

This is a “preproof” accepted article for *Quantitative Plant Biology*.

This version may be subject to change during the production process.

10.1017/qpb.2025.10021

1 **Small language models enable rapid and accurate extraction of structured data from unstructured text: an**
2 **example with plants and their specialized metabolites**

3
4
5 Lucas Busta^{1,*} and Alan R. Oyler¹

6
7 ¹Department of Chemistry and Biochemistry, University of Minnesota Duluth, USA

8
9 * To whom correspondence should be addressed: bust0037@d.umn.edu
10
11
12

13 **Abstract**

14
15 Transformer-based large language models are receiving considerable attention because of their ability to analyze
16 scientific literature. Small language models (SLMs), however, also have potential in this area, have smaller compute
17 footprints, and allow users to keep data in-house. Here, we quantitatively evaluate the ability of SLMs to: (i) score
18 references according to project-specific relevance and (ii) extract and structuring data from unstructured sources
19 (scientific abstracts). By comparing SLMs’ outputs against those of a human on hundreds of abstracts, we found that
20 (i) SLMs can effectively filter literature and extract structured information relatively accurately (error rates as low as
21 10%), but not with perfect yield (as low as 50% in some cases), (ii) that there are tradeoffs between accuracy, model
22 size, and computing requirements, and (iii) that clearly written abstracts are needed to support accurate data
23 extraction. We recommend advanced prompt engineering techniques, full-text resources, and model distillation as
24 future directions.
25

This is an Open Access article, distributed under the terms of the Creative Commons Attribution-NonCommercial-NoDerivatives licence (<http://creativecommons.org/licenses/by-nc-nd/4.0/>), which permits non-commercial re-use, distribution, and reproduction in any medium, provided the original work is unaltered and is properly cited. The written permission of Cambridge University Press must be obtained for commercial re-use or in order to create a derivative work.

26 **1. Introduction**

27

28 Language models are emerging as powerful tools for a wide array of tasks, with a particularly promising role in
29 processing scientific literature (Agathokleous et al. 2024; Jin et al. 2024; Lam et al. 2024; Simon et al. 2024; Busta
30 et al. 2024b; Knapp et al. 2024b). Scientific articles compile results from decades, if not centuries, of effort by
31 scientists worldwide. However, the automation of classification, summarization, and data extraction tasks related to
32 this literature remains a challenge because natural language is a complex data type. In other fields with intricate
33 data, such as image and sound, a proven strategy is to build mathematical models of the input data type that can then
34 be leveraged to summarize, classify, or otherwise manipulate the input. Modeling natural language is a long-
35 standing field of study, but recently, the development and increase in accessibility of transformer-based language
36 models have led to substantial advances in our language processing ability. Perhaps we can solve some of the many
37 challenges with automated processing of scientific literature by applying transformer-based language models.

38

39 A considerable number of recent investigations are focused on applying large language models to scientific literature
40 (Jin et al. 2024; Busta et al. 2024b; Shiu and Lehti-Shiu 2024; Sarumi and Heider 2024; Knapp et al. 2024b). For
41 example, large language models have been utilized to perform tasks such as text classification, text summarization,
42 and question answering (Dalal et al. 2024; Riordan 2024; Shiu and Lehti-Shiu 2024; Guo et al. 2023; Yin et al.
43 2019). Generally, these large models require significant memory—hundreds of gigabytes—to store high billions or
44 trillions of parameters required at runtime. However, a diverse range of language models exists beyond the popular
45 large models from, for example, OpenAI, Anthropic, Google, and Mistral. In particular, small language models
46 (SLMs) have gained attention due to their smaller sizes (low billions or even just millions of parameters) and thus
47 reduced computing requirements. Furthermore, though the small models are not as general purpose as the large
48 models, the emerging evidence suggesting the small models are effective in various, albeit specific natural language
49 processing tasks (Lepagnol et al. 2024; Guo et al. 2023; Zhu 2024; Lewis et al. 2019). Thus, these small language
50 models are intriguing because they suggest that individual scientists could use them on ordinary personal computing
51 devices to potentially enhance scientific literature processing tasks. Importantly, running the small models on local
52 hardware also avoids passing private and/or copyrighted content to large language model companies, which is
53 prohibited by many research institutions and industrial organizations.

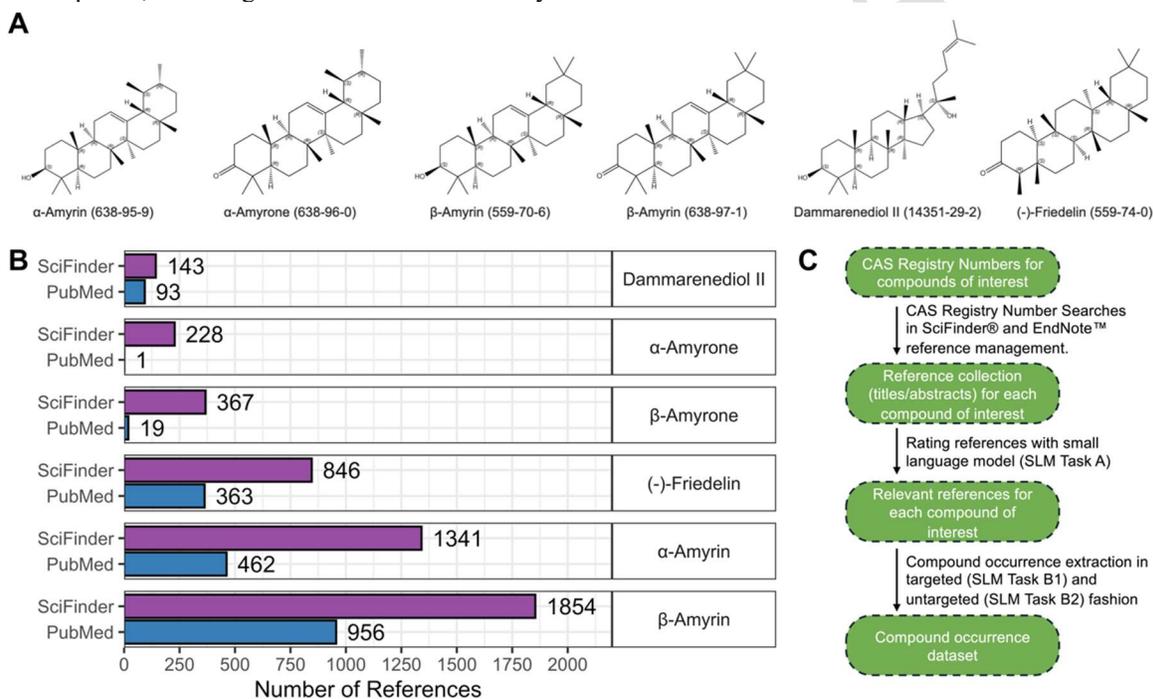
54

55 In the present work, we aimed to develop and evaluate a proof-of-concept small language model processes to
56 support the expansion of databases that document plants and the specialized metabolites that each may produce.
57 Other databases have been created in the past to document this same type of information (Zeng et al. 2024; Tay et al.
58 2023; Gallo et al. 2023; Rutz et al. 2022; Sorokina and Steinbeck 2020; Nguyen-Vo et al. 2020; Yang et al. 2019;
59 Chen et al. 2017; Xie et al. 2015), but these databases, so far, do not leverage the potential provided by language
60 models. We experimented with models to conduct two major tasks: (i) scoring articles based on their relevance to
61 need-specific criteria (in this case, whether they contained reports of a specific plant making a specific chemical
62 compound) and (ii) extracting and structuring information on the occurrence of specific chemical compounds in
63 specific plant species. We tested a dozen language models' abilities on these tasks by manually reading, labelling,
64 and extracting data from more than 100 to more than 1000 scientific abstracts, depending on the task, then measured
65 the models' ability to perform those same tasks. Overall, our findings indicate that small language models, while not
66 perfect, effectively aid in filtering scientific literature references and in extracting data. We recommend that
67 researchers both experiment with these models and monitor for updates in literature processing software that
68 incorporate language model-enabled features.

69

70 **2. Results and Discussion**

71
 72 To develop and evaluate a potential role for small language models in creating a phytochemical occurrence database,
 73 we assessed such models' abilities with regard to two tasks: (i) to quickly score references according to whether the
 74 reference reports the occurrence of a specific compound in a specific plant species (Task 1, Section 2.1), and (ii) to
 75 evaluate language models' ability to extract an experimentally-supported compound occurrence dataset (Task 2,
 76 Section 2.2). For these investigations, we chose to use six triterpenoid compounds as test cases (**Fig. 1A**). The six
 77 triterpenoid test cases under study here presented a challenge because they have been mentioned in the literature
 78 (going back to the 1960s) by many names. Indeed, CAS SciFinder® indicates that a total of more than 52 names has
 79 been associated with these six compounds, potentially complicating efforts to retrieve references describing the
 80 occurrence of specific plant chemicals. Fortunately, triterpenoids (and the vast majority of all other chemical
 81 entities) are identified explicitly by their CAS Registry® numbers (**Fig. 1A**), which means that references to a given
 82 compound that use varied nomenclature can be collected simultaneously and non-ambiguously when using CAS
 83 Registry® number-based search strategies. While other identification number systems exist, such as PubChem® and
 84 LOTUS numbers, these alternate systems are not as comprehensive as CAS Registry® numbers. Thus, where
 85 possible, searches with identification numbers, as opposed to common names, are preferred because this approach
 86 ensures not only that a broader array of references is retrieved, but also that those reference relate to one and the
 87 same compound, including the correct stereochemistry.



