# Outbreak.info Research Library: A standardized, searchable platform to discover and explore COVID-19 resources

Laura Hughes ( ✉ lhughes@scripps.edu )

Scripps Research    https://orcid.org/0000-0003-1718-6676

Ginger Tsueng

Scripps Research Institute

Julia Mullen

The Scripps Research Institute

Manar Alkuzweny

University of Notre Dame    https://orcid.org/0000-0002-6069-5778

Marco Cano

The Scripps Research Institute

Benjamin Rush

Ocuvera

Emily Haag

The Scripps Research Institute

Alaa Abdel Latif

Invitae    https://orcid.org/0000-0002-3713-8420

Xinghua Zhou

Scripps Research    https://orcid.org/0000-0002-9119-3906

Zhongchao Qian

Scripps Research    https://orcid.org/0000-0001-8334-9467

Emory Hufbauer

The Scripps Research Institute

Mark Zeller

The Scripps Research Institute

Kristian Andersen

Scripps Research

Chunlei Wu

The Scripps Research Institute    https://orcid.org/0000-0002-2629-6124

Andrew Su

The Scripps Research Institute

Karthik Gangavarapu

1  **Outbreak.info Research Library: A standardized, searchable platform to discover and**
2  **explore COVID-19 resources**

3  Ginger Tsueng[1]*, Julia L. Mullen[1], Manar Alkuzweny[2], Marco Cano[1], Benjamin Rush[3], Emily
4  Haag[1], Outbreak Curators, Alaa Abdel Latif[4], Xinghua Zhou[1], Zhongchao Qian[1], Emory
5  Hufbauer[4], Mark Zeller[4], Kristian G. Andersen[4,5], Chunlei Wu[1,5,6], Andrew I. Su[1,5,6], Karthik
6  Gangavarapu[7], Laura D. Hughes[1]

7

8  [1]Department of Integrative, Structural and Computational Biology, The Scripps Research
9  Institute, La Jolla, CA 92037, USA
10 [2] Department of Biological Sciences, University of Notre Dame, Notre Dame, IN 46556, USA
11 [3] Ocuvera, Lincoln, NE 68512, USA
12 [4] Department of Immunology and Microbiology, The Scripps Research Institute, La Jolla, CA
13 92037, USA
14 [5] Scripps Research Translational Institute, La Jolla, CA 92037, USA
15 [6] Department of Molecular Medicine, The Scripps Research Institute, La Jolla, CA 92037, USA
16 [7] Department of Human Genetics, David Geffen School of Medicine, University of California
17 Los Angeles, Los Angeles, CA 90095, USA
18 *Corresponding authors: gtsueng@scripps.edu, lhughes@scripps.edu

19  **Abstract**

20  To combat the ongoing COVID-19 pandemic, scientists have been conducting research at
21  breakneck speeds, producing over 52,000 peer-reviewed articles within the first year. To
22  address the challenge in tracking the vast amount of new research located in separate
23  repositories, we developed outbreak.info Research Library, a standardized, searchable
24  interface of COVID-19 and SARS-CoV-2 resources. Unifying metadata from fourteen
25  repositories, we assembled a collection of over 270,000 publications, clinical trials, datasets,
26  protocols, and other resources as of May 2022. We used a rigorous schema to enforce
27  consistency across different sources and resource types and linked related resources.
28  Researchers can quickly search the latest research across data repositories, regardless of
29  resource type or repository location, via a search interface, public API, and R package. Finally,
30  we discuss the challenges inherent in combining metadata from scattered and
31  heterogeneous resources and provide recommendations to streamline this process to aid
32  scientific research.

## Introduction

In early January 2020, SARS-CoV-2 was identified as the virus responsible for a series of pneumonia cases with unknown origin in Wuhan, China[1]. As the virus quickly spread all over the world, the global scientific community began to study the new virus and disease, resulting in the rapid release of research outputs (such as publications, clinical trials, datasets) and resources (i.e. research outputs, websites, portals and more). The frequently uncoordinated generation and curation of resources by different types of resource generators (such as government agencies, NGOs, research institutes, etc.) exacerbate four factors that make finding and using resources a challenge: volume, fragmentation, variety, and standardization (**Figure 1**). These four factors hamper the ability for researchers to discover these resources, and consequently, impede the translation of these protocols, data, and insights into a synthesized understanding of the virus to help combat the pandemic.

For example, the volume of peer-reviewed articles from a single resource (LitCovid) has grown from about 52,000 published within the first twelve months to over 250,000 as of June 2022.  Since April 2020, over 1,000 different research outputs have been published on a weekly basis, spanning new protocols, datasets, clinical trials, as well as publications. The rapid proliferation of resources could be manageable if there were a centralized repository for finding them, but none exists. In addition to research outputs like scientific literature, researchers, public health officials, media outlets, and concerned communities independently developed websites providing highly localized or specialized information on infection rates[2,3], prevention policies[4,5,6], and travel restrictions[7] resulting in a fragmented landscape of very different types of resources (**Figure 1**).

The volume and fragmentation issues were immediately obvious. Lacking alternate solutions for addressing these issues, individual and community efforts for curating these resources were created via shared Google spreadsheets[8,9,10] to aid in discoverability. However, the sheets were not a scalable solution and usually lacked sufficient metadata for describing resources, with the exception of Navarro and Capdarest-Arest. Several projects have attempted to address the volume and fragmentation issues, but were most often focused on a single type of resource. For example, NIH's iSearch COVID-19 portfolio[11] and the Kaggle COVID-19 Open Research Dataset Challenge (CORD-19)[12] aggregate scholarly articles, but do not include clinical trials, datasets, or other types of resources.

