

1 **HydDB: A web tool for hydrogenase**

2 **classification and analysis**

3 **Dan Søndergaard^a, Christian N. S. Pedersen^a, Chris Greening^{b, c*}**

4

5 ^a Aarhus University, Bioinformatics Research Centre, C.F. Møllers Allé 8, Aarhus
6 DK-8000, Denmark

7 ^b The Commonwealth Scientific and Industrial Research Organisation, Land and
8 Water Flagship, Clunies Ross Street, Acton, ACT 2060, Australia

9 ^c Monash University, School of Biological Sciences, Clayton, VIC 2800, Australia

10

11 **Correspondence:**

12 Dr Chris Greening (chris.greening@monash.edu), Monash University, School of
13 Biological Sciences, Clayton, VIC 2800, Australia

14 Dan Søndergaard (das@birc.au.dk), Aarhus University, Bioinformatics Research
15 Centre, C.F. Møllers Allé 8, Aarhus DK-8000, Denmark

16 **Abstract**

17 H₂ metabolism is proposed to be the most ancient and diverse mechanism of
18 energy-conservation. The metalloenzymes mediating this metabolism,
19 hydrogenases, are encoded by over 60 microbial phyla and are present in all major
20 ecosystems. We developed a classification system and web tool, HydDB, for the
21 structural and functional analysis of these enzymes. We show that hydrogenase
22 function can be predicted by primary sequence alone using an expanded
23 classification scheme (comprising 29 [NiFe], 8 [FeFe], and 1 [Fe] hydrogenase
24 classes) that defines 11 new classes with distinct biological functions. Using this
25 scheme, we built a web tool that rapidly and reliably classifies hydrogenase primary
26 sequences using a combination of *k*-nearest neighbors' algorithms and CDD
27 referencing. Demonstrating its capacity, the tool reliably predicted hydrogenase
28 content and function in 12 newly-sequenced bacteria, archaea, and eukaryotes.
29 HydDB provides the capacity to browse the amino acid sequences of 3248
30 annotated hydrogenase catalytic subunits and also contains a detailed repository of
31 physiological, biochemical, and structural information about the 38 hydrogenase
32 classes defined here. The database and classifier are freely and publicly available at
33 <http://services.birc.au.dk/hyddb/>

34

35 **Introduction**

36 Microorganisms conserve energy by metabolizing H₂. Oxidation of this high-energy
37 fuel yields electrons that can be used for respiration and carbon-fixation. This
38 diffusible gas is also produced in diverse fermentation and anaerobic respiratory
39 processes ¹. H₂ metabolism contributes to the growth and survival of microorganisms

40 across the three domains of life: chemotrophs and phototrophs, lithotrophs and
41 heterotrophs, aerobes and anaerobes, mesophiles and extremophiles alike ^{1,2}. On
42 the ecosystem scale, H₂ supports microbial communities in most terrestrial, aquatic,
43 and host-associated ecosystems ^{1,3}. It is also proposed that H₂ was the primordial
44 electron donor ^{4,5}. In biological systems, metalloenzymes known as hydrogenases
45 are responsible for oxidizing and evolving H₂ ^{1,6}. Our recent survey showed there is a
46 far greater number and diversity of hydrogenases than previously thought ². It is
47 predicted that over 55 microbial phyla and over a third of all microorganisms harbor
48 hydrogenases ^{2,7}. Better understanding H₂ metabolism and the enzymes that
49 mediate it also has wider implications, particularly in relation to human health and
50 disease ^{3,8}, biogeochemical cycling ⁹, and renewable energy ^{10,11}.

51

52 There are three types of hydrogenase, the [NiFe], [FeFe], and [Fe] hydrogenases,
53 that are distinguished by their metal composition. Whereas the [Fe]-hydrogenases
54 are a small methanogenic-specific family ¹², the [NiFe] and [FeFe] classes are widely
55 distributed and functionally diverse. They can be classified through a hierarchical
56 system into different groups and subgroups/subtypes with distinct biochemical
57 features (e.g. directionality, affinity, redox partners, and localization) and
58 physiological roles (i.e. respiration, fermentation, bifurcation, sensing) ^{1,6}. It is
59 necessary to define the subgroup or subtype of the hydrogenase to predict
60 hydrogenase function. For example, while Group 2a and 2b [NiFe]-hydrogenases
61 share > 35% sequence identity, they have distinct roles as respiratory uptake
62 hydrogenases and H₂ sensors respectively ^{13,14}. Likewise, discrimination between
63 Group A1 and Group A3 [FeFe]-hydrogenases is necessary to distinguish
64 fermentative and bifurcating enzymes ^{2,15}. Building on previous work ^{16,17}, we

65 recently created a comprehensive hydrogenase classification scheme predictive of
66 biological function ². This scheme was primarily based on the topology of
67 phylogenetic trees built from the amino acid sequences of hydrogenase catalytic
68 subunits/domains. It also factored in genetic organization, metal-binding motifs, and
69 functional information. This analysis identified 22 subgroups (within four groups) of
70 [NiFe]-hydrogenases and six subtypes (within three groups) of [FeFe]-hydrogenases,
71 each proposed to have unique physiological roles and contexts ².

72

73 In this work, we build on these findings to develop the first web database for the
74 classification and analysis of hydrogenases. We developed an expanded
75 classification scheme that captures the full sequence diversity of hydrogenase
76 enzymes and predicts their biological function. Using this information, we developed
77 a classification tool based on the *k*-nearest neighbors' (*k*-NN) method. HydDB is a
78 user-friendly, high-throughput, and functionally-predictive tool for hydrogenase
79 classification that operates with precision exceeding 99.8%.

80

81 **Results and Discussion**

82 **A sequence-based classification scheme for hydrogenases**

83 We initially developed a classification scheme to enable prediction of hydrogenase
84 function by primary sequence alone. To do this, we visualized the relationships
85 between all hydrogenases in sequence similarity networks (SSN) ¹⁸, in which nodes
86 represent individual proteins and the distances between them reflect BLAST *E*-
87 values. As reflected by our analysis of other protein superfamilies ^{19,20}, SSNs allow
88 robust inference of sequence-structure-function relationships for large datasets
89 without the problems associated with phylogenetic trees (e.g. long-branch attraction).

90 Consistent with previous phylogenetic analyses ^{2,16,17}, this analysis showed the
91 hydrogenase sequences clustered into eight major groups (Groups 1 to 4 [NiFe]-
92 hydrogenases, Groups A to C [FeFe]-hydrogenases, [Fe]-hydrogenases), six of
93 which separate into multiple functionally-distinct subgroups or subtypes at narrower
94 logE filters (**Figure 1; Figure S1**). The SSNs demonstrated that all [NiFe]-
95 hydrogenase subgroups defined through phylogenetic trees in our previous work ²
96 separated into distinct clusters, which is consistent with our evolutionary model that
97 such hydrogenases diverged from a common ancestor to adopt multiple distinct
98 functions ². The only exception were the Group A [FeFe]-hydrogenases, which as
99 previously-reported ^{2,17}, cannot be classified by sequence alone as they have
100 principally diversified through changes in domain architecture and quaternary
101 structure. It remains necessary to analyze the organization of the genes encoding
102 these enzymes to determine their specific function, e.g. whether they serve
103 fermentative or electron-bifurcating roles.

