

UC Santa Cruz

UC Santa Cruz Previously Published Works

Title

SnoVault and encodeD: A novel object-based storage system and applications to ENCODE metadata.

Permalink

<https://escholarship.org/uc/item/5740q1gk>

Journal

PLoS ONE, 12(4)

Authors

Hitz, Benjamin
Rowe, Laurence
Podduturi, Nikhil
[et al.](#)

Publication Date

2017

DOI

10.1371/journal.pone.0175310

Copyright Information

This work is made available under the terms of a Creative Commons Attribution License, available at <https://creativecommons.org/licenses/by/4.0/>

Peer reviewed

RESEARCH ARTICLE

SnoVault and encodeD: A novel object-based storage system and applications to ENCODE metadata

Benjamin C. Hitz¹✉, Laurence D. Rowe¹✉, Nikhil R. Poddaturi¹, David I. Glick¹, Ulugbek K. Baymuradov¹, Venkat S. Malladi¹, Esther T. Chan¹, Jean M. Davidson¹, Idan Gabdank¹, Aditi K. Narayana¹, Kathrina C. Onate¹, Jason Hilton¹, Marcus C. Ho¹, Brian T. Lee², Stuart R. Miyasato¹, Timothy R. Dreszer¹, Cricket A. Sloan¹, J. Seth Strattan¹, Forrest Y. Tanaka¹, Eurie L. Hong¹, J. Michael Cherry¹*

1 Stanford University School of Medicine, Department of Genetics, Stanford, California, United States of America, **2** University of California Santa Cruz, Baskin School of Engineering, Center for Biomolecular Science and Engineering, Santa Cruz, California, United States of America

✉ These authors contributed equally to this work.

* cherry@stanford.edu



OPEN ACCESS

Citation: Hitz BC, Rowe LD, Poddaturi NR, Glick DI, Baymuradov UK, Malladi VS, et al. (2017) SnoVault and encodeD: A novel object-based storage system and applications to ENCODE metadata. PLoS ONE 12(4): e0175310. <https://doi.org/10.1371/journal.pone.0175310>

Editor: Zhang Zhang, Beijing Institute of Genomics Chinese Academy of Sciences, CHINA

Received: October 15, 2016

Accepted: March 23, 2017

Published: April 12, 2017

Copyright: © 2017 Hitz et al. This is an open access article distributed under the terms of the [Creative Commons Attribution License](https://creativecommons.org/licenses/by/4.0/), which permits unrestricted use, distribution, and reproduction in any medium, provided the original author and source are credited.

Data Availability Statement: The software and database models can be found at <https://github.com/ENCODE-DCC/encodeD> and <https://github.com/ENCODE-DCC/snovault>. The production website (generated by the above software) for the ENCODE portal is <http://www.encodeproject.org>.

Funding: This work was funded by the National Human Genome Resource Institute of the National Institute of Health (<https://www.genome.gov/>) under the following U41 grant and supplement: U41HG006992. The funders had no role in study

Abstract

The Encyclopedia of DNA elements (ENCODE) project is an ongoing collaborative effort to create a comprehensive catalog of functional elements initiated shortly after the completion of the Human Genome Project. The current database exceeds 6500 experiments across more than 450 cell lines and tissues using a wide array of experimental techniques to study the chromatin structure, regulatory and transcriptional landscape of the *H. sapiens* and *M. musculus* genomes. All ENCODE experimental data, metadata, and associated computational analyses are submitted to the ENCODE Data Coordination Center (DCC) for validation, tracking, storage, unified processing, and distribution to community resources and the scientific community. As the volume of data increases, the identification and organization of experimental details becomes increasingly intricate and demands careful curation. The ENCODE DCC has created a general purpose software system, known as SnoVault, that supports metadata and file submission, a database used for metadata storage, web pages for displaying the metadata and a robust API for querying the metadata. The software is fully open-source, code and installation instructions can be found at: <http://github.com/ENCODE-DCC/snovault/> (for the generic database) and <http://github.com/ENCODE-DCC/encodeD/> to store genomic data in the manner of ENCODE. The core database engine, SnoVault (which is completely independent of ENCODE, genomic data, or bioinformatic data) has been released as a separate Python package.

Introduction

The Encyclopedia of DNA Elements (ENCODE) project (<https://www.encodeproject.org/>) is an international consortium with a goal of annotating regions of the genome[1–6] initiated as a follow up to the Human Genome Project[7]. The Data Coordination Center (DCC)[8–10] is

design, data collection and analysis, decision to publish, or preparation of the manuscript.

Competing interests: The authors have declared that no competing interests exist.

charged with validating, tracking, storing, processing, visualizing and distributing these data files and their metadata to the scientific community. During 6 years of the pilot and initial scale-up phase, the project surveyed the landscape of the *H. sapiens* and *M. musculus* genomes using over 20 high-throughput genomic assays in more than 350 different cell and tissue types, resulting in over 3000 datasets. ENCODE employs a wide range of genomics techniques[11,12]; [13]. To determine the transcribed elements including mRNAs, the GENCODE (<http://www.genencodegenes.org>) project employs a team of biocurators who synthesize data from a wide range of assays to build accurate gene models that address the full complexity of vertebrate transcription including alternative promoters, alternative splicing, and RNA editing[14]. These assays include sequencing selected RNAs for promoter characterization, next-generation RNA sequencing and assembly methods for detecting alternative splicing and for quantizing the level of RNA in different cellular contexts[15]. To determine transcriptional regulatory regions ENCODE uses DNase I[15,16] and nucleosome positioning assays to locate regions of the chromosome accessible to regulatory elements, DNA methylation assays and ChIP-seq of modified histones to define the overall chromatin architecture [17][18], and ChIP-seq of transcription factors to determine the players involved in the regulatory interactions[19]. To determine which regulatory elements interact with each other, ENCODE uses ChIA-PET[19,20] and Hi-C[21]. To identify regulatory elements that operate on the RNA rather than the DNA level, ENCODE employs RIP-seq[22] as well as a variety of computational assays [23,24]; [25][23,24][26]. Over the current phase of ENCODE that began in 2012, the diversity and volume of data has increased as new genomic assays are added to the project. In addition to data from experiments performed in diversity of biological samples, data from additional species (*D. melanogaster* and *C. elegans*) [27] and other projects have been incorporated. Experimental data are validated and analyzed using new methods, and higher level analyses have been included to form a true Encyclopedia of genomic elements. The DCC currently houses experimental data and metadata from ENCODE3 (the current ENCODE), modENCODE[28], ENCODE2[29](the phase of ENCODE ending in 2012), mouse-ENCODE[29,30], Roadmap for Epigenomics Mapping (REMC)[29–31], modERN (fly and worm experiments associated with ENCODE3) and Genomics for Gene Regulation (GGR) (<https://www.genome.gov/27561317>). These contributions are summarized on the ENCODE Portal at: <https://www.encodeproject.org/about/data-access/>.

The primary goals of the DCC are to track and compile the experimental metadata for each experiment produced by the consortium. This metadata is an absolutely critical resource both for internal consortium metrics and bookkeeping and to provide the best possible user experience for scientists and educators wishing to make sense of such a large corpus of datafiles, experiments, and outputs[8–10].

One of the challenges in this field is keeping a flexible data model without sacrificing data continuity and integrity. To this end, the SnoVault system includes an infrastructure to build (quality or validity) data audits and reports to maintain data integrity. To resolve these conflicting design criteria, we have developed a hybrid relational/object database using JSON, JSON-LD (<http://json-ld.org/>), JSON-SCHEMA (<http://json-schema.org/>), and PostgreSQL (<http://www.postgresql.org>). For efficient searching and display, the data are de-normalized and indexed in Elasticsearch (<https://www.elastic.co/>). These software components are wrapped in a Python Pyramid web framework application which executes the business logic of the application and provides the RESTful web API. The website itself, also called the “front-end” is constructed using the ReactJS framework (<https://facebook.github.io/react/>).

