


SOFTWARE

Open Access



Schema Playground: a tool for authoring, extending, and using metadata schemas to improve FAIRness of biomedical data

Marco A. Cano¹, Ginger Tsueng^{1*} , Xinghua Zhou¹, Jiwen Xin¹, Laura D. Hughes¹, Julia L. Mullen¹, Andrew I. Su¹ and Chunlei Wu¹

*Correspondence:
gtsueng@scripps.edu

¹The Scripps Research Institute,
San Diego, USA

Abstract

Background: Biomedical researchers are strongly encouraged to make their research outputs more Findable, Accessible, Interoperable, and Reusable (FAIR). While many biomedical research outputs are more readily accessible through open data efforts, finding relevant outputs remains a significant challenge. Schema.org is a metadata vocabulary standardization project that enables web content creators to make their content more FAIR. Leveraging Schema.org could benefit biomedical research resource providers, but it can be challenging to apply Schema.org standards to biomedical research outputs. We created an online browser-based tool that empowers researchers and repository developers to utilize Schema.org or other biomedical schema projects.

Results: Our browser-based tool includes features which can help address many of the barriers towards Schema.org-compliance such as: The ability to easily browse for relevant Schema.org classes, the ability to extend and customize a class to be more suitable for biomedical research outputs, the ability to create data validation to ensure adherence of a research output to a customized class, and the ability to register a custom class to our schema registry enabling others to search and re-use it. We demonstrate the use of our tool with the creation of the Outbreak.info schema—a large multi-class schema for harmonizing various COVID-19 related resources.

Conclusions: We have created a browser-based tool to empower biomedical research resource providers to leverage Schema.org classes to make their research outputs more FAIR.

Keywords: Metadata, Schema, Standardization, FAIR resources

Background

Funding agencies, international consortia, institutional policies, and publisher requirements have helped promote the adoption of the FAIR (Findability, Accessibility, Interoperability, and Reusability) guiding principles [4, 41] for biomedical research data sharing to varying degrees of success. While it is now standard to make datasets accessible and potentially reusable via deposition of the dataset in a repository, metadata



© The Author(s) 2023. **Open Access** This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons licence, and indicate if changes were made. The images or other third party material in this article are included in the article's Creative Commons licence, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons licence and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this licence, visit <http://creativecommons.org/licenses/by/4.0/>. The Creative Commons Public Domain Dedication waiver (<http://creativecommons.org/publicdomain/zero/1.0/>) applies to the data made available in this article, unless otherwise stated in a credit line to the data.

standardization issues (i.e.—lack of standardization in how datasets are described) continue to make it challenging for researchers to make datasets findable, interoperable, and reusable. To address these issues, domain experts and data stewards have been inspecting the gap between principle and practice [23]; extending [19], adapting [15], and adopting the principles [12]; creating their own metadata standards [6] and data schemas [12, 16, 29]. However a large gap remains between the communities that develop standards and the adoption of these standards by data and resource providers due to issues in communication, education/training, incentives, and the availability of supportive tools [14, 17]. For example, the Dublin Core Metadata Initiative (DCMI) provides a metadata ontology (i.e.—a structured vocabulary for classifying and describing metadata): terms and data elements (Dublin Core Metadata Initiative [9]), two general-use schema classes (i.e.—sets of metadata vocabulary used to describe a conceptual entity): core and qualified, and a thorough guide for utilizing their ontology with their model-based framework for creating schemas: the Dublin Core Application Profile (DCAP) guide [8]. The DCAP guide was intended to empower data providers to mix and match Dublin Core (and other) metadata terms/elements (properties) to create new application profiles (schemas) to suit their needs. While the core (data element) schema has been widely-adopted, the lack of authoring tools to help create more type/concept-specific schemas and the lack of tools for transforming schemas into working formats for consumption and implementation has hampered the adoption and implementation of DCAP [1]. Even after standardization communities successfully introduce standards, their adoption, modification, and implementation are frequently defined by widely used tools or repositories within a specific community [12].

Schema.org is a metadata vocabulary standardization project founded by the major search engine companies such as Google, Microsoft, Yahoo, and Yandex. It is an open source, collaborative initiative that develops metadata standards for improved searchability. While domain-neutral, Schema.org welcomes proposals and discussions of new properties and classes from anyone, including domain-specific ontology or schema development groups, via participation in their W3C group [38]. Members of metadata ontology development communities (including the aforementioned DCMI, as well as LRMI, and other W3C groups) [3, 39] have been involved with, have influenced, and have successfully integrated some of their vocabulary into Schema.org [2, 32]. Schema.org already includes some biomedically relevant classes (i.e.—conceptual entities) like Datasets and Medical Study, and applying Schema.org classes to biomedical research resources would improve interoperability, enabling researchers readily ingest existing resources and to leverage search engine-based solutions (like Google Dataset Search) to find resources of interest. Furthermore, the hierarchical nature of schemas from Schema.org allows for inheritance of vocabulary sets (sets of properties) from parent schemas. Although there have been some efforts to leverage Schema.org to improve findability of scientific research data [20, 29, 31] and many generic repositories (like Figshare and Zenodo) are compliant, Schema.org remains largely underutilized by the biomedical research community. Bioschemas is an open and collaborative effort that has been actively promoting the use of Schema.org in the life sciences by serving as a hub for researchers to create new biomedically relevant classes with the goal of refining and

proposing these classes to Schema.org [11], Profiti et al. [30], and by raising awareness about the usefulness of metadata schemas. The Bioschemas community has also identified the need for easy-to-use tools to help improve public accessibility and participation in the schema development process.

Here, we describe the Data Discovery Engine's (DDE) Schema Playground, a web-based tool that improves the ease of using any registered schema or Schema.org classes. Our tool allows users to easily find and visualize relevant classes from Schema.org, Bioschemas, BioLink [5], and others, extend them; create JSON schema validation rules [22]; and save/share the newly created classes for others to reuse. Our tool expresses schemas in JSON-LD format, improving interoperability of schemas which might normally otherwise be viewed as HTML tables. Our tool also includes a framework for building data registries and creating guides for data submission; however, the implementation and integration of these features on our site is restricted to partner organizations. We introduce the features of this tool, review its value to different types of users, demonstrate its application towards the creation of a new schema for COVID-19-related resources, and discuss its adoption by the Bioschemas metadata standardization community.

