined entity types:

• {labels}

- For each sentence:
  - Label each word in the text with the appropriate entity type if it matches the specified categories.
  - Extract **multiple entities** of the same class if they exist.

The output should be in valid JSON format, with each word and its corresponding label as shown below.

Follow the structure strictly and do not add any other explanation.

In entities, label the word exactly as in the text. All the text is case-sensitive.

## Input

{input\_text}

#### **Context Prompt Structure**

#### **Task Description**

You are an advanced Named-Entity Recognition (NER) system.

Your task is to analyze the given sentence or passage, identify, extract, and classify specific named entities according to the following predefined entity types:

{labels}

For each sentence:

• Label each word in the text with the appropriate entity type if it matches the specified categories.

• Extract **multiple entities** of the same class if they exist.

The output should be in **valid JSON format**, with each word and its corresponding label as shown below. Follow the structure strictly and do not add any other explanation.

In entities, label the word exactly as in the text. All the text is case-sensitive.

#### Examples

{context\_examples}

#### Input

{input\_text}

#### **H** Examples

This section provides examples of prompts from the training data for different datasets used in this study. For visual purposes, we used only only *top5* examples in context. Follows several prompt examples for the: (*i*) CoNLL-2003, (*ii*) WNUT-17, (*iii*) SKILLSPAN datasets, and (*iv*) GUM datasets.

Example 1-CoNLL-2003

#### **Task Description**

You are an advanced Named-Entity Recognition (NER) system. Your task is to analyze the given sentence or passage, identify, extract, and classify specific named entities according to the following predefined entity types:

['PER', 'ORG', 'LOC', 'MISC']

For each sentence:

- Label each word in the text with the appropriate entity type if it matches the specified categories.
- Extract **multiple entities** of the same class if they exist.

The output should be in **valid JSON format**, with each word and its corresponding label as shown below. Follow the structure strictly and do not add any other explanation. In entities, label the word exactly as in the text. All the text is case-sensitive.

## Examples

["A South African boy is writing back to an American girl whose message in a bottle he found washed up on President Nelson Mandela 's old prison island .", [{'Entity': 'South African', 'Label': 'MISC'}, {'Entity': 'American',

1003

1004

1005

1006

1007

1008

'Label': 'MISC'}, {'Entity': 'Nelson Mandela', 'Label': 'PER'}]]

- ['A rottweiler dog belonging to an elderly South African couple savaged to death their two-year-old grandson who was visiting, police said on Thursday.', [{'Entity': 'South African', 'Label': 'MISC'}]]
- ['The princess, who has carved out a major role for herself as a helper of the sick and needy, is said to have turned to Mother Teresa for guidance as her marriage crumbled to heir to the British throne Prince Charles .', [{'Entity': 'Mother Teresa', 'Label': 'PER'}, {'Entity': 'British', 'Label ': 'MISC'}, {'Entity': 'Prince Charles', 'Label': 'PER'}]]
- ['South African answers U.S. message in a bottle .', [{'Entity': 'South African', 'Label': 'MISC'}, {'Entity': 'U.S.', 'Label': 'LOC'}]]
- ["But Carlo Hoffmann , an 11-year-old jailer 's son who found the bottle on the beach at Robben Island off Cape Town after winter storms , will send his letter back by ordinary mail on Thursday , the post office said .", [{'Entity': 'Carlo Hoffmann', 'Label': 'PER'}, {'Entity': 'Robben Island', 'Label': 'LOC'}, {'Entity': 'Cape Town', 'Label': 'LOC'}]]

#### Input

Revered skull of S. Africa king is Scottish woman 's .

#### Response

[Entity: S. Africa, Label: LOC, Entity: Scottish, Label: MISC]

#### Example 2–CoNLL-2003

#### **Task Description**

You are an advanced Named-Entity Recognition (NER) system. Your task is to analyze the given sentence or passage, identify, extract, and classify specific named entities according to the following predefined entity types:

#### ['PER', 'ORG', 'LOC', 'MISC']

For each sentence:

- Label each word in the text with the appropriate entity type if it matches the specified categories.
- Extract **multiple entities** of the same class if they exist.

The output should be in **valid JSON format**, with each word and its corresponding label as shown below. Follow the structure strictly and do not add any other explanation. In entities, label the word exactly as in the text. All the text is case-sensitive.