88
 89 **Figure 1. Comparison of SciFinder® versus PubMed® as a data source and schematic of the small language**
 90 **model workflow for retrieving compound-species associations from literature. A. Structures, common names, and**
 91 **CAS Registry® Numbers for the six triterpenoid compounds used as test cases in our small language model**
 92 **development and evaluation work. B. Bar plot comparing the number of references (x-axis) found by SciFinder®**
 93 **and PubMed® (y-axis) for the six different triterpenoids (vertically arranged panels) studied in this work. Each bar**
 94 **represents the number of references found by the indicated search tool for a particular triterpenoid. The absolute**
 95 **number of references found is shown in text to the right of each bar. Bars are color coded according to search tool**
 96 **(SciFinder® in purple and PubMed® in blue). SciFinder® searches were conducted using CAS Registry® Numbers,**
 97 **while PubMed® (which does not generally use these registry numbers) searches were conducted using compound**
 98 **common names. C. Schematic for the workflow we developed to extract compound occurrence data from information**
 99 **in the literature. Files or information are shown in green bubbles, while steps or actions are shown as arrows. The**
 100 **workflow consists of searching the literature with SciFinder® based on CAS Registry® numbers then creating a**
 101 **repository of references and associated full text PDF files in an EndNote™ database; then filtering references for**
 102 **those of highest task-specific relevance (SLM Task A) and finally extracting compound occurrence data in either a**
 103 **targeted (SLM Task B1) or untargeted (SLM Task B2) fashion. Abbreviations: SLM: small language model.**

104 To obtain references describing our six triterpenoids of interest, we used CAS Registry® numbers to search
105 SciFinder®, which, although requiring a subscription, allows the user to enter a CAS Registry® number and then
106 navigate directly to literature references that relate to that specified compound. PubMed®, although providing open
107 access, does not generally support searches based upon CAS Registry® numbers or PubChem ID numbers, so we
108 conducted searches in PubMed using compound common names. We first considered the two most common
109 compounds in our case study set, α and β -amyrin. In SciFinder®, we found over 1,340 and more than 1,850 hits for
110 these two compounds, respectively, compared to fewer than 500 and 1,000 hits in PubMed® (Fig. 1B). Results were
111 similar for the other four triterpenoid test cases (Fig. 1B). In total, ~3,200 SciFinder® references were retrieved
112 using our searches, while ~1,500 references were retrieved by PubMed®. Therefore, we used SciFinder®-retrieved
113 references to develop and evaluate small language model-based reference ranking and occurrence dataset extraction
114 processes (Fig. 1C).

115 116 2.1 SLM Task A: Rating references according to relevance with a small language model

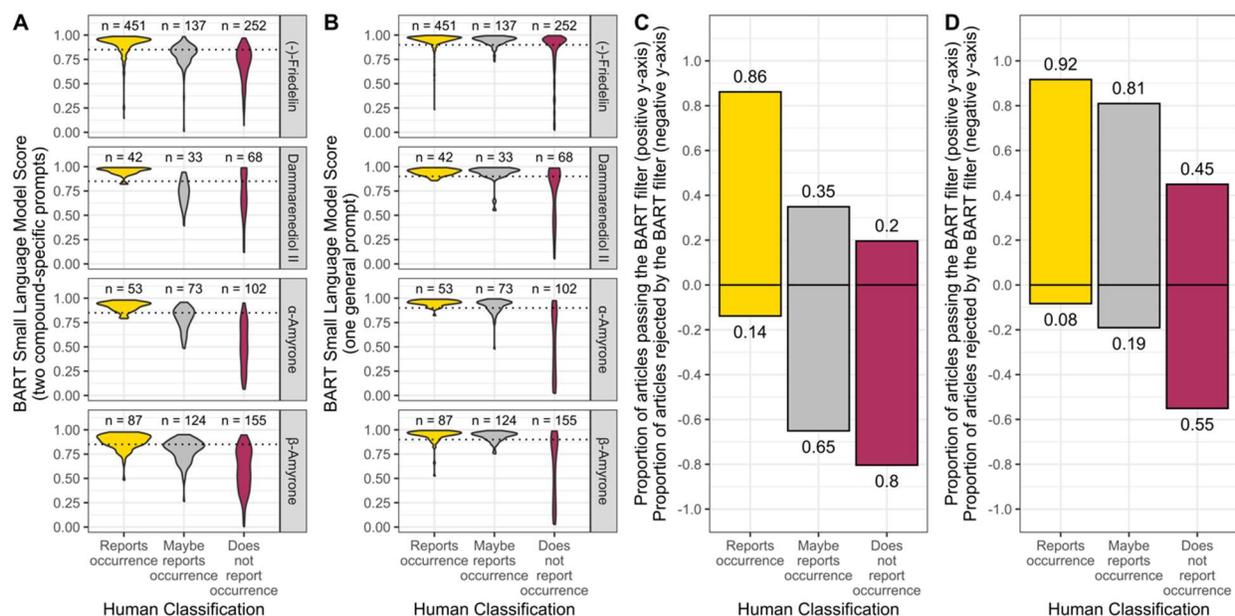
117
118 At this stage in the present work, we had used SciFinder® to collect more than 3,000 references associated with one
119 or more of the six triterpenoids that comprised our test cases for compound occurrence data collection. Our first aim
120 was to determine the efficacy of small language models with respect to filtering the references for articles of interest.
121 In this case, our interest was in articles that reported phytochemical occurrences (i.e., evidence for a specific plant
122 species producing a specific chemical compound). To establish a benchmark against which to evaluate small
123 language model performance we read more than 1,500 of the references in our collection, including their titles and
124 abstracts, and classified each as “reporting an occurrence”, “maybe reporting an occurrence”, or “not reporting an
125 occurrence” (Supplemental File 1). These human-read citations included all the reference citations for α -amyrone,
126 β -amyrone, dammarenediol II, as well as (–)-friedelin. For an article to be considered as “reporting an occurrence”
127 its title or abstract needed to indicate that the article in question provided experimental evidence for the presence of
128 a particular plant chemical in a particular plant species. Articles whose titles or abstracts merely contained co-
129 occurrences of a plant chemical name and a plant species name without indicating that there was experimental
130 evidence for an association between the two were classified as “not reporting an occurrence”. Citations that did not
131 explicitly indicate that their articles contained experimental evidence for a compound’s occurrence but instead
132 implied that such evidence might be present in the full text (to which we did not have access) were classified as
133 “maybe reporting an occurrence”. Of the 1,558 references that we read, 720 were classified as “reporting an
134 occurrence” (46%), 332 were classified as “maybe reporting an occurrence”, (21%) and 506 were classified as “not
135 reporting an occurrence” (33%).

136
137 We next evaluated how well language models could classify references according to whether they reported the
138 occurrence of a phytochemical using the 1,558 manually classified references as a ground-truth set. We used the
139 bart-large-mnli model, selected because it is one of the most downloaded on Huggingface.co, a major hub for open-
140 source language development, largely due to its versatility and high speed – we found that it could process 45,000
141 articles / hour, a desirable characteristic for a model that will be used to filter inputs into a multi-step processing
142 pipeline. This small language model is employed by providing it with a body of text and then one or more classifier
143 phrases. The model then assigns a score to each phrase to indicate how closely that phrase relates to the provided
144 text. The bart-large-mnli model card (i.e., the instruction manual) suggests presenting the model with a classifier
145 phrase framed as a hypothesis (e.g., “This text is about politics”). Accordingly, we investigated phrases such as
146 “Amyrin is present in plants” as well as paired phrases in which a hypothesis was matched with the exact negative
147 (i.e., “Amyrin is present in plants” and “Amyrin is not present in plants”). Our early experiments showed that
148 composite scores derived from the pairs’ individual scores improve the signal-to-noise ratio in the classification task.
149 Furthermore, we noted that multiple compound names could be included in these positive and negative phrases (for
150 instance, “friedelin, friedooleanan-3-one, friedelan-3-one, friedelanone, or friedeline is found in plants”; full
151 classifier phrase details are provided in Supplemental File 2). In future large-scale operations, a single, general
152 classifier phrase, which is not based on compound names, would be preferred if the performance was comparable to
153 that of our specific classifier phrase system, which is based on compound names. Therefore, we also tested the more
154 general classifier phrase, “The text discusses plants that contain specific compounds.”

155
156 Using the two classifier phrase approaches described in the previous paragraph, we instructed the bart-large-mnli
157 model to assign two scores to each of the 1,558 references, one composite score from the binary/two classifier
158 phrase system, as well as a score for the general classifier phrase. Composite scores (means and standard deviations)
159 for, respectively, references that reported occurrences / maybe reported occurrences / did not report occurrences for

160 (-)-friedelin were 0.9 ± 0.1 , 0.8 ± 0.1 , and 0.7 ± 0.1 (Fig. 2A, top panel). Results were similar for the other three
 161 triterpenoids (Fig. 2A). Scores from the general classifier phrase for, respectively, references that reported
 162 occurrences / maybe reported occurrences / did not report occurrences for (-)-friedelin were 0.9 ± 0.06 , 0.9 ± 0.04 ,
 163 and 0.8 ± 0.2 (Fig. 2B, top panel), and again, results were similar for the other three triterpenoids (Fig. 2B). This
 164 illustrates that the two-classifier phrase system and the general classifier phrase system both worked comparably
 165 well among references describing four different triterpenoid compounds and may also to a similar extent for
 166 compounds other than triterpenoids.