Implementation

The Data Discovery Engine's Schema Playground is a browser-based tool built with Vue.js [43], Python/Tornado [35], and the BioThings Software Development Toolkit [25]. Schemas from Schema.org and other consortia/projects are stored and made searchable using MongoDB [27] and Elasticsearch [34]. The code for the Schema Playground can be found at <https://github.com/biothings/discovery-app> and is free to use under the Apache License 2.0. The schema generated by the DDE are exported as JSON-LD files [21], following RDF schema specifications [40] with embedded JSON Schema metadata validation rules [22]. The COVID-19 Outbreak.info resource schemas were developed by comparing metadata properties across multiple type-specific repositories to identify properties in common. For example, metadata from LitCovid/PubMed, BioRxiv/MedRxiv, various journals like JAME, NEJM and others, and the metadata from publications found on Zenodo, Figshare and others were compared in order to identify a suitable schema for COVID-19-related publications. Similarly, protocols from protocols.io and the Bioschemas *LabProtocol* class were compared to develop a schema for COVID-19-related protocols. Once the desired properties and structure for each class of COVID-19-related resource was identified, the schemas were created by extending existing Schema.org classes using the DDE Schema Playground.

Results

The DDE Schema Playground consists of two standard (and fully-accessible) components and two related, custom (limited-access) components (Fig. 1). The standard components improve the ease of use of schemas and classes, while the custom components help communities to reap the benefits of their use. The Schema Editor allows users to import community standard schemas like Schema.org and customize them for biomedical purposes. These extended schemas can then be shared in the Schema Registry, which allows users to view the schemas and reuse them. When used in conjunction with Data

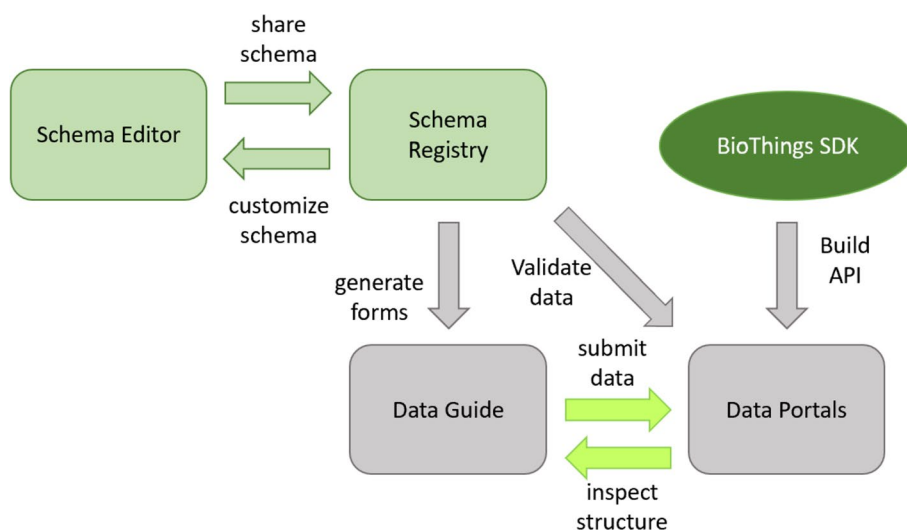


Fig. 1 Components of the DDE Schema Playground and how they work together

		Standard Components		Custom Components	
		Schema Registry	Schema Editor	Data Portals	Data Guides
Types of Users	Data Providers	Find schemas for new classes of data Register in-use schema for FAIRness	Extend schemas to ensure maximum compatibility and minimize duplicate effort	Collaboratively collect and share data	Submit compliant data via forms
	Data Consumers	Download schemas to understand how to consume and validate data	Easily create schema validation for structured data if schema validation not available	Explore and reuse standardized data	Understand structure of available data
	Data Standardizers	Explore community needs for new classes & properties Share schemas for community uptake	Easily create schema from existing classes and validation once consensus on structure has been established	Identify data structuring needs and potential structuring improvements	Understand Structure Burden

Fig. 2 The value and utility of DDE Schema Playground components to different types of users

Portals built with BioThings SDK [25], The DDE Schema Playground can automatically generate data submission forms known as Data Guides.

To understand how the Schema Playground might help to bridge the gap between data standardization communities and data resource providers, we identified potential utility and value of each of the DDE Schema Playground components for different types of users in our partner communities (Fig. 2).

Any data portals and guides can be used by anyone with sufficient access rights, but the creation of a data portal or data guide requires partnerships with our team to actualize. For the outbreak portal, data submission via the guide is open to all and utilizes GitHub for authentication. For other portals, access may be restricted as required by the responsible partner organization. The data portal and data guides allow data providers and data consumers to collect, share, and use data. Since the data guide converts a

custom schema into a web-based data submission form, it enables data consumers and data standardizers to visually inspect and understand the burden of structure.

The schema registry and editor allows data providers and/or standardization communities to find, customize, and share schemas. Sharing schemas via the registry will make it easier for data consumers to understand how to consume data from a data provider and to create data validation if one is not available from the data provider. For example, data-use restrictions usually require a data consumer to create an account with the data provider in order to access data. However, data consumers cannot easily determine whether the data locked behind the account-creation process will actually be useful prior to creating an account. Sharing the data schema via the DDE registry could address this issue by allowing data consumers to understand what's available without actually displaying any restricted-access data. Having a central location for schemas submitted by data providers will also make it easier for data standardization communities to evaluate the needs of the biomedical research community. To our knowledge, the DDE's schema registry is the only crowd-sourced registry for type/concept-specific schemas created specifically for the biological and biomedical research space. To further illustrate the value of the schema registry and editor, we compare and detail the features of the DDE Schema Playground with available tools for creating, applying, and consuming other major schemas such as Schema.org and Bioschemas.

