Assessing the quality of information extraction

Anonymous Author(s) Affiliation Address email

Abstract

Advances in large language models have notably enhanced the efficiency of infor-1 mation extraction from unstructured and semi-structured data sources. As these 2 technologies become integral to various applications, establishing an objective 3 4 measure for the quality of information extraction becomes imperative. However, 5 the scarcity of labeled data presents significant challenges to this endeavor. In this paper, we introduce an automatic framework to assess the quality of the in-6 formation extraction/retrieval and its completeness. The framework focuses on 7 information extraction in the form of entity and its properties. We discuss how to 8 handle the input/output size limitations of the large language models and analyze 9 their performance when extracting the information. In particular, we introduce 10 11 scores to evaluate the quality of the extraction and provide an extensive discussion on how to interpret them. 12

13 1 Introduction

In the domain of natural language processing (NLP), information extraction (IE) stands as a critical task, transforming unstructured or semi-structured data into a structured format conducive to indexing, exploration, and further analysis. The increasing amount of data across digital platforms underscores the urgency for sophisticated IE techniques that can parse through volumes of information with precision. An extensive survey about IE is provided by [1], where the authors highlight the complexity of processing and analyzing text to derive meaningful information, given the heterogeneity and volume of such data.

21 Large language models (LLMs) have revolutionized IE by introducing generative methods for structuring knowledge from text. LLMs excel across diverse domains without extensive task-specific 22 training. A survey by [9] details the progress of LLMs on IE tasks. Here, the authors address specific 23 aspects of information extraction, including entity recognition, relation extraction, event detection, 24 and universal IE. They review the existing models and their efficiency on a comprehensive collection 25 of annotated benchmarks. Nonetheless, the challenge of quantitatively assessing the quality and 26 completeness of extracted information persists, particularly in the absence of labeled datasets for 27 benchmarking. Before conducting the experiments introduced in this paper, we perform IE on a vast 28 corpus of business documents utilizing LLMs. While the extraction process is beyond the scope of 29 30 this paper, some details about the extraction are given in Section 3.

To measure the quality of extraction, we propose an evaluation framework that relies on artificially 31 generated complex information which is infused into the document to test the efficiency of LLMs in 32 IE tasks. This paper introduces an iterative extraction process and a novel score, MINEA (Multiple 33 Infused Needle Extraction Accuracy), to address the critical need for objective quality assessment 34 35 measures. By inserting artificial information ("needles") into the data, the proposed method creates 36 a synthetic ground truth for evaluation, enabling the measurement of extraction quality in various specific domains even without manually labeled data. The empirical analysis demonstrates the 37 utility of MINEA for evaluating LLM-based IE in scenarios where ground truth is unavailable. By 38

³⁹ automating the quality assessment of information extraction, the framework could reduce the need

40 for manual review by experts, saving time and resources and thus enhance the efficiency and accuracy

41 of information extraction from large volumes of unstructured data.

The paper is organized as follows: Section 2 presents a related work that inspired us when developing our IE quality assessment method; Section 3 sketch a way in which structured information is obtained using LLMs; Section 4 deals with shortcomings arising when treating long contexts by LLMs; finally Section 5 introduces the novel method to access the quality of IE and provide the reader with practical tips; Sections 4 and 5 are supplemented by numerical studies. The data used in these studies are an internal set of documents related to a business case in the healthcare industry.

48 2 Related work

A common practice in many specialized IE tasks is that well-trained experts review what was extracted
 and provide ground truth as done in [5]. Such an approach is relatively reliable, however, it is manual
 and very time-consuming.

In [4] they suggest *summary score without reference* (SUSWIR), a score to evaluate the quality of text summaries without the need for human annotations. The SUSWIR score can be used for IE tasks where the extracted information is viewed as a compression of original data. The score compares the original text with its summary. From its nature, it is very useful when comparing the outputs of extraction tasks among themselves, i.e., the best extraction/summary has the highest score value. On the other hand, its ability to provide an objective absolute evaluation of a single extraction is disadvantaged because the desirable output is not known.

59 Recently, an effort to eliminate the requirement for human involvement relies on LLMs. These prove

themselves as highly cost-effective data creators, either by labeling unlabeled data or generating data
 given the labels, see [7]. Therefore they may substitute human experts providing the ground truth by

62 doing their work in an automatic way.

Needle In A Haystack (NIAH)¹ evaluation is a tool designed to evaluate the performance of LLMs in retrieval across different sizes of context. Short targeted information, the 'needle', is inserted into a large, more complex text body, the 'haystack'. The goal is to test an LLM's ability to find and make use of this piece of information.

Our method builds on LLMs acting as data creators, but instead of annotating the complete data, it only automatizes the process of creating the needle. I.e., given an original text, an LLM generates the needle. The needle then substitutes the ground truth.

70 **3** Capturing the structure

The form of needles depends on a form of data, on structure capturing the information and on the task being solved. The needles can be short paragraphs of text, account records, graph nodes as you extract information from continuous text, table, graph, respectively. The structured arrangement of information is beneficial for consecutive processing and analysis. It helps to highlight relationships among distinct information pieces. There are countless ways to impose a structure on unstructured data in order to capture the relevant information. To demonstrate our methodology for measuring the quality of information extraction, we specify a particular structure and tailor the needles to it.

78 **3.1 Schema**

To impose a structure on the data, we adopt the idea of schema markup [3] which is used to communicate the content of a web page to the search tool. The schema markup is in the form of structured data and can be viewed as a compression of the essential information. The structure is defined by Schema.org² vocabulary which is a set of entity types, each associated with a set of properties and hierarchically arranged. Figure 1 shows an example of structured information inspired by Schema.org. It describes three entities of types 'Insight', 'Person' and 'Organization'. Each

¹https://github.com/gkamradt/LLMTest_NeedleInAHaystack

²https://schema.org

- ⁸⁵ type has its own set of properties, e.g., an entity of type 'Person' is described by 'type', 'name',
- ⁸⁶ 'birthDate', 'worksFor', and 'jobTitle'. In other words, each entity is a set of key-value pairs, e.g.,
- ⁸⁷ 'name' is the key and 'AI Enthusiast' is the value.

```
E
  {
   "@type": "Insight",
   "name": "Information exctraction tested by Needle in a Haystack test",
   "description": "A short targeted information pieces, the 'needles', are inserted to
       a large, more complex text body, the 'haystack'. The quality of information
       extraction task is measured by ratio of succesfully extracted needles.",
   "keywords": "information extraction (9), large language models (8), quality evaluation
       (10), needle in a hayastack (8), named entity recogniction (7), schema.org (6)",
   "author": {
      "@type": "Person",
      "name": "AI Enthusiast",
      "birthDate": " ",
      "worksFor": {
         "@type": "Organization",
         "name": "Creative Dock"
         "description": "Creative Dock builds and scales disruptive tech companies, as
               a startup and corporate venture builder. The company provides end-to-end
               venture-building, from idea to building and scaling.",
         "keywords": "tech company (8), venture builder (9), AI (7), startup growth (8)"
        }
      "jobTitle": "Data Scientist"
     }
 }
]
```

Figure 1: Toy example: structured information encapsulating three entities using schema.org.