A hybrid relational-object data store

Traditionally, biological databases have been implemented using a relational database model. Relational database software, both open-source (e.g., MySQL or PostgreSQL) and commercial

(e.g. Oracle), is ubiquitous on all current hardware. It is a robust system for maintaining transactional integrity and keeping concepts normalized (storing each data only once, and connecting them via foreign-key relations). In the SnoVault system, each major category of metadata is a JSON object or “document” typically corresponding to an experimental component of a genomic assay, such as an RNA-seq or ChIP-seq assay. For example, there are JSON objects that represent a specific human donor, which can be linked to 1 or more biological samples (biosamples) from that donor as well as objects that represent key reagents in an assay such as an antibody lot that is used in a ChIP-seq assay or an shRNA used to knockdown a gene target. In addition, objects represent the files and computational analyses required to generate these files. In addition to objects that represent experimental components and reagents of a genomic assay, there are objects that allow the grouping of these objects to represent a single biological replicate and a grouping of biological replicates to represent a replicated assay.

These JSON objects are stored in our transactional implementation of an document store in PostgreSQL. This low-level database is analogous to a commercial or OSS object/document store like MongoDB (<http://www.mongodb.com>) or couchDB (<http://couchdb.apache.org>) with additional transactional and rollback capabilities. The other critical features we have added to “classical” object storage methods are JSON-SCHEMA and JSON-LD. JSON-SCHEMA is simply a way to template JSON objects to control the allowed fields in each object and their allowed values in each field. There are separate schema files for each object type, which are analogous to tables in a traditional relational database. JSON-LD is a data standards format that provides a unified and straightforward way of linking JSON objects via pointers or “links” (analogous to foreign keys in a traditional relational database). These links are used in the metadata storage in a manner analogous to foreign key relationships in a traditional relational model. These links allow our software to “upgrade” objects to a new version of a schema when the schema changes. Our upgrading system allows us to return objects that are compatible with the most current schema without necessarily reindexing them in Elasticsearch; they are permanently upgraded when the objects are PUT or PATCHed. Herein lies the critical advance we have made—JSON-LD and JSON-SCHEMA allow us to strike a balance between the flexibility of an unstructured object database and the data integrity and normalization features of a relational database.

To date, the encodeD system (built on SnoVault) contains 75 object types that are able to handle over 40 different classes and flavors of genomic assays. The addition of new assays during the course of the ENCODE project has resulted in minimal changes that often include the addition of a new experimental properties in an existing object or a new object to allow new relationships to be created among existing objects. The flexibility of the SnoVault system allows the addition of these new properties or objects with minimal intervention from software engineers. Data wranglers or other data scientists are perfectly capable of extending or modifying the JSON-SCHEMA model, with only a small amount of training. The specific object types and properties as defined in the schema are available directly from the API via the /profiles/ endpoint. For convenience, each object type has a profile endpoint as well, for example: <https://www.encodeproject.org/profiles/experiment.json>.

We use the open-source python web framework Pyramid (<http://www.pylonsproject.org/>) as the base python layer of our back end. Pyramid’s excellent buildout system (similar in principle to “make”) allows us to manage our dependencies for python, node.js, and ruby and makes installation on an arbitrary UNIX machine very straight forward.

One of the features of ENCODE DCC website and RESTful API is the ability to set permissions (view, edit, or delete) on a single object independent of the objects that might refer to or embed it. Data that are released are publicly viewable without needing a login at all, while metadata that is still incomplete or in progress is available to view by consortium members. Viewing permissions are defined by a combination of the release status of an object. Editing

permissions are defined by a combination of the release status of an object and a membership in a user group. An “admin” user group has permission to edit all objects, regardless of release status. Individual users may belong to specific groups, as defined by labs, to allow editing of data that are generated by that lab. An individual user may be authorized to edit only data for 1 lab or as many labs as appropriate. These user groups are only allowed to edit data that is not released. In this way, we distinguish between “admin” level users (such as an ENCODE data wrangler) who have full rights to all data; a “submitter” level user (someone from a production lab) who can modify the objects of a lab that they submit for; a “consortium” user, who can view all submitted data, and a “base” user, who can only view data that is released. Pyramid has excellent support for detailed “Action Control Lists” (ACLs) which were utilized in our application. The ability to edit and view can be customized by defining the set of permissions that are granted to the specific group. Specifically, submitters from the ENCODE and GGR projects are distinguished without our system; members of those groups may neither view or edit unreleased data from another group.

Indexing and performance optimization. One of the challenges in web database development is to create a system that optimized both for WRITE access and storage and for READ access. Storing data requires transactional integrity, flexibility of schema and some level of denormalization. Our JSON-SCHEMA document model using links to other objects for references fulfills this goal, but retrieving complex linked objects using multiple GETs for each linked object (and so on down through the object tree) is very inefficient. We take the classical approach of indexing denormalized data that we need to retrieve rapidly, but add some modern twists. Each object or document in our system has a uuid that acts similar to a primary key in a relational database. Using JSON-LD, we reference other objects by their object type and uuid, in the parlance of JSON-LD these are called links. Lists of objects or more complex data structures can be used as well. Objects can be returned from the RESTful API either with or without embedded linked objects using a URL query parameter (frame = object, frame = embedded, frame = raw or frame = page). An embedded JSON-LD object is one in which it's links have been substituted with “copies” (or denormalization) of the linked objects (often just a subsection of the linked objects; the details are specified in the python code for the views of each object type).

To make the database performant, and to enable rapid and faceted searching of the metadata database implemented in encodeD, we index all frames of all objects in Elasticsearch (<https://www.elastic.co/products/elasticsearch>). Elasticsearch is a robust search and live indexing wrapper based around Apache Lucene (<https://lucene.apache.org/core/>). As new objects are added to the database via POST or modified via PATCH, they are indexed in real time by the indexer process.

The real-time indexing occurs according to the following rubric:

- When rendering a response, we record the set of embedded_uuids and linked_uuids used.
- The embedded_uuids are those objects embedded into the response or whose properties have been consulted in rendering of the response.
- Any change to one of these objects should cause an invalidation.
- Linked_uuids are the objects linked to in the response. Only changes to their URL (i.e, the canonical name of the link) need trigger an invalidation.
- When modifying objects, event subscribers keep track of which objects are updated and their resource paths before and after the modification. This is used to record the set of ‘updated_uuids’ and ‘renamed_uuids’ in the transaction log.
- The indexer process listens for notifications of new transactions.

- With the union of updated_uuids and union of renamed_uuids across each transaction in the log since its last indexing run, the indexer performs a search for all objects where embedded_uuids intersect with the updated_uuids or linked_uuids intersect with the renamed_uuids.
- The result is the set of invalidated objects which must be reindexed in addition to those that were modified (recorded in updated_uuids.)

RESTful API. Programmatic interaction with the ENCODE DCC metadata database is typically done through scripts written and executed on a local computer. These scripts interact with the database through an industry-standard, Hypertext-Transfer-Protocol-(HTTP)-based, Representational-state-transfer (RESTful) application programming interface (API). Data objects exchanged with the server conform to the standard JavaScript Object Notation (JSON) format. Each object type in the encodeD system (or any like system built with SnoVault) can be accessed as an individual object or a collection (list) of objects. For example, to return a specific ENCODE experiment with the accession ENCSR107SLP, you can send a “GET” request to the URL: <https://www.encodeproject.org/experiments/ENCSR107SLP/>. Accessions and unique names are aliased so that objects can be returned without referring to their object type in the URL. For example, <https://www.encodeproject.org/ENCSR107SLP> returns the same object. If the request is a browser, the JSON object returned is transformed into the webpage HTML by the ReactJS front end (Fig 1).