Schema.org, Bioschemas and other data standardization efforts have built strong communities to generate consensus on data modeling for the creation of new schemas or the improvement of existing schemas. Hence, there are extensive processes in place (but few tools) for the creation of a new schema based on Schema.org or any other schemas. Because of its widespread adoption, there are third party tools available for utilizing and consuming markup from Schema.org. The Bioschemas community has developed a process for defining new classes and has a set of tools which cover both the creation of a new schema (google spreadsheet conversion), utilization of a schema (markup generation), and evaluation of use (markup validation, scraping), but these tools vary in usability based on the users programming experience. In contrast to Schema.org, Bioschemas also defines cardinality (allowable number of values per property) and marginality (optional vs required value) in its profile schemas as these are important to the life sciences research community. Although the DDE schema playground was developed independently from the Bioschemas community, our interests aligned and we sought to provide complimentary schema tools to facilitate biomedical schema development and adoption. To do this, we identified schema tools and features available directly from the Schema.org and Bioschemas communities. We expanded the list of tools by searching for "Schema.org tools", "schema generation tools", "schema creation tools", "schema editing tools", "schema validation tools", "bioschemas tools" in Google and in bio.tools). Bio.tools yielded two relevant results (biovalidator and ObjTables), while Google yielded multiple tool reviews/lists which generally featured similar sets of tools. Most user-friendly tools were aimed towards the generation, extraction, or validation of schema-compatible markup rather than the development of schemas themselves (Additional file 2: Supplemental Table 1). The Bioschemas community has a few well-documented tools for schema

Table 1 Comparison of Schema.org, Bioschemas, and DDE Schema Playground

	Schema.org	Bioschemas	DDE schema Playground
Search for classes from own community	✓	✓	✓
Search for classes from other communities			✓
Visualize any schema in json			✓
Create a schema	✓*	✓*	✓
Extend a schema			✓
Create schema validation			✓
Export/Save a schema		✗	✓
View example markup using a schema	✓		
Create markup using a schema	✗	✗	✓**
Account for cardinality		✓	✓
Account for marginality		✓	✓
Google doc integration		✗	
github integration			✓
Markup validation	✗	✗	
Markup scraper	✗	✗	

✓ available, ✗ separate tools available, *process in place, **feature exists, but not generally available

development, but many of those tools were only available as source code and required basic programming experience. We focused our efforts on features for which user-friendly tools for schema creation and reuse, resulting in a web-based application that empowers individual data resource providers to utilize and customize existing schemas from Schema.org and other similar efforts. As seen in Table 1, these features include:

1. Searching and viewing schemas from Schema.org and other metadata standardization efforts

The DDE Schema Playground allows for the visualization of JSON-LD-formatted schemas hosted online either on GitHub or elsewhere (Additional file 1: Supplemental Figure 1A). This allows users who are familiar with Schema.org to review their compliant schema in a more human readable format. The DDE Schema Playground also has a searchable registry of classes from Schema.org, BioLink, BioThings, Bioschemas, and others. Users may browse and visualize the schemas for various classes from these sources to identify the classes of most interest to them (Additional file 1: Supplemental Figure 1B). If a community like Bioschemas or consortia like the National COVID Cohort Collaborative (N3C) [13] is interested in making a new schema available for searching and viewing, they can import and register their JSON-LD-formatted schema. The DDE Schema Playground also enables users to compare up to four schemas. For example, there are multiple Dataset schemas available in the registry, and users can compare them to see what properties are unique to each and what properties they share (Additional file 1: Supplemental Figure 1C).

2. Extending and customizing a pre-existing schema for a particular use

The ability to browse and inspect pre-existing schemas makes it easier for a user to customize or extend the schema to suit their own purpose. All the properties from the pre-existing schema will be inherited in the extended schema; however, the user may select properties for which validation is desirable. The user can also create new properties to be included in the extended schema. For example, the Dataset class from Schema.org serves as a potential foundation, but a schema focused on COVID-19-related datasets may need additional fields (e.g., *infectiousAgent*). To tailor the Dataset schema, we find and extend it from the registry (Additional file 1: Supplemental Figure 2A). After we create a name for our schema (the namespace) and the class, we can customize it. We can select to include any property that is available from the schema we are extending (Additional file 1: Supplemental Figure 2B), and we can create new properties (e.g., *infectiousAgent*) that are tailored to our needs (Additional file 1: Supplemental Figure 2C). This feature also serves as an easy way to maintain Bioschemas profiles as users can update a registered profile by extending from it, making the necessary changes, and pushing them back to Bioschemas. Outside the tool, there is only manual writing/editing of JSON-LD, YAML, TTL, SHACL, ShEx, or other file types and running command-line tools for customizing an existing schema in an interoperable format and making it human-friendly viewable online.

3. Creating validation for the schema for data quality enforcement

Marginality (whether a property is required or not) and cardinality (whether a property can have one or multiple values) are two aspects of schema properties that are not expressed well by Schema.org but are desirable to biomedical researchers (Additional file 1: Supplemental Figure 3A and 3C). In the DDE Schema Playground, this is handled via the creation of JSON Schema validation rules, and the DDE's Schema Validation Editor provides a simple drag and drop mechanism to create straightforward validations (Additional file 1: Supplemental Figure 3B). For slightly more complex validations, the user can edit the validation rule before dragging and dropping it into the property of interest. In our example *Dataset* schema, we may want to restrict the values for our new property (*infectiousAgent*) such that they map to and are standardized by an ontology. We edit the example JSON schema validation rules for an ontology to tailor it to the NCBI Taxon ontology (Additional file 1: Supplemental Figure 3D). Schema development working groups often leverage the work done by ontology groups to ensure that the values of a property are standardized. Once these JSON schema validation rules have been created, they can be used to test the validity of JSON-LD-formatted metadata using any of the many third party metadata validation tools and program libraries that are already available. Future features of the DDE will include a built-in metadata validation tool.

4. Exporting and saving a schema generated by the Schema Playground editor

The DDE Schema Playground allows you to export/download your newly created schema locally and it is also integrated with GitHub, allowing users to save to their GitHub repository (Additional file 1: Supplemental Figure 4A-C). The integration with GitHub allows the edits to the schema to be made by multiple parties and provides the schema owner the option of pulling changes to the schema. Additionally,

the schema can be forked and edited/customized allowing for re-use of the schemas which in turn improves findability and reusability of resources which follow the schemas.

5. Registering a newly created schema in the DDE schema registry to facilitate its extension and re-use

Once saved in GitHub, users can review their schema with the schema viewer and add it to the registry to enable others to easily re-use it (Additional file 1: Supplemental Figure 1A). This provides a user-friendly interface for editing, customizing, and re-using schemas for those who prefer not to manually edit text and format in JSON-LD.