Similarly, we extract and compress the relevant information contained in data using an LLM. 88 Schema.org presents a clear basis for the categorization of various entities contained in data. In the 89 rest of the paper, by schema we mean a predetermined set of types, such as {'Person', 'Project', 90 'Product', 'Legislation', 'Event', 'OpportunityArea', 'Insight', 'Substance', 'Thing', 'BioChemEn-91 tity', 'MedicalCondition'}, together with their properties. The schema is set at the beginning and 92 the information to be extracted depends on it. Therefore the schema has to be tailored to a particular 93 scope of the (proprietary) knowledge and application. If a more complex or uncommon entity needs 94 to be captured, it is natural and very easy to extend the set of core types by more detailed descriptive 95 and custom vocabulary. E.g., 'Insight' and 'OpportunityArea' are not native Schema.org types, but 96 we will use them in our study. The usage of suitably tailored schema is beneficial for specialized 97 applications since it narrows the information to the relevant core and hence potentially improves the 98 overall performance. On the other hand, the usage of schemata is not restrictive as the scope can be 99 always extended by using a broader set of types. 100

101 3.2 The role of LLMs

LLMs are rather effective in the creation of structured data, cf. [9]. Using dedicated prompts, we get 102 103 a structured text file describing entities found in the documents and matching types of predefined 104 schema. The predefined schema (types and properties) is given to an LLM within the prompt. The LLM is asked to analyse the document, identify an information relevant to the mentioned types of 105 entities and populate the schema with this information. It is asked to be attentive to nested entities, 106 maintain consistency and uniqueness of extracted entities. Indeed, LLM is not prohibited from 107 extracting entities whose types do not appear in the predefined schema. It is worthy to note, that 108 LLMs are known to inherit biases present in their training data. If not carefully managed, these biases 109 could lead to unfair or inaccurate information extraction, impacting decision-making processes. 110 Besides the information extraction task, LLMs can be used to suggest suitable Schema.org types for 111

a particular document. An example together with a prompt is shown in Appendix B1.

113 4 Length aspects

When focusing on the quality of IE performed by an LLM, several limitations that LLM presents in terms of the length of data to be extracted from must be considered. Each LLM has a maximal content limit it can process, both on the input and the output. The limit on the output is typically much more strict. When trying to use the maximal possible input another issue may appear – the *Lost in the middle* phenomenon [8] says that the ability of LLMs to retrieve information from a long context declines and that the attention focuses on the beginning and the end of the context while it tends to attenuate information in the middle.

To demonstrate shortcomings arising from these limitations numerically we use *gpt-4-1106-preview* model.³ The model is limited by 4095 tokens on the output and by 128000 tokens on the input (context window limit). The following sections present two major LLM limitations we have to consider before performing IE, namely length restrictions in Section 4.1 and *Lost in the middle* problem in Section 4.2.

126 **4.1 Length restrictions**

- 127 Long data are difficult to process because of the restrictions posed by the maximum amount of:
- (O) output tokens: The restriction on output tokens means that there is some maximal length of
 data from which most entities can be efficiently extracted. If the length of the text exceeds
 this maximum, there would be no tokens for extra entities.
- (I) input tokens: Maximal size of context window (input) prohibits the extraction of data
 exceeding the specific token limit.

Another difficulty regarding the output is the tendency of LLMs to generate rather brief responses which do not use the allowed maximal number of tokens. This unwillingness of models can be circumvented by prompting. Even so, the limited number of output tokens is typically too low and prevents effective extraction from long texts.

With a more sophisticated approach, the restriction (O) becomes irrelevant and only the restriction (I) 137 will apply. The issue imposed by (O) is overcome by splitting the source document into smaller pieces 138 which are extracted independently. A significant drawback is that the extracted information can be 139 easily duplicated – extracted independently from multiple text pieces. Iterating the calls to the LLM 140 with instruction to continue with already started extraction, i.e., continuing with the extraction in a 141 single thread, helps to extract more information and avoid duplication. As we insist on continuation, 142 more and more information is added and the extraction is more thorough, at least to some point – this 143 will be addressed in detail in Section 5.1. Further, a lower number of duplicates is found due to the 144 extraction history, i.e., all information extracted until present, which is kept within the thread. 145

The combination of both improvements – text splitting and iterated calls, has proven itself to perform the best. We split the document into distinct text pieces which we extract sequentially. Extraction from each text piece is carried out by several iterated LLM calls while taking into account the extraction history from previously extracted text pieces. Once the sum of the lengths of the text pieces and the extraction history exceeds the context window limit, i.e., restriction (I) applies, a new independent extraction starts. A single structured output, per document or once (I) is applied, is created by appending all entities identified from each text piece.

153 **4.2 Lost in the middle**

In the case of long documents, whose extraction consumes almost the whole context window, 154 LLMs are giving more inconsistent results and we can observe a presence of the Lost in the middle 155 phenomenon, see [8]. We extract information from several long documents from our business case 156 which are each split into 15 pieces and its processing consumes almost the whole context window. 157 We add the sixteenth piece identical to one of the fifteen that are already extracted and measure a 158 redundancy score, for details see Appendix A. Each column of Table 1 then states the redundancy of 159 the newly extracted information with the information that was already extracted from the same piece 160 of the text before. The table presents mean values per four distinct documents. We can notice that 161

³https://platform.openai.com/docs/models/overview

for the parts 'in the middle' the proportion of redundantly extracted entities (entities with the same 'name' attribute) is higher than for those at the beginning and the end.

Table 1: Are we lost in the middle? After finishing the extraction of a whole document (consisting of fifteen pieces), we re-extract the information from each of its pieces. Columns 1-15 then compare the re-extracted information with the information that was extracted from the same piece of the text before. The pieces in the middle of the document contain more duplicated entities then those at the beginning and the end.

| redunda | part ncy (key = | - 'name') | 1 0 | 2 0 | 3 0.2266 | 4 0.1150 | 5 0.1482 | 6 0.3816 |
|---------|--------------------|-----------|--------|--------|-------------|-------------|-------------|-------------|
| 7 | 8 | 9 | 10 | 11 | | 13 | 14 | 15 |
| 0.3334 | 0.4643 | 0.7398 | 0.5152 | 0.6672 | | 0.3820 | 0.4473 | 0.4086 |

164 5 Quality of extraction

165 Once the information is extracted from data into a structured form defined by the chosen schema, 166 e.g., Figure 1, the quality of such extraction is important to evaluate. In practice, it is very rare to be equipped with ground truth and its human generation requires vast expertise in the scope of data and a 167 ridiculous amount of time. Therefore we adopt methods from [4]. They examine semantic similarity, 168 relevance, redundancy, and bias and compound these into a single score called SUSWIR, for details 169 see Appendix A. The score and its subparts are very useful when comparing distinct extractions 170 among themselves, e.g., we can use it to find an optimal number of iterated LLM calls. Unfortunately, 171 the score does not represent an absolute way of evaluation. It does not provide a complete insight into 172 the task – some information (= entities) can be missing, misclassified or their properties not filled 173 in correctly. To come up with a robust and general solution we generalize the NIAH test, which is 174 commonly used to measure the ability of LLMs to process long documents, cf. [6]. 175

176 5.1 Iterated LLM calls

177 Since the first LLM extraction is typically not exhaustive, iterating the extraction process helps with 178 the completeness of extraction. To improve the quality of extraction, we ask LLM to process the document again and search for other entities which were not extracted vet. A question arises: What is 179 the optimal number of iterations? It is desirable to stop when additional LLM call will return no or 180 only a few new entities. The answer however depends heavily on the text being extracted and on the 181 chosen schema. Below, we present a small comparative study regarding the contribution of iterated 182 extraction to its quality. We interpret the extracted structured data, e.g., Figure 3, as a summary of 183 the original text document. To measure the quality of the summary we adopt the scores from [4] (a 184 convex combination of these scores creates the overall SUSWIR metric), namely *semantic similarity*, 185 relevance, and redundancy avoidance. We use a modified bias avoidance score from [4] and add two 186 new scores, *relevance spread*, and *incompleteness score*. See Appendix A for more details. 187

