

1 Unveiling the Synthesis Patterns of Nanomaterials: A Text Mining and Meta-Analysis Approach 2 with ZIF-8 as a Case Study

3 Joseph R.H. Manning^{*a} and Lev Sarkisov^{*a}

4 (a) Department of Chemical Engineering, University of Manchester, UK, M13 9PL

5 Joseph.Manning@manchester.ac.uk

6 Lev.Sarkisov@manchester.ac.uk

7 8 **Abstract**

9 With the continuously growing number of scientific articles on synthesis of nanomaterials, it
10 becomes impossible for researchers to grasp and comprehend the landscape of synthetic protocols
11 available for a particular material. The aim of this study is to explore the feasibility of extracting the
12 collective knowledge on synthesis of a particular material accumulated over the years from the
13 published corpus of articles and organizing it in a systematic manner. Accordingly, we developed
14 methods to perform detailed text mining on a single nanomaterial target for the purposes of
15 methodology optimisation. Taking the common material ZIF-8 as a case study, we analysed 1600
16 synthesis protocols to identify trends in parameters, such as reagents, concentrations, and reaction
17 time/temperature. We used this information to find the distribution of synthesis parameters and
18 their relationships to one another, identifying the limits of common reaction parameters and
19 revealing subtle details, such as insolubility of metal acetate reagents in alcoholic solvents, or the
20 occurrence of amorphous oxides at low stoichiometric ratios. We then clustered similar synthesis
21 protocols together, using their relative popularity to identify promising regions of the synthesis
22 phase space for optimisation, reducing the need for brute force synthesis optimisation. The
23 techniques developed here are a general tool accelerating the synthesis development of a wide
24 range of nanomaterials by aggregating existing research trends, averting the need for laborious
25 manual comparison of existing synthesis protocols or repetition of previously-developed techniques.

26 **Introduction**

27 The number of chemical syntheses reported is large and growing exponentially.¹ While naturally
28 indicative of greater scientific progress, this leads to two significant challenges. Firstly, researchers
29 are confronted with the growing difficulty of maintaining a comprehensive overview and
30 understanding of the diverse landscape of synthetic routes and conditions accessible for a particular
31 group of compounds. Secondly, although the repository of published synthesis data contains an
32 immense wealth of information, its potential for systematic development of new synthesis protocols
33 remains largely untapped and underutilized. In response to this, various informatics approaches
34 have been adopted to standardise the data produced during chemical research. For example, the
35 creation of chemical synthesis ontologies²⁻⁴ and automated reactionware^{5,6} has enabled new
36 procedures to be directly compared against previously-published data or shared openly through
37 chemical “programming languages”.^{7,8} However, the nature of reporting synthesis protocols – as
38 unformatted prose in a written report – has remained largely unchanged.

39 As a result, most new publications and the entire body of prior chemical synthesis reports remains
40 unlabelled, with the potential for far broader data mining and informatics research if these reports
41 could be standardised. Accordingly, with the advent of text mining methods and natural language
42 processing (NLP),⁹ software has been developed to interpret chemical details from the plain text
43 within chemistry publications^{10,11} including compound structure,¹² reaction stoichiometry,¹³ and
44 performance.¹⁴ Using these tools, large databases of organic^{14,15} and inorganic¹⁶⁻¹⁹ chemicals and
45 reactions have been developed and used for novel materials discovery. For example, Cole and co-

46 workers created a database of organic dyes to identify ideal mixtures for broad-spectrum light
47 absorption in dye-sensitized solar cells, regardless of the intension of the original studies.¹⁵ Similar
48 strategies have been used by Olivetti and co-workers to analyse how synthesis gel composition and
49 organic structure directing agent (OSDA) can dictate crystal polymorphs for a range of zeolite
50 syntheses.¹⁶

51 One weakness of these text mining approaches is their reliance on unambiguous identification of the
52 chemical entities in question, using named-entity recognition (NER)^{9,20} and the programmatic
53 naming conventions defined by IUPAC²¹ to succeed. In the absence of such well-accepted naming
54 schemes – as is the case for a variety of emerging nanomaterial families like porous silicas, polymers
55 of intrinsic microporosity, and covalent organic framework materials – large scale data mining
56 becomes far less practical. An excellent example of this is metal-organic framework (MOF) materials
57 - infinite condensation polymers of various organic ligands and metal ions or clusters. There are
58 millions of possible MOFs,^{22–25} and hundreds of thousands of frameworks already synthesized,^{26–29}
59 necessitating data-driven approaches to accelerate progress in the field. However, unambiguous
60 naming conventions for MOFs have yet to be fully adopted,³⁰ frustrating text-mining of the primary
61 publications themselves. Instead, informatics methods have largely been driven by the creation of a
62 subset of the Cambridge Structural Database (CSD)³¹ focused on MOF materials,²⁸ as these resources
63 allow researchers to analyse the full range of experimentally known MOF structures, identifying the
64 best experimentally-realised materials for future research and development.

65 To accelerate development of experimental procedures to make MOFs, however, data-mining
66 approaches must look beyond structure into the synthesis protocols leading to different
67 frameworks. By understanding the relationships between protocol and eventual material, new
68 synthesis methods can be digitally generated, obviating the need for arduous trial-and-error or
69 intuition-based approaches.⁶ To this end, large-scale post-hoc analyses of experimental MOF
70 synthesis protocols have recently been developed.^{32,33} These studies apply NLP to the underlying
71 publications in the CSD MOF subset to interpret their synthesis protocols, identifying such details as
72 solvents used, specific reagents, solvents, and reaction parameters. As a result, broad descriptive
73 statistics about the synthesis strategies to produce MOFs have been developed,³³ and even
74 predictive models to suggest synthesis parameters for novel MOF materials when given a
75 hypothetical structure.³²

76 While these approaches give an excellent overview of the field of MOFs in general, they are
77 vulnerable to bias in the papers submitting to the CSD. As the database focuses on chemical structure
78 rather than synthesis protocols, only 1-2 synthesis examples of each framework are included.
79 Further, the synthesis protocols are generally submitted from initial studies reporting the discovery
80 of a material, rather than exploring the full range of potential approaches to a single target, meaning
81 that only a very vague understanding of any individual MOF can be generated with this approach.
82 For example, while candidate solvents and reaction parameters can be suggested, other salient
83 parameters such as reagent ratios, product isolation methods, and alternative synthesis strategies
84 (e.g. hydrothermal or mechanochemical versus solvent crystallisation) cannot. Deeper insight into
85 individual MOFs and the peculiarities of their synthesis protocols can be gained through targeted
86 meta-analysis of studies focusing on that particular material,³⁴ enabling regression of product
87 properties like defect density against synthesis details. However, challenges of manually comparing
88 synthesis protocols against one another severely limit the scale of such meta-analyses, preventing
89 their widespread use.

90 To address these issues, in this article we pose the following questions: can we leverage previously-
91 developed chemistry text mining tools to analyse the synthesis protocols for a single target

92 nanomaterial? If so, can we develop methods to process the extracted information on a uniform
93 basis, enabling like-for-like comparison regardless of original format? Finally, can we harness this
94 information to accelerate synthesis refinement of the material e.g. by generating proposed synthesis
95 conditions correlated to high material quality and yield?

96 As a case study, we consider ZIF-8, a commonly synthesized MOF material which has been
97 extensively studied within the literature. ZIF-8 is constructed from a combination of zinc ions and 2-
98 methylimidazole in the sodalite topology, held together with metal-amine bonds rather than the
99 more common metal-carboxylate bonds, thus rendering the material both hydrophobic and water-
100 stable.^{35,36} Accordingly, ZIF-8 has garnered significant interest in the literature for applications
101 including gas storage and separation, adsorptive refrigeration,³⁷ biomolecule encapsulation,³⁸
102 catalysis,³⁹ and sensing.⁴⁰ Further, ZIF-8 can be synthesized from a number of strategies – for
103 example using protic or aprotic solvents,⁴¹ a range of temperatures,⁴² reagent concentrations,⁴³
104 modulators and crystal growth modifiers,⁴⁴ and acid/base conditions.⁴⁵ In sum, over 7500 papers
105 have been published regarding ZIF-8 to date. Given the breadth of synthesis protocols established
106 for ZIF-8, it is practically impossible to manually compare all possible synthesis methodologies to one
107 another. Applying text mining methods to automatically and quantitatively analyse ZIF-8 synthesis
108 protocols would enable larger-scale analysis and the identification of promising synthesis strategies.

109 In this study we developed methods to extract and aggregate synthesis protocols in a uniform
110 format. We studied 1600 synthesis protocols of ZIF-8 and related materials from 3197 original
111 articles, performing an automated meta-analysis of the synthesis methods contained. We analysed
112 the chemical identities used alongside quantities and reaction conditions to provide a systematic
113 design space for ZIF-8, identifying key trends in the approaches used. Finally, we group similar
114 synthesis protocols together with unsupervised clustering techniques, identifying hidden patterns in
115 the data.

116 Methods development

117 The workflow of extracting and analysing synthesis protocols was split into four overarching steps:
118 text collection, where a corpus of research papers is identified and downloaded; paragraph
119 identification, where raw synthesis protocols are identified within the prose; grammar parsing,
120 where the natural language is converted into hierarchical data for later interpretation; and synthesis
121 protocol extraction, where the extracted data is standardised to produce a structured “recipe” for
122 each synthesis protocol. Key steps in the workflow are depicted in Figure 1. The first three steps
123 have been widely described elsewhere, and only a brief description is provided in this section (with
124 associated code provided by the authors on GitHub at [https://github.com/SarkisovTeam/SynOracle-
125 preprocessing](https://github.com/SarkisovTeam/SynOracle-preprocessing)). The final stage of the workflow was developed in this study using python 3.9,⁴⁶ and is
126 made freely available by the authors on GitHub at
127 <https://github.com/SarkisovTeam/SyntheticOracle>.

