

---

# Automated, LLM enabled extraction of synthesis details for reticular materials from scientific literature

---

**Viviane Torres da Silva**  
IBM Research  
vivianet@br.ibm.com

**Alexandre Rademaker**  
IBM Research  
alexrad@br.ibm.com

**Krystelle Lioni**  
IBM Research  
klioni@us.ibm.com

**Ronaldo Giro**  
IBM Research  
rgiro@br.ibm.com

**Geisa Lima**  
IBM Research  
Geisa.Lima@ibm.com

**Sandro Fiorini**  
IBM Research  
srfiorini@ibm.com

**Marcelo Archanjo**  
IBM Research  
marcelo.archanjo@ibm.com

**Breno W. Carvalho**  
IBM Research  
brenow@ibm.com

**Rodrigo Neumann**  
IBM Research  
rneumann@br.ibm.com

**Anaximandro Souza**  
IBM Research  
anaximandrosouza@ibm.com

**João Pedro Souza**  
Idiap Research Institute  
joao.gandarela@idiap.ch

**Gabriela de Valnísio**  
IBM Research  
gvalnísio@ibm.com

**Carmen Nilda Paz**  
IBM Research  
cpaz@br.ibm.com

**Renato Cerqueira**  
IBM Research  
rcerq@br.ibm.com

**Mathias Steiner**  
IBM Research  
mathiast@br.ibm.com

## Abstract

1 Automated knowledge extraction from scientific literature can potentially accelerate  
2 materials discovery. We have investigated an approach for extracting synthesis  
3 protocols for reticular materials from scientific literature using large language  
4 models (LLMs). To that end, we introduce a Knowledge Extraction Pipeline (KEP)  
5 that automatizes LLM-assisted paragraph classification and information extraction.  
6 By applying prompt engineering with in-context learning (ICL) to a set of open-  
7 source LLMs, we demonstrate that LLMs can retrieve chemical information from  
8 PDF documents, without the need for fine-tuning or training and at a reduced risk  
9 of hallucination. By comparing the performance of five open-source families of  
10 LLMs in both paragraph classification and information extraction tasks, we observe  
11 excellent model performance even if only few example paragraphs are included in  
12 the ICL prompts. The results show the potential of the KEP approach for reducing  
13 human annotations and data curation efforts in automated scientific knowledge  
14 extraction.

## 15 1 Introduction

16 Reticular materials are a class of crystalline, porous materials made of molecular building blocks  
17 that are linked by strong chemical bonds [1]. They exhibit exceptional properties due to their highly  
18 porous structure, high surface area, tunable pore sizes and morphologies [2]. Their versatility is  
19 evidenced by a broad range of industrial applications, among them heterogeneous catalysis [3], energy

20 storage [4], water treatment [5], chemical sensing [6], heat transfer [7], gas capture [8] and drug  
21 delivery [9].

22 Following recent advances in generative AI, several models have been proposed to explore the large  
23 chemical space covered by reticular materials [10–14]. These models aim to generate reticular  
24 structures with optimized properties. Such structures are hypothetical as they have not yet been  
25 synthesised and tested in the lab. Devising a synthesis protocol for computationally generated  
26 structures requires a subject matter expert (SME). This is, however a challenging task given the large  
27 number of possible structures. An AI model that correlates a computationally discovered material  
28 with a lab synthesis protocol is, therefore, highly desirable. A first step towards the creation of such a  
29 model is building a database of existing synthesis protocols.

30 One approach for creating such database is applying information extraction techniques to the existing  
31 body of scientific literature. A large number of reticular materials have been reported in the literature  
32 alongside their respective synthesis protocols [15, 16]. It is worth noting, however, that overlapping  
33 discoveries are common, given that the same material can be produced by means of different synthesis  
34 protocols [17]. Transfer learning has been suggested as means to improve information extraction on  
35 existing corpora of scientific texts related to materials [18]. For example, fine-tuning techniques allow  
36 for adapting existing general-purpose AI models to specific tasks in domains for which comparatively  
37 little data exists. However, recent developments in LLMs have enabled information extraction based  
38 on prompt engineering and few-shot learning tasks [19].

39 In this paper, we propose using large language models (LLMs), without the need for additional training  
40 or fine-tuning, for extracting synthesis protocols of reticular materials from scientific literature, i.e.,  
41 unstructured PDF documents. We use prompt engineering with in-context learning (ICL) [20] for  
42 providing in the prompt all the context needed by the LLM to process the instructions. Together  
43 with instructions and input data, we provide examples that guide the LLM output production. This  
44 technique reduces the risk of hallucination, since all the context needed to execute the instruction is  
45 provided within the prompt. Also, it accelerates the process of information extraction because it does  
46 not require SME-based annotation of thousands of sentences/paragraphs for fine-tuning the models.

47 Our domain-independent Knowledge Extraction Pipeline (KEP) uses LLMs for extracting relevant  
48 information from PDF documents. The pipeline is composed of four main modules: (i) *PDF extractor*:  
49 processes the PDF to extract the text; (ii) *Paragraph classification*: processes the text in order to select  
50 only the relevant paragraphs (i.e., paragraphs that have the information the user is interested in); (iii)  
51 *Information extraction*: processes the relevant paragraphs and extract the relevant information; and  
52 (iv) *Knowledge representation*: interprets and assigns meaning to the information while representing  
53 the related knowledge. The pipeline uses LLMs with prompt-engineering and ICL in two modules,  
54 namely *paragraph classification* and *information extraction*, which are the focus of this paper. In  
55 addition, for identifying the best set of examples to be used in the prompts of these two modules,  
56 we propose the *Examples selection* phase. This phase measures the performance of the LLMs in a  
57 given task and, by using different sets of examples, identifies the set to be used for optimal LLM  
58 performance.

59 We have used five families of LLMs in both *paragraph classification* and *information extraction*  
60 modules and have compared their performance. We note that these open-source LLMs are not  
61 domain-specific and were not fine-tuned for our tasks. Our experiments indicate that: (i) even without  
62 fine-tuning or training, some of these models have achieved high performance in case ICL was used  
63 to provide examples in the prompt; (ii) the examples used in the prompt affect model performance  
64 and, hence, must be chosen carefully; and (iii) the same set of examples may lead to varying results if  
65 used in different models.

66 Some recent papers share our work’s objectives, however, they differ methodologically [19, 21–23].  
67 For example, Polak *et al.* (2024) [19] reported a pipeline for extracting information from unstructured  
68 text in the material discovery domain using language models. However, the cited work focused  
69 on simple extraction tasks, e.g., *material, value and unit*, while our pipeline is aimed at complex  
70 information associated with synthesis protocols that require additional classification. Unlike in our  
71 approach which is based on few-shot prompts providing examples for facilitating the information  
72 extraction, the cited work applies zero-shot methods for determining the relevance of sentences or  
73 paragraphs. Huo *et al.* (2019) [21] introduced a semi-supervised machine learning approach for  
74 classifying inorganic materials synthesis steps in scientific papers. The authors used the Latent  
75 Dirichlet Allocation (LDA) unsupervised topic modeling algorithm for clustering terms that are

76 typically used in synthesis descriptions. A random forest classifier, based on annotations of hundreds  
77 of paragraphs, categorized the occurring synthesis types. This approach also used a Markov chain for  
78 modeling the sequence of steps, creating flowcharts of synthesis procedures.

79 In Kononova *et al.* (2019) [22], the authors generated a dataset with “codified recipes” for solid-state  
80 synthesis which was automatically extracted from scientific publications using traditional text mining  
81 and natural language processing approaches. The authors used the two-step paragraph classification  
82 approach described in Huo *et al.* (2019) [21] for finding paragraphs on solid-state synthesis. The ex-  
83 traction pipeline consisted of several algorithms (BiLSTM-CRF, Material Parser, etc.) for identifying  
84 materials related information, including synthesis steps and conditions. Compared to our method, the  
85 cited work required considerable annotation effort and employed a less straightforward extraction  
86 pipeline. We note that our method relies primarily on the LLM capabilities for text understanding,  
87 without specialized tokenizers or entity recognizers. Finally, Park *et al.* (2022) [23] created a four-step  
88 pipeline, with text extraction from XML/HTML or PDF files and classifying relevant paragraphs,  
89 performing named entity recognition and, a fully connected multi-layer with dropout as classifier.

90 Another promising, less related approach is using “AI chatbot agents” for assisting materials scientists  
91 in specific pipeline tasks. In reference [24], the authors used prompt engineering for guiding  
92 a ChatGPT-based bot to extract MOF synthesis information from various sources. The authors  
93 leveraged a bot-like interface for answering questions about synthesis procedures and chemical  
94 reactions. In reference [25], the authors leveraged multiple AI assistants, such as LLMs and specific  
95 ML algorithms, as lab assistants to support a human SME, enabling productivity levels similar  
96 to those of an entire research team. While the approach was not fully automated, it provided a  
97 proof-of-concept of how language models can be leveraged for accelerating materials discovery.

98 The remainder of this paper is organized as follows. Section 2 introduces the use case, Section 3  
99 describes in details the pipeline applied to the use case and Section 4 presents our experiments.  
100 Section 5 concludes and presents some future work.