104

105 The SSN analysis revealed that several branches that clustered together on the
106 phylogenetic tree analysis ² in fact separate into several well-resolved subclades
107 (**Figure 1**). We determined whether this was significant by analyzing the taxonomic
108 distribution, genetic organization, metal-binding sites, and reported biochemical or
109 functional characteristics of the differentiated subclades. On this basis, we concluded
110 that 11 of the new subclades identified are likely to have unique physiological roles.
111 We therefore refine and expand the hydrogenase classification to reflect the
112 hydrogenases are more diverse in both primary sequence and predicted function
113 than accounted for by even the latest classification scheme ². The new scheme

114 comprises 38 hydrogenase classes, namely 29 [NiFe]-hydrogenase subclasses, 8
115 [FeFe]-hydrogenase subtypes, and the monophyletic [Fe]-hydrogenases (**Table 1**).
116
117 Three lineages originally classified as Group 1a [NiFe]-hydrogenases were
118 reclassified as new subgroups, namely those affiliated with Coriobacteria (Group 1i),
119 Archaeoglobi (Group 1j), and Methanosarcinales (Group 1i). Cellular and molecular
120 studies show these enzymes all support anaerobic respiration of H₂, but differ in the
121 membrane carriers (methanophenazine, menaquinone) and terminal electron
122 acceptors (heterodisulfide, sulfate, nitrate) that they couple to^{21,22}. The previously-
123 proposed 4b and 4d subgroups² were dissolved, as the SSN analysis confirmed
124 they were highly polyphyletic. These sequences are reclassified here into five new
125 subgroups: the formate- and carbon monoxide-respiring Mrp-linked complexes
126 (Group 4b)²³, the ferredoxin-coupled Mrp-linked complexes (Group 4d)²⁴, the well-
127 described methanogenic Eha (Group 4h) and Ehb (Group 4i) supercomplexes²⁵,
128 and a more loosely clustered class of unknown function (Group 4g). Enzymes within
129 these subgroups, with the exception of the uncharacterized 4g enzymes, sustain
130 well-described specialist functions in the energetics of various archaea²³⁻²⁵. Three
131 crenarchaeotal hydrogenases were also classified as their own family (Group 2e);
132 these enzymes enable certain crenarchaeotes to grow aerobically on O₂^{26,27} and
133 hence may represent a unique lineage of aerobic uptake hydrogenases currently
134 underrepresented in genome databases. The Group C [FeFe]-hydrogenases were
135 also separated into three main subtypes given they separate into distinct clusters
136 even at relatively broad log*E* values (**Figure 1**); these subtypes are each transcribed
137 with different regulatory elements and are likely to have distinct regulatory roles^{2,17,28}
138 (**Table 1**).

139

140 **HydDB reliably predicts hydrogenase class using the k -NN method and CDD**
141 **referencing**

142 Using this information, we built a web tool to classify hydrogenases. Hydrogenase
143 classification is determined through a three-step process following input of the
144 catalytic subunit sequence. Two checks are initially performed to confirm if the
145 inputted sequence is likely to encode a hydrogenase catalytic subunit/domain. The
146 Conserved Domain Database (CDD) ²⁹ is referenced to confirm that the inputted
147 sequence has a hydrogenase catalytic domain, i.e. “Complex1_49kDa superfamily”
148 (cl21493) (for NiFe-hydrogenases), “Fe_hyd_lg_C superfamily” (cl14953) (for FeFe-
149 hydrogenases), and “HMD” (pfam03201) (for Fe-hydrogenases). A homology check
150 is also performed that computes the BLAST E -value between the inputted sequence
151 and its closest homolog in HydDB. HydDB classifies any inputted sequence that
152 lacks hydrogenase conserved domains or has low homology scores (E -value $> 10^{-5}$)
153 as a non-hydrogenase (**Table S1**).

154

155 In the final step, the sequence is classified through the k -NN method that determines
156 the most similar sequences listed in the HydDB reference database. To determine
157 the optimal k for the dataset, we performed a 5-fold cross-validation for $k = 1 \dots 10$
158 and computed the precision for each k . The results are shown in **Figure 2**. The
159 classifier predicted the classes of the 3248 hydrogenase sequences with 99.8%
160 precision and high robustness when performing a 5-fold cross-validation (as
161 described in the Methods section) for $k = 4$. The six sequences where there were
162 discrepancies between the SSN and k -NN predictions are shown in **Table S2**. The
163 classifier has also been trained to detect and exclude protein families that are

164 homologous to hydrogenases but do not metabolize H₂ (Nuo, Ehr, NARF, HmdII^{1,2})
165 using reference sequences of these proteins (**Table S1**).

166

167 Sequences of the [FeFe] Group A can be classified into functionally-distinct subtypes
168 (A1, A2, A3, A4) based on genetic organization². The classifier can classify such
169 hydrogenases if the protein sequence immediately downstream from the catalytic
170 subunit sequence is provided. The classifier references the CDD to search for
171 conserved domains in the downstream protein sequence. A sequence is classified
172 as [FeFe] Group A2 if one of the domains “GltA”, “GltD”, “glutamate synthase small
173 subunit” or “putative oxidoreductase”, but not “NuoF”, is found in the sequence.

174 Sequences are classified as [FeFe] Group A3 if the domain “NuoF” is found and
175 [FeFe] Group A4 if the domain “HycB” is present. If none of the domains are found,
176 the sequence is classified as A1. These classification rules were determined by
177 collecting 69 downstream protein sequences. The sequences were then submitted to
178 the CDD and the domains which most often occurred in each subtype were
179 extracted.

180

181 In addition to its precision, the classifier is superior to other approaches due to its
182 usability. It is accessible as a free web service at <http://services.birc.au.dk/hyddb/>
183 HydDB allows the users to paste or upload sequences of hydrogenase catalytic
184 subunit sequences in FASTA format and run the classification (**Figure S2**). When
185 analysis has completed, results are presented in a table that can be downloaded as
186 a CSV file (**Figure S3**). This provides an efficient and user-friendly way to classify
187 hydrogenases, in contrast to the previous standard which requires visualization of
188 phylogenetic trees derived from multiple sequence alignments³⁰.

189

190 **HydDB infers the physiological roles of H₂ metabolism**

191 As summarized in **Table 1**, hydrogenase class is strongly correlated with
192 physiological role. As a result, the classifier is capable of predicting both the class
193 and function of a sequenced hydrogenase. To demonstrate this capacity, we used
194 HydDB to analyze the hydrogenases present in 12 newly-sequenced bacteria,
195 archaea, and eukaryotes of major ecological significance. The classifier correctly
196 classified all 24 hydrogenases identified in the sequenced genomes, as validated
197 with SSNs (**Table 2**). On the basis of these classifications, the physiological roles of
198 H₂ metabolism were predicted (**Table 2**). For five of the organisms, these predictions
199 are confirmed or supported by previously published data ^{27,31–34}. Other predictions
200 are in line with metabolic models derived from metagenome surveying ^{35–37}. In some
201 cases, the capacity for organisms to metabolize H₂ was not tested or inferred in
202 previous studies despite the presence of hydrogenases in the sequenced genomes
203 ^{32,38–40}.

204

205 While HydDB serves as a reliable initial predictor of hydrogenase class and function,
206 further analysis is recommended to verify predictions. Hydrogenase sequences only
207 provide organisms with the genetic capacity to metabolise H₂; their function is
208 ultimately modulated by their expression and integration within the cell ^{1,41}. In
209 addition, some classifications are likely to be overgeneralized due to lack of
210 functional and biochemical characterization of certain lineages and sublineages. For
211 example, it is not clear if two distant members of the Group 1h [NiFe]-hydrogenases
212 (*Robiginitalea biformata*, *Sulfolobus islandicus*) perform the same H₂-scavenging
213 functions as the core group ⁹. Likewise, it seems probable that the Group 3a [NiFe]-

214 hydrogenases of Thermococci and Aquificae use a distinct electron donor to the
215 main class⁴². Prominent cautions are included in the enzyme pages in cases such
216 as these. HydDB will be updated when literature is published that influences
217 functional assignments.

