

Text-mining, metagenomics, names and phenotype: A high-level view.

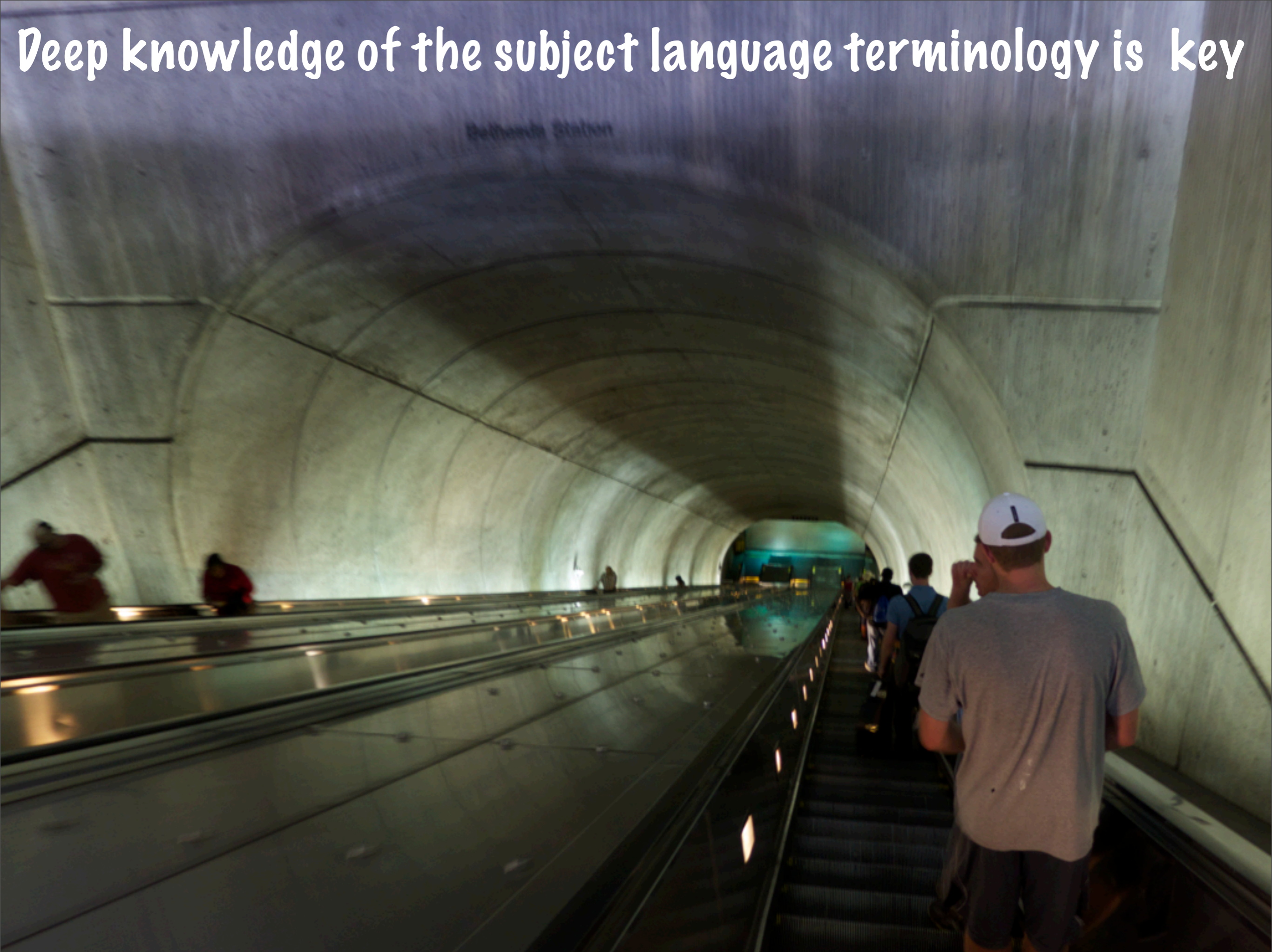
George Garrity
Michigan State University
NamesforLife, LLC





Zipf's law applies
...but so do the rules of nomenclature
...and communities of practice
...and the principles of semiotics and taxonomy

Deep knowledge of the subject language terminology is key



“Kbase also should incorporate metabolic, physiological, and morphological phenotypes used to identify species (e.g., from Bergey’s “differential characteristics” tables).”

**From the DOE Systems Biology
Knowledgebase Implementation Plan**



Description of *Legionella steelei* sp. nov.

Legionella steelei (steel'e.i. N.L. masc. gen. n. *steelei* of Steele, pertaining to the Australian microbiologist Trevor Steele, who performed pioneering work on the ecology and pathogenesis of *L. longbeachae* infection).

Gram-negative-staining rod. Grows on BCYE α agar, but not on tryptic soy blood agar or BCYE α agar without L-cysteine. Optimal growth temperature is 35 °C, with no or very poor growth at ≥ 37 °C. Produces large amounts of diffusible yellow-green fluorescent pigment. Variable production of colony-bound blue-white fluorescence. Positive in tests for activities oxidase, catalase, gelatinase, hippurate hydrolysis and β -lactamase. Negative result for glucose utilization. Motile, with monopolar flagella. Dominant cellular fatty acids are anteiso15:0, 16:1 and 16:0. Virulent in *Acanthamoeba castellanii*. 16S rRNA, *mip*, *rpoB*, *rnpB* and *proA* gene sequences differ significantly from all recognized species of the genus *Legionella*.

The type strain, IMVS-3376^T (=IMVS 3113^T=ATCC BAA-2169^T), was isolated from a human respiratory tract specimen.

Major features included in the NamesforLife Phenotypic Index

Strain metadata

- N4L Exemplar ID
- Isolation source
- Isolation method
- Isolation substrate
- Geographic location
- Environmental information
- Host
- Strain designation
- Collection ID(s)
- Taxon status (type/non-type)

Genotypic

- 16S rRNA sequence
- Other marker genes
- % DNA-DNA similarity
- % G+C composition
- Whole genome

Morphology

Micromorphology

- Cell size
- Cell shape
- Flagellation
- Sporulation
- Staining characteristics
- Other characteristics
- Intracellular inclusions
- Extracellular features
- Life cycle

Macromorphology

- Growth on solid surfaces
 - Colony morphology
- Growth in liquid
- Pigment production
- Other features

Chemotaxonomy

- Fatty acids
- Polar Lipids
- Mycolic Acids
- Respiratory quinones
- Peptidoglycan composition
- Polyamines

Physiological

- terminal e- acceptor
- substrate utilization
- metabolic end-products
- sensitivity/tolerance to
 - chemical and physical agents

Technical and non-technical barriers

National Park Service
U.S. Department of the Interior

Because of the
Federal Government SHUTDOWN,
**All National Parks
Are CLOSED.**



N4L DOI mediated semantic services

- Terminology development tools
- Annotation tools
- Nomenclature, data and ontologies



INTERNATIONAL
COMMITTEE ON
SYSTEMATICS OF
PROKARYOTES

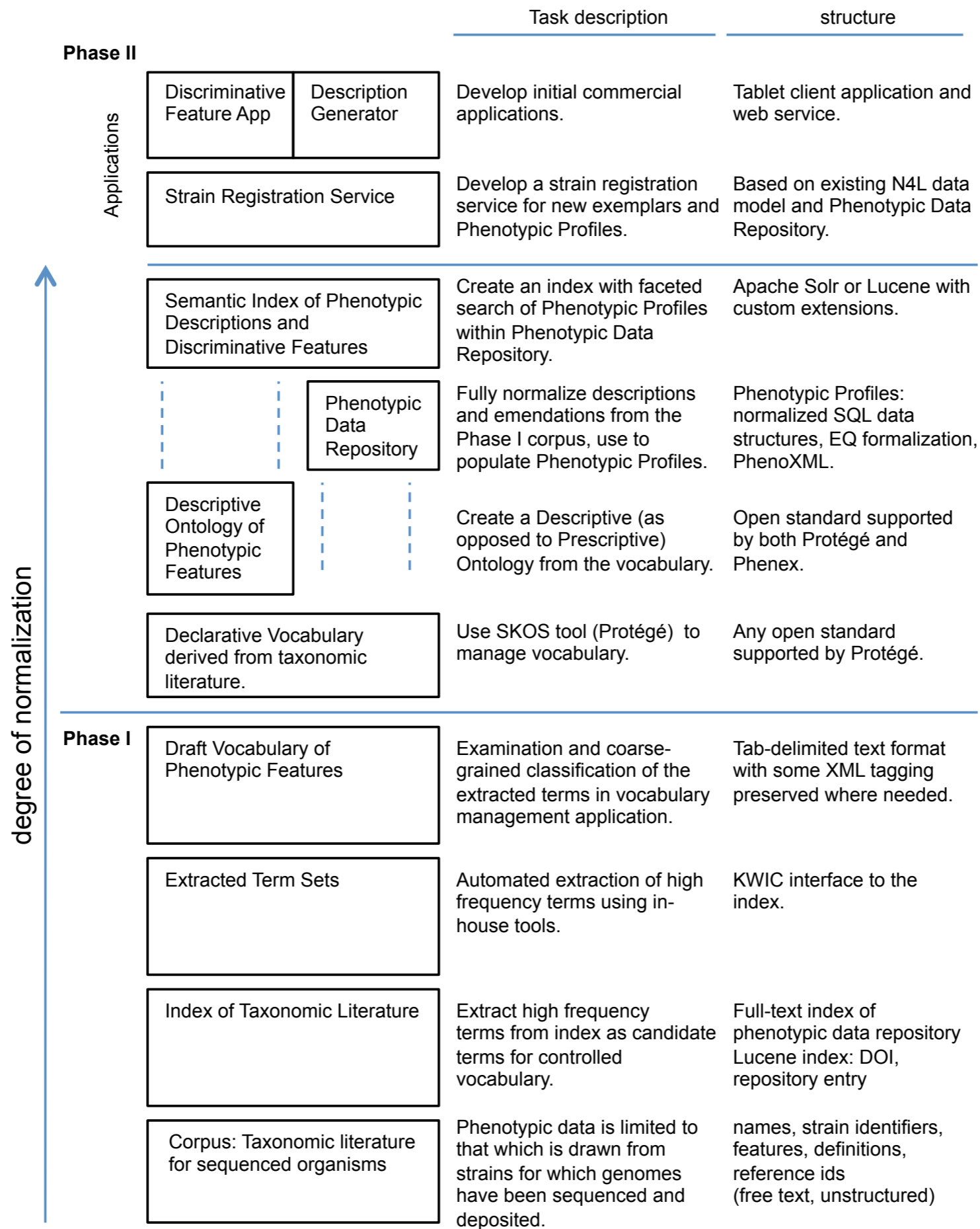


Project personnel
Chuck Parker
Nenad Krzavac
Kevin Petersen
Grace Rodriquez
Amber Roberts