Consider document which length is approximately 12k chars. Table 2 compares the content of the 188 document with extracted information created iteratively by succeeding LLM calls. Each iteration 189 enriches the extracted information, but the benefit decreases. From the third iteration, i.e., after 190 four LLM calls, the majority of scores in Table 2 are either getting worse or stagnating (the arrows 191 following the score name indicate the direction in which the score improves). It is obvious that shorter 192 and longer text will require less or more iterations to extract majority of information without reducing 193 its semantic and factually relevant meaning, respectively. Further, the risk that the LLM will suffer 194 from hallucinations increases as we observe a growth of bias. In the rest of the paper we use three 195 iterations to extract documents of approximate length 12k chars within all extractions (if not stated 196 otherwise). 197

198 5.2 Test the quality

This section introduces a robust and versatile score to objectively measure the quality of IE. Assuming the structure is imposed by some schema, see Section 3.1, we would like to measure the IE quality as

Table 2: Quality of extraction depends on a number of calls to LLM. The first iterated call is the most beneficial one. From some point (bold) the scores stagnate or even deteriorate. All scores have values between 0 and 1, the arrows indicate whether lower (\downarrow) or higher (\uparrow) values are desired.

| # iterations | 0 | 1 | 2 | 3 | 4 | 5 |
|---|--------|--------|--------|--------|--------|--------|
| semantic similarity \uparrow | 0.5416 | 0.6316 | 0.6899 | 0.7572 | 0.7540 | 0.7685 |
| relevance ↑ | 0.3409 | 0.4396 | 0.4449 | 0.4746 | 0.4522 | 0.4445 |
| relevance spread \downarrow | 0.3364 | 0.2493 | 0.2350 | 0.1445 | 0.1428 | 0.1368 |
| redundancy avoidance (0.2) \uparrow | 0.7727 | 0.8670 | 0.8810 | 0.9257 | 0.9251 | 0.9307 |
| redundancy avoidance (0.1) \uparrow | 0.4697 | 0.5936 | 0.6854 | 0.8002 | 0.7972 | 0.8119 |
| redundancy avoidance | 0.8182 | 0.9163 | 0.9422 | 0.9650 | 0.9699 | 0.9726 |
| $(0.5, \text{key='name'}) \uparrow$ | | | | | | |
| bias avoidance \uparrow | 0.5614 | 0.5515 | 0.4925 | 0.4559 | 0.4447 | 0.4247 |
| incompleteness ↓ | 0. | 0.5862 | 0.6735 | 0.4217 | 0.5413 | 0.4615 |

a portion of successfully extracted entities, i.e., the accuracy of name entity recognition (NER) task
 taking into account even the context captured by entity properties. Unfortunately, such an experiment
 unfeasible without labeled data. As a consequence, it is unfeasible in many specialized tasks
 because of the absence of suitable labeled data unseen by LLM models. This can be the case with
 very recent datasets as well as proprietary datasets. To overcome this issue we use inspiration by
 NIAH test to build up an automatic and general procedure to access the quality of IE tasks.

207 **5.2.1 Needles**

A 'needle' in our context represents an entity. It is created according to the chosen schema, i.e., 208 a list of types we want to extract from the document. We use an LLM to generate a short paragraph 209 introducing a new original (not appearing in the document) entity, but still relevant to the scope of the 210 document, for an example see Figure 2, and for more details on generation process see Appendix B2. 211 This artificial paragraph, the needle, is then placed into the document body at random (taking into 212 the account natural units within the text as sentences, paragraphs, etc. if applicable). Moreover, 213 the needle is accompanied with several properties, namely we assign to the needle a name, short 214 description and keywords, see Figure 2. This additional properties are assigned to the needle by the 215 LLM. 216

217 5.2.2 Multiple infused needle extraction accuracy

To measure the quality of extraction we propose a *multiple infused needle extraction accuracy* 218 (MINEA) score. Its computation combines the approach of NIAH evaluation and NER task. We 219 scatter several needles at random over the text document body (such that the inserted needles fill 10 220 to 30% of the enriched text) and measure how many of them were successfully extracted. Since we 221 know what exactly was inserted, we know what should be extracted. Then we can objectively measure 222 the quality of extraction on these new entities and moreover, we can compare extracted information 223 from the document with and without needles. Table 3 shows extraction accuracy - MINEA score 224 - total and per schema type - measured on a vast corpus of business documents with predefined 225 schema consisting of types 'BioChemEntity', 'Event', 'Insight', 'Legislation', 'MedicalCondition', 226 'OpportunityArea', 'Person', 'Product', 'Project', 'Substance' and 'Thing'. 227

228 5.2.3 Identification of needles

Matching the generated needles with extracted entities imposes a challenge and mostly depends on the formulation of needles. If the needles are too complex or too vague, the straightforward identification changes into a serious problem. For this reason, we equip the needles with additional properties which are then used to compare the needles with extracted entities and to decide whether the needles were extracted successfully or not.

²³⁴ We present several alternative ways how to measure whether the extraction of a needle is successful:

- n an entity with a name perfectly matching the needle name is found;
- ns the needle name is found among the extracted information;

```
[
 {
"@type": "Event",
"The A
  "needle": "The AI Clan Meeting on Thursday aims to bring together a diverse team for
       collaboration and knowledge sharing. It is a hybrid event, with team members
       gathering in person at the office while also connecting online via video
       conferencing. The meeting will feature discussions on recent AI projects,
      updates on upcoming initiatives, and collaborative brainstorming sessions.",
  "name": "AI Clan Meeting"
  "description": "The aim of hybrid event AI Clan Meeting happening on Thursday is to
       foster collaboration and engagement among the team. The agenda includes project
  discussions, updates on upcoming initiatives, and brainstorming sessions.
"keywords": "AI (9), AI projects (9), project updates (7), team collaboration (6),
       knowledge sharing (7), hybrid event (4)",
 },
{
   "@type": "Product",
  "needle": "Graph Index (GRIX) is a cutting-edge retrieval-augmented generation model
       that is based on a knowledge graph. A graph representation of the knowledge base
       enhances effectiveness and ability to answer complex user queries. It is end-to-
       end solution for question-answering task dealing with the knowledge graph construction from and the retrieval of a relevant information from it.",
  "name": "Graph Index"
  "description": "GRIX is an innovative retrieval-augmented generation model based on
       a knowledge graph. A great focus is laid on proper extraction of information from
       data, its composition into the graph and retrieval of a relevant subgraph.'
  "keywords": "retrieval-augmented generation (9), knowledge graph (8), information
       extraction (6), product innovation (7), graph index (8), question-answering (8)"
}
]
```

Figure 2: Toy example: two needles, highlighted by blue color, accompanied by additional information described by 'name', 'description', and 'keywords'.

Table 3: Quality of extraction – MINEA score – total and per schema type. Entity types are grouped into five classes - 1. three most frequent schema.org types in the documents; 2. med-bio-chem entities, somewhat interchangeable types; 3. best distinguishable types; 4. custom (non Schema.org) types; 5. Schema.org types related to documents, but not stated in the chosen schema. Note: an entity is assumed to be extracted if it is contained within the extracted information - often its type can be misclassified (Project-Product-OpportunityArea, Substance-Thing-BioChemEntity) or sometimes it can be mentioned indirectly (Organization is related to a Person by property 'works for').

| class | entity type | extraction accuracy | # entities used for evaluation |
|-------|------------------|---------------------|--------------------------------|
| | Person | 0.884 | 69 |
| 1 | Project | 0.702 | 47 |
| | Product | 0.750 | 52 |
| | Substance | 0.822 | 45 |
| 2 | Thing | 0.739 | 46 |
| | BioChemEntity | 0.674 | 43 |
| | MedicalCondition | 0.636 | 44 |
| 3 | Legislation | 0.942 | 52 |
| | Event | 0.915 | 47 |
| 4 | OpportunityArea | 0.671 | 73 |
| | Insight | 0.747 | 91 |
| 5 | Organization | 0.907 | 43 |
| | Place | 0.767 | 43 |
| | overall | 0.780 | 695 |

237 238 239

k an entity with some number of keywords perfectly matching the needle keywords is found, the number is determined by the threshold parameter determining the percentage of keywords to be matched;

```
Γ
   {
    "@type": "Insight",
    "name": "Information exctraction tested by Needle in a Haystack test",
   },
   {
    "@type": "Event"
    "name": "AI Meeting"
    "description": "A hybrid event bringing together a diverse team for collaboration and
        knowledge sharing.",
    "keywords": "AI Clan Meeting (9), collaboration (8), knowledge sharing (8), hybrid
event (7), team gathering (7), video conferencing (6)"
  },
   {
    "@type": "Product",
    "name": "GRIX",
                   "Cutting-edge retrieval-augmented generation model based on a knowledge
    "description":
       graph".
    "keywords": "GRIX (10), retrieval-augmented generation (9), knowledge graph (10),
        question-answering (8), graph construction (6), information extraction (7)"
   }
 1
```

Figure 3: Toy example: extracted information from the data infused by needles from Figure 2.