167



168

169

170 **Figure 2: Performance of small language models on a reference relevance ranking task. A and B.** Violin plot
 171 showing the score (BART Small Model Language Score, y-axis) assigned to references by the bart-large-mnli small
 172 language model. Scores range from zero (low relevance) to one (high relevance) and indicate the relevance of a
 173 given reference to a user-defined natural language criterion. In panel A, the score is derived from two, chemical
 174 compound-specific criteria (full details in methods section), while in panel B, the score is derived from a single,
 175 generic criterion ("chemical compounds are found in plants"). In both panels, scores are broken out according to
 176 whether the reference was labeled by a human as "reporting an occurrence", "maybe reporting an occurrence",
 177 "not reporting an occurrence" of a specific chemical compound in a specific species (x-axis). The number of
 178 references belonging to each group are shown above each violin. In panel A, the dotted line represents a threshold of
 179 0.85 and in panel B, the dotted line represents a threshold of 0.9; details of thresholds discussed in main text). **C and**
 180 **D.** Column plot showing the proportion of references (y-axis) from each human labeled category ("reporting an
 181 occurrence", "maybe reporting an occurrence", or "not reporting an occurrence"; x-axis) that would be retained if
 182 a threshold small language model score was used for filtering references. The proportion of each column in the
 183 positive y space indicates the fraction of references that would pass the filter and be retained, while the proportion
 184 of each column in the negative y space indicates the fraction of references that would be rejected by the filter and
 185 eliminated. Exact proportions are shown in numbers above and below each column. In panel C the threshold is 0.85,
 186 based on two-prompt scoring, while in panel D the threshold is 0.9, based on single, general prompt scoring (details
 187 in main text and methods section). For example, if a score of 0.85 were used as a threshold with which to filter
 188 references that had been scored using the two-prompt small language model scoring system, then 86% of references
 189 reporting occurrences would be retained while 14% of such references would be rejected, 35% of references maybe
 190 reporting occurrences would be retained, while 65% of such references would be rejected, and 20% of references not
 191 reporting occurrences would be retained while 80% of such references would be rejected. In all panels A-D, colors
 192 correspond to the three human label categories ("reporting an occurrence", "maybe reporting an occurrence", "not
 193 reporting an occurrence"). BART stands for the bart-large-mnli small language model.

194

195 Next, we investigated the ability of these scores to act as a filter to separate articles of interest that report chemical
196 occurrences from those that did not report such occurrences. Thus, we examined the proportion of the former type
197 articles that would be retained if a threshold score were to be used as a filtering criterion for the reference collection
198 (i.e., if references with a score higher than a threshold were to be retained and those with a score lower than the
199 threshold were to be eliminated from the collection). Based on the distribution of scores assigned to articles that
200 reported chemical occurrences versus those that did not (Fig. 2A and B), we selected 0.85 as a threshold for the
201 specific two-prompt scores and 0.90 as a threshold for the general prompt-derived scores. With these thresholds, the
202 specific two-prompt scoring system acting as a filter would have retained 86% of the references that report
203 phytochemical occurrences (the references of interest in our study), and rejected 80% of the references that did not
204 report an occurrence (Fig. 2C). The general prompting system, with a 0.90 filtering threshold, would have retained
205 92% of the references reporting phytochemical occurrences and eliminated 55% of the references that did not report
206 occurrences (Fig. 2D). While both the two-prompt and general prompt filtering approach led to the retention and
207 rejection, respectively, of article of interest and not of interest, the two approaches handled articles we had labelled
208 as “maybe reports an occurrence” differently: the two-prompt approach kept only 35% of these, while the general
209 approach kept 81%. To learn more about these “maybe” references, we obtained and read 100 full text articles for
210 these references (those related to α -amyrone and dammarenediol II, Supplemental File 3). This manual inspection
211 revealed that approximately 65% percent of these “maybe” references contained reports of compound occurrence
212 data, which suggested that access to full text information will help create more comprehensive chemical occurrence
213 datasets. After manual re-annotation of the 100 articles based on full texts, we tested to see if the scores of
214 occurrence-reporting articles differed from articles that did not report occurrences, but there was no significant
215 difference in the scores. However, regardless of whether full texts are available or not, our results show that small
216 language model relevance scores provide a means to quickly (~45,000 references / hr.) and accurately (~80%
217 relevant articles kept, ~80% of irrelevant articles rejected) identify references that are most likely to provide the
218 information that a user might be seeking. This ability will be highly useful when dealing with many thousands of
219 references. Our data also indicate that there will likely be a benefit to developing more nuanced filtering approaches
220 to handle edge cases like the ‘maybe’ articles we identified here.

221

222 2.2: SLM Task B: Extracting compound occurrence data with language models

223

224 After filtering our collection of references to include only entries with high scores concerning phytochemical
225 occurrence data, we evaluated the ability of language models to extract experimentally supported compound
226 presence details. In this task, two steps can be envisioned: (i) a first step in which a model receives a body of text
227 including the title and abstract of a scientific article and (ii) a second step in which a model receives a query about
228 compound occurrences. For example, in the second step, we might ask the model: “Does the provided text offer
229 experimental evidence that *Arabidopsis thaliana* produces the chemical compound thalanelol?” This mode of
230 operation represents a **targeted approach**. A second mode of operation (for the second step) could be to pass a
231 language model a text passage containing the title and abstract of a scientific article and pose an open-ended query
232 such as: “List all of the plant species mentioned in the provided text and indicate which chemical compounds were
233 reported from each one as part of the experimental investigation described in the passage.” This second mode
234 represents an **untargeted approach**. Several advantages and disadvantages of each approach can be imagined from
235 the outset. For example, an untargeted approach does not require a preconceived set of chemical compounds or plant
236 species of interest about which to query the model, and a single untargeted query can potentially extract multiple
237 compound-occurrence data simultaneously. In contrast, one benefit of the targeted approach is the relative simplicity
238 of creating human-labeled data. Thus, a true/false answer about one plant/compound occurrence can be supplied by
239 the human or model instead of meticulously generating a complete list of such occurrences. Furthermore, a model's
240 rate of detecting true negative associations can be measured directly by comparing the model's response to plant and
241 compound names appearing in an abstract, without experimental association data, to the corresponding human
242 response. Thus, the targeted and untargeted approaches each offer distinct benefits, so we tested and herein present
243 results from both approaches. For either approach, a model must correctly distinguish between characters used in
244 chemical names (in the present study, especially Greek letters like α and β) and recognize the synonymous nature of
245 certain symbols and words (for example, that α -amyrin and alpha-amyrin are the same compound). Previous studies
246 have shown that Greek letters occupy their own positions in language model input spaces (Stevenson et al. 2025)
247 and that such models can reason over diverse alphabets (Maronikolakis et al. 2021), suggesting that modern
248 language models, in this regard at least, may be suited to the above-described approaches. We conducted preliminary
249 tests by asking each model (see model details below) a series of six questions like “are α -amyrin and alpha-amyrin
the same compound?”, “are α -amyrin and beta-amyrin the same compound?”, “are β -amyrin and beta-amyrin the

250

251 same compound”, and so forth. All but the smallest two models (gemma-3-1B-instruct and qwen-2.5-0.5B-instruct)
252 answered these questions with 100% accuracy, illustrating that detailed investigations of task/project-specific
253 assumptions should be empirically tested during the model selection step of a language model-based investigation.

254

255 **2.2.1 SLM Task B1: Targeted compound occurrence data extraction**

256

257 To evaluate the ability of large language models to extract compound occurrence data from scientific abstracts, we
258 first prepared and manually evaluated a set of candidate occurrences. For this effort, we used regular expression-
259 based pattern matching to identify accepted plant species names in the abstracts associated with the six triterpenoids
260 that comprised the present test case. We then compiled a data set containing three columns: the title and abstract of
261 each reference, the chemical compound linked to it (the SciFinder® search compound that retrieved that reference in
262 the first place), and accepted plant species name(s) found in that title or abstract. We manually evaluated 500
263 candidate associations and annotated each occurrence as positive (the abstract described experimental support for
264 the occurrence of that compound in that plant species) or a negative (the abstract did not provide such support). We
265 found that roughly 350 (71%) of the candidate associations were negatives, while around 150 (29%) were positives
266 (**Supplemental File 4**). With a set of human-labeled compound species or candidate compound species associations
267 in hand, we next turned to evaluating whether open-source language models could perform the same task. For this
268 task, we used open-source language models that accepted two types of prompts. The first prompt was a system
269 prompt that contained detailed instructions on how the model should generate an output. The second prompt (also
270 called user text) delivered content from which the model generated that output. We used the second prompt to
271 supply information on the candidate compound species association (title/abstract, compound name, and species
272 name) and the system prompt to convey detailed instructions on how the model was supposed to evaluate this given
273 information (full details in Methods).

274

275 Past research has shown that language models of different sizes vary in their ability to perform natural language
276 processing tasks (Brown et al. 2020; Kaplan et al. 2020), including tasks related to chemical occurrence data
277 extraction (Busta et al. 2024a). Accordingly, in evaluating their capacity for the present targeted occurrence
278 extraction task, we tested 12 language models of various scales, spanning 0.5 billion to 32 billion parameters (often
279 denoted 0.5B to 32B, **Fig. 3A**). These models included variants of different sizes from the Qwen family (Qwen: An
280 et al. 2025) (32B, 14B, 7B, and 0.5B), the Gemma family (Gemma et al. 2025) (27B, 12B, 4B, and 1B), and the Phi-
281 4 family (phi-4 14B and phi4-mini-instruct 4B) (Abdin et al. 2024). Each model was given the same system prompt
282 and all 500 candidate occurrences that had been previously examined manually. During these assessments, all
283 models were run at 16-bit precision, except gemma-3-27B-it-unsloth and phi-4-unsloth-bnb-4bit, which are
284 dynamically quantized instances operating at 4-bit precision (**Fig. 3A**). When reviewing the 500 candidate
285 associations, run times generally varied in direct proportion with size; qwen-2.5-32B-instruct handled about 200
286 references per hour, while qwen-2.5-0.5B-instruct surpassed 32,000 per hour (**Fig. 3A**). Notably, the quantized
287 variants processed references at speeds only slightly higher than their full-resolution counterparts (for example, the
288 4-bit phi-4-unsloth at 1,500 references / hr. and the 16-bit phi-4 at 1,200 per hour). These speeds will be important
289 when applying language model-based approaches to larger projects or the assembly of databases.

