

Database and ontologies

Stitcher: An entity resolution framework for comprehensive data integration of approved drugs

Dac-Trung Nguyen,^{1,*} Ivan Grishagin,¹ Daniel Katzel,¹ Tyler Peryea,^{1,2} Ajit Jadhav,¹ and Noel Southall^{1,*}

¹Division of Pre-clinical Innovation, National Center for Advancing Translational Sciences (NCATS), National Institutes of Health, USA

²Present address: Office of Health Informatics, Office of Chief Scientist, Food and Drug Administration, USA

*To whom correspondence should be addressed.

Associate Editor: XXXXXXXX

Received on XXXXX; revised on XXXXX; accepted on XXXXX

Abstract

Motivation: As biomedical data continues to grow at a rapid rate, the need to integrate diverse data across scales and modalities has never been more urgent. In the context of drug data, the data integration challenges are considerably more complex due to the lack of clarity around the term “drug.” What is considered as “drug” can, therefore, vary significantly between data sources. The data integration challenges and the lack of a comprehensive resource of drugs that been marketed or approved in the United States for human use are the primary goals behind our current work.

Results: Through the combination of a reference data source (G-SRS) and entity resolution strategies tailored specifically for drug entity type, we have developed a comprehensive data integration pipeline (**Stitcher**) for drug data. The resource <https://drugs.ncats.io> is a comprehensive drug resource that is supported by **Stitcher**.

Availability: The complete source code along with data and build instructions for **Stitcher** is readily available on Github <https://github.com/ncats/stitcher>. The latest version of the stitch multigraph is available as a Neo4j database at <https://stitcher.ncats.io/browser> (use [stitcher.ncats.io:80](https://drugs.ncats.io) as the hostname). And the **InXight Drugs** resource is accessible at <https://drugs.ncats.io>.

Contact: southalln@mail.nih.gov

Supplementary information: Supplementary data are available at *Bioinformatics* online.

1 Introduction

As the volume of biological data continues to grow at an unprecedented rate, data de-duplication—also commonly known as record linkage or *entity resolution*—is proportionally playing a prominent role in data integration. From the construction of training data for machine learning to building knowledge graphs as epistemological frameworks for artificial intelligence, proper entity resolution is essential in generating ground-truth data. The core challenge of entity resolution is in establishing *uniqueness*. For well-defined entity types (e.g., gene, tissue, cell line), uniqueness is determined solely based on established identifiers and nomenclature; for other entity types (e.g., drug, disease, phenotype), however, uniqueness

is not as well-established due to conceptual ambiguities in how entities are defined and represented. Take the disease entity type as an example. The discrepancy between the theoretical concept of “disease entity” from its clinical nosology (Hucklenbroich, 2014) is what makes disease entity resolution extremely challenging.

Herein we report on our recent data integration effort to build a comprehensive resource of drugs that have either been marketed or approved in the United States for human use. Such a resource is not only instrumental for drug repurposing but also serves as a valuable tool to further our understanding of the mechanistic properties of molecular targets (Huang *et al.*, 2011, 2019). To the best of our knowledge, **InXight Drugs** is currently the most comprehensive resource of its kind. In the remainder of this paper, we discuss data integration challenges associated with drug data, conceptually as well as technically. This discussion serves

as the backdrop for the development of **Stitcher**, an entity resolution framework that we have developed to address the shortcomings of traditional approaches.