#### Examples

- ['Rwanda said on Saturday that Zaire had expelled 28 Rwandan Hutu refugees accused of being " trouble-makers " in camps in eastern Zaire .', [{' Entity': 'Rwanda', 'Label': 'LOC'}, {'Entity': 'Zaire', 'Label': 'LOC'}, {'Entity': 'Rwandan', 'Label': 'MISC'}, {'Entity': 'Hutu', 'Label': 'MISC '}, {'Entity': 'Zaire', 'Label': 'LOC'}]]
- ['Repatriation of 1.1 million Rwandan Hutu refugees announced by Zaire and Rwanda on Thursday could start within the next few days, an exiled Rwandan Hutu lobby group said on Friday .', [{'Entity': 'Rwandan Hutu', ' Label': 'MISC'}, {'Entity': 'Zaire', 'Label': 'LOC'}, {'Entity': 'Rwanda', 'Label': 'LOC'}, {'Entity': 'Rwandan Hutu', 'Label': 'MISC'}]]
- ['Innocent Butare, executive secretary of the Rally for the Return of Refugees and Democracy in Rwanda (RDR) which says it has the support of Rwanda \'s exiled Hutus, appealed to the international community to deter the two countries from going ahead with what it termed a "forced and inhuman action ".', [{'Entity': 'Innocent Butare', 'Label': 'PER'}, {' Entity': 'Rally for the Return of Refugees and Democracy in Rwanda', ' Label': 'ORG'}, {'Entity': 'RDR', 'Label': 'ORG'}, {'Entity': 'Rwanda', ' Label': 'LOC'}, {'Entity': 'Hutus', 'Label': 'MISC'}]]

['Rwanda says Zaire expels 28 Rwandan refugees .', [{'Entity': 'Rwanda', ' Label': 'LOC'}, {'Entity': 'Zaire', 'Label': 'LOC'}, {'Entity': 'Rwandan', 'Label': 'MISC'}]]

['Rwandan group says expulsion could be imminent .', [{'Entity': 'Rwandan', ' Label': 'MISC'}]]

## Input

Captain Firmin Gatera, spokesman for the Tutsi-dominated Rwandan army, told Reuters in Kigali that 17 of the 28 refugees handed over on Friday from the Zairean town of Goma had been soldiers in the former Hutu army which fled to Zaire in 1994 after being defeated by Tutsi forces in Rwanda 's civil war.

#### Response

[Entity: Captain Firmin Gatera, Label: PER, Entity: Rwandan, Label: MISC, Entity: Reuters, Label: ORG, Entity: Kigali, Label: LOC, Entity: Zairean, Label: MISC, Entity: Goma, Label: LOC, Entity: Hutu, Label: MISC, Entity: Zaire, Label: LOC, Entity: Tutsi, Label: MISC, Entity: Rwanda, Label: LOC]

#### Example 3–WNUT-17

#### **Task Description**

You are an advanced Named-Entity Recognition (NER) system. Your task is to analyze the given sentence or passage, identify, extract, and classify specific named entities according to the following predefined entity types:

['corporation', 'creative-work', 'group', 'location', 'person', 'product']

For each sentence:

- Label each word in the text with the appropriate entity type if it matches the specified categories.
- Extract **multiple entities** of the same class if they exist.

The output should be in **valid JSON format**, with each word and its corresponding label as shown below. Follow the structure strictly and do not add any other explanation. In entities, label the word exactly as in the text. All the text is case-sensitive.

## Examples

- ['@justinbieber i just wanna say you make me smile everyday :) thanks for smiling because when u smile i smile ! :)', []]
- ["@joeymcintyre I heart you . Even if I haven't seen u in months ... SEND A PIC !", []]
- ['@lovable\_sin OMG OMG OMG ! Thank you for " tumblring " it to me , I so wasn \'t expecting them today . OMG !', []]
- ['RT @aplusk : This made me laugh today http://bit.ly/bjOhom < --- courtesy of splurb . What made you laugh ?', []]
- ['RT @Sn00ki : Haha yes !!! I love that you knew that :) RT @trishamelissa @Sn00ki Is phenomenal the word of the day ?', []]

## Input

@jimmyfallon is following me ! OMG ! My life is now complete ! I heart you JF and have for years ! Thank you for making me laugh everyday !