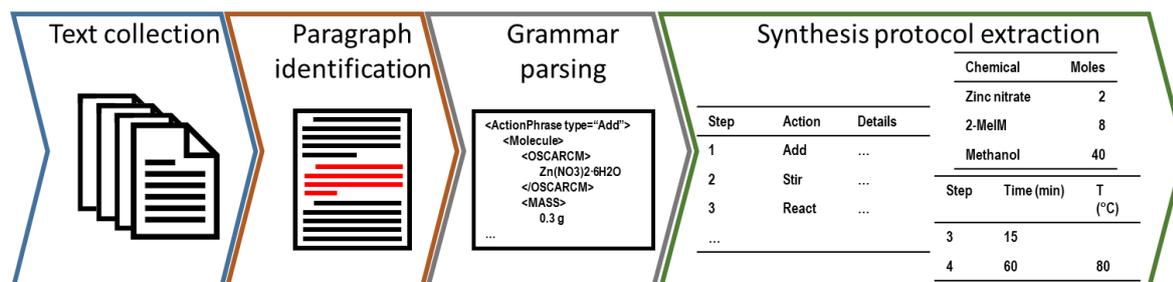


Figure 1 – Scheme of the data processing pipeline used in this study.

130 *Text collection, paragraph identification, and grammar parsing*

131 To produce a corpus of ZIF-8 synthesis protocols, we initially followed established methods to
132 download collections of papers and identify synthesis protocols within them.^{32,33} Synthesis papers
133 were identified by searching the SCOPUS database using Elsevier's elsapy software
134 (<https://github.com/ElsevierDev/elsapy>). Papers were identified using the search term "ZIF OR
135 zeolitic imidazol* AND synthesis," returning 4198 results. These were then categorised by publisher,
136 from which the three largest groups were targeted for downloading (ACS, RSC, and Elsevier),
137 reducing the total corpus to 3179 papers. XML or HTML versions of each paper were then
138 downloaded according to their publisher's specifications – using elsapy in the case of Elsevier, web
139 scraping in the case of the RSC, and through the text and data mining service at the ACS.

140 Once downloaded, synthesis paragraphs were identified using ChemDataExtractor2.1¹⁰ according to
141 previously developed protocols for identifying MOF synthesis methods.^{32,33} In this procedure,
142 chemical named entity recognition was performed using BERT⁴⁷ to identify potential reagents, and
143 part-of-speech (POS) tagging was carried out on the remaining tokens to interpret sentence
144 grammar. Chemical quantities were identified from the POS tags as CD-NN bigrams (phrases
145 consisting of a cardinal number followed by a noun), and regex matching of the noun against a
146 library of SI units. Synthesis paragraphs were identified as containing three or more chemical named
147 entities and three or more chemical quantities, after which each paragraph was extracted as plain
148 text for manual confirmation and later analysis.

149 Once confirmed that each extracted paragraph contained a synthesis procedure, hierarchical
150 grammar parsing was performed in the ChemicalTagger software¹¹ to associate chemical named
151 entities with quantities and specific synthesis actions (termed *ActionPhrases*). These were stored as
152 nested tags within an XML document.

153 *Synthesis protocol extraction*

154 To interpret and compare synthesis protocols against one another, data about synthesis steps,
155 conditions, and chemicals involved had to be converted from nested XML data into useful
156 information using the software developed in this study. To perform this, XML data extracted from
157 ChemicalTagger was recursively parsed into strings within a pandas^{48,49} DataFrame object such that
158 each row consisted of a single *ActionPhrase*, its associated time and temperature, and details of any
159 chemical entity involved.

160 Chemical identities were first confirmed by cross-referencing identified chemical names against the
161 PubChem database⁵⁰ using the pubchempy python library
162 (<https://pubchempy.readthedocs.io/en/latest/index.html>). From this, a unique identifier for each
163 individual chemical was generated, enabling extraction of key information about each chemical and
164 summation of identical chemicals together. To prevent semantically identical reagents from being
165 considered separately (e.g. zinc nitrate and their hydrates), PubChem identifiers were supplemented
166 with structural information gathered from the cheminformatics tool RDkit.⁵¹ Specifically, chemicals
167 whose formulae contained the elements zinc or cobalt, as well as the nitrate, acetate, sulfate, and
168 imidazole substructures were separately identified.

169 Then, numerical quantities associated with each chemical were calculated. To do this, chemical
170 quantities were categorised by type from the structured XML output of ChemicalTagger (e.g. by
171 volume, moles, mass etc.), and parsed into physically meaningful units with the pint python library
172 (<https://pint.readthedocs.io/en/0.20.1/index.html>). To prevent double-counting in situations where
173 two units were mentioned, e.g. by the common phrase "5 g of [reagent] (0.8 mmol)," only a single

174 unit type was considered for each chemical entity according to the priority list (moles > mass >
 175 volume). These units were then converted into moles using the molecular mass identified from the
 176 PubChem identity. In the case of converting volume to moles, densities were estimated from the
 177 ChEDL database of critical point properties⁵² using the COSTALD method.⁵³ Once chemical identities
 178 and quantities had been fully converted, these were aggregated into a single bill of materials for
 179 each synthesis (visualised in Table 1). Conditions (i.e. time and temperature values) were similarly
 180 parsed from strings into meaningful units using the pint python library, and stored as minutes and
 181 degrees Kelvin, respectively.

182 *Table 1 – example of a synthesis protocol bill of materials taken from reference* ⁵⁴

PubChem Identifier	Chemical name	Original quantities	Amount (millimoles)
12749	2-methylimidazole	0.24 g, 3.4 mmol	3.4
15865313	Zn(NO3)2.6H2O	0.956 g, 3.2 mmol	3.2
6212	Chloroform	40 mL	500
6228	DMF	70 mL	1210

183
 184 Finally, to reduce semantically meaningless differences between different synthesis sequences,
 185 synthesis actions were grouped using a similar technique to the recently developed ULSA for
 186 inorganic nanomaterials syntheses.⁵⁵ Synthesis actions were categorised as either being related to
 187 “addition,” “extraction,” “reaction,” or “other” (Table 2) and collocated steps of the same kind were
 188 grouped together. A fifth category, “start,” was used to signify opening statements of synthesis
 189 protocols (e.g. “ZIF-8 was produced by our previously published method”), which would otherwise
 190 be miscategorised as an “extraction” or “other” action. “Start” actions were then excluded from
 191 further analysis.

192 *Table 2 – Relationship between ChemicalTagger-identified ActionPhrase types and aggregated action types used here.*

Action type	ActionPhrase
“addition”	Add, Dissolve, Stir
“reaction”	ApparatusAction, Synthesize, Wait
“extraction”	Degass, Dry, Extract, Filter, Partition, Precipitate, Purify, Quench, Recover, Remove, Yield
“other”	Concentrate, Cool, Heat

193
 194 *Grouping similar synthesis protocols together*

195 To group synthesis protocols together, we related individual syntheses to one another by the
 196 identity of the reagents used only. To calculate the mathematical relationship between different
 197 synthesis protocols the list of chemicals was first vectorised, creating a numerical representation of
 198 the chemical combination used in each synthesis. Briefly, an $M \times N$ matrix was created, where M is
 199 the number of synthesis protocols, and N is the number of unique chemicals present across all of
 200 synthesis protocols studied. To reduce noise in the data, only synthesis protocols containing 2-

201 methylimidazole were considered, and metal sources were grouped by chemical substructures as
202 described previously. In total, 139 unique chemicals were identified across 1134 synthesis protocols.

203 For each synthesis protocol, a vector was generated using the term frequency–inverse document
204 frequency algorithm (TF-IDF), a commonly used text mining method to estimate the importance of
205 words in a group of documents.⁵⁶ The TF-IDF algorithm weights the frequency of a word used in each
206 document against its frequency across the group of documents – words present in many documents
207 are given a low weight, while words occurring in only rarely are given a high weight. This is shown in
208 Equation 1, which calculates the weight of word t in the individual document d as part of the group
209 of documents D , where f is the frequency the word occurs. As in this study the “words” are chemical
210 names, common chemicals like methanol are afforded a low weight, while rarer chemicals like CTAB
211 are afforded a relatively higher weight.

212

$$213 \quad tfidf(t, d, D) = f_{t,d} \cdot \log_{10} \left(\frac{1 + n}{1 + f_{t,D}} \right)$$

214 Once the chemical identities had been vectorised, similarity was calculated by the DBSCAN clustering
215 method.⁵⁷ DBSCAN calculates the local density of data points in Euclidean space (synthesis protocols
216 in the case of this study), defined as the number of neighbours closer than a threshold distance from
217 each data point. Clusters are identified as disconnected regions containing a high density of data
218 points, while isolated data points with no connection to a larger cluster as identified as noise.

219 To visualise the results of the clustering analysis, the high dimensional data were projected into two
220 dimensions using the t-distributed stochastic neighbour embedding (t-SNE) method.⁵⁸ To do this the
221 algorithm calculates the distances between each datapoint in high dimensional space, and estimates
222 low-dimensional coordinates for each datapoints which preserves the distance between each point
223 and its neighbours.

