# NamesforLife Semantic Resolution Services for the Life Sciences

Charles T. Parker[1], Dorothea Rohlfs[1], Kara Mannor[1], Sarah Wigley[1], Nicole Osier[1], Catherine Lyons[1], George M. Garrity[1,2]*(garrity@namesforlife.com)

[1]NamesforLife, LLC, East Lansing, MI and [2]Michigan State University, East Lansing, MI
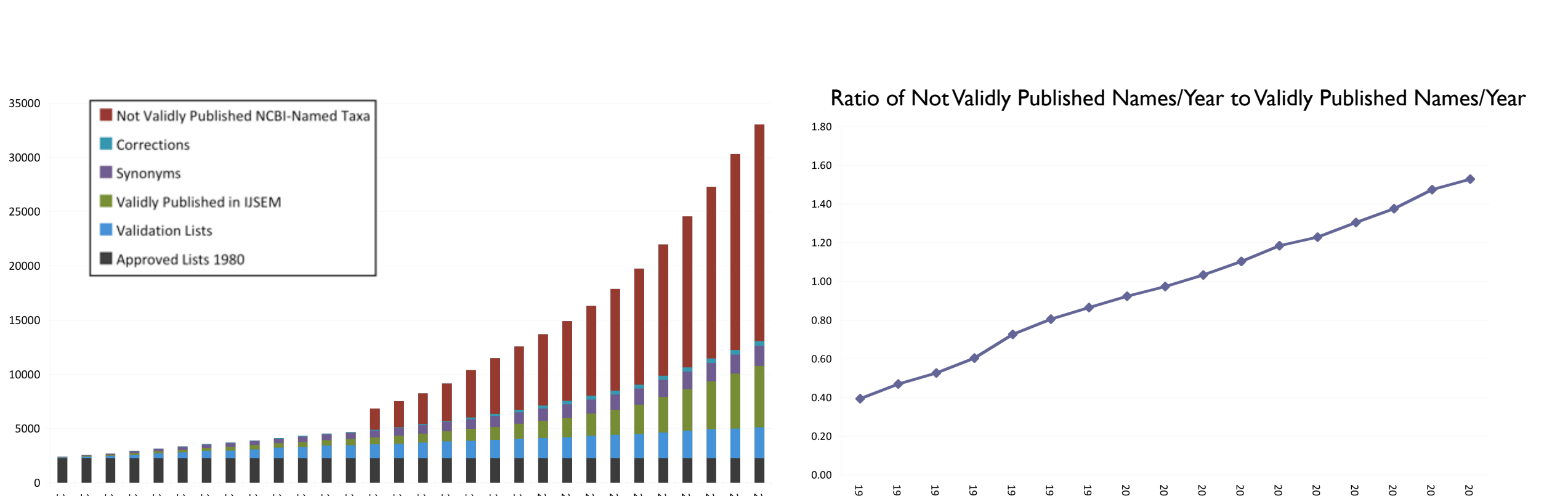
MICHIGAN STATE UNIVERSITY

## Abstract

Within the Genomes-to-Life Roadmap, the DOE states that a significant barrier to effective communication in the life sciences is a lack of standardized semantics that accurately describe data objects and persistently express knowledge change over time. As research methods and biological concepts evolve, certainty about correct interpretation of prior data and published results decreases because both become overloaded with synonymous and polysemous terms. Ambiguity in rapidly evolving terminology is a common and chronic problem in science and technology. N4L is a novel technology designed to solve this problem.