1.1 What is a “drug”?

While the word is included within the name of the organization, the U.S. Food and Drug Administration (FDA) does not have a straightforward definition of the word “drug.” The Federal Food Drug and Cosmetic Act (FD&C Act) and FDA regulations define the term drug, in part, by reference to its intended use, as “articles intended for use in the diagnosis, cure, mitigation, treatment, or prevention of disease and articles (other than food) intended to affect the structure or any function of the body of man or other animals.” (FDA, 2020a) More practically, the agency defines “drug substance” and “drug product” respectively as the physical ingredients found in marketed products. Others use the word “drug” to sometimes refer to “drug substances” and sometimes to “drug products” as convenient, and this causes a great deal of semantic confusion within drug data found on the web. The National Library of Medicine produces a semantic product, RxNorm, that provides a variety of precise semantic types for ingredients, tradenames, dose forms, semantic clinical drug components, semantic clinical drug forms, and semantic clinical drugs which facilitate working with drug data, but its terminology is unfortunately limited to commonly used prescription drugs, “clinically significant ingredients,” and adoption of this complex semantic scheme is limited (NLM, 2005).

There is a third definition of the word drug that is commonly used in the literature and used by the FDA when it refers to an active moiety and a new molecular entity. In this case, ingredients whose pharmacological effect occurs through the same molecular entity are considered the same drug. This holds for different salt forms such as sumatriptan succinate and sumatriptan hemisulfate, but it also holds for prodrugs and their metabolized active forms such as brincidofovir and cidofovir (FDA, 2020b). *An active moiety is a molecule or ion, excluding those appended portions of the molecule that cause the drug to be an ester, salt (including a salt with hydrogen or coordination bonds), or other noncovalent derivative (such as a complex, chelate, or clathrate) of the molecule, responsible for the physiological or pharmacological action of the drug substance* (FDA, 2012). Under the Food and Drug Administration Amendments Act of 2007, all newly introduced active moieties must first be reviewed by an advisory committee before the FDA can approve these products.

As in other information domains, the names used to refer to drug substances and products are particularly problematic because their definitions change as a function of location or jurisdiction, time and context. FDA and other national regulators of medicines have collaborated to produce ISO 11238 (ISO, 2018) which endeavors to define an information scheme for the unambiguous identification of all ingredients found in medicinal products, and FDA uses an implementation of ISO 11238 as the backbone of its information systems within the agency (FDA, 2019a). While this facilitates data exchange within the FDA and with other national authorities, the task still remains to be able to map other, external data sources into this rigorously-defined scheme using whatever names and data are at hand.

2 Approach

2.1 Preliminary concepts

The conceptual data model underlying **Stitcher** is a *multigraph*. Within this multigraph, a node can either be a *stitch node* or *data node*. Each data node represents a “raw” entity as ingested from the data source; its corresponding stitch node is a *standardized* representation that is used for *stitching*. An edge between two stitch nodes can either be a *stitch key*

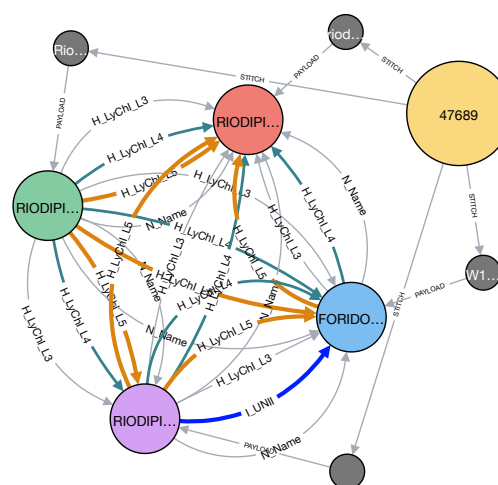


Fig. 1. A connected component in the stitch multigraph with four stitch nodes (medium) and corresponding data nodes (small). Each stitch value forms a clique within this connected component. The edge labels between stitch nodes are the stitch keys. The large node is the derived entity (i.e., sgroup node) generated from entity resolution.

(undirected) or *relationship* (directed). A unique *stitch value* is associated with each stitch key such that it forms a clique. Figure 1 shows an instance of a connected component of a stitch multigraph with overlapping cliques.