218

219 **HydDB contains interfaces for hydrogenase browsing and analyzing**

220 In addition to its classification function, HydDB is designed to be a definitive
221 repository for hydrogenase retrieval and analysis. The database presently contains
222 entries for 3248 hydrogenases, including their NCBI accession numbers, amino acid
223 sequences, hydrogenase classes, taxonomic affiliations, and predicted behavior
224 **(Figure S4)**. To enable easy exploration of the data set, the database also provides
225 access to an interface for searching, filtering, and sorting the data, as well as the
226 capacity to download the results in CSV or FASTA format. There are individual
227 pages for the 38 hydrogenase classes defined here **(Table 1)**, including descriptions
228 of their physiological role, genetic organization, taxonomic distribution, and
229 biochemical features. This is supplemented with a compendium of structural
230 information about the hydrogenases, which is integrated with the Protein Databank
231 (PDB), as well as a library of over 500 literature references **(Figure S5)**.

232

233 **Conclusions**

234 To summarize, HydDB is a definitive resource for hydrogenase classification and
235 analysis. The classifier described here provides a reliable, efficient, and convenient
236 tool for hydrogenase classification and functional prediction. HydDB also provides
237 browsing tools for the rapid analysis and retrieval of hydrogenase sequences.

238 Finally, the manually-curated repository of class descriptions, hydrogenase

239 structures, and literature references provides a deep but accessible resource for
240 understanding hydrogenases.

241

242 **Methods**

243 **Sequence datasets**

244 The database was constructed using the amino acid sequences of all curated non-
245 redundant 3248 hydrogenase catalytic subunits represented in the NCBI RefSeq
246 database in August 2014 ² (**Dataset S1**). In order to test the classification tool,
247 additional sequences from newly-sequenced archaeal and bacterial phyla were
248 retrieved from the Joint Genome Institute's Integrated Microbial Genomes database
249 ⁴³.

250

251 **Sequence similarity networks**

252 Sequence similarity networks (SSNs) ¹⁸ constructed using Cytoscape 4.1 ⁴⁴ were
253 used to visualize the distribution and diversity of the retrieved hydrogenase
254 sequences. In this analysis, each node represents one of the 3248 hydrogenase
255 sequences in the reference database (**Dataset S1**). Each edge represents the
256 sequence similarity between them as determined by *E*-values from all-vs-all BLAST
257 analysis, with all self and duplicate edges removed. Three networks were
258 constructed, namely for the [NiFe]-hydrogenase large subunit sequences (**Dataset**
259 **S2**), [FeFe]-hydrogenase catalytic domain sequences (**Dataset S3**), and [Fe]-
260 hydrogenase sequences (**Dataset S4**). To control the degree of separation between
261 nodes, $\log E$ cutoffs that were incrementally decreased from -5 to -200 until no major
262 changes in clustering was observed. The $\log E$ cutoffs used for the final
263 classifications are shown in **Figure 1** and **Figure S1**.

264

265 **Classification method**

266 The k -NN method is a well-known machine learning method for classification ⁴⁵.

267 Given a set of data points x_1, x_2, \dots, x_N (e.g. sequences) with known labels y_1, y_2, \dots, y_N

268 (e.g. type annotations), the label of a point, x , is predicted by computing the distance

269 from x to x_1, x_2, \dots, x_N and extracting the k labeled points closest to x , i.e. the

270 neighbors. The predicted label is then determined by majority vote of the labels of

271 the neighbors. The distance measure applied here is that of a BLAST search. Thus,

272 the classifier corresponds to a homology search where the types of the top k results

273 are considered. However, formulating the classification method as a machine

274 learning problem allows the use of common evaluation methods to estimate the

275 precision of the method and perform model selection. The classifier was evaluated

276 using k -fold cross-validation. The dataset is first split in to k parts of equal size. $k - 1$

277 parts (the *training set*) are then used for training the classifier and the labels of the

278 data points in the remaining part (the *test set*) are then predicted. This process,

279 called a *fold*, is repeated k times. The predicted labels of each fold are then

280 compared to the known labels and a precision can be computed.

281

282 **Acknowledgements**

283 We thank A/Prof Colin J. Jackson, Dr Hafna Ahmed, Dr Andrew Warden, Dr Stephen

284 Pearce, and the two anonymous reviewers for their helpful advice and comments

285 regarding this manuscript. This work was supported by a PUMPkin Centre of

286 Excellence PhD Scholarship awarded to DS and a CSIRO Office of the Chief

287 Executive Postdoctoral Fellowship awarded to CG.

288

289 Author Contributions

290 CG and DS designed experiments. DS and CG performed experiments. CG, DS,
291 and CNSP analyzed data. CNSP supervised students. CG and DS wrote the paper.

292

293 Competing financial interests

294 The authors declare no competing financial interests.

295

296 References

- 297 1. Schwartz, E., Fritsch, J. & Friedrich, B. *H₂-metabolizing prokaryotes*. (Springer Berlin
298 Heidelberg, 2013).
- 299 2. Greening, C. *et al.* Genome and metagenome surveys of hydrogenase diversity indicate H₂ is
300 a widely-utilised energy source for microbial growth and survival. *ISME J.* **10**, 761–777 (2016).
- 301 3. Cook, G. M., Greening, C., Hards, K. & Berney, M. in *Advances in Bacterial Pathogen Biology*
302 (ed. Poole, R. K.) **65**, 1–62 (Academic Press, 2014).
- 303 4. Lane, N., Allen, J. F. & Martin, W. How did LUCA make a living? Chemiosmosis in the origin of
304 life. *BioEssays* **32**, 271–280 (2010).
- 305 5. Weiss, M. C. *et al.* The physiology and habitat of the last universal common ancestor. *Nat.*
306 *Microbiol.* **1**, 16116 (2016).
- 307 6. Lubitz, W., Ogata, H., Rüdiger, O. & Reijerse, E. Hydrogenases. *Chem. Rev.* **114**, 4081–148
308 (2014).
- 309 7. Peters, J. W. *et al.* [FeFe]- and [NiFe]-hydrogenase diversity, mechanism, and maturation.
310 *Biochim. Biophys. Acta - Mol. Cell Res.* (2014). doi:10.1016/j.bbamcr.2014.11.021
- 311 8. Carbonero, F., Benefiel, A. C. & Gaskins, H. R. Contributions of the microbial hydrogen
312 economy to colonic homeostasis. *Nat Rev Gastroenterol Hepatol* **9**, 504–518 (2012).
- 313 9. Greening, C. *et al.* Atmospheric hydrogen scavenging: from enzymes to ecosystems. *Appl.*

- 314 *Environ. Microbiol.* **81**, 1190–1199 (2015).
- 315 10. Levin, D. B., Pitt, L. & Love, M. Biohydrogen production: prospects and limitations to practical
316 application. *Int. J. Hydrogen Energy* **29**, 173–185 (2004).
- 317 11. Cracknell, J. A., Vincent, K. A. & Armstrong, F. A. Enzymes as working or inspirational
318 catalysts for fuel cells and electrolysis. *Chem. Rev.* **108**, 2439–2461 (2008).
- 319 12. Shima, S. *et al.* The crystal structure of [Fe]-Hydrogenase reveals the geometry of the active
320 site. *Science* **321**, 572–575 (2008).
- 321 13. Lenz, O. & Friedrich, B. A novel multicomponent regulatory system mediates H₂ sensing in
322 *Alcaligenes eutrophus*. *Proc. Natl. Acad. Sci. U. S. A.* **95**, 12474–12479 (1998).
- 323 14. Greening, C., Berney, M., Hards, K., Cook, G. M. & Conrad, R. A soil actinobacterium
324 scavenges atmospheric H₂ using two membrane-associated, oxygen-dependent [NiFe]
325 hydrogenases. *Proc. Natl. Acad. Sci. U. S. A.* **111**, 4257–4261 (2014).
- 326 15. Schuchmann, K. & Müller, V. A bacterial electron-bifurcating hydrogenase. *J. Biol. Chem.* **287**,
327 31165–31171 (2012).
- 328 16. Vignais, P. M., Billoud, B. & Meyer, J. Classification and phylogeny of hydrogenases. *FEMS*
329 *Microbiol. Rev.* **25**, 455–501 (2001).
- 330 17. Calusinska, M., Happe, T., Joris, B. & Wilmotte, A. The surprising diversity of clostridial
331 hydrogenases: a comparative genomic perspective. *Microbiology* **156**, 1575–1588 (2010).
- 332 18. Atkinson, H. J., Morris, J. H., Ferrin, T. E. & Babbitt, P. C. Using sequence similarity networks
333 for visualization of relationships across diverse protein superfamilies. *PLoS One* **4**, e4345
334 (2009).
- 335 19. Ahmed, F. H. *et al.* Sequence-structure-function classification of a catalytically diverse
336 oxidoreductase superfamily in mycobacteria. *J. Mol. Biol.* **427**, 3554–3571 (2015).
- 337 20. Ney, B. *et al.* The methanogenic redox cofactor F₄₂₀ is widely synthesized by aerobic soil
338 bacteria. *ISME J.* In press (2016).
- 339 21. Stetter, K. O. *Archaeoglobus fulgidus* gen. nov., sp. nov.: a new taxon of extremely
340 thermophilic archaeobacteria. *Syst. Appl. Microbiol.* **10**, 172–173 (1988).