The DDE Schema Playground offers any user the ability to reuse and extend existing schemas. This tool is primarily to assist in the authoring of schemas for use in other applications. In addition, we have converted three *Dataset* schemas into "guides", which are web-based forms for annotating resources using schemas authored in the DDE Schema Playground. Annotations created using these guides are stored within a resource registry hosted within the DDE. There are currently three public guides based on the *Dataset* schemas for the Outbreak.info web application [28], the N3C initiative, and the CD2H consortium [7]. While the creation of guides from schemas is not a fully-automated feature that is available to all users, most of the underlying components are reusable, additional guides can be constructed and hosted within the DDE through collaboration. The Bioschemas community has integrated the DDE schema playground as part of its schema creation and update process to improve participation by members who lack the programming expertise needed to participate via their previous pipeline.

Creating the COVID-19 Outbreak schema using the Schema Playground

Schema.org classes are often simultaneously too broad (lacking properties needed) and too narrow (including too many irrelevant properties) for a specific research purpose. For this reason, it becomes necessary to adapt schemas to suit needs of a biomedical research project. Outbreak.info is a project from the Su, Wu, and Andersen labs at Scripps Research to unify COVID-19 and SARS-CoV-2 epidemiology and genomic data, published research, and other resources [10, 37]. The standardization of published research and other resources was accomplished by creating a single, multiclass schema to harmonize the metadata: The COVID-19 Outbreak schema. This schema can be found in the DDE registry at <https://discovery.biothings.io/view/outbreak/> and was built via the DDE Schema Playground with some manual editing (for merging all the classes into a single schema). There are six principal classes in the Outbreak schema (*Analysis*, *Dataset*, *ClinicalTrial*, *ComputationalTool*, *Protocol*, *Publication*) and many subclasses to support the principal classes. As seen in Table 2, the classes in the Outbreak schema were extended from related Schema.org classes (whenever available) and were created based on metadata comparisons from a variety of related sources. By extending from existing schemas, we reuse existing metadata properties when appropriate, and create new properties only when necessary.

Table 2 Classes in the Outbreak schema and how they were created and used

Outbreak schema class	Source Schema.org class	Inspired by	Used in
Analysis	CreativeWork	Outbreak Protocol, Dataset Class, various COVID-19 Analyses	–
ClinicalTrial	MedicalStudy	NCT Clinical Trials PRS, WHO ClinicalTrials	–
ComputationalTool	Software	Bioschemas, NIAID schema, Bio.tools schema	–
DataDownload	DataDownload	NIAID schema	ComputationalTool
Dataset	Dataset	NIAID schema	–
Protocol	HowTo	Bioschemas LabProtocol, protocols.io	–
Publication	MedicalScholarlyArticle	NLM Medline, biorxiv, various journals	–
ArmGroup	Thing	NCT Clinical Trials PRS, WHO ClinicalTrials	ClinicalTrial
CitationObject	CreativeWork	Outbreak primary classes	All
Correction	CreativeWork	CitationObject	Publication
DataDownload	DataDownload	NIAID schema	Dataset
Eligibility	Thing	NCT Clinical Trials PRS, WHO ClinicalTrials	ClinicalTrial
Instrument	Product	Bioschemas LabProtocol, protocols.io	Protocol
Intervention	Thing	NCT Clinical Trials PRS, WHO ClinicalTrials	ClinicalTrial
MonetaryGrant	MonetaryGrant	NIAID schema	All
Organization	Organization	NIAID schema	All
Outcome	Thing	NCT Clinical Trials PRS, WHO ClinicalTrials	ClinicalTrial
Person	Person	NIAID schema	All
Product	Product	Bioschemas LabProtocol, protocols.io	Protocol
StudyEvent	Thing	NCT Clinical Trials PRS, WHO ClinicalTrials	ClinicalTrial
StudyStatus	Thing	NCT Clinical Trials PRS, WHO ClinicalTrials	ClinicalTrial
StudyDesign	Thing	NCT Clinical Trials PRS, WHO ClinicalTrials	ClinicalTrial
StudyLocation	Place	NCT Clinical Trials PRS, WHO ClinicalTrials	ClinicalTrial

For example, the level of detail provided by Protocol Registration System (PRS) schema used by the National Clinical Trial (NCT) registry is more granular than Schema.org's *MedicalStudy* class, but broad enough that it encompasses properties from both child classes of *MedicalStudy* (*MedicalTrial* and *MedicalObservationalStudy*). The child classes of *MedicalStudy* only differ in the property name for the study design (*trialDesign* vs *studyDesign*), and this property is not delineated in PRS. Further, the PRS includes many properties not currently available in any of these Schema.org classes. Adopting the PRS directly was also problematic as we planned to ingest records from other registries like the World Health Organization's Clinical Trial registry (WHOCT), and the PRS was also more granular than WHOCT. For this reason, the Outbreak.info *ClinicalTrial* class was created by using the DDE to extend from Schema.org, leveraging the PRS-WHO crosswalk [42], and creating properties that could help with issues previously identified [26].

In addition to adapting Schema.org classes to normalize record data from multiple sources within a class, Outbreak.info needed to normalize common metadata properties between different classes. The hierarchical nature of Schema.org classes simplified this process, as many derivative classes inherit properties from the *Thing* class. For example, the *Protocol* class in the Outbreak schema was extended from the *HowTo* class in Schema.org and was based on properties identified from available metadata in protocols.io and the *LabProtocol* profile from Bioschemas. Since both the Schema.org classes, *MedicalStudy* and *HowTo*, are derivatives of *Thing*, the Outbreak schema naturally has properties in common across multiple classes and can normalize the metadata across these classes allowing for cleaner query design and improved search functionality. This schema is currently used to harmonize and improve FAIRness of metadata from over 300,000 resource entries in the Outbreak.info research library at <https://outbreak.info/resources>.