A connected component in the stitch multigraph represents the basic unit of work for entity resolution. While the majority of connected components are of reasonable sizes (e.g., 20 to 50 stitch nodes), the real challenges center around effective strategies for handling very large connected components—or also commonly known as *hairballs* (Croset et al., 2016). For example, the current version of the **InXight Drugs** resource has an hairball close to 30,000 stitch nodes spanning across 15 data sources. We discuss our strategies in detail for untangling through such an hairball in Section 3.3.

The primary goal of entity resolution is to determine the number of unique entities in a connected component. These derived entities are represented as *sgroup nodes* in the stitch multigraph. There can be multiple instances of sgroup nodes for any given set of stitch nodes, with each instance reflects a specific algorithmic strategy or version. Figure 1 shows that there is only one unique entity as determined by the entity resolution algorithm for the given connected component.

2.2 Stitch keys

Stitch key is a core concept in **Stitcher**. It defines how entities are matched, which, in turn, determines how cliques and connected components are formed. By virtue of its importance, the stitch key should reflect the true identity of the entity as much as possible. Depending on the entity type, the stitch key can be generic (e.g., synonym) or very specific (e.g., molecular hash key). For drug entity type, **Stitcher** relies on the following stitch keys for each entity:

N_Name. This is the most generic stitch key available. Stitch values associated with this stitch key can be any established names or nomenclature; e.g., tradenames, INN (International Nonproprietary Names), USAN (United States Adopted Names), IUPAC (International Union of Pure and Applied Chemistry).

I_UNII, I_CAS, I_CID, I_CODE. These stitch keys represent (i) unique identifiers assigned to the entity by a well-known registrar (e.g., the U.S. Food and Drug Administration in the case of UNII) or (ii) internal

company code. `I_UNII`, `I_CAS`, and `I_CID` are specific to drug (or substance in general) entity type, whereas `I_CODE` can be used for any type of identifiers. The decision to use specific stitch keys over generic ones ultimately rests on the strategies used for entity resolution.

`H_LyChI_L5`, `H_LyChI_L4`, `H_LyChI_L3`. For the small molecule class of drugs, perhaps more important than any identifiers is the underlying chemical structure definition. These stitch keys are hash values derived from the molecular structure at different resolutions (NCATS, 2013). Section 3.1 discusses in detail how these derived stitch values are generated.

`R_activeMoiety`. Technically not a stitch key, the active moiety relationship between two drugs provides a strong evidence of equivalence. While this relationship can be inferred directly from the chemical structures (e.g., freebase and salt forms, with and without esters), there is some level of curation needed to handle structures with metal complex.

Table 1 shows an example of stitch keys and stitch values for the drug entity *imatinib mesylate*. In this example, the `R_activeMoiety` relationship specifies the UNII of the freebase form of imatinib mesylate.

Table 1. Stitch keys and stitch values for the drug imatinib mesylate

Stitch key	Stitch value
<code>N_Name</code>	IMATINIB MESYLATE; GLEEVEC; GLIVEC
<code>I_UNII</code>	8A1O1M485B
<code>I_CAS</code>	220127-57-1
<code>I_CID</code>	5291
<code>I_CODE</code>	STI-571; CHEMBL941
<code>H_LYCHI_L5</code>	7S4GKGNQ6N3X-N
<code>H_LYCHI_L4</code>	VLU17BQBSGWU-N; K83X3L3XSSHK-S
<code>H_LYCHI_L3</code>	VL3FPUQ59CU-N; K846NBMB7T3-S
<code>R_activeMoiety</code>	BKJ8M8G5HI

2.3 Data sources

Stitcher utilizes a number of diverse data sources for the **InXight Drugs** resource. Among the data sources, of particular importance is the public G-SRS data source from the FDA (FDA, 2019b). This data source is well-curated and contains over 100K substances across six different classes: chemical, structurally diverse, protein, mixture, polymer, and nucleic acid. As a data source derived from the FDA’s internal substance registry system (FDA, 2019a), the G-SRS data source naturally forms the basis of our data integration effort. Using this data source as the “seed” from which other data sources can map onto has the following benefits:

- Since the G-SRS data source implements the ISO 11238 standard (ISO, 2018) for defining medicinal substances, it serves as an ideal starting point for what constitutes a “drug.”
- The data is a public version of the internal substance registry within the FDA; as such, it is well-curated and up-to-date.