Otherwise, Accept-headers (<https://www.w3.org/Protocols/rfc2616/rfc2616-sec14.html>) or the URL query parameter format = json can be used to request that the system return text JSON objects instead of HTML pages; but the URLs are precisely the same (Figs 2 and 3). Because Elasticsearch indexes JSON objects, we can expose parts of the search API to allow command-line or programmatic access to our full facet search capability.

Members of the ENCODE consortium use either the REST API or HTML forms to submit data files and metadata to the DCC. New objects are created in the database by a POST request to the collection URL (for example, POSTing a valid biosample JSON object to the biosamples collection URL (<https://www.encodeproject.org/biosamples/>) creates a new biosample). Existing records are modified by a PATCH request to the object URL (for example, posting {"assay_term_name": "RNA-seq", "assay_term_id": "OBI:0001271"} to <https://www.encodeproject.org/ENCSR107SLP> would change it from a ChIP-seq experiment to an RNA-seq experiment). Connection between experiments objects and biosample objects is maintained by the “linking” objects for Replicate and Library. A Replicate is assigned to one and only one Experiment, and references a specific Library created from (and linked to) the specific Biosample from which nucleic acid was extracted. Most submitters have several hundred experiments, and use simple python scripts to post their metadata and upload their datafiles to the DCC. For cases where a low throughput solution is expedient, a logged-in submitter can edit the properties of objects that belong to a lab they submit for, using a simple HTML form interface.

When new objects are added or edited, the objects are compared against the schema of the object to determine validity. A valid object (either POST, PUT, or PATCH) returns a “200” and an invalid one returns “422” with a brief description of the error. Our JSON Schema objects are, by design, quite permissive in the fields that they require. This approach was chosen to encourage ENCODE submitters to provide as much metadata as possible as early in the process, in effect “registering” experiments, biosamples, and antibodies before the complete metadata is marshalled and even before the experiments themselves have been completed. The earlier metadata is submitted to the DCC, the earlier we can catch irregularities, errors or data

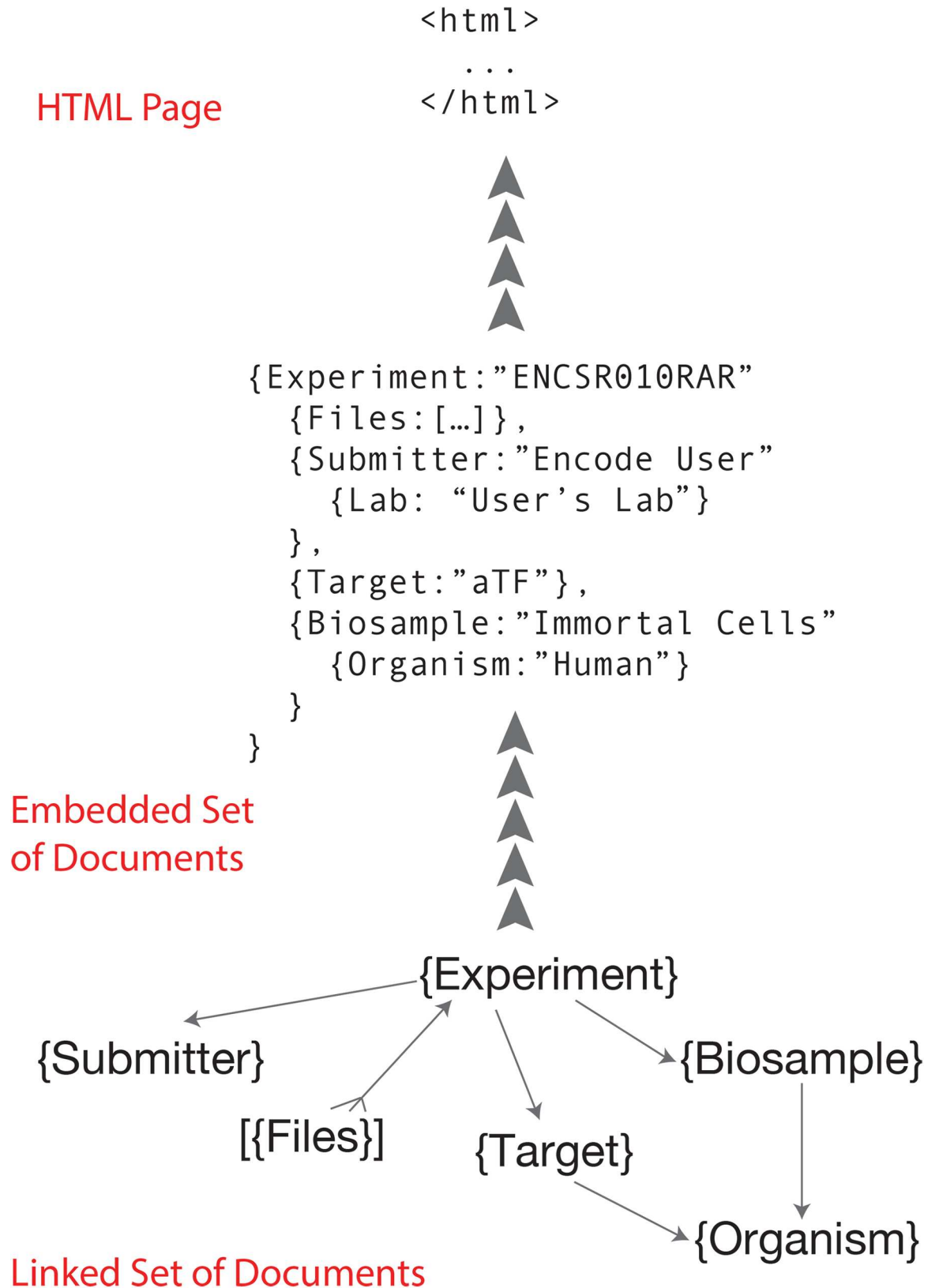


Fig 1. Object graph and embedding. The normalized object graph (linked set of documents) stored in PostgreSQL is framed as an expanded JSON-LD document (embedded set of documents) before indexing into Elasticsearch and rendering in JavaScript.