224 **Results and discussion**

225 *Validation against manually-extracted information*

226 To perform a quantitative meta-analysis of ZIF-8 synthesis, we first demonstrate the validity of the
227 information extracted by comparing the performance of our text mining approach against a
228 manually identified “ground truth” from a small number of papers sourced from the NIST database
229 of emerging adsorbent materials. Using this database served two purposes: it was sufficiently small
230 to provide a tractable number of articles for high-fidelity analysis, and each synthesis report was
231 confirmed to contain ZIF-8 by the isotherm data provided. Overall, 44 publications describing ZIF-8
232 synthesis were identified, of which full information could be extracted for 42. The manuscripts were
233 downloaded from their publisher, synthesis paragraphs manually identified, and synthesis
234 information extracted both manually and using our software. In all cases, data reported within the
235 paper and manually collated were considered as the ground truth.

236 From these paragraphs, three key parameters were extracted: a sequence of synthesis actions
237 taken, a table of constituent chemicals, and the reaction conditions (i.e. temperatures and quoted
238 times). For each parameter, the F1-score was calculated providing a numeric score for each text
239 mining task compared against the manually-extracted ground truth. Extracted chemical identities
240 were cross-referenced against the PubChem database of compounds to act as both a unique
241 identifier and source of key information about each species. Finally, physical quantities – the values
242 of time, temperature, and chemical quantity – were converted from plain text to numerical units

243 using the pint python library and compared against their manually extracted counterparts. These
244 data are summarised in Table 3.

245 *Table 3 – Parsing fidelity metrics as a percentage for manually-labelled quantities in the NIST ISODB corpus of ZIF-8*
246 *synthesis procedures.*

<i>Metric</i>	<i>Precision</i>	<i>Recall</i>	<i>F1-score</i>	<i>Matching quantities</i>
<i>Synthesis actions</i>	79	53	63	-
<i>Aggregated actions</i>	94	94	93	-
<i>Reagent identification</i>	60.2	77.4	66.7	81.7
<i>Temperature parsing</i>	74.4	72.5	73.4	68.3
<i>Time parsing</i>	73.8	91.2	81.6	73.8

247

248 Individual synthesis actions were relatively poorly identified with text mining, with a F1-score of ca.
249 60%. This was primarily due to the low recall (i.e. true positive) rate of action parsing. Inspection of
250 the synthesis paragraphs themselves showed that actions that were implicitly repeated, for example
251 in the phrase “washed with water and methanol subsequently for 3 times”,⁵⁹ were not captured by
252 ChemicalTagger thereby leading to lower scores. Conversely, when synthesis actions were converted
253 to their conceptual types and aggregated, the F1-score increased significantly to over 90% indicating
254 that all synthesis stages were identified even if the specific *ActionPhrases* themselves were not.
255 Therefore, we conclude that the text mining captures the essence of the synthesis protocol, but is
256 unable to fully summarise the semantics of synthesis due to “linguistic noise” i.e. variability between
257 different authors writing styles.

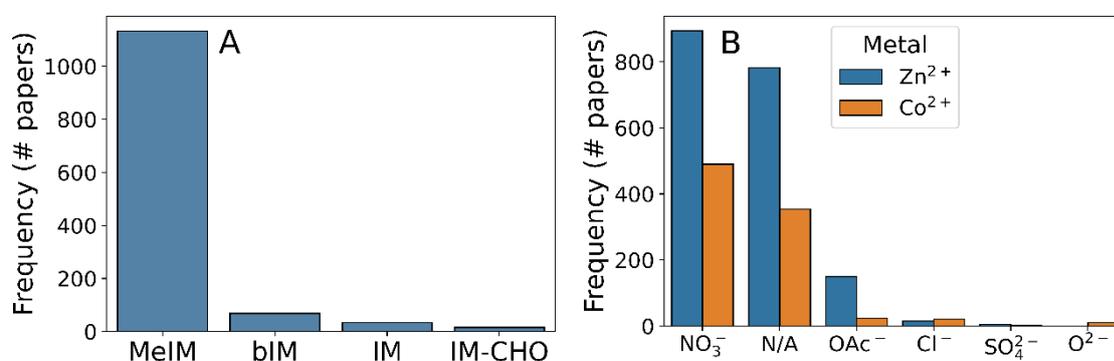
258 In terms of synthesis parameters, F1-scores and quantity matching were between 60-80% in all
259 cases. These range of scores are slightly lower than previous text-mining efforts, which generally
260 score between 60-98%.^{1,60} We ascribe this relatively low score to more stringent criteria used in this
261 study: as we define true positive to be the successful identification of a PubChem database entry,
262 precision is lowered when cross-referencing fails. This is further exacerbated by the presence of
263 typographical errors and colloquial chemical names which are not recognised by an automated
264 PubChem database search (e.g. 2-methylimidazole or 2-MeIM, rather than 2-methylimidazole).
265 Failure to successfully convert numerical quantities similarly reduced the F1-score during time and
266 temperature parsing.

267 In sum, while individual synthesis features could be reliably extracted using the methods developed
268 here, it is currently impossible to reliably reproduce the entirety of any specific synthesis protocol.
269 To achieve such high-fidelity reproduction, methods would have to be developed to estimate the
270 completeness of a synthesis protocol, requiring a much larger set of manually-labelled synthesis
271 sequences, similar to that developed by Wang *et al.* for individual synthesis actions.⁵⁵ Efforts to
272 create such a dataset are ongoing in our research group. Instead, further analysis in this study is
273 performed by compiling a group of similar synthesis protocols to extract a representative aggregate
274 of synthesis details, hence enabling quantitative meta-analysis.

275 *Interpreting ZIF-8 synthesis strategies*

276 Given the effectiveness of our text mining methods to extract synthesis information from text, we
277 progressed to a larger dataset of 3179 experimental synthesis reports of ZIF-8. From this dataset we
278 processed 1600 synthesis protocols, enabling strong statistical analysis of the synthesis options
279 which have been explored.

280 We first analysed the reagent compounds used during synthesis, which should consist of 2-
281 methylimidazole and Zn salts only. As can be seen in Figure 2, this is not the case: while
282 methylimidazole was by far the most common linker molecule mentioned (Figure 2A), 34% of the
283 synthesis protocols mentioned cobalt salts. In fact, 32% of the synthesis protocols omitted zinc
284 entirely, indicating that these were synthesis protocols of ZIF-67 instead – the cobalt equivalent of
285 ZIF-8. The remaining cobalt-mentioning synthesis protocols also contained zinc, indicating that they
286 may be mixed-metal systems. This ambiguity highlights some of the key nomenclature issues with
287 MOF materials – ZIF-8 and -67 are practically the same material in terms of synthesis protocol but
288 this proximity is not reflected in the common name. The use of tools such as MOFid³⁰ can avoid this
289 linguistic ambiguity, even accurately describing the continuous transition between the two
290 frameworks.



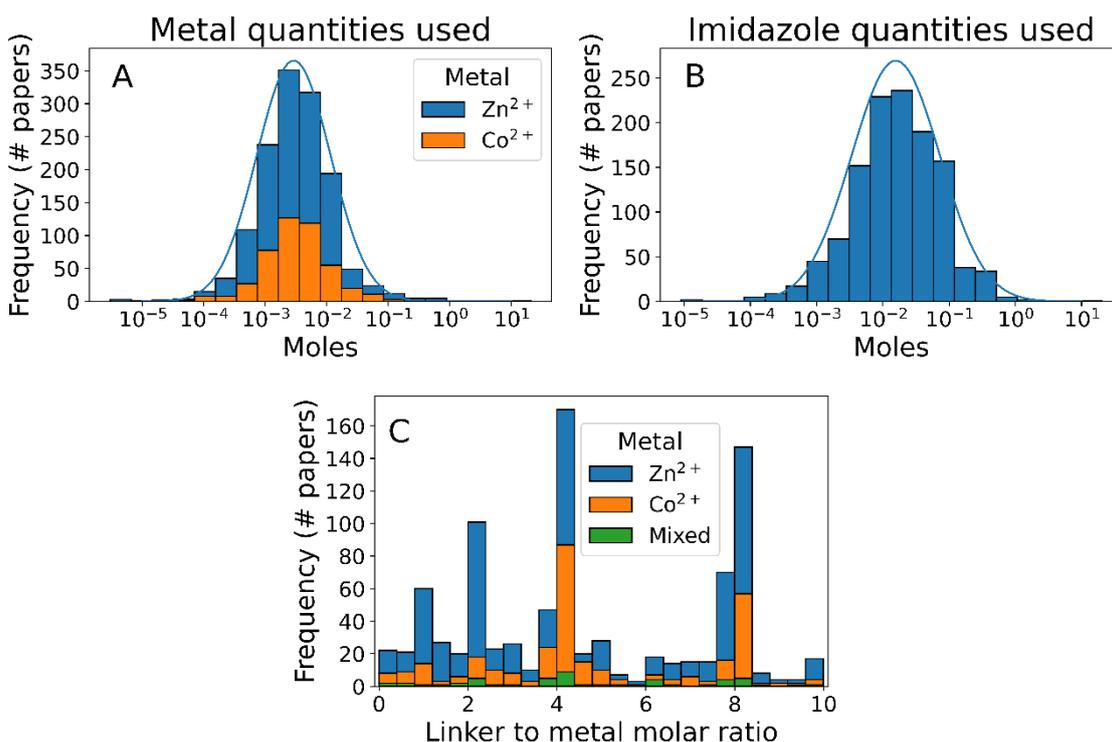
291

292 *Figure 2 – Histograms of reagent compound frequency in ZIF-8 syntheses, broken down by (A) linker choice and (B) metal*
293 *choice. Abbreviated chemical names refer to: MeIM – 2-methylimidazole; bIM – 2-benzylimidazole; IM – imidazole; IM-CHO*
294 *– imidazole-2-carbaldehyde*

295 To further analyse the reagents used we grouped the metal salts used by anion type (Figure 2B),
296 assuming that there was no consequence of using anhydrous versus hydrated salts. Nitrate was the
297 most commonly used counterion, being present in 75% of syntheses. Ambiguous mentions of zinc
298 and cobalt compounds were present in 17.2% of the 1600 protocols, encompassing minor zinc salts
299 (e.g. Zn(OH)₂ in the case of reference 61), indirect reference to zinc precursors in synthesis (e.g. “The
300 sample obtained with Zn”⁶²), or mis-identified zinc compounds due to word tokenisation errors (e.g.
301 “Firstly, 645 mg (2.469 mmol) of Zn (NO₃)₂ .4H₂O was dissolved”⁶³, where the space character
302 between “Zn” and its counterions causes incorrect chemical parsing). Aside from nitrates and
303 ambiguous mentions, the only other commonly-mentioned metal salt was zinc acetate (present in
304 11.5% of synthesis protocols). The presence of chloride, acetate, and oxide precursors indicate that
305 the synthesis is compatible to a range of electrolyte environments, agreeing with experimental
306 reports which have shown that counterion choice significantly alters crystal nucleation and growth
307 rates.^{64,65} Despite the utility of these other salts, the overwhelming popularity of nitrate counterions
308 found during our analysis indicates that other factors e.g. cost may have been prohibitive to their
309 widespread adoption.