The complete list of data sources currently used by **Stitcher** is shown in Table 2.

2.4 Overall strategy

The basic premise behind **Stitcher** is that data integration is often done within the context of a specific data source. This is a reasonable assumption given the data quality varies when integrating across disparate sources.

Table 2. Data sources used in the current version of Stitcher.

Data source	Size
G-SRS, April 2019	105,019
Withdrawn and Shortage Drugs List Feb 2018	674
Broad Institute Drug List 2018-09-07	6,125
NCATS Pharmaceutical Collection, April 2012	14,814
Rancho BioSciences, March 2019	51,591
Pharmaceutical Manufacturing Encyclopedia (Third Edition)	2,268
DailyMed Rx, January 2019	74,850
DrugBank, December 2018	11,922
DailyMed Other, January 2019	13,393
DailyMed OTC, January 2019	79,448
DrugsFDA & Orange Book, July 2019	28,256
ClinicalTrials, December 2017	305,833
OTC Monographs, December 2018	2,713
FDA NADA and ANADAs, December 2018	554
FDA Excipients, December 2018	10,212

Furthermore, by establishing a reference data source for data integration, we have finer control over the following:

Data quality. A reference data source is typically selected such that it is of high quality. Here, we can also impose other data quality constraints (e.g., no synonyms can span multiple entities) to guide entity resolution.

Data resolution. Entity resolution is particularly challenging when data integration involves ontologies. A reference data source can serve as the anchor ontology from which other ontologies can be mapped. As with data quality, we can also impose any additional semantic constraints; e.g., prostate cancer is not one of the diagnoses for a female patient in an electronic health record.

Data curation. Generating ground-truth data is more manageable with a single data source than across multiple data sources. This is particularly important due to the iterative feedback between data curation and data integration.

The G-SRS data source serves as an ideal reference data source. Its rigorous substance models and well-structured data elements give us a good starting point for drug data integration. In the next section, we discuss our strategies in utilizing the G-SRS reference data source to address entity resolution for drug data.

3 Methods

In general, data integration with **Stitcher** consists of four basic steps: *ingestion*, *stitching*, *entity resolution*, and *entity normalization*. With the exception of *entity resolution*, all other steps—as they are currently implemented in **Stitcher**—are generic and can be applied to a wide range of entity types.

3.1 Data ingestion

Stitcher is capable of ingesting data in a wide variety of sources and formats. Semantic formats such as OWL, RDF, and Turtle are supported as are JSON, delimiter separated text, and custom formats. For non-semantic format, a separate configuration file is required to map properties to stitch keys.

An important step in data ingestion is the standardization and validation of stitch values. For `N_Name` stitch key, the standardization procedure is simply to convert the input string to uppercase; no validation is performed.

For `I_UNII` and `I_CAS` stitch keys, no standardization is required, and validation is a simple checksum calculation to ensure the stitch value is proper. Depending on the input format, **Stitcher** also provides basic utilities (e.g., regular expression) to help with data transformation during ingestion.

Perhaps the most unique feature of **Stitcher** is its ability to incorporate knowledge of chemical structures into entity resolution. Whereas traditional approaches rely on names and identifiers to determine equivalence substances, **Stitcher** goes a step further and utilizes the underlying chemical structures to infer equivalence. This is particularly relevant when the drug is a mixture, prodrug, or active moiety with complex excipient (or derivative thereof). As an example, consider the drug entity *IMATINIB MESYLATE* and its active ingredient *IMATINIB*. Here, it is obvious that the two entities cannot be matched by name alone. Instead, having structural information by way of molecular hash keys for each molecular component allows us to determine equivalence from the common active moiety *IMATINIB* between the two entities. This trivial example might suggest that, instead of comparing names exactly, we find the longest common substring of the names. The approach would certainly work in this example, but to make it work in general would require very specialized parsing rules and dictionaries.