²⁴⁰ **IIm** an entity matching the needle according to LLM is found.

Table 4: Toy example: fulfillment of the conditions. The text enriched by two needles from Figure 2 was extracted into the form shown in Figure 3.

| entity type | сс | condition for needle identification | | | | | | |
|-------------|----|-------------------------------------|------|------|------|-----|--|--|
| | n | ns | k0.5 | k0.6 | k0.7 | llm | | |
| Event | 0 | 1 | 1 | 0 | 0 | 1 | | |
| Product | 0 | 0 | 1 | 1 | 0 | 1 | | |

Note that other conditions can be constructed, e.g., based on the short description instead of keywords, 241 242 etc. Table 4 shows whether the conditions are fulfilled in the example illustrated by Figures 2 and 3. Namely, the condition **n** is not satisfied ('AI Clan Meeting' \neq 'AI Meeting', 'Graph Index' \neq 243 'GRIX'). Condition **ns** is satisfied only for needle representing an entity of type 'Event' ('AI Clan 244 Meeting' can be found in the extracted information). There are three keywords out of the six assigned 245 to the needle representing the entity of type 'Event' which match the keywords of an extracted entity, 246 hence **k0.5** is, and **k0.6**, **k0.7** are not satisfied (there is an entity within the extracted information 247 with 50% of keywords being the same as the keywords of the needle). In the case of the second 248 needle, there are four such keywords, therefore **k0.5** and **k0.6** are satisfied. Finally, both needles are 249 identified within the extracted information by an LLM. 250

Table 5 shows scores (ratios of successfully extracted entities) based on the above criteria in the case 251 of our business documents. The types of inserted needles are 'BioChemEntity', 'Country', 'Event', 252 'Insight', 'Legislation', 'Person', 'Product', 'Project' and 'Substance'. Matching the needle and 253 entity name usually does not perform well if the name is prone to modification (e.g., person name 254 with and without title), or if the entity is easy to be misclassified (an entity of type 'Country' was 255 often extracted as 'Place' whose name did not match the country name). Searching for a needle name 256 in all extracted information gives very accurate results if the entities are well characterized by their 257 name (compare for example types 'Person' and 'Legislation' with type 'Insight' where the name is 258 not a natural attribute). Matching the needle and entity keywords depends on the threshold parameter 259 - with a lower proportion of keywords that have to match the score value increases and the reliability 260 of the entity identification decreases. An LLM performs well the entity identification and it is an 261 important criterion in the case of more creative types such as 'Insight'. Finally, the MINEA score for 262 each type is taken as the maximum of the scores (the values are highlighted). 263

Table 5: The decision about the success of needle extraction can be made based on several criteria: comparing the corresponding needle and entity properties (columns **n** and **k0.5-k0.7** compare name and keywords, respectively), full-text search (column **ns** search for the needle name in extracted information), comparison of needles and entities using LLM (column **llm**).

| entity type | | condition for needle identification | | | | | # entities used |
|---------------|-------|-------------------------------------|-------|-------|-------|-------|-----------------|
| | n | ns | k0.5 | k0.6 | k0.7 | llm | for evaluation |
| Person | 0.594 | 0.884 | 0.652 | 0.362 | 0.232 | 0.826 | 69 |
| Project | 0.170 | 0.702 | 0.638 | 0.234 | 0.085 | 0.681 | 47 |
| Product | 0.596 | 0.712 | 0.462 | 0.192 | 0.135 | 0.750 | 52 |
| Country | 0 | 0.765 | 0.412 | 0.294 | 0.059 | 0.471 | 17 |
| Legislation | 0.635 | 0.942 | 0.365 | 0.269 | 0.096 | 0.942 | 52 |
| Event | 0.830 | 0.851 | 0.638 | 0.511 | 0.149 | 0.915 | 47 |
| Insight | 0.176 | 0.187 | 0.714 | 0.418 | 0.088 | 0.747 | 91 |
| BioChemEntity | 0.116 | 0.605 | 0.651 | 0.581 | 0.488 | 0.674 | 43 |
| Substance | 0.289 | 0.578 | 0.822 | 0.644 | 0.222 | 0.800 | 45 |

264 5.2.4 Model comparison

MINEA score can be used to compare the performance of distinct LLMs, see Table 6. A corpus 265 of documents is infused by needles representing entities whose types match the schema introduced 266 in Section 5.2.2. Three OpenAI LLMs⁴ are used to extract a relevant information under the same 267 setting (the same model parameters such as temperature, the same number of iterations, the same 268 prompting, etc.). Model gpt-3.5-turbo is outperformed by gpt-4-turbo by almost 15% and gpt-4-turbo 269 is outperformed by gpt-40 model by another 12%. Note that the achieved accuracy is lower than 270 presented in Table 3, since only one iteration instead of three was performed in order to reduce the 271 computational time. 272

 Table 6:
 LLMs comparison using MINEA score.

| model | gpt-3.5-turbo | gpt-4-turbo | gpt-40 |
|-------|---------------|-------------|----------|
| MINEA | 0.449198 | 0.593583 | 0.716578 |

273 Conclusions

In this paper, we focused on quality evaluation of information extraction (IE) performed by large 274 language models (LLMs). First, we delved into the technical limitations of LLMs complicating the 275 extraction of information from a long context. To extract reasonable information from data it is 276 needed to take into the account features such as context window limits, iterated extractions, extraction 277 history recording and Lost in the middle phenomenon. Once the extraction is performed, assessing its 278 quality is essential. However in many customized tasks, a truly objective method is missing, because 279 of the lack of labeled data fitting the scope of the application. The versatile method presented in this 280 paper overcomes the issue by adjustment of the data by insertion of an artificial information, a needle, 281 into it. The artificial information created to this purpose is application and data-specific, but the 282 method itself is applicable generally across the field of IE. By controlling the generation process of 283 the needles, we created a synthetic ground truth that enables us to absolutely measure the extraction 284 quality even when no labeled data is available. We introduced a MINEA score to measure the quality 285 of extraction. The key part is a decision rule on whether a needle was successfully extracted or not. 286 MINEA possibly combines several decision rules into one final score. Our empirical analysis of the 287 MINEA score on a specialized dataset demonstrated its utility for evaluation of LLM-based IE tasks 288 when ground truth is unavailable. 289