<https://doi.org/10.1371/journal.pone.0175310.g001>

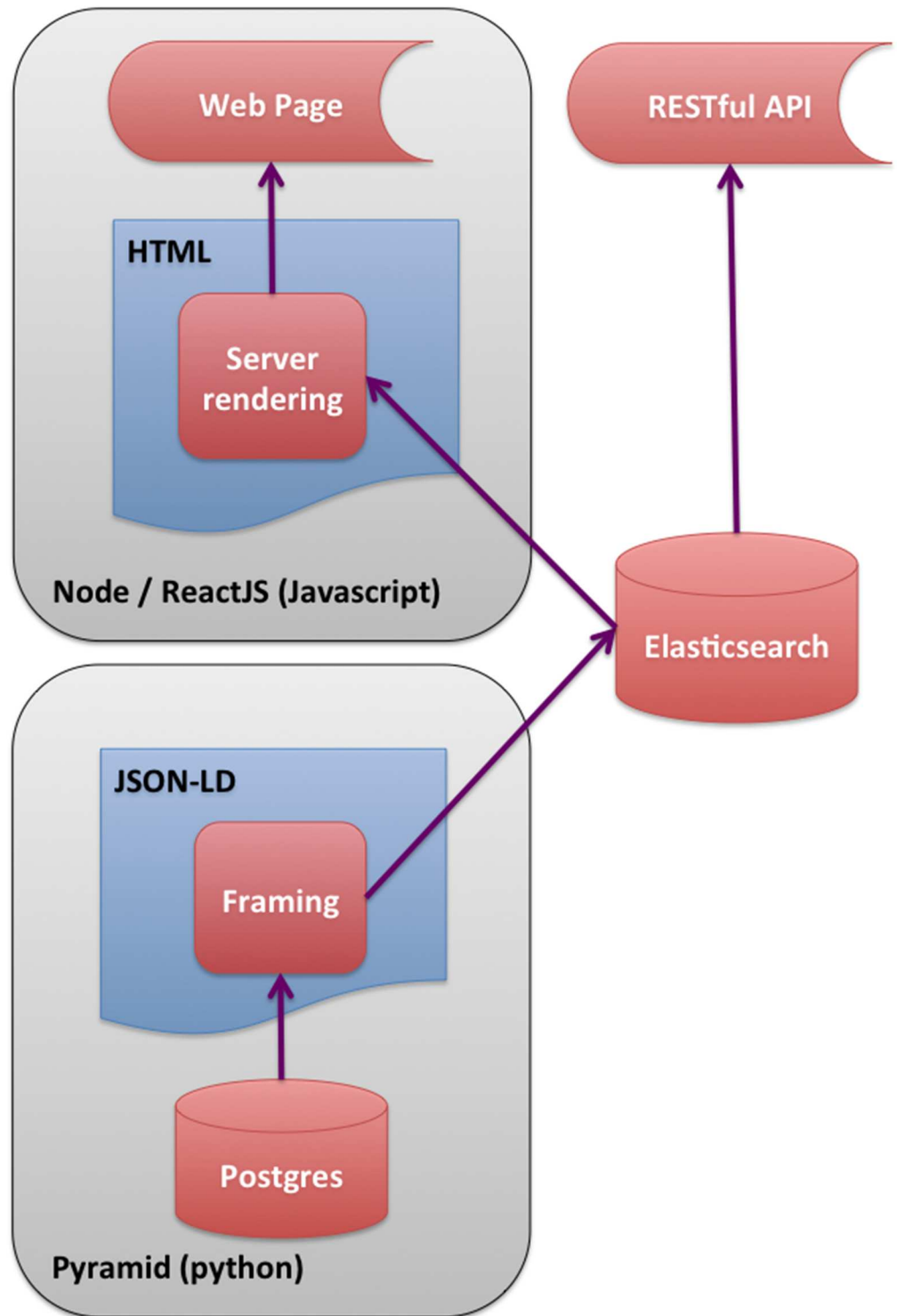


Fig 2. Software Stack. Schematic diagram of software stack showing different paths for page rendering. The most efficient rendering is with the HTML rendered on the server and the embedded documents indexed in Elasticsearch.

<https://doi.org/10.1371/journal.pone.0175310.g002>

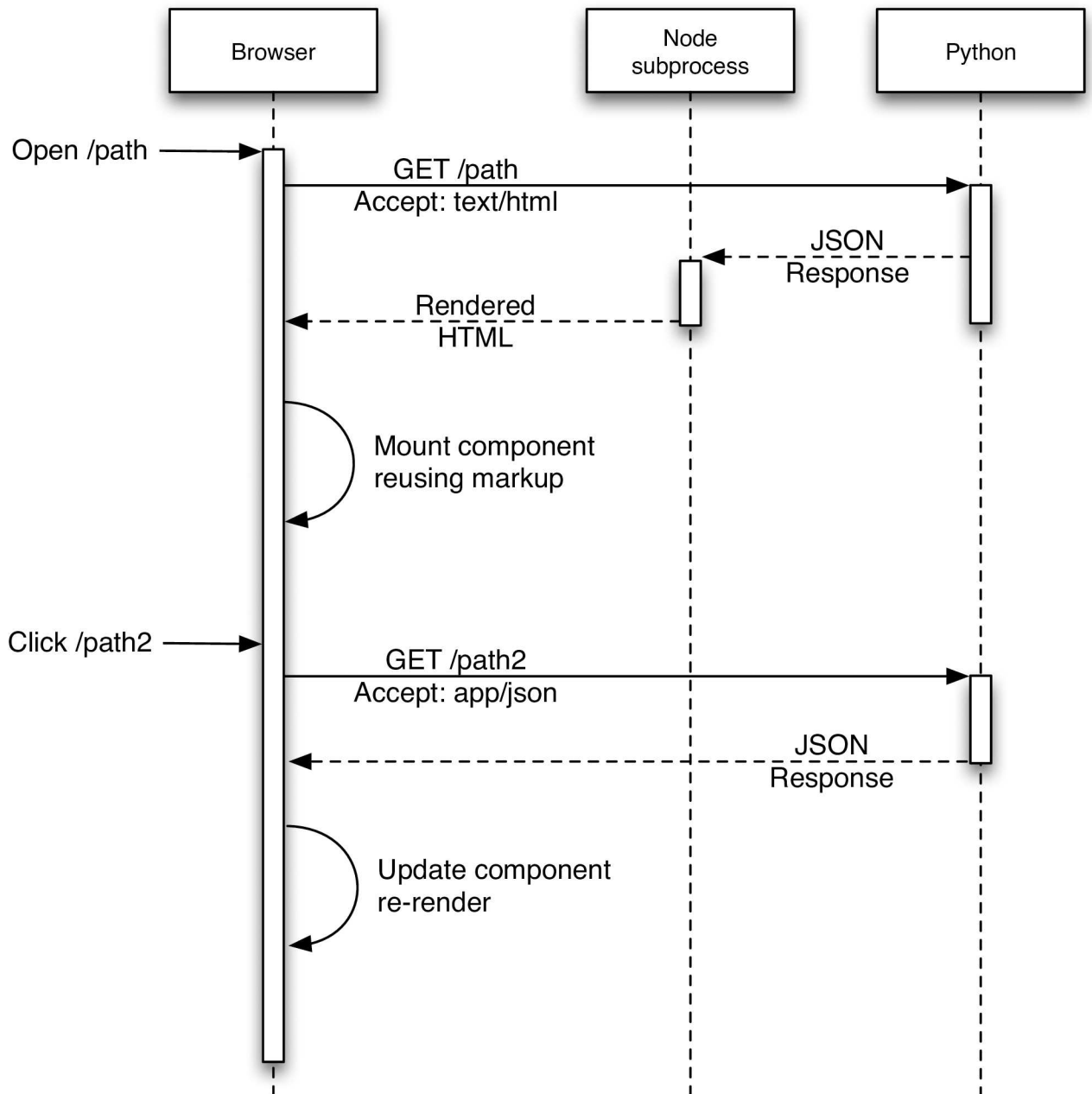
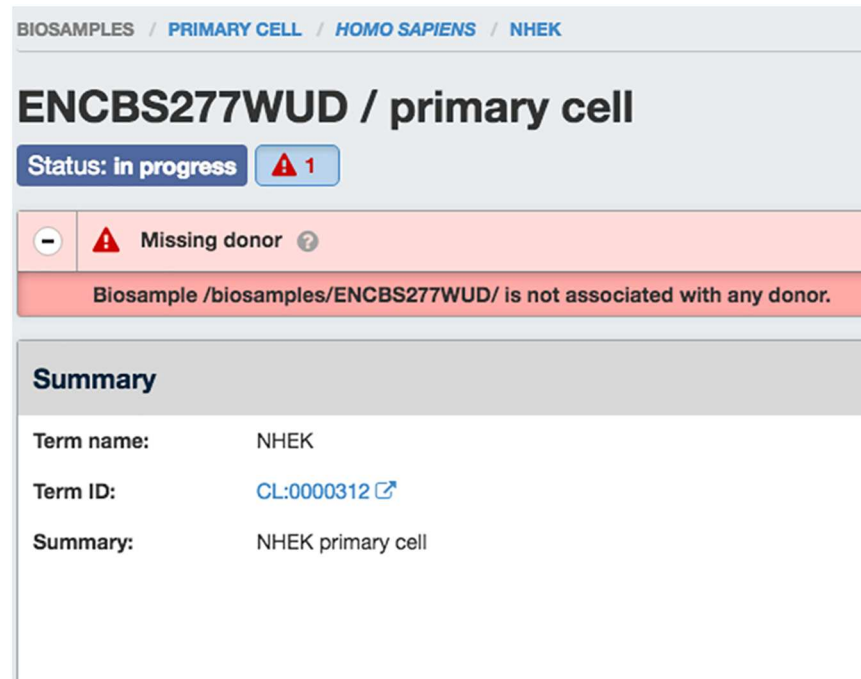


Fig 3. Rendering Overview. Initial page loads are rendered to HTML on the server for immediate display on the client. Once the JavaScript is fully loaded, subsequent page loads can be rendered on the client.