290

291 Alongside measuring how quickly various models processed 500 candidate associations, we also examined model
292 accuracy. To gauge that accuracy, we compared whether each model labeled every candidate association as positive
293 or negative against the corresponding human label. The results let us classify each model output as a true positive
294 (when the model labeled a candidate association as positive, matching the human label), a true negative (when both
295 the model and the human labeled it negative), a false positive (when the model labeled it positive but the human did
296 not), or a false negative (when the model labeled it negative but the human did not). Because 71% of the 500
297 candidate associations were negative, a high-performing model would have a true negative rate approaching 71%.
298 The true negative rates for the models tested ranged from 52% to 67%, with models containing more parameters
299 generally showing higher percentages (**Fig. 3B**). One exception was qwen-2.5-0.5B-instruct, which had a 0% true
300 negative rate, as it labeled all candidates occurrences as positive. These differences in true negative rates came with
301 parallel differences in false positive rates, since false positives arise when a model incorrectly labels a negative
302 result as positive. The false positive rate is one of the most important metrics for this task because those errors
303 represent fabricated occurrence data. In our experiments, larger models achieved lower false positive rates overall,
304 with qwen-2.5-32B-instruct and phi-4 showing the lowest values at 4% and 5%, respectively (**Fig. 3B**). Because
305 both were also the slowest and largest, there is a clear trade-off between parameter count and computational
306 requirements on one hand and task-specific accuracy on the other.

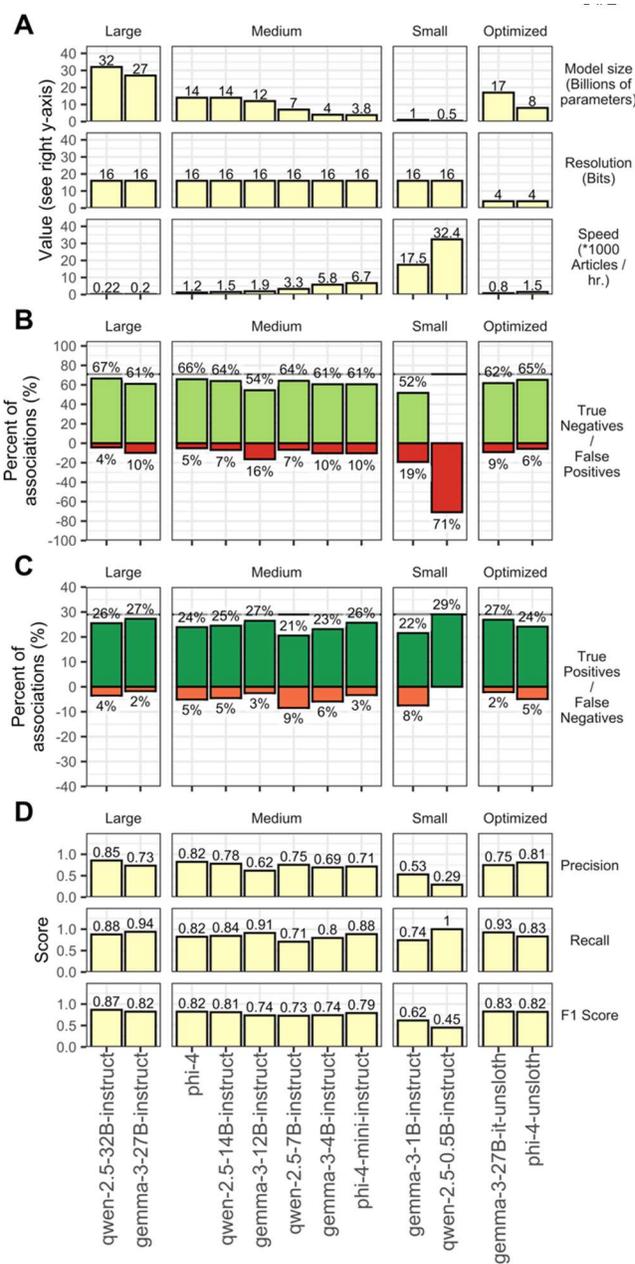


Figure 3: Performance of language models on a targeted compound occurrence data extraction task.

A. Bar plot showing various metrics (y axes in each row of panels) for different language models (x axis). The first row shows model size in billions of parameters, the second row shows model resolution in bits, the third row shows the speed with which a model processes references (using the prompt shown in the methods section) in units of 1000 references per hour. **B.** Bar plot showing the raw performance metrics of each model (false negative, false positive, true negative, and true positive rates). False negatives arise when a model erroneously marks a real compound occurrence as not being true. False positives arise when a model erroneously marks a simple textual occurrence of a compound name and species name as an occurrence data point. True negatives arise when a model correctly marks a simple textual occurrence of a compound name and species name as such, and not as an occurrence data point. True positives arise when a model correctly marks a compound occurrence as such. According to human evaluation of the 500 putative occurrences used to test the models, 71% of the putative occurrences were real (i.e. "positives"), and 29% of the putative occurrences were just textual co-occurrence (i.e. "negatives"). Thus, a perfect model would have, in this experiment, a 71% true negative rate and a 29% true positive rate. Bars are colored according to true/false positive/negative. **C.** Bar plot showing the processed performance metrics of each model. In the first row, the precision of each model is shown (the ratio of true positives to the sum of true positives and false positives). In the second row, the recall of each model is shown (the ratio of true positives to the sum of true positives and false negatives). In the third row, the F1 score is shown, which is the harmonic mean of the precision and recall. In A–C, models are organized into columns of panels by type (large: > 20 B parameters, medium: 1–20 B parameters, small: 0–1 B parameters, and optimized: 4-bit resolution models).

349 The models we tested here did not only vary in their (true negative)/(false positive) rates, but also in their (true
350 positive)/(false negative) rates. Since positive associations comprised 29% of the 500 candidate associations, a
351 perfect model in our experiment would have a 29% true positive rate. True positive rates among the models tested
352 here generally ranged from 21% to 27% (**Fig. 3B**). This variability did not correlate as strongly with model size as
353 did the (true negative)/(false positive) rates. For example, a large model (qwen-2.5-32B-instruct), two medium
354 models (gemma-3-12B-instruct and phi-4-mini-instruct (4B)), and one of the quantized models (gemma-3-27B-it-
355 unsloth) all had very similar true positive rates (26% or 27%, **Fig. 3B**). Note that the perfect true positive rate of
356 qwen-2.5-0.5B-instruct is a misleading statistic, since this model simply labeled all associations with which it was
357 presented as positive. To account for such potentially misleading rates, we computed precision and recall statistics.
358 Precision is calculated as the number of true positive results divided by the sum of true positive and false positive
359 results, which indicates how reliable the model is when it marks an association as positive. Recall is calculated as
360 the number of true positive results divided by the sum of true positive and false negative results, which reflects the
361 model's ability to correctly identify all actual positive associations. Excluding qwen-2.5-0.5B-instruct, precision
362 varied from 0.5 to as high as 0.85 and recall varied from 0.71 to 0.94 (**Fig. 3C**). We also computed F1 scores, which
363 are the harmonic mean of precision and recall, to provide a single metric to balance both reliability (precision) and
364 completeness (recall). F1 scores (excluding qwen-2.5-0.5B-instruct) ranged from 0.62 (gemma-3-1B-instruct) to
365 0.87 (qwen-2.5-32B-instruct) and varied, again, according to model size, which reinforced the importance of that
366 parameter in task-specific accuracy.