⁴https://platform.openai.com/docs/models

290 **References**

- [1] Kiran Adnan and Rehan Akbar. Limitations of information extraction methods and techniques for
 heterogeneous unstructured big data. *International Journal of Engineering Business Management*,
 11:1847979019890771, 2019.
- [2] Satanjeev Banerjee and Alon Lavie. METEOR: An automatic metric for MT evaluation with
 improved correlation with human judgments. In *Proceedings of the acl workshop on intrinsic and extrinsic evaluation measures for machine translation and/or summarization*, pages 65–72,
 2005.
- [3] Matthew Edgar. Schema and structured data markup. In *Tech SEO Guide: A Reference Guide for Developers and Marketers Involved in Technical SEO*, pages 67–78. Springer, 2023.
- [4] Abdullah Al Foysal and Ronald Böck. Who Needs External References?—Text Summarization
 Evaluation Using Original Documents. *AI*, 4(4):970–995, 2023.
- [5] Neil Jethani, Simon Jones, Nicholas Genes, Vincent J Major, Ian S Jaffe, Anthony B Cardillo,
 Noah Heilenbach, Nadia Fazal Ali, Luke J Bonanni, Andrew J Clayburn, et al. Evaluating
 ChatGPT in Information Extraction: A Case Study of Extracting Cognitive Exam Dates and
 Scores. 2023.
- [6] Yuri Kuratov, Aydar Bulatov, Petr Anokhin, Dmitry Sorokin, Artyom Sorokin, and Mikhail
 Burtsev. In Search of Needles in a 11M Haystack: Recurrent Memory Finds What LLMs Miss.
 arXiv preprint arXiv:2402.10790v2, 2024.
- [7] Dong-Ho Lee, Jay Pujara, Mohit Sewak, Ryen W White, and Sujay Kumar Jauhar. Making large
 language models better data creators. *arXiv preprint arXiv:2310.20111*, 2023.
- [8] Nelson F Liu, Kevin Lin, John Hewitt, Ashwin Paranjape, Michele Bevilacqua, Fabio Petroni,
 and Percy Liang. Lost in the middle: How language models use long contexts. *arXiv preprint arXiv:2307.03172*, 2023.
- [9] Derong Xu, Wei Chen, Wenjun Peng, Chao Zhang, Tong Xu, Xiangyu Zhao, Xian Wu, Yefeng
 Zheng, and Enhong Chen. Large language models for generative information extraction: A
 survey. *arXiv preprint arXiv:2312.17617*, 2023.

317 Appendix A

To measure the quality of the summary we adopt the methods from [4]: *semantic similarity* combines latent semantic similarity and cosine similarity; *relevance* is measured using METEOR score, see [2], without chunk penalty; *redundancy avoidance* compares extracted entities among themselves using a threshold parameter – entities with a higher cosine similarity are assumed to be redundant; redundancy avoidance can be focused on a single particular property of entities (we use 'name' as this pivotal property).

We modify the *bias avoidance* score from [4] to be $J^*(A, B) = \frac{|A \cap B|}{|B|}$, where *A* represents the entities in the original text document and we normalize by a number of entities that were extracted, |B|. The score controls how much information in the structured file is not present in the original text, i.e., a potential hallucination of an LLM.

We add two new scores: the *relevance spread* is the standard deviation of relevance over the text pieces to which the document is split and normalized by the mean value, its higher values indicate that the extraction from distinct text pieces is unbalanced; the *incompleteness score* just measures the proportion of entities with incomplete information (at least one property value missing or unfilled), e.g., the entity 'AI Enthusiast' in Figure 1 has an unknown 'birthDate'.

333 Appendix B

Except for the IE task, LLMs are used in several subtasks within the paper, namely to determine schema types appearing in the document, to create a suitable needles fitting contextually to the document and to identify whether a needle was extracted or not. In the following, we provide the reader with prompts and examples of these subtasks.

338 **B1 Discovering a schema**

Figure 4 shows a prompt to obtain the Schema.org types from the attached text – Wikipedia article about IE.⁵ An LLM is asked to assign relevance to the types to distinguish the most important ones.

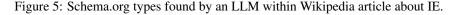
Figure 5 shows the entity types that were deduced from the text, together with their relevance and reasoning for why they were chosen. The most relevant types are those directly mentioned – 'Article', as the webpage content itself is represented as an article, 'SoftwareApplication', and 'WebSite' (all with maximal relevance). The least relevant identified types are generic – 'Thing', as a parent type of many directly mentioned types, and 'LearningResource', as a categorization of the article style.

| SYSTEM | |
|--|--|
| You are skilled proof reader. Analyze the attached text and entities it contains carefully and determine all schema.org types being present. | |
| Output just a list of all TYPES found. Provide rating of relevance (1-10) for each identified type and url to website where the type is defined. | |
| *Attached text: https://en.wikipedia.org/wiki/Information_extraction | |

Figure 4: Prompt to determine a possible suitable schema from a given text – Wikipedia article about IE.

ASSISTANT Based on the attached text, the following schema.org types can be identified:

1. **Article**: This type is clearly represented through the mention of "Information extraction" as an article within the content. Relevance: 10 - [Article](https://schema.org/Article) 2. **EducationalOrganization** (due to references to universities and research institutes): Given the academic citations and references, organizations involved in education and research are implicitly present. Relevance: 8 -[EducationalOrganization](https://schema.org/EducationalOrganization) 3. **Person** (authors of the referenced works): The references section mentions multiple authors, indicating their presence in a structured form. Relevance: 9 - [Person](https://schema.org/Person) 4. **Event** (Message Understanding Conferences): The text mentions specific conferences, which are events. Relevance: 8 - [Event](https://schema.org/Event) 5. **SoftwareApplication** (GATE, Apache OpenNLP, DBpedia Spotlight, etc.): The section on free or open-source software indicates these entities. Relevance: 10 - [SoftwareApplication](https://schema.org/SoftwareApplication) 6. **CreativeWork**: Given the references to published papers, books, and academic outputs, this type is represented. Relevance: 9 - [CreativeWork](https://schema.org/CreativeWork) 7. **WebSite** (Wikipedia itself and mentioned web services like OpenCalais): The context of the text being on Wikipedia and references to other web services. Relevance: 10 - [WebSite](https://schema.org/WebSite) 8. **Organization** (in references to DARPA, Reuters, etc.): Various organizations are mentioned throughout the text Relevance: 8 - [Organization](https://schema.org/Organization) 9. **Thing** (as a generic type for topics covered, e.g., "natural language processing", "information extraction"): Given that 'Thing' is a top-level Schema.org type under which many topics can be categorized. Relevance: 7 - [Thing] (https://schema.org/Thing) 10. **LearningResource** (implied by the educational and instructional nature of the text): Considering the text's aim to inform and educate about Information Extraction. Relevance: 7 - [LearningResource] (https://schema.org/LearningResource) 11. **DataFeed** (considering the reference to news data feeds for IE applications): Specifically mentioned in the context of information extraction from news feeds. Relevance: 8 - [DataFeed](https://schema.org/DataFeed) 12. **ProgrammingLanguage** (when referring to usage of Java, Python, etc. in software and toolkits): Implicit through mentioning specific programming languages used in the context of NLP and IE. Relevance: 9 -[ProgrammingLanguage](https://schema.org/ProgrammingLanguage)



⁵https://en.wikipedia.org/wiki/Information_extraction

346 **B2** Creating needles

A needle, i.e., a text paragraph fitting thematically to the document, but being new and unique to it, is generated by an LLM using the prompt in Figure 6. The prompt specifies the type of entity that the

needle should represent. Multiple needles of the same type can be obtained easily within a single

350 LLM call.

Figure 7 shows ten needles representing the entities of type 'Person' generated based on a Wikipedia

- article about IE. In the next step properties such as a name, description and keywords can be generated
- 353 by an LLM.

| SYSTEM | |
|--|---|
| You have be readers. | een created to test attentiveness and mindfulness of proof |
| document t completely to it THEMA | to suggest a piece of information to be inserted to a hey will carefully examine. The information needs to be NEW and UNIQUE, not appearing in the document, but FIT VTICALLY. A proof reader will pass the test only, if he discover the newly added text. |
| particular e about is 'Pe | eaders are tested repeatedly. Each test focuses on a ntity type. The type you are asked to create a piece of text rson'. Stick only to this type!!! You will be severely or each created entity classified to a wrong schema.org |
| DIFFERENT proof reade always at th | self sufficient information pieces to be inserted. Use format and style for each piece in order to confuse the rs. Avoid patterns! E.g., placing the name of the entity e beginning is a typical pattern behavior. DO NOT the name of the entity. |
| *Type: Pers | on |
| *Output: tex | cts to insert |
| Doodinoni | to insert to: ikipedia.org/wiki/Information_extraction |
| | |

Figure 6: Prompt to generate needles. Given a Wikipedia article about IE, the LLM is asked to think out 10 relevant persons.