310 In addition to reagent identity, our text mining method provides information about the quantity of
311 each reagent used, enabling analysis of synthesis protocol scale and reaction stoichiometry (Figure
312 3). The scale of ZIF-8 synthesis follows approximately a log-normal distribution, with 95% of

313 synthesis using 0.18-46 millimoles of metal ions and 0.73-330 millimoles of 2-methylimidazole
 314 (Figure 3A and B, respectively), demonstrating the flexibility of ZIF-8 synthesis with respect to scale.
 315 In terms of reaction stoichiometry, most synthesis protocols use an excess of linkers compared to
 316 the stoichiometric ratio of 2:1 (Figure 3C). This excess has been shown to control particle sizes by
 317 slowing the rate of crystal growth,^{38,66-68} although few synthesis protocols use a higher ratio than
 318 8:1. Interestingly, despite clear evidence that excess concentration of metal ions forms undesired by-
 319 products such as $\text{Zn}(\text{OH})(\text{NO}_3)(\text{H}_2\text{O})$,^{43,68-70} 6% of the synthesis protocols analysed used a molar ratio
 320 of 1:1 or lower.



321

322

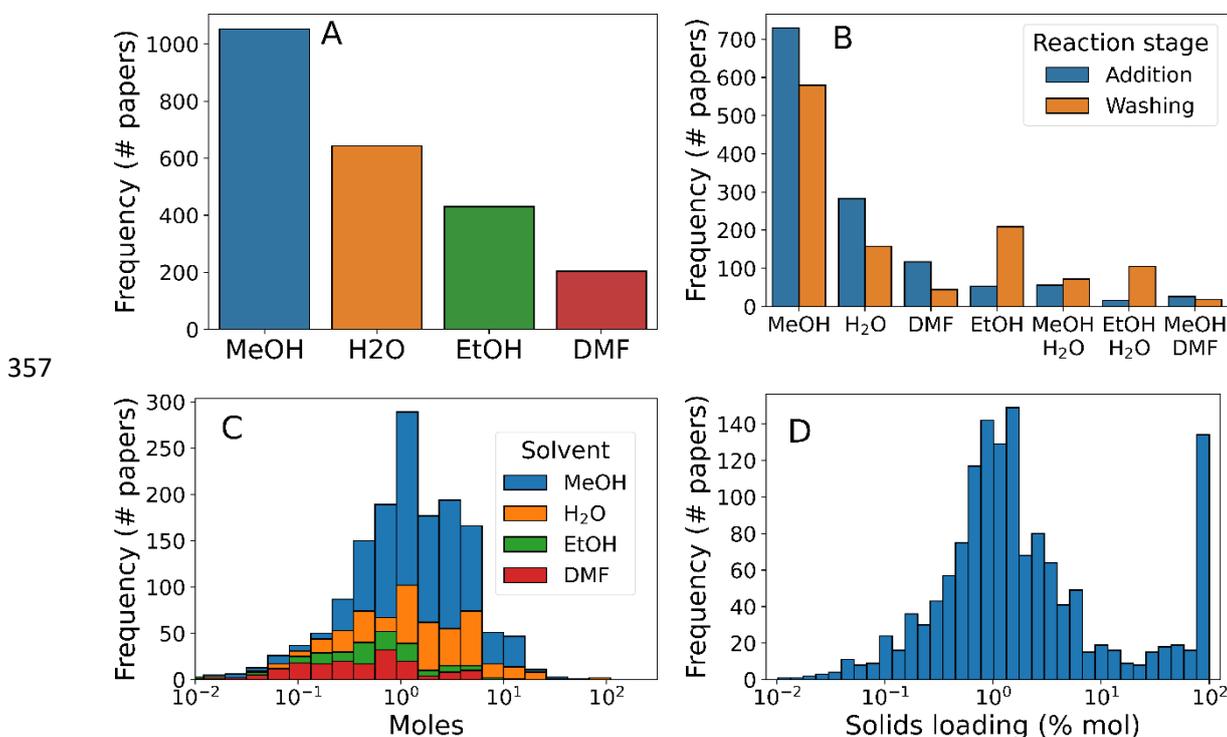
323 *Figure 3 - Histograms of reagent quantities used. (A) metals, (B) linkers, and (C) metal/linker ratios broken down by*
 324 *synthesis metal. Where multiple variables are plotted in A and C, data bars are stacked on top of one another.*

325 After considering reagents, the next most important aspect of a synthesis protocol lies in the choice of
 326 solvent environment for the reaction. Solvent choice has ramifications on the reaction mixture
 327 dielectric constant, in turn dictating factors such as reagent solubility and reaction kinetics. Further,
 328 the choice between protic and aprotic solvents, can accelerate reaction mechanisms relying on
 329 proton transfer, such as the linker deprotonation present during ZIF-8 synthesis.⁶⁶ Finally, overall
 330 reaction concentration is critical for determining whether the reaction mixture will act as an ideal
 331 solution, and in terms of the relative mass efficiency of the synthesis, both of which have
 332 consequences in terms of synthesis protocol viability in terms of scaleup to process-level
 333 manufacture.

334 The vast majority of synthesis protocols studied here contain one of methanol, ethanol, water, and
 335 DMF. Methanol was by far the most frequently mentioned solvent, present in 66% of synthesis
 336 protocols (Figure 4A), followed by water (40% of synthesis protocols), ethanol (27%), and finally DMF
 337 (12%). Less frequently used solvents included chloroform (1.4%), toluene (1.0%), and ethylene glycol
 338 (0.88%). To analyse the usage of each solvent present, we separated them by “synthesis” and
 339 “workup” procedure steps, as well as incorporating binary solvent mixtures (Figure 4B). This analysis
 340 revealed that, while ethanol was the third most prevalent solvent overall, it was the second most

341 common solvent used for washing and purification (and the fifth most common reaction solvent).
 342 Mixed solvent systems, primarily methanol-water, were present in 8% of syntheses presumably to
 343 tune the reaction dielectric and proton transfer catalysis rate.⁷¹

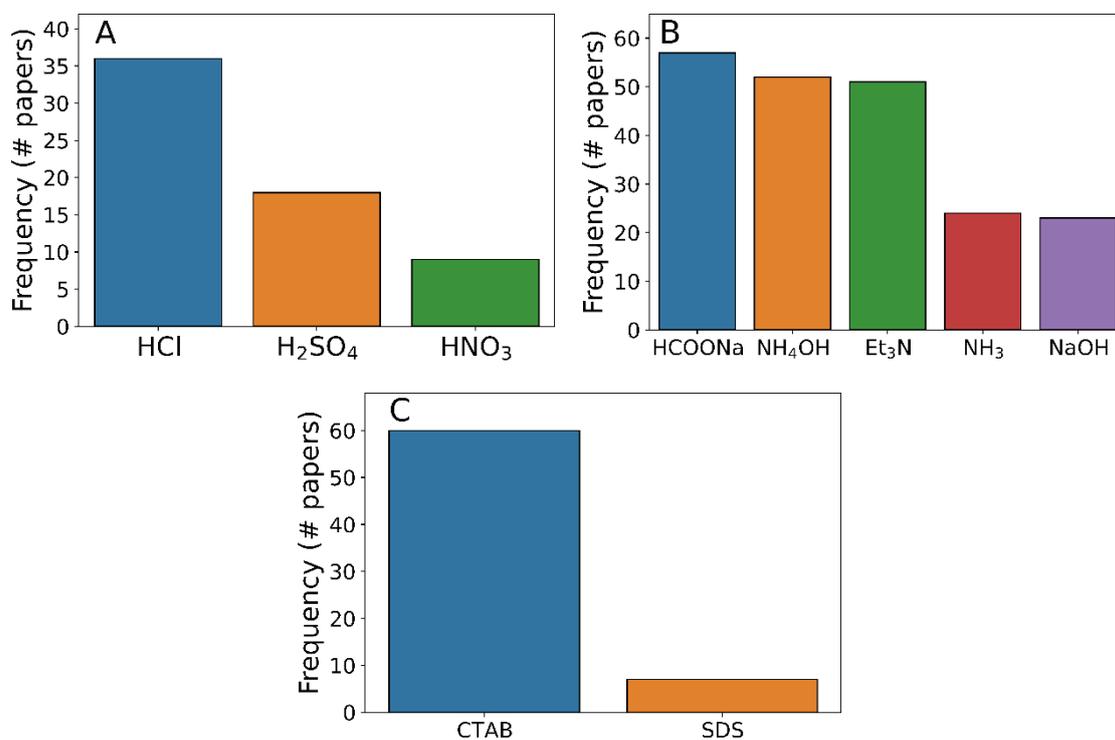
344 The distribution of solvent quantities used within the syntheses studied (Figure 4C) showed that
 345 each solvent followed approximately lognormal distributions. Both DMF and ethanol were used in
 346 smaller quantities than methanol or water (means of 0.4, 0.6, 1.4, and 1.6 moles per synthesis,
 347 respectively), indicating that the latter two solvents were more appropriate for scaling up the
 348 synthesis. Finally, we analysed the total solids concentration of synthesis protocols by dividing total
 349 reagent amounts by the solvent amounts used (Figure 4D). As with individual reagent
 350 concentrations, the total solids concentration followed an approximately log-normal distribution
 351 between 0.1-10 %mol. Separately, 7.7% of synthesis protocols had a solids loading of approximately
 352 100 %mol - signifying mechanochemical synthesis protocols. Although mechanochemistry is a
 353 promising synthesis route due to its high yields⁷² and low environmental impact⁷³ compared to
 354 conventional solvent synthesis methods, the relatively low popularity may be explained due to
 355 practical difficulties of mechanochemical synthesis e.g. prevention of hot-spot formation in the
 356 reaction vessel.⁷⁴