## 101 2 Use Case: Synthesis Protocols of Reticular Materials

102 With the goal of extracting knowledge about the synthesis of reticular materials, i.e., MOFs, ZIFs,  
103 COFs and zeolites, we have searched the scientific literature by using Elsevier’s API<sup>1</sup> and downloaded  
104 full-text PDFs from the SCOPUS database.<sup>2</sup> Our approach is based on extracting information from  
105 PDFs, and not XMLs, since not always a XML file will be available for a given document. Notice  
106 that our extraction pipeline (see Section 3) was not created to manipulate only documents available in  
107 Elsevier, where their XML files are also provided, but to process any PDF document (including those  
108 that are images).

109 Our search employed the following keywords and wildcard terms to capture relevant references:  
110 ‘MOF’, ‘metal organic framework’, ‘metal-organic framework’, ‘metal-organic-framework’, ‘COF’,  
111 ‘covalent organic’, ‘covalent-organic’, ‘ZIF’, and ‘zeolit\* imidazol\*’. We further limited the search  
112 to articles published in journals within Chemistry, Chemical Engineering, Materials Science, Energy,  
113 Engineering, Environmental Science, Physics and Astronomy, and Biochemistry, Genetics, and  
114 Molecular Biology, retrieving 6,669 articles.

115 The results were then filtered, by using the filter provided in the Elsevier API, to include only  
116 open-access articles with DOI identifiers from the following publishers: Elsevier (10.1016), Wiley  
117 Blackwell (10.1002), The Royal Society of Chemistry (10.1039), American Chemical Society  
118 (10.1021), Springer-Verlag (10.1007), Nature Publishing Group (10.1038), and MDPI (10.3390).

119 To create a public dataset, we finally kept only articles under the CC-BY-4.0 or CC-BY-3.0 licenses,  
120 resulting in 2,032 CC-BY-4.0 articles and 255 CC-BY-3.0 articles. These CCBY license papers  
121 were selected by performing web-scraping from the list of DOIs provided by the Elsevier API. Since  
122 we are considering only papers with CCBY 3.0 and 4.0 licenses, everyone can retrieve the PDFs.

---

<sup>1</sup><https://github.com/ElsevierDev/elsapy>

<sup>2</sup><https://www.scopus.com>

123 After collecting the data, we randomly selected 305 articles in PDF format <sup>3</sup>. We then extracted  
124 from these PDFs 188 paragraphs describing synthesis protocols, and 137 examples of paragraphs not  
125 describing synthesis protocols (a total of 325 paragraphs). This curated set of paragraphs constitutes  
126 our golden collection of classified paragraphs. For details about how those paragraphs were extracted,  
127 see Section 3.

128 Subsequently, a team of eleven research scientists (composed of 2 SMEs) annotated each of the  
129 synthesis-related paragraphs on a case-by-case basis for extracting the following information: (i) the  
130 description of the synthesis product; (ii) the equipment used as an energy source; (iii) the conditions  
131 under which the synthesis occurred (e.g., reaction time, reaction temperature, current density); and  
132 (iv) the reactants and solvents used, including their descriptions, quantities, and units of measurement.

133 Intentionally, some paragraphs were selected for annotation by multiple SMEs, leading to some  
134 inconsistencies. These inconsistencies were then used to refine the annotation guidelines. The data  
135 was reviewed on a case-by-case basis by SMEs using a custom-built graphical interface and compiled  
136 in a final set of 131 syntheses descriptions encoded in a JSON format, thereby creating our golden  
137 dataset of annotated synthesis information. Table 1 summarizes the data in our golden dataset.

Table 1: Overview of golden dataset

	Synthesis	Not Synthesis
paragraphs classified	188	137
annotated paragraphs	131	-

### 138 3 Knowledge Extraction Pipeline (KEP)

139 KEP is a domain-independent pipeline that helps extract knowledge from unstructured data. It is  
140 composed of four main modules: *PDF extractor*, *Paragraph classification*, *Information extraction*  
141 and *Knowledge representation*, as shown in Figure 1. The *PDF extractor* processes the PDF to extract  
142 paragraphs, since we assume that SMEs are interested in paragraphs containing specific information.  
143 The *Paragraph classification* classifies the extracted paragraphs into *relevant* or *irrelevant*, according  
144 to the task the SME is interested in. When applying this module to our use case, *relevant* paragraphs  
145 are those describing synthesis protocols of reticular materials.

146 *Information extraction* processes the relevant paragraphs and extracts the relevant information. When  
147 applying this module to our use case, the relevant information is the synthesis details such as the de-  
148 scription of the synthesis product, the experimental conditions (such as reaction time and temperature),  
149 and the reagents and solvents used in the synthesis. The final module, *Knowledge representation*,  
150 interprets and assigns meaning to the extracted information while creates the knowledge represen-  
151 tation. In the synthesis protocol use case, the knowledge representation is characterized by (i) the  
152 normalization of the unities; (ii) by the instantiation of entities of different kinds (such as productions,  
153 reactants and solvents), and (iii) by the instantiation of the relationships (such as used-reactant and  
154 used-solvent) that link the entities to the synthesis where they take part. For instance, it is possible to  
155 represent that the same reactant is being used in syntheses of two different products and that same  
156 product can be synthesized by two different synthesis.

157 The *PDF extractor* was implemented using the DS4SD open-source tool<sup>4</sup> that converts unstructured  
158 PDF documents into JSON files containing the document elements such as section titles, paragraphs,  
159 footnotes, headers, figure captions and tables, etc. DS4SD is also able to process PDFs that are  
160 indeed images since it uses an OCR engine to extract text-snippets from those images. The *Para-*  
161 *graph classification* and *Information extraction* modules, which are the focus of this paper, were  
162 implemented by using open source LLMs of the Flan, Granite, LLaMa, Mistral and Mixtral families.  
163 As detailed in Section 4, we compare the performance of these five families of LLMs when used in  
164 both the *Paragraph classification* and *Information extraction* modules. The LLMs were used without  
165 fine-tuning or training for the extraction of synthesis related information or on any task defined  
166 specifically for the Material Discovery domain. We only used prompt-engineering and ICL.

<sup>3</sup>171 articles with at least one paragraph describing a synthesis protocol and 134 articles without any synthesis protocol description.

<sup>4</sup><https://ds4sd.github.io/>

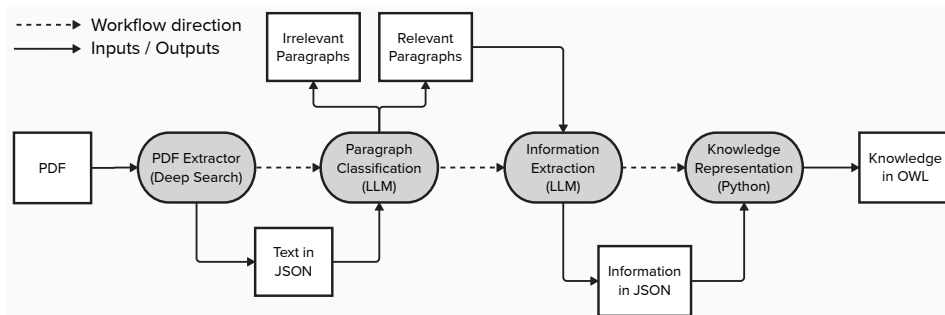


Figure 1: Knowledge Extraction Pipeline (KEP) with the four KEP modules highlighted in gray color. Also shown are the respective inputs and outputs

167 To select the best set of examples to be provided in the prompt, the pipeline adds an additional step  
 168 to each of the LLM’s modules, namely *Paragraph classification* and *Information extraction*. The  
 169 *Examples selection* step aims to select the best set of examples to be used in each tested LLMs for  
 170 each one of the tasks, *paragraph classification* and *information extraction*, see Section 3.3.

### 171 3.1 Paragraph classification

172 Since the goal of this module is the classification of paragraphs as *relevant* or *irrelevant*, the prompt  
 173 to be used in this model should describe the difference between a relevant and an irrelevant paragraph.  
 174 In addition, a sentence explicitly instructing the LLM that it should not provide an explanation  
 175 together with the classification may be required.

176 Since we are not using zero-shot prompting but ICL prompting, we not only provide the LLM with the  
 177 aforementioned instructions, but also give it several examples of paragraphs and their corresponding  
 178 classifications. In Section 4 we demonstrate that, by providing just a few examples in the prompt, the  
 179 performance of the LLMs tends to increase significantly. Below is an example of instructions used,  
 180 along with an example of paragraph<sup>5</sup> and its corresponding classification, also provided in the prompt.  
 181 This paragraph was classified as "S" meaning it is a paragraph describing a synthesis protocol.