## Background

Understanding the correct meaning of a biological name in the appropriate context is important. It is not a trivial task, and the number of individuals with expertise in biological nomenclature is limited. Such knowledge can, however, be accurately modeled and delivered through a networked semantic resolution service that provides end-users of biological nomenclature or other dynamic terminologies with the appropriate information, in proper context, on demand. Such a service can also be used by database owners, publishers, or other information providers to semantically enable them more readily discoverable, even when the definition of a name or term has changed.
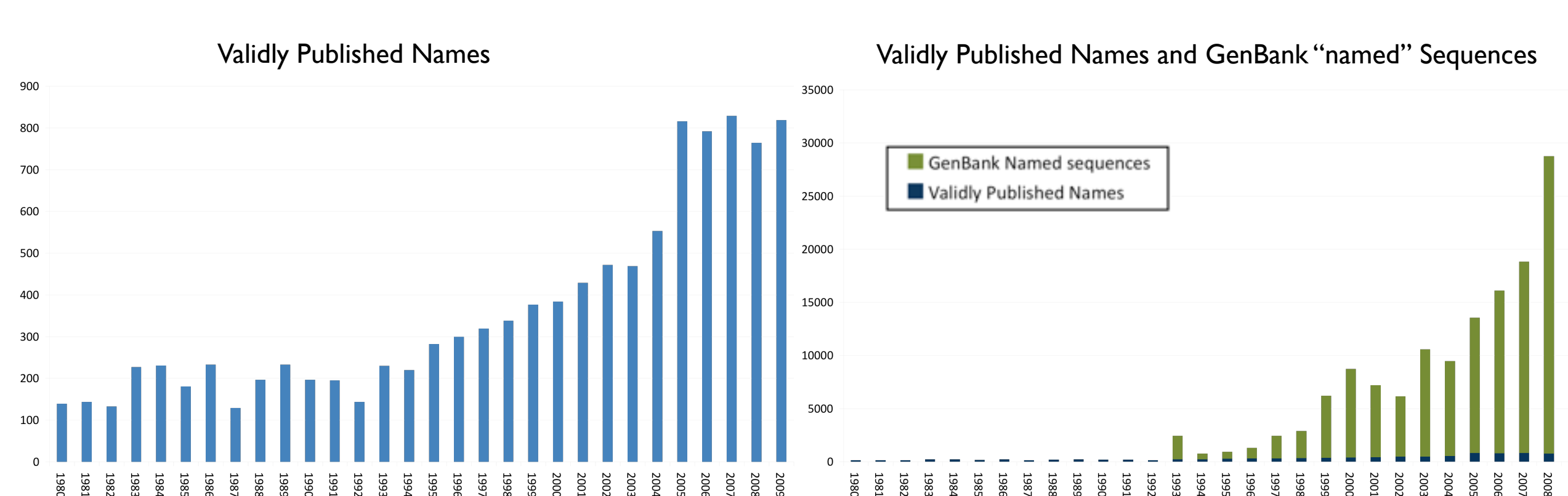
The figure below shows the cumulative growth in the number of validly published names since 1980. At the beginning of 1Q 2010, there were 13,069 validly published names (subspecies through classes) and 1,854 synonyms for which there were publicly available type strains. As seen in the chart below, the rate at which new validly published names appear in the literature plateaued in 2006, at around 800 names per year (excluding taxonomic rearrangements and emendations). GenBank began cataloging taxonomic information for sequences derived from Bacteria and Archaea in 1993, and introduced the practice of tagging data with trivial or invalid names. By 2002, the number of invalid names associated with GenBank records exceeded the number of validly published names. At present, invalidly published names outnumber validly published names (3:2). For both tracking and identification purposes, this adds significantly to the burden of those who must rely on this information for various purposes by adding an additional 4.8 names/day, none of which are governed by any rules or principles for assignment, usage, or uniqueness.



The chart above shows that the ratio of not validly published names to validly published names has been increasing at a linear rate since 1993.

## A looming problem

The adoption of DNA sequencing as the preferred method of rapidly characterizing *Bacteria* and *Archaea* has tremendously accelerated during the past five years, with the expected consequences. At present, the rate at which "named" sequences are added to the GenBank taxonomy exceeds the rate at which validly published names appear in the taxonomic record by a factor of approximately 35. This confounds the retrieval of related information from various databases and the scientific, technical and medical literature as many of these invalidly named species can not be readily tracked over time, nor can relationships be inferred to those species for which at least one genome sequence is available. This disconnect between the knowledge contained in the literature and the accumulated genomic data is likely to grow as faster and cheaper sequencing methods come into the market place.