Visit us at <http://services.namesforlife.com>





Phase I

| | | |
|--|---|--|
| Draft Vocabulary of Phenotypic Features | Examination and coarse-grained classification of the extracted terms in vocabulary management application. | Tab-delimited text format with some XML tagging preserved where needed. |
| Extracted Term Sets | Automated extraction of high frequency terms using in-house tools. | KWIC interface to the index. |
| Index of Taxonomic Literature | Extract high frequency terms from index as candidate terms for controlled vocabulary. | Full-text index of phenotypic data repository Lucene index: DOI, repository entry |
| Corpus: Taxonomic literature for sequenced organisms | Phenotypic data is limited to that which is drawn from strains for which genomes have been sequenced and deposited. | names, strain identifiers, features, definitions, reference ids (free text, unstructured) |

Name Information

phylum class subclass order suborder family **genus** species subspecies
Collimonas
 10.1601/nm.8548

Article Information

Int J Syst Evol Microbiol 2004; **54**(3):857.
 10.1099/ijls.0.02920-0

Description of *Collimonas* gen. nov.

Cells are strictly aerobic, straight or slightly curved, Gram-negative rods, 0.3–0.5×1.0–2.0 µm. They occur singly and possess flagella (mostly one to three polar, but in some cases, several lateral) and pili when cultured in liquid media. Oxidase activity is positive; catalase activity is negative or weakly positive. Maximal growth is observed between 20 and 30 °C, without a sharp optimum. Maximum temperature that supports growth is approximately 35 °C. Optimal growth occurs at pH 6.5. Cells are able to hydrolyse colloidal chitin and milk proteins, but not lichenan or cellulose. Isolates vary in their ability to hydrolyse colloidal chitosan (deacetylated chitin). On minimal colloidal chitin/yeast agar, circular cleared haloes (final diameter, 4–10 mm) that support little translucent biomass are produced. Halo formation is completely repressed in chitin agar that contains TSB or glucose. On 0.1× TSB agar, colony morphology is variable. A wide range of sugars, alcohols, organic acids and amino acids can be metabolized. Several di- and trisaccharides cannot be used as substrates. In purified sand, cells proliferate upon introduction of intact, living fungal hyphae of various species of soil fungi. Major cellular fatty acids are C_{16:0} and C_{16:1ω7c}. DNA G+C content is 57–62 mol%. A member of the β-Proteobacteria, related most closely to the genera *Herbaspirillum* and *Janthinobacterium* in the family 'Oxalobacteraceae', order 'Burkholderiales'. Characteristics useful to differentiate the genus *Collimonas* from these related genera are given in Tables 1 and 2, as well as in the Results and Discussion. So far, the genus *Collimonas* is only known to occur in slightly acidic sandy dune soils. The type species is *Collimonas fungivorans*.

Etymology

Collimonas (Col.li.mo'nas, L. masc. n. *collis* hill; Gr. n. *monas* a unit, monad; N.L., fem. n. *Collimonas* cell from the hill).

KWIC Index : 714 occurrences of "c16:1ω7c"

haemolytic on sheep blood. The predominant fatty acid is C_{16:1ω7c}. The G+C content is 64.5 mol% (sd 0.1) (det
 talase is not produced. The major cellular fatty acid is C_{16:1ω7c}. Cell-wall murein is type A4α containing 1
 fungi. Major cellular fatty acids are C_{16:0} and C_{16:1ω7c}. DNA G+C content is 57–62 mol%. A member of the β-
 llow, flexirubin-like pigment and contain C_{16:1ω5c} and C_{16:1ω7c} as the major fatty acids. They do not hydro
 acturonate or l-arabitol. Major cellular fatty acids are C_{16:1ω7c} and C_{16:0}. Q-8 is the predominant respira
 acturonate or l-arabitol. Major cellular fatty acids are C_{16:1ω7c} and C_{16:0}. Q-8 is the predominant respira
 acturonate or l-arabitol. Major cellular fatty acids are C_{16:1ω7c} and C_{16:0}. Q-8 is the predominant respira

| Filters (inclusive) | invert | Unfiltered Terms (749) | frequency |
|--|-------------------------------------|------------------------|-----------|
| <input type="checkbox"/> Extracted Terms dictionary (4,763 terms) | <input checked="" type="checkbox"/> | c16:0 | 1375 |
| <input type="checkbox"/> Bacterial Names dictionary (15,346 terms) | <input checked="" type="checkbox"/> | iso-c15:0 | 979 |
| <input type="checkbox"/> Authorities dictionary (3,906 terms) | <input checked="" type="checkbox"/> | c8 | 976 |
| <input type="checkbox"/> PubChem CID-MeSH dictionary (63,702 terms) | <input checked="" type="checkbox"/> | c4 | 965 |
| <input type="checkbox"/> KEGG Lipids ontology (2,739 terms) | <input checked="" type="checkbox"/> | c14 | 816 |
| <input type="checkbox"/> PubChem MeSH-Pharm dictionary (6,133 terms) | <input checked="" type="checkbox"/> | anteiso-c15:0 | 770 |
| <input type="checkbox"/> American National Corpus dictionary (237,267 terms) | <input checked="" type="checkbox"/> | c16:1ω7c | 714 |
| <input type="checkbox"/> Webster's 1934 Dictionary dictionary (310,536 terms) | <input checked="" type="checkbox"/> | iso-c16:0 | 663 |
| <input type="checkbox"/> Geographic Regions dictionary (24,734 terms) | <input checked="" type="checkbox"/> | c18:1ω7c | 602 |
| <input type="checkbox"/> Culture Collections dictionary (349 terms) | <input checked="" type="checkbox"/> | anteiso-c17:0 | 447 |
| <input type="checkbox"/> Nucleotides RegEx | <input checked="" type="checkbox"/> | summed feature 3 | 419 |
| <input type="checkbox"/> Numbers RegEx | <input checked="" type="checkbox"/> | iso-c15:0 2-oh | 384 |
| <input type="checkbox"/> Term Length 3 ≤ length ≤ any | <input checked="" type="checkbox"/> | c18:0 | 375 |
| <input type="checkbox"/> Term Frequency 5 ≤ frequency ≤ 1000 | <input checked="" type="checkbox"/> | c14:0 | 348 |
| | | c15:0 | 271 |
| | | iso-c17:0 | 263 |
| | | iso-c17:0 3-oh | 253 |
| | | c18:1ω9c | 239 |
| | | iso-c14:0 | 236 |
| | | c17:0 | 211 |
| | | c17:1ω8c | 153 |
| | | iso-c17:1ω9c | 140 |
| | | c18:1 | 138 |
| | | c12:0 | 123 |
| | | c16:1 | 117 |
| | | c12:0 3-oh | 109 |
| | | iso-c15:0 3-oh | 99 |
| | | c19:0 cyclo ω8c | 97 |
| | | c10:0 3-oh | 91 |
| | | summed feature 4 | 88 |
| | | c16:1ω5c | 86 |
| | | c17:0 cyclo | 80 |
| | | c17:1ω6c | 78 |
| | | iso-c15:1 | 74 |
| | | c16:1ω9c | 69 |
| | | c16:0 3-oh | 64 |
| | | i-c15:0 | 64 |
| | | iso-c15:1 g | 63 |
| | | c20 | 62 |
| | | c14:0 3-oh | 61 |
| | | 10-methyl c18:0 | 59 |
| | | 11-methyl c18:1ω7c | 58 |