- 341 22. Deppenmeier, U. & Blaut, M. Analysis of the *vhoGAC* and *vhtGAC* operons from
342 *Methanosarcina mazei* strain Gö1, both encoding a membrane-bound hydrogenase and a
343 cytochrome *b*. *Eur. J. Biochem.* **269**, 261–269 (1995).
- 344 23. Kim, Y. J. *et al.* Formate-driven growth coupled with H₂ production. *Nature* **467**, 352–5 (2010).
- 345 24. McTernan, P. M. *et al.* Intact functional fourteen-subunit respiratory membrane-bound [NiFe]-
346 hydrogenase complex of the hyperthermophilic archaeon *Pyrococcus furiosus*. *J. Biol. Chem.*
347 **289**, 19364–19372 (2014).
- 348 25. Lie, T. J. *et al.* Essential anaplerotic role for the energy-converting hydrogenase Eha in
349 hydrogenotrophic methanogenesis. *Proc. Natl. Acad. Sci. U. S. A.* **109**, 15473–8 (2012).
- 350 26. Auernik, K. S. & Kelly, R. M. Physiological versatility of the extremely thermoacidophilic
351 archaeon *Metallosphaera sedula* supported by transcriptomic analysis of heterotrophic,
352 autotrophic, and mixotrophic growth. *Appl. Environ. Microbiol.* **76**, 931–935 (2010).
- 353 27. Giaveno, M. A., Urbietta, M. S., Ulloa, J. R., González Toril, E. & Donati, E. R. Physiologic
354 versatility and growth flexibility as the Main characteristics of a novel thermoacidophilic
355 *Acidianus* strain isolated from Copahue geothermal area in Argentina. *Microb. Ecol.* **65**, 336–
356 346 (2012).
- 357 28. Poudel, S. *et al.* Unification of [FeFe]-hydrogenases into three structural and functional groups.
358 *Biochim. Biophys. Acta (BBA)-General Subj.* doi:10.1016/j.bbagen.2016.05.034 (2016).
- 359 29. Marchler-Bauer, A. & Bryant, S. H. CD-Search: protein domain annotations on the fly. *Nucleic*
360 *Acids Res.* **32**, W327–31 (2004).
- 361 30. Berney, M., Greening, C., Hards, K., Collins, D. & Cook, G. M. Three different [NiFe]
362 hydrogenases confer metabolic flexibility in the obligate aerobe *Mycobacterium smegmatis*.
363 *Environ. Microbiol.* **16**, 318–330 (2014).
- 364 31. Greening, C. *et al.* Persistence of the dominant soil phylum *Acidobacteria* by trace gas
365 scavenging. *Proc. Natl. Acad. Sci.* **112**, 10497–10502 (2015).
- 366 32. Chen, Z. *et al.* *Phaeodactylibacter xiamenensis* gen. nov., sp. nov., a member of the family
367 *Saprospiraceae* isolated from the marine alga *Phaeodactylum tricornutum*. *Int. J. Syst. Evol.*

- 368 *Microbiol.* **64**, 3496–3502 (2014).
- 369 33. Koch, H. *et al.* Growth of nitrite-oxidizing bacteria by aerobic hydrogen oxidation. *Science* **345**,
370 1052–1054 (2014).
- 371 34. Carere, C. R. *et al.* Growth and persistence of methanotrophic bacteria by aerobic hydrogen
372 respiration. *Proc. Natl. Acad. Sci. U. S. A.* (2016).
- 373 35. Haroon, M. F. *et al.* Anaerobic oxidation of methane coupled to nitrate reduction in a novel
374 archaeal lineage. *Nature* **500**, 567–70 (2013).
- 375 36. Evans, P. N. *et al.* Methane metabolism in the archaeal phylum *Bathyarchaeota* revealed by
376 genome-centric metagenomics. *Science* **350**, 434–438 (2015).
- 377 37. Brown, C. T. *et al.* Unusual biology across a group comprising more than 15% of domain
378 *Bacteria*. *Nature* **523**, 208–211 (2015).
- 379 38. Spang, A. *et al.* Complex archaea that bridge the gap between prokaryotes and eukaryotes.
380 *Nature* **521**, 173–179 (2015).
- 381 39. Eloë-Fadrosch, E. A. *et al.* Global metagenomic survey reveals a new bacterial candidate
382 phylum in geothermal springs. *Nat Commun* **7**, (2016).
- 383 40. Wilson, M. C. *et al.* An environmental bacterial taxon with a large and distinct metabolic
384 repertoire. *Nature* **506**, 58–62 (2014).
- 385 41. Greening, C. & Cook, G. M. Integration of hydrogenase expression and hydrogen sensing in
386 bacterial cell physiology. *Curr. Opin. Microbiol.* **18**, 30–8 (2014).
- 387 42. Greening, C. *et al.* Physiology, biochemistry, and applications of F₄₂₀- and F_o-dependent redox
388 reactions. *Microbiol. Mol. Biol. Rev.* **80**, 451–493 (2016).
- 389 43. Markowitz, V. M. *et al.* IMG: the integrated microbial genomes database and comparative
390 analysis system. *Nucleic Acids Research* **40**, D115–22 (2012).
- 391 44. Shannon, P. *et al.* Cytoscape: a software environment for integrated models of biomolecular
392 interaction networks. *Genome Res.* **13**, 2498–2504 (2003).
- 393 45. Cover, T. & Hart, P. Nearest neighbor pattern classification. *IEEE Trans. Inf. Theory* **13**,

394 (1967).

395 46. Constant, P., Chowdhury, S. P., Pratscher, J. & Conrad, R. Streptomyces contributing to
396 atmospheric molecular hydrogen soil uptake are widespread and encode a putative high-
397 affinity [NiFe]-hydrogenase. *Environ. Microbiol.* **12**, 821–829 (2010).

398 47. Hamann, E. *et al.* Environmental *Breviatea* harbour mutualistic *Arcobacter* epibionts. *Nature*
399 **534**, 254–258 (2016).

400 48. Sousa, F. L., Neukirchen, S., Allen, J. F., Lane, N. & Martin, W. F. Lokiarchaeon is hydrogen
401 dependent. *Nat. Microbiol.* **1**, 16034 (2016).

402

403 **Figure Legends**

404 **Figure 1.** Sequence similarity network of hydrogenase sequences. Nodes represent
405 individual proteins and the edges show the BLAST E -values between them at the
406 $\log E$ filter defined at the bottom-left of each panel. The sequences are colored by
407 class as defined in the legends. **Figure S1** shows the further delineation of the
408 encircled [NiFe] hydrogenase classes.

409

410 **Figure 2.** Evaluating the k -NN classifier for $k = 1 \dots 10$. For each k , a 5-fold cross-
411 validation was performed. The mean precision \pm two standard deviations of the folds
412 is shown in the figure (note the y -axis). $k = 1$ provides the most accurate classifier.
413 However, $k = 4$ provides almost the same precision and is more robust to errors in
414 the training set (reflected by the lower standard deviation). In general, the standard
415 deviation is very small, indicating that the predictions are robust to changes in the
416 training data.