For data sources with chemical structures, the most computationally demanding step in data ingestion is the generation of molecular hash keys. Hash keys are generated for each component of a chemical structure in three different structural levels: L5, L4, and L3, which correspond to stitch keys `H_LYCHI_L5`, `H_LYCHI_L4`, and `H_LYCHI_L3`, respectively. Level L5 is the most specific; it represents the chemical structure as-is, i.e., without structure normalization and standardization. With the exception of the relation `R_activeMoiety`, a match at this level has higher priority over other stitch keys. The next level L4 represents the structure after normalization and standardization per the LyChI software package (NCATS, 2013). A match at this level implies that two structures are equivalent in terms of stereochemistry, resonance, and tautomerism. And the last level L3 is the same as L4 but without stereochemistry. A match at this level is considered weak and does not constitute equivalence without other significant supporting evidence. The purpose for L3 is in anticipation of incorrect or missing stereo information, which is one of the most common type of errors associated with chemical structures. For each hash key, a suffix `-M`, `-S`, or `-N` is also assigned to designate the molecular component as either a *metal*, *salt*, or *neither*, respectively. Table 1 illustrates all three representations for the drug *IMATINIB MESYLATE*. Note that the cardinality for L5 is always one, whereas for L4 and L3 the cardinality is equal to the number of non-hydrate molecular components. (Hydrate components are removed prior to processing.)

3.2 Data stitching

Stitching is the process by which the stitch multigraph is incrementally constructed as data is ingested. Algorithm 1 describes the basic stitching algorithm of **Stitcher**. This algorithm is applied to each data source, and upon its completion produced a stitch multigraph such that any stitch value that spans N stitch nodes is an induced clique, i.e., a complete subgraph of N nodes and $\frac{N(N-1)}{2}$ edges. Overlapping induced cliques form the basis for the proposed entity resolution approach discussed in the next section. As a side-effect, the stitching algorithm also utilizes the union-find algorithm (Cormen et al., 2001) to efficiently track connected components.

3.3 Entity resolution

After all data sources have been stitched together, the next step is to identify *unique* entities from the stitch multigraph. Formally, this step is known as *entity resolution* and is the only step within **Stitcher** that is specific to the drug entity type. (This is to be expected: Given that entity resolution

Algorithm 1: Entity stitching algorithm

```

Let  $W$  denote the set of stitch nodes created in the data ingestion
step for a given data source  $D$ .
Let  $\langle k, v \rangle$  be the tuple of stitch key and value, respectively,
defined for a stitch node  $w$ .
 $G = (V, E)$  is the current stitch multigraph.
 $\text{Find}(k, v)$  is a function that returns all stitch nodes in  $V$ 
containing stitch key  $k$  and stitch value  $v$ .
 $\text{Union}(w, z)$  is a union-find algorithm for tracking disjoint sets
(i.e., connected components).
for  $w \in W$  do
  for  $\langle k_i, v_i \rangle \in w$  do
    for  $z \in \text{Find}(k_i, v_i)$  do
       $E \leftarrow E \cup z \sim w$ 
       $\text{Union}(w, z)$ 
    end
  end
   $V \leftarrow V \cup w$ 
end

```

is about adjudicating the splitting and merging of entities, a reasonable amount of knowledge of the entity type is required for the adjudication to be effective.) For a given connected component, the iterative process of assigning equivalence labels to stitch nodes is known as *untangling*. Algorithm 2 gives a high level outline of the untangling process.