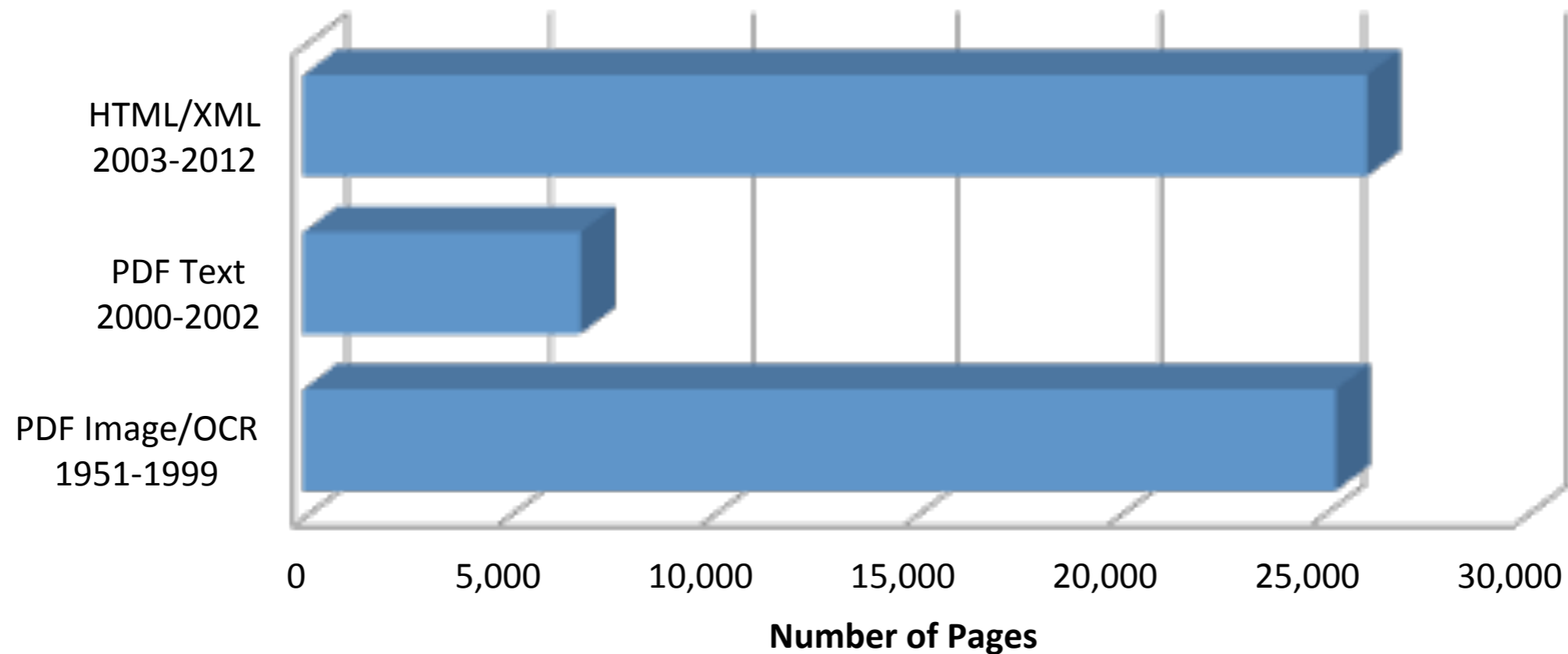
"c16:1ω7c"

Phase II

Applications

| | | Task description | structure |
|--------------|---|--|--|
| Applications | Discriminative Feature App | Develop initial commercial applications. | Tablet client application and web service. |
| | Description Generator | | |
| | Strain Registration Service | Develop a strain registration service for new exemplars and Phenotypic Profiles. | Based on existing N4L data model and Phenotypic Data Repository. |
| | Semantic Index of Phenotypic Descriptions and Discriminative Features | Create an index with faceted search of Phenotypic Profiles within Phenotypic Data Repository. | Apache Solr or Lucene with custom extensions. |
| | Phenotypic Data Repository | Fully normalize descriptions and emendations from the Phase I corpus, use to populate Phenotypic Profiles. | Phenotypic Profiles: normalized SQL data structures, EQ formalization, PhenoXML. |
| | Descriptive Ontology of Phenotypic Features | Create a Descriptive (as opposed to Prescriptive) Ontology from the vocabulary. | Open standard supported by both Protégé and Phenex. |
| | Declarative Vocabulary derived from taxonomic literature. | Use SKOS tool (Protégé) to manage vocabulary. | Any open standard supported by Protégé. |

Scoring functions



The corpus

pre-2000 descriptions [corrected OCR, from 2,964 descriptions
 2003-2009 freely available on web (3,338 articles)
 2010-2012 provided via partnership with SGM (1,080 articles)

Yielded 5,750 descriptions and emendations

41.full.pdf - Adobe Reader

File Edit View Window Help

Obviously if *Pseudomonas*
Migula 1894 is to be recognized
as a genus it must be associated
with the names of the species
designated and described by
Migula in 1895. Alternatively
the generic name *Pseudomonas*
might be recognized as having
been validly published in 1895.

Migula (1900 p. 884) (2) re-
cognized that the name *Bacterium*
aequalium Schroter 1872 (3)
antedated *Bacterium pyroganum*
Gessard 1882 (4) and published
the name as *Pseudomonas aequali-*
nosa (Schroter) Migula 1900.