#### Response

[Entity: @jimmyfallon, Label: person, Entity: JF, Label: person]

#### Example 4–WNUT-17

#### **Task Description**

You are an advanced Named-Entity Recognition (NER) system. Your task is to analyze the given sentence or passage, identify, extract, and classify specific named entities according to the following predefined entity types:

['corporation', 'creative-work', 'group', 'location', 'person', 'product']
<ul> <li>For each sentence:</li> <li>Label each word in the text with the appropriate entity type if it matches the specified categories.</li> <li>Extract multiple entities of the same class if they exist.</li> <li>The output should be in valid JSON format, with each word and its corresponding label as shown below.</li> <li>Follow the structure strictly and do not add any other explanation. In entities, label the word exactly as in the text. All the text is case-sensitive.</li> </ul>
Examples
['We are one step closer to our new kitchens . We chose a maker and had official measurements taken today !', []]
['We were all enjoying a glass of wine in the office when a fudge delivery showed up . I love my job . And I love Fridays .', []]
['800 miles to see clients , 3 ACC candidate/commissioner meetings , big press release , making it to Friday PRICELESS !', []]
["I hope the weeks keep flying. It 's actually fantastic the way none of the days dragged this week like NONE. :D", []]
['Feeling really good after great week in our SF and LA offices . Glad to kick back on AMerican flight back to NYC', [{'Entity': 'SF', 'Label': ' location'}, {'Entity': 'LA', 'Label': 'location'}, {'Entity': 'AMerican', 'Label': 'corporation'}, {'Entity': 'NYC', 'Label': 'location'}]]

#### Input

Great week in the Optimise office, another new client on board and we are close to signing a new team member

#### Response

[Entity: Optimise, Label: corporation]

#### Example 5–SKILLSPAN

#### **Task Description**

You are an advanced Named-Entity Recognition (NER) system. Your task is to analyze the given sentence or passage, identify, extract, and classify specific named entities according to the following predefined entity types:

#### ['Skill']

For each sentence:

- Label each word in the text with the appropriate entity type if it matches the specified categories.
- Extract **multiple entities** of the same class if they exist.

The output should be in **valid JSON format**, with each word and its corresponding label as shown below. Follow the structure strictly and do not add any other explanation. In entities, label the word exactly as in the text. All the text is case-sensitive.

#### Examples

['Hands on experience with automated testing using Java .', []]

['Experience with automation systems framework design/use and deployment .', []]

['Good understanding of Agile methodologies and Continuous Delivery .', []]

- ['Demonstrate clear understanding of automation and orchestration principles .', []]
- ['Good exposure to UI Frameworks like Angular Proficiency in SQL and Database development .', []]

["Ability to understand and use efficient Defect management regular view of test coverage to identify gaps and provide improvements Personal Specification 5+ years of relevant IT/quality assurance work experience Bachelor's degree in Computer Science or related field of study or equivalent relevant experience; demonstrated experience within the quality assurance / testing arena; demonstrated skills in quality assurance methods/processes and practices .", [{'Entity': 'understand and use efficient Defect management', 'Label': 'Skill'}, {'Entity': 'identify gaps ', 'Label': 'Skill'}]]

## Input

Very good understanding of test automation frameworks.

#### Response

[Entity: test automation frameworks, Label: Skill]

#### Example 6–SKILLSPAN

#### **Task Description**

You are an advanced Named-Entity Recognition (NER) system. Your task is to analyze the given sentence or passage, identify, extract, and classify specific named entities according to the following predefined entity types:

#### ['Skill']

For each sentence:

- Label each word in the text with the appropriate entity type if it matches the specified categories.
- Extract multiple entities of the same class if they exist.

The output should be in **valid JSON format**, with each word and its corresponding label as shown below. Follow the structure strictly and do not add any other explanation. In entities, label the word exactly as in the text. All the text is case-sensitive.