182 **Instruction:** *You are assisting a chemist in classifying paragraphs from scientific articles. Mark*  
 183 *the paragraph as 'S' if it describes the components of synthesis protocols for reticular materials,*  
 184 *or 'I' if it does not include a synthesis description. After reviewing the examples, classify the given*  
 185 *paragraph. Do not add any information or explanation besides 'S' or 'I' in the answer.*

186 **Example:** *"Synthesis of Zn-MOF: Bis(imidazole-1-yl)methane was synthesized analogously to a*  
 187 *the procedure reported in [43]. All other materials were obtained from commercial sources and were*  
 188 *used as received. {[Zn(bim)(bdc)]0.8DMF0.4EtOH0.1H<sub>2</sub>O}<sub>n</sub> (Zn-MOF). Bis(imidazol-1-yl)methane*  
 189 *(bim) (3.0 mg, 0.02 mmol), terephthalic acid (6.6 mg, 0.04 mmol), and Zn(NO<sub>3</sub>)<sub>2</sub>·6H<sub>2</sub>O (7.6 mg, 0.02*  
 190 *mmol) were dissolved in DMF/EtOH/H<sub>2</sub>O (2:1:1, vol.) mixture (1 mL), placed in a 4 mL screw-cap*  
 191 *vial, and heated to 100 °C for 24 h."*

192  
 193 *Classification: S*

### 194 3.2 Information extraction

195 The prompt used in the *Information extraction* module should inform to the LLM the kind of  
 196 knowledge that should be extracted. In case of a complex structure, the prompt should suggest to the  
 197 LLM to represent the extracted information following a given schema in well-known format, such  
 198 as JSON [27]. It is reasonable to assume that the LLM will be able to parse this format since it is  
 199 a commonly used data format that appeared in several documents used to train the LLM. In order  
 200 to exemplify, find below the instruction we used and the JSON annotation related to the synthesis  
 201 paragraph presented in Section 3.1.

<sup>5</sup>Paragraph extracted from [26].

202 **Instruction:** *You are assisting a chemist in identifying and extracting descriptions of the synthesis*  
203 *of reticular materials from paragraphs. For each synthesis described in a paragraph, your task*  
204 *is to produce a JSON object that encodes the components involved in the synthesis, following the*  
205 *format provided in the examples. After reviewing the examples, carefully analyze the last paragraph*  
206 *and create a JSON object for each synthesis you find, ensuring that it adheres to the structure and*  
207 *conventions demonstrated.*

208 **Example:** *"Synthesis of Zn-MOF: Bis(imidazole-1-yl)methane was synthesized analogously to a*  
209 *... screw-cap vial, and heated to 100 °C for 24 h."*

```
210 {"output": {  
211   "product": {  
212     "description": "Zn-MOF",  
213     "material_type": "MOF",  
214     "conditions": [  
215       {"description": "reaction temperature", "value": 100, "unit": "oC"},  
216       {"description": "reaction time", "value": 24, "unit": "h"}  
217     ]  
218   },  
219   "reactants": [  
220     {"description": "Bis(imidazol-1-yl)methane (bim)", "value": 0.02, "unit": "mmol"},  
221     {"description": "terephthalic acid", "value": 0.04, "unit": "mmol"},  
222     {"description": "Zn(NO3)2-6H2O", "value": 0.02, "unit": "mmol"}  
223   ],  
224   "solvents": [  
225     {"description": "DMF/EtOH/H2O (2:1:1, vol)", "value": 1.0, "unit": "mL"}  
226   ]  
227 }}
```

### 228 3.3 Examples selection

229 It is well-known that the performance of LLMs to execute a given task is significantly influenced by  
230 the set of examples provided in the prompt. In addition, due to the different characteristics of how  
231 the LLMs were trained, it is expected that different LLMs will require different sets of examples to  
232 achieve their highest performance when executing the same task.

233 Therefore, the *Examples selection* step was included and associated with each KEP module that uses  
234 LLMs to help on the selection of the best set of prompt examples to be used. *Examples selection*  
235 receives as input the model to be tested, a golden dataset and the number of examples to be selected as  
236 examples. It randomly selects from the dataset some instances to be used as examples in the prompt,  
237 and all other instances are used to measure the performance of the model. This step is executed for  
238 all possible combinations of examples or until the user is satisfied with the performance of the model  
239 in one of the executions. The set of examples that leads the LLM to achieve the highest performance  
240 is the one selected to be used in the associated KEP module.

## 241 4 Experiments

242 This section presents the experiments we ran with 5 families of open-source LLMs. None of them  
243 were trained or fine-tuned to extract synthesis details from paragraphs or to execute any specific task  
244 in the Material Discovery domain. We selected 2 models of each family<sup>6</sup>, prioritized the models  
245 that have been fine-tuned using a collection of instructions (not related to our tasks) and chosen the  
246 last released ones<sup>7</sup>. Ultimately, the selected models were: (i) flan: flan-t5-xxl-11b, flan-ul2-20b;  
247 (ii) granite: granite-20b-code-instruct, granite-34b-code-instruct; (iii) llama: llama-3-70b-instruct,  
248 llama-3.1-405b-instruct; (iv) mistral: mistral-large; and (v) mixtral: mixtral-8x7b-instruct-v01. See  
249 the description of each model in Appendix A.

<sup>6</sup>Exceptions: mistral and mixtral

<sup>7</sup>Exception: llama-3-70b-instruct selected instead of llama-3-1-70b-instruct since it has a highest performance in the tasks we are testing.

250 **4.1 Examples selection**

251 **Paragraph classification:** From the original set of 325 classified paragraphs, we reduced the golden  
 252 dataset by downselecting only 50 paragraphs to demonstrate that, even when testing the prompt  
 253 examples selection in a small dataset, it is possible to achieve a good performance on a majority of  
 254 the tested models. In addition, the use of a small dataset helps demonstrate that the approach does not  
 255 require the manual classification/annotation of thousands or hundreds of examples.

256 In the set of 50 paragraphs we ensure that 25 paragraphs are relevant (i.e., classified with "S" and  
 257 mentioning synthesis protocol) and 25 are irrelevant (i.e., classified with "I" and not mentioning  
 258 synthesis protocols). We fixed the number of examples to be provided in the prompt to 5, since  
 259 paragraphs describing synthesis protocols are typically very large and the prompts have a limited  
 260 number of tokens. Our goal is to find the best set of 5 examples used in the prompt that helps the  
 261 models achieve their highest performance. The accuracy of each model was measured by using the  
 262 F1 metric.

263 For each model, we executed 100 runs by providing in the prompt the instruction mentioned in  
 264 Section 3.1 and 5 examples randomly selected from 50 possibilities. We tested the output with the  
 265 remaining 45 paragraphs not provided in the prompt. Table 2 presents the result of our experiments.  
 266 For each one of the models, the table indicates the number of paragraphs mentioning synthesis  
 267 protocols ("S") and the number of irrelevant paragraphs ("I") used in both the worst and best prompt  
 268 together with the F1 value for each case.

Table 2: The best-case (highlighted in bold) and worst-case (underlined) scenarios in the selection of examples to be used in the prompt of the *Paragraph classification* module.

Model	Worst			Best		
	#S	#I	F1	#S	#I	F1
flan-t5-xxl-11b	1	4	0.93	3	2	<b>1.0</b>
flan-ul2-20b	3	2	<u>0.0</u>	1	4	0.98
granite-34b-code-instruct	1	4	<u>0.30</u>	2	3	0.92
granite-20b-code-instruct	2	3	0.32	2	3	0.74
llama-3-70b-instruct	1	4	0.71	4	1	<b>1.0</b>
llama-3.1-405b-instruct	3	2	<u>0.0</u>	3	2	0.95
mistral-large	2	3	<u>0.76</u>	4	1	<b>1.0</b>
mixtral-8x7b-instruct-v01	3	2	0.61	3	2	<b>1.0</b>

269 The models with highest performance were flan-t5-xxl-11b, llama-3-70b-instruct, mistral-large and  
 270 mixtral-8x7b-instruct-v01. Although llama-3-70b-instruct and mistral-large used the same number of  
 271 relevant paragraphs and the same number of irrelevant paragraphs in their best cases, their prompts  
 272 share only one paragraph (see Table 6 in Appendix B). When testing the best prompt for mistral-large  
 273 in llama-3-70b-instruct by using the same 45 testing examples, the performance of the model did not  
 274 achieve F1=1.0, but F1=0.98. Although it is a small difference, it demonstrate that, different LLMs  
 275 may need different examples in their prompts to achieve their highest performance. The models  
 276 with worst performance were flan-ul2-20b and llama-3.1-405b-instruct. Although we included in the  
 277 prompt a sentence stating that the answer should only include "S" or "I", their answers often also  
 278 include an explanation; which we considered to be a hallucination and, thus, an incorrect answer.

279 **Information extraction:** The golden dataset used in this step is the 25 paragraphs mentioning  
 280 synthesis protocols used in the previous step together with their corresponding JSON annotations.  
 281 Different from the previous step, here we fixed the number of examples used in the prompt to 2, since  
 282 the JSON annotation is being provided together with the paragraph, which significantly increases the  
 283 number of tokens. Even with only 2 examples, flan-t5-xxl and flan-ul2 could not be tested since their  
 284 prompt+result do not accept so many tokens <sup>8</sup>.