<https://doi.org/10.1371/journal.pone.0175310.g003>

model confusions. We use PUT to create a wholly new object, POST to replace an existing object in its entirety (except for uuid), and PATCH to update only specific properties of an existing object.

One extension we have made to the JSON-SCHEMA representation of objects is to add what we call “calculated properties”. The JSON-SCHEMA definition for a particular object or document type is the minimum specification for properties that define the POSTing or creation of an instance of the object. For display and reporting purposes, the python code extends objects by properties that are calculated from one or more existing object properties or even



BIOSAMPLES / PRIMARY CELL / HOMO SAPIENS / NHEK

ENCBS277WUD / primary cell

Status: in progress ▲ 1

▲ Missing donor ?

Biosample /biosamples/ENCBS277WUD/ is not associated with any donor.

Summary

Term name:	NHEK
Term ID:	CL:0000312
Summary:	NHEK primary cell

Fig 4. Example of the “missing donor” audit. This audit appears on relevant web pages (and in JSON return from the API), alerting submitters and DCC data wranglers that this epidermal keratinocyte biosample is missing the (human) donor object from whom it was derived.

<https://doi.org/10.1371/journal.pone.0175310.g004>

linked properties. For example, this is how we show the files that belong to an experiment. In our system, files have a link to their “home” Dataset (of which Experiment is a sub-class). The “Files” of an Experiment are a property that appears in the return JSON objects but is not specified in the schema. Instead, they are calculated as a “reverse link” from the file objects themselves. Note that when calculated properties are indexed, updating a “child” object like File, it will result in the invalidation of the parent Experiment (or other home Dataset object) and trigger a re-indexing of the parent in Elasticsearch.

A similar feature we have added are something we call Audits. Audits were developed to act as cross-object type data integrity measures (in a way similar to triggers can be used by RDBS). Audits are small snippets of python code that are checked whenever an object added or updated (via POST, PUT, or PATCH). Because our system considers each object POSTed individually, we cannot use the JSON-SCHEMA representation to validate properties across linked data types. Generally speaking, a criteria of the object and its associated linked objects are compared to look for violations (e.g., human donor object should not have a “lifestage” property) or to track metadata inconsistencies (both replicates of an experiment should have the same biosample ontology term) or quality metric (bam files are typically expected to have a minimum number of uniquely mapped reads). Audits are classified as ERROR (Red Flag), WARNING (Orange Flag), NON-COMPLAINT (Yellow Flag), or INTERNAL ACTION (Grey Ambulance) and are indexed with the page frame so that they can be displayed in the user interface (Fig 4). These audits help us retain a relational-like integrity across our sophisticated data graph, and give valuable feedback to submitter and wrangler alike. Finally, we use the audit system to report discrepancies between the ENCODE submission standards and quality control metrics and the true properties of the data and metadata. For example, if a particular experiment has low read mapped

read depth (relative to standards set by the consortium for that particular class of experiment, it will be reported on the user interface and API as a “Warning” with a yellow flag.

One of the drawbacks of not using a fully relational database backend is that one cannot query with arbitrary joins across data types. For web applications, this is generally not a significant drawback, because each web page or API end point has to be specified in the source code anyway. Specifically, in our system we must either embed properties from linked object type B into object type A, or perform multiple queries (GET requests). There are some purposes for which this is inconvenient; for example creating database reports that span many objects. We outline here three different approaches that can be used if more arbitrary queries are needed. One approach is to simply use the python command line interpreter and the excellent requests module (<http://docs.python-requests.org/en/master/>) to download collections of objects in JSON and traverse them programmatically. This method is simple and straightforward but can be a little tedious if multiple queries across many different object types is wanted. A second approach is to query the PostgreSQL database directly with psql. The objects are stored as JSON-type fields in PostgreSQL, and starting with version 9.4 the individual properties of the JSON documents can be queried using the JSON operators implemented by PostgreSQL. This method is not suitable if you are interested in calculated or embedded properties of the objects. Finally, we have added a utility to the SnoVault package that allows one to convert the JSON-LD object graph into Resource Description Format or RDF (<https://www.w3.org/RDF/>). A single API call to fetch all Objects (`/search/?type = Item&limit = all&format = json`) can be dumped to a file locally. Our script, `jsonld_rdf.py` uses `rdflib` (<https://rdflib.readthedocs.org/en/stable/>) converts objects to RDF in any standard format (xml, turtle, n3, trix and others). Once in RDF, it can be converted to a SPARQL (<https://www.w3.org/TR/rdf-sparql-query/>) store using a software package such as Virtuoso (<http://semanticweb.org/wiki/Virtuoso.html>). It could also be integrated with the EBI RDF Platform[32,33]. In this way, the ENCODE meta-data database (and any database created using SnoVault) can be made fully available to the semantic web.

Front end. One of the advantages of using an object store where the documents are JSON objects, is that JSON (JavaScript Object Notation) is the native format used in JavaScript front end frameworks. No translations need to be made between the backend data structures and the front end. We have implemented the web front end to encoded in ReactJS (<http://facebook.github.io/react/>) using JSX (<http://jsx.github.io>). ReactJS is a clean and efficient library for user interface programming in JavaScript, and JSX allows us to use XML-like stanzas in place of an HTML or XML templating system. We further optimized the performance of ReactJS by compiling (from JavaScript to XHTML) the pages on the server using a node (<https://nodejs.org>) process (Figs 2 and 3). Most object types have useful landing pages in HTML that provide a human-friendly version of the raw JSON data, while collections of objects are accessed via the faceted search interface (<https://www.encodeproject.org/search/>). We configure the facets for each object type in the JSON schema objects. These are then translated into Elasticsearch aggregations providing rapid filtering and windowing through data. We have recently implemented an Experiment matrix view (<https://www.encodeproject.org/matrix/?type=Experiment> and Fig 5) and a spreadsheet report view (<https://www.encodeproject.org/report/?type=Experiment> Fig 6) within the search, expanding the modes in which our users can process our data. Individual experiments or sets of experiments with browser-viewable files can be directly viewed at the UCSC genome browser[34] via a track hub[32] created on the fly (Visualize Button), and we provide a metadata driven method for downloading all data files associated with a specific search using the “Download” button. From searching, one can drill down to a specific object page to get detailed information on all the metadata collected for an object.

Experiment Matrix - ENCO x
<https://www.encodeproject.org/matrix/?type=Experiment>

ENCODE Data Methods About Help Search...

Experiment Matrix

Click or enter search terms to filter the experiments included in the matrix.

Enter search term(s)

A.

Assay	Assay category	Target of assay
ChIP-seq 5054	DNA binding 5054	histone 2673
RNA-seq 1359	Transcription 2536	histone modification 2673
DNase-seq 729	DNA accessibility 768	broad histone mark 1451
shRNA/RNA-seq 381	DNA methylation 672	transcription factor 1450
RNA microarray 293	RNA binding 334	control 1136

+ See more... + See more... + See more...

B.