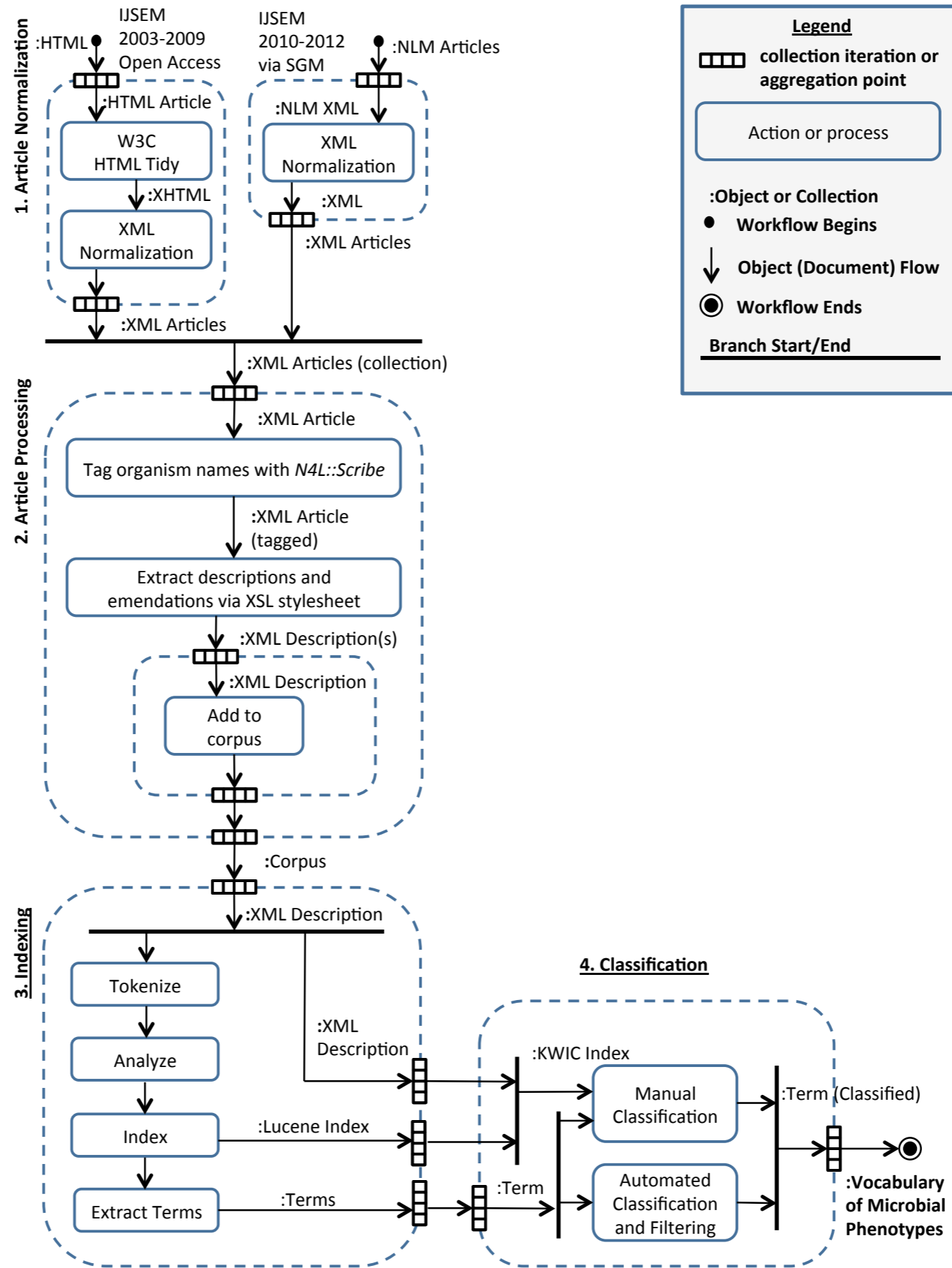
Euchanan (1918 p. 18) desig-

6.50 x 8.50 in

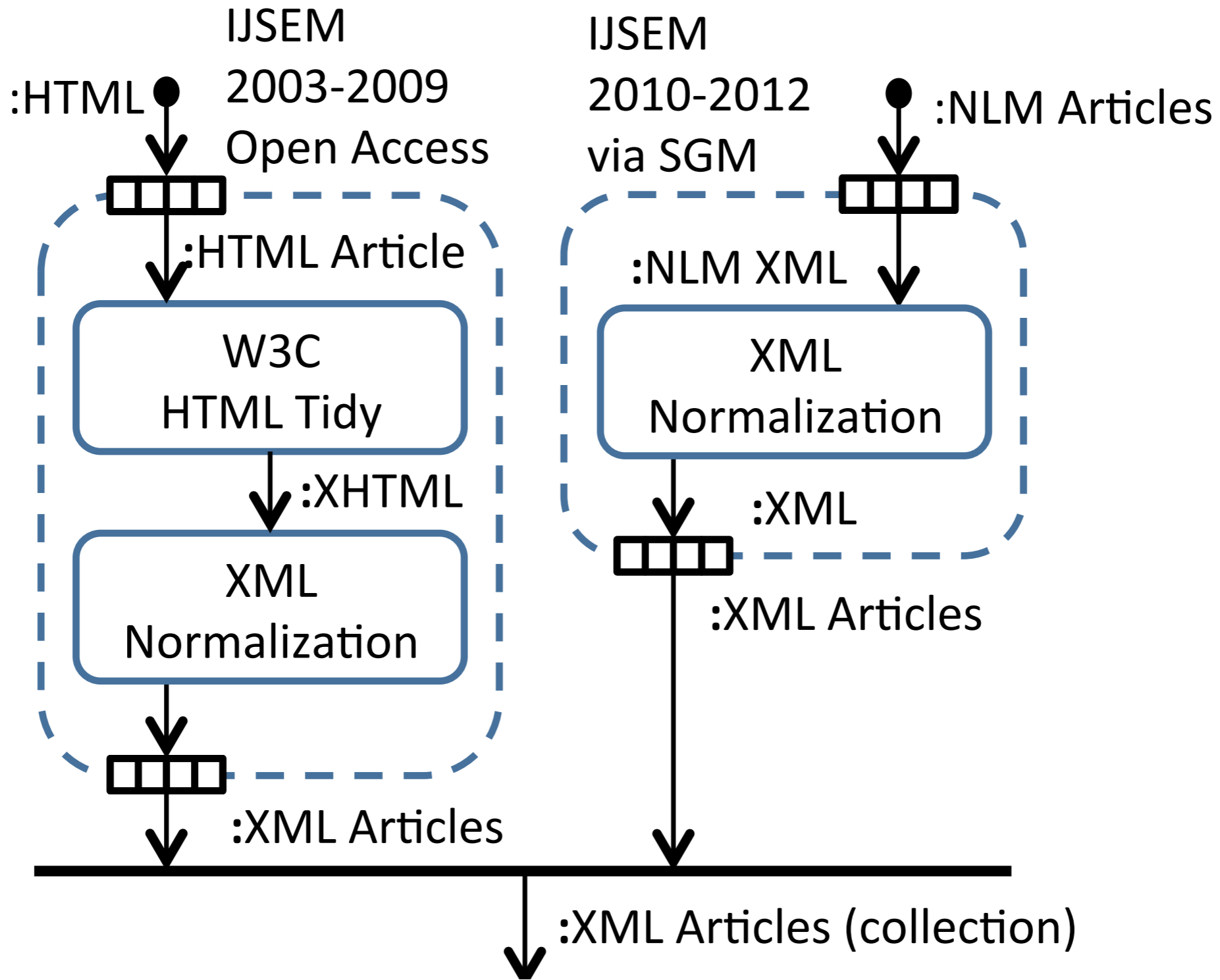
IJSB_1_41.txt - Notepad

File Edit Format View Help

Obviously if
Migula 1894 is to be recognized
as a genus it must be associated
with the names of the species
designated and described by
Migula in 1895. Alternatively
might be recognized as havin
k e n validly published i n 188s
the generic name *Pseudomonas*
Gessard 1882 (4) and published
(1) Schreter J Schizomycetes
i n C o b . F Kryptogamen FloRa
von Schlesien 3 136 174, Qe6
(2) Migula. W System der
Bakterien 2 884 -3Uc
6 urch Bakterien gebildete P i g