285 The experiment begun by randomly selecting 2 paragraphs+JSON to be used in the prompt for  
 286 each one of the 6 models. For each model, we executed 100 runs by providing in the prompt the  
 287 instructions mentioned in Section 3.2 and the 2 examples of paragraph+JSON randomly selected  
 288 from 25 possibilities. We tested the performance of the model with each prompt by using the 23

<sup>8</sup>Both flan models accept only 4,096 when comparing to llama that accepts 8,192

289 paragraphs that were not provided as examples in the prompt. The results are presented in Table 3.  
 290 To compare the JSON annotations provided by the LLM with the JSON annotations included in the  
 291 golden dataset, a structure analysis based on each JSON key (i.e., name/value pair) was defined<sup>9</sup>.

Table 3: The best-case and worst-case scenarios in the selection of examples of the *Information extraction* module. The best results are highlighted in bold and the worst results are underlined.

Model	Worst accuracy	Best accuracy
granite-34b-code-instruct	0.70	0.93
granite-20b-code-instruct	0.65	0.84
llama-3-70b-instruct	0.54	<b>0.95</b>
llama-3.1-405b-instruct	0.53	0.94
mistral-large	<u>0.22</u>	0.94
mixtral-8x7b-instruct-v01	0.70	0.93

292 The models that achieved the highest accuracy were llama-3-70b-instruct, llama-3.1-405b-instruct  
 293 and mistral-large. However, it is important to notice that all of them achieved an accuracy higher than  
 294 **0.84** even using only two examples in the prompt. Similar to what happened in the previous step, the  
 295 experiments illustrate the influence of the examples in the accuracy of the model (E.g. llama-3.1-405b-  
 296 instruct worst case was 0.53 and best case was 0.94). In addition, one of the paragraphs presented in  
 297 the worst case of mistral-large appeared in the best case of mixtral-8x7b-instruct-v01 (see Table 7 in  
 298 Appendix B). Two related models that have the same example in opposite scenarios. Moreover, it is  
 299 important to highlight that the two granites, the two llamas, and mixtral-8x7b-instruct-v01 included  
 300 in their worst scenarios the same paragraph (see Table 7 in Appendix B). It may indicate that there  
 301 are examples that really do not help the models on executing their tasks.

## 302 4.2 Paragraph classification

303 After selecting the final set of five examples that maximize the performance of each model, the  
 304 *paragraph classification* module was tested by using the entire golden dataset of 275 paragraphs  
 305 (325 minus the 50 used for prompt selection). For each model, the prompt was composed of the  
 306 instructions mentioned in Section 3.1 and the best set of examples selected for that model, as presented  
 307 in Section 4.1. Table 4 summarizes the results for each model in terms of Precision, Recall, and F1  
 308 achieved with the best prompt. Llama-3-70b-instruct and mistral-large achieved **F1=0.98**. Although  
 309 llama-3.1-405b-instruct and flan-ul2-20b have more parameters than the other model of their families,  
 310 their performances were worse. It occurred due the hallucination mentioned in the *Example section*  
 311 step. Excluding granite-20b-code-instruct, all the models achieved **F1>0.84**, which is very good  
 312 accuracy given that only five examples were provided in the prompt to these models.

Table 4: Experiments for the *Paragraph classification* module (best results highlighted in bold).

Model	Precision	Recall	F1
flan-t5-xxl-11b	<b>0.98</b>	0.96	0.97
flan-ul2-20b	0.96	0.96	0.96
granite-34b-code-instruct	0.87	0.83	0.84
granite-20b-code-instruct	0.75	0.70	0.72
llama-3-70b-instruct	<b>0.98</b>	<b>0.98</b>	<b>0.98</b>
llama-3.1-405b-instruct	<b>0.98</b>	0.83	0.88
mistral-large	<b>0.98</b>	<b>0.98</b>	<b>0.98</b>
mixtral-8x7b-instruct-v01	0.95	0.93	0.94

## 313 4.3 Information extraction

314 This module was tested by using the golden dataset of 106 annotated paragraphs (131 minus the 25  
 315 used for prompt selection). For each model, the prompt was composed of the instructions mentioned  
 316 in Section 3.2 and the best set of examples selected for that model, as presented in Section 4.1.

<sup>9</sup>To create a more fine-grained comparison between the JSONs, it would be necessary to compare their semantics and not only their structures, as different structures could have the same meaning.



317 Table 5 summarizes the results of our experiments. The model that achieved the highest accuracy  
318 (**0.96**) was llama-3.1-405b-instruct, which is the biggest one. Other four models also achieved a  
319 very similar and high performance (mixtral-8x7b-instruct-v01, mistral-large, llama-3-70b-instruct  
320 and granite-34b-code-instruct). Notice that the smallest model (granite-20b-code-instruct) was the  
321 one that achieved the lower performance. The high accuracy achieved by the biggest models when  
322 compared to the smallest one is expected due to the complex of the task that involves the creation  
323 of a correct JSON. When considering both the *Paragraph classification* and *Information extraction*  
324 modules, the three models with highest performance and, thus, those that should be considered to  
325 be used in KEP to process all the selected papers mention in Section 2 are: llama-3-70b-instruct,  
326 mistral-large and mixtral-8x7b-instruct-v01.

Table 5: Experiments for the *Information extraction* module (best results highlighted in bold).

Model	Accuracy
granite-34b-code-instruct	0.93
granite-20b-code-instruct	0.84
llama-3-70b-instruct	0.93
llama-3.1-405b-instruct	<b>0.96</b>
mistral-large	0.95
mixtral-8x7b-instruct-v01	0.94

## 327 5 Conclusions and Future Research

328 In summary, we present a knowledge extraction pipeline for synthesis protocols of reticular materials  
329 that significantly reduces SME based classification and annotation tasks related to the training or fine-  
330 tuning of machine learning models. Our experimental results indicate that LLMs can achieve high  
331 performance with a limited set of examples within the prompt, even without training or fine-tuning  
332 the models for the specific domain. For example, by including five representative paragraphs in the  
333 prompt, we have reproducibly achieved F1=**0.98** in paragraph classification tasks. In information  
334 extraction tasks, we have used two paragraphs + JSON and llama-3.1-405b-instruct for achieving  
335 Accuracy=**0.96**.

336 Our results highlight the necessity of testing different examples to be used in the prompt as this  
337 variation strongly influences model performance. For instance, in the *Paragraph classification*  
338 module, the performance of mixtral-8x7b-instruct-v01, one of the best models in our study, ranges  
339 from F1=**0.61** to F1=**1.0**. In addition, the experiments show that different LLMs may require different  
340 sets of examples for achieving top performance. Although both llama-3-70b-instruct and mistral-large  
341 included four synthesis paragraphs and one irrelevant paragraph in their best set of examples, llama-  
342 3-70b-instruct has not achieved its highest performance with the best prompt chosen for mistral-large.  
343 Finally, a huge number of parameters in the model does not necessarily guarantee a superior model  
344 performance. Both flan-ul2 and llama-3.1-405b-instruct failed to achieve top performance in the  
345 classification of paragraphs if compared to other models of the same family.

346 Future research work should include comparative analyses with nonLLM methods in view of extrac-  
347 tion time and quality, as well as measuring LLMs’ performance for different materials applications.  
348 For creating a dataset of synthesis protocols for reticular materials, future research should address the  
349 following: (i) **refine JSONs comparison**: The creation of metrics for semantically comparing JSONs  
350 is needed to validate if the output of the model is structurally comparable with the golden dataset, and  
351 if it should be considered a valid JSON; (ii) **workflow extraction**: The extension of the *Information*  
352 *extraction* module for extracting the synthesis workflow step-by-step; and (iii) **increase use case**  
353 **coverage**: The application of KEP to all paragraphs extracted from the selected 2,287 papers. Once  
354 processed, the resulting data set should be explored for analyzing the distributions of synthesis details  
355 made available in the scientific literature.

## 356 References

357 [1] O. M. Yaghi, M. O’Keeffe, N. W. Ockwig, H. K. Chae, M. Eddaoudi, and J. Kim, “Reticular  
358 synthesis and the design of new materials,” *Nature*, vol. 423, no. 6941, pp. 705–714, 2003.