Project	Count
ENCODE	5664
Roadmap	3128
modENCODE	874

Organism

<i>Homo sapiens</i>	7236
<i>Mus musculus</i>	1425
<i>Drosophila melanogaster</i>	481
<i>Caenorhabditis elegans</i>	368
<i>Drosophila pseudoobscura</i>	10

+ See more...

Biosample type

immortalized cell line	3344
tissue	2712
primary cell	1773
stem cell	616
whole organisms	550

+ See more...

Organ

brain	613
skin of body	285
lung	238
muscle organ	237
kidney	235

+ See more...

ASSAY

9666 results

BIOSAMPLE	ChIP-seq	RNA-seq	DNase-seq	shRNA/RNA-seq	RNA microarray	RRBS	WGBS	eCLIP	CAGE	DNAse array	RNA Bind-n-Seq	Replic-Seq	single cell RNA-seq	genotyping array	MRE-seq	RAMPAGE	Replic-chip	siRNA/RNA-seq	MeDIP-seq	microRNA-seq	...and 16 more	
immortalized cell line																						
K562	331	32	16	191	12	1	98	9	2	6	1	2	1					1	1			
HepG2	146	13	3	190	7	2	1	66	6	2	6	1	1									
GM12878	172	27	2		8	2	2	6	2	6	11	1	2	1						1	1	
HeLa-S3	102	11	4		4	1		6	1	6	1	1			1							
A549	110	19	2		2	1		3	1													
...and 165 more																						
tissue																						
liver	113	27	4		1	9		1			3	3	2								5	
heart	82	19	17		8	1	7								2						4	
lung	79	13	14		3	3	6		2					1							2	
kidney	66	13	14		1	5	4		2												2	
stomach	60	16	17			1	8						1		1						2	
...and 141 more																						
primary cell																						
CD4-positive helper T cell	77	4	8		1	3						1	3								3	
common myeloid progenitor, CD34-positive	68	1	15		5	9	1															
trophoblast cell	58	5	2			2																
endothelial cell of umbilical vein	35	9	2		2			5	1	6		1										
Purkinje cell	3														58							
...and 141 more																						
stem cell																						
H1-hESC	188	21	3		1	8	4		4	2			2	4		1		1	1			
H9	60	1	2		1	1	2						1	1		1						
neuronal stem cell	43																					
mesenchymal stem cell	37		2																			
ES-E14	25	1	1																			
...and 35 more																						
whole organisms																						
multi-cellular organism	374	123							15													
carcass		16							4													
in vitro differentiated cells																						
mesendoderm	42	4	2				2															
mesenchymal stem cell	32	4					2	2													1	
chondrocyte	32						3															
neural cell	21	5																			1	
myocyte	26	1																				
...and 39 more																						
induced pluripotent stem cell line																						
iPS-20b	17						1	5														
induced pluripotent stem cell	3	3	5		4																1	2
iPS DF 19.11	4	2	2		1		2						1									
iPS-15b	8						3															
iPS-18c	7						1	3														
...and 12 more																						

Fig 5. The ENCODE portal data matrix. This page shows all experiments released at the ENCODE Portal, including the Roadmap for Epigenomic Mapping Consortium (REMC). Experiments are organized by their biosample (tissue, cell or cell line) on the Y axis and by Assay type on the X-axis (A) Facets select specific properties, such as target (histone, transcription factor) or experimental type (ChIP-seq, RNA-seq, etc.) (B) Facets apply specifically to the biosample, including organism (human, mouse, fly, worm), type (tissue, immortalized cell line, stem cell, etc.) or organ system (as inferred from ontological relations).

<https://doi.org/10.1371/journal.pone.0175310.g005>

Using Github and the cloud for rapid development and deployment. To enable rapid software development and robust deployment, the ENCODE DCC has implemented cloud dev-ops (development operations) code directly into the encodedD system. Using the cloud effectively is a combination of a few small python and well-defined coding practices. Our coding practices start with git, and github. Git is an open-source distributed version control software that is the state of the art in current software engineering. We manage bug requests and feature requests (“Issues”) in tracking software, where each issue is assigned a number and a developer. The developer creates a git branch specific to this issue and works on his or her local version of the software (including a small text database or “fixture”). The branch code is pushed daily to github so that the team is always kept up to date. Github is a web-based Git

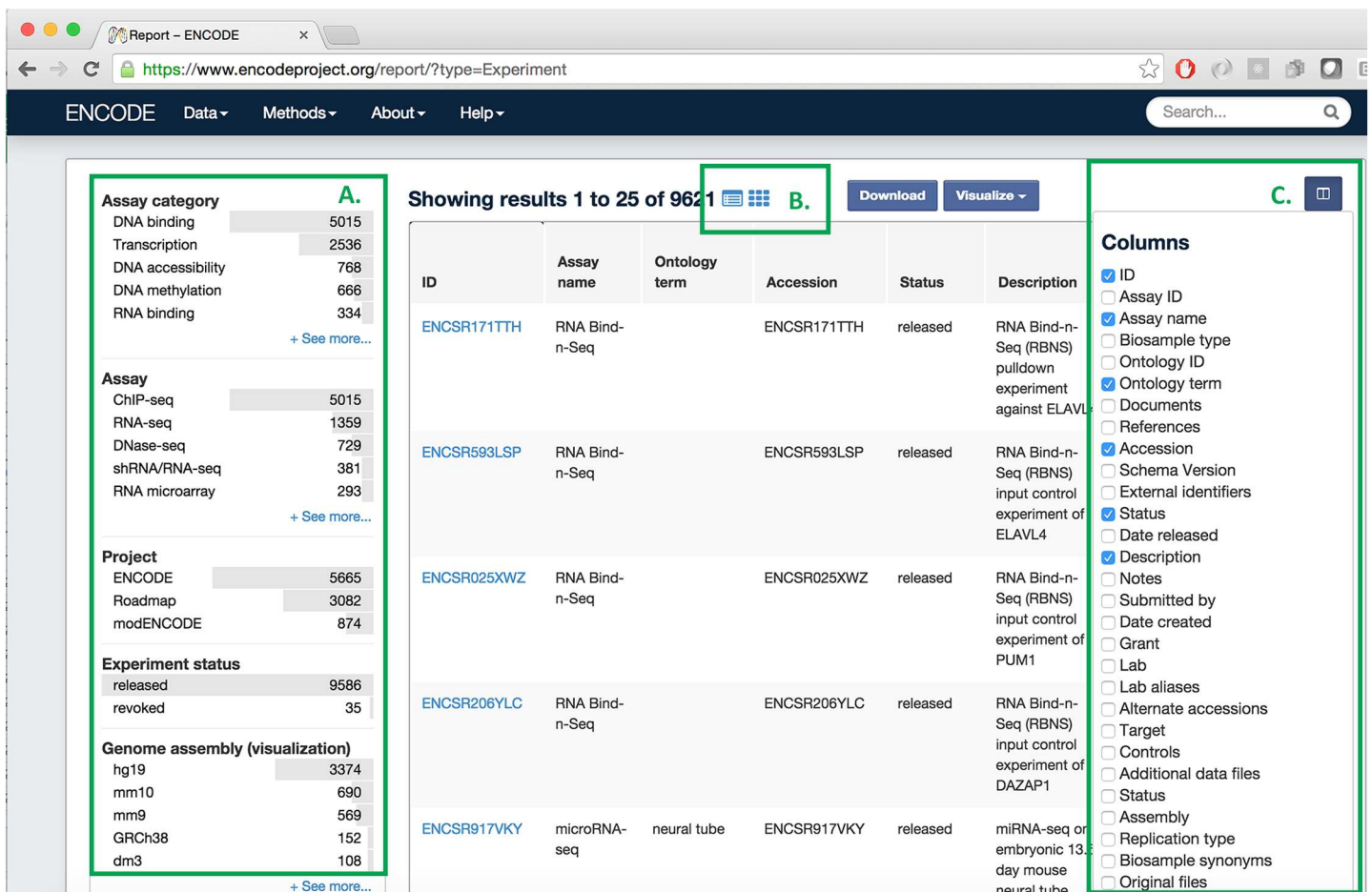


Fig 6. The ENCODE portal report page. This page shows where metadata can be downloaded in spreadsheet format (csv). Report views exist for all collection searches, for example, Experiment, Biosample or Antibody. (A) Standard ENCODE Assay facets to filter the rows of interest (in this case Experiments). (B) Toggle between the report, matrix, and standard search output. (C) Select the columns (individual properties) that will appear in the columns of your spreadsheet.