Compounding search issues caused by the variety of resource types, there has been a long-standing lack of standardization even *within* a particular type of resource. Existing resource repositories which were able to pivot quickly and curate COVID-19 content from their
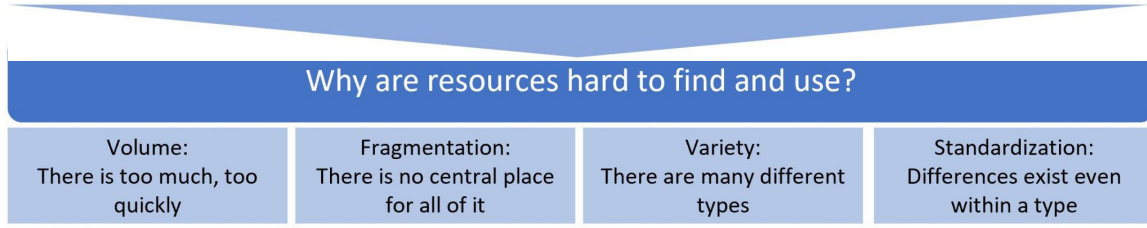
70   collections utilized pre-existing metadata standards. For example, researchers involved in
71   PubMed, which uses Medline citation standards, shifted quickly to create LitCovid[13] which
72   follows the same standard. Similarly, the National Clinical Trials Registry has their own
73   custom list of COVID-19 Clinical Trials which follows their own Protocol Registration and
74   Results System (PRS) schema[14], but these conventions are not followed by the WHO
75   International Clinical Trials Registry Platform. Zenodo[15] and Figshare[16] , which both enable
76   export to multiple open data formats including schema.org, do not completely agree on the
77   marginality, cardinality, and selection of the properties in profiles they use[17,18,19].
78

## What are resources?

- Websites, dashboards, portals
- Research outputs (Publications, Datasets, etc.)
- Lists and collections of other resources
- Guides, policies, recommend-dations
- Tools and code
- Other (learning materials, videos, etc.)

79

## Who contributes to resource proliferation during an outbreak?

| | Existing resources | Medical / research institute / association | Governmental Agencies / Trade associations / NGOs | News and Media Outlets | Citizen Scientists | Search engines and others |
|---|---|---|---|---|---|---|
| **Examples of resource generators** | PubMed, Figshare | JHU, Boston Children's Hospital | CDC, IATA, WHO, ACAPS, Dept. of Public Health | The New York Times, The Atlantic | Phildini, COVID-19 LST | google |
| **What kind of resources do they generate?** | Subsets of existing resources | Research Outputs, Websites, portals, tools | Research outputs, portals, lists, guides, policies, and more | Dashboards, Datasets and other Research outputs, websites, apps | Lists, collections, tools, websites, portals, dashboards | Dashboard-like, specialized search results |
| **What are some examples of resources they generate?** | LitCovid, Covid Figshare | JHU COVID cases and deaths data, VaccineFinder | WHO Reports, Local county dashboards, CDC, State Declarations | The Atlantic's Covid Tracker, New York Time's data | stayinghomeclub, Daily COVID-19 LST Reports | Google's covid19 search including maps with testing sites, etc. |
| **What is the volume and fragmentation of the resources they generate?** | Volume and fragmentation depend on resource | High collective volume but very fragmented | High collective volume, fragmentation dependent on data infrastructure | Varies—depends on syndication and shared infrastructure | Volume unknown due to high level of fragmentation | Low--single resources (search engines) aggregating metadata from other resources |
| **What is the level of standardization for the resource data or metadata?** | High within a resource, Low across different resources | Varies– can be high within a resource type or conform to repository requirements. Low across types /repositories | Varies depending on existing policies and infrastructure—lack of guidance/policy: less standardization | High within organization, low between organizations | Low across individual efforts, maybe higher in collective efforts | Prefers high–specialized views require standardized resources |

80
81

## Why are resources hard to find and use?

| Volume: There is too much, too quickly | Fragmentation: There is no central place for all of it | Variety: There are many different types | Standardization: Differences exist even within a type |
|---|---|---|---|

## How can we support the use of standardized, centralized resources?

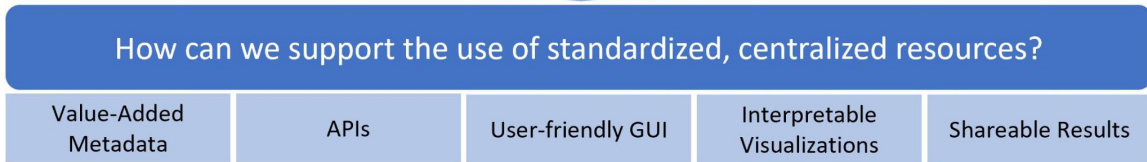| Value-Added Metadata | APIs | User-friendly GUI | Interpretable Visualizations | Shareable Results |
|---|---|---|---|---|

82
83 **Figure 1.** What are resources, who contributes to the proliferation of resources, why are resources
84 difficult to find and use, and how can we support their use?