- 359 [2] H. Lyu, Z. Ji, S. Wuttke, and O. M. Yaghi, "Digital reticular chemistry," *Chem*, vol. 6, no. 9, pp.  
360 2219–2241, 2020.
- 361 [3] A. Bavykina, N. Kolobov, I. S. Khan, J. A. Bau, A. Ramirez, and J. Gascon, "Metal–organic  
362 frameworks in heterogeneous catalysis: recent progress, new trends, and future perspectives,"  
363 *Chemical reviews*, vol. 120, no. 16, pp. 8468–8535, 2020.
- 364 [4] R. Zhao, Z. Liang, R. Zou, and Q. Xu, "Metal-organic frameworks for batteries," *Joule*, vol. 2,  
365 no. 11, pp. 2235–2259, 2018.
- 366 [5] R. Li, N. N. Adarsh, H. Lu, and M. Wriedt, "Metal-organic frameworks as platforms for the  
367 removal of per- and polyfluoroalkyl substances from contaminated waters," *Matter*, vol. 5, no. 10,  
368 pp. 3161–3193, 2022.
- 369 [6] L. E. Kreno, K. Leong, O. K. Farha, M. Allendorf, R. P. Van Duyne, and J. T. Hupp, "Metal–  
370 organic framework materials as chemical sensors," *Chemical reviews*, vol. 112, no. 2, pp.  
371 1105–1125, 2012.
- 372 [7] M. Islamov, H. Babaei, R. Anderson, K. B. Sezginel, J. R. Long, A. J. McGaughey, D. A.  
373 Gomez-Gualdrón, and C. E. Wilmer, "High-throughput screening of hypothetical metal-organic  
374 frameworks for thermal conductivity," *npj Computational Materials*, vol. 9, no. 1, p. 11, 2023.
- 375 [8] P. Z. Moghadam, Y. G. Chung, and R. Q. Snurr, "Progress toward the computational discovery  
376 of new metal–organic framework adsorbents for energy applications," *Nature Energy*, vol. 9,  
377 no. 2, pp. 121–133, 2024.
- 378 [9] P. Horcajada, T. Chalati, C. Serre, B. Gillet, C. Sebrie, T. Baati, J. F. Eubank, D. Heurtaux,  
379 P. Clayette, C. Kreuz *et al.*, "Porous metal–organic-framework nanoscale carriers as a potential  
380 platform for drug delivery and imaging," *Nature materials*, vol. 9, no. 2, pp. 172–178, 2010.
- 381 [10] Z. Yao, B. Sánchez-Lengeling, N. S. Bobbitt, B. J. Bucior, S. G. H. Kumar, S. P. Collins,  
382 T. Burns, T. K. Woo, O. K. Farha, R. Q. Snurr *et al.*, "Inverse design of nanoporous crystalline  
383 reticular materials with deep generative models," *Nature Machine Intelligence*, vol. 3, no. 1, pp.  
384 76–86, 2021.
- 385 [11] H. Park, S. Majumdar, X. Zhang, J. Kim, and B. Smit, "Inverse design of metal–organic  
386 frameworks for direct air capture of CO<sub>2</sub> via deep reinforcement learning," *Digital Discovery*,  
387 vol. 3, no. 4, pp. 728–741, 2024.
- 388 [12] H. Park, X. Yan, R. Zhu, E. A. Huerta, S. Chaudhuri, D. Cooper, I. Foster, and E. Tajkhorshid,  
389 "A generative artificial intelligence framework based on a molecular diffusion model for the  
390 design of metal-organic frameworks for carbon capture," *Communications Chemistry*, vol. 7,  
391 no. 1, p. 21, 2024.
- 392 [13] Y. Kang and J. Kim, "Chatmof: an artificial intelligence system for predicting and generating  
393 metal-organic frameworks using large language models," *Nature Communications*, vol. 15,  
394 no. 1, p. 4705, 2024.
- 395 [14] F. Cipcigan, J. Booth, R. N. B. Ferreira, C. R. dos Santos, and M. Steiner, "Discovery of novel  
396 reticular materials for carbon dioxide capture using gflownets," *Digital Discovery*, vol. 3, no. 3,  
397 pp. 449–455, 2024.
- 398 [15] P. Z. Moghadam, A. Li, S. B. Wiggin, A. Tao, A. G. Maloney, P. A. Wood, S. C. Ward, and  
399 D. Fairen-Jimenez, "Development of a Cambridge structural database subset: a collection of  
400 metal–organic frameworks for past, present, and future," *Chemistry of Materials*, vol. 29, no. 7,  
401 pp. 2618–2625, 2017.
- 402 [16] P. Z. Moghadam, A. Li, X.-W. Liu, R. Bueno-Perez, S.-D. Wang, S. B. Wiggin, P. A. Wood,  
403 and D. Fairen-Jimenez, "Targeted classification of metal–organic frameworks in the Cambridge  
404 structural database (csd)," *Chemical science*, vol. 11, no. 32, pp. 8373–8387, 2020.
- 405 [17] Y. G. Chung, E. Haldoupis, B. J. Bucior, M. Haranczyk, S. Lee, H. Zhang, K. D. Vo-  
406 giatzis, M. Milisavljevic, S. Ling, J. S. Camp *et al.*, "Advances, updates, and analytics for  
407 the computation-ready, experimental metal–organic framework database: Core mof 2019,"  
408 *Journal of Chemical & Engineering Data*, vol. 64, no. 12, pp. 5985–5998, 2019.

- 409 [18] E. A. Olivetti, J. M. Cole, E. Kim, O. Kononova, G. Ceder, T. Y.-J. Han, and A. M. Hiszpanski,  
410 “Data-driven materials research enabled by natural language processing and information  
411 extraction,” *Applied Physics Reviews*, vol. 7, no. 4, p. 041317, 12 2020. [Online]. Available:  
412 <https://doi.org/10.1063/5.0021106>
- 413 [19] M. P. Polak and D. Morgan, “Extracting accurate materials data from research papers with  
414 conversational language models and prompt engineering,” *Nature Communications*, vol. 15,  
415 no. 1, p. 1569, 2024. [Online]. Available: <https://doi.org/10.1038/s41467-024-45914-8>
- 416 [20] N. Wies, Y. Levine, and A. Shashua, “The learnability of in-context learning,” in  
417 *Advances in Neural Information Processing Systems*, A. Oh, T. Naumann, A. Globerson,  
418 K. Saenko, M. Hardt, and S. Levine, Eds., vol. 36. Curran Associates, Inc., 2023, pp.  
419 36 637–36 651. [Online]. Available: [https://proceedings.neurips.cc/paper\\_files/paper/2023/file/73950f0eb4ac0925dc71ba2406893320-Paper-Conference.pdf](https://proceedings.neurips.cc/paper_files/paper/2023/file/73950f0eb4ac0925dc71ba2406893320-Paper-Conference.pdf)
- 421 [21] H. Huo, Z. Rong, O. Kononova, W. Sun, T. Botari, T. He, V. Tshitoyan, and  
422 G. Ceder, “Semi-supervised machine-learning classification of materials synthesis procedures,”  
423 *npj Computational Materials*, vol. 5, no. 1, p. 62, Jul 2019. [Online]. Available:  
424 <https://doi.org/10.1038/s41524-019-0204-1>
- 425 [22] O. Kononova, H. Huo, T. He, Z. Rong, T. Botari, W. Sun, V. Tshitoyan, and G. Ceder,  
426 “Text-mined dataset of inorganic materials synthesis recipes,” *Scientific Data*, vol. 6, no. 1, p.  
427 203, Oct 2019. [Online]. Available: <https://doi.org/10.1038/s41597-019-0224-1>
- 428 [23] H. Park, Y. Kang, W. Choe, and J. Kim, “Mining insights on metal–organic framework synthesis  
429 from scientific literature texts,” *Journal of Chemical Information and Modeling*, vol. 62, no. 5,  
430 pp. 1190–1198, 2022.
- 431 [24] Z. Zheng, O. Zhang, C. Borgs, J. T. Chayes, and O. M. Yaghi, “Chatgpt chemistry assistant for  
432 text mining and the prediction of mof synthesis,” *Journal of the American Chemical Society*,  
433 vol. 145, no. 32, pp. 18 048–18 062, 2023.
- 434 [25] Z. Zheng, O. Zhang, H. L. Nguyen, N. Rampal, A. H. Alawadhi, Z. Rong, T. Head-Gordon,  
435 C. Borgs, J. T. Chayes, and O. M. Yaghi, “Chatgpt research group for optimizing the crystallinity  
436 of mofs and cofs,” *ACS Central Science*, vol. 9, no. 11, pp. 2161–2170, 2023.
- 437 [26] V. V. Matveevskaya, D. I. Pavlov, A. A. Ryadun, V. P. Fedin, and A. S.  
438 Potapov, “Synthesis, crystal structure, and luminescent sensing properties of a  
439 supramolecular 3d zinc(ii) metal–organic framework with terephthalate and bis(imidazol-  
440 1-yl)methane linkers,” *Inorganics*, vol. 11, no. 7, p. 264, Jun. 2023. [Online]. Available:  
441 <http://dx.doi.org/10.3390/inorganics11070264>
- 442 [27] *ECMA-404: The JSON data interchange syntax*, ECMA International Std. 404, 2017. [Online].  
443 Available: <https://ecma-international.org/publications-and-standards/standards/ecma-404/>
- 444 [28] H. W. Chung, L. Hou, S. Longpre, B. Zoph, Y. Tay, W. Fedus, E. Li, X. Wang, M. De-  
445 hghani, S. Brahma *et al.*, “Scaling instruction-finetuned language models,” *arXiv preprint*  
446 *arXiv:2210.11416*, 2022.
- 447 [29] Y. Tay, H. W. Chung, L. Hou, B. Zoph, S. Borgeaud, P. He, S. Narang, W. Fedus, and D. G. Patil,  
448 “UI2 20b: An open-source unified language learner model,” *arXiv preprint arXiv:2301.07520*,  
449 2023.
- 450 [30] I. Research, “Granite code models: A family of open foundation models for code intelligence,”  
451 *IBM Documentation*, 2024. [Online]. Available: <https://www.ibm.com/granite/playground/code/>
- 452 [31] *Llama-3-70B-Instruct*, Meta Std., 2024. [Online]. Available: <https://huggingface.co/meta-llama/Meta-Llama-3-70B-Instruct>
- 454 [32] *Llama-3.1-405B*, Meta Std., 2024. [Online]. Available: <https://huggingface.co/meta-llama/Llama-3.1-405B>  
455