ASSISTANT 1. Arnold Jennings is a renowned expert in the field of Information Extraction and has published several groundbreaking research papers on the topic in prestigious journals.

2. In 1995, Professor Elena Martinez pioneered a new approach to Named Entity Recognition, revolutionizing the way entities are identified in natural language processing.

3. Dr. James Carter, a leading figure in the field of Information Extraction, will be delivering a keynote speech at the upcoming International Conference on NLP and Machine Learning.

4. Vanessa Rodriguez, a rising star in the world of Information Extraction, recently developed a cutting-edge algorithm that significantly enhances the accuracy of event extraction from texts.

5. The collaboration between Dr. Michael Nguyen and Dr. Sophia Lee resulted in the creation of a powerful tool for Relationship Extraction, which has been widely acclaimed in the NLP community.

6. Professor Samuel Brown received the prestigious Excellence in Information Extraction Award for his groundbreaking research on Coreference Resolution in complex textual data.

7. In her latest research project, Dr. Lily Martinez successfully applied Information Extraction techniques to analyze and extract valuable insights from a large dataset of social media posts.

8. The innovative work of Dr. Oliver Walker in Template Filling has paved the way for more efficient extraction of structured information from unstructured documents.

9. Christina Lee, a leading researcher in the field of Named Entity Recognition, has developed a state-of-the-art system that can accurately detect and classify various types of entities in textual data.

10. Professor Nathan Turner's contributions to Semi-structured Information Extraction have been instrumental in advancing the capabilities of automated systems in processing and interpreting diverse types of data sources.

Figure 7: Needles generated by an LLM and representing ten entities of type 'Person'.

354

B3 Identifying needles 355

356

- The quality of extraction is evaluated based on the proportion of successfully extracted needles. An LLM can be used to decide whether the needle was extracted or not using the prompt presented in 357
- Figure 8. 358

| SYSTEM |
|--|
| You are an expert in named entity recognition (NER). |
| Given a list of entities your goal is to decide for each of them whether |
| in the attached structured list is an element directly matching to it (the |
| element and the entity have to be one to one correspondence). Each |
| entity has a type, name, description and keywords. Take care about |
| these attributes when making the decision. |
| For EACH entity in the list you are supposed state a simple YES or NO |
| and in the case of positive answer output the matching element, in the |
| case of negative answer output by meaning the closest element. |
| Rule: number of decisions = number of entities |
| You will be strictly penalized for doing decisions upon entities which |
| are not mentioned!!! |
| |
| List of entities: {list_ents} |
| |
| Attached structured list: {extracted} |

Figure 8: Prompt to identify whether the needles were extracted or not.

NeurIPS Paper Checklist

| 1. | Claims |
|----|--------|
| | 1. |

Question: Do the main claims made in the abstract and introduction accurately reflect the paper's contributions and scope?

363 Answer: [Yes]

Justification: The abstract and introduction clearly state the development of an automatic framework to assess the quality of information extraction (IE), which is the main contribution of the paper. This is supported by the introduction of the MINEA score and the discussion on handling input/output size limitations of large language models (LLMs).

368 Guidelines:

- The answer NA means that the abstract and introduction do not include the claims made in the paper.
- The abstract and/or introduction should clearly state the claims made, including the contributions made in the paper and important assumptions and limitations. A No or NA answer to this question will not be perceived well by the reviewers.
- The claims made should match theoretical and experimental results, and reflect how much the results can be expected to generalize to other settings.
 - It is fine to include aspirational goals as motivation as long as it is clear that these goals are not attained by the paper.

2. Limitations

Question: Does the paper discuss the limitations of the work performed by the authors?

Answer: [Yes]

Justification: The paper discusses the limitations related to the complexity or vagueness of the needles, dependence on the chosen schema and criteria for needle identification (Section 5). Further the paper focuses on limitations of LMMs in IE tasks such as input/output size constraints, lost in the middle phenomenon, bias and hallucinations (Section 4).

- Guidelines:
 - The answer NA means that the paper has no limitation while the answer No means that the paper has limitations, but those are not discussed in the paper.
 - The authors are encouraged to create a separate "Limitations" section in their paper.
 - The paper should point out any strong assumptions and how robust the results are to violations of these assumptions (e.g., independence assumptions, noiseless settings, model well-specification, asymptotic approximations only holding locally). The authors should reflect on how these assumptions might be violated in practice and what the implications would be.
 - The authors should reflect on the scope of the claims made, e.g., if the approach was only tested on a few datasets or with a few runs. In general, empirical results often depend on implicit assumptions, which should be articulated.
 - The authors should reflect on the factors that influence the performance of the approach. For example, a facial recognition algorithm may perform poorly when image resolution is low or images are taken in low lighting. Or a speech-to-text system might not be used reliably to provide closed captions for online lectures because it fails to handle technical jargon.
 - The authors should discuss the computational efficiency of the proposed algorithms and how they scale with dataset size.
 - If applicable, the authors should discuss possible limitations of their approach to address problems of privacy and fairness.
- While the authors might fear that complete honesty about limitations might be used by reviewers as grounds for rejection, a worse outcome might be that reviewers discover limitations that aren't acknowledged in the paper. The authors should use their best judgment and recognize that individual actions in favor of transparency play an important role in developing norms that preserve the integrity of the community. Reviewers will be specifically instructed to not penalize honesty concerning limitations.

| 412 | 3. Theory Assumptions and Proofs |
|------------|--|
| 413 | Question: For each theoretical result, does the paper provide the full set of assumptions and |
| 414 | a complete (and correct) proof? |
| 415 | Answer: [NA] |
| 416 | Justification: The paper does not include theoretical results that require formal proofs. |
| 417 | Guidelines: |
| 418 | • The answer NA means that the paper does not include theoretical results. |
| 419 | • All the theorems, formulas, and proofs in the paper should be numbered and cross- |
| 420 | referenced. |
| 421 | • All assumptions should be clearly stated or referenced in the statement of any theorems. |
| 422 | • The proofs can either appear in the main paper or the supplemental material, but if |
| 423 | they appear in the supplemental material, the authors are encouraged to provide a short proof sketch to provide intuition. |
| | Inversely, any informal proof provided in the core of the paper should be complemented |
| 425 426 | by formal proofs provided in appendix or supplemental material. |
| 420 | Theorems and Lemmas that the proof relies upon should be properly referenced. |
| | |
| 428 | 4. Experimental Result Reproducibility |
| 429 | Question: Does the paper fully disclose all the information needed to reproduce the main ex- |
| 430 | perimental results of the paper to the extent that it affects the main claims and/or conclusions of the paper (regardless of whether the code and data are provided or pot)? |
| 431 | of the paper (regardless of whether the code and data are provided or not)? |
| 432 | Answer: [Yes] |
| 433 | Justification: The paper provides detailed descriptions of the experimental setup, including |
| 434 | the use of LLMs for IE and the creation of synthetic ground truth data. This is detailed in |
| 435 | Sections 3 and 5. |
| 436 | Guidelines: |
| 437 | The answer NA means that the paper does not include experiments. |
| 438 | • If the paper includes experiments, a No answer to this question will not be perceived |
| 439 | well by the reviewers: Making the paper reproducible is important, regardless of |
| 440 | whether the code and data are provided or not. |
| 441 | • If the contribution is a dataset and/or model, the authors should describe the steps taken |
| 442 | to make their results reproducible or verifiable. |
| 443 444 | • Depending on the contribution, reproducibility can be accomplished in various ways. For example, if the contribution is a novel architecture, describing the architecture fully |
| 445 | might suffice, or if the contribution is a specific model and empirical evaluation, it may |
| 446 | be necessary to either make it possible for others to replicate the model with the same |
| 447 | dataset, or provide access to the model. In general, releasing code and data is often |
| 448 | one good way to accomplish this, but reproducibility can also be provided via detailed |
| 449 | instructions for how to replicate the results, access to a hosted model (e.g., in the case |
| 450 | of a large language model), releasing of a model checkpoint, or other means that are |
| 451 | appropriate to the research performed. |
| 452 | • While NeurIPS does not require releasing code, the conference does require all submis- |
| 453 | sions to provide some reasonable avenue for reproducibility, which may depend on the |
| 454 | nature of the contribution. For example |
| 455 | (a) If the contribution is primarily a new algorithm, the paper should make it clear how to reproduce that algorithm |
| 456 | to reproduce that algorithm. (b) If the contribution is primarily a new model architecture, the paper should describe |
| 457 458 | (b) If the contribution is primarily a new model arcmeeture, the paper should describe the architecture clearly and fully. |
| 459 | (c) If the contribution is a new model (e.g., a large language model), then there should |
| 460 | either be a way to access this model for reproducing the results or a way to reproduce |
| 461 | the model (e.g., with an open-source dataset or instructions for how to construct |
| 462 | the dataset). |
| 463 | (d) We recognize that reproducibility may be tricky in some cases, in which case |
| 464 | authors are welcome to describe the particular way they provide for reproducibility. |
| | |

