

1 **An atomic-resolution view of neofunctionalization in the**
2 **evolution of apicomplexan lactate dehydrogenases**

3
4 Jeffrey I. Boucher¹, Joseph R. Jacobowitz¹, Brian C. Beckett¹, Scott Classen² and Douglas L.
5 Theobald¹

6
7 ¹*Brandeis University, Department of Biochemistry, Waltham MA 02454, USA.*

8 ²*Physical Biosciences Division, Lawrence Berkeley National Laboratory, Berkeley, CA 94720, USA.*

9
10 **Malate and lactate dehydrogenases (MDH and LDH) are homologous, core metabolic**
11 **enzymes that share a fold and catalytic mechanism yet possess strict specificity for their**
12 **substrates. In the Apicomplexa, convergent evolution of an unusual LDH from MDH**
13 **resulted in a difference in substrate preference exceeding 12 orders of magnitude. The**
14 **molecular and evolutionary mechanisms responsible for this extraordinary functional**
15 **shift are currently unknown. Using ancestral sequence reconstruction, we find that the**
16 **evolution of pyruvate specificity in apicomplexan LDHs arose through a classic**
17 **neofunctionalization mechanism characterized by long-range epistasis, a promiscuous**
18 **intermediate, and relatively few gain-of-function mutations of large effect. Residues far**
19 **from the active site determine specificity, as shown by the crystal structures of three**
20 **ancestral proteins that bracket the key gene duplication event. This work provides an**
21 **unprecedented atomic-resolution view of evolutionary trajectories resulting in the *de novo***
22 **creation of a nascent enzymatic function.**

24 **Background**

25 The common ancestor of the eukaryotic Apicomplexa evolved nearly one billion years
26 ago (1), and its modern descendants comprise a large phylum of intracellular parasites that are
27 currently responsible for numerous devastating metazoan diseases, including malaria
28 (*Plasmodium*), toxoplasmosis (*Toxoplasma*), cryptosporidiosis (*Cryptosporidium*),
29 cyclosporiasis (*Cyclospora*), and babesiosis (*Babesia*). A key event in the early evolution of the
30 Apicomplexa was the acquisition of a malate dehydrogenase (MDH) via lateral gene transfer
31 from α -proteobacteria (2-4). Following a gene duplication event roughly 700-900 Mya, one
32 copy of this MDH evolved a novel substrate specificity to become a highly specific lactate
33 dehydrogenase (LDH) that is now essential to the life cycle of many modern apicomplexans (5).
34 As a core metabolic enzyme that evolved independently of metazoan LDH, the unique
35 apicomplexan LDH has attracted significant attention as a potential drug target (6-9). However,
36 the molecular and evolutionary mechanisms that drove this switch in substrate specificity are
37 currently unknown.

38 LDH and MDH are homologous, 2-ketoacid oxidoreductases that share both a protein
39 fold (10) (**Figure 1-figure supplement 1**) and a common catalytic mechanism (11-16)
40 (**Figure 1**). Both enzymes are found in central metabolism: MDH catalyzes the interconversion
41 of oxaloacetate and malate in the citric acid cycle, and LDH converts pyruvate to lactate in the
42 final step of anaerobic glycolysis. Despite their structural and catalytic similarities, modern
43 apicomplexan LDHs and MDHs have extraordinarily strict substrate specificity. For example,
44 *Plasmodium falciparum* (*Pf*) MDH and LDH each prefer their respective substrates by over six
45 orders of magnitude. The biophysical basis for this extraordinary substrate preference is
46 presently an unresolved question.