Algorithm 2: An algorithm to untangle a connected component

```

Let  $U$  be the disjoint set data structure for all entities.
Let  $S$  denote the set of unlabeled entities (i.e., singletons).
 $C = (V, E)$  is the connected component.
 $\text{MergeNodes}(U, r)$  is a function that performs transitive
closure on stitch nodes in  $C$  which are connected by a relation
 $r \in E$ . The results are accumulated in  $U$ .
 $\text{MergeCliques}(U, K)$  is a function that takes a set of stitch
keys  $K$ , finds overlapping cliques that span two or more stitch
keys, and performs transitive closure on the entities.
 $\text{MergeSingletons}(U, S, K)$  is a function that also takes in a
set of stitch keys  $K$ , a set of singleton stitch node  $S$ , and find the
best mapping to an already labeled stitch node.
 $\text{MergeNodes}(U, R_{\text{activeMoiety}})$ 
 $\text{MergeNodes}(U, I_{\text{UNII}})$ 
 $\text{MergeNodes}(U, H_{\text{LyChI\_L4}})$ 
 $\text{MergeCliques}(U, N_{\text{Name}}, I_{\text{CAS}}, I_{\text{CID}}, H_{\text{LyChI\_L4}})$ 
 $\text{MergeSingletons}(U, S, N_{\text{Name}}, I_{\text{CAS}})$ 

```

At the core of the algorithm is the implied priorities associated with the stitch keys. The relation `R_activeMoiety` has the highest priority as it is manually generated and is a special relation that only available to the G-SRS data source. As an example, consider the entities *acetylsalicylic acid* (or also commonly known as *aspirin*) and *ethyl acetylsalicylate* shown in Figures 2 and 3, respectively. While the two entities have nothing in common, in G-SRS *acetylsalicylic acid* is annotated as being an *active moiety* of *ethyl acetylsalicylate*. Further examination of the structural differences shows that only an *ester* separates the two entities; this falls well within what the FDA considers as equivalent drugs. This example also highlights a quagmire: Computationally, there is nothing to prevent us from imputing *active moiety* relationships through efficient (sub-) graph isomorphisms. This, however, is a very tempted trap that we have

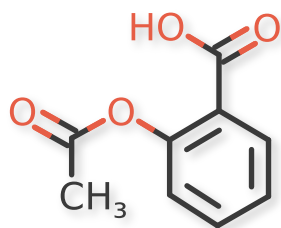


Fig. 2. Chemical structure for acetylsalicylic acid.

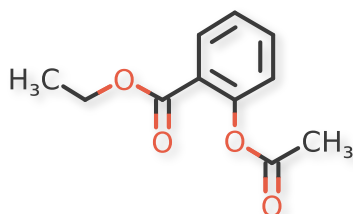


Fig. 3. Chemical structure for ethyl acetylsalicylate.

thus far resisted due to other forms of *active moiety* relationships—e.g., metabolites and metals—that would require considerable investment of effort.

The next priority is the stitch key `I_UNII`. As UNII is the primary identifier for the G-SRS data source, any data source that provides mapping based on this identifier implies that the data source has sufficient knowledge of G-SRS (i.e., guilt by association). For entities that can be represented by chemical structures, the stitch key `H_LyChI_L4` has the next level of priority. The complexity required for two entities to have the same stitch values means that the entities are less likely to match by errors. The rest of the stitch keys (i.e., `N_Name`, `I_CAS`, `I_CID`) all have the lowest priority.

At the completion of Algorithm 2, the disjoint set data structure U contains all equivalence entity classes such that each class is represented by an *sgroup* node in the stitch multigraph. The *sgroup* nodes are the *resolved entities*.

3.4 Entity normalization

The last step in the data integration pipeline is to decide how the resolved entities are defined. This step is referred to as *entity normalization* and its goals are to have (i) clear and consistent strategies for merging properties and (ii) conflict resolution (semantic as well as self consistency). While this step can be quite trivial if the properties are mutually exclusive across all data sources, to address this in a general setting will require considerable efforts in terms of understanding the data source and its metadata. Within the context of the current work, we resort to a simple strategy: When merging properties, we preferentially use those that come from G-SRS with a basic consistency constraint that a synonym is never associated with more than one resolved entity. It also helps that many of the data sources in Table 2 have mutually exclusive properties; e.g., the property “drug approval year” in the Drugs@FDA data source is not a property of G-SRS.