## Database Statistics

The NamesforLife Nomenclature Database (N4LDB) contains 13,534 distinct names, 13,117 of which are validly published, 111 Candidatus, and 306 that are illegitimate but relevant to Bacterial Nomenclature. N4LDB also contains 10,293 exemplars, 9,005 of which represent distinct type strains for 10,558 taxa and 10,942 names. The remaining 2,592 names are associated with higher taxa that do not have exemplars. The breakdown is shown in Table 1 to the right.
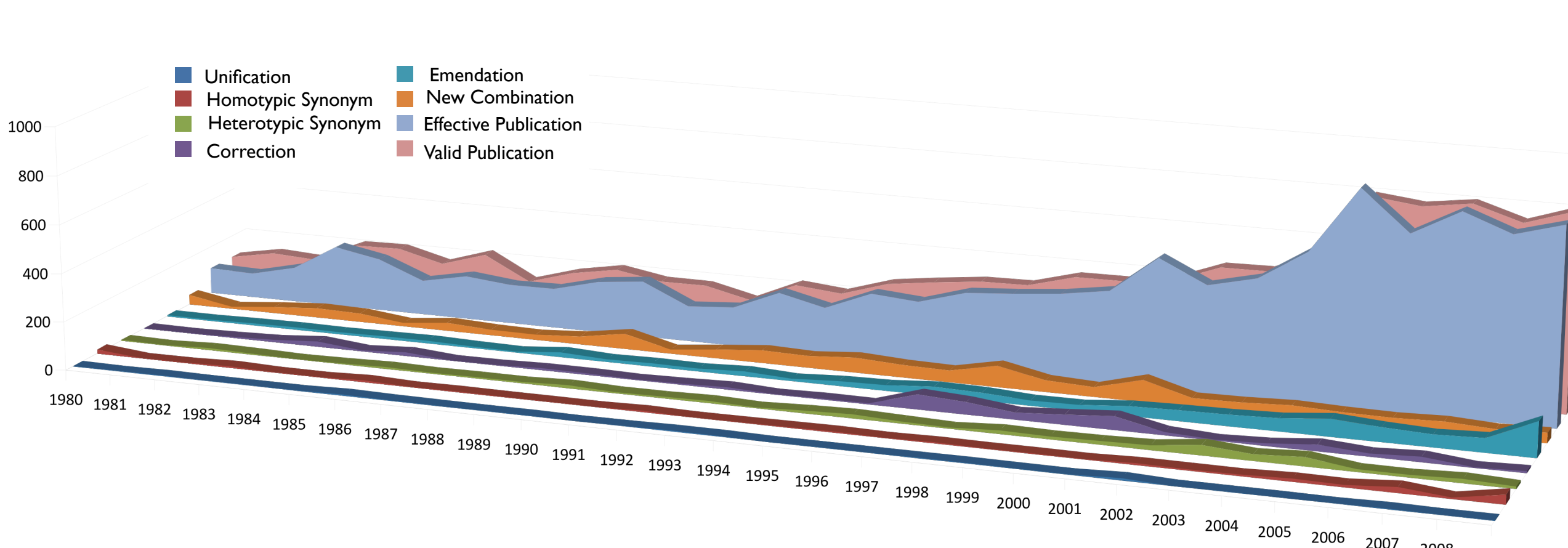
Table 2 shows a breakdown of the nomenclatural events by event type since the publication of the Approved Lists in 1980. Some of the less common events (Judicial Opinions, etc) are not shown here.

### Table 1: N4LDB Records by Rank

| Rank | Taxa | Names |
|---|---|---|
| Domain | 2 | 2 |
| Phylum | 34 | 35 |
| Class | 62 | 63 |
| Subclass | 5 | 5 |
| Order | 119 | 123 |
| Suborder | 20 | 21 |
| Family | 317 | 322 |
| Subfamily | 1 | 1 |
| Genus | 1965 | 1996 |
| Subgenus | 5 | 5 |
| Species | 10058 | 10403 |
| Subspecies | 519 | 558 |
| Total | 13107 | 13534 |