359 *Figure 4 - histograms of solvent usage in ZIF-8 synthesis. (A) frequency of solvent mentions in all synthesis procedures, (B)*
 360 *frequency of solvent usage broken down by stage of the procedure, (C) quantity of solvent used, broken down by solvent*
 361 *type, and (D) total solids loading. Where multiple variables are plotted in C, data bars are stacked on top of one another.*

362 In addition to reagents and solvents, ancillary chemicals such as surfactants, pH modifiers, and
 363 modulators are often key to ensure the success of MOF syntheses as well as dictating secondary
 364 particle characteristics such as size and crystal form. Three chemical types were prevalent within the
 365 synthesis protocols studied: acids, bases, and surfactant compounds. Unlike solvents and reagents,
 366 no individual ancillary chemical was identified in more than 3.5% of synthesis protocols (Figure 5).
 367 However, bases were present in 18% of all the synthesis protocols analysed, carrying out the
 368 important role of deprotonating the linker molecule in the reaction mixture. From the variety of

369 distinct molecules used for this role, it appears that no molecular recognition occurs, simply pH
 370 control. Despite the requirement for methylimidazole deprotonation for the reaction to progress,
 371 acids were detected in 6.3% of syntheses, however from inspection of the individual synthesis
 372 protocols acids only appeared during post-synthetic modification of the ZIF-8 materials e.g. after
 373 carbonisation⁷⁵ or impregnation into silicas.⁷⁶ Finally, surfactants like cetyltrimethylammonium
 374 bromide (CTAB) or sodium dodecylsulfate (SDS) were present in 4.6% of synthesis protocols, being
 375 used to slow the growth of individual ZIF-8 crystals and therefore control the particle shape.^{59,77}

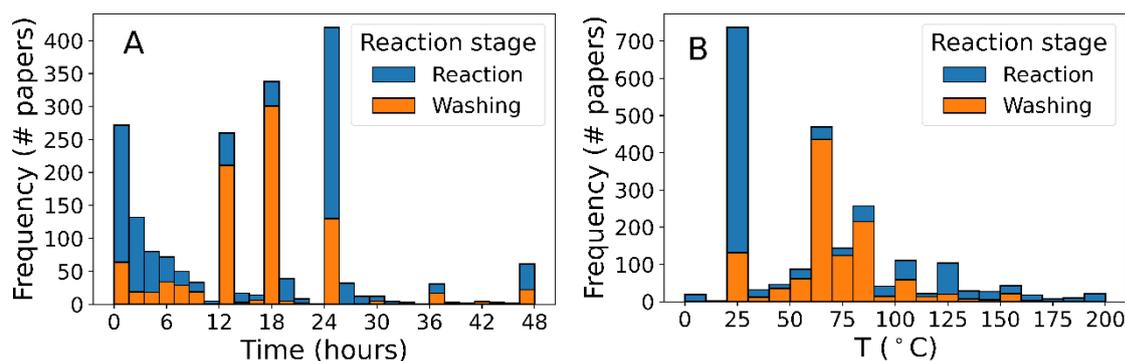


376

377

378 *Figure 5 -Histograms of ancillary chemical prevalence in ZIF-8 synthesis. (A) acids, (B) bases, and (C) surfactants.*

379 While it is possible to identify broad differences in synthesis strategy from feedstock compounds
 380 alone, it is impossible to understand why one chemical is chosen over another without further detail
 381 about the synthesis protocol being described. For example, the modulator sodium formate has been
 382 shown to perform different roles in room-temperature syntheses compared to hydrothermal
 383 alternatives.^{44,78} In the first instance, we also consider the conditions (i.e. time and temperature)
 384 during the process. These are shown in Figure 6, demonstrating that the majority of protocols have
 385 synthesis times under six hours. Even after disregarding protocols with a reported synthesis time of
 386 0 minutes as being spurious, it is clear that synthesis can be completed very quickly. In terms of
 387 synthesis temperature, the majority of the extracted temperatures were found to be room
 388 temperature indicating that thermal driving forces were not necessary for the formation of ZIF-8.
 389 This is further corroborated by the relative lack of procedures mentioning heated reaction
 390 conditions compared to heated drying conditions (Figure 6B).



391

392 *Figure 6 – Histograms of conditions during ZIF-8 synthesis processes. (A) total time elapsed and (B) temperatures used*
 393 *during synthesis. Annotational on (B) indicate the boiling points of the four most common solvents identified. Data are*
 394 *broken down by reaction step type as defined in Table 2. Where multiple variables are plotted, data bars are stacked on top*
 395 *of one another.*

396 Overall, the tools developed in this study provide wide-ranging descriptive statistics of various ZIF-8
 397 synthesis routes. The data generated are an excellent addition to existing literature review methods,
 398 facilitating the interpretation of different synthesis aspects e.g. reagent choices, stoichiometric
 399 ratios and reaction conditions. From these data we are able to identify gaps in the existing literature
 400 or synthesis conditions most likely to succeed, as well as providing useful input data for later
 401 techno-economic analysis.

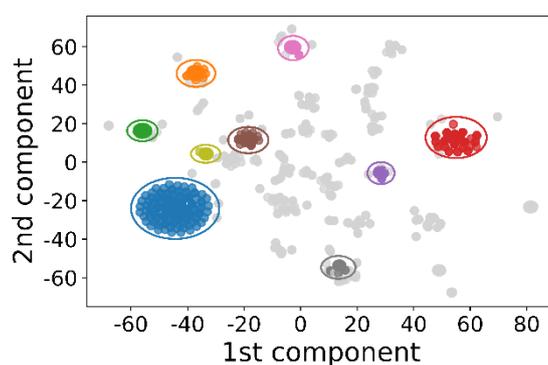
402 *Harnessing synthesis information for accelerated methodology development*

403 While the analysis performed is useful as a means of understanding the ZIF-8 reaction system, a key
 404 aim of this study was to systematize the collective synthesis knowledge for the material, thereby
 405 connecting synthesis protocols to some key performance indicators either of the synthesis (e.g.
 406 yield) or material (e.g. crystal form, surface area). One crucial barrier to this goal was the correlation
 407 of material performance data with synthesis protocol information: research papers are inconsistent
 408 in reporting of material properties (primarily as different quality metrics are used depending on the
 409 motivation of the original research), and the sample naming conventions used within research
 410 articles prevent unambiguous linking between the described protocols and materials produced. For
 411 example, while a synthesis paragraph might detail the synthesis of “nano-sized ZIF-8,” later mentions
 412 in the text may be labelled differently e.g. “ZIF-8_{nano},”⁷⁹ confounding attempts for automated
 413 identification of reaction products using regular expressions.³³ While this issue will undoubtedly be
 414 resolved by the adoption of transformer-based language models such as BERT⁸⁰ and GPT-4,⁸¹ such
 415 models became available only recently and the scientific community, including our group, is in the
 416 process of probing their extension to scientific data mining. In fact, the current study highlighted a
 417 number of issues with the current structure and completeness of reported synthetic protocols,
 418 understanding of which will be very helpful in engineering and fine-tuning GPT-based models.

419 As a result, the analysis performed in this study can only provide insight into how the MOF material
 420 is made rather than linking different synthesis features to specific outcomes like yield or quality. In
 421 the absence of such synthesis outcome information, we instead focus on how best to prepare the
 422 information gathered in this study for the generation of predictive models for ZIF-8 materials quality.
 423 A key challenge when attempting to optimise synthesis protocols either through systematic
 424 experimentation⁵ or by training machine learning models⁶ is the high dimensionality of the
 425 information contained in each synthesis. For example, 8 unique reagent chemicals were discussed in
 426 the previous section – 3 metal sources, 1 linker, and 4 solvents – meaning that to fully explore the 8-
 427 dimensional chemical space alone, N^8 experiments would be required (where N is the number of

428 quantity values tested for each variable). Even when limits on the complexity of each individual
 429 reaction are used – i.e. to contain a maximum of two metal salts and solvents – the dimensionality is
 430 only reduced to $12(N^5)$ experiments. While theoretically this dimensionality would scale with the
 431 number of synthesis steps used, we were unable to identify meaningfully distinct groups of synthesis
 432 actions (data not shown here, for brevity) and hence did not consider the sequence as impacting the
 433 synthesis outcome.