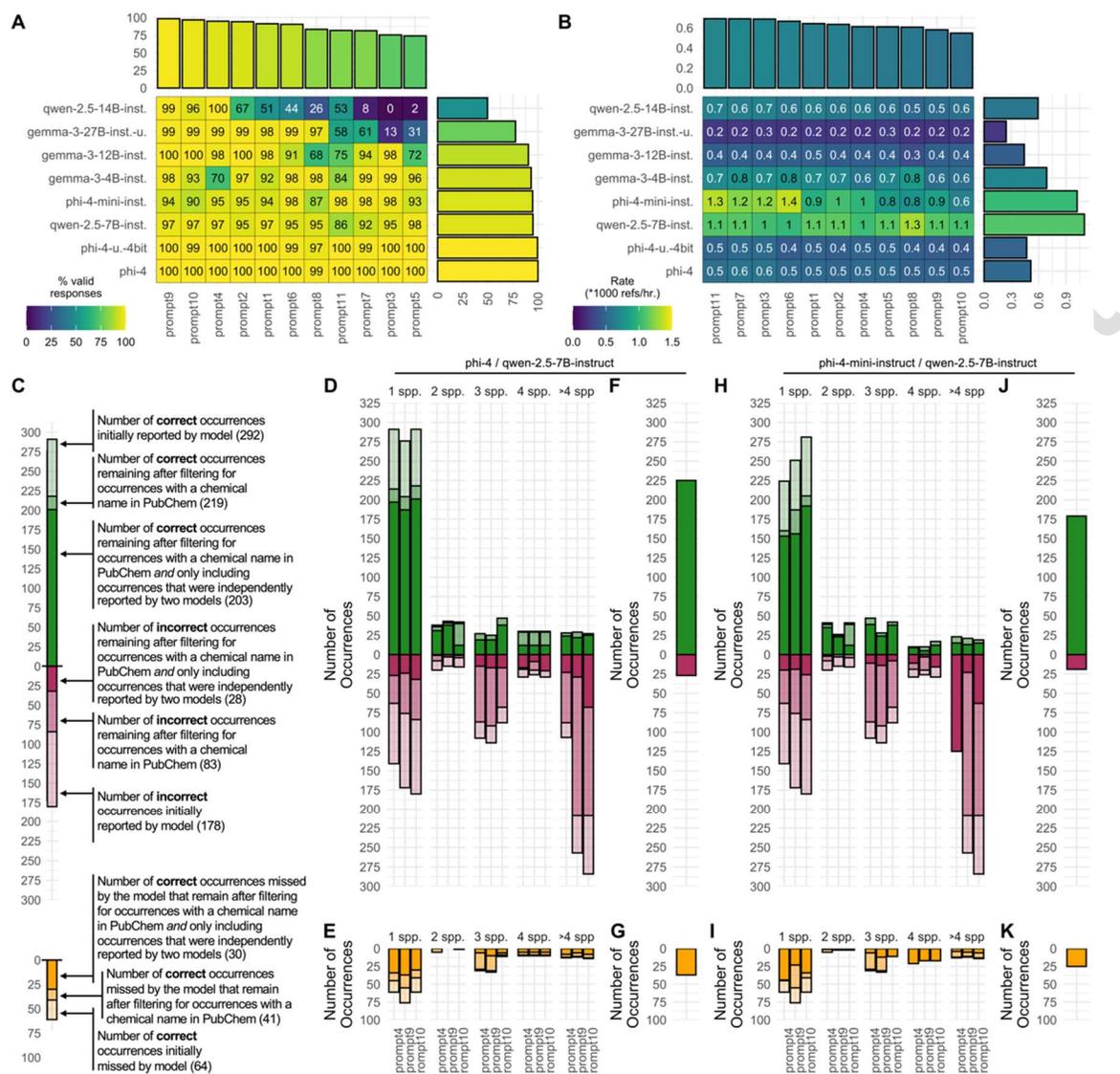
367
368 So far, our results indicated that language models can assess whether an abstract describes experimental support for
369 a particular compound, but no model was entirely accurate in performing this task. Accordingly, we next turned our
370 attention to a detailed examination of the candidate associations that were frequently labeled incorrectly by the
371 language models. Specifically, we reviewed the incorrect answers generated by the phi-4 model. First, we focused
372 on references in which no experimental support for a compound's occurrence was provided, yet the model
373 (erroneously) indicated such support was presented (i.e., false positives). Among these occurrences, two main text
374 structures appeared to "confuse" the model. The first scenario involved abstracts where occurrence data were not
375 presented in separate sentences but instead merged with multiple data types. For example, some passages combined
376 information from authentic standards and plant extracts, or from sediments and plant extracts, or listed multiple
377 compounds from several species in a single statement. The second scenario leading to false positives involved
378 abstracts that failed to provide clear statements about plant/compound occurrences, even to a human reader. As an
379 example, one such abstract stated "beta-sitosterol and alpha-amyrin were isolated from unsaponifiable fractions of
380 mature seeds of solanaceae plants" and mentioned the solanaceous species *Hyoscyamus muticus*, which caused the
381 model to label alpha-amyrin as present in *Hyoscyamus muticus*, even though this link was not explicitly supported
382 by the text. Finally, we examined references where positive associations were mistakenly labeled by the models as
383 negative (i.e., false negatives). We identified three main cases: (i) abstracts that were written in confusing ways,
384 which lead the model to produce an incorrect result, (ii) abstracts that contained an alternative spelling or
385 abbreviation for a compound or species name, and (iii) clearly written abstracts in which the model nevertheless
386 failed to provide the correct answer. These scenarios appeared in roughly equal proportions among phi-4's false
387 negatives. To summarize, the model sometimes makes clear mistakes, but, just as often, the model produces
388 incorrect answers because of inconsistencies or unclear information in the input data. Finally, we also examined the
389 performance of the models when alternative spellings of compound names were present in abstracts. Across the 500
390 candidate associations we manually evaluated there were 28 instances where alternative spellings were used in the
391 abstract (amyrin/amirine, friedelin/friedeline, amyron/amyrenone). Evaluating these candidate associations, the
392 highest performing models were correct ~50% of the time, which is lower than model performance across the entire
393 dataset (~10% overall error rate). Thus, we conclude that these alternative spellings do impact model performance
394 and strategies to deal with such should be included in the design of small language model-based pipelines.

395
396 Several conclusions arise from our work with targeted compound-occurrence data set extraction. First, models with
397 more parameters ("larger" models) appear to perform the task with higher accuracy, though that improvement comes
398 alongside increased computational demands and time requirements. To balance speed and performance, architectures
399 such as phi-4 stand out from those evaluated in this study. Next, the abilities of systems like phi-4 to accurately
400 detect true negatives indicate that they are distinguishing references with textual co-occurrence of plant and
401 compound names from references that present experimental evidence for a plant producing a given compound.
402 Finally, our examination of the underlying reasons for incorrect answers revealed many errors arise from
403 inconsistencies or unclear information in the input data, which suggests that using full-text articles instead of titles
404 and abstracts may improve results beyond the approach described here.

2.2.2 SLM Task B2: Untargeted compound occurrence data extraction

After assessing the extent to which language models can classify compound occurrences in a targeted manner, we next examined these systems' abilities with the same task in an untargeted way. For this process we used models that, as before, accept a system prompt with detailed instructions and a second prompt containing content with which to work. Our general approach was to provide a system prompt directing the model to read the input text (title/abstract) and write all experimentally supported compound occurrences in a Python dictionary format (for example: {"Arabidopsis thaliana": ["arabidiol", "beta-sitosterol"], "Brassica oleracea": ["beta-sitosterol", "alpha-amyrin"]}). Thus, this task is considerably more complicated than the targeted approach. Due to this complexity, we conducted some preliminary tests to determine which of our 12 models might be suitable for this task. We found that the two large models and the two small ones were, respectively, too slow and too inaccurate to be feasible. For this reason, we proceeded with the six medium models as well as the two quantized 4-bit variants described in the previous section. Previous work has shown that the exact phrasing of system prompts can have substantial impacts on the accuracy of language model outputs (Razavi et al. 2025; Sclar et al. 2024), which included the context of phytochemical data processing (Knapp et al. 2024a). This phenomenon is the basis for prompt engineering. This untargeted task was inherently more complicated than the above-described targeted approach, but further complications arose because we wanted a specific output format (the Python dictionary). We investigated a variety of prompts to determine how they might impact results from each model. As in the previous sections, to benchmark the ability of the models to perform this task, we again began by performing this task manually. We read 100 abstracts and wrote out the compound species associations reported in each in the JSON, or Python dictionary, format. This led to the identification of just over 400 compound occurrences across the 100 abstracts (**Supplemental File 5**). Below, we describe the performance of the 8 models and the 11 prompts on this untargeted compound occurrence extraction task with the 100 manually evaluated abstracts.

To begin, we carefully created a detailed system prompt and then employed a commercial large language model to produce 10 additional prompt variants that contained the same instructions but with different phrasings (all prompts included in **Supplemental File 6**). We then used each of the eleven prompts to instruct each of the eight models to write out all experimentally supported occurrences in each of the 100 manually evaluated abstracts. Next, we examined the ability of each model/prompt combination to provide results in a valid Python dictionary (the structure of the response needed to perform this data extraction task) and the speed at which each model/prompt combination could process the 100 abstracts. The percentage of responses from each model in answer to each prompt varied considerably, with some model/prompt combinations producing zero valid dictionaries and others generating 100% valid dictionaries (**Fig. 4A**). Most model-prompt combinations produced >90% correctly structured responses, with some notable exceptions. Interestingly, qwen-2.5-14B-instruct struggled to consistently produce valid dictionary outputs, while its smaller sibling, qwen-2.5-7B-instruct, yielded over 90% valid dictionaries in most cases. This result breaks the trend of larger models being more proficient, as described in the previous section of this report. Phi-4 was the best model tested at this task since it returned 100% valid Python dictionaries, except for one response to prompt 8 (**Fig. 4A**). We also observed variation among the prompts tested, with prompts 9, 10, and 4 eliciting higher proportions of valid responses across all the models than other prompts. We also examined the rate at which each model and prompt pairing could process queries. Rates ranged from about 200 references per hour to almost 1,500 references per hour, with model size as the primary determinant of speed (**Fig. 4B**). Different prompts sometimes caused variability in processing times for the same model, though these shifts were negligible compared to those driven by scale. Overall, the largest model, gemma-3-27B-instruct-unsloth, was the slowest. Meanwhile, phi-4-mini-instruct and qwen-2.5-7B-instruct performed the fastest, at rates around 1,000 articles or references per hour. Altogether, the outcomes suggested that the phi-4 family models, along with qwen-2.5-7B-instruct combined with prompts 9, 10, and 4, were the most accurate for further detailed investigation. The four best performing models for producing valid Python dictionaries included the two fastest frameworks (phi-4-mini-instruct and qwen-2.5-7B-instruct), which showed that larger models do not always perform more proficiently than smaller versions.



453
 454 **Figure 4: Performance of language models on a targeted compound occurrence data extraction task.**
 455 *Heat map showing the percent of outputs that contain valid python dictionaries (encoded with color and written inside*
 456 *each box) from each language model (y-axis) in response to each prompt (x-axis). The marginal (i.e. top and right)*
 457 *plots show the mean percent valid responses across all models for each prompt or across all prompts for each*
 458 *model. **B.** *Heat map showing the rate (in 1000 references per hour) of processing by each language model (y-axis)*
 459 *in response to each prompt (x-axis). The marginal (i.e. top and right) plots show the mean percent valid responses*
 460 *across all models for each prompt or across all prompts for each model.* **C.** *Guide describing how to interpret panels*
 461 *D-K.* **D-K.** *Evaluation of occurrence data reported by language models (D/E/F/G: phi-4 and, in darkest bars, phi-4*
 462 *in agreement with qwen-2.5-7B-instruct; H/I/J/K: phi-4-mini-instruct, and, in darkest bars, phi-4-mini-instruct in*
 463 *agreement with qwen-2.5-7B-instruct). D and H show the number of correct occurrences (true positives, positive y-*
 464 *axis) and incorrect occurrences (false positives, negative y-axis) reported, as indicated in panel C. E and I show the*
 465 *number of correct occurrences (false negatives, negative y-axis) reported, as indicated in panel C. F and J show the*
 466 *number of correct occurrences (true positives, positive y-axis) and incorrect occurrences (false positives, negative y-*
 467 *axis) reported after filtering for occurrences whose compounds are in PubChem and were agreed upon by the two*
 468 *models. G and K show the number of correct occurrences missed by the models after PubChem and agreement*
 469 *filtering (false negatives, negative y-axis). In D-K, bar orientation emphasizes desired model behavior: bars*
 470 *pointing upwards indicate correct model responses (desired behavior), while bars pointing down indicate incorrect*
 471 *model responses or correct answers not reported by the model (undesired behavior).**