417

418 **Tables**

419 **Table 1.** Expanded classification scheme for hydrogenase enzymes. The majority of
420 the classes were defined in previous work ^{2,16,17,46}. The [NiFe] Group 1i, 1j, 1j, 2e,
421 4d, 4g, 4h, and 4i enzymes and [FeFe] Groups C1, C2, and C3 enzymes were
422 defined in this work based on their separation into distinct clusters in the SSN
423 analysis (**Figure 1**). HydDB contains detailed information on each of these classes,
424 including their taxonomic distribution, genetic organization, biochemistry, and
425 structures, as well a list of primary references.

[NiFe] Group 1: Respiratory H₂-uptake [NiFe]-hydrogenases			
1a	Periplasmic	Electron input for sulfate, metal, and organohalide respiration. [NiFeSe] variants.	2
1b	Prototypical	Electron input for sulfate, fumarate, metal, and nitrate respiration.	2
1c	Hyb-type	Electron input for fumarate, nitrate, and sulfate respiration. Physiologically reversible.	2
1d	Oxygen-tolerant	Electron input for aerobic respiration and oxygen-tolerant anaerobic respiration.	2
1e	Isp-type	Electron input primarily for sulfur respiration. Physiologically reversible.	2
1f	Oxygen-protecting	Unresolved role. May liberate electrons to reduce reactive oxygen species.	2
1g	Crenarchaeota-type	Electron input primarily for sulfur respiration.	2
1h	Actinobacteria-type	Electron input for aerobic respiration. Scavenges electrons from atmospheric H ₂ .	2,46
1i	Coriobacteria-type (putative)	Undetermined role. May liberate electrons for anaerobic respiration.	This work
1j	Archaeoglobi-type	Electron input for sulfate respiration ²¹ .	This work
1k	Methanophenazine-reducing	Electron input for methanogenic heterodisulfide respiration ²² .	This work
[NiFe] Group 2: Alternative and sensory uptake [NiFe]-hydrogenases			
2a	Cyanobacteria-type	Electron input for aerobic respiration. Recycles H ₂ produced by other cellular processes.	16
2b	Histidine kinase-linked	H ₂ sensing. Activates two-component system controlling hydrogenase expression.	16
2c	Diguanylate cyclase-linked (putative)	Undetermined role. May sense H ₂ and regulate processes through cyclic di-GMP production.	2
2d	Aquificae-type	Unresolved role. May generate reductant for carbon fixation or have a regulatory role.	2
2e	Metallosphaera-type (putative)	Undetermined role. May liberate electrons primarily for aerobic respiration ²⁶ .	This work
[NiFe] Group 3: Cofactor-coupled bidirectional [NiFe]-hydrogenases			
3a	F ₄₂₀ -coupled	Couples oxidation of H ₂ to reduction of F ₄₂₀ during methanogenesis. Physiologically reversible. [NiFeSe] variants.	16
3b	NADP-coupled	Couples oxidation of NADPH to evolution of H ₂ . Physiologically reversible. May have sulfhydrogenase activity.	16
3c	Heterodisulfide reductase-linked	Bifurcates electrons from H ₂ to heterodisulfide and Fd _{ox} in methanogens. [NiFeSe] variants.	16
3d	NAD-coupled	Interconverts electrons between H ₂ and NAD depending on cellular redox state.	16
[NiFe] Group 4: Respiratory H₂-evolving [NiFe]-hydrogenases			
4a	Formate hydrogenlyase	Couples formate oxidation to fermentative H ₂ evolution. May be H ⁺ -translocating.	2
4b	Formate-respiring	Respires formate or carbon monoxide using H ⁺ as electron acceptor. Na ⁺ -translocating via Mrp ²³ .	This work
4c	Carbon monoxide-respiring	Respires carbon monoxide using H ⁺ as electron acceptor. H ⁺ -translocating.	2

4d	Ferredoxin-coupled, Mrp-linked	Couples Fd _{red} oxidation to H ⁺ reduction. Na ⁺ -translocating via Mrp complex ²⁴ .	This work
4e	Ferredoxin-coupled, Ech-type	Couples Fd _{red} oxidation to H ⁺ reduction. Physiologically reversible via H ⁺ /Na ⁺ translocation.	2
4f	Formate-coupled (putative)	Undetermined role. May couple formate oxidation to H ₂ evolution and H ⁺ translocation.	2
4g	Ferredoxin-coupled (putative)	Undetermined role. May couple Fd _{red} oxidation to proton reduction and H ⁺ /Na ⁺ translocation.	This work
4h	Ferredoxin-coupled, Eha-type	Couples Fd _{red} oxidation to H ⁺ reduction in anaplerotic processes. H ⁺ /Na ⁺ -translocating ²⁵ .	This work
4i	Ferredoxin-coupled, Ehb-type	Couples Fd _{red} oxidation to H ⁺ reduction in anabolic processes. H ⁺ /Na ⁺ -translocating ²⁵ .	This work
[FeFe] Hydrogenases			
A1	Prototypical	Couples ferredoxin oxidation to fermentative or photobiological H ₂ evolution.	2,17
A2	Glutamate synthase-linked (putative)	Undetermined role. May couple H ₂ oxidation to NAD reduction, generating reductant for glutamate synthase.	2,17
A3	Bifurcating	Reversibly bifurcates electrons from H ₂ to NAD and Fd _{ox} in anaerobic bacteria.	2,17
A4	Formate dehydrogenase-linked	Couples formate oxidation to H ₂ evolution. Some bifurcate electrons from H ₂ to ferredoxin and NADP.	2,17
B	Colonic-type (putative)	Undetermined role. May couple Fd _{red} oxidation to fermentative H ₂ evolution.	17
C1	Histidine kinase-linked (putative)	Undetermined role. May sense H ₂ and regulate processes via histidine kinases ² .	This work
C2	Chemotactic (putative)	Undetermined role. May sense H ₂ and regulate processes via methyl-accepting chemotaxis proteins ² .	This work
C3	Phosphatase-linked (putative)	Undetermined role. May sense H ₂ and regulate processes via serine/threonine phosphatases ² .	This work
[Fe] Hydrogenases			
All	Methenyl-H ₄ MPT dehydrogenase	Reversibly couples H ₂ oxidation to 5,10-methenyltetrahydromethanopterin reduction.	16

426

427 **Table 2.** Predictive capacity of the HydDB. HydDB accurately determined hydrogenase content and predicted the physiological
428 roles of H₂ metabolism in 12 newly-sequenced archaeal and bacterial species.