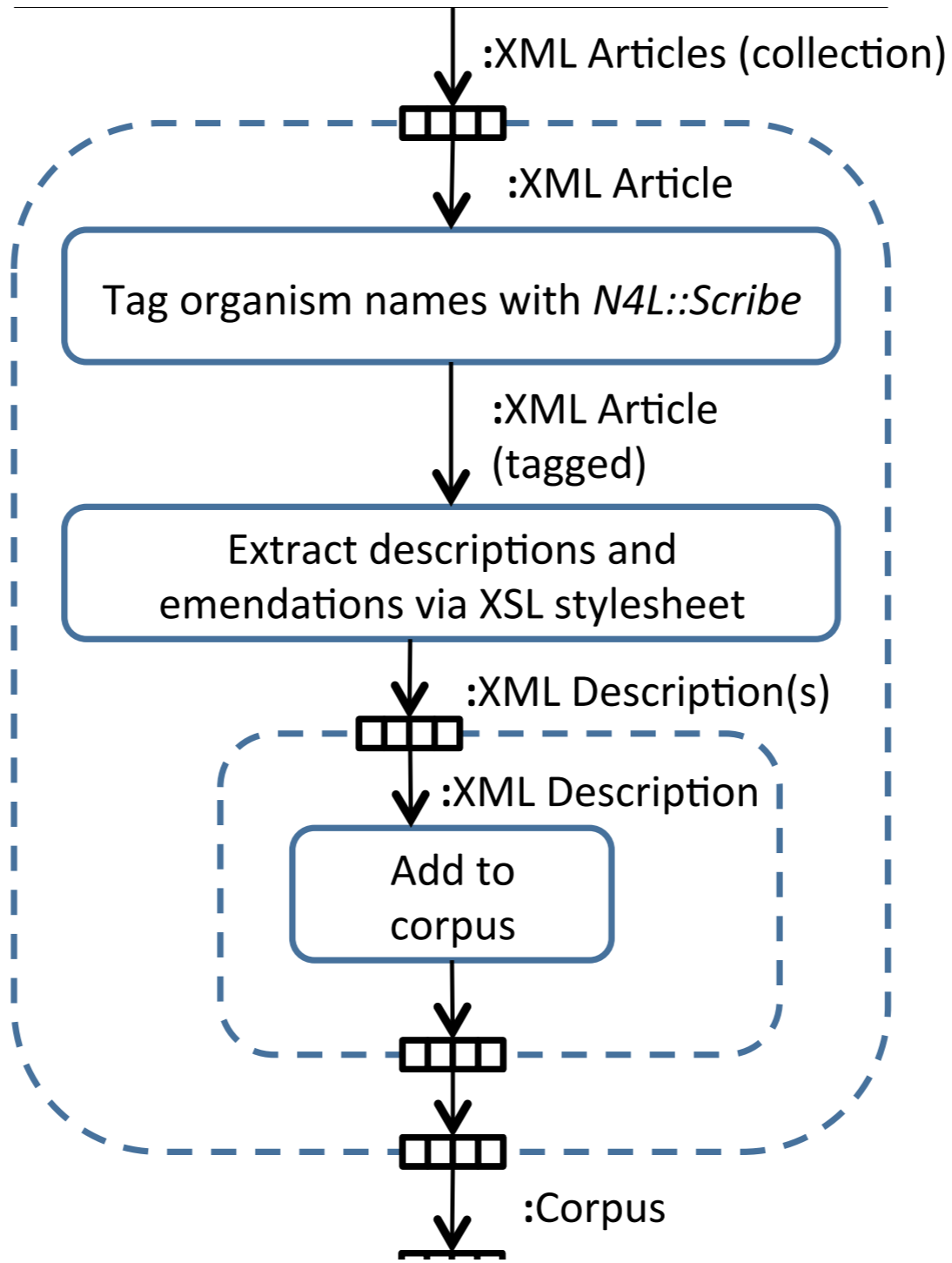


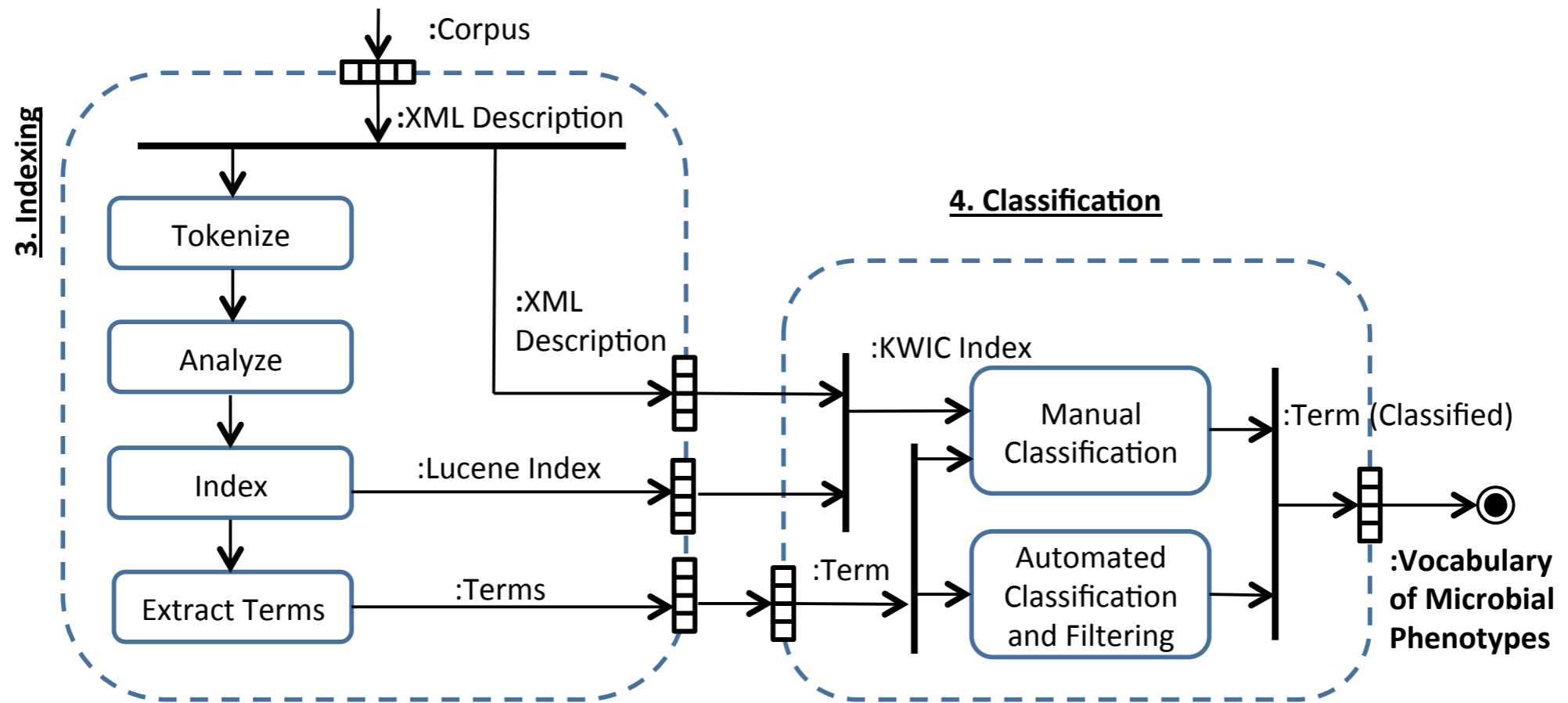
1. Article Normalization

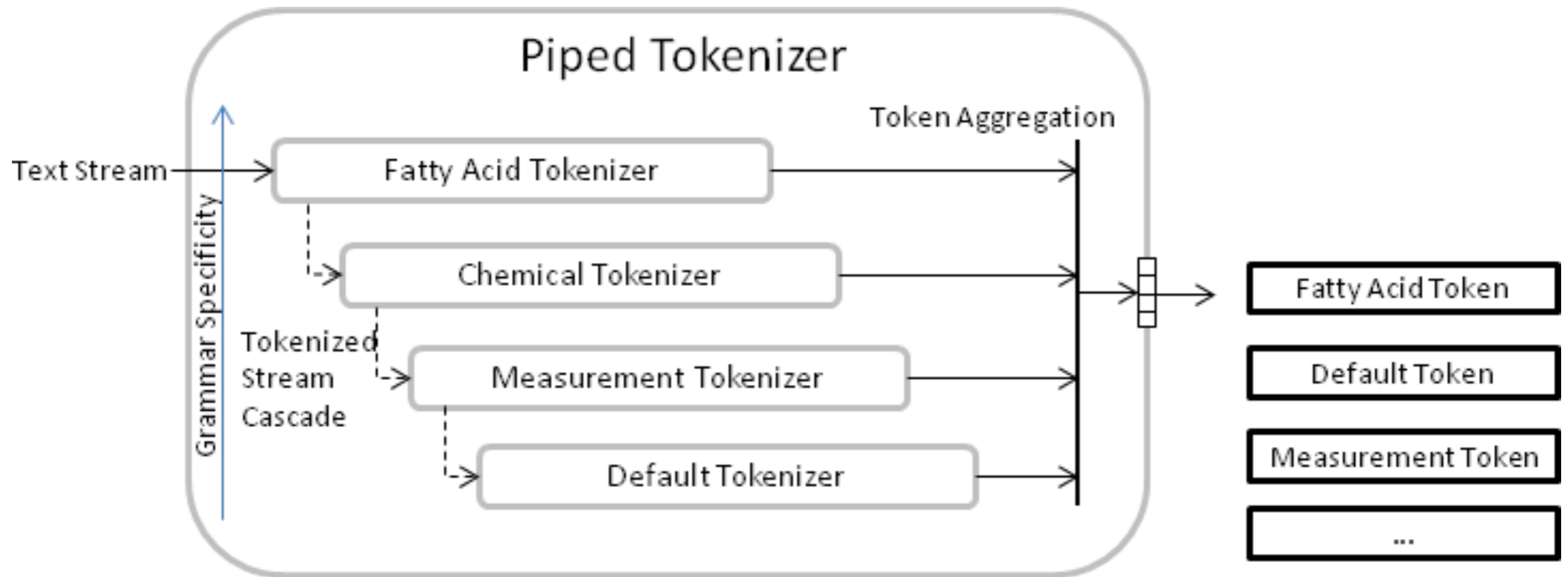


3

2. Article Processing





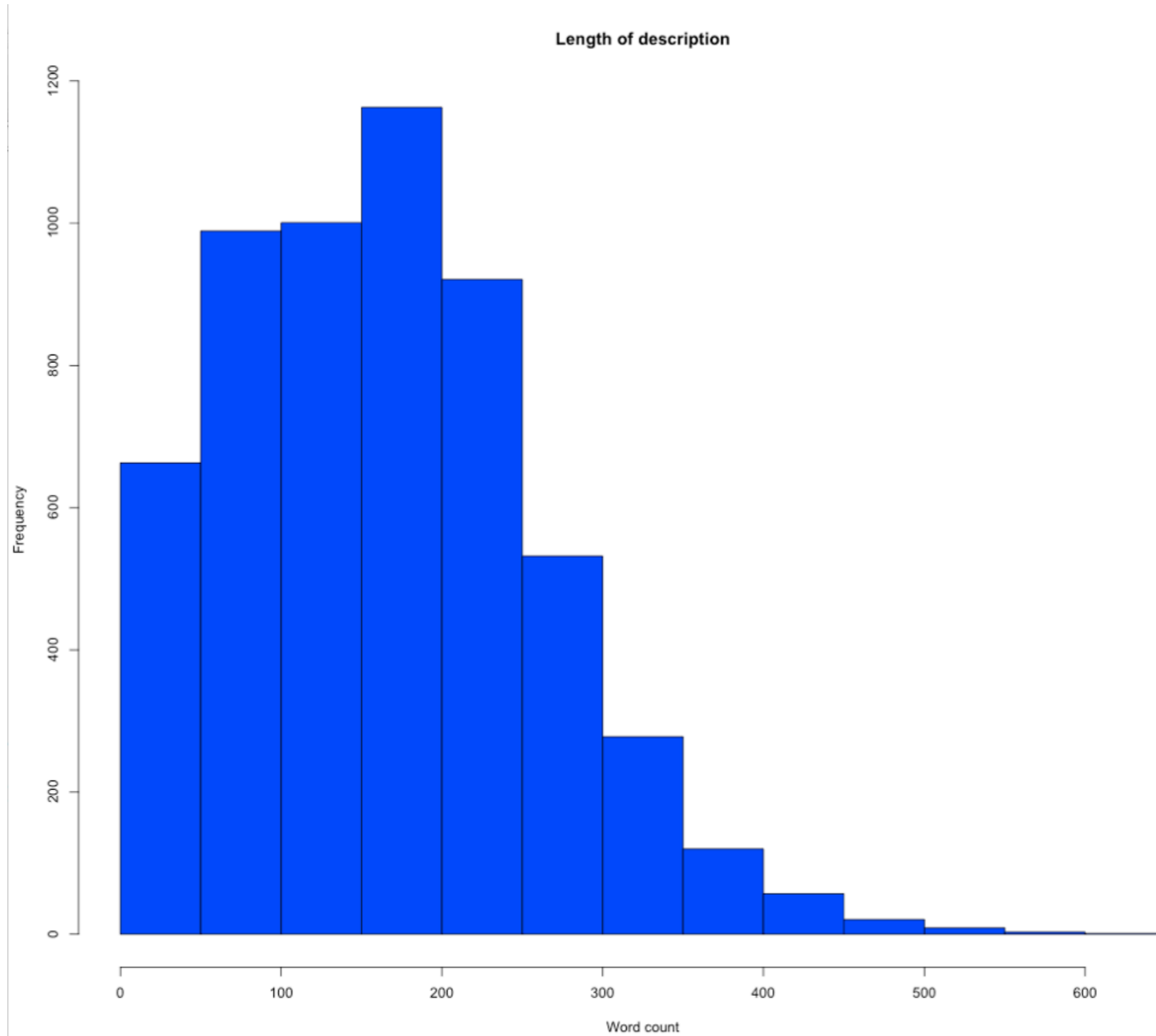


Taxonomic composition of the corpus

| Rank | Frequency |
|-------------|------------------|
| phylum | 4 |
| class | 17 |
| subclass | 4 |
| order | 31 |
| suborder | 19 |
| family | 131 |
| genus | 1095 |
| species | 4382 |
| subspecies | 75 |

Corpus metrics

| Phylum | n | total | mean | Phylum | n | total | mean |
|----------------------------|------|--------|------|-----------------------|----|-------|------|
| <i>Proteobacteria</i> | 2220 | 272688 | 123 | Chlorobi | 15 | 1691 | 113 |
| <i>Actinobacteria</i> | 1255 | 152022 | 121 | Synergistetes | 15 | 951 | 63 |
| <i>Firmicutes</i> | 994 | 120469 | 121 | Planctomycetes | 13 | 1604 | 123 |
| <i>Bacteroidetes</i> | 765 | 105384 | 138 | Deferribacteres | 12 | 803 | 67 |
| <i>Euryarchaeota</i> | 157 | 16990 | 108 | Armatimonadetes | 11 | 459 | 42 |
| <i>Verrucomicrobia</i> | 48 | 5499 | 115 | Caldiserica | 6 | 219 | 37 |
| <i>Deinococcus-Thermus</i> | 38 | 4594 | 121 | Fusobacteria | 6 | 488 | 81 |
| <i>Chloroflexi</i> | 34 | 1845 | 54 | Gemmatimonadetes | 6 | 208 | 35 |
| <i>Aquificae</i> | 28 | 1953 | 70 | Thermodesulfobacteria | 5 | 423 | 85 |
| <i>Tenericutes</i> | 27 | 2118 | 78 | Nitrospirae | 4 | 323 | 81 |
| <i>Spirochaetes</i> | 24 | 2309 | 96 | Chlamydiae | 3 | 277 | 92 |
| <i>Crenarchaeota</i> | 18 | 1416 | 79 | Cyanobacteria | 2 | 92 | 46 |
| <i>Acidobacteria</i> | 16 | 1305 | 82 | Lentisphaerae | 2 | 64 | 32 |
| <i>Thermotogae</i> | 16 | 1626 | 102 | Chrysiogenetes | 1 | 75 | 75 |



Not all descriptions are equally informative

Description of *Kribbella koreensis* (Lee *et al.* 2000) comb. nov.

Kribbella koreensis (ko.re.en'sis. N.L. fem. adj. *koreensis* pertaining to Korea, the location of the soil sample from which the type strain was isolated).