472 In the previous section, we identified that the results from prompts 9, 10, and 4, in conjunction with phi-4, phi-4-
473 mini-instruct, and qwen-2.5-7B-instruct warranted further scrutiny. Therefore, we next examined the accuracy of
474 occurrences generated by those models in response to those prompts. In contrast to our quantitative assessment of
475 the models' ability to evaluate targeted compound instances, this broader approach allowed for quantifying only
476 three response types: true positives (correct occurrences reported by a model), false positives (incorrect occurrences
477 reported), and false negatives (correct occurrences missed by the model but found during manual evaluation, **Fig.**
478 **4C**). Note that true negatives are not present in this untargeted analysis since the model is only asked to report
479 existing occurrences, not to classify candidate occurrences. We quantified the number and category of each
480 occurrence identified by each model in response to prompts 4, 9, and 10. We observed that using different system
481 prompts led to only minor variations in the total correct versus incorrect instances flagged by a given model, but,
482 interesting, that correct versus incorrect outputs varied greatly with respect to the number of species described in a
483 given abstract (**Fig. 4D, 4E, 4H, and 4I**). Specifically, references involving more than four species appeared
484 "confusing" to the models, resulting in large numbers of inaccuracies from those sources (**Fig. 4D and 4H**), while
485 abstracts focused on one or two species typically yielded substantially more correct instances compared to incorrect
486 ones (**Fig. 4D and 4H**). Even so, the ratio of correct to incorrect responses typically generated from articles
487 reporting on one or two species was roughly 2:1, an approximately 30% false positive rate.
488

489 To reduce the false positive rate observed during this untargeted compound occurrence extraction task, we
490 introduced two types of filters. For the first filter, we programmatically compared the compound name reported in
491 each occurrence against the PubChem database to check if it appeared among the entries. We removed all reported
492 occurrences describing compounds missing from PubChem, which generally produced a bigger drop in incorrect
493 results than in correct ones. The second filter relied on two language models identifying the same occurrence from a
494 given abstract. Only those occurrences found by both, working independently, were kept, while partial matches
495 (instances flagged by a single model but not recognized by another) were excluded. We tested this two-part filtering
496 approach with two pairs of models: (i) one containing the most advanced model: phi-4 + qwen-2.5-7B-instruct, and
497 (ii) another featuring the two fastest options: phi-4-mini-instruct + qwen-2.5-7B-instruct. In both scenarios, the
498 agreement filter yielded a marked decrease in inaccurate entries in the final dataset and only a small decline in valid
499 ones (**Fig. 4D and H**). Finally, to produce a dataset that reflects the lowest likely false positive rate for these models
500 on the untargeted task at hand, we combined three filtering strategies: we restricted data to abstracts mentioning one
501 or two species, retained only occurrences describing chemicals found in the PubChem database, and kept only those
502 occurrences that were independently detected from the same abstract by two different language models. Using this
503 threefold approach, phi-4 + qwen-2.5-7B-instruct produced about 225 accurate occurrences and 25 inaccurate ones
504 (an 11% error rate and ~55% yield, relative to the 400 occurrences found during manual inspection of the 100
505 abstracts, **Fig. 4F**). Meanwhile, phi-4-mini-instruct + qwen-2.5-7B-instruct yielded 175 valid occurrences and
506 around 20 erroneous findings (also an 11% error rate and ~44% yield/recall, **Fig. 4J**). Thus, pairing two fastest
507 models led to a dataset that was less comprehensive but maintained similar accuracy as a pair that contained a
508 considerably larger and more sophisticated model.
509

3. Conclusions and Future Directions

Here, we evaluated the ability of small language models to perform two major tasks: to numerically score references based on their relevance to a given topic (SLM Task A) and to extract structured data from unstructured inputs in both a targeted (SLM Task B1) and untargeted fashion (SLM Task B2). Our efforts showed that a small language model could rapidly and effectively score references in a set so that a threshold score could be used as a filter to substantially enrich that set for articles of interest (Section 2.1). Limitations arose when handling edge cases, with highly tailored, task-specific prompts emerging as a possible approach to address those shortcomings. When using small language models to classify candidate compound occurrences as true or false, we observed that if an abstract directly reports the detection of a particular compound in a specific plant species, the models nearly always label the candidate occurrence correctly (Section 2.2.1). In this task however, a trade-off did appear between accuracy and model size (parameter count) and compute requirements. Among the mistakes noted (false positives as low as 5% and false negatives as low as 2% for certain models), these misclassifications were as often tied to convoluted and unclear writing in the input abstract as they were to outright model errors. For extracting compound occurrence information from unstructured text in an untargeted manner, we found small language models to be effective, though choosing a suitable model and pipeline strategy proved more challenging than earlier tasks (Section 2.2.2). We discovered that prompt engineering, selecting a model, and filtering reported detections by cross-referencing chemical databases, along with requiring two small language models to independently agree on an occurrence, yielded the best reporting statistics (~10% false positives and ~50% yield). Of note, this relatively low yield arises because many correct associations are filtered out, essentially sacrificed, to lower the false positive rate. Regarding all tasks considered, more advanced prompting techniques (e.g., chain-of-thought prompting (Wei et al. 2022) or model distillation (Hinton et al. 2015; Sanh et al. 2020) could reduce error rates further and improve yield/recall. In addition, future model releases, including small reasoning models, may also address these limitations. Finally, we will note that many abstracts we worked with here presented problems for humans and language models alike by failing to contain clear and concise information. We read hundreds of abstracts for the present project. Fully understanding many abstracts in a timely fashion was extremely difficult due to long, convoluted sentences, the presentation of connected data types (e.g., plants and compounds) in multiple sentences spread throughout a long abstract, the use of compound numbers or abbreviations instead of compound names, poor grammar, and so forth. In a variety of cases, we were surprised that the language models performed reasonably well while humans needed considerable time to understand the same abstracts.

Overall, though the approaches here represent a considerable advance over manual curation (at least, with respect to the creation of large databases, where speed is a prime consideration), a substantial amount of plant chemical occurrence data will still not be retrieved from the literature using the techniques presented here. One important step forward will be the development of pipelines that can handle articles reporting occurrence data from dozens of species, including in tabular format. In addition, further attempts towards occurrence databases, and in fact scientific endeavors in general, need literature databases that include the full text files along with reference citations and abstracts. The separation of the full text from the citations seems to be a systematic and legal barrier that needs to be overcome. The expanded posting of pre-prints is suggested as a potential, albeit partial, solution to this issue. In addition to the tasks we quantitatively evaluated here, we also experimented with several versions of Microsoft's Phi-4 model to conduct multiple activities related to reference citations (e.g., species name extraction, compound number or plant number extraction, etc.) and found that the models could perform a range of additional functions, suggesting versatility and application in order domains. In our case, these functions have allowed us to identify publications that most likely contain extensive tabular data in the full text, flagging them for analysis by a pipeline suitable for such reports. Finally, in our efforts, we found that filtering capabilities such as those provided by SciFinder® and EndNote™ showed usefulness in a somewhat orthogonal way to the value of the small language model scores. For example, in our case, we were able to eliminate many articles of low relevance to our case studies using EndNote™ keyword filters. As these commercial software tools and other related programs are outfitted with language model ("artificial intelligence") capabilities, it will be important to evaluate and incorporate those features into discipline-specific workflows. We strongly encourage the scientific community to look for new versions of their favorite research tools that incorporate language model features, and to experiment and empirically test and report on such functionality in field-specific tasks as they emerge.

563 4. Methods:

564
565 Literature searches were conducted with CAS SciFinder®. SciFinder® searches were conducted by entering the
566 compound CAS Registry® number from the SUBSTANCE menu and then working with all the references that were
567 assigned to this Registry number. SciFinder® references were downloaded as “tagged” text files. The “tagged” text
568 file selection provides numerous fields including the CAS Registry numbers for all compounds discussed in a given
569 article. Multiple tagged files were downloaded for each compound (according to year ranges) since the SciFinder®
570 software limits an individual tagged export file to 100 citations. SciFinder limits the number of citations that can be
571 exported in one file to 100. Thus, for a compound such as alpha-amyrin with 4,344 SciFinder references, the
572 downloading of all references was not possible. If the number of filtered references was greater than 400, the word
573 “plant” was entered into the “search within results.” Thus, only English-language journal references that
574 corresponded to the “search within results” term “plant” were downloaded (1,744 references, in the example of
575 alpha-amyrin). PubMed® searches for the six triterpenoids were also conducted based on their major common
576 names (not all synonyms were used). These PubMed® searches were conducted with the compound names shown at
577 the top of Table 1 since PubMed® does not generally recognize CAS Registry® numbers. PubMed® files were
578 downloaded as PubMed (NLM) files. Of note is that PubMed® provides automated access to its search and abstract
579 download services through a REST API and various language-specific packages like R and Trez. These tools could
580 be leveraged in the future to further streamline literature analysis projects and automate data extraction and
581 tabulation.

582
583 EndNote™ Version 21.5 (<https://endnote.com/>) was used to import and combine the sets of “tagged” SciFinder®
584 export text files for each compound into an individual EndNote™ compound folders (with the “discard duplicate”
585 feature turned on). Furthermore, EndNote™ “Smart Groups” were set up for each of the six triterpenoids, which
586 included the CAS Registry® number and multiple names for each compound (i.e., synonyms). The references in
587 each of the six Smart Groups were then added to the corresponding original six triterpenoid folders (with automatic
588 elimination of duplicates). As noted above, some plants contained more than one of the six triterpenoids. These
589 EndNote™ operations ensured that any references that might have been missed in a given SciFinder® compound
590 search, but included in another compound search, would end up in the appropriate folders (i.e., one reference might
591 be in more than one compound folder). In EndNote®, the user can select scores of references and then right-click on
592 “Find full text.” EndNote will then automatically download the PDF files for each reference that cites a journal for
593 which the user's institution has a subscription or an open-source journal. However, in some cases, software blocks
594 (e.g., the “Are you a human filter?”) prevent the downloading of some files. In our case at our institution, EndNote
595 downloads approximately 40-50% of the PDF files for the selected references.