4

85 Once the issues of volume, fragmentation, variety, and standardization of resources are
86 addressed, accessibility of the resulting resources for reuse must be addressed.
87 Standardized, centralized resources are of no value if researchers are not able to leverage
88 them. Researchers seeking to process information *en masse* will need an API, while
89 researchers seeking to browse and explore will prefer a user-friendly interface. APIs
90 themselves are less useful without a means of understanding the underlying metadata/data
91 (such as documentation or a GUI), and a user-friendly search portal will be less useful without
92 the inclusion of value-added metadata (such as ones supporting search/filter, linkage and
93 exploration, or qualitative evaluation) for improving resource discovery and interpretation.
94 Interpretability of metadata/data is influenced by the order in which information is
95 presented. To address this challenge, the user interface must encourage exploration which
96 gives users control over the information flow to suit their needs. Lastly, if a user has been
97 able to successfully leverage the standardized, centralized resources, they should be able to
98 easily save and share the results of their efforts.
99
100 We address the aforementioned challenges inherent in combining metadata from disparate
101 and heterogeneous resources and making information more interpretable by building
102 outbreak.info, a website which integrates a searchable interface for a diverse,
103 heterogeneous resources which we have collected and standardized (metadata) with
104 surveillance reports on SARS-CoV-2 variants and mutants (data). Following implementation
105 considerations for FAIRness[20], our website includes programmatic access via APIs and a
106 standardized metadata interface built off schema.org. Daily updates ensure that site users
107 have up-to-date information, essential in the midst of a constantly changing research
108 landscape. Based on our experience unifying metadata across repositories, we will discuss
109 issues with centralizing, standardizing, and returning resource metadata, epidemiological
110 data, and supporting the use of the metadata/data. In a companion piece, we present our
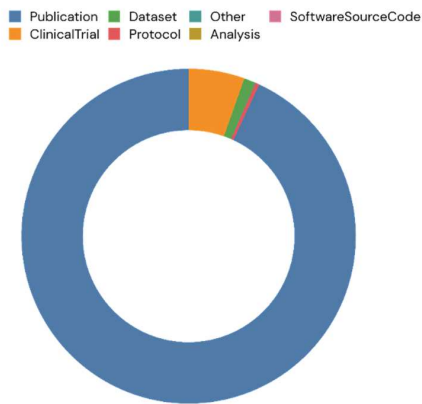111 efforts to develop genomic reports to scalably and dynamically track SARS-CoV-2 variants[21].

112 **Results**

113 **Standardizing metadata through a schema harmonizing a variety of resource types.**
114 We address issues with metadata variety, standardization, and fragmentation by developing
115 a harmonized schema. Schema.org provides a framework to standardize metadata for many
116 different types of data found on the world wide web. However, these standards are not
117 preserved across different types of data. For example, publication providers like PubMed
118 typically use the 'author' property in their metadata, while dataset providers like Figshare
119 and Zenodo are compliant with the DataCite schema and typically prefer 'creator'. Although
120 both properties are valid for their respective schema.org classes, we normalized our schema

5

121 to use 'author' for all 5 of our classes since we expected the volume of publications to dwarf
122 all other classes of resources. We developed a schema that encompassed five types of
123 resources based on their proliferation at the beginning of the pandemic and their
124 importance to the research community: Publications, Datasets, Clinical Trials, Analysis, and
125 Protocols. We added this schema to the Schema Registry of the Data Discovery Engine
126 (DDE)[22], a project to share and reuse schemas and register datasets according to a particular
127 schema. Using this schema we ingested and harmonized metadata from an initial set of
128 fourteen key resources: LitCovid (Publications), bioRxiv and medRxiv (Publications), COVID-
129 19 Literature Surveillance Team (COVID-19 LST) (Publications), ClinicalTrials.gov (NCT)
130 (ClinicalTrials), WHO International Clinical Trials Registry Platform (WHO ICTRP)
131 (ClinicalTrials), Figshare (Datasets, Publications, and more), Zenodo (Datasets, Publications,
132 and more), MRC Centre for Global Infectious Disease Analysis (Analyses, Publications, and
133 More), Protocols.io (Protocols), Protein Data Bank (PDB) (Datasets), Data Discovery Engine
134 (Datasets), Harvard Dataverse (Datasets), and ImmPort (Datasets) (**Figure 2a**).

**Figure 2.** Supporting resource centralization and standardization by developing a harmonizing schema. **a,** Distribution of resources by resource type and source. Note that the x-axis in the bar graphs have different scales. **b,** Heterogeneous and filterable resources (i.e. publications, clinical trials, datasets, etc.) resulting from a single search of the phrase "Delta Variant"**.**

Sources of certain metadata did not map readily to existing schema.org classes. For example, clinical trials registries like NCT have one general schema for both observational and interventional studies, while schema.org provides separate classes for each of these types of studies. Since NCT was a primary source of clinical trials metadata for our research library, we tailored the Outbreak schema based on the combined general NCT schema. Fortunately, many dataset repositories offered schema.org-compliant metadata, even if the repositories differed in the metadata fields that were available. Dataset metadata was harvested from Zenodo, Figshare, PDB, and Harvard Datasets, while Protocols were imported from

7

148   Protocols.io and NCT Protocols. Once our schema was developed, we created parsers (data
149   plugins) to import and standardize metadata from our initial set of resources. We assembled
150   the data plugins into a single API via BioThings SDK[23], and scheduled them to update on a
151   daily basis to ensure up-to-date information. By leveraging the BioThings SDK, we developed
152   a technology stack that addresses the fragmentation issue by easily integrating metadata
153   from different pre-existing resources. With a unified schema that harmonizes information
154   across heterogeneous resource types, a single search (for example "delta variant") to our API
155   can return relevant publications, datasets, clinical trials, and more (**Figure 2b**).
156
157
158   **Enabling community curation and metadata submission to address fragmentation**
159   **and standardization issues.**
160   At the start of the pandemic many curation efforts were neither coordinated, standardized,
161   or easy to find; however, these efforts served an important role in organizing information
162   early on. Given the highly-fragmented, diffuse and frequently changing nature inherent to
163   biomedical resources, we built outbreak.info with the idea that it should be expanded with
164   the participation of the community. Not only is finding and adding resources to the collection
165   an onerous process, it also requires us to know the full landscape of resources on the
166   internet. Furthermore, many resources do not collect metadata useful for linkage,
167   exploration, and evaluation in machine-readable formats. We enabled community-based
168   contributions of resource metadata in a variety of ways (**Figure 3a**).
169