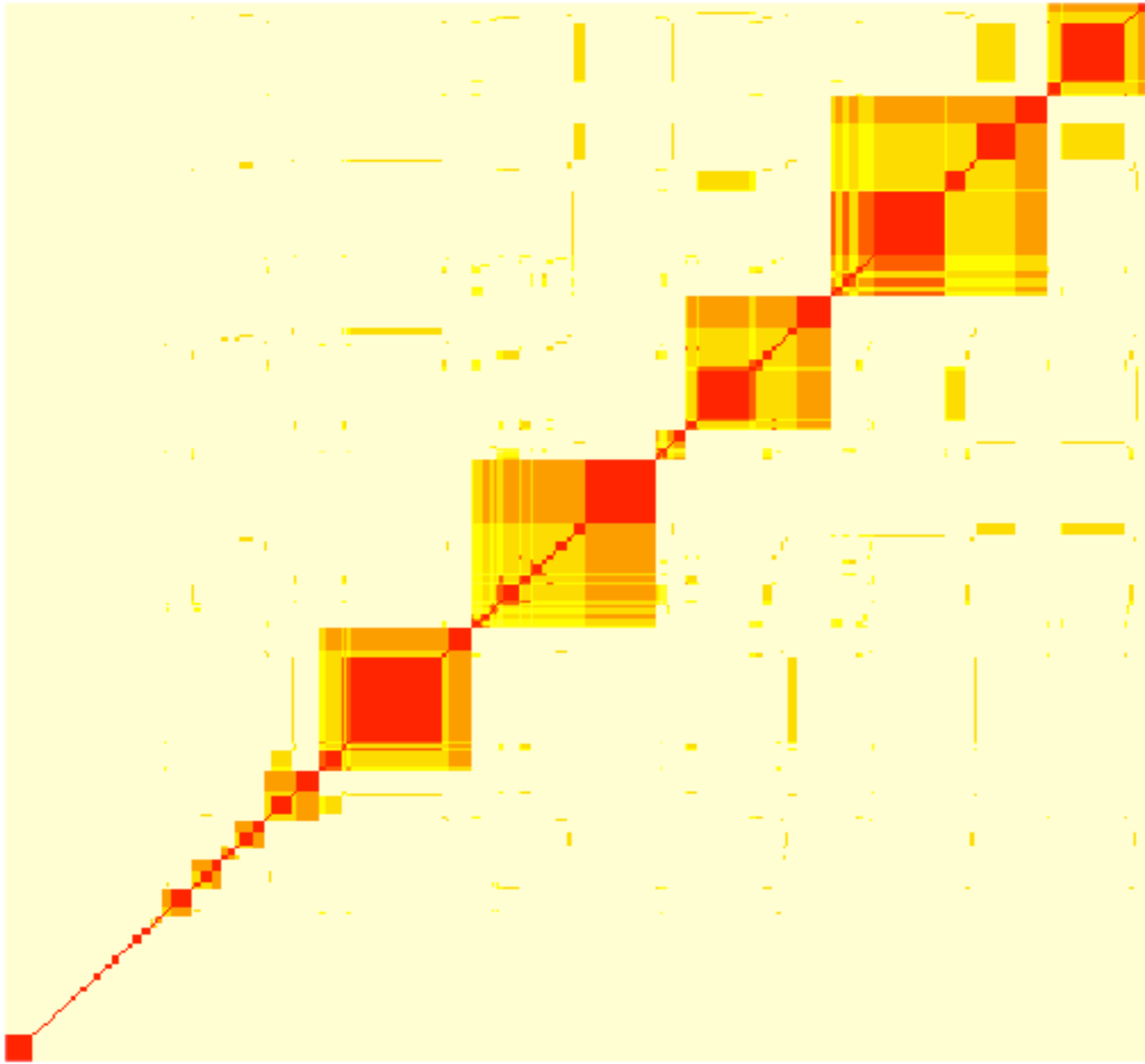
The description is identical to that of Lee *et al.* (2000).

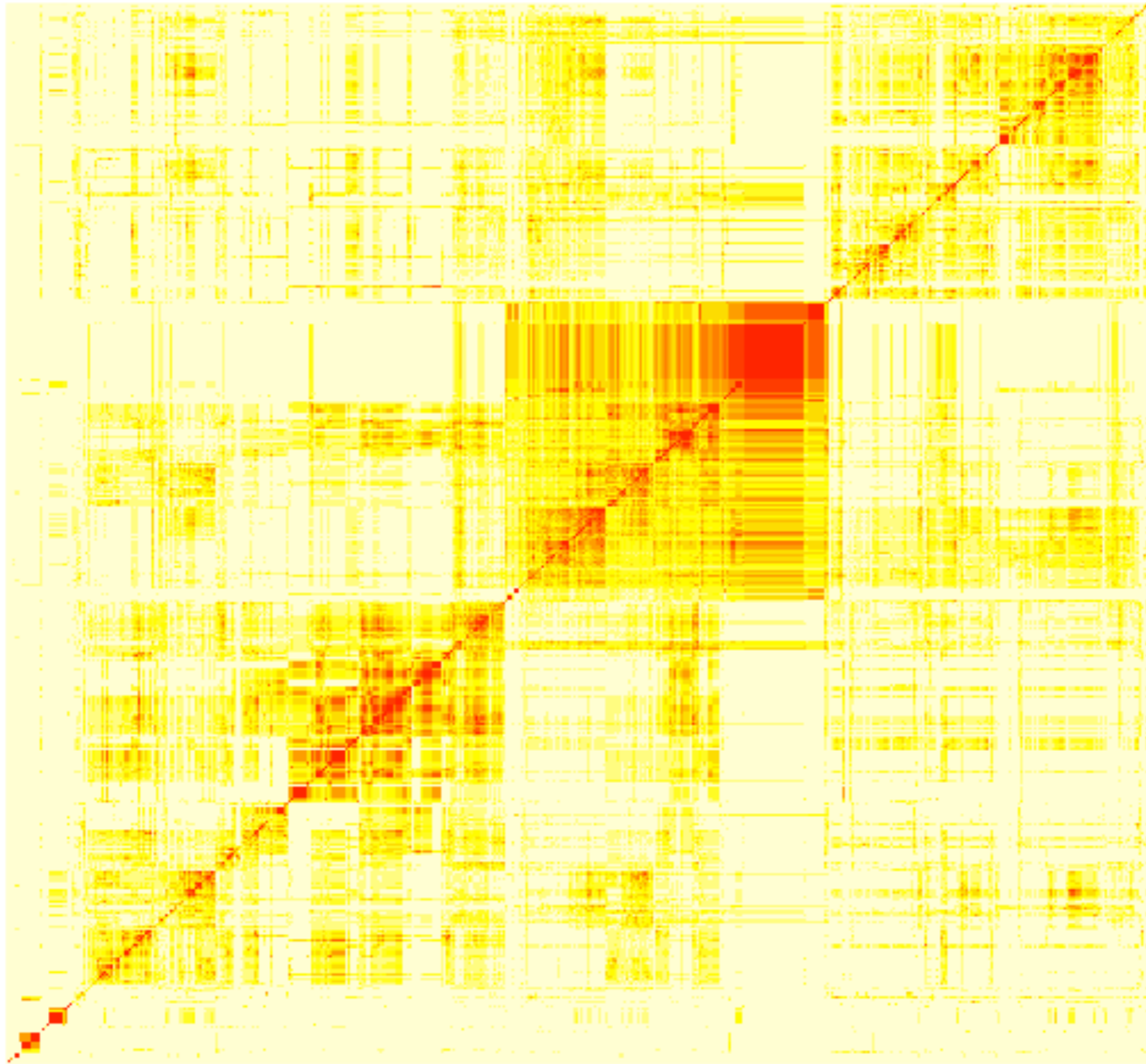
Emended description of *Anoxybacillus pushchinoensis* corrig. Pikuta *et al.* 2000