596
597 All manual evaluation of reference relevance (“reports an occurrence”, “maybe reports an occurrence”, “does not
598 report an occurrence”), manual evolution of candidate occurrences (targeted) and manual extraction of associations
599 (untargeted) was performed by opening the list of references in Microsoft Excel and entering the manual annotations
600 into a new column. References were labeled as “maybe reports an occurrence” if they mentioned specific plant
601 species and the isolation of multiple compounds from the species but did not mention the specific compounds’
602 names in the abstract. While the likelihood of a plant/compound association appearing in the full article was high,
603 we nevertheless conservatively chose to label these types of citations as “maybe reports an occurrence” until the full
604 text article file could be evaluated. An “maybe” example is: "Medicinal attributes of *Solanum capsicoides* All.: an
605 antioxidant perspective. *Int. J. Pharm. Sci. Res.* 12(5): 2810-2817. The study evaluates the medicinal efficacy of
606 *Solanum capsicoides* fruits as an antioxidant. Fruit extracts were prepared using acetone, ethanol, HCl, and water
607 [...] A neg. correlation was observed between the pigments, anthocyanins, and carotenoids, with DPPH and
608 CUPRAC activity. [...] From this study, it can be considered that the phenolics present in the fruits contribute to the
609 characteristic antioxidant property."

610
611 The Facebook BART-Large-MNLI zero-shot classification model (<https://huggingface.co/facebook/bart-large-mnli>)
612 was applied to the individual sets of compound reference citations in the EndNote™ database. The model was run
613 on a single NVIDIA GV100GL [Quadro GV100] GPU. First, the set of references in the curated EndNote™ folder
614 for a given compound was selected and exported from this folder to a text file (with the “annotated” style selected).
615 This text file was then imported into an Excel file (e.g., with the legacy “get text from file” Excel wizard. The
616 resulting Excel sheet was then modified so that each reference citation (author/year/journal/abstract) was contained
617 in one cell and all cells resided in one column. This Excel sheet, which contained all the reference citations for a
618 given compound, was then saved as a CSV UTF-8 (Comma delimited) file. This CSV file was used via JupyterLab

619 (<https://jupyter.org/>, operating in a WINDOWS 11 environment) and a custom Python program (full code in
620 Supplemental File 8). System prompt-accepting chat language models were downloaded from HuggingFace.co and
621 run on a single NVIDIA GV100GL [Quadro GV100] GPU using custom code (full code provided in Supplemental
622 File 8). Calculation of precision, recall, and F1 scores as well as plotting were performed in R. Additional system
623 prompts for the prompt engineering reported in Section 2.2.2 were generated by OpenAI's o4-mini-high language
624 model using the ChatGPT browser interface.

625

626 **5. Supplemental Materials:**

627

628 Supplemental File 1: 1,558 references manually scored for relevance to compound occurrence.

629 Supplemental File 2: Details of classifier phrases.

630 Supplemental File 3: Details of maybe references.

631 Supplemental File 4: 500 manually evaluated candidate occurrences.

632 Supplemental File 5: 100 abstracts from which untargeted occurrence data was manually extracted.

633 Supplemental File 6: Prompts that were used in small language model untargeted occurrence data extraction.

634 Supplemental File 7: Schematic of reference acquisition process.

635 Supplemental File 8: Code used in this work.

636

637 **6. Acknowledgements**

638

639 The authors wish to acknowledge the support the University of Minnesota Duluth Chemistry and Biochemistry
640 Department. During the writing of this manuscript, the large language model o4-mini-high from OpenAI was
641 utilized for copy editing to suggest alternative sentence structures and word choices that, in our view, enhanced
642 grammar and readability. Finally, we collectively acknowledge that the University of Minnesota Duluth is located on
643 the traditional, ancestral, and contemporary lands of Indigenous people. The University resides on land that was
644 cared for and called home by the Ojibwe people, before them the Dakota and Northern Cheyenne people, and other
645 Native peoples from time immemorial. Ceded by the Ojibwe in an 1854 treaty, this land holds great historical,
646 spiritual, and personal significance for its original stewards, the Native nations, and peoples of this region. We
647 recognize and continually support and advocate for the sovereignty of the Native nations in this territory and
648 beyond. By offering this land acknowledgment, we affirm tribal sovereignty and will work to hold the University of
649 Minnesota Duluth accountable to American Indian peoples and nations.

650

651 **7. Data Availability Statement**

652

653 All data and code used in this study are available for free as a Supplement to this document.

654

655 **8. Financial Support**

656

657 This work was funded by startup funds granted to Lucas Busta from the University of Minnesota Duluth Swenson
658 College of Science and Engineering.

659

660 **9. Author Contributions**

661

662 LB and ARO conceived and designed the study, gathered data, and wrote the manuscript.

663

664 **10. Conflicts of Interest:**

665

666 None

667

668

669 **11. References:**

670

671 Abdin M, Eldan R, Javaheripi M, Li Y, Price E, Shah S, Yu D, Aneja J, Gunasekar S,
672 Kauffmann P, Liu W, Rosa Gd, Wang X, Zhang C, Behl H, Harrison M, Lee JR,
673 Mendes CCT, Saarikivi O, Ward R, Zhang Y, Bubeck Se, Hewett RJ, Lee YT,
674 Nguyen A, Salim A, Wu Y (2024) Phi-4 Technical Repor. arXiv:241208905v1
675 [csCL] 12 Dec 2024

676 Agathokleous E, Rillig MC, Penuelas J, Yu Z (2024) One hundred important questions
677 facing plant science derived using a large language model. *Trends in Plant*
678 *Science* 29 (2):210-218. doi:10.1016/j.tplants.2023.06.008

679 Brown TB, Mann B, Ryder N, Subbiah M, Kaplan J, Dhariwal P, Neelakantan A, Shyam
680 P, Sastry G, Askell A, Agarwal S, Herbert-Voss A, Krueger G, Henighan T, Child
681 R, Ramesh A, Ziegler DM, Wu J, Winter C, Hesse C, Chen M, Sigler E, Litwin M,
682 Gray S, Chess B, Clark J, Berner C, McCandlish S, Radford A, Sutskever I,
683 Amodei D (2020) Language Models are Few-Shot Learners. arXiv,

684 Busta L, Hall D, Johnson B, Schaut M, Hanson CM, Gupta A, Gundrum M, Wang Y, A.
685 Maeda H (2024a) Mapping of specialized metabolite terms onto a plant
686 phylogeny using text mining and large language models. *The Plant Journal*
687 (2024-7-8). doi:10.1111/tpj.16906

688 Busta L, Hall D, Johnson B, Schaut M, Hanson CM, Gupta A, Gundrum M, Wang Y, H
689 AM (2024b) Mapping of specialized metabolite terms onto a plant phylogeny
690 using text mining and large language models. *Plant J* 120 (1):406-419.
691 doi:10.1111/tpj.16906

692 Chen Y, de Bruyn Kops C, Kirchmair J (2017) Data Resources for the Computer-Guided
693 Discovery of Bioactive Natural Products. *Journal of Chemical Information and*
694 *Modeling* 57 (9):2099-2111. doi:10.1021/acs.jcim.7b00341

695 Dalal A, Ranjan S, Bopaiah Y, Chembachere D, Steiger N, Burns C, Daswani V (2024)
696 Text summarization for pharmaceutical sciences using hierarchical clustering with
697 a weighted evaluation methodology. *Scientific Reports* 14 (1):20149.
698 doi:10.1038/s41598-024-70618-w

699 Gallo K, Kemmler E, Goede A, Becker F, Dunkel M, Preissner R, Banerjee P (2023)
700 SuperNatural 3.0-a database of natural products and natural product-based
701 derivatives. *Nucleic Acids Res* 51 (D1):D654-D659. doi:10.1093/nar/gkac1008

702 Gemma T, Kamath A, Ferret J, Pathak S, Vieillard N, Merhej R, Perrin S, Matejovicova
703 T, Ramé A, Rivière M, Rouillard L, Mesnard T, Cideron G, Grill J-b, Ramos S,
704 Yvinec E, Casbon M, Pot E, Penchev I, Liu G, Visin F, Kenealy K, Beyer L, Zhai
705 X, Tsitsulin A, Busa-Fekete R, Feng A, Sachdeva N, Coleman B, Gao Y, Mustafa
706 B, Barr I, Parisotto E, Tian D, Eyal M, Cherry C, Peter J-T, Sinopalnikov D,
707 Bhupatiraju S, Agarwal R, Kazemi M, Malkin D, Kumar R, Vilar D, Brusilovsky I,
708 Luo J, Steiner A, Friesen A, Sharma A, Sharma A, Gilady AM, Goedeckemeyer A,
709 Saade A, Feng A, Kolesnikov A, Bendebury A, Abdagic A, Vadi A, György A,
710 Susano Pinto A, Das A, Bapna A, Miech A, Yang A, Paterson A, Shenoy A,
711 Chakrabarti A, Piot B, Wu B, Shahriari B, Petrini B, Chen C, Le Lan C,
712 Choquette-Choo CA, Carey CJ, Brick C, Deutsch D, Eisenbud D, Cattle D,
713 Cheng D, Paparas D, Shivakumar Sreepathihalli D, Reid D, Tran D, Zelle D,
714 Noland E, Huizenga E, Kharitonov E, Liu F, Amirkhanyan G, Cameron G,
715 Hashemi H, Klimczak-Plucińska H, Singh H, Mehta H, Lehri HT, Hazimeh H,