Adoption of the Schema Playground into the Bioschemas schema development and maintenance pipeline

Previously, the pipeline for updating a Bioschemas specification involved the use of a google spreadsheet for attaining community consensus, a command-line tool for converting the CSV from the spreadsheet to yaml, cloning the Bioschemas website repository and copying/editing HTML and YAML files, running Jekyll to test the changes, editing example files in the Bioschemas specifications repository, and creating pull requests for the Bioschemas website repository once everything had performed as tested. The level of expertise needed in order to update a specification has been discussed in multiple Bioschemas community calls as a potential barrier to participation. After initial tests during and after Biohackathon 2021, the Bioschemas community has decided to adopt the DDE into its schema development and maintenance pipeline. Manuals for using the DDE to create or update Bioschemas specifications have been developed, and automated scripts using GitHub actions have been developed to more tightly integrate the tool into the pipeline. As seen in Fig. 3, the process for updating

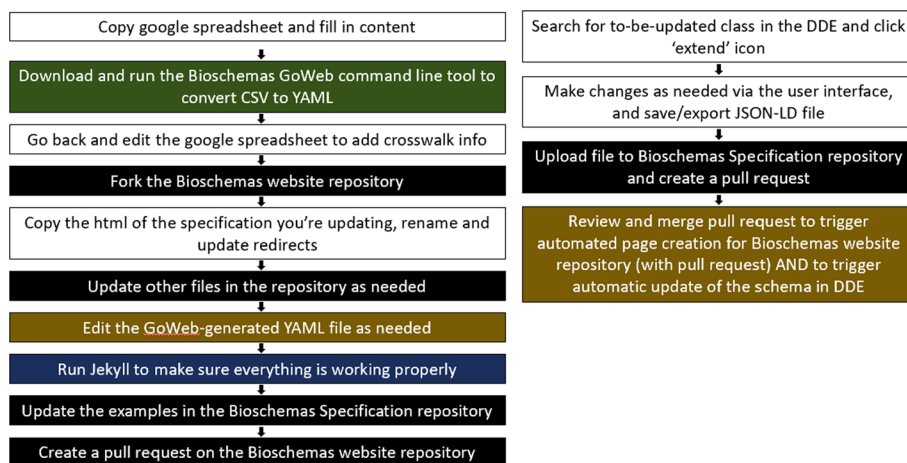


Fig. 3 The Bioschemas profile update process before (left) and after (right) the inclusion of the DDE

a Bioschemas profile requires less technical expertise after the integration of the DDE. While the process prior to and after the DDE still requires the ability to edit a YAML/JSON file (brown) and the ability to use GitHub (black), the DDE-based process does not require the user to have the technical knowledge needed to run tools via the command line (green), or to use Jekyll (blue).

Discussion

In an effort to make scientific resources more FAIR, communities in the biological sciences (Bioschemas), earth sciences (ESIP's Science on Schema.org cluster), and more are working diligently to align and influence Schema.org to suit the needs of the scientific research community [11, 33]. These communities serve as an important bridge between domain-specific ontology development groups and the domain-neutral Schema.org by introducing Schema.org to the scientific research resource providers, identifying existing ontologies to leverage, and creating tailored schemas more suitable for the research community. For example, ontology development groups, like PPEO/MIAPPE [24] and EDAM [18], have been consulted or have participated in the development of the *Sample* and *ComputationalWorkflow* schemas by the corresponding Bioschemas working groups. A term from PPEO/MIAPPE was included as a property in the *Sample* schema, while JSON schema validation rules enforce the use of terms from EDAM as values for certain properties in the *ComputationalWorkflow* schema.

Although communities like Bioschemas have helped to create more relevant classes or improve existing classes, it is difficult to push these suggestions to Schema.org without compelling use cases or widespread adoption of these tailored classes. For example, the Bioschemas community first developed the *Gene* class (with input from gene resource providers, Gene Ontology proponents and gene resource consumers) in 2018. However, it was not included as a pending class in Schema.org until 2021 due to a lack of widespread adoption. The Bioschemas community spent considerable time and effort on education and training in order to increase the adoption of Bioschemas classes; however, participation in the development of the classes was hampered by the technological expertise needed in order to update a Bioschemas class. The availability of user-friendly tools can make it easier to find and use Schema.org and other community-driven schema classes, and empower data providers and researchers to engage in schema authoring and sharing.

Most tools for utilizing existing Schema.org classes focus on the utilization of an existing schema (such as markup generation) and lack the ability to customize the schema in a Schema.org-compliant way. Tools that do allow customizing/creating a schema (e.g., Bioschemas GoWeb) often require some degree of programming. The DDE Schema Playground is a browser-based tool that enables members of the research community to easily adapt schemas to suit their need and to enable community re-use of their schemas through the DDE schema registry. This encourages and empowers researchers to structure their data in a Schema.org-compliant fashion earlier on in the scientific research process rather than as an afterthought. The schema authoring by the research community, for the research community will encourage the creation and adoption of new classes and properties, which may have previously been neglected due to the absence of representation (e.g.,

volunteers with subject matter expertise) in data standardization communities. In this fashion, the DDE Schema Playground allows for researchers to express and share their data structuring needs with the data standardization community without diverting attention away from their primary research efforts. Data standardization communities also benefit because their volunteer time can be concentrated on classes already in use by researchers (but could benefit from some standardization), and diverted away from classes which lack interest/support from the research community at large.

There are many ways to express schemas (i.e., SHACL, ShEx), but the DDE only supports the expression of schemas in JSON-LD/JSON Schema format due to the widespread adoption of the JSON-LD format by resource providers and library/tool developers. In addition to this restriction, there are important limitations as to what can or cannot be registered into the DDE schema registry. Schema registration in the DDE is currently limited class-based schema (i.e., classes described by sets of properties) rooted in Schema.org, while many well-used, domain-neutral metadata ontologies (such as DCMI) and schema have properties that are not necessarily tied to any class. These classless metadata vocabularies intentionally do not group the properties into classes in order to encourage the mix-and-match of properties. Although classless metadata vocabularies cannot be registered in their entirety as classless properties in the DDE at this time, the DDE can flexibly ingest properties from any metadata vocabulary (whether or not they are class-based) as long as it is properly formatted (i.e., conforms to JSON-LD/JSON Schema formatting). This means that users can build their schema by extending from Schema.org, Bioschemas, or any registered schema, and incorporate properties from OWL, DCMI, or any other accessible vocabulary as needed. For example, all Bioschemas profile classes also include the *conformsTo* property from DCMI, and the NIAID *Dataset* schema [36] also leverages properties from OWL. In theory, classes inheriting just a single property from a Schema.org class, but otherwise built entirely from other metadata ontologies can be viewed and registered in the DDE.