| 465 466 467 | In the case of closed-source models, it may be that access to the model is limited in some way (e.g., to registered users), but it should be possible for other researchers to have some path to reproducing or verifying the results. |
|-------------------|--|
| 468 | 5. Open access to data and code |
| 469 | Question: Does the paper provide open access to the data and code, with sufficient instruc- |
| 470 | tions to faithfully reproduce the main experimental results, as described in supplemental |
| 471 | material? |
| 472 | Answer: [No] |
| 473 | Justification: The paper does not provide open access to the data and code due to the |
| 474 475 | proprietary nature of the business documents used in the experiments. However, it provides detailed instructions on how to replicate the methodology. |
| 476 | Guidelines: |
| 477 | • The answer NA means that paper does not include experiments requiring code. |
| 478 | • Please see the NeurIPS code and data submission guidelines (https://nips.cc/ |
| 479 | public/guides/CodeSubmissionPolicy) for more details. |
| 480 | • While we encourage the release of code and data, we understand that this might not be possible, so "No" is an acceptable answer. Papers cannot be rejected simply for not |
| 481 482 | including code, unless this is central to the contribution (e.g., for a new open-source |
| 483 | benchmark). |
| 484 | • The instructions should contain the exact command and environment needed to run to |
| 485 | reproduce the results. See the NeurIPS code and data submission guidelines (https: |
| 486 | <pre>//nips.cc/public/guides/CodeSubmissionPolicy) for more details.</pre> |
| 487 | • The authors should provide instructions on data access and preparation, including how to access the raw data, preprocessed data, intermediate data, and generated data, etc. |
| 488 | • The authors should provide scripts to reproduce all experimental results for the new |
| 489 490 | proposed method and baselines. If only a subset of experiments are reproducible, they |
| 491 | should state which ones are omitted from the script and why. |
| 492 | • At submission time, to preserve anonymity, the authors should release anonymized |
| 493 | versions (if applicable). |
| 494 495 | • Providing as much information as possible in supplemental material (appended to the paper) is recommended, but including URLs to data and code is permitted. |
| 496 | 6. Experimental Setting/Details |
| 497 | Question: Does the paper specify all the training and test details (e.g., data splits, hyper- |
| 498 | parameters, how they were chosen, type of optimizer, etc.) necessary to understand the |
| 499 | results? |
| 500 | Answer: [Yes] |
| 501 | Justification: The paper specifies the use of LLMs, the schema used for structuring data, and |
| 502 | the process of generating needles for evaluation. These details are provided in Sections 3, 4 and 5. |
| 503 | |
| 504 | Guidelines: |
| 505 | • The answer NA means that the paper does not include experiments. |
| 506 | • The experimental setting should be presented in the core of the paper to a level of detail that is necessary to appreciate the results and make sense of them. |
| 507 | • The full details can be provided either with the code, in appendix, or as supplemental |
| 508 509 | material. |
| 510 | 7. Experiment Statistical Significance |
| 511 | Question: Does the paper report error bars suitably and correctly defined or other appropriate |
| 512 | information about the statistical significance of the experiments? |
| 513 | Answer: [No] |
| 514 | Justification: The paper does not include experiments that require statistical significance |
| 515 516 | testing or error bars. The experiments in Sections 4 and 5 present mean values of reasonably large samples. The experiments are not repeated, each of them is carried once on a set of |
| 516 | ange samples. The experiments are not repeated, each of them is carried once on a set of |

| 517 518 | distinct documents containing a large amount of entities. In Section 5, a vast set of unique needles (with repeating types) is used to infuse the documents. |
|------------|--|
| 519 | Guidelines: |
| 520 | • The answer NA means that the paper does not include experiments. |
| 521 | • The authors should answer "Yes" if the results are accompanied by error bars, confi- |
| 522 | dence intervals, or statistical significance tests, at least for the experiments that support |
| 523 | the main claims of the paper. |
| 524 | • The factors of variability that the error bars are capturing should be clearly stated (for |
| 525 | example, train/test split, initialization, random drawing of some parameter, or overall |
| 526 | run with given experimental conditions). |
| 527 | • The method for calculating the error bars should be explained (closed form formula, |
| 528 | call to a library function, bootstrap, etc.) |
| 529 | • The assumptions made should be given (e.g., Normally distributed errors). |
| | • It should be clear whether the error bar is the standard deviation or the standard error |
| 530 531 | of the mean. |
| | |
| 532 | • It is OK to report 1-sigma error bars, but one should state it. The authors should preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis |
| 533 534 | of Normality of errors is not verified. |
| 535 | • For asymmetric distributions, the authors should be careful not to show in tables or |
| 536 | figures symmetric error bars that would yield results that are out of range (e.g. negative |
| 537 | error rates). |
| 538 | • If error bars are reported in tables or plots, The authors should explain in the text how |
| 539 | they were calculated and reference the corresponding figures or tables in the text. |
| 540 | 8. Experiments Compute Resources |
| | Question: For each experiment, does the paper provide sufficient information on the com- |
| 541 542 | puter resources (type of compute workers, memory, time of execution) needed to reproduce |
| 543 | the experiments? |
| 544 | Answer: [No] |
| 545 | Justification: The paper does not provide detailed information on the compute resources used |
| 546 | for the experiments. The requirements such as time of execution are determined especially |
| 547 | by used LLMs. |
| 548 | Guidelines: |
| 549 | • The answer NA means that the paper does not include experiments. |
| | The paper should indicate the type of compute workers CPU or GPU, internal cluster, |
| 550 551 | or cloud provider, including relevant memory and storage. |
| 552 | • The paper should provide the amount of compute required for each of the individual |
| 553 | experimental runs as well as estimate the total compute. |
| 554 | • The paper should disclose whether the full research project required more compute |
| 555 | than the experiments reported in the paper (e.g., preliminary or failed experiments that |
| 556 | didn't make it into the paper). |
| 557 | 9. Code Of Ethics |
| 559 | Question: Does the research conducted in the paper conform, in every respect, with the |
| 558 | NeurIPS Code of Ethics https://neurips.cc/public/EthicsGuidelines? |
| 559 | |
| 560 | Answer: [Yes] |
| 561 | Justification: The research adheres to the NeurIPS Code of Ethics, ensuring that the methods |
| 562 | and data used do not violate ethical guidelines. The proprietary data used is handled with |
| 563 | confidentiality and integrity. |
| 564 | Guidelines: |
| 565 | • The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics. |
| 566 | • If the authors answer No, they should explain the special circumstances that require a |
| 567 | deviation from the Code of Ethics. |
| | |