**Figure 3.** Aggregating resource metadata by leveraging community contributions. **a,** The community contribution pipeline and technology stack for outbreak.info's Research Library. Curators may submit dataset metadata using the DDE built-in guide or from GitHub via the DDE/BioThings SDK. Python-savvy contributors can create parsers to contribute even more metadata via the BioThings SDK plugin architecture. A resource plugin allows the site to automatically ingest and update metadata from the corresponding external resource. Blue arrows indicate manual steps, yellow arrows indicate automatable steps after an initial set up, green arrows indicate completely automated steps. **b,** An example of a detailed metadata record manually-curated by volunteers as it appears in the Research Library.

180  For single datasets, contributors can submit the metadata via outbreak.info's dataset
181  submission guide on the Data Discovery Engine, which ensures that the curated metadata
182  conforms to our schema. From there, it can be saved to GitHub, where it can be improved
183  by other contributors via forking and pull requests. The DDE automatically passes the
184  information to the outbreak.info Resources API where it is made discoverable with the
185  Research Library. We demonstrated its utility by asking two volunteers to annotate metadata
186  from thirty different individual resources from across the internet and submitted the
187  metadata for integration via the DDE. As seen in **Figure 3b**, community-contributed
188  metadata using the DDE is standardized and can be exhaustively detailed. Although both of
189  our volunteers provided values for many of the available metadata properties (name,
190  description, topicCategories, keywords, etc.), one provided an extensive list of authors. Using
191  the BioThings SDK in conjunction with the DDE allows us to centralize and leverage
192  individualized curation efforts that often occur at the start of a pandemic. Additionally,
193  collections of standardized datasets, publications and other resources can be submitted to
194  the Outbreak Resources API by contributing a resource plugin. Resource plugins are
195  BioThings-compatible Python scripts to harvest metadata from a source and standardize it
196  to our schema; these parsers can be submitted by anyone with Python coding skills[23]. Our
197  community contribution pipeline allows us to quickly and flexibly integrate the
198  uncoordinated data curation efforts, particularly apparent at the start of the pandemic
199  (**Supplemental Figure 1**).
200
201  **Improving searching, linkage and evaluation of resources to support exploration**
202  Centralizing and standardizing the resources does not automatically make the resources
203  explorable to a user. While centralizing and standardizing allows for search, aggregation and
204  some filtering; additional metadata and a user-friendly interface is needed to allow thematic
205  browsing/filtering and to enable iterative traversal from query to search result to refined
206  query and vice versa. To support resource exploration and interpretation, we added
207  properties (value-added metadata) to every class in our schema that would support
208  searching/filtering/browsing (topicCategories), linkage/exploration (correction, citedBy,
209  isBasedOn, isRelatedTo), and interpretation (qualitative evaluations) of resources.
210
211  We selected these properties based on pre-existing citizen science and resource curation
212  activities, suggesting their value in promoting discoverability. For example, citizen scientists
213  categorized resources in their lists/collections by type (Dataset, Clinical Trials, etc.) in their
214  outputs[10] or area of research (Epidemiological, Prevention, etc.)[24] as they found these
215  classifications helpful for searching, filtering, and browsing their lists/collections. They also
216  evaluated the level of evidence provided by these resources in order to improve its

217 interpretability (i.e. understanding the credibility/quality of the resource)[24]. Existing
218 repositories such as LitCovid also organized information to enhance browsability, but these
219 efforts were often not captured in the metadata. For instance, LitCovid organized
220 publications into eight research areas such as Treatments or Prevention, but these
221 classifications are not available in the actual metadata records for each publication. To
222 obtain these classifications from LitCovid, subsetted exports of identifiers must be
223 downloaded from LitCovid and then mapped to the metadata records from PubMed.
224
225 To classify resources by topicCategory and improve search/browse/filtering capabilities in
226 our user interface, we used a combination of existing work (LitCovid) and human curation to
227 augment that categorization to provide higher specificity of topics and to extend to new
228 types of data (datasets, clinical trials). We applied out-the-box logistic regression,
229 multinomial naive bayes, and random forest algorithms to create models for classifying each
230 resource as belonging or not belonging to each topic. These three algorithms were found to
231 perform best on this binary classification task using out-the-box tests. For example, if a user
232 wants to browse for all resources (or filter down search results) related to the prevention of
233 COVID-19, they can select the appropriate topicCategory in the search/search results view of
234 the resources (**Figure 4a**). Users can also easily traverse from a view of a resource record to
235 start a new search by clicking on a topicCategory of interest (**Figure 4b**). We further enable
236 exploration by populating the linkage properties (corrections, citedBy, isBasedOn,
237 isRelatedTo) from citation metadata (whenever possible), corrections metadata (from
238 LitCovid, when available), and via an algorithm for matching peer-reviewed papers in LitCovid
239 with their corresponding preprints in bioRxiv/medRxiv. Together with the corrections
240 metadata from LitCovid, the algorithm has matched over 2,600 peer-reviewed articles with
241 their corresponding preprints, enabling users to follow from a publication record from
242 LitCovid to a publication record in bioRxiv/MedRxiv (**Figure 4b**).
243
244 Once a user has found a record of interest, they might wonder about the credibility of the
245 resource. To populate resource evaluations so that users can assess the quality of a resource
246 and tailor their interpretation accordingly, we leveraged the Oxford 2011 Levels of Evidence
247 annotations generated by the COVID-19 Literature Surveillance (COVID-19 LST) team[24] as well
248 as Digital Science's Altmetrics[25]. These evaluations are currently visible in the search results,
249 and in the future, we will enable users to further filter or sort search results by some
250 measurement of quality (i.e. Altmetrics: degree of access, or COVID-19 LST: level of evidence).
251 Lastly, we integrated resources with data and analyses we curated to track SARS-CoV-2
252 variants[21]. Researchers can seamlessly traverse from a specific variant report like Omicron
253 to resources on that variant to help understand its behavior (**Figure 4c**). In the absence of a

centralized search interface with linked records, a similar attempt to explore resources outside the outbreak.info portal would require extensive manual searching from multiple different sites (**Supplemental Figure 2**), each with their own interfaces and corresponding search capabilities.