- 456 [33] H.-C. Tsai, Y.-F. Huang, and C.-W. Kuo, “Comparative analysis of automatic literature review  
457 using mistral large language model and human reviewers,” *Sciety*, 2024. [Online]. Available:  
458 <https://sciety.org/articles/activity/10.21203/rs.3.rs-4022248/v1>
- 459 [34] M. AI, “Sparse mixture of experts in large language models: Mixtral 8x7b,” *arXiv preprint*  
460 *arXiv:2401.04088*, 2024. [Online]. Available: <https://arxiv.org/abs/2401.04088>

## 461 **A Models Description**

462 Flan-T5 [28] is a variant of the T5 (Text-to-Text Transfer Transformer) model, further fine-tuned using  
463 a mixture of instruction-based learning tasks. Like the original T5, Flan-T5 leverages a transformer  
464 architecture, specifically designed for text-to-text tasks, which means it treats both the input and  
465 output as text sequences, regardless of the task (e.g., translation, summarization, question-answering).  
466 The “Flan” component (Fine-tuned LAngeuage Net) introduces instruction tuning, where the model  
467 is exposed to a variety of natural language instructions during its fine-tuning phase. This method  
468 allows the model to generalize better across tasks by learning to follow explicit human instructions.  
469 In essence, Flan-T5 adapts the standard pre-training and fine-tuning methods of T5 but adds an  
470 additional layer of task diversity through its instruction-based training. This approach enhances its  
471 performance on zero-shot and few-shot learning tasks, making it more versatile for a wide range of  
472 NLP applications.

473 Flan-UL2 (Unified Language Learner) [29] is a variant of the UL2 architecture, designed for improved  
474 instruction-based fine-tuning similar to Flan-T5. UL2 is an advanced architecture that introduces a  
475 novel pre-training method utilizing a mixture of denoising tasks with different difficulty levels. This  
476 approach allows the model to adapt to a wider range of NLP tasks by balancing between simple and  
477 complex learning objectives. In the case of Flan-UL2, this model takes UL2 and further enhances it  
478 with instruction tuning, similar to the Flan-T5 approach. It is trained on a large variety of instruction  
479 tasks, making it highly proficient at zero-shot and few-shot learning across many tasks, such as  
480 summarization, translation, and question answering. The model’s ability to generalize across these  
481 tasks is further improved by the fine-tuning process with diverse datasets of instructions, allowing it to  
482 better understand human-like queries and execute complex tasks. This makes Flan-UL2 particularly  
483 powerful for applications requiring high versatility and adaptability in natural language understanding.

484 Granite-20B-Code-Instruct and Granite-34B-Code-Instruct [30] are part of the Granite family of  
485 large language models (LLMs) designed specifically for code-related tasks. Both models are fine-  
486 tuned versions of their respective base models, Granite-20B-Code-Base and Granite-34B-Code-Base,  
487 using instruction-based datasets to improve their ability to follow natural language instructions.  
488 These models, developed by IBM Research, are built for tasks such as code generation, bug fixing,  
489 code explanation, and translation across a wide range of programming languages, making them  
490 versatile tools for code-centric applications. Granite-20B-Code-Instruct, with 20 billion parameters,  
491 was trained on trillions of tokens from various sources, including high-quality code, mathematical  
492 data, and instructional prompts. Its fine-tuning emphasizes logical reasoning and problem-solving,  
493 with a focus on generating and explaining code, alongside supporting tasks like API calling and  
494 debugging . Granite-34B-Code-Instruct, with 34 billion parameters, extends these capabilities by  
495 being a more computationally powerful model, trained on a larger and more diverse dataset of code  
496 instructions. It can handle more complex coding tasks and demonstrates state-of-the-art performance  
497 across benchmarks for code synthesis, explanation, and debugging . Both models are decoder-only  
498 architectures, optimized for generating human-readable code outputs from natural language inputs,  
499 and are trained with instruction tuning to improve their accuracy in code-based applications.

500 Llama-3-70B-Instruct [31] is part of Meta’s Llama 3 family of large language models, specifically  
501 designed for instruction-following tasks. The model contains 70 billion parameters and is optimized  
502 for generating text in response to user prompts. It is a decoder-only model, which uses an optimized  
503 transformer architecture. The instruction-tuned version of Llama-3-70B benefits from Supervised  
504 Fine-Tuning (SFT) and Reinforcement Learning with Human Feedback (RLHF) to align its outputs  
505 with human preferences for helpfulness and safety. This fine-tuning process makes it particularly  
506 suitable for assistant-like applications, such as chatbots and task-oriented dialogue systems. Llama-3-  
507 70B-Instruct was trained on an extensive corpus of 15 trillion tokens from publicly available datasets  
508 and supports a wide range of use cases, including multilingual text generation and code-related  
509 tasks. It incorporates improvements like Grouped-Query Attention (GQA) for faster inference and an

510 expanded 8,192 token context window, allowing it to handle longer inputs effectively. The model  
511 has been tested extensively for safety, and Meta has integrated safeguards to limit misuse, including  
512 rigorous red teaming and cybersecurity assessments. The model is available under the Meta Llama 3  
513 Community License for both commercial and research applications. It’s praised for outperforming  
514 other models in several benchmarks, demonstrating significant advancements in multilingual dialogue  
515 capabilities and code generation.

516 Llama 3.1-405B-Instruct [32] is the largest model in the Llama 3.1 series by Meta, designed to provide  
517 state-of-the-art performance in multilingual dialogue and complex instruction-following tasks. With  
518 405 billion parameters, it utilizes a transformer-based, decoder-only architecture optimized for  
519 extensive text generation tasks. It introduces enhancements in context handling, supporting up  
520 to 128,000 tokens, which makes it ideal for tasks like document summarization and long-context  
521 conversation . This model is fine-tuned using a combination of Supervised Fine-Tuning (SFT) and  
522 Reinforcement Learning with Human Feedback (RLHF), enabling it to align better with human  
523 preferences and improve the safety and helpfulness of its outputs . Llama 3.1-405B was trained on a  
524 mixture of publicly available datasets containing approximately 15 trillion tokens, and its fine-tuning  
525 included more than 25 million synthetically generated instruction-based examples . Furthermore,  
526 it offers improved multilingual support beyond English, covering languages like German, French,  
527 Italian, Portuguese, Hindi, Spanish, and Thai . The model is open-source and available under Meta’s  
528 custom open model license, encouraging use in both research and commercial applications .

529 Mistral AI’s large language models, particularly Mistral Large 2 [33], represent significant advance-  
530 ments in both computational efficiency and reasoning capabilities. This model, featuring 123 billion  
531 parameters, is designed for tasks that require extensive reasoning, such as multilingual text processing,  
532 code generation, and mathematical problem-solving. With support for over 80 coding languages and  
533 a context window of 128,000 tokens, it excels in handling large documents and long, complex inputs.  
534 Mistral Large 2 is particularly strong in benchmarks like MMLU (Massive Multitask Language  
535 Understanding), where it achieves an accuracy of 84

536 Mixtral-8x7B-Instruct-v0.1 [34] is an advanced sparse mixture-of-experts (SMoE) model developed  
537 by Mistral AI. It incorporates a unique architecture where each layer contains eight experts (feedfor-  
538 ward blocks), but only two are activated for each token during inference. This selective processing  
539 allows the model to manage a large number of parameters—47 billion in total—while only using 13  
540 billion active parameters per token, which significantly reduces computation costs during inference.  
541 Mixtral-8x7B-Instruct has been fine-tuned for instruction-following tasks through a combination  
542 of supervised fine-tuning (SFT) and Direct Preference Optimization (DPO). This model excels in  
543 benchmarks such as MMLU and GSM8K, matching or outperforming larger models like GPT-3.5  
544 Turbo and Llama 2 70B in several areas, particularly code generation, reasoning, and multilingual  
545 tasks. Its ability to handle long sequences with a 32k token context window makes it highly effective  
546 for long-range information retrieval and complex prompts.

## 547 **B Examples selection**

548 **Paragraph classification** Table 6 shows excerpts of JSONs with best and worst paragraphs selected  
549 as examples for each model. It is possible to see that few paragraphs appear in more than one prompt.

Table 6: Partial JSONs with the best and worse examples for *Paragraph classification* module.