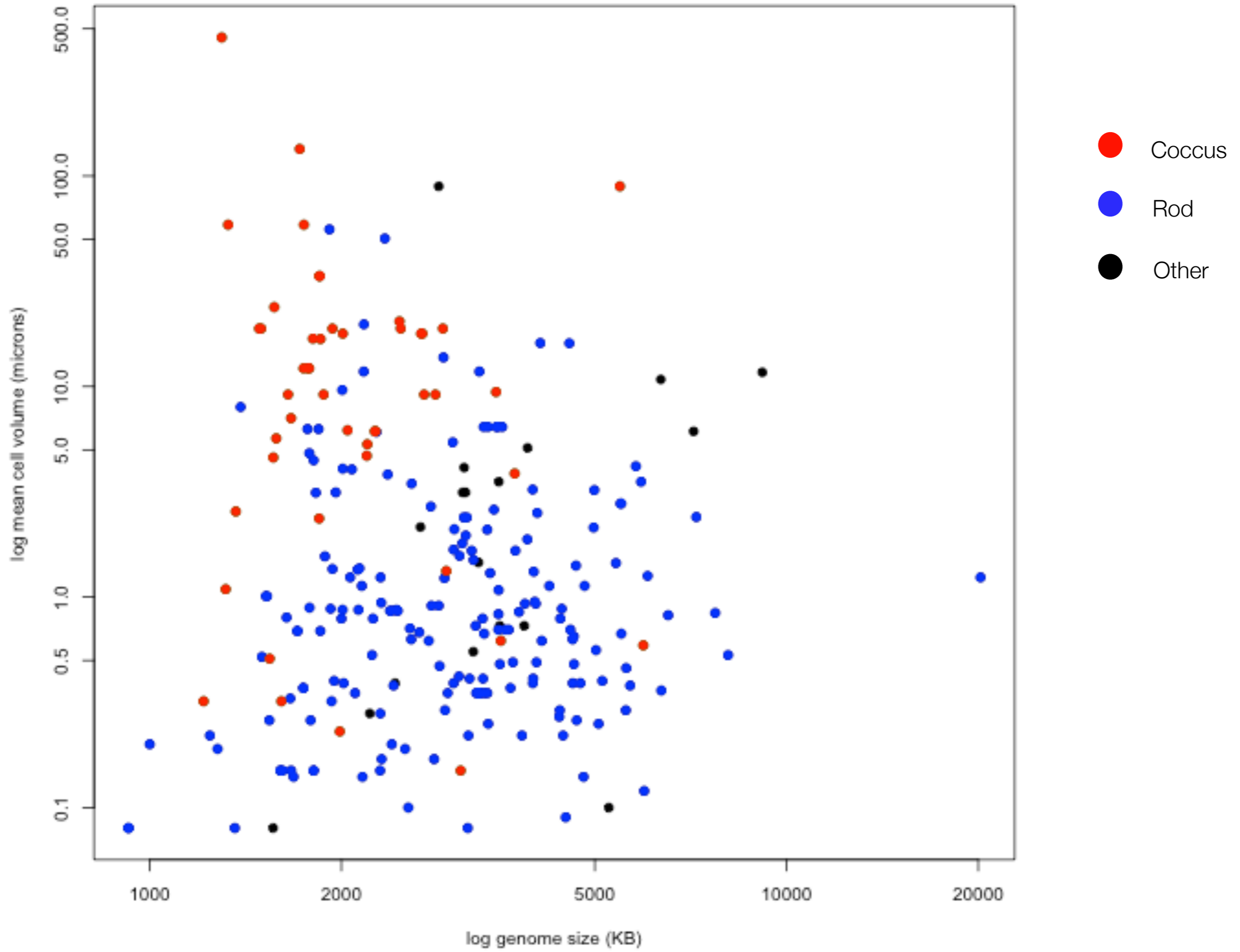
Aerotolerant anaerobe.

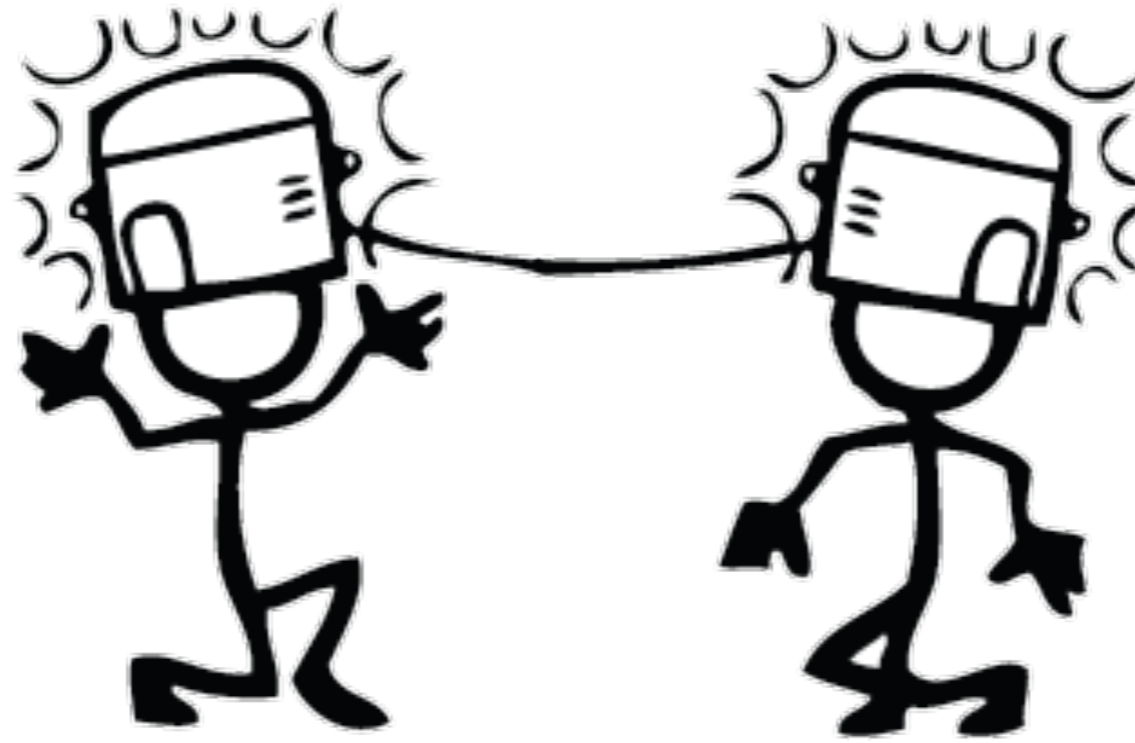
| Fatty Acid | Extracted |
|-------------------|------------------|
| Tokenized | 733 |
| Reduced | 675 |
| Synonyms | 412 |
| Straight chain | 399 |
| Unsaturated | |
| mono | 255 |
| di | 12 |
| tri | 4 |
| tetra-hexa | 6 |
| branched | 315 |
| hydroxy | 141 |
| DMA | 13 |
| cyclo | 4 |

| | |
|----------------------|-----|
| Matched in LipidMaps | 205 |
| Unidentified | 529 |
| Reduced set | 675 |
| Normalized | 412 |
| synonyms (11-21) | 6 |
| synonyms (2-9) | 76 |
| unique | 330 |