4 Results

Stitcher and the data integration pipeline developed for the **InXight Drugs** resource are available in source form at <https://github.com/ncats/stitcher>. The stitch multigraph built with data sources listed in Table 2 is also available as a Neo4j database at <https://stitcher.ncats.io/browser>. (While no credentials are needed for the database,

the web interface requires that the string `stitcher.ncats.io:80` is entered into the field `Host`.) This database currently contains 192,413 stitch nodes and 11,948,470 edges (relationships and stitch keys). Tables 3 and 4 give a breakdown of the stitch- keys and values, respectively, in the stitch multigraph. All figures used throughout this paper have been generated directly from this database.

Table 3. Distribution of edge types for the stitch multigraph.

Stitch key	Size
<code>H_LyChI_L3</code>	6,140,942
<code>H_LyChI_L4</code>	5,300,078
<code>N_Name</code>	177,446
<code>H_LyChI_L5</code>	162,160
<code>I_UNII</code>	139,176
<code>R_activeMoiety</code>	14,940
<code>I_CAS</code>	11,684
<code>I_CID</code>	2,044

Table 4. Top stitch value for each stitch key.

Stitch key	Stitch value	Size
<code>H_LyChI_L3</code>	VUSPQLGXN18-M	1,344,440
<code>H_LyChI_L4</code>	VU8BQZFPPYTZ-M	1,307,592
<code>N_Name</code>	ROFECOXIB	72
<code>H_LyChI_L5</code>	9DKQLD7D29DN-N	162,160
<code>I_UNII</code>	UNKNOWN	210
<code>R_activeMoiety</code>	2M83C4R6ZB	106
<code>I_CAS</code>	25322-68-3	1,806
<code>I_CID</code>	121225712	380

The LyChI hash keys L3 and L4 correspond to the potassium ion (K+).

In the remainder of this section, we provide examples that highlight different entity resolution challenges.

4.1 ASPIRIN

ASPIRIN is a versatile drug that can be used alone or in combination with other drugs. Shown in Figure 4 is the induced subgraph of the much larger *ASPIRIN* connected component that forms the *ASPIRIN* entity. This example demonstrates **Stitcher**’s ability to tease out only the relevant stitch nodes for which *ASPIRIN* is likely to be the active moiety for the underlying substance.

4.2 LEVOMETHADYL

LEVOMETHADYL and its derivative *LEVACETYLMETHADOL* are often considered as two separate drugs. This is readily apparent in Figure 5, which shows that there are two distinct “clusters” in the stitch multigraph. If entity resolution is based on graph metrics (e.g., betweenness centrality), it is likely that this connected component will yield two drugs instead of one. Here, the priority of the stitch key allows the two clusters to be merged to indicate that there is only one drug.

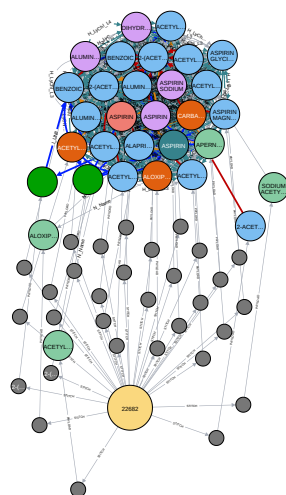


Fig. 4. A connected component for ASPIRIN.

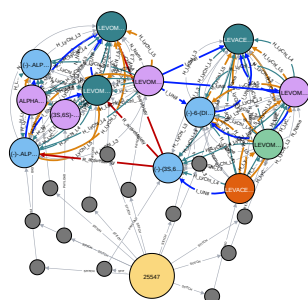


Fig. 5. A connected component for LEVOMETHADYL that clearly shows two distinct clusters.