**Figure 4.** Enabling exploration of the resources. **a,** Selectable options for filtering results by topic category or other facets enhance searchability and exploration from the search results view. **b,** Links to other records or to additional potential searches of interest enabling further exploration from a record view. **c,** Links from the Omicron Variant report to related resources.

## Discussion

Over the course of the COVID-19 outbreak, researchers have shared the results of their work at unprecedented levels – exacerbating existing issues in resource volume, fragmentation, variety, and standardization. These issues make it challenging to assemble, traverse, and maintain up-to-date resources. Further, the urgency of a pandemic requires that these issues be addressed quickly, and in a scalable manner to be able to accommodate more data flexibly. We launched outbreak.info within 2 months of the start of the COVID-19 pandemic to address these issues and to highlight barriers in rapidly sharing research outputs in the midst of a pandemic.

272

273 To address the structure and standardization issue, we developed a standardized schema,
274 integrated metadata from different resources into an accessible API, and created a user-
275 friendly search-and-filter, web-based interface. In addition to difficulties standardizing
276 inconsistent metadata models between resources, it is also challenging to maintain a
277 resource library that imports metadata from so many sources, particularly when the
278 metadata updates daily and is prone to change structure. Any changes to the upstream
279 metadata offered by an external site necessitates a change in the parser which imports
280 them. The resource API utilizes the BioThings SDK plugin architecture to handle errors in
281 individual parsers without affecting the availability of the API itself. Using the plugin
282 architecture also allows the creation and maintenance of the individual resource parsers to
283 be crowdsourced to anyone with basic Python knowledge and a GitHub account. Although
284 resource plugins allow outbreak.info to ingest large amounts of standardized metadata,
285 there are still many individual datasets and research outputs scattered throughout the web
286 which are not located in large repositories. Since it is not feasible for one team to locate,
287 identify, and collect standardized metadata from these individual datasets and research
288 outputs, we leveraged the Data Discovery Engine to enable crowdsourcing and citizen
289 science participation in the curation of individual resource metadata.

290

291 At the onset of our data harvesting and harmonization efforts, we focused on creating a
292 unified search interface backed by a common schema.org-based schema. With an
293 extendable pipeline in place, we focused next on augmenting the existing metadata by
294 adding properties to help researchers find information more quickly: topic categorization to
295 group related research, resource linking to connect related entities along the data lifecycle
296 from data generation through publication, and integrating external evaluations of the
297 research trustworthiness using a combination of human curation and automated methods.

298

299 Citizen scientists have played an active role in data collection[26,27] and making information
300 more accessible[12,24] throughout the current pandemic. Given their ability to perform
301 information extraction[28] and their immense contributions to classification tasks[29], we
302 incorporated citizen science contributions into the training data for classifying resources into
303 our topic categories. Some resource aggregators have used clustering algorithms to
304 categorize the entries in their resource libraries, though many only aggregate resources of a
305 single type (i.e. publications). We employed a different approach due to the heterogeneity of
306 our resources, but our API is openly accessible, so anyone is welcome to apply clustering
307 approaches to classify the entries.

308

309    In addition to generating metadata values for improved searching and filtering, we enabled
310    linkages between resources in our schema. For instance, ideally a publication about a clinical
311    trial would link to its clinical trial record, protocols used to collect the data, datasets used in
312    their analyses, and software code underlying the analyses to enable a more meaningful
313    understanding of this trial. However, these connections rarely exist within the metadata; as
314    a result, we have generated linkages between preprints and peer-reviewed publications, and
315    plan to create more linkages between other resource types.  Challenges to include these
316    linkages included: the lack of unique identifiers, inconsistent use of citation metadata fields
317    between resources, and the lack of structured linkage metadata. For example, the ONS
318    Deaths Analysis does not have a unique identifier as assigned by Imperial College London,
319    lacks any citation metadata fields, and instead mentions a potential linkage to an Imperial
320    College London report in its mention of limitations[30]. Although preprints from bioRxiv[31] and
321    medRxiv may have links to the corresponding peer-reviewed manuscript on the bioRxiv site,
322    this information is not accessible via their API, necessitating the use of algorithms to
323    generate these links.
324
325    As a result of this centralization, standardization, and linkage, the outbreak.info Research
326    Library and resources API has been widely used by the external community, including
327    journalists, members of the medical and public health communities, students, and
328    biomedical researchers[32]. For instance, the Radx-Rad Data Coordination Center
329    (https://www.radxrad.org) is utilizing the Outbreak API to collect articles for customized
330    research digests for its partners. Using the Radx-Rad SearchOutbreak app
331    (https://searchoutbreak.netlify.app), users select topics based on information submitted
332    from partners. These are turned into queries for the Outbreak API, and every week, new
333    articles are added to the digests which are available at the website. A workflow sends an
334    email to subscribed users. These digests are not currently available to the public but are
335    expected to be released publicly in the future[33]. Overall, the site receives over a thousand
336    hits per day on average and its visualizations are shared frequently across social media
337    platforms like Twitter.
338
339    While we have developed a framework for addressing resource volume, fragmentation, and
340    variety that can be applicable to future pandemics, our efforts during this framework
341    exposed additional limitations in how data and metadata are currently collected and shared.
342    Researchers have embraced pre-publications, but resources (especially datasets and
343    computational tools) needed to replicate and extend research results are not linked in ways
344    that are discoverable. Although many journals and funders have embraced dataset and
345    source code submission requirements, the result is that the publication of datasets and