434 To enable faster and more efficient searching of the synthesis phase space, we used clustering to
 435 identify lower-dimensional sub-regions of the synthesis phase space which have been widely
 436 researched in experimental papers – essentially using a chemical combination’s popularity as a proxy
 437 for its importance. The chemical identities used were encoded using TF-IDF vectorisation, then
 438 similar synthesis protocols were grouped by their density in the encoded space. The outcome of this
 439 clustering analysis is visualised using a 2-d projection in Figure 7 and summarised in Table 4, where
 440 the distance between points is indicative of each protocol’s similarity to its neighbours. Eight clusters
 441 of reagent combinations were identified each containing 2-4 chemicals of a total of 6 reagents. We
 442 posit that these clusters represent well defined strategies to synthesize ZIF-8, which can be explored
 443 separately, therefore reducing the total amount of information required to explore these regions of
 444 the synthesis space.



445

446 *Figure 7 – 2-dimensional representation of the chemical combination space for ZIF-8 synthesis, generated using the t-SNE*
 447 *algorithm. Major synthesis pathways are identified using the DBSCAN clustering method and colour coded, while noise data*
 448 *is shown in light grey. Clusters are circled and described in Table 4.*

449 *Table 4 - Cluster labels and common features from Figure 7. N.B. all synthesis protocols included 2-methylimidazole, which*
 450 *was omitted for brevity.*

Cluster number (colour)	Common chemicals	Protocols in cluster
1 (blue)	zinc, nitrate, methanol	225
2 (red)	cobalt, nitrate, methanol	147
3 (brown)	zinc, nitrate, water	50
4 (orange)	Zinc, nitrate	39
5 (green)	Zinc, cobalt, nitrate, methanol	31
6 (pink)	cobalt, nitrate, water	25
7 (purple)	Zinc, acetate, water	22
8 (olive)	Zinc, nitrate, methanol, water	20
9 (grey)	Zinc, nitrate, DMF	17

451

452 The well-defined synthesis strategies clustered in Figure 7 are notably different from the analysis
453 performed in the previous section. In the first instance, ethanol was fully absent signifying its
454 insignificance as a reaction solvent and matching the earlier analyses. Separately, acetate salts are
455 only identified in one cluster and only associated with water. This association is due to the lack of
456 solubility of zinc acetate in methanol (ca. 15 g/L cf. 430 g/L in water), information which can only
457 otherwise be gained by specific knowledge of the chemistry of zinc acetate. While obvious to those
458 who already are aware of the system, this information may otherwise be overlooked by chemists
459 naive to the intricacies of ZIF-8 synthesis – an example of chemical intuition.⁶ Therefore, clustering of
460 similar synthesis protocols together can help users to avoid some common pitfalls when planning
461 experiments for the first time.

462 Finally, to demonstrate the benefit of this approach towards synthesis optimisation, we consider the
463 reduction in experiments that would be required to explore the identified popular sub-regions of the
464 synthesis space. From the clustering analysis, we identified 6 sub-regions with only 3 chemicals of
465 interest – clusters 1, 2, 3, 6, 7, and 8 in Table 4, containing only a single metal salt, 2-
466 methylimidazole, and a single solvent – and a further 2 sub-regions with 4 chemicals of interest:
467 clusters 5 and 8 containing either mixed salts or solvents. Accordingly, rather than requiring N^8 or
468 $12(N^5)$ experiments, full exploration would only require $6(N^3) + 2(N^4) \approx N^{4.4}$ experiments. To
469 illustrate the extent of dimensionality reduction in real terms, the number of experiments required
470 to explore the synthesis space are shown in Table 5 for various values of N . In combination with the
471 quantity distributions shown in Figure 3 and Figure 4, text mining and data reduction tools
472 demonstrated in this paper will provide excellent initial values for efficient searching of chemical
473 synthesis space, thereby accelerating methodology refinement for a range of nanomaterials.

474 *Table 5 – Approximate number of experiments required to fully characterise the synthesis space, for various values of N .*

	Full exploration (N^8)	Limited experimental complexity ($12(N^5)$)	Identified clusters only ($6(N^3) + 2(N^4)$)
$N = 3$	6,500	2,900	320
$N = 5$	390,000	37,500	2,000
$N = 10$	1×10^8	1.2×10^6	2.6×10^4

475

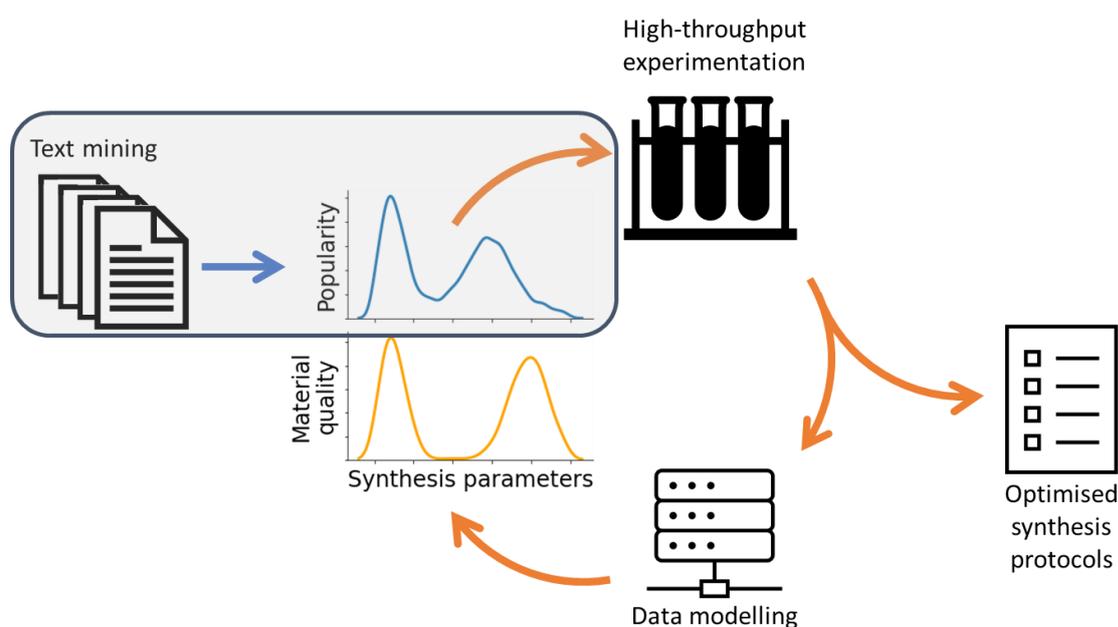
476

477 Conclusions

478 In this study, we applied text mining to the problem of synthesis methodology optimisation,
479 exploring to what extent the previously accumulated collective knowledge of a particular
480 nanomaterial can accelerate the development of reliable and scalable synthesis protocols. As the
481 first step toward this objective, in this study, we posed three research questions: first, is it possible
482 to use text mining tools to provide deep insight into a single synthetic target, rather than a
483 comprehensive overview of a family of materials? Second, is it possible to standardise the synthesis
484 details extracted as a means of performing like-for-like comparison between different studies?
485 Finally, is it possible to use this analysis to suggest optimal synthesis conditions, thereby accelerating
486 methodology development?

487 To this end, we developed software to systematically analyse nanomaterials synthesis methods
488 based on established text mining protocols. We extracted structured data to describe the details of

489 each synthesis protocol, enabling large-scale statistical analysis of the synthesis parameter space and
 490 clustering of similar methods together to identify well-explored regions of the synthesis space. We
 491 believe that this progress represents the first step in creating a closed feedback loop for the
 492 automated optimisation of experimental nanomaterials synthesis, visualised in Figure 8. In this
 493 feedback loop text mined information can identify common limits to parameters as well as low-
 494 dimensional sub-regions of interest in the synthesis space. By using this information as initial
 495 conditions for iterative high-throughput experimentation, the search for synthesis protocols
 496 optimised against any target material quality metric can be greatly accelerated.



497

498 *Figure 8 - Scheme of a synthesis protocol optimisation feedback loop. Work carried out in this study is shaded in grey.*

499 As a case study to demonstrate the utility of this approach, we performed a quantitative meta-
 500 analysis of 1600 synthesis methods for the common MOF ZIF-8. Using this framework, we identified
 501 key aspects of the synthesis including the range of chemicals used as reagents, solvents, and
 502 ancillary modulators/pH modifiers. We extracted information about the quantity of each reagent
 503 used during the synthesis, enabling us to identify the distribution of synthesis scales, reagent ratios,
 504 and reaction mixture solids concentration, as well as reaction times and temperatures. Further
 505 insight was gathered by cross-referencing chemicals mentioned against the stage they were
 506 introduced into the synthesis protocol – for example identifying that ethanol is primarily used as a
 507 washing solvent rather than in the reaction medium. We demonstrated how the quantitative meta-
 508 analysis performed here can assist in systematic searches of the synthesis phase space by identifying
 509 both low-dimensional regions of interest and the distribution of synthesis parameters. As a result,
 510 we were able to reduce the number of hypothetical experiments required to optimise ZIF-8
 511 significantly. Notably, while we considered MOF materials as a case study in this work, the methods
 512 developed here are general to any synthesis type. Particularly, we envisage they will be useful the
 513 systematising understanding of other emerging nanomaterial systems such as mesoporous
 514 (organo)silicas, covalent organic frameworks, and polymers of intrinsic microporosity.