### Table 2: Nomenclatural Events Recorded in N4LDB

| Event | Count |
|---|---|
| Corrections | 431 |
| New Combinations | 1243 |
| Heterotypic Synonyms | 343 |
| Homotypic Synonyms | 171 |
| Unifications | 46 |
| Automatically created names via rule 40d | 65 |
| Emendations | 1104 |
| Validation List events | 2841 |
| Valid Publication | 10372 |

The graph below shows nomenclatural activity from the Approved Lists through the end of 2009. Clearly visible is the slight lag between effective and valid publication, as well as periods of higher activity for taxonomic rearrangements and emendations. The total number of references in the N4L database describing the events is 9,812.
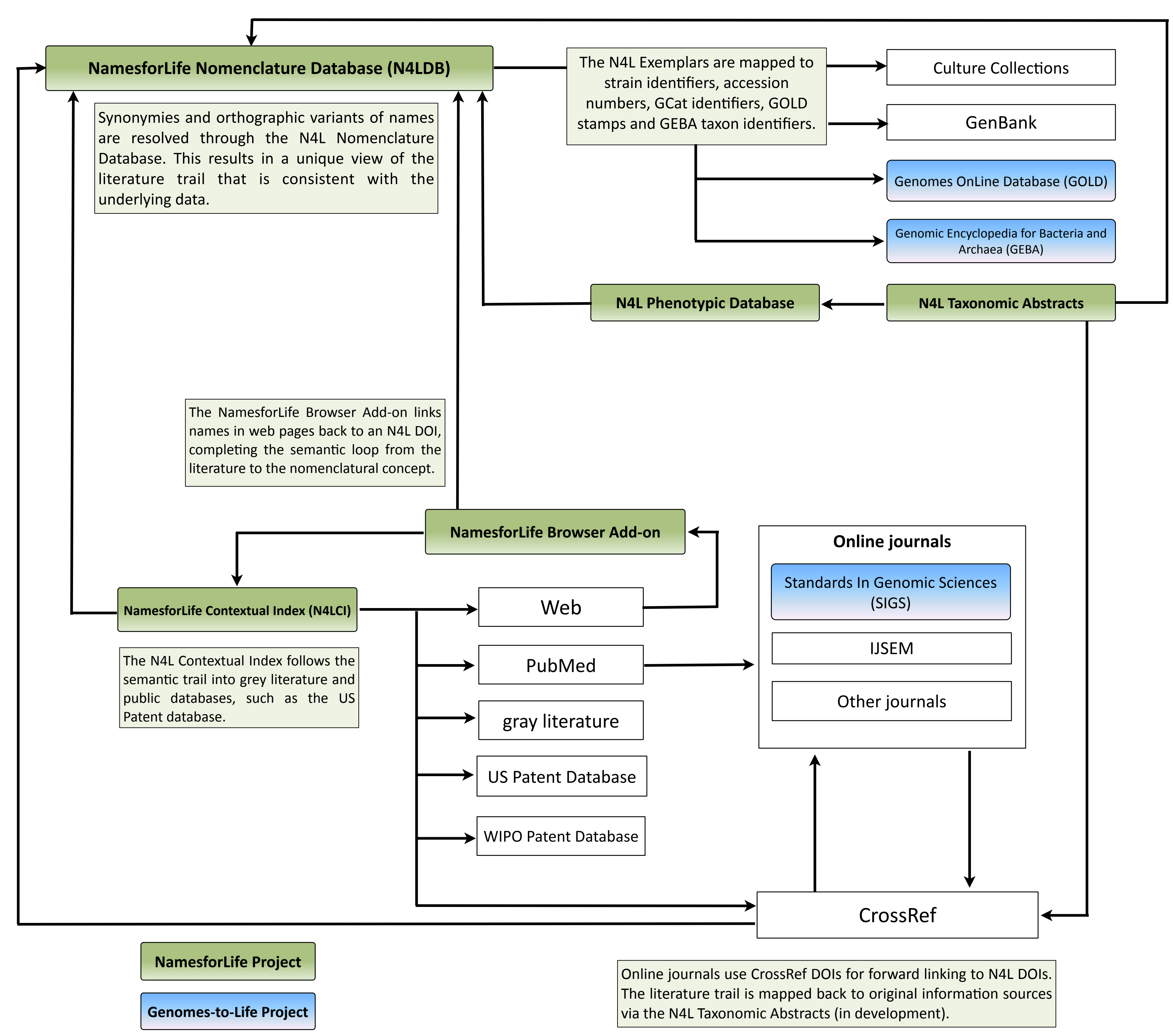