#### Examples

- ['Strong communication skills including the ability to express complex technical concepts to different audiences in writing and conference calls .', [{'Entity': 'communication skills', 'Label': 'Skill'}, {'Entity': ' express complex technical concepts to different audiences', 'Label': ' Skill'}]]
- ['Excellent organizational verbal and written communication skills .', [{' Entity ': 'organizational verbal and written communication skills', 'Label ': 'Skill '}]]
- ['Excellent organizational verbal and written communication skills .', [{' Entity ': 'organizational verbal and written communication skills ', 'Label ': 'Skill '}]]
- ['The ability to work within a team and in collaboration with others is critical to this position and excellent communication skills verbal and written are essential .', [{'Entity': 'work within a team and in collaboration with others', 'Label': 'Skill'}, {'Entity': 'communication skills', 'Label': 'Skill'}]]
- ['This role requires a wide variety of strengths and capabilities including Ability to work collaboratively in teams and develop meaningful relationships to achieve common goals Strong organizational skills Ability to multi-task and deliver to a tight deadline Excellent written and verbal communication skills Experience developing UI components in Angular Good experience in using design patterns UML OO concepts .', [{'Entity': 'work collaboratively in teams', 'Label': 'Skill'}, {'Entity': 'develop meaningful relationships', 'Label': 'Skill'}, {'Entity': 'achieve common goals', 'Label': 'Skill'}, {'Entity': 'organizational skills', 'Label': ' Skill'}, {'Entity': 'deliver to a tight deadline', 'Label': 'Skill'}, {'Entity': 'communication skills ', 'Label': 'Skill'}, {'Entity': 'developing UI components', 'Label': '

Skill'}, {'Entity': 'using design patterns', 'Label': 'Skill'}]]

#### Input

Must have excellent verbal and written skills being able to communicate effectively on both a technical and business level Ability to work under pressure to resolve issues affecting the production services .

#### Response

[Entity: verbal and written skills, Label: Skill, Entity: communicate effectively on both a technical and business level, Label: Skill, Entity: work under pressure, Label: Skill, Entity: resolve issues affecting the production services, Label: Skill]

Example 7-GUM

#### **Task Description**

You are an advanced Named-Entity Recognition (NER) system. Your task is to analyze the given sentence or passage, identify, extract, and classify specific named entities according to the following predefined entity types:

['abstract', 'animal', 'event', 'object', 'organization', 'person', 'place', 'plant', 'quantity', 'substance', 'time']

For each sentence:

- Label each word in the text with the appropriate entity type if it matches the specified categories.
- Extract **multiple entities** of the same class if they exist.

The output should be in **valid JSON format**, with each word and its corresponding label as shown below. Follow the structure strictly and do not add any other explanation. In entities, label the word exactly as in the text. All the text is case-sensitive.

#### Examples

- [ 'The 131-page document was found on Castlefrank Road in Kanata , Ontario in a rain-stained , tire-marked brown envelope by a passerby', 'Entities': [{'Entity': 'The 131-page document was found on Castlefrank Road in Kanata , Ontario in a rain-stained , tire-marked brown envelope by a passerby', 'Label': 'event'}]]
- ['Also the language is important in writing and in literature', 'Entities': [{'Entity': 'the language', 'Label': 'abstract'}, {'Entity': 'writing', Label': 'abstract'}, {'Entity': 'literature', 'Label': 'abstract'}]]

['Ingredients Basil comes in many different varieties, each of which have a unique flavor and smell', 'Entities': [{'Entity': 'Ingredients', 'Label': 'object'}, {'Entity': 'Basil', 'Label': 'plant'}, {'Entity': 'many different varieties', 'Label': 'abstract'}, {'Entity': 'each of which', ' Label': 'abstract'}, {'Entity': 'a unique flavor and smell', 'Label': ' abstract'}]]

['We do not want to just traffic in the same 24 hour news cycle', 'Entities': [{'Entity': 'We do not want to just traffic in the same 24 hour news cycle ', 'Label': 'abstract'}]]

#### Input

If you are just visiting York for the day, using a Park and Ride [1] costs a lot less than trying to park in or near the city centre, and there are five sites dotted around the Outer Ring Road

#### Response

['Entity': 'York', 'Label': 'place', 'Entity': 'the day', 'Label': 'time', 'Entity': 'a Park and Ride', 'Label': 'object', 'Entity': 'the city centre', 'Label': 'place', 'Entity': 'five sites', 'Label': 'quantity', 'Entity': 'the Outer Ring Road', 'Label': 'place']

<sup>[&#</sup>x27;You go through quite a bit', 'Entities': [{'Entity': 'You', 'Label': 'person '}, {'Entity': 'quite a bit', 'Label': 'quantity'}]]

#### Example 8–GUM

#### **Task Description**

You are an advanced Named-Entity Recognition (NER) system. Your task is to analyze the given sentence or passage, identify, extract, and classify specific named entities according to the following predefined entity types:

#### ['abstract', 'animal', 'event', 'object', 'organization', 'person', 'place', 'plant', 'quantity', 'substance', 'time']

For each sentence:

- Label each word in the text with the appropriate entity type if it matches the specified categories.
- Extract multiple entities of the same class if they exist.