• The authors should make sure to preserve anonymity (e.g., if there is a special consid-568 eration due to laws or regulations in their jurisdiction). 569 10. Broader Impacts 570 Question: Does the paper discuss both potential positive societal impacts and negative 571 societal impacts of the work performed? 572 573 Answer: [Yes] Justification: The paper is primarily concerned with the technical methodology, the intro-574 575 duction of the MINEA score, and the empirical analysis of the framework's performance. The potential positive impacts are mentioned in Introduction: by automating the quality 576 assessment of information extraction, the framework could reduce the need for manual 577 review by experts, saving time and resources and thus enhance the efficiency and accuracy 578 of information extraction from large volumes of unstructured data. The negative aspects of 579 using LLMs for IE tasks such as inherited bias and potential hallucinations are mentioned 580 especially in Sections 4.2 (Lost in the middle problem) and 5.1 (bias avoidance score). 581 Guidelines: 582 • The answer NA means that there is no societal impact of the work performed. 583 • If the authors answer NA or No, they should explain why their work has no societal 584 impact or why the paper does not address societal impact. 585 Examples of negative societal impacts include potential malicious or unintended uses 586 (e.g., disinformation, generating fake profiles, surveillance), fairness considerations 587 (e.g., deployment of technologies that could make decisions that unfairly impact specific 588 groups), privacy considerations, and security considerations. 589 · The conference expects that many papers will be foundational research and not tied 590 to particular applications, let alone deployments. However, if there is a direct path to 591 any negative applications, the authors should point it out. For example, it is legitimate 592 to point out that an improvement in the quality of generative models could be used to 593 generate deepfakes for disinformation. On the other hand, it is not needed to point out 594 that a generic algorithm for optimizing neural networks could enable people to train 595 models that generate Deepfakes faster. 596 • The authors should consider possible harms that could arise when the technology is 597 being used as intended and functioning correctly, harms that could arise when the 598 technology is being used as intended but gives incorrect results, and harms following 599 from (intentional or unintentional) misuse of the technology. 600 • If there are negative societal impacts, the authors could also discuss possible mitigation 601 strategies (e.g., gated release of models, providing defenses in addition to attacks, 602 mechanisms for monitoring misuse, mechanisms to monitor how a system learns from 603 feedback over time, improving the efficiency and accessibility of ML). 604 11. Safeguards 605 Question: Does the paper describe safeguards that have been put in place for responsible 606 release of data or models that have a high risk for misuse (e.g., pretrained language models, 607 image generators, or scraped datasets)? 608 Answer: [NA] 609 Justification: The paper does not release any data or models that pose a high risk for misuse. 610 Guidelines: 611 612 • The answer NA means that the paper poses no such risks. • Released models that have a high risk for misuse or dual-use should be released with 613 necessary safeguards to allow for controlled use of the model, for example by requiring 614 that users adhere to usage guidelines or restrictions to access the model or implementing 615 safety filters. 616 · Datasets that have been scraped from the Internet could pose safety risks. The authors 617 should describe how they avoided releasing unsafe images. 618 • We recognize that providing effective safeguards is challenging, and many papers do 619 not require this, but we encourage authors to take this into account and make a best 620 faith effort. 621

| 622 | 12. | Licenses for existing assets |
|------------|-----|--|
| 623 | | Question: Are the creators or original owners of assets (e.g., code, data, models), used in |
| 624 | | the paper, properly credited and are the license and terms of use explicitly mentioned and |
| 625 | | properly respected? |
| 626 | | Answer: [Yes] |
| 627 | | Justification: All existing models are properly referenced and credit to their creators is given. |
| 628 | | These are either LLMs or metrics such as SUSWIR and METEOR (Section 5 and Appendix |
| 629 | | A). |
| 630 | | Guidelines: |
| 631 | | • The answer NA means that the paper does not use existing assets. |
| 632 | | • The authors should cite the original paper that produced the code package or dataset. |
| 633 | | • The authors should state which version of the asset is used and, if possible, include a |
| 634 | | URL. |
| 635 | | • The name of the license (e.g., CC-BY 4.0) should be included for each asset. |
| 636 637 | | • For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided. |
| 638 | | • If assets are released, the license, copyright information, and terms of use in the |
| 639 | | package should be provided. For popular datasets, paperswithcode.com/datasets |
| 640 | | has curated licenses for some datasets. Their licensing guide can help determine the |
| 641 | | license of a dataset. |
| 642 | | • For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided. |
| 643 644 | | If this information is not available online, the authors are encouraged to reach out to |
| 644 645 | | the asset's creators. |
| 646 | 13. | New Assets |
| 647 | | Question: Are new assets introduced in the paper well documented and is the documentation |
| 648 | | provided alongside the assets? |
| 649 | | Answer: [NA] |
| 650 | | Justification: The paper does not introduce new assets that require documentation. |
| 651 | | Guidelines: |
| 652 | | • The answer NA means that the paper does not release new assets. |
| 653 | | • Researchers should communicate the details of the dataset/code/model as part of their |
| 654 | | submissions via structured templates. This includes details about training, license, |
| 655 | | limitations, etc. |
| 656 | | • The paper should discuss whether and how consent was obtained from people whose asset is used. |
| 657 | | At submission time, remember to anonymize your assets (if applicable). You can either |
| 658 659 | | create an anonymized URL or include an anonymized zip file. |
| 660 | 14. | Crowdsourcing and Research with Human Subjects |
| 661 | | Question: For crowdsourcing experiments and research with human subjects, does the paper |
| 662 | | include the full text of instructions given to participants and screenshots, if applicable, as |
| 663 | | well as details about compensation (if any)? |
| 664 | | Answer: [NA] |
| 665 | | Justification: The paper does not involve crowdsourcing or research with human subjects. |
| 666 | | Guidelines: |
| 667 | | • The answer NA means that the paper does not involve crowdsourcing nor research with |
| 668 | | human subjects. |
| 669 | | • Including this information in the supplemental material is fine, but if the main contribu- tion of the paper involves human subjects, then as much detail as possible should be |
| 670 671 | | tion of the paper involves human subjects, then as much detail as possible should be included in the main paper. |
| 071 | | menadou in the main puper. |

| 672 673 674 | • According to the NeurIPS Code of Ethics, workers involved in data collection, curation, or other labor should be paid at least the minimum wage in the country of the data collector. |
|--------------------------|--|
| 675 676 | 15. Institutional Review Board (IRB) Approvals or Equivalent for Research with Human Subjects |
| 677 678 679 680 | Question: Does the paper describe potential risks incurred by study participants, whether such risks were disclosed to the subjects, and whether Institutional Review Board (IRB) approvals (or an equivalent approval/review based on the requirements of your country or institution) were obtained? |
| 681 | Answer: [NA] |
| 682 683 | Justification: The paper does not involve research with human subjects that would require IRB approval. |
| 684 | Guidelines: |
| 685 686 | The answer NA means that the paper does not involve crowdsourcing nor research with human subjects. |
| 687 688 689 | • Depending on the country in which research is conducted, IRB approval (or equivalent) may be required for any human subjects research. If you obtained IRB approval, you should clearly state this in the paper. |
| 690 691 692 | • We recognize that the procedures for this may vary significantly between institutions and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the guidelines for their institution. |
| 693 694 | • For initial submissions, do not include any information that would break anonymity (if applicable), such as the institution conducting the review. |