<pre> "flan-t5-xxl":{   "best":{     "prompt":"You are assisting a chemist in classifying paragraphs from     \"Paragraph1\": \"DFT simulations on CD-MOF-2 were conducted to understa     \"Label1\": \"I\",     \"Paragraph2\": \"2.2. The Construction of Cu-As MOF The process of cre     \"Label2\": \"S\",     \"Paragraph3\": \"2.4. Synthesis of Cu<sub>3</sub>(BTC)<sub>2</sub> and Zn<sub>3</sub>(B)     \"Label3\": \"S\",     \"Paragraph4\": \"2.2. Synthesis of {[Tb<sub>5</sub>SL<sub>6</sub>(OH)<sub>3</sub>](HS<sub>2</sub>)     \"Label4\": \"S\",     \"Paragraph5\": \"Two porous organic polymers (POPs) incorporating a fe     \"Label5\": \"I\",     \"accuracy\":1.0,     \"f1_score\":1.0,     \"run_number\":41,     \"correct_items\":43,     \"incorrect_items_count\":0,     \"S_exemples\":3,     \"I_exemples\":2   },   \"worst\":{     \"prompt\":\"You are assisting a chemist in classifying paragraphs from     \"Paragraph1\": \"Preparation of CdMOF-1. A mixture of Cd(NOS<sub>3</sub>)<sub>2</sub>(     \"Label1\": \"S\",     \"Paragraph2\": \"Since amines readily adsorb pollutants, boosting the     \"Label2\": \"I\",     \"Paragraph3\": \"The potential for using MOF materials to remove fluor     \"Label3\": \"I\",     \"Paragraph4\": \"Stability tests upon light illumination with differer     \"Label4\": \"I\",     \"Paragraph5\": \"Researchers have conducted extensive studies on wette     \"Label5\": \"I\",     \"accuracy\":0.9333333333333333,     \"f1_score\":0.9327380158829563,     \"run_number\":28,     \"correct_items\":42,     \"incorrect_items_count\":3,     \"S_exemples\":1,     \"I_exemples\":4   } } </pre>	<pre> "flan-u12":{   "best":{     "prompt":"You are assisting a chemist in classifying paragraphs from     \"Paragraph1\": \"Despite obvious advantages, the systematic study of     \"Label1\": \"I\",     \"Paragraph2\": \"Jiang' group reported a series of 2D COFs with 1D of     \"Label2\": \"I\",     \"Paragraph3\": \"Firstly, 9.08 g of 2-methylimidazole was dissolved i     \"Label3\": \"S\",     \"Paragraph4\": \"As mentioned above, COFs are a class of crystalline     \"Label4\": \"I\",     \"Paragraph5\": \"Two porous organic polymers (POPs) incorporating a 1     \"Label5\": \"I\",     \"accuracy\":0.9777777777777777,     \"f1_score\":0.977733532437365,     \"run_number\":4,     \"correct_items\":44,     \"incorrect_items_count\":1,     \"S_exemples\":1,     \"I_exemples\":4   },   \"worst\":{     \"prompt\":\"You are assisting a chemist in classifying paragraphs from     \"Paragraph1\": \"Electrically conductive metal organic frameworks (MC     \"Label1\": \"I\",     \"Paragraph2\": \"Post-synthetic modification has also shown potential     \"Label2\": \"I\",     \"Paragraph3\": \"Preparation of CdMOF-1. A mixture of Cd(NOS<sub>3</sub>)<sub>2</sub>(     \"Label3\": \"S\",     \"Paragraph4\": \"The synthesis of NOS<sub>2</sub>-MIL-53(Cu) was carried out     \"Label4\": \"S\",     \"Paragraph5\": \"Synthesis of [Co(H<sub>2</sub>SO<sub>4</sub>)<sub>2</sub>](HCOO)<sub>2</sub> \\u00b     \"Label5\": \"S\",     \"accuracy\":0,     \"f1_score\":0,     \"run_number\":21,     \"correct_items\":0,     \"incorrect_items_count\":0,     \"S_exemples\":3,     \"I_exemples\":2   } } </pre>
<pre> "granite-20b-code-instruct":{   "best":{     "prompt":"You are assisting a chemist in classifying paragraphs from s     \"Paragraph1\": \"Syntheses of {[Sr<sub>2</sub>(MTA)(H<sub>2</sub>O)]u00b7H<sub>2</sub>O]u00b74DMF)n (UP     \"Label1\": \"S\",     \"Paragraph2\": \"Since amines readily adsorb pollutants, boosting the p     \"Label2\": \"I\",     \"Paragraph3\": \"4.1. Synthesis of UiO-66-(COOH)<sub>2</sub> \$ 10 mmol 1,2,4,5-     \"Label3\": \"S\",     \"Paragraph4\": \"Two porous organic polymers (POPs) incorporating a fer     \"Label4\": \"I\",     \"Paragraph5\": \"The potential for using MOF materials to remove fluori     \"Label5\": \"I\",     \"accuracy\":0.7333333333333333,     \"f1_score\":0.7473698338833187,     \"run_number\":28,     \"correct_items\":33,     \"incorrect_items_count\":12,     \"S_exemples\":2,     \"I_exemples\":3   },   \"worst\":{     \"prompt\":\"You are assisting a chemist in classifying paragraphs from s     \"Paragraph1\": \"Post-synthetic modification has also shown potential:     \"Label1\": \"I\",     \"Paragraph2\": \"Preparation of MWCNTx@ZIF-67 precursor: first, differe     \"Label2\": \"S\",     \"Paragraph3\": \"Powder XRD patterns were recorded on samples obtained     \"Label3\": \"I\",     \"Paragraph4\": \"3.3. Preparation of T-Ni(OH)<sub>2</sub>@TiO<sub>2</sub> and Ni(OH     \"Label4\": \"S\",     \"Paragraph5\": \"Proteins' secondary structure elucidates key character     \"Label5\": \"I\",     \"accuracy\":0.2666666666666666,     \"f1_score\":0.3226337448559671,     \"run_number\":92,     \"correct_items\":12,     \"incorrect_items_count\":33,     \"S_exemples\":2,     \"I_exemples\":3   } } </pre>	<pre> "granite-34b-code-instruct":{   "best":{     "prompt":"You are assisting a chemist in classifying paragraphs from     \"Paragraph1\": \"4.1. Synthesis of UiO-66-(COOH)<sub>2</sub> \$ 10 mmol 1,2,4,     \"Label1\": \"S\",     \"Paragraph2\": \"Proteins' secondary structure elucidates key charact     \"Label2\": \"I\",     \"Paragraph3\": \"To study the role of reactive oxygen species (ROS) c     \"Label3\": \"I\",     \"Paragraph4\": \"In the first instance the synthesis of Cu-BDC was ce     \"Label4\": \"S\",     \"Paragraph5\": \"Despite obvious advantages, the systematic study of     \"Label5\": \"I\",     \"accuracy\":0.9111111111111111,     \"f1_score\":0.9218290079766318,     \"run_number\":11,     \"correct_items\":41,     \"incorrect_items_count\":4,     \"S_exemples\":2,     \"I_exemples\":3   },   \"worst\":{     \"prompt\":\"You are assisting a chemist in classifying paragraphs from     \"Paragraph1\": \"Previously, it was noted in [29,31,32,48] that the s     \"Label1\": \"I\",     \"Paragraph2\": \"Since amines readily adsorb pollutants, boosting the     \"Label2\": \"I\",     \"Paragraph3\": \"2.2. Synthesis of {[Tb<sub>5</sub>SL<sub>6</sub>(OH)<sub>3</sub>](HS<sub>2</sub>)     \"Label3\": \"S\",     \"Paragraph4\": \"Nanocellulose can be subdivided into cellulose nanof     \"Label4\": \"I\",     \"Paragraph5\": \"Powder XRD patterns were recorded on samples obtaine     \"Label5\": \"I\",     \"accuracy\":0.3333333333333333,     \"f1_score\":0.3026573426573427,     \"run_number\":61,     \"correct_items\":15,     \"incorrect_items_count\":30,     \"S_exemples\":1,     \"I_exemples\":4   } } </pre>

Table 6: (continued)