## N4L Technology

A full-text literature search can quickly reach a dead end if the terminology has changed. N4L, the technology developed by NamesforLife, resolves changes in Bacterial Nomenclature via a fully referenced model of the synonymies and corrections (see nomenclature example below).
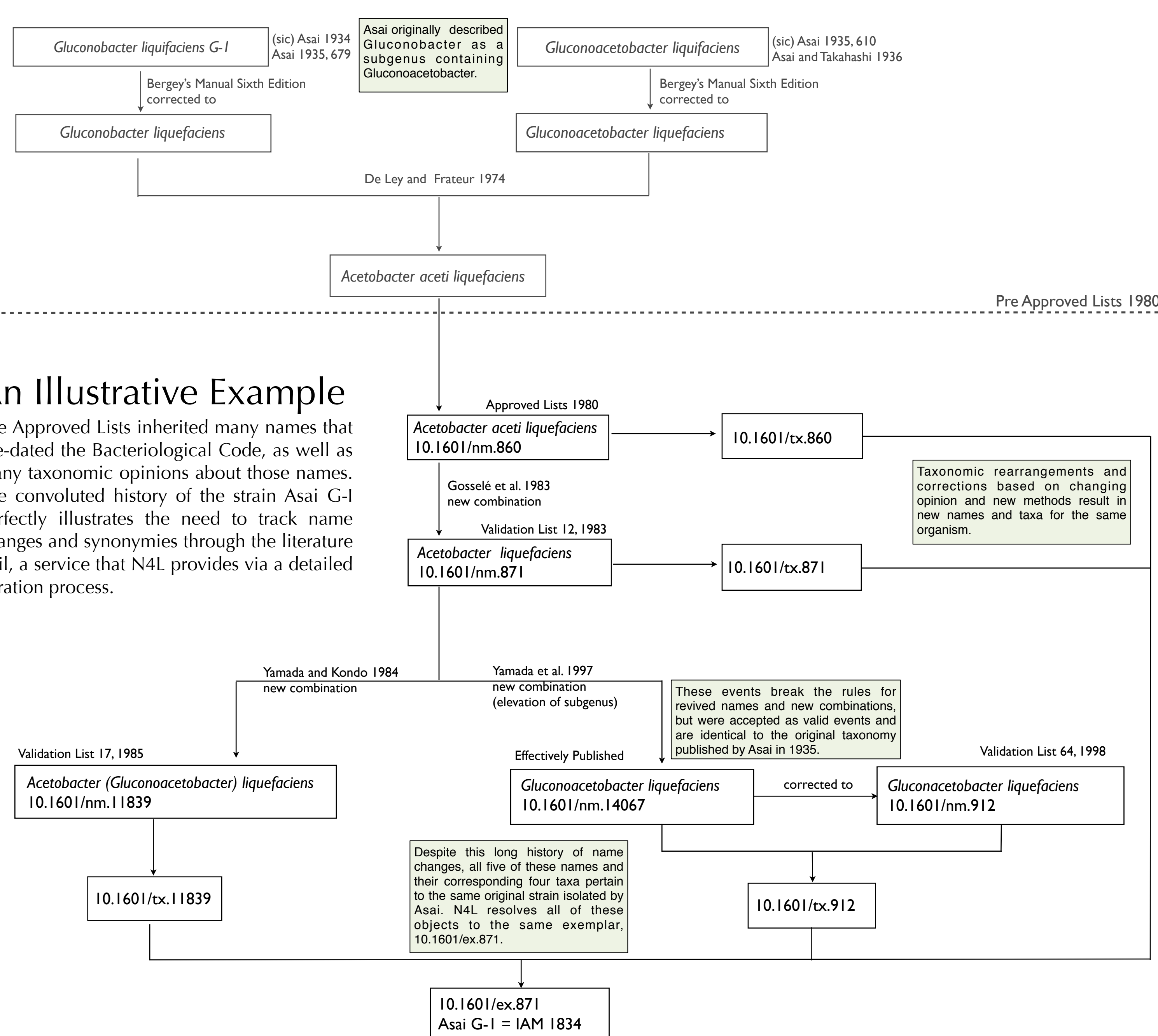
The power of the N4L model is that it does not aggregate data, but rather metadata; it links existing unique identifiers of literature and data together through the terminology itself (a bacterial name or strain), via a single persistent identifier: the N4L Digital Object Identifier (DOI). As shown in the figure below, several once-disparate information sources are now linked together through terminology in a way that is traversable.