515 Despite the deep insight we were able to gain into the synthesis system of ZIF-8, the current study
 516 also identified significant challenges associated with developing a true “synthetic oracle” for
 517 predicting the ideal synthesis parameters for any given material. While we were able to identify and
 518 extract information about the synthesis, we were unable to reliably connect the quality of the

519 material produced to the methods themselves (e.g. by identifying specific yield or surface area). A
520 crucial next step is therefore to adopt state of the art transformer-based methods e.g. BERT or GPT-
521 4 to better interpret the entire research article as a single unit and therefore identify implicitly
522 described synthesis protocols (e.g. tabulated changes to individual synthesis parameters). A second
523 challenge lies in the estimating the viability of synthesis parameters extracted during text mining or
524 proposed by generative models, preventing automated reproduction of a synthesis protocol without
525 human oversight and validation. Finally, as has been discussed elsewhere, the synthesis protocol
526 extraction methods developed here can only build from published information, which is biased
527 towards the most successful synthesis methods only. More comprehensive reporting of synthesis
528 information using structured formats akin to the crystallographic information file format would
529 enable far more wide-reaching analysis to be performed.

530 In summary, the methods developed in this study acts as a preliminary approach for the large-scale
531 standardisation and analysis of experimental synthesis data, representing the first step in creating a
532 closed feedback loop for the automated optimisation of experimental nanomaterials synthesis. By
533 interfacing with automated and high throughput reactionware e.g. through integration of the XDL
534 chemical programming language, methodology development will be significantly accelerated
535 thereby easing the adoption of nanomaterials at larger scales and in new settings.

536 **References**

- 537 1. Kononova, O. *et al.* Opportunities and challenges of text mining in materials research.
538 *iScience* **24**, 102155 (2021).
- 539 2. Kim, E., Huang, K., Kononova, O., Ceder, G. & Olivetti, E. Distilling a Materials Synthesis
540 Ontology. *Matter* **1**, 8–12 (2019).
- 541 3. Baer, D. R. & Gilmore, I. S. Responding to the growing issue of research reproducibility. *J. Vac.*
542 *Sci. Technol. A* **36**, 068502 (2018).
- 543 4. Baer, D. R., Munusamy, P. & Thrall, B. D. Provenance information as a tool for addressing
544 engineered nanoparticle reproducibility challenges. *Biointerphases* **11**, 04B401 (2016).
- 545 5. Burger, B. *et al.* A mobile robotic chemist. *Nature* **583**, 237–241 (2020).
- 546 6. Moosavi, S. M. *et al.* Capturing chemical intuition in synthesis of metal-organic frameworks.
547 *Nat. Commun.* **10**, 539 (2019).
- 548 7. Wilbraham, L., Mehr, S. H. M. & Cronin, L. Digitizing Chemistry Using the Chemical Processing
549 Unit: From Synthesis to Discovery. *Acc. Chem. Res.* **54**, 253–262 (2021).
- 550 8. Hammer, A. J. S., Leonov, A. I., Bell, N. L. & Cronin, L. Chemputation and the Standardization
551 of Chemical Informatics. *JACS Au* **1**, 1572–1587 (2021).
- 552 9. Jessop, D. M., Adams, S. E., Willighagen, E. L., Hawizy, L. & Murray-Rust, P. OSCAR4: A flexible
553 architecture for chemical textmining. *J. Cheminform.* **3**, 1–12 (2011).
- 554 10. Mavračić, J., Court, C. J., Isazawa, T., Elliott, S. R. & Cole, J. M. ChemDataExtractor 2.0:
555 Autopopulated Ontologies for Materials Science. *J. Chem. Inf. Model.* **61**, 4280–4289 (2021).
- 556 11. Hawizy, L., Jessop, D. M., Adams, N. & Murray-Rust, P. ChemicalTagger: A tool for semantic
557 text-mining in chemistry. *J. Cheminform.* **3**, 1–13 (2011).
- 558 12. Vazquez, M., Krallinger, M., Leitner, F. & Valencia, A. Text mining for drugs and chemical
559 compounds: Methods, tools and applications. *Mol. Inform.* **30**, 506–519 (2011).

- 560 13. Guo, J. *et al.* Automated Chemical Reaction Extraction from Scientific Literature. *J. Chem. Inf. Model.* **62**, 2035–2045 (2022).
561
- 562 14. Beard, E. J., Sivaraman, G., Vázquez-Mayagoitia, Á., Vishwanath, V. & Cole, J. M. Comparative
563 dataset of experimental and computational attributes of UV/vis absorption spectra. *Sci. Data*
564 **6**, 1–11 (2019).
- 565 15. Cooper, C. B. *et al.* Design-to-Device Approach Affords Panchromatic Co-Sensitized Solar
566 Cells. *Adv. Energy Mater.* **9**, 1–10 (2019).
- 567 16. Jensen, Z. *et al.* Discovering Relationships between OSDAs and Zeolites through Data Mining
568 and Generative Neural Networks. *ACS Cent. Sci.* **7**, 858–867 (2021).
- 569 17. Beard, E. J. & Cole, J. M. Perovskite- and Dye-Sensitized Solar-Cell Device Databases Auto-
570 generated Using ChemDataExtractor. *Sci. Data* **9**, 1–19 (2022).
- 571 18. Zhao, J. & Cole, J. M. A database of refractive indices and dielectric constants auto-generated
572 using ChemDataExtractor. *Sci. Data* **9**, 1–11 (2022).
- 573 19. Kononova, O. *et al.* Text-mined dataset of inorganic materials synthesis recipes. *Sci. Data* **6**,
574 1–11 (2019).
- 575 20. Isazawa, T. & Cole, J. M. Single Model for Organic and Inorganic Chemical Named Entity
576 Recognition in ChemDataExtractor. *J. Chem. Inf. Model.* **62**, 1207–1213 (2022).
- 577 21. Klinger, R., Kolářik, C., Fluck, J., Hofmann-Apitius, M. & Friedrich, C. M. Detection of IUPAC
578 and IUPAC-like chemical names. *Bioinformatics* **24**, 268–276 (2008).
- 579 22. Majumdar, S., Moosavi, S. M., Jablonka, K. M., Ongari, D. & Smit, B. Diversifying Databases of
580 Metal Organic Frameworks for High-Throughput Computational Screening. *ACS Appl. Mater.*
581 *Interfaces* **13**, 61004–61014 (2021).
- 582 23. Boyd, P. G. & Woo, T. K. A generalized method for constructing hypothetical nanoporous
583 materials of any net topology from graph theory. *CrystEngComm* **18**, 3777–3792 (2016).
- 584 24. Wilmer, C. E. *et al.* Large-scale screening of hypothetical metal-organic frameworks. *Nat.*
585 *Chem.* **4**, 83–89 (2012).
- 586 25. Lee, S. Y. S. *et al.* Computational Screening of Trillions of Metal-Organic Frameworks for High-
587 Performance Methane Storage. *ACS Appl. Mater. Interfaces* **13**, 23647–23654 (2021).
- 588 26. Moosavi, S. M. *et al.* Understanding the diversity of the metal-organic framework ecosystem.
589 *Nat. Commun.* **11**, 4068 (2020).
- 590 27. Moghadam, P. Z. *et al.* Targeted classification of metal-organic frameworks in the Cambridge
591 structural database (CSD). *Chem. Sci.* **11**, 8373–8387 (2020).
- 592 28. Moghadam, P. Z. *et al.* Development of a Cambridge Structural Database Subset: A Collection
593 of Metal-Organic Frameworks for Past, Present, and Future. *Chem. Mater.* **29**, 2618–2625
594 (2017).
- 595 29. Ongari, D., Talirz, L. & Smit, B. Too Many Materials and Too Many Applications: An
596 Experimental Problem Waiting for a Computational Solution. *ACS Cent. Sci.* **6**, 1890–1900
597 (2020).
- 598 30. Bucior, B. J. *et al.* Identification Schemes for Metal–Organic Frameworks To Enable Rapid
599 Search and Cheminformatics Analysis. *Cryst. Growth Des.* **19**, 6682–6697 (2019).
- 600 31. Groom, C. R. & Allen, F. H. The Cambridge Structural Database in retrospect and prospect.