14

346  software code are still heavily based in publications instead of in community repositories
347  with well-described metadata to promote discoverability and reuse. In the outbreak.info
348  Research Library, the largest research output by far is publications, while dataset submission
349  lags in standardized repositories encouraged by the NIH such as ImmPort, Figshare, and
350  Zenodo. We hypothesize that this disparity between pre-print and data sharing reflects the
351  existing incentive structure, where researchers are rewarded for writing papers and less for
352  providing good, reusable datasets. Ongoing efforts to improve metadata standardization
353  and encourage schema adoption (such as the efforts in the Bioschemas community) will help
354  make resources more discoverable in the future – provided researchers adopt and use them.
355  For this uptake to happen, fundamental changes in the incentive structure for sharing
356  research outputs may be necessary.
357
358  Within the eighteen months since SARS-CoV-2 was first identified as the infectious agent of
359  the COVID-19 pandemic, there have been over 170 million cases and nearly 4 million deaths.
360  As those numbers continue to grow, so too does the research and understanding of the
361  causes and consequences of the spread of this virus. Given that there will be other
362  pandemics in the future, we demonstrate how we built and launched an extendable and
363  searchable platform for exploring COVID-19 research outputs and genomics data within two
364  months of the pandemic. We address many of the challenges faced when assembling a
365  collection of heterogeneous research outputs and data into a searchable platform. Our
366  platform, outbreak.info, seeks to make COVID-19 data more findable, accessible,
367  interoperable, reusable and interpretable by addressing many data management issues
368  exposed by an urgent and frequently-changing situation. Our site is used by a wide variety
369  of professionals including journalists, members of the medical and public health
370  communities, students, and biomedical researchers[32]. On average, the site receives over a
371  thousand hits per day and its visualizations are shared frequently across social media
372  platforms like Twitter. This platform is also easily extensible to add new metadata sources,
373  allowing the Research Library to grow with the pandemic as research changes.
374

## Methods

### Schema development

377  The development of the schema for standardizing our collection of resources is as previously
378  described[22]. Briefly, we prioritized five classes of resources which had seen a rapid expansion
379  at the start of the pandemic due to their importance to the research community:
380  Publications, Datasets, Clinical Trials, Analysis, and Protocols. We identified the most closely
381  related classes from schema.org and mapped their properties to available metadata from 2-

382    5 of the most prolific sources. Additionally, we identified subclasses which were needed to
383    support our main five classes and standardized the properties within each class. In addition
384    to standardizing ready-to-harvest metadata, we created new properties which would
385    support the linkage, exploration, and evaluation of our resources. Our schema was then
386    refined as we iterated through the available metadata when assembling COVID-19
387    resources.          The          Outbreak          schema          is          available          at
388    https://discovery.biothings.io/view/outbreak.
389
390    **Assembly of COVID-19 resources**
391    The resource metadata pipeline for outbreak.info includes two ways to ingest metadata.
392    First, metadata can be ingested from other resource repositories or collections using the
393    BioThings SDK data plugins. For each resource repository/collection, a parser/data plugin
394    enables automated import and updates from that resource. Second, metadata for individual
395    resources can be submitted via an online form. To assemble the outbreak.info collection of
396    resources, we collected a list of over a hundred separate resources on COVID-19 and SARS-
397    CoV-2. This list ([Supplemental Table 1](#)) included generalist open data repositories,
398    biomedical-specific data projects including those recommended by the NIH[34] and NSF[36] to
399    house open data, and individual websites we came across through search engines and other
400    COVID-19 publications. Prioritizing those resources which had a large number of resources
401    related to COVID-19, we selected an initial set of 2-3 sources per resource type to import into
402    our collection. Given the lack of widespread repositories for Analysis Resources, only one
403    source would be included in our initial import (Imperial College London). An Analysis
404    resource is defined as a frequently-updated, web-based, data visualization, interpretation,
405    and/or analysis resource.
406
407    **Community curation of resource metadata**
408    Resource plugins such as those used in the assembly of COVID-19 resources do not
409    necessarily have to be built by our own team. We used the BioThings SDK[23] and the Data
410    Discovery Engine[22] so that individual resource collections can be added by writing BioThings
411    plugins that conform to our schema. Expanding available classes of resources can be done
412    easily by extending other schema.org classes via the DDE Schema Playground at
413    [https://discovery.biothings.io/schema-playground](https://discovery.biothings.io/schema-playground). Community contributions of resource
414    plugins can be done via GitHub. In addition to contributing resource plugins for
415    collections/repositories of metadata, users can enter metadata for individual resources via
416    the automatic guides created by the Data Discovery Engine. To investigate potential areas of
417    community contribution, we asked two volunteers to inspect 30 individual datasets sprinkled
418    around the web and collect the metadata for these datasets. We compared the results