<https://doi.org/10.1371/journal.pone.0175310.g006>

repository hosting service, which is free to open-source projects. We have written hooks to a web-based continuous integration service called travis-ci (<http://travis-ci.com>), such that when any branch is pushed to github, our testing suite is initiated on a virtual machine hosted by travis-ci. This service is also free for open-source projects. Every 2–3 weeks, completed features are collected into a “release candidate”, merged together, put through a manual Quality Assurance (QA) process, and released to production. Release candidates or even specific feature branches that require full production databases or deeper QA, can be installed on specific demonstration machines via a script that creates and “names” a virtual machine in the Amazon Cloud (AWS).

We use PostgreSQL Write-Ahead Logging (WAL;<http://www.postgresql.org/docs/9.4/static/wal-intro.html>) and the software package WAL-E (<https://github.com/wal-e/wal-e>) to maintain a consistent database across our production, test, and various “feature-demo” instances created by engineers and data wranglers for testing and QA. A simple flag to the deploy script configures each instance to be a production candidate, test (global testing) or demo (local, in-house testing machine), and they immediately create a local copy of the PostgreSQL database via the WAL logs. The difference between a production instance and any other instance, is that the production instance acts as the PostgreSQL “master” and is the only machine authorized to *write* to the WAL logs or to upload files to the production Amazon S3 bucket.

When the ENCODE DCC project moved to Stanford University in 2012, we had to develop a software system that could handle all existing data and metadata from ENCODE, incorporate distinct and additional metadata from modENCODE and REMC, as well as all future ENCODE metadata. In many cases, the experimental assays, reagents, and protocols were not yet fully defined when we began our work on encodeD, while experimental data was being continuously created from existing ENCODE production labs. In short we had to develop a system flexible enough to handle nearly arbitrary experimental definitions within the field of genomics and epigenomics, yet still maintain strong data integrity and control the input specification to preserve univocity in our data and metadata descriptions. Our schema changes with nearly every software release, and the software is released in place every 3–4 weeks, with almost zero disruption to submitting or viewing users.

Using encodeD or SnoVault in your project

Open source repository. All of the source code created by the ENCODE DCC is available from GitHub: <http://github.com/ENCODE-DCC/>. SnoVault is a completely generic hybrid object database with Elasticsearch interface. It is completely data agnostic and could be used for any web database purpose. Our system used for the ENCODE DCC metadata database is at <http://github.com/ENCODE-DCC/encodeD>, which could be adapted to other projects that specific store genomic data and metadata. Scripts and python modules specific to retrieving and submitting metadata to the ENCODE DCC can be found here: <https://github.com/ENCODE-DCC/pyencoded-tools>; we plan on adding a command-line interface in the future.

PyPi package. We have extracted the main Pyramid/JSON-SCHEMA/JSON-LD/Elasticsearch framework to a separate repository, and created a python package called SnoVault (<http://github.com/ENCODE-DCC/snovault>). SnoVault is available from PyPi or other python distribution systems. The main object-relational backend can be used for any project regardless of content. Creating application specific code starts with modeling your data in JSON using JSON-SCHEMA and JSON-LD. A small amount of custom code is required to get the bare bones application running for item and collection pages, and from there frontend components using ReactJS or any other system can be used to give your application a unique

user experience. We have used the encodeD system to create a curation application and portal for the ClinGen project (<http://clinicalgenome.org>; <http://github.com/ClinGen/clincodeD>), and it is currently being evaluated by other genomics Data Coordination Centers at Washington University (Saint Louis) and Harvard Medical School.

Acknowledgments

We wish to thank all participants in the ENCODE Consortium for our collaborations that have enhanced the metadata definitions, as well as ENCODE Portal users for their useful feedback.

Author Contributions

Conceptualization: BCH LDR ELH.

Data curation: VSM ETC JMD IG AKN KCO JH MCH CAS JSS ELH.

Funding acquisition: JMC.

Methodology: BCH ELH.

Project administration: ELH BCH CAS.

Software: LDR NRP DIG UKB TRD FYT BCH SRM BTL IG.

Supervision: JMC BCH ELH CAS.

Validation: BTL.

Visualization: DIG FYT BCH.

Writing – original draft: BCH.

Writing – review & editing: JH.

References

1. International Human Genome Sequencing Consortium. Finishing the euchromatic sequence of the human genome. *Nature*. 2004; 431: 931–945. <https://doi.org/10.1038/nature03001> PMID: 15496913
2. Rosenbloom KR, Sloan CA, Malladi VS, Dreszer TR, Learned K, Kirkup VM, et al. ENCODE data in the UCSC Genome Browser: year 5 update. *Nucleic Acids Res*. 2013; 41: D56–63. <https://doi.org/10.1093/nar/gks1172> PMID: 23193274
3. Ho JWK, Jung YL, Liu T, Alver BH, Lee S, Ikegami K, et al. Comparative analysis of metazoan chromatin organization. *Nature*. 2014; 512: 449–452. <https://doi.org/10.1038/nature13415> PMID: 25164756
4. Boyle AP, Araya CL, Brdlik C, Cayting P, Cheng C, Cheng Y, et al. Comparative analysis of regulatory information and circuits across distant species. *Nature*. 2014; 512: 453–456. <https://doi.org/10.1038/nature13668> PMID: 25164757
5. Gerstein MB, Rozowsky J, Yan K-K, Wang D, Cheng C, Brown JB, et al. Comparative analysis of the transcriptome across distant species. *Nature*. 2014; 512: 445–448. <https://doi.org/10.1038/nature13424> PMID: 25164755
6. Yue F, Cheng Y, Breschi A, Vierstra J, Wu W, Ryba T, et al. A comparative encyclopedia of DNA elements in the mouse genome. *Nature*. 2014; 515: 355–364. <https://doi.org/10.1038/nature13992> PMID: 25409824
7. Lander ES, Linton LM, Birren B, Nusbaum C, Zody MC, Baldwin J, et al. Initial sequencing and analysis of the human genome. *Nature*. 2001; 409: 860–921. <https://doi.org/10.1038/35057062> PMID: 11237011
8. Malladi VS, Erickson DT, Podduturi NR, Rowe LD, Chan ET, Davidson JM, et al. Ontology application and use at the ENCODE DCC. *Database*. 2015; 2015.