- 601 *Angew. Chemie - Int. Ed.* **53**, 662–671 (2014).
- 602 32. Luo, Y. *et al.* MOF Synthesis Prediction Enabled by Automatic Data Mining and Machine
603 Learning**. *Angew. Chemie Int. Ed.* **61**, (2022).
- 604 33. Gubsch, K., Bence, R., Glasby, L. T. & Moghadam, P. Z. DigiMOF: A Database of MOF Synthesis
605 Information Generated via Text Mining. *ChemRxiv* 1–26 (2022) doi:10.26434/chemrxiv-2022-
606 41t70.
- 607 34. Cox, C. S., Slavich, E., Macreadie, L. K., McKemmish, L. K. & Lessio, M. Understanding the Role
608 of Synthetic Parameters in the Defect Engineering of UiO-66: A Review and Meta-analysis.
609 *Chem. Mater.* **35**, 3057–3072 (2023).
- 610 35. Zhang, H. & Snurr, R. Q. Computational Study of Water Adsorption in the Hydrophobic Metal-
611 Organic Framework ZIF-8: Adsorption Mechanism and Acceleration of the Simulations. *J.*
612 *Phys. Chem. C* **121**, 24000–24010 (2017).
- 613 36. Bhattacharyya, S. *et al.* Acid gas stability of zeolitic imidazolate frameworks: Generalized
614 kinetic and thermodynamic characteristics. *Chem. Mater.* **30**, 4089–4101 (2018).
- 615 37. De Lange, M. F. *et al.* Metal-Organic Frameworks in Adsorption-Driven Heat Pumps: The
616 Potential of Alcohols as Working Fluids. *Langmuir* **31**, 12783–12796 (2015).
- 617 38. Kinoshita, M., Yanagida, S., Gessei, T. & Monkawa, A. Precursor concentration effects on
618 crystallite size and enzyme immobilization efficiency of Enzyme@ZIF-8 composite. *J. Cryst.*
619 *Growth* **600**, 126877 (2022).
- 620 39. Lee, Y. R., Do, X. H., Hwang, S. S. & Baek, K. Y. Dual-functionalized ZIF-8 as an efficient acid-
621 base bifunctional catalyst for the one-pot tandem reaction. *Catal. Today* **359**, 124–132
622 (2021).
- 623 40. Paul, A., Banga, I. K., Muthukumar, S. & Prasad, S. Engineering the ZIF-8 Pore for
624 Electrochemical Sensor Applications-A Mini Review. *ACS Omega* **7**, 26993–27003 (2022).
- 625 41. Lewis, A. *et al.* Crystallization and phase selection of zeolitic imidazolate frameworks in
626 aqueous cosolvent systems: The role and impacts of organic solvents. *Results Eng.* **17**, 100751
627 (2023).
- 628 42. Tsai, C.-W. & Langner, E. H. G. The effect of synthesis temperature on the particle size of
629 nano-ZIF-8. *Microporous Mesoporous Mater.* **221**, 8–13 (2016).
- 630 43. Kida, K., Okita, M., Fujita, K., Tanaka, S. & Miyake, Y. Formation of high crystalline ZIF-8 in an
631 aqueous solution. *CrystEngComm* **15**, 1794–1801 (2013).
- 632 44. Cravillon, J. *et al.* Formate modulated solvothermal synthesis of ZIF-8 investigated using time-
633 resolved in situ X-ray diffraction and scanning electron microscopy. *CrystEngComm* **14**, 492–
634 498 (2012).
- 635 45. Nordin, N. A. H. M. *et al.* Aqueous room temperature synthesis of zeolitic imidazole
636 framework 8 (ZIF-8) with various concentrations of triethylamine. *RSC Adv.* **4**, 33292–33300
637 (2014).
- 638 46. Van Rossum, G. & Drake, F. L. *Python 3 Reference Manual*. (CreateSpace, 2009).
- 639 47. Devlin, J., Chang, M.-W., Lee, K. & Toutanova, K. BERT: Pre-training of Deep Bidirectional
640 Transformers for Language Understanding. (2018) doi:10.48550/arXiv.1810.04805.
- 641 48. McKinney, W. Data Structures for Statistical Computing in Python. in *Proceedings of the 9th*

- 642 *Python in Science Conference* (eds. van der Walt, S. & Millman, J.) vol. 1 56–61 (2010).
- 643 49. The pandas development team. pandas-dev/pandas: Pandas. (2023)
644 doi:10.5281/zenodo.3509134.
- 645 50. Kim, S. *et al.* PubChem 2023 update. *Nucleic Acids Res.* **51**, D1373–D1380 (2023).
- 646 51. Landrum, G. RDKit: Open-source cheminformatics. (2010).
- 647 52. Bell, C. & Cortes-Pena, Y. R. Chemicals: Chemical properties component of Chemical
648 Engineering Design Library (ChEDL). (2016).
- 649 53. Thomson, G. H., Hankinson, R. W. & Thomson, G. H. A new correlation for saturated densities
650 of liquids and their mixtures. *AIChE J.* **25**, 653–663 (1979).
- 651 54. Zhang, Z. *et al.* Enhancement of CO₂ Adsorption and CO₂/N₂ Selectivity on ZIF-8 via
652 Postsynthetic Modification. *AIChE J.* **59**, 2195–2206 (2013).
- 653 55. Wang, Z. *et al.* ULSA: unified language of synthesis actions for the representation of inorganic
654 synthesis protocols. *Digit. Discov.* **1**, 313–324 (2022).
- 655 56. Sparck Jones, K. A Statistical Interpretation of Term Specificity And Its Application In Retrieval.
656 *J. Doc.* **28**, 11–21 (1972).
- 657 57. Esther, M., Kriegel, H.-P., Sander, J. & Xu, X. A density-based algorithm for discovering
658 clusters in large spatial databases with noise. in *Proceedings of the Second International*
659 *Conference on Knowledge Discovery and Data Mining* 226–231 (1996).
- 660 58. Maaten, L. van der & Hinton, G. Visualizing High-Dimensional Data Using t-SNE. *J. Mach.*
661 *Learn. Res.* **9**, 2579–2605 (2008).
- 662 59. Pan, Y. *et al.* Tuning the crystal morphology and size of zeolitic imidazolate framework-8 in
663 aqueous solution by surfactants. *CrystEngComm* **13**, 6937 (2011).
- 664 60. Swain, M. C. & Cole, J. M. ChemDataExtractor: A Toolkit for Automated Extraction of
665 Chemical Information from the Scientific Literature. *J. Chem. Inf. Model.* **56**, 1894–1904
666 (2016).
- 667 61. Liu, J. *et al.* The active roles of ZIF-8 on the enhanced visible photocatalytic activity of
668 Ag/AgCl: Generation of superoxide radical and adsorption. *J. Alloys Compd.* **693**, 543–549
669 (2017).
- 670 62. Butova, V. V. *et al.* Modification of ZIF-8 with triethylamine molecules for enhanced iodine
671 and bromine adsorption. *Inorganica Chim. Acta* **509**, 119678 (2020).
- 672 63. Samadi-Maybodi, A., Ghasemi, S. & Ghaffari-Rad, H. Ag-doped zeolitic imidazolate
673 framework-8 nanoparticles modified CPE for efficient electrocatalytic reduction of H₂O₂.
674 *Electrochim. Acta* **163**, 280–287 (2015).
- 675 64. Schejn, A. *et al.* Controlling ZIF-8 nano- and microcrystal formation and reactivity through zinc
676 salt variations. *CrystEngComm* **16**, 4493–4500 (2014).
- 677 65. Jian, M. *et al.* Water-based synthesis of zeolitic imidazolate framework-8 with high
678 morphology level at room temperature. *RSC Adv.* **5**, 48433–48441 (2015).
- 679 66. Öztürk, Z., Filez, M. & Weckhuysen, B. M. Decoding Nucleation and Growth of Zeolitic
680 Imidazolate Framework Thin Films with Atomic Force Microscopy and Vibrational
681 Spectroscopy. *Chem. - A Eur. J.* **23**, 10915–10924 (2017).

- 682 67. Hamidon, N. F., Tahir, M. I. M., Latif, M. A. M. & Abdul Rahman, M. B. Effect of altering linker
683 ratio on nano-ZIF-8 polymorphisms in water-based and modulator-free synthesis. *J. Coord.*
684 *Chem.* **75**, 1180–1192 (2022).
- 685 68. Yamamoto, D. *et al.* Synthesis and adsorption properties of ZIF-8 nanoparticles using a
686 micromixer. *Chem. Eng. J.* **227**, 145–150 (2013).
- 687 69. Chen, B., Bai, F., Zhu, Y. & Xia, Y. A cost-effective method for the synthesis of zeolitic
688 imidazolate framework-8 materials from stoichiometric precursors via aqueous ammonia
689 modulation at room temperature. *Microporous Mesoporous Mater.* **193**, 7–14 (2014).
- 690 70. Zhang, Y., Jia, Y., Li, M. & Hou, L. Influence of the 2-methylimidazole/zinc nitrate hexahydrate
691 molar ratio on the synthesis of zeolitic imidazolate framework-8 crystals at room
692 temperature. *Sci. Rep.* **8**, 1–7 (2018).
- 693 71. Albright, P. S. & Gosting, L. J. Dielectric Constants of the Methanol-Water System from 5 to
694 55°. *J. Am. Chem. Soc.* **68**, 1061–1063 (1946).
- 695 72. Lee, Y. R. *et al.* ZIF-8: A comparison of synthesis methods. *Chem. Eng. J.* **271**, 276–280 (2015).
- 696 73. Xia, W., Lau, S. K. & Yong, W. F. Comparative life cycle assessment on zeolitic imidazolate
697 framework-8 (ZIF-8) production for CO₂ capture. *J. Clean. Prod.* **370**, 133354 (2022).
- 698 74. Tricker, A. W., Samaras, G., Hebisch, K. L., Realff, M. J. & Sievers, C. Hot spot generation,
699 reactivity, and decay in mechanochemical reactors. *Chem. Eng. J.* **382**, 122954 (2020).
- 700 75. Torad, N. L. *et al.* Nanoarchitected porous carbons derived from ZIFs toward highly
701 sensitive and selective QCM sensor for hazardous aromatic vapors. *J. Hazard. Mater.* **405**,
702 124248 (2021).
- 703 76. Zhou, C. *et al.* Epoxy composite coating with excellent anticorrosion and self-healing
704 performances based on multifunctional zeolitic imidazolate framework derived
705 nanocontainers. *Chem. Eng. J.* **385**, (2020).
- 706 77. Zheng, G. *et al.* Shape control in ZIF-8 nanocrystals and metal nanoparticles@ZIF-8
707 heterostructures. *Nanoscale* **9**, 16645–16651 (2017).
- 708 78. Cravillon, J. *et al.* Controlling zeolitic imidazolate framework nano- and microcrystal
709 formation: Insight into crystal growth by time-resolved in situ static light scattering. *Chem.*
710 *Mater.* **23**, 2130–2141 (2011).
- 711 79. Tuffnell, J. M. *et al.* Comparison of the ionic conductivity properties of microporous and
712 mesoporous MOFs infiltrated with a Na-ion containing IL mixture. *Dalt. Trans.* **49**, 15914–
713 15924 (2020).
- 714 80. Huang, S. & Cole, J. M. BatteryDataExtractor: battery-aware text-mining software embedded
715 with BERT models. *Chem. Sci.* **13**, (2022).
- 716 81. Jablonka, K. M., Schwaller, P. & Ortega-guerrero, A. Is GPT-3 all you need for low-data
717 discovery in chemistry? *ChemRxiv* 1–32 (2023).
- 718