419  between the two volunteers and their combined results were subsequently submitted into
420  the collection via the Data Discovery Engine's Outbreak Data Portal Guide at
421  https://discovery.biothings.io/guide/outbreak/dataset. Improvements or updates for
422  manually curated metadata can be submitted via GitHub pull requests.
423
424  **Community curation of searching, linkage, and evaluation metadata and scaling with**
425  **machine learning**
426  In an effort to enable improved searching and filtering, we developed a nested list of
427  thematic or topic-based categories based on an initial list developed by LitCovid[13] with input
428  from the infectious disease research community and volunteer curators. The list consists of
429  11 broad categories and 24 specific child categories. Whenever possible, sources with
430  thematic categories were mapped to our list of categories in order to develop a training set
431  for basic binary (in group/out group) classifications of required metadata fields such as (title,
432  abstract and/or description). If an already-curated training set could not be found for a broad
433  category, it would be created via an iterative process involving term/phrase searching on
434  LitCovid, evaluating the specificity of the results, identifying new search terms by keyword
435  frequency, and repeating the process. To generate training data for classifying resources into
436  specific topic categories, the results from several approaches were combined. These
437  approaches include direct mapping from LitCovid research areas, keyword mapping from
438  LitCovid, logical mapping from NCT Clinical Trials metadata, the aforementioned terms
439  search iteration, and citizen science curation of Zenodo and Figshare datasets.
440
441  The efforts of our two volunteers suggested that non-experts were capable of thematically
442  categorizing datasets, so we built a simple interface to allow citizen scientists to thematically
443  classify the datasets that were available in our collection at that point in time. Each dataset
444  was assigned up to 5 topics by at least three different citizen scientists. Citizen scientists were
445  asked to prioritize specific topic categories over broader ones. 90 citizen scientists
446  participated in classifying 500 datasets pulled from Figshare and Zenodo. The citizen science
447  curation site was originally hosted at https://curate.outbreak.info, the code for the site can
448  be found at https://github.com/outbreak-info/outbreak.info-
449  resources/tree/master/citsciclassify and the citizen science classifications at
450  https://github.com/outbreak-
451  info/topic_classifier/blob/main/data/subtopics/curated_training_df.pickle. These
452  classifications have been incorporated into the appropriate datasets in our collection, and
453  have been used to build our models for topic categorization. Basic in-group/out-group
454  classification models were developed for each category using out-the-box logistic regression,

455 multinomial naive bayes, and random forest algorithms available from SciKitLearn. The topic
456 classifier can be found at https://github.com/outbreak-info/topic_classifier.

457

458 In addition to community curation of topic categorizations, we identified a citizen science
459 effort, the COVID-19 Literature Surveillance Team (COVID-19 LST), that was evaluating the
460 quality of COVID-19 related literature. The COVID-19 LST consists of medical students,
461 practitioners and researchers who evaluate publications on COVID-19 based on the Oxford
462 Levels of Evidence criteria and write Bottom Line, Up Front summaries[24]. With their
463 permission, we integrated their outputs (daily reports/summaries, and evaluations) into our
464 collection.

465

466 We further integrated our publications by adding structured linkage metadata, connecting
467 preprints and their peer-reviewed versions. We performed separate Jaccard's similarity
468 calculations on the title/text and authors for preprint vs LitCovid Publications. We identified
469 thresholds with high precision, low sensitivity and binned the matches into (expected match
470 vs needs review). We also leveraged NLM's pilot preprint program to identify and incorporate
471 additional matches. The code used for the preprint-matching can be found at
472 https://github.com/outbreak-info/outbreak_preprint_matcher. Expected matches were
473 linked via the `correction` property in our schema.

474

475

476 **Harmonization and integration of resources and genomics data**
477 The integration of genomics data from GISAID is discussed elsewhere[21]. We built separate
478 API endpoints for our resources (metadata resources API) and genomics (genomics data API)
479 using the BioThings SDK[23]. Data is available via our API at api.outbreak.info and through our
480 R package, as described in Gangavarapu et al.

481 **Acknowledgements**

498 **Conflicts of Interest**
499 KGA has received consulting fees and/or compensated expert testimony on SARS-CoV-2 and
500 the COVID-19 pandemic.

## References

1. World Health Organization. Novel Coronavirus (2019-nCoV): situation report, 1. *World Health Organization*. https://apps.who.int/iris/handle/10665/330760 (2020)

2. Dong, E. *et al*. An interactive web-based dashboard to track COVID-19 in real time. *The Lancet Infectious Diseases*, 20(5), 533-534 https://doi.org/10.1016/S1473-3099(20)30120-1 (2020)

3. Kaiser, J. 'Every day is a new surprise.' Inside the effort to produce the world's most popular coronavirus tracker. *Science* https://doi.org/10.1126/science.abc1085 (2020)

4. Noren, L.E. *et al*. Institutional Response to Covid. https://docs.google.com/spreadsheets/d/1IbF_wlmldVssG5spcmNE82nR9btcbF7rUlEqtcXW03o/edit#gid=0 (2020)

5. Morris, A. & citizen scientists. USA COVID-19 K-12 School Closures, Quarantines, and/or Deaths https://docs.google.com/spreadsheets/d/e/2PACX-1vQSD9mm5HTXhxAiHabZA6BPUByWBlP5HZ2jfOPEeGZkMB0ZFsmFBL5orqjIq22mjFNZ7n-11ObCylGn/pubhtml?fbclid=IwAR2tJ8yDVehGpxoP97Cco5HYAxoN014opwwm6uYt4s3E2xDr_8u9KF_LlgI# (2020)

6. James, P. & citizen scientists. Staying Home Club https://github.com/phildini/stayinghomeclub (2020)

7. Pogkas, D. *et al*. The Airlines Halting Flights as Virus Outbreak Spreads. *Bloomberg* https://www.bloomberg.com/graphics/2020-china-coronavirus-airlines-business-effects/ (2020)