9. Sloan CA, Chan ET, Davidson JM, Malladi VS, Strattan JS, Hitz BC, et al. ENCODE data at the ENCODE portal. *Nucleic Acids Res.* 2016; 44: D726–32. <https://doi.org/10.1093/nar/gkv1160> PMID: [26527727](https://pubmed.ncbi.nlm.nih.gov/26527727/)
10. Hong EL, Sloan CA, Chan ET, Davidson JM, Malladi VS, Strattan JS, et al. Principles of metadata organization at the ENCODE data coordination center. *Database.* 2016; 2016.
11. ENCODE Project Consortium. A user's guide to the encyclopedia of DNA elements (ENCODE). *PLoS Biol.* 2011; 9: e1001046. <https://doi.org/10.1371/journal.pbio.1001046> PMID: [21526222](https://pubmed.ncbi.nlm.nih.gov/21526222/)
12. Pazin MJ. Using the ENCODE Resource for Functional Annotation of Genetic Variants. *Cold Spring Harb Protoc.* 2015; 2015: 522–536. <https://doi.org/10.1101/pdb.top084988> PMID: [25762420](https://pubmed.ncbi.nlm.nih.gov/25762420/)
13. Diehl AG, Boyle AP. Deciphering ENCODE. *Trends Genet.* 2016;
14. Harrow J, Frankish A, Gonzalez JM, Tapanari E, Diekhans M, Kokocinski F, et al. GENCODE: the reference human genome annotation for The ENCODE Project. *Genome Res.* 2012; 22: 1760–1774. <https://doi.org/10.1101/gr.135350.111> PMID: [22955987](https://pubmed.ncbi.nlm.nih.gov/22955987/)
15. Haase K, Mösch A, Frishman D. Differential expression analysis of human endogenous retroviruses based on ENCODE RNA-seq data. *BMC Med Genomics.* 2015; 8: 71. <https://doi.org/10.1186/s12920-015-0146-5> PMID: [26530187](https://pubmed.ncbi.nlm.nih.gov/26530187/)
16. Vierstra J, Stamatoyannopoulos JA. Genomic footprinting. *Nat Methods.* 2016; 13: 213–221. <https://doi.org/10.1038/nmeth.3768> PMID: [26914205](https://pubmed.ncbi.nlm.nih.gov/26914205/)
17. Siggens L, Ekwall K. Epigenetics, chromatin and genome organization: recent advances from the ENCODE project. *J Intern Med.* 2014; 276: 201–214. <https://doi.org/10.1111/joim.12231> PMID: [24605849](https://pubmed.ncbi.nlm.nih.gov/24605849/)
18. Ma Y, Smith CE, Lai C-Q, Irvin MR, Parnell LD, Lee Y-C, et al. The effects of omega-3 polyunsaturated fatty acids and genetic variants on methylation levels of the interleukin-6 gene promoter. *Mol Nutr Food Res.* 2016; 60: 410–419. <https://doi.org/10.1002/mnfr.201500436> PMID: [26518637](https://pubmed.ncbi.nlm.nih.gov/26518637/)
19. Fan Y, Zhang Y, Xu S, Kong N, Zhou Y, Ren Z, et al. Insights from ENCODE on Missing Proteins: Why β -Defensin Expression Is Scarcely Detected. *J Proteome Res.* 2015; 14: 3635–3644. <https://doi.org/10.1021/acs.jproteome.5b00565> PMID: [26258396](https://pubmed.ncbi.nlm.nih.gov/26258396/)
20. Li G, Fullwood MJ, Xu H, Mulawadi FH, Velkov S, Vega V, et al. ChIA-PET tool for comprehensive chromatin interaction analysis with paired-end tag sequencing. *Genome Biol.* 2010; 11: R22. <https://doi.org/10.1186/gb-2010-11-2-r22> PMID: [20181287](https://pubmed.ncbi.nlm.nih.gov/20181287/)
21. Lieberman-Aiden E, van Berkum NL, Williams L, Imaikaev M, Ragoczy T, Telling A, et al. Comprehensive mapping of long-range interactions reveals folding principles of the human genome. *Science.* 2009; 326: 289–293. <https://doi.org/10.1126/science.1181369> PMID: [19815776](https://pubmed.ncbi.nlm.nih.gov/19815776/)
22. Zhao J, Ohsumi TK, Kung JT, Ogawa Y, Grau DJ, Sarma K, et al. Genome-wide identification of polycomb-associated RNAs by RIP-seq. *Mol Cell.* 2010; 40: 939–953. <https://doi.org/10.1016/j.molcel.2010.12.011> PMID: [21172659](https://pubmed.ncbi.nlm.nih.gov/21172659/)
23. Hoffman MM, Ernst J, Wilder SP, Kundaje A, Harris RS, Libbrecht M, et al. Integrative annotation of chromatin elements from ENCODE data. *Nucleic Acids Res.* 2013; 41: 827–841. <https://doi.org/10.1093/nar/gks1284> PMID: [23221638](https://pubmed.ncbi.nlm.nih.gov/23221638/)
24. Yip KY, Cheng C, Bhardwaj N, Brown JB, Leng J, Kundaje A, et al. Classification of human genomic regions based on experimentally determined binding sites of more than 100 transcription-related factors. *Genome Biol.* 2012; 13: R48. <https://doi.org/10.1186/gb-2012-13-9-r48> PMID: [22950945](https://pubmed.ncbi.nlm.nih.gov/22950945/)
25. Gerstein MB, Kundaje A, Hariharan M, Landt SG, Yan K-K, Cheng C, et al. Architecture of the human regulatory network derived from ENCODE data. *Nature.* 2012; 489: 91–100. <https://doi.org/10.1038/nature11245> PMID: [22955619](https://pubmed.ncbi.nlm.nih.gov/22955619/)
26. Cheng Y, Ma Z, Kim B-H, Wu W, Cayting P, Boyle AP, et al. Principles of regulatory information conservation between mouse and human. *Nature.* 2014; 515: 371–375. <https://doi.org/10.1038/nature13985> PMID: [25409826](https://pubmed.ncbi.nlm.nih.gov/25409826/)
27. Boley N, Nathan B, Wan KH, Bickel PJ, Celniker SE. Navigating and mining modENCODE data. *Methods.* 2014; 68: 38–47. <https://doi.org/10.1016/j.ymeth.2014.03.007> PMID: [24636835](https://pubmed.ncbi.nlm.nih.gov/24636835/)
28. Brown JB, Celniker SE. Lessons from modENCODE. *Annu Rev Genomics Hum Genet.* 2015; 16: 31–53. <https://doi.org/10.1146/annurev-genom-090413-025448> PMID: [26133010](https://pubmed.ncbi.nlm.nih.gov/26133010/)
29. ENCODE Project Consortium. An integrated encyclopedia of DNA elements in the human genome. *Nature.* 2012; 489: 57–74. <https://doi.org/10.1038/nature11247> PMID: [22955616](https://pubmed.ncbi.nlm.nih.gov/22955616/)
30. Mouse ENCODE Consortium, Stamatoyannopoulos JA, Snyder M, Hardison R, Ren B, Gingeras T, et al. An encyclopedia of mouse DNA elements (Mouse ENCODE). *Genome Biol.* 2012; 13: 418. <https://doi.org/10.1186/gb-2012-13-8-418> PMID: [22889292](https://pubmed.ncbi.nlm.nih.gov/22889292/)

31. Roadmap Epigenomics Consortium, Kundaje A, Meuleman W, Ernst J, Bilenky M, Yen A, et al. Integrative analysis of 111 reference human epigenomes. *Nature*. 2015; 518: 317–330. <https://doi.org/10.1038/nature14248> PMID: 25693563
32. Raney BJ, Dreszer TR, Barber GP, Clawson H, Fujita PA, Wang T, et al. Track data hubs enable visualization of user-defined genome-wide annotations on the UCSC Genome Browser. *Bioinformatics*. 2014; 30: 1003–1005. <https://doi.org/10.1093/bioinformatics/btt637> PMID: 24227676
33. Jupp S, Malone J, Bolleman J, Brandizi M, Davies M, Garcia L, et al. The EBI RDF platform: linked open data for the life sciences. *Bioinformatics*. 2014; 30: 1338–1339. <https://doi.org/10.1093/bioinformatics/btt765> PMID: 24413672
34. Speir ML, Zweig AS, Rosenbloom KR, Raney BJ, Paten B, Nejad P, et al. The UCSC Genome Browser database: 2016 update. *Nucleic Acids Res*. 2016; 44: D717–25. <https://doi.org/10.1093/nar/gkv1275> PMID: 26590259