47 A conspicuous structural difference between apicomplexan MDHs and LDHs is an
48 insertion within the active site loop of the LDHs (17, 18) (**Figure 2**). In the LDH/MDH

49 superfamily, closure of this loop over the active site is rate-limiting during catalysis (16), and
50 mutations within this loop have a large effects on activity and substrate specificity (19). For
51 example, simply mutating Gln102 to Arg in the specificity loop of *Bacillus stearothermophilus*
52 (*Bs*) LDH converts the enzyme into an MDH, shifting specificity from a 10³-fold preference for
53 pyruvate to a 10⁴-fold preference for oxaloacetate (19) (**Figure 1**, residue numbering is based on
54 the dogfish LDH convention (20)). In fact, all known MDHs have an Arg at position 102, while
55 canonical LDHs have a Gln, and consequently residue 102 has been called the “specificity
56 residue” (21). Residue 102 is thought to contribute to substrate discrimination by balancing the
57 substrate charge within the active site: the positively charged Arg in MDHs forms a salt bridge
58 with the C4 carboxylate of oxaloacetate, whereas the neutral Gln in canonical LDHs packs with
59 the C3 methyl of pyruvate (**Figure 1**). Yet, attempts to convert an MDH into an LDH by
60 mutating Arg102 to Gln have met with limited success (22, 23). In the apicomplexan LDHs,
61 residue 102 is not a Gln but a Lys, a relatively conservative substitution compared to the MDH
62 Arg. It is currently not understood why *Plasmodium* LDHs lack activity towards oxaloacetate,
63 despite having a positively charged sidechain at residue 102 similar to MDHs (8, 24-28).

64 Apicomplexan LDH evolved from the duplication of an ancestral MDH gene (3, 4). Gene
65 duplication is widely considered the major force that has driven the evolutionary diversity of
66 protein functions (29). There are three general ways duplicated genes can be fixed in a
67 population by selection: 1) “dosage selection”, beneficial increase in dosage due to multiple
68 copies, 2) “subfunctionalization”, specialization of previously existing functions, or 3)
69 “neofunctionalization”, creation of a novel function through the accumulation of beneficial,
70 gain-of-function mutations (30). Most mutations, however, are either neutral or detrimental. A
71 new duplicated gene typically degrades to a crippled pseudogene before it can acquire the rare
72 beneficial mutations needed to confer a selectable function (31, 32). Hence, classical
73 neofunctionalization has fallen out of favor in preference for models that begin with the
74 duplication of a multifunctional protein, such as “specialization” and “subfunctionalization”

75 models. Currently the molecular and evolutionary mechanisms that create novel functions in
76 gene duplicates are fiercely debated (29, 33-35), and there are no clear examples of classic
77 neofunctionalization or gain-of-function mutations (36-38).

78 The apicomplexan LDH and MDH enzyme family provides an exceptional model system
79 for investigating several long-standing questions in molecular evolution, including the
80 mechanisms available to convergent evolution, the number of mutations required to produce a
81 nascent function, the role of promiscuous intermediates during evolution of function, and the
82 effects of epistasis on evolutionary irreversibility. In order to identify the biophysical and
83 evolutionary mechanisms responsible for pyruvate specificity in apicomplexan LDHs, we have
84 reconstructed ancestral proteins along the evolutionary trajectories leading to modern
85 apicomplexan MDHs and LDHs (**Figure 3B**). We kinetically and structurally characterized the
86 ancestral proteins together with multiple evolutionary intermediates. This work provides a clear
87 example of neofunctionalization in protein evolution and the first crystal structures
88 documenting the evolution of a new enzyme. We show that apicomplexan LDHs evolved as the
89 result of few mutations of large effect via the classic neofunctionalization of a duplicated MDH
90 gene.

91

92 **Results**

93 ***LDH enzymes have evolved independently at least four times***

94 A maximum likelihood phylogeny of representatives of all known LDH and MDH
95 proteins provides strong support for five distinct protein clades (**Figure 3A, Figure 3-figure**
96 **supplement 1**): canonical LDHs, “LDH-like” MDHs, mitochondrial-like MDHs, cytosolic-like
97 MDHs, and the poorly characterized HicDHs (hydroxyisocaproate-related dehydrogenases),
98 confirming previous phylogenetic analyses (2-4, 39).

99 The HicDH clade are close sequence homologs of a known hydroxyisocaproate
100 dehydrogenase. They all possess a residue other than a Gln or an Arg at the “specificity”
101 position 102, as well as insertions of varying lengths within the catalytic loop between residues
102 102 and 109. Despite these alterations within the catalytic loop, all other catalytic residues
103 (Arg109, Asp168, Arg171, and His195) are conserved. Only one taxon within the HicDH clade
104 has been functionally characterized, DHL2_LACCO, which is a specific hydroxyisocaproate
105 dehydrogenase (85). These observations suggest that the clade features dehydrogenases with
106 altered substrate specificity.