We tested the use of the DDE Schema Playground to create customized Schema.org-compliant classes that could be used to normalize metadata between multiple types (datasets, clinical trials, publications, etc.) of COVID-19-related resources and applied these schemas towards a searchable resource site (<https://outbreak.info>). The Outbreak resource schema is available in the DDE schema registry which also includes schemas from Schema.org, Bioschemas, BioLink, the National COVID Cohort Collaborative (N3C), the National Institute of Allergy and Infectious Diseases (NIAID) and more. We hope others will join us in making their open data more interpretable, interoperable, and reusable by adding their schemas to the schema registry.

Conclusion

We have created a user-friendly browser-based tool which facilitates the application of Schema.org towards biomedical research outputs. We demonstrate its use with the creation of the Outbreak.info schema, its adoption into the Bioschemas schema development pipeline, and we encourage others to register and reuse Schema.org-compliant schemas. We welcome user feedback which has and continues to help identify desirable new features and tools (i.e., metadata validation tools) which will be added in the near future.

Availability and requirements

Project name: Data Discovery Engine Schema Playground.

Project home page: <https://discovery.biothings.io/schema-playground>

Project source code: <https://github.com/biothings/discovery-app>

Operating system(s): Web-based, Platform independent.

Programming language: JavaScript, Python.

Other requirements: GitHub account for schema editing.

License: Creative Commons Attribution 4.0 International license (Content), Apache 2.0 (source code).

Any restrictions to use by non-academics: No.

Abbreviations

CSV	Comma-Separated Values
DCAP	Dublin Core Application Profile
DCMI	Dublin Core Metadata Initiative
DDE	Data Discovery Engine
EDAM	EMBRACE Data and Methods
ESIP	Earth Science Information Partners
FAIR/FAIRness	Findability, Accessibility, Interoperability, Reusability
HTML	HyperText Markup Language
JAMA	Journal of the American Medical Association
JSON-LD	JavaScript Object Notation-Linked Data
LRMI	Learning Resource Metadata Innovation
MIAPPE	Minimal Information About Plant Phenotyping Experiment
N3C	National COVID Cohort Collaborative
NCT	National Clinical Trials
NEJM	New England Journal of Medicine
NIAID	National Institute of Allergy and Infectious Diseases
OWL	Web Ontology Language
PPEO	Plant Phenotyping Experiment Ontology
PRS	Protocol Registration System
RDF	Resource Description Framework
SDK	Software Development Kit
SHACL	Shapes Constraint Language
ShEx	Shape Expressions
TTL	Terse RDF Triple Language
W3C	World Wide Web Consortium
WHOCT	World Health Organization Clinical Trials
YAML	YAML ain't markup language

Supplementary Information

The online version contains supplementary material available at <https://doi.org/10.1186/s12859-023-05258-4>.

Additional files 1: Supplemental Figures. The interface and features of the Schema Playground.

Additional file 2. Supplemental Table 1. Comparison of Schema Playground with other tools.

Acknowledgements

We thank Ben Rush for his suggestions early on in the development of the Outbreak.info schema. We thank the Bioschemas community, especially Nick Juty and Alasdair Gray for their feedback, suggestions, and assistance in improving the DDE and integrating it into the Bioschemas development and maintenance pipeline.

Author contributions

MC developed the front-end of the DDE Schema Playground with feedback from GT, LDH, and JM. XZ and JX developed the backend of the DDE Schema Playground with feedback from MC. CW and AIS guided the overall development of the DDE Schema Playground. GT developed the Outbreak.info schema with feedback from LDH and JM. GT wrote the manuscript with feedback from LDH, AIS, and CW. All authors read and approved the final manuscript.

Funding

The development of the Data Discovery Engine was supported by the National Center for Advancing Translational Sciences, as part of the National Center for Data to Health (5 U24 TR002306) award to CW. Work on Outbreak.info was supported by National Institute for Allergy and Infectious Diseases (5 U19 AI135995-02), the National Center for Data to Health (5 U24 TR002306) and Centers for Disease Control and Prevention (75D30120C09795).

Availability of data and materials

Data sharing is not applicable to this article as no datasets were generated or analyzed during the current study; however, the Outbreak.info schema generated is available at <https://discovery.biothings.io/view/outbreak> and from GitHub at <https://github.com/outbreak-info/outbreak.info-resources/tree/master/yaml>. The primary classes of the Outbreak schema may be viewed at: Analysis: <https://discovery.biothings.io/view/outbreak/outbreak:Analysis>. ClinicalTrial: <https://discovery.biothings.io/view/outbreak/outbreak:ClinicalTrial>. ComputationalTool: <https://discovery.biothings.io/view/outbreak/outbreak:ComputationalTool>. Dataset: <https://discovery.biothings.io/view/outbreak/outbreak:Dataset>. Protocol: <https://discovery.biothings.io/view/outbreak/outbreak:Protocol>. Publication: <https://discovery.biothings.io/view/outbreak/outbreak:Publication>.

Declarations

Ethics approval and consent to participate

Not applicable.

Consent for publication

Not applicable.

Competing interests

The authors declare that they have no competing interests.