4.3 BENOXAPROFEN

Figure 6 shows the connected component for *BENOXAPROFEN*, a nonsteroidal antiinflammatory drug approved in 1982. This drug is a racemic mixture. The density of this connected component is a reflection of the lack of specified stereocenter that caused many spurious stitch keys. *Stitcher* is able to disambiguate the connected component into three distinct entities that represent the mixture, *R*-, and *S*-.

Acknowledgments

We thank our colleagues, Mark Williams and Tyler Beck, for their valuable proof-reading of early drafts of the manuscript. We also thank our colleague, Tongan Zhao, for his help in developing a prototype curation user interface for *Stitcher*. We are particularly grateful to Alexey Zakharov and Tim Sheils for their constant encouragement and support.

References

Cormen, T. et al. (2001). *Introduction to algorithms*. MIT Press. Second edition.

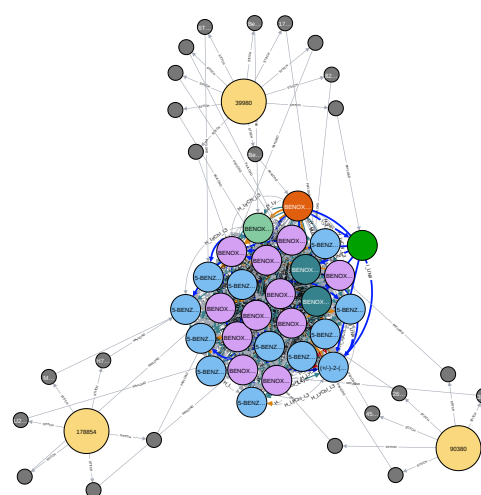


Fig. 6. A dense connected component for BENOXAPROFEN that resolved to three unique entities.

- Croset, S. et al. (2016). Flexible data integration and curation using a graph-based approach. *BMC Bioinformatics*, **32**, 918–925.
- FDA (2012). Drugs for human use. [<https://www.govinfo.gov/content/pkg/CFR-2012-title21-vol5/pdf/CFR-2012-title21-vol5-chapI-subchapD.pdf>; accessed January, 2020].
- FDA (2019a). Fda’s global substance registration system. [<https://www.fda.gov/industry/fda-resources-data-standards/fdas-global-substance-registration-system>; accessed January, 2020].
- FDA (2019b). G-srs substance database. [<https://tripod.nih.gov/ginas>; accessed January, 2020].
- FDA (2020a). Human drugs. [<https://www.fda.gov/industry/regulated-products/human-drugs>; accessed January, 2020].
- FDA (2020b). New drugs at fda. [<https://www.fda.gov/drugs/development-approval-process-drugs>; accessed January, 2020].
- Huang, R. et al. (2011). The NCGC Pharmaceutical Collection: A Comprehensive Resource of Clinically Approved Drugs Enabling Repurposing and Chemical Genomics. *Science Translational Medicine*, **3**, 80ps16.
- Huang, R. et al. (2019). The NCATS Pharmaceutical Collection: a 10-year update. *Drug Discovery Today*, **24**, 2341–2349.
- Hucklenbroich, P. (2014). “Disease Entity” as the Key Theoretical Concept of Medicine. *Journal of Medicine and Philosophy*, **39**, 609–633.
- ISO (2018). Health informatics—identification of medicinal products—data elements and structures for the unique identification and exchange of regulated information on substances. [<https://www.iso.org/standard/69697.html>; accessed January, 2020].
- NCATS (2013). Layered chemical identifier. [<https://github.com/ncats/lychi>; accessed October, 2019].
- NLM (2005). Rxnorm: Prescription for electronic drug information exchange. [<https://www.nlm.nih.gov/research/umls/rxnorm/RxNorm.pdf>; accessed January, 2020].