<pre> "llama-3-405b-instruct":{   "best":{     "prompt":"You are assisting a chemist in classifying paragraphs fr     \"Paragraph1\": 3.3. Preparation of T-Ni(OH)_{2}S_{2} and N     \"Label1\": \"S\",     \"Paragraph2\": First, 25 mg of Ni(Ac)_{2}S_{2} was     \"Label2\": \"S\",     \"Paragraph3\": Researchers have conducted extensive studies on we     \"Label3\": \"I\",     \"Paragraph4\": Stability tests upon light illumination with diffe     \"Label4\": \"I\",     \"Paragraph5\": Synthesis of [Co(HS_{2}SO)_{2}(HCOO)_{2}] \u00     \"Label5\": \"S\",     \"accuracy\":0.9111111111111111,     \"f1_score\":0.9534872534872535,     \"run_number\":82,     \"correct_items\":41,     \"incorrect_items_count\":4,     \"S_examples\":3,     \"I_examples\":2   },   \"worst\":{     \"prompt\":\"You are assisting a chemist in classifying paragraphs fr     \"Paragraph1\": Post-synthetic modification has also shown potenti     \"Label1\": \"I\",     \"Paragraph2\": 2.2. Synthesis of {[Tb_{5}SL_{6}(OH)_{3}(HS_{2}     \"Label2\": \"S\",     \"Paragraph3\": 3.3. Preparation of T-Ni(OH)_{2}S_{2} and N     \"Label3\": \"S\",     \"Paragraph4\": Electrically conductive metal organic frameworks (     \"Label4\": \"I\",     \"Paragraph5\": {[Zn(bim)(bdc)] \u00b7 0.8DMF \u00b7 0.4EtOH \u00b     \"Label5\": \"S\",     \"accuracy\":0.0,     \"f1_score\":0.0,     \"run_number\":77,     \"correct_items\":0,     \"incorrect_items_count\":45,     \"S_examples\":3,     \"I_examples\":2   }, } </pre>	<pre> "llama-3-70b-instruct":{   "best":{     "prompt":"You are assisting a chemist in classifying paragraphs from     \"Paragraph1\": 2.2. Synthesis of the Five MOFs For the Cu-BTC synth     \"Label1\": \"S\",     \"Paragraph2\": Two porous organic polymers (POPs) incorporating a f     \"Label2\": \"I\",     \"Paragraph3\": 2.2. Synthesis of {[Tb_{5}SL_{6}(OH)_{3}(HS_{2}     \"Label3\": \"S\",     \"Paragraph4\": Preparation of MIL-68 (In). First, 0.24 g of CS_{8}S     \"Label4\": \"S\",     \"Paragraph5\": First, 25 mg of Ni(Ac)_{2}S_{2} was di     \"Label5\": \"S\",     \"accuracy\":1.0,     \"f1_score\":1.0,     \"run_number\":1,     \"correct_items\":45,     \"incorrect_items_count\":0,     \"S_examples\":4,     \"I_examples\":1   },   \"worst\":{     \"prompt\":\"You are assisting a chemist in classifying paragraphs from     \"Paragraph1\": Nanocellulose can be subdivided into cellulose nanof     \"Label1\": \"I\",     \"Paragraph2\": Syntheses of {[Sr2(MTA)(H2O)]\u00b7H2O\u00b74DMF)n (     \"Label2\": \"S\",     \"Paragraph3\": Jiang' group reported a series of 2D COFs with 1D op     \"Label3\": \"I\",     \"Paragraph4\": Rechargeable (also called 'secondary') batteries con     \"Label4\": \"I\",     \"Paragraph5\": However, unlike sMMO, many properties of pMMO are no     \"Label5\": \"I\",     \"accuracy\":0.7333333333333333,     \"f1_score\":0.7066666666666668,     \"run_number\":33,     \"correct_items\":33,     \"incorrect_items_count\":12,     \"S_examples\":1,     \"I_examples\":4   }, } </pre>
<pre> "mistral-large":{   "best":{     "prompt":"You are assisting a chemist in classifying paragraphs for     \"Paragraph1\": 2.2. Fabrication of SL Ce-BTC MOF NS The method prop     \"Label1\": \"S\",     \"Paragraph2\": Stability tests upon light illumination with differ     \"Label2\": \"I\",     \"Paragraph3\": Preparation of CdMOF-1. A mixture of Cd(NO_{3})_{2}S_{2}     \"Label3\": \"S\",     \"Paragraph4\": Synthesis of [Co(HS_{2}SO)_{2}(HCOO)_{2}] \u00b7     \"Label4\": \"S\",     \"Paragraph5\": 2.2. Synthesis of the Five MOFs For the Cu-BTC synt     \"Label5\": \"S\",     \"accuracy\":1.0,     \"f1_score\":1.0,     \"run_number\":82,     \"correct_items\":45,     \"incorrect_items_count\":0,     \"S_examples\":4,     \"I_examples\":1   },   \"worst\":{     \"prompt\":\"You are assisting a chemist in classifying paragraphs for     \"Paragraph1\": Oxidative stress is an imbalance between reactive o     \"Label1\": \"I\",     \"Paragraph2\": 2.2. Synthesis of {[Tb_{5}SL_{6}(OH)_{3}(HS_{2}     \"Label2\": \"S\",     \"Paragraph3\": Proteins' secondary structure elucidates key charact     \"Label3\": \"I\",     \"Paragraph4\": Preparation of MWCNTx@ZIF-67 precursor: first, diffe     \"Label4\": \"S\",     \"Paragraph5\": In 2017, Campagna, Hanan and their colleagues report     \"Label5\": \"I\",     \"accuracy\":0.7777777777777778,     \"f1_score\":0.7649393999569561,     \"run_number\":45,     \"correct_items\":35,     \"incorrect_items_count\":10,     \"S_examples\":2,     \"I_examples\":3   }, } </pre>	<pre> "mistral-8x7b-instruct-v01":{   "best":{     "prompt":"You are assisting a chemist in classifying paragraphs 1     \"Paragraph1\": The construction of MIL-101@SiO_{2} was carried     \"Label1\": \"S\",     \"Paragraph2\": Preparation of CdMOF-1. A mixture of Cd(NO_{3})_{2}S_{2}     \"Label2\": \"S\",     \"Paragraph3\": The sonochemical method is a synthesis procedure     \"Label3\": \"I\",     \"Paragraph4\": Despite obvious advantages, the systematic study     \"Label4\": \"I\",     \"Paragraph5\": 2.2. Synthesis of {[Tb_{5}SL_{6}(OH)_{3}(HS_{2}     \"Label5\": \"S\",     \"accuracy\":1.0,     \"f1_score\":1.0,     \"run_number\":30,     \"correct_items\":45,     \"incorrect_items_count\":0,     \"S_examples\":3,     \"I_examples\":2   },   \"worst\":{     \"prompt\":\"You are assisting a chemist in classifying paragraphs 1     \"Paragraph1\": The construction of MIL-101@SiO_{2} was carried     \"Label1\": \"S\",     \"Paragraph2\": Previously, it was noted in [29,31,32,48] that th     \"Label2\": \"I\",     \"Paragraph3\": Firstly, 9.08 g of 2-methylimidazole was dissolve     \"Label3\": \"S\",     \"Paragraph4\": First, 25 mg of Ni(Ac)_{2}S_{2} was     \"Label4\": \"S\",     \"Paragraph5\": Proteins' secondary structure elucidates key char     \"Label5\": \"I\",     \"accuracy\":0.5333333333333333,     \"f1_score\":0.6137467700258398,     \"run_number\":27,     \"correct_items\":24,     \"incorrect_items_count\":21,     \"S_examples\":3,     \"I_examples\":2   }, } </pre>

550 **Information extraction** Table 7 shows the JSONs that include the best and worst paragraphs  
 551 selected as examples for each model.

Table 7: Partial JSON with the best and worse examples for *Information extraction* module.

```

"meta-llama/llama-3-405b-instruct":{ ☐
  "best":{ ☐
    "Paragraph1":"Synthesis of UPJS-15 and UPJS-16 Syntheses of {[Sr2(MTA)(H2O)]-H2O
    "Paragraph2":" The construction of MIL-101@SiO2 was carried out according to the
  },
  "worst":{ ☐
    "Paragraph1":" The Construction of Cu-As MOF The process of creating a copper as
    "Paragraph2":" Synthesis of the Five MOFs: For the Cu-BTC synthesis, 1.925 g Cu(
  }
},
"mistralai/mistral-large":{ ☐
  "best":{ ☐
    "Paragraph1":" NiDMOF was synthesized by a solvothermal reaction according to th
    "Paragraph2":" 3.3. Preparation of T-Ni(OH)2@TiO2 and Ni(OH)2@TiO2 Photoanodes:
  },
  "worst":{ ☐
    "Paragraph1":"Preparation of MWCNTx@ZIF-67 precursor: first, different masses of
    "Paragraph2":"In the first instance the synthesis of Cu-BDC was carried out foll
  }
},
"meta-llama/llama-3-70b-instruct":{ ☐
  "best":{ ☐
    "Paragraph1":" Fabrication of SL Ce-BTC MOF NS The method proposed by Liu and co
    "Paragraph2":" NiDMOF was synthesized by a solvothermal reaction according to th
  },
  "worst":{ ☐
    "Paragraph1":" The Construction of Cu-As MOF The process of creating a copper as
    "Paragraph2":" Synthesis of ZIF-67 All the chemicals utilized in this study were
  }
},
"ibm/granite-34b-code-instruct":{ ☐
  "best":{ ☐
    "Paragraph1":" Synthesis of Zn-MOF Bis(8midazole-1-yl)methane was synthesized
    "Paragraph2":" Synthesis of poly(1,10-ferrocenediyl-bis(metylphosphinate) Zn(
  },
  "worst":{ ☐
    "Paragraph1":" 3.3. Preparation of T-Ni(OH)2@TiO2 and Ni(OH)2@TiO2 Photoanode
    "Paragraph2":" The Construction of Cu-As MOF The process of creating a copper
  }
},
"mistralai/mixtral-8x7b-instruct-v01":{ ☐
  "best":{ ☐
    "Paragraph1":" Preparation of MWCNTx@ZIF-67 precursor: first, different masse
    "Paragraph2":" Synthesis of Zn-MOF Bis(8midazole-1-yl)methane was synthesized
  },
  "worst":{ ☐
    "Paragraph1":" 3.3. Preparation of T-Ni(OH)2@TiO2 and Ni(OH)2@TiO2 Photoanode
    "Paragraph2":" The Construction of Cu-As MOF The process of creating a copper
  }
},
"ibm/granite-20b-code-instruct":{ ☐
  "best":{ ☐
    "Paragraph1":" Preparation of MWCNTx@ZIF-67 precursor: first, different masse
    "Paragraph2":" In the first instance the synthesis of Cu-BDC was carried out
  },
  "worst":{ ☐
    "Paragraph1":" The conventional laccase-ZIF-8 biocomposites (Lac@ZIF-8) were
    "Paragraph2":" The Construction of Cu-As MOF The process of creating a copper
  }
}

```