Received: 7 June 2022 Accepted: 27 March 2023

Published online: 20 April 2023

References

- Baker T. Dublin Core Application Profiles at eleven years (2011). Dublincore.org. 2019. Available at: https://www.dublincore.org/blog/2011/application_profile/. Accessed 19 Apr 2022.
- Barker P. examples for LRMI properties—fixes #912. github.com/schemaorg. 2015. Available at: <https://github.com/schemaorg/schemaorg/pull/902>. Accessed 7 Feb 2023.
- Barker P, Campbell LM. LRMI, Learning Resource Metadata on the Web. WWW '15 Companion: Proceedings of the 24th International Conference on World Wide Web. 2015. P687. <https://doi.org/10.1145/2740908.2741745>
- Boeckhout M, Zielhuis GA, Bredenoord AL. The FAIR guiding principles for data stewardship: Fair enough? *Eur J Hum Genet.* 2018;26(7):931–6. <https://doi.org/10.1038/s41431-018-0160-0>.
- Bruskiewich R, Deepak, Moxon S, Mungall C. Biolink Model. Biolink Model. 2021 <https://doi.org/10.5281/zenodo.1242670>. Available at: <https://biolink.github.io/biolink-model/>. Accessed 1 Sept 2021.
- Canham S, Ohmann C. A metadata schema for data objects in clinical research. *Trials.* 2016;17(1):557. <https://doi.org/10.1186/s13063-016-1686-5>. PMID:27881150;PMCID:PMC5122021.
- CD2H. Center for Data to Health. 2021 Available at: <https://cd2h.org/>. Accessed 10 Aug 2021.
- Coyle K, Baker T. Guidelines for Dublin Core Application Profiles. Dublincore.org. 2009. Available at: <https://www.dublincore.org/specifications/dublin-core/profile-guidelines/#>. Accessed 19 April 2022.
- Dublin Core Metadata Initiative. DCMI Metadata Terms. Dublincore.org. 2020. Available at: <https://www.dublincore.org/specifications/dublin-core/dcmi-terms/>. Accessed 19 April 2022.
- Gangavarapu K, Latif A, Mullen J, Alkuzweny M, Hufbauer E, Tsueng G, Haag E, Zeller M, Aceves C, Zaiets K, Cano M, Zhou J, Qian Z, Sattler R, Matteson N, Levy J, Lee R, Freitas L, Maurer-Stroh S, Suchard M, Wu C, Su A, Andersen K and Hughes L. Outbreak.info genomic reports: scalable and dynamic surveillance of SARS-CoV-2 variants and mutations. *Nat Methods* 2023. <https://doi.org/10.1038/s41592-023-01769-3>.
- Gray AJ, Gobel C, Jimenez RC, Bioschemas Community. Bioschemas: From Potato Salad to Protein Annotation. ISWC 2017, CEUR-WS.org. 2017. Available at: <http://ceur-ws.org/Vol-1963/paper579.pdf>
- Gundersen S, Boddu S, Capella-Gutierrez S, Drabløs F, Fernández JM, Kompova R, Taylor K, Titov D, Zerbino D, Hovig E. Recommendations for the FAIRification of genomic track metadata. *F1000Res.* 2021;10:ELIXIR-268. doi: <https://doi.org/10.12688/f1000research.28449.1>
- Haendel M, Chute C, Bennett T, Eichmann D, Guinney J, Kibbe W, Payne P, Pfaff E, Robinson P, Saltz J, Spratt H, Suver C, Wilbanks J, Wilcox A, Williams A, Wu C, Blacketer C, Bradford R, Cimino J, Clark M, Colmenares E, Francis P, Gabriel D, Graves A, Hemadri R, Hong S, Hripscak G, Jiao D, Klann J, Kostka K, Lee A, Lehmann H, Lingrey L, Miller R, Morris M, Murphy S, Natarajan K, Palchuk M, Sheikh U, Solbrig H, Visweswaran S, Walden A, Walters K, Weber G, Zhang X, Zhu R, Amor B, Girvin A, Manna A, Qureshi N, Kurilla M, Michael S, Portilla L, Rutter J, Austin C, Gersing K. The National COVID Cohort Collaborative (N3C): Rationale, design, infrastructure, and deployment. *J Am Med Inf Assoc.* 2020;28(3):427–43.
- Hollmann S, Kremer A, Baebler Š, Trefois C, Gruden K, Rudnicki WR, Tong W, Gruca A, Bongcam-Rudloff E, Evelo CT, Nechyporenko A, Frohme M, Šafránek D, Regierer B, D'Elia D. The need for standardisation in life science research—an approach to excellence and trust. *F1000Res.* 2020;9:1398. <https://doi.org/10.12688/f1000research.27500.2>.
- Holub P, Kohlmayer F, Prasser F, Mayrhofer MT, Schlünder I, Martin GM, Casati S, Koumakis L, Wutte A, Kozera Ł, Strapagiel D, Anton G, Zanetti G, Sezerman OU, Mendy M, Valik D, Lavitrano M, Dagher G, Zatloukal K, van Ommen GB, Litton JE. Enhancing reuse of data and biological material in medical research: from FAIR to FAIR-Health. *Biopreserv Biobank.* 2018;16(2):97–105. <https://doi.org/10.1089/bio.2017.0110>.
- Hruby GW, Hoxha J, Ravichandran PC, Mendonça EA, Hanauer DA, Weng C. A data-driven concept schema for defining clinical research data needs. *Int J Med Inform.* 2016;91:1–9. <https://doi.org/10.1016/j.ijmedinf.2016.03.008>.