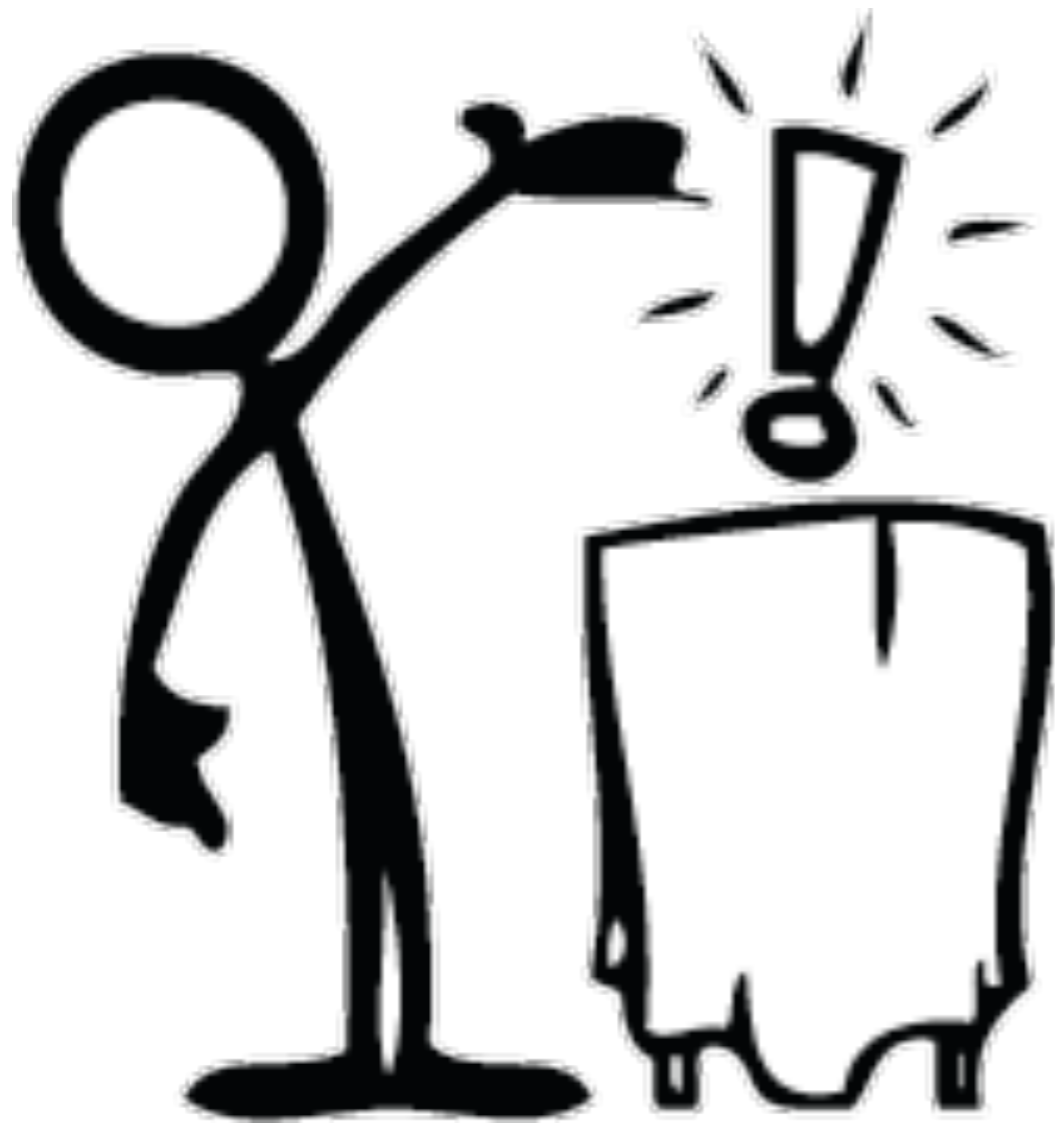




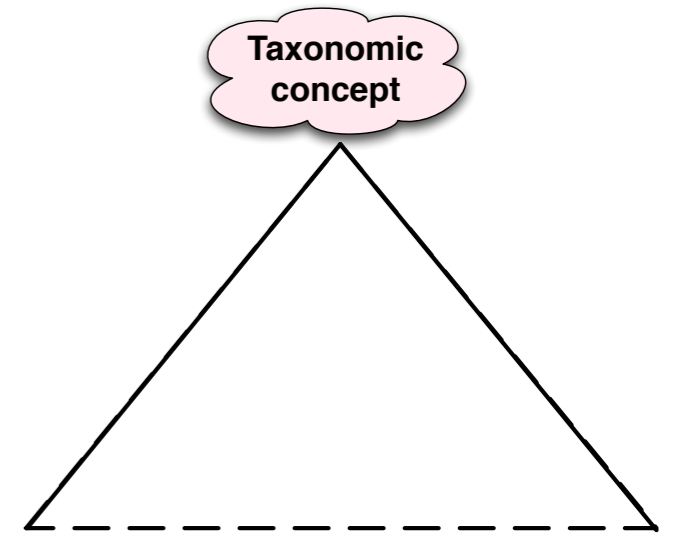
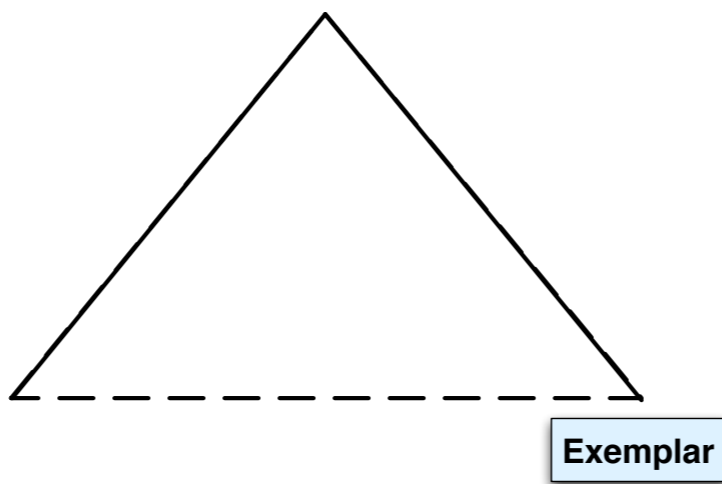
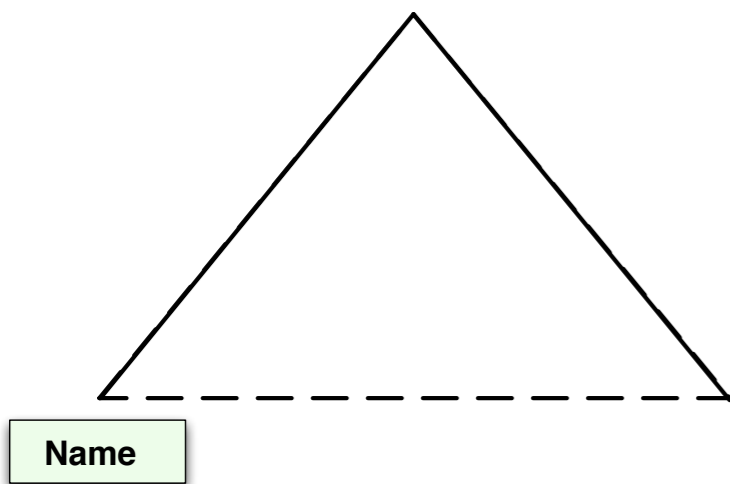
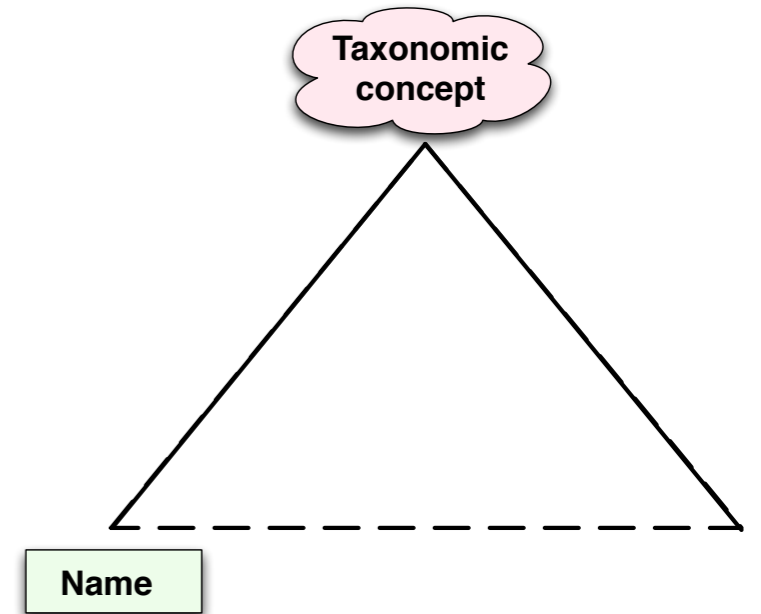
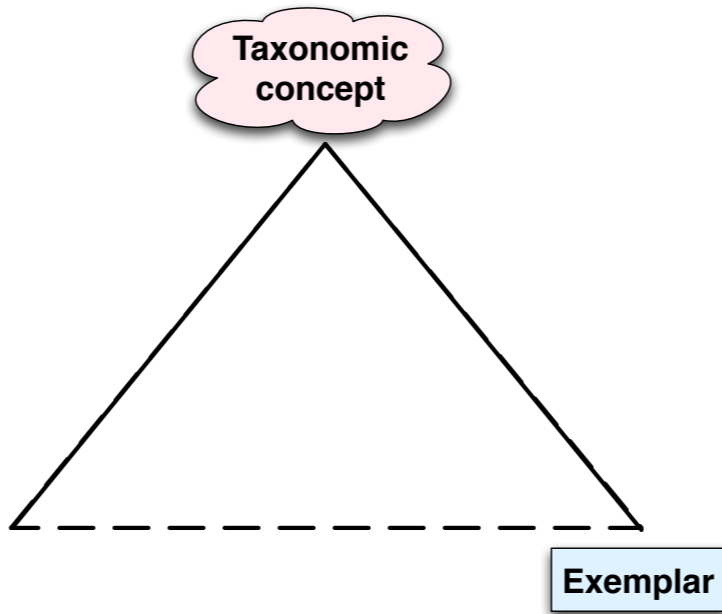
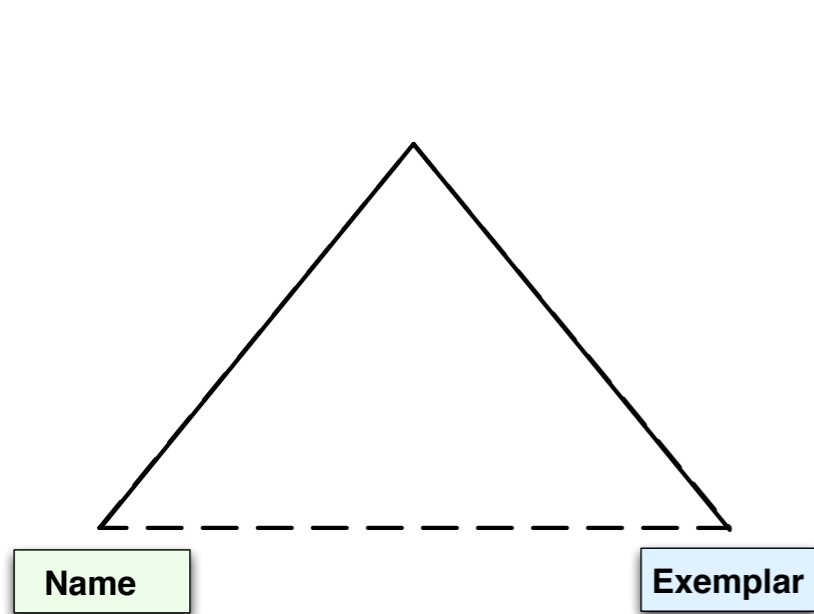
Our goal

To provide on-demand access to the **correct** information so that end-users can interpret like a subject matter expert.

Following up on Folker's bright idea



The reality of modern microbiology





Task description

structure

Phase II

Applications

| | |
|-----------------------------|-----------------------|
| Discriminative Feature App | Description Generator |
| Strain Registration Service | |

Develop initial commercial applications.

Tablet client application and web service.

Develop a strain registration service for new exemplars and Phenotypic Profiles.

Based on existing N4L data model and Phenotypic Data Repository.

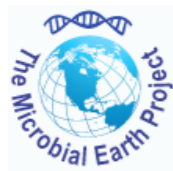
The parts are falling into place,
so what would it take?



Catherine Lyons
Charles T. Parker
Dorothea Taylor
Sarah Wigley
Nicole Osier
Kara Mannor
Grace Rodreguez

**INTERNATIONAL
COMMITTEE ON
SYSTEMATICS OF
PROKARYOTES**

Brian Tindall
Aharon Oren



Nikos Kyrpides
Hans-Peter Klenk

SGM

Ron Fraser
Robin Dunford
Karen Rowlett

ATCC™

Timothy G. Lilburn

SIGS

Oranmiyan Nelson

Beta-testers

MICHIGAN STATE
UNIVERSITY

James R. Cole
Jordan Fish
Xiong Wang
Donna MacGarell

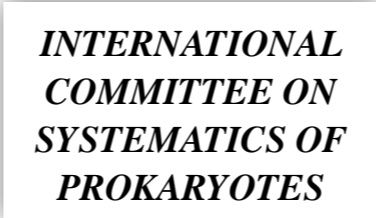
UNIVERSITY OF
Nebraska
Lincoln

Khalid Sayood
Ufuk Nalbantoglu
Sam Way

DTU

Dave Ussery

Thanks to our partners, collaborators and sponsors



Your logo could go here

The NamesforLife system and methods of semantic resolution are covered under US Patent 7,925,444. Other US and PCT patents pending.