Now that the Bacterial Nomenclature database is complete and updated in synchrony with the valid publication of nomenclatural changes, NamesforLife is in the process of linking together Bacterial Nomenclature, technical literature, and the various projects of the Genomes-to-Life program. In N4L, each individual organism is represented by a metadata object (an N4L Exemplar), which is identified by a DOI.

An N4L Exemplar aggregates what is known about an individual organism. The Genomes OnLine Database (GOLD), Standards in Genomic Sciences (SIGS), Genomic Encyclopedia for Bacteria and Archaea (GEBA) and Genomes and Metagenomes Catalogue (GEM) all use unique identifiers that link to each other in some way; via the GCat identifier, GOLD stamp, and GEBA Taxon Identifier. However, there is no single common link to the literature. NamesforLife is closing this gap by tying these disparate sources of information together via N4L Exemplars, which are integrated with the N4L Nomenclature Database and N4L Contextual Index.



### An Illustrative Example

The Approved Lists inherited many names that pre-dated the Bacteriological Code, as well as many taxonomic opinions about those names. The convoluted history of the strain Asai G-I perfectly illustrates the need to track name changes and synonymies through the literature trail, a service that N4L provides via a detailed curation process.



## Current Work

The NamesforLife curators are currently reviewing references for the last groups of type exemplars and finishing the links into the Genomes-to-Life projects, including non-type genomes. We expect this work to be completed by April of this year. The N4L Nomenclature Database will continue to be updated as nomenclatural events occur.

The development of the N4L Contextual Index is nearing completion. The first component of this system will come online this spring and will create a semantic path into US Patents and Patent Applications. Additional resources will become available via the Contextual Index throughout the year.

The N4L Taxonomic Abstracts are currently in development. These will provide a snapshot of Bacterial Nomenclature in the form of a citable micro-publication, and will serve to link existing literature to nomenclature via CrossRef.

The Beta release of the N4L Browser Add-on is officially scheduled to coincide with the Society for General Microbiology conference at the end of March 2010, but it is already available for early testing. Instructions on installation and use can be found at the NamesforLife services website, https://services.namesforlife.com. This Firefox Add-on detects and links bacterial names to the N4LDB, providing up-to-date nomenclature, strain and genome information, and a full bibliography. The screenshots below demonstrate the use of this tool on a Standards in Genomic Sciences article.



## Future Work

The Taxonomic Outline of Bacteria and Archaea (TOBA) is an open publication that has served the microbiology community for several years. NamesforLife will continue to support its ongoing publication. Release 8.0 of TOBA is planned following the completion of the current curatorial efforts.

In tandem with the new release of TOBA, NamesforLife will begin to publish the N4L Taxonomic Abstracts, described above. These will leverage the existing citation system developed by CrossRef to complete the semantic link between existing publications as shown in the diagram to the left.

An N4L Phenotypic Database is currently in the planning stages. A prototype of this database will be available in the summer of 2010 and will be incorporated into the N4L Taxonomic Abstracts. The goal is to provide a source of fully-referenced MIGS-compliant data for bacterial strains based on historical literature.

Additional tools and delivery methods for N4L content are also planned, including authoring and research tools, and alert services. New services will be announced as they become available.

## Acknowledgments