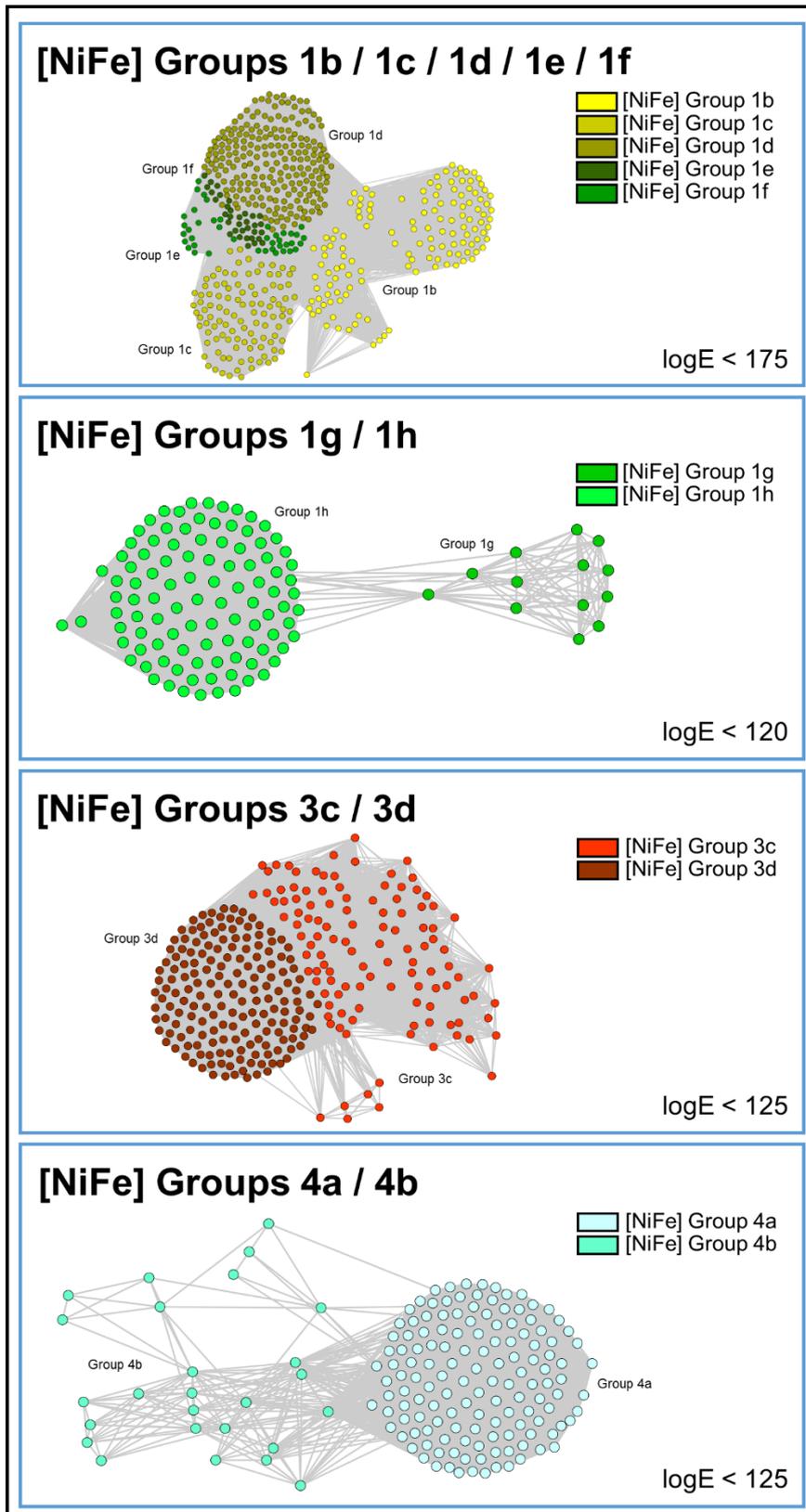
429

430

Organism	Phylum	Hydrogenase accession no.	HydDB classification	SSN classification	Predicted H ₂ metabolism	Confirmed H ₂ metabolism
<i>Pyrinomonas methylaliphatogenes</i>	Acidobacteria	WP_041979300.1	[NiFe] Group 1h	[NiFe] Group 1h	Persistence by aerobic respiration of atmospheric H ₂	Confirmed experimentally ³¹
<i>Phaeodactylibacter xiamenensis</i>	Bacteroidetes	WP_044227713.1 WP_044216927.1 WP_044227053.1	[NiFe] Group 1d [NiFe] Group 2a [NiFe] Group 3d	[NiFe] Group 1d [NiFe] Group 2a [NiFe] Group 3d	Chemolithoautotrophic growth by aerobic H ₂ oxidation	Bacterium grows aerobically, but H ₂ oxidation untested ³²
<i>Bathyarchaeota archaeon BA1</i>	Bathyarchaeota	KPV62434.1 KPV62673.1 KPV62298.1	[NiFe] Group 3c [NiFe] Group 3c [NiFe] Group 4g	[NiFe] Group 3c [NiFe] Group 3c [NiFe] Group 4g	Couples Fd _{red} oxidation to H ₂ evolution in energy-conserving and bifurcating processes	Unconfirmed but consistent with metagenome-based models ³⁶
<i>Lenisia limosa</i>	Obozoa (Breviatea class)	LenisMan28	[FeFe] Group A1	[FeFe] Group A	Fermentative evolution of H ₂	Confirmed experimentally ⁴⁷
<i>Acidianus copahuensis</i>	Crenarchaeota	WP_048100721.1 WP_048100713.1 WP_048100378.1 WP_048100359.1	[NiFe] Group 1g [NiFe] Group 1g [NiFe] Group 1h [NiFe] Group 2e	[NiFe] Group 1g [NiFe] Group 1g [NiFe] Group 1h [NiFe] Group 2e	Chemolithoautotrophic growth by H ₂ oxidation using O ₂ or S ₀ as electron acceptors	Partially confirmed experimentally ²⁷
<i>Arcobacter</i> sp. E1/2/3	Proteobacteria (Epsilon class)	Arc.peg.2312	[NiFe] Group 1b	[NiFe] Group 1b	Chemolithoautotrophic growth by anaerobic H ₂ oxidation	Confirmed experimentally ⁴⁷
<i>Methanoperedens nitroreducens</i>	Euryarchaeota (ANME)	WP_048088262.1 WP_048090768.1	[NiFe] Group 3b [NiFe] Group 3b	[NiFe] Group 3b [NiFe] Group 3b	Secondary role for H ₂ metabolism limited to fermentative evolution of H ₂	Unconfirmed but consistent with metagenome-based models ³⁵
<i>Kryptonium thompsoni</i>	Kryptonionia	CUU03002.1 CUU06124.1	[NiFe] Group 1d [NiFe] Group 3b	[NiFe] Group 1d [NiFe] Group 3b	Chemolithoautotrophic growth by aerobic H ₂ oxidation, fermentative evolution of H ₂ .	Untested, candidate phylum identified by metagenomics ³⁹
<i>Lokiarchaeum</i> sp. GC14_75	Lokiarchaeota	KKK40681.1	[NiFe] Group 3c	[NiFe] Group 3c	Bifurcates electrons between H ₂ , heterodisulfide, and ferredoxin	Unconfirmed but consistent with metagenome-based models ⁴⁸
<i>Nitrospira moscoviensis</i>	Nitrospirae	WP_053379275.1	[NiFe] Group 2a	[NiFe] Group 2a	Chemolithoautotrophic growth by aerobic H ₂ oxidation	Confirmed experimentally ³³
<i>Bacterium</i> GW2011_GWE1_35_17	Moranbacteria	KKQ46070.1 KKQ45273.1	[NiFe] Group 1a [NiFe] Group 3b	[NiFe] Group 1a [NiFe] Group 3b	Chemolithoautotrophic growth by anaerobic H ₂ oxidation, fermentative evolution of H ₂ .	Unconfirmed but consistent with metagenome-based models ³⁷
<i>Bacterium</i> GW2011_GWA2_33_10	Peregrinibacteria	KKP36897.1	[FeFe] Group A3	[FeFe] Group A	Bifurcates electrons between H ₂ , NADH, and ferredoxin	Unconfirmed but consistent with metagenome-based models ³⁷
<i>Entotheonella</i> sp. TSY1	Tectomicrobia	ETW97737.1 ETW94065.1	[NiFe] Group 1h [NiFe] Group 3b	[NiFe] Group 1h [NiFe] Group 3b	Persistence by aerobic respiration of atmospheric H ₂ , fermentative evolution of H ₂	Untested, candidate phylum identified by metagenomics ⁴⁰

432 **Supporting Information**

433 **Figure S1.** Sequence similarity networks showing the relationships between closely
434 related subgroups of [NiFe]-hydrogenases as narrow $\log E$ filters.



435

Figure S2. Screenshot showing interface of HydDB classification page.