107 Except for the HicDHs, which are exclusively eubacterial, both eukaryotic and
108 eubacterial enzymes are found in all major clades (**Figure 3-figure supplement 2**). The
109 “LDH-like” MDH clade additionally contains archaeal dehydrogenases, which are basal and
110 group to the exclusion of the bacterial MDHs.

111 Intriguingly, three different groups of LDH proteins cluster with high confidence outside
112 of the canonical LDH clade. A set of trichomonad LDHs found in the cytosolic-like MDH clade
113 are thought to have evolved from a recent gene duplication of an MDH (40). The
114 Trichomonads appear to lack a canonical LDH. A prominent eukaryotic group of LDH and
115 MDH proteins from the Apicomplexa nests deep within the bacterial “LDH-like” MDHs, sister to
116 many Rickettsiales sequences, signifying a horizontal gene transfer event from α -proteobacteria

117 to the eukaryotic Apicomplexa. We find no evidence that the Apicomplexa have canonical LDH
118 or conventional eukaryotic-type MDH (either cytosolic- or mitochondrial-like MDHs), despite
119 searching in many available complete apicomplexan genomes (multiple *Eimeria*, *Neospora*,
120 *Toxoplasma*, *Plasmodium*, and *Cryptosporidium* species) (41-43). In the Apicomplexa, LDH
121 activity has apparently evolved independently twice (**Figure 3B, Figure 3-figure**
122 **supplement 3**), once in a lineage leading to *Plasmodium*-related species and once in
123 *Cryptosporidium*. The apicomplexan portion of the LDH/MDH gene phylogeny is consistent
124 with recent apicomplexan species phylogenies constructed from concatenated protein sequences
125 (44).

126 We rooted the MDH/LDH phylogeny using the Rossmann fold domain of the distantly
127 related α/β -glucosidases and aspartate dehydrogenases as outgroups. The ML root position
128 apparently splits the tree into two large groups: (1) one which contains the cytosolic- and
129 mitochondrial-like MDHs, which are largely dimeric, and (2) another which contains the
130 canonical LDHs, “LDH-like” MDHs, and HicDHs, which are primarily tetrameric (**Figure 3-**
131 **figure supplement 1**). While the ML root position is robust to variation in taxon coverage,
132 the exact location is poorly supported. Nevertheless, there is strong support for a root position
133 within the central MDH section of the tree and outside of the five identified clades, including the
134 canonical LDH clade (confidence level > 0.99985 according to the aLRT), indicating that the
135 canonical LDHs evolved from an ancestral MDH. The global rooting and the location of the
136 three separate LDH groups, deep within MDH clades, indicate that LDH enzymes have evolved
137 convergently from MDHs at least four times in the superfamily.

138 ***An insertion in the catalytic loop of apicomplexan LDHs***

139 In the present work, our focus is on the convergent evolution of the unusual
140 apicomplexan LDHs. With the α -proteobacterial “LDH-like” MDHs as the closest outgroup, the
141 apicomplexan enzymes are split into two main groups: (1) LDHs belonging to *Toxoplasma*,

142 *Plasmodium*, and related protists, and (2) MDHs belonging to *Plasmodium* and
143 *Cryptosporidium*. Apart from their atypical phylogenetic position, the apicomplexan MDHs
144 appear as typical α -proteobacterial “LDH-like” MDHs, containing all the key catalytic residues
145 including Arg102. The *Cryptosporidium* LDHs are an exception, being nested within the
146 apicomplexan MDH clade partitioned from the rest of the apicomplexan LDHs.
147 *Cryptosporidium* LDHs have a Gly at position 102 and are thought to be a product of an
148 independent, convergent duplication event (39).