17. Hughes LD, Tsueng G, DiGiovanna J, Horvath TD, Rasmussen LV, Savidge TC, Stoeger T, Turkarslan S, Wu Q, Wu C, Su AI, Pache L and the NIAID Systems Biology Data Dissemination Working Group, Addressing barriers in FAIR data practices for biomedical data. *Scientific Data*. 2023. <https://doi.org/10.1038/s41597-023-01969-8>.
18. Ison J, Kalas M, Jonassen I, Bolser D, Uludag M, McWilliam H, Malone J, Lopez R, Pettifer S, Rice P. EDAM: an ontology of bioinformatics operations, types of data and identifiers, topics and formats. *Bioinformatics*. 2013;29(10):1325–32. <https://doi.org/10.1093/bioinformatics/btt113>.
19. Jauer ML, Deserno TM. Data Provenance Standards and Recommendations for FAIR Data. *Stud Health Technol Inform*. 2020;16(270):1237–8. <https://doi.org/10.3233/SHTI200380>.
20. Jones MB, Richard S, Vieglais D, Shepherd A, Duerr R, Fils D, McGibbney L. Science-on-Schema.org v1.2.0 (Version 1.2.0). Zenodo. 2021. <https://doi.org/10.5281/zenodo.4477164>
21. Jsdn.org. JSON-LD—JSON for Linking Data. 2022. Available at: <https://json-ld.org/>. Accessed 6 June 2022.
22. JSON Schema. Specification. 2022. Available at: <https://json-schema.org/specification.html>. Accessed 6 June 2022.
23. Koesten L, Vougiouklis P, Simperl E, Groth P. Dataset reuse: toward translating principles to practice. *Patterns* (N.Y). 2020;1(8):100136. <https://doi.org/10.1016/j.patter.2020.100136>.
24. Larmande P, Costa BV, Cwiek-Kupczynska H, Cornut G, Neveu P, Chaves I, Pommier C, Papoutsoglou E, Ruiz M, Faria D, Laporte MA. Plant phenotype experiment ontology. *AgroPortal*. 2019. Available at: <https://agroportal.lirmm.fr/ontologies/PPEO>. Accessed 8 Feb 2023.
25. Lelong S, Zhou X, Afrasiabi C, Qian Z, Cano MA, Tsueng G, Xin J, Mullen J, Yao Y, Avila R, Taylor G, Su AI, Wu C. BioThings SDK: a toolkit for building high-performance data APIs in biomedical research. *Bioinformatics*. 2022;38(7):2077–9. <https://doi.org/10.1093/bioinformatics/btac017>.
26. Miron L, Gonçalves RS, Musen MA. Obstacles to the reuse of study metadata in ClinicalTrials.gov. *Sci Data*. 2020 **7**, 443. <https://doi.org/10.1038/s41597-020-00780-z>
27. MongoDB. What Is MongoDB? 2022. Available at: <https://www.mongodb.com/what-is-mongodb>. Accessed 06 June 2022.
28. Outbreak.info. Research Library. 2020. Available at: <https://outbreak.info/resources>. Accessed 10 Aug 2021.
29. Papadiamantis AG, Klaessig FC, Exner TE, Hofer S, Hofstaetter N, Himly M, Williams MA, Doganis P, Hoover MD, Afantitis A, Melagraki G, Nolan TS, Rumble J, Maier D, Lynch I. Metadata Stewardship in Nanosafety Research: Community-Driven Organisation of Metadata Schemas to Support FAIR Nanoscience Data. *Nanomaterials* (Basel). 2020;10(10):2033. <https://doi.org/10.3390/nano10102033>.
30. Profiti G, Jimenez RC, Zambelli F et al. Using community events to increase quality and adoption of standards: the case of Bioschemas [version 1; not peer reviewed]. *F1000Research* 2018, 7(ELIXIR):1696 (document). <https://doi.org/10.7490/f1000research.1116233.1>
31. Sansone SA, Gonzalez-Beltran A, Rocca-Serra P, Alter G, Grethe JS, Xu H, Fore IM, Lyle J, Gururaj AE, Chen X, Kim HE, Zong N, Li Y, Liu R, Ozyurt IB, Ohno-Machado L. DATS, the data tag suite to enable discoverability of datasets. *Sci Data*. 2017;4:170059. <https://doi.org/10.1038/sdata.2017.59>.
32. Schema.org. Releases. 2022. Available at: <https://schema.org/docs/releases.html>. Accessed 06 June 2022.
33. Shepherd A, Jones MB, Richard S, Jarboe N, Vieglais S, Fils D, Duerr R, Verhey C, Minch M, Mecum B, Bentley N. Science-on-Schema.org v1.3.0. Zenodo. 2022. <https://doi.org/10.5281/zenodo.6502539>
34. Tong Z. What is an Elasticsearch Index? | Elastic. Elastic.co. 2013. Available at: <https://www.elastic.co/blog/what-is-an-elasticsearch-index>. Accessed 06 June 2022.
35. Tornadoweb.org. Tornado Web Server—Tornado 6.1 documentation. 2022. Available at: <https://www.tornadoweb.org/en/stable/>. Accessed 06 June 2022.
36. Tsueng G, Cano MA, Bento J, Czech C, Kang M, Pache L, Rasmussen LV, Savidge TC, Starren J, Wu Q, Xin J, Yeaman MR, Zhou X, Su AI, Wu C, Brown L, Shabman RS, Hughes LD, and the NIAID Systems Biology Data Dissemination Working Group. Developing a standardized but extendable framework to increase the findability of infectious disease datasets. *Sci Data*. 2023a. <https://doi.org/10.1038/s41597-023-01968-9>
37. Tsueng G, Mullen J, Alkuzweny M, Cano M, Rush B, Haag E; Lin J, Welzel DJ, Zhou X, Qian Z, Latif AA, Emory H, Zeller M, Andersen KG, Wu C, Su AI, Gangavarapu K, Hughes LD. Outbreak.info Research Library: a standardized, searchable platform to discover and explore COVID-19 resources. *Nat Methods*. 2023b. <https://doi.org/10.1038/s41592-023-01770-w>.
38. W3C Team. Call for Participation in Schema.org Community Group. 2015. Available at: <https://www.w3.org/community/schemaorg/2015/03/31/call-for-participation-in-schema-org-community-group/>. Accessed 8 Feb 2023.
39. W3.org. About W3C groups. n.d. Available at: <https://www.w3.org/groups/>. Accessed 8 February 2023.
40. W3.org. W3C RDF Core Working Group (Closed). 2004. Available at: <https://www.w3.org/2001/sw/RDFCore/#documents>. Accessed 6 June 2022.
41. Wilkinson MD, Dumontier M, Aalbersberg IJ, Appleton G, Axton M, Baak A, Blomberg N, Boiten JW, da Silva Santos LB, Bourne PE, Bouwman J, Brookes AJ, Clark T, Crosas M, Dillo I, Dumon O, Edmunds S, Evelo CT, Finkers R, Gonzalez-Beltran A, Gray AJ, Groth P, Goble C, Grethe JS, Heringa J, 't Hoen PA, Hooft R, Kuhn T, Kok R, Kok J, Lusher SJ, Martone ME, Mons A, Packer AL, Persson B, Rocca-Serra P, Roos M, van Schaik R, Sansone SA, Schultes E, Sengstag T, Slater T, Strawn G, Swertz MA, Thompson M, van der Lei J, van Mulligen E, Velterop J, Waagmeester A, Wittenburg P, Wolstencroft K, Zhao J, Mons B. The FAIR Guiding Principles for scientific data management and stewardship. *Sci Data*. 2016;3:160018. <https://doi.org/10.1038/sdata.2016.18>. Erratum in: *Sci Data*. 2019 Mar 19;6(1):6. PMID: 26978244; PMCID: PMC4792175.
42. World Health Organization (WHO)/International Committee of Medical Journal Editors (ICMJE)-ClinicalTrials.gov Cross Reference 2019. Available at: <https://prsinfo.clinicaltrials.gov/trainTrainer/WHO-ICMJE-ClinTrialsgov-Cross-Ref.pdf>. Accessed 10 Apr 2022.
43. You E. Vue.js—The Progressive JavaScript Framework. 2014. Available at: <https://vuejs.org/#is-vue-lightweight>. Accessed 6 June 2022.

Publisher's Note

Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.