The output should be in **valid JSON format**, with each word and its corresponding label as shown below. Follow the structure strictly and do not add any other explanation. In entities, label the word exactly as in the text. All the text is case-sensitive.

#### Examples

- ['" NASA Administrator Charles Bolden announces where four space shuttle orbiters will be permanently displayed at the conclusion of the Space Shuttle Program during an event commemorating the 30th anniversay of the first shuttle launch on April 12, 2011', 'Entities': [{'Entity': 'NASA Administrator Charles Bolden', 'Label': 'person'}, {'Entity': 'where four space shuttle orbiters will be permanently displayed', 'Label': 'place'}, {'Entity': 'the conclusion of the Space Shuttle Program', 'Label': 'event '}, {'Entity': 'an event', 'Label': 'event'}, {'Entity': '30th anniversay of the first shuttle launch', 'Label': 'event'}, {'Entity': 'April 12, 2011', 'Label': 'time'}]]
- ['NASA celebrated the launch of the first space shuttle Tuesday at an event at the Kennedy Space Center (KSC) in Cape Canaveral, Florida', 'Entities ': [{'Entity': 'NASA', 'Label': 'organization'}, {'Entity': 'the launch of the first space shuttle', 'Label': 'event'}, {'Entity': 'Tuesday', 'Label ': 'time'}, {'Entity': 'an event', 'Label': 'event'}, {'Entity': 'Kennedy Space Center', 'Label': 'place'}, {'Entity': 'KSC', 'Label': 'place'}, {' Entity': 'Cape Canaveral, Florida', 'Label': 'place'}]
- ['Looking back : Space Shuttle Columbia lifts off on STS-1 from Launch Pad 39A at the Kennedy Space Center on April 12, 1981', 'Entities': [{'Entity': 'Space Shuttle Columbia', 'Label': 'object'}, {'Entity': 'STS-1', 'Label': 'event'}, {'Entity': 'Launch Pad 39A', 'Label': 'place'}, {'Entity': ' Kennedy Space Center', 'Label': 'place'}, {'Entity': 'April 12, 1981', ' Label': 'time'}]]
- ['At the ceremony, NASA Administrator Charles Bolden announced the locations that would be given the three remaining Space Shuttle orbiters following the end of the Space Shuttle program', 'Entities': [{'Entity': 'the ceremony', 'Label': 'event'}, {'Entity': 'NASA Administrator Charles Bolden', 'Label': 'person'}, {'Entity': 'the locations', 'Label': 'place '}, {'Entity': 'the three remaining Space Shuttle orbiters', 'Label': ' object'}, {'Entity': 'the end of the Space Shuttle program', 'Label': ' event'}]]
- ['On April 12, 1981, Space Shuttle Columbia lifted off from the Kennedy Space Center on STS-1, the first space shuttle mission', 'Entities': [{' Entity': 'April 12, 1981', 'Label': 'time'}, {'Entity': 'Space Shuttle Columbia', 'Label': 'object'}, {'Entity': 'Kennedy Space Center', 'Label': 'place'}, {'Entity': 'STS-1', 'Label': 'event'}, {'Entity': 'the first space shuttle mission', 'Label': 'event'}]]

#### Input

Tuesday, September 22, 2015 Discovery is undergoing decommissioning and currently being prepped for display by removing toxic materials from the orbiter

#### Response

['Entity': 'Tuesday', 'Label': 'time', 'Entity': 'September 22, 2015', 'Label': 'time', 'Entity': 'Discovery', 'Label': 'object', 'Entity': 'decommissioning', 'Label': 'event', 'Entity': 'display', 'Label': 'event', 'Entity': 'toxic materials',

'Label': 'substance', 'Entity': 'the orbiter', 'Label': 'object']