8. Joachimiak, M. *et al*. SARS-COV-2 and COVID-19 datasets https://docs.google.com/spreadsheets/d/1eMhot7MjusyM7_2IBnzqi7RlzWWoYnfheWhMgDlPToQ/edit#gid=0 (2020)

9. Skenderi, J. *et al*. COVID-19 Resource Library. https://docs.google.com/spreadsheets/u/2/d/1cqxDAg4jMHXI6gHOnoV8HqDdRHnmxEJRl-bhhpe1HEo/htmlview# (2020)

10. Navarro, C. & Capdarest-Arest, N. COVID-19 Open Dataset Sources https://docs.google.com/spreadsheets/d/10t3vtULr3nTz7mrlKj0rldUys47wsIfOVReHnx3Xu18/edit#gid=0 (2020)

11. NIH OPA. iSearch COVID-19 Portfolio https://icite.od.nih.gov/covid19/search (2020)

12. Allen Institute For AI. COVID-19 Open Research Dataset Challenge (CORD-19) https://www.kaggle.com/allen-institute-for-ai/CORD-19-research-challenge (2020)

13. Chen, Q. et al. LitCovid: an open database of COVID-19 literature. *Nucleic Acids Research*, 49(D1), pp.D1534-D1540 (2020)

536    14. ClinicalTrials.gov. Protocol Record Schema - XML Schema for electronic transfer of
537    protocol information into the ClinicalTrials.gov Protocol Registration System.
538    https://prsinfo.clinicaltrials.gov/ProtocolRecordSchema.xsd (2018)

539    15. Fava, I. et al. Coronavirus Disease Research Community - COVID-19. *Zenodo.org*.
540    https://zenodo.org/communities/covid-19/?page=1&size=20 (2020)

541    16. Hyndman, A. A Figshare COVID-19 Research Publishing Portal. *COVID19.figshare.com*.
542    https://figshare.com/blog/A_Figshare_COVID-19_Research_Publishing_Portal/558 (2020)

543    17. European Organization for Nuclear Research. Zenodo FAIR Principles. *Zenodo.org*.
544    https://about.zenodo.org/principles/ (2013)

545    18. Hahnel, M. What Google Dataset Search means for academia. *Figshare Blog*.
546    https://figshare.com/blog/What_Google_Dataset_Search_means_for_academia/422 (2018)

547    19. Schema.org. About Schema.org. https://schema.org/docs/about.html (2015)

548    20. Jacobsen, A. *et al*. FAIR Principles: Interpretations and Implementation Considerations.
549    *Data Intelligence* 2(1-2), pp.10-29. https://doi.org/10.1162/dint_r_00024 (2020)

550    21. Gangavarapu, K. *et al*. Outbreak.info: Real-time surveillance of SARS-CoV-2 mutations
551    and variants. *medRxiv*. https://doi.org/10.1101/2022.01.27.22269965 (2022)

552    22. Cano, M. et al. Schema Playground: A tool for authoring, extending, and using metadata
553    schemas to improve FAIRness of biomedical data. *bioRxiv*.
554    https://doi.org/10.1101/2021.09.02.458726 (2021)

555    23. Lelong, S. *et al*. BioThings SDK: a toolkit for building high-performance data APIs in
556    biomedical research. *Bioinformatics*. https://doi.org/10.1093/bioinformatics/btac017 (2021)

557    24. Rah, J. *et a*l. COVID-19 Literature Surveillance Team. *Covid19lst.org*.
558    https://www.covid19lst.org/copy-of-about (2020)

559    25. Digital Science. About Us. *Altmetric*. https://www.altmetric.com/about-us/ (Accessed 07
560    March 2022)

561    26. Birkin, L.J. et al. Citizen science in the time of COVID-19. *Thorax* 76(7), 636-637.
562    https://doi.org/10.1136/thoraxjnl-2020-216673 (2021)

563    27. Rohwer Lab at San Diego State University. #swab4corona. CoVID-19 Citizen Science.
564    https://covidsample.org/ (Accessed 17 September 2021)

565    28. Tsueng, G. *et al*. Applying citizen science to gene, drug and disease relationship
566    extraction from biomedical abstracts. *Bioinformatics* 36(4), 1226-1233.
567    10.1093/bioinformatics/btz678 (2020)

568    29. Blickhan, S. *et al*. Transforming research (and public engagement) through citizen
569    science. *Proceedings of the International Astronomical Union*. Cambridge University Press,
570    14(A30), pp. 518–523. https://doi.org/10.1017/S174392131900526X (2018)

571    30. Imperial College COVID-19 Response Team. ONS excess deaths. *Imperial College London*.
572    http://www.imperial.ac.uk/mrc-global-infectious-disease-analysis/covid-19/covid-19-
573    reports/ (Accessed 21 October 2021)

574    31. bioRxiv. COVID-19 SARS-CoV-2 preprints from medRxiv and bioRxiv.
575    https://connect.biorxiv.org/relate/content/181 (2021)

576    32. Haag, E. User Stories Outbreak.info Blog https://blog.outbreak.info/?tag=user_stories
577    (Accessed 06 January 2022)

578    33. Valentine, D. & Radx. SearchOutbreak. Radical Data Coordination Center.
579    https://searchoutbreak.netlify.app (2021)

580    34. BioMedical Informatics Coordinating Committee. Data Sharing Resources. National
581    Institutes of Health. https://www.nlm.nih.gov/NIHbmic/nih_data_sharing_repositories.html
582    (2020)

583    36. National Science Foundation. Open Data at NSF. https://www.nsf.gov/data/ (2013)

# Supplementary Files

This is a list of supplementary files associated with this preprint. Click to download.

- outbreak.inforesourcessupfile2.xlsx
- outbreak.inforesourcessupfile1.pdf