HydDB  Browse  Information Pages 

Classify

HydDB provides access to an accurate classifier for hydrogenase sequences and a curated database of hydrogenases by known type. The service is provided by the School of Biological Sciences, Monash University and the Bioinformatics Research Centre, Aarhus University.

Classify

Sequences

Sequences File

No file chosen

Check sequences using CDD?
If enabled, HydDB will use CDD to check whether the submitted sequences encode catalytic subunits of putative before classification. Since this step is time-consuming, you may want to uncheck this option if you are certain your sequences encode hydrogenase catalytic subunits.

Mail

If an e-mail address is provided, a mail will be sent when the job succeeds or fails.

Instructions

To use the classifier to predict the type of one or more hydrogenases from sequence, either:

- paste your FASTA-formatted protein sequences into the text area, or
- upload a FASTA-formatted file with your protein sequences.

Press the "Submit" button to upload the sequences and begin the classification.

If you provided an e-mail address you will receive an e-mail when your job finishes or fails including a link to the results. You will also be able to download the results as a CSV file.

Only sequences encoding the catalytic subunits of hydrogenases will be classified, i.e. those binding the [NiFe]-centre (NiFe-hydrogenases), [FeFe]-centre (FeFe-hydrogenases), or [Fe]-centre (Fe-hydrogenases). Electron-transfer subunits, accessory proteins, and maturation factors cannot be classified by this service.

Limits

A job can at most run for 2 hours. This should be enough for about 2500 sequences to be classified. Results will be stored for 2 weeks. However, we recommend to download the results as they may be deleted due to the rare event of a power outage or server crash.

Statistics

Jobs completed in total	40
Sequences classified in total	232
Jobs completed in the last 24 hours	0
Sequences classified in the last 24 hours	0

Figure S3. Screenshot showing the information provided in the data entry pages for 3248 individual hydrogenases in HydDB.

HydDB
Classify
Browse
Information Pages ▾

Entry WP_004030875.1

Phylum	Euryarchaeota
Order	Methanobacteriales
Organism	Methanobacterium formicicum
Hydrogenase	[Fe]
Activity (Predicted)	Bidirectional
Oxygen Tolerance (Predicted)	Tolerant
Subunits (Predicted)	1
Metal Centres (Predicted)	Fe ion
Accessory Subunits (Predicted)	None

```

MKLAILGAGCYRTHAASGITNFSRACEVAEQVGKPEIAMTHSTIAMGAEKELAGIDEIVVSDPVFDNDFTVIDDFEYEAVIEAHKDPESIMPQIREKVNVAKDLKPKPPKG
AIHFTHPEDLGFVETTTDDNEAVQDADWMTWFPKGMQMGIIKEFADNLKEGAILTHACTVPTTFQKIFEDLSSDEMNIAPKVNVSYPHGAVPEMKGGVYIAEGYASEDAI
CKLVDWGVAAARGDAFKLPAELLPVCDMCSALTAITYAGILSYRDSVMNIIIGAPAGFAQWIAKESLTQVTDLMNSVGDHMEKLDPGALLGTADSMNFGAAADVLPVLEVL
ENRKGKGP TCNI
                    
```

Figure S4. Screenshot showing the capacity for browsing hydrogenase data entries in HydDB.

NCBI Accession	Organism	Hydrogenase Class	Phylum	Order	Activity (Predicted)	Oxygen Tolerance (Predicted)	Subunits (Predicted)	Metal Centres (Predicted)	Accessory Subunits (Predicted)
WP_004030875.1	Methanobacterium formicum	[Fe]	Euryarchaeota	Methanobacteriales	Bidirectional	Tolerant	1	Fe ion	None
WP_012955328.1	Methanobrevibacter ruminantium	[Fe]	Euryarchaeota	Methanobacteriales	Bidirectional	Tolerant	1	Fe ion	None
WP_019263574.1	Methanobrevibacter smithii	[Fe]	Euryarchaeota	Methanobacteriales	Bidirectional	Tolerant	1	Fe ion	None
WP_016357634.1	Methanobrevibacter sp. AbM4	[Fe]	Euryarchaeota	Methanobacteriales	Bidirectional	Tolerant	1	Fe ion	None
WP_013296316.1	Methanothermobacter marburgensis	[Fe]	Euryarchaeota	Methanobacteriales	Bidirectional	Tolerant	1	Fe ion	None
WP_010876766.1	Methanothermobacter thermoautotrophicus	[Fe]	Euryarchaeota	Methanobacteriales	Bidirectional	Tolerant	1	Fe ion	None
WP_013413799.1	Methanothermus fervidus	[Fe]	Euryarchaeota	Methanobacteriales	Bidirectional	Tolerant	1	Fe ion	None

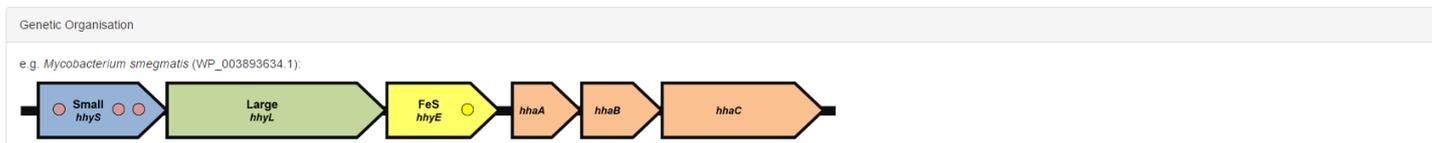
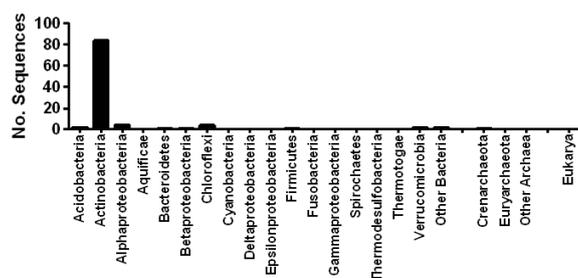
Figure S5. Screenshot showing the detailed content of the information pages about each hydrogenase class on HydDB. Equivalent information pages are available for all 38 hydrogenase classes defined in this work (**Table 1**).

[NiFe] Group 1h-hydrogenase

This entry was last updated at: June 13, 2016, 11:11 a.m.

Properties	
Group	[NiFe] Group 1: Respiratory H ₂ -uptake [NiFe] hydrogenases
Subgroup	[NiFe] Group 1h: Actinobacteria-type
Function	Hydrogenotrophic respiration using O ₂ as terminal electron acceptor. Enzyme scavenges electrons from atmospheric H ₂ to fuel respiratory chain during carbon-starvation. Route of electron transfer unresolved.
Activity	H ₂ -uptake (unidirectional, high-affinity)
Oxygen tolerance	O ₂ -tolerant or O ₂ -insensitive
Localisation	Membrane-associated?

Distribution	
Ecosystem distribution	Upland soils, plant tissues, possibly surface waters
Taxonomic distribution	Widespread among obligately aerobic soil bacteria, especially Actinobacteria, Acidobacteria, and Chloroflexi



Architecture	
Structures	5AAS (<i>Ralstonia eutropha</i> , 2.5 Å resolution, active)
Subunits	3?
Subunit description	HhyL (hydrogenase large subunit) HhyS (hydrogenase small subunit) HhyE (putative iron-sulfur protein and proposed physiological electron acceptor)
Catalytic site	[NiFe]-centre
FeS clusters	Proximal: 3Cys1Asp[4Fe4S] Medial: 4Cys[4Fe4S] Distal: 3Cys1His[4Fe4S]

Important Notes

The *Robiginitalea biformata* and *Sulfolobus islandicus* enzymes are relatively to distantly related to the main group. No studies have yet tested whether these enzymes have a H₂-scavenging role like other Group 1h [NiFe]-hydrogenases. They may instead represent founding members of a functionally-distinct lineage.