716 Ballantyne I, Szpektor I, Nardini I, Pouget-Abadie J, Chan J, Stanton J, Wieting J,
717 Lai J, Orbay J, Fernandez J, Newlan J, Ji J-y, Singh J, Black K, Yu K, Hui K,
718 Vodrahalli K, Greff K, Qiu L, Valentine M, Coelho M, Ritter M, Hoffman M, Watson
719 M, Chaturvedi M, Moynihan M, Ma M, Babar N, Noy N, Byrd N, Roy N, Momchev
720 N, Chauhan N, Sachdeva N, Bunyan O, Botarda P, Caron P, Rubenstein PK,
721 Culliton P, Schmid P, Sessa PG, Xu P, Stanczyk P, Tafti P, Shivanna R, Wu R,
722 Pan R, Rokni R, Willoughby R, Vallu R, Mullins R, Jerome S, Smoot S, Girgin S,
723 Iqbal S, Reddy S, Sheth S, Pöder S, Bhatnagar S, Raghuram Panyam S, Eiger
724 S, Zhang S, Liu T, Yacovone T, Liechty T, Kalra U, Evcı U, Misra V, Roseberry V,
725 Feinberg V, Kolesnikov V, Han W, Kwon W, Chen X, Chow Y, Zhu Y, Wei Z,
726 Egyed Z, Cotruta V, Giang M, Kirk P, Rao A, Black K, Babar N, Lo J, Moreira E,
727 Martins LG, Sanseviero O, Gonzalez L, Gleicher Z, Warkentin T, Mirrokni V,
728 Senter E, Collins E, Barral J, Ghahramani Z, Hadsell R, Matias Y, Sculley D,
729 Petrov S, Fiedel N, Shazeer N, Vinyals O, Dean J, Hassabis D, Kavukcuoglu K,
730 Farabet C, Buchatskaya E, Alayrac J-B, Anil R, Lepikhin D, Borgeaud S, Bachem
731 O, Joulin A, Andreev A, Hardin C, Dadashi R, Hussenot L (2025) Gemma 3
732 Technical Report. arXiv,
733 Guo Z, Wnag Y, Wang P, Yu P (2023) Improving Small Language Models on PubMedQA
734 via Generative Data Augmentation. arXiv:230507804v4 [csCL] 1 Aug 2023
735 Hinton G, Vinyals O, Dean J (2015) Distilling the Knowledge in a Neural Network. vol
736 1503.02531.
737 Jin Q, Leaman R, Lu Z (2024) PubMed and beyond: biomedical literature search in the
738 age of artificial intelligence. EBioMedicine 100:104988.
739 doi:10.1016/j.ebiom.2024.104988
740 Kaplan J, McCandlish S, Henighan T, Brown TB, Chess B, Child R, Gray S, Radford A,
741 Wu J, Amodei D (2020) Scaling Laws for Neural Language Models. arXiv,
742 Knapp R, Johnson B, Busta L (2024a) Advancing Plant Metabolic Research By Using
743 Large Language Models To Expand Databases And Extract Labelled Data.
744 (2024-11-6). doi:10.1101/2024.11.05.622126
745 Knapp R, Johnson B, Busta L (2024b) Advancing plant metabolic research by using
746 large language models to expand databases and extract labelled data. BioRxiv.
747 doi:10.1101/2024.11.05.622126
748 Lam HYI, Ong XE, Mutwil M (2024) Large language models in plant biology. Trends in
749 Plant Science 29 (10):1145-1155. doi:10.1016/j.tplants.2024.04.013
750 Lepagnol P, Gerald T, Ghannay S, Seran C, Rosset S (2024) Small Language Models
751 are Good Too: An Empirical Study of Zero-Shot Classification.
752 arXiv:240411122v1 [csAI] 17 Apr 2024.
753 doi:<https://arxiv.org/abs/2404.11122?form=MG0AV3>
754 Lewis M, Liu Y, Goyal N, Ghazvininejad M, Mohamed A, Levy O, Stoyanov S,
755 Zettlemoyer L (2019) BART: Denoising Sequence-to-Sequence Pre-training for
756 Natural Language Generation, Translation, and Comprehension.
757 arXiv:191013461v1 [csCL] 29 Oct 2019
758 Maronikolakis A, Dufter P, Schütze H (2021) BERT Cannot Align Characters.
759 Nguyen-Vo T-H, Nguyen L, Do N, Nguyen T-N, Trinh K, Cao H, Le L (2020) Plant
760 Metabolite Databases: From Herbal Medicines to Modern Drug Discovery.

761 Journal of Chemical Information and Modeling 60 (3):1101-1110.
762 doi:10.1021/acs.jcim.9b00826
763 Qwen: An Y, Baosong Y, Beichen Z, Binyuan H, Bo Z, Bowen Y, Chengyuan L,
764 Dayiheng L, Fei H, Haoran W, Huan L, Jian Y, Jianhong T, Jianwei Z, Jianxin Y,
765 Jiayi Y, Jingren Z, Junyang L, Kai D, Keming L, Keqin B, Kexin Y, Le Y, Mei L,
766 Mingfeng X, Pei Z, Qin Z, Rui M, Runji L, Tianhao L, Tianyi T, Tingyu X,
767 Xingzhang R, Xuancheng R, Yang F, Yang S, Yichang Z, Yu W, Yuqiong L, Zeyu
768 C, Zhenru Z, Zihan Q (2025) Qwen2.5 Technical Report
769 :contentReference[oaicite:1]{index=1}. arXiv. doi:10.48550/arXiv.2412.15115
770 :contentReference[oaicite:3]{index=3}
771 Razavi A, Soltangheis M, Arabzadeh N, Salamat S, Zihayat M, Bagheri E (2025)
772 Benchmarking Prompt Sensitivity in Large Language Models. arXiv,
773 Riordan BB, Abraham Albert; He, Zhen; Nibali, Aiden; Anderson-Luxford, Dan;
774 Kuntsche, Emmanuel (2024) How to apply zero-shot learning to text data in
775 substance use research: An overview and tutorial with media data. . Addiction
776 (Abingdon, England) 119 (5):951-959. doi:10.1111/add.16427
777 Rutz A, Sorokina M, Galgonek J, Mietchen D, Willighagen E, Gaudry A, Graham JG,
778 Stephan R, Page R, Vondrášek J, Steinbeck C, Pauli GF, Wolfender J-L, Bisson
779 J, Allard P-M (2022) The LOTUS initiative for open knowledge management in
780 natural products research. eLife 11. doi:10.7554/eLife.70780
781 Sanh V, Debut L, Chaumond J, Wolf T (2020) DistilBERT, a distilled version of BERT:
782 smaller, faster, cheaper and lighter. vol 1910.01108.
783 Sarumi OA, Heider D (2024) Large language models and their applications in
784 bioinformatics. Computational and Structural Biotechnology Journal 23:3498-
785 3505. doi:10.1016/j.csbj.2024.09.031
786 Sclar M, Choi Y, Tsvetkov Y, Suhr A (2024) Quantifying Language Models' Sensitivity to
787 Spurious Features in Prompt Design or: How I learned to start worrying about
788 prompt formatting. arXiv,
789 Shiu S-H, Lehti-Shiu MD (2024) Assessing the evolution of research topics in a
790 biological field using plant science as an example. PLoS Biology 22
791 (5):e3002612. doi:10.1371/journal.pbio.3002612
792 Simon E, Swanson K, Zou J (2024) Language models for biological research: a primer.
793 Nat Methods 21 (8):1422-1429. doi:10.1038/s41592-024-02354-y
794 Sorokina M, Steinbeck C (2020) Review on natural products databases: where to find
795 data in 2020. J Cheminform 12 (1):20. doi:10.1186/s13321-020-00424-9
796 Stevenson CE, Pafford A, van der Maas HLJ, Mitchell M (2025) Can Large Language
797 Models generalize analogy solving like people can?
798 Tay DWP, Yeo NZX, Adaikkappan K, Lim YH, Ang SJ (2023) 67 million natural product-
799 like compound database generated via molecular language processing. Sci Data
800 10 (1):296. doi:10.1038/s41597-023-02207-x
801 Wei J, Wang X, Schuurmans D, Bosma M, Ichter B, Xia F, Chi E, Le QV, Zhou D (2022)
802 Chain-of-Thought Prompting Elicits Reasoning in Large Language Models. Paper
803 presented at the Advances in Neural Information Processing Systems,
804 Xie T, Song S, Li S, Ouyang L, Xia L, Huang J (2015) Review of natural product
805 databases. Cell Prolif 48 (4):398-404. doi:10.1111/cpr.12190

- 806 Yang B, Mao J, Gao B, Lu X (2019) Computer-Assisted Drug Virtual Screening Based
807 on the Natural Product Databases. *Current Pharmaceutical Biotechnology* 20
808 (4):293-301. doi:10.2174/1389201020666190328115411
- 809 Yin W, Hay J, Rother D (2019) Benchmarking Zero-shot Text Classification: Datasets,
810 Evaluation and Entailment Approach. *Proceedings of the 2019 Conference on*
811 *Empirical Methods in Natural Language Processing and the 9th International*
812 *Joint Conference on Natural Language Processing*, pages 3914–3923, Hong
813 Kong, China, November 3–7, 2019
- 814 Zeng T, Li J, Wu R (2024) Natural product databases for drug discovery: Features and
815 applications. *Pharmaceutical Science Advances* 2.
816 doi:10.1016/j.pscia.2024.100050
- 817 Zhu XL, Jian; Liu, Yong; Ma, Can; Wang, Weiping (2024) Distilling mathematical
818 reasoning capabilities into Small Language Models. . *Neural networks : the*
819 *official journal of the International Neural Network Society* 179:106594.
820 doi:10.1016/j.neunet.2024.106594
821

Accepted Manuscript