NCBI Accession	Organism	Hydrogenase Class	Phylum	Order	Activity (Predicted)	Oxygen Tolerance (Predicted)	Subunits (Predicted)	Metal Centres (Predicted)	Accessory Subunits (Predicted)
WP_014267363.1	<i>Granulicella mallensis</i>	[NiFe] Group 1h	Acidobacteria	Acidobacteriales	Aerobic Uptake	Tolerant	3	[NiFe]-centre, 3 x [4Fe4S] clusters	[FeS] protein
WP_011688202.1	<i>Soilbacter usitatus</i>	[NiFe] Group 1h	Acidobacteria	Solibacterales	Aerobic Uptake	Tolerant	3	[NiFe]-centre, 3 x [4Fe4S] clusters	[FeS] protein
WP_021597135.1	<i>Actinomadura madurae</i>	[NiFe] Group 1h	Actinobacteria	Actinomycetales	Aerobic Uptake	Tolerant	3	[NiFe]-centre, 3 x [4Fe4S] clusters	[FeS] protein
WP_026402909.1	<i>Actinomadura rifamycinii</i>	[NiFe] Group 1h	Actinobacteria	Actinomycetales	Aerobic Uptake	Tolerant	3	[NiFe]-centre, 3 x [4Fe4S] clusters	[FeS] protein
WP_018330638.1	<i>Actinomycetospira changmaiensis</i>	[NiFe] Group 1h	Actinobacteria	Actinomycetales	Aerobic Uptake	Tolerant	3	[NiFe]-centre, 3 x [4Fe4S] clusters	[FeS] protein
WP_007735075.1	<i>Rhodococcus qingshengii</i>	[NiFe] Group 1h	Actinobacteria	Actinomycetales	Aerobic Uptake	Tolerant	3	[NiFe]-centre, 3 x [4Fe4S] clusters	[FeS] protein
WP_003935326.1	<i>Rhodococcus ruber</i>	[NiFe] Group 1h	Actinobacteria	Actinomycetales	Aerobic Uptake	Tolerant	3	[NiFe]-centre, 3 x [4Fe4S] clusters	[FeS] protein
WP_005443931.1	<i>Saccharomonospora azurea</i>	[NiFe] Group 1h	Actinobacteria	Actinomycetales	Aerobic Uptake	Tolerant	3	[NiFe]-centre, 3 x [4Fe4S] clusters	[FeS] protein

« 1 2 3 »

Literature

Genetics:

- Berney, M., Greening, C., Hards, K., Collins, D., and Cook, G.M. (2014) Three different [NiFe] hydrogenases confer metabolic flexibility in the obligate aerobic *Mycobacterium smegmatis*. *Environ. Microbiol.* **16**: 318-330.
- Constant, P., Chowdhury, S.P., Hesse, L., and Conrad, R. (2011) Co-localization of atmospheric H₂ oxidation activity and high affinity H₂-oxidizing bacteria in non-axenic soil and sterile soil amended with *Streptomyces* sp. PCB7. *Soil Biol. Biochem.* **43**: 1888-1893.
- Constant, P., Chowdhury, S.P., Hesse, L., Pratscher, J., and Conrad, R. (2011) Genome data mining and soil survey for the novel group 5 [NiFe]-hydrogenase to explore the diversity and ecological importance of presumptive high-affinity H₂-oxidizing bacteria. *Appl. Environ. Microbiol.* **77**: 6027-6035.
- Greening, C., Bliswas, A., Carere, C.R., Jackson, C.J., Taylor, M.C., Stott, M.B., Cook, G.M., and Morales, S.E. (2016) Genomic and metagenomic surveys of hydrogenase distribution indicate H₂ is a widely utilised energy source for microbial growth and survival. *ISME J.* **10**: 761-777.
- Khdir, M., Hesse, L., Popa, M.E., Quiza, L., Lalonde, I., Meredith, L.K., Röckmann, T., and Constant, P. (2015) Soil carbon content and relative abundance of high affinity H₂-oxidizing bacteria predict atmospheric H₂ soil uptake activity better than soil microbial community composition. *Soil Biol. Biochem.* **85**: 1-9.

Physiology:

- Berney, M., Greening, C., Conrad, R., Jacobs, W.R., and Cook, G.M. (2014) An obligately aerobic soil bacterium activates fermentative hydrogen production to survive reductive stress during hypoxia. *Proc. Natl. Acad. Sci. U. S. A.* **111**: 11479-11484.
- Constant, P., Chowdhury, S.P., Pratscher, J., and Conrad, R. (2010) *Streptomyces* contributing to atmospheric molecular hydrogen soil uptake are widespread and encode a putative high-affinity [NiFe]-hydrogenase. *Environ. Microbiol.* **12**: 821-829.

Table S1. Validation that HydDB classifies only hydrogenase catalytic subunit sequences. HydDB excludes non-hydrogenase sequences through a combination of homology checks (sequences are only classified as hydrogenases if BLAST *E*-value of the closest hit in HydDB is less than 10^{-5}) and CDD checks (sequences are only classified as hydrogenases if signature conserved domains are found). In addition, the classifier has been specifically trained to exclude four protein families that are homologous to hydrogenase catalytic subunits (HmdII, Her, NuoD, NARF) but lack hydrogenase activity.

NCBI Accession	Sequence type	Homology check	CDD check	Final result
WP_041979300.1	Validated hydrogenase catalytic subunit	Highest sequence homology with [NiFe] Group 1h (E = 0)	Ni,Fe-hydrogenase I large subunit (COG0374)	Hydrogenase
WP_011729412.1	P-type ATPase (unrelated to hydrogenases)	Low sequence homology with hydrogenases (E = 5.6)	Non-hydrogenase	Non-hydrogenase
WP_003895387.1	Chaperone (unrelated to hydrogenases)	Low sequence homology with hydrogenases (E = 3.8)	Non-hydrogenase	Non-hydrogenase
WP_013295714.1	HmdII (homologous with [Fe]-hydrogenases)	Highest sequence homology with HmdII (E = 0)	HMD (pfam03201)	Non-hydrogenase
WP_003901794.1	Ehr (homologous with [NiFe]-hydrogenases)	Highest sequence homology with Ehr (E = 0)	Complex1_49kDa superfamily (cl21493)	Non-hydrogenase
WP_003901553.1	NuoD (homologous with [NiFe]-hydrogenases)	Highest sequence homology with NuoD (E = 0)	NuoD (COG0649)	Non-hydrogenase
NP_114174.1	NARF (homologous with [FeFe]-hydrogenases)	Highest sequence homology with NARF (E = 0)	Fe_hyd_Ig_C (pfam02906)	Non-hydrogenase

Table S2. Hydrogenase sequences where there is disagreement between classification by SSN and *k*-NN methods. These sequences represent six out of the total 3248 sequences analyzed, i.e. 0.0018%.

NCBI Accession	Organism	<i>k</i>-NN Classification	SSN Classification
WP_027414715.1	Aneurinibacillus terranovensis	[NiFe] Group 1e	[NiFe] Group 1d
WP_027358538.1	Desulforegula conservatrix	[NiFe] Group 3d	[NiFe] Group 3c
WP_012532312.1	Geobacter bemidjiensis	[NiFe] Group 3d	[NiFe] Group 3c
WP_012469611.1	Geobacter lovleyi	[NiFe] Group 3d	[NiFe] Group 3c
WP_004512544.1	Geobacter metallireducens	[NiFe] Group 3d	[NiFe] Group 3c
WP_015839165.1	Geobacter sp. M21	[NiFe] Group 3d	[NiFe] Group 3c

Dataset S1. Excel spreadsheet listing the sequence, taxonomy, and hydrogenase class of all 3248 hydrogenase catalytic subunit sequences listed in HydDB.

Dataset S2. Zip file containing the Cytoscape network for [NiFe]-hydrogenases.

Dataset S3. Zip file containing the Cytoscape network for [FeFe]-hydrogenases.

Dataset S4. Zip file containing the Cytoscape network for [Fe]-hydrogenases.