149 In contrast, the large apicomplexan LDH clade is demarcated by a unique, conserved
150 five-residue insertion in the active site loop. While the apicomplexan LDH and MDH proteins
151 are moderately divergent, with about 45% sequence identity, the differences are largely confined
152 to exterior residues removed from the active sites. One important difference is that the
153 apicomplexan LDHs have Lys102 for the “specificity residue”, rather than a Gln as found in the
154 canonical LDHs (**Figure 2-figure supplement 1**). Apicomplexan proteins frequently contain
155 numerous insertions relative to proteins from other species (45, 46), a characteristic thought to
156 result from various factors, including high AT genome content, DNA strand slippage, double
157 strand break repair, high recombination rates, and selection pressure for parasite antigenic
158 variation. Except for Met106, the amino acid and coding sequence immediately flanking the
159 apicomplexan LDH loop insertion is largely conserved with α -proteobacterial MDHs (**Figure 2-**
160 **figure supplement 1**). It is therefore likely that a mutation “expanded” the Met106 codon to
161 code for six residues, resulting in the observed five-residue insertion and the Met106Lys
162 mutation. Henceforth we will refer to this expansion mutation as the “six-residue loop
163 insertion”.

164 ***Trp107f is the modern apicomplexan LDH specificity residue***

165 In the modern apicomplexan enzymes, the six-residue insertion in the LDH specificity
166 loop (positions 99-112) induces two significant structural changes relative to MDH (**Figure 2**).

167 First, LDH residue Lys102 is excluded from of the active site, unlike the corresponding Arg102
168 in MDH, which is enclosed within the active site and participates in functionally important
169 interactions with the substrate. Second, LDH Trp107f, which is part of the novel insertion,
170 occupies the same space as Arg102 in MDH (by convention, residues in the insertion are labeled
171 using numbers and letters to maintain consistency with homologous positions in the dogfish
172 LDH, **Figure 2**).

173 The only prominent structural difference between the active sites of the LDH and MDH
174 proteins is the replacement of MDH Arg102 with LDH Trp107f. Trp107f is positioned where it
175 could presumably interact with the distinguishing C3 methyl of the pyruvate substrate, while
176 MDH Arg102 interacts with the C4 carboxylate of oxaloacetate (21). As a bulky, hydrophobic
177 residue, Trp107f could recognize pyruvate in preference to oxaloacetate by two mechanisms: (1)
178 a hydrophobic interaction with the pyruvate C3 methyl vs the negatively charged oxaloacetate
179 methylene carboxylate and (2) steric occlusion of the methylene carboxylate of oxaloacetate.
180 Furthermore, Trp107f is conserved in all apicomplexan LDHs (**Figure 2-figure supplement**
181 **1**), suggesting negative selection and functional importance. We therefore hypothesized that
182 Trp107f plays an important role in pyruvate recognition.

183 We tested the functional importance of residues in the specificity loop in *Pf*LDH with an
184 “alanine scan” by individually mutating each residue in positions 101-108 to an alanine (**Figure**
185 **2-figure supplement 2**, note Ala103 was mutated to a serine). We assessed the activity of the
186 mutants using k_{cat}/K_m , a measure of enzymatic specificity and catalytic efficiency, as determined
187 from steady state kinetic assays. Mutating Trp107f to Ala reduced pyruvate activity by five
188 orders of magnitude, whereas mutations at all other positions had effects less than a single order
189 of magnitude, including the canonical specificity residue at position 102. The Trp107fAla
190 mutation affects both k_{cat} (1500-fold decrease) and K_m (50-fold increase).

191 To assess the effects of Trp107fAla mutation on the specificity loop conformation, we
192 solved the crystal structure of *Pf*LDH-W107fA (1.1Å) in the presence of oxamate and NADH.

193 The protein crystallizes in the same space group as the wild-type *Pf*LDH, with nearly identical
194 cell dimensions (**Figure 2-figure supplement 3A**). In the W107fA mutant, the specificity
195 loop is disordered between residues Leu86 and Arg99, as is often seen in structures in which the
196 loop is in the open conformation. In the mutant, residues 102-105 are in a linear α -helical
197 conformation, in contrast to the wild-type *Pf*LDH closed state which has a very prominent 60°
198 kink in the α -helix at Pro104. Thus, the only significant difference between the wild-type and
199 mutant structures is that the *Pf*LDH-W107fA specificity loop is found in the open conformation,
200 consistent with weaker binding of substrate (**Figure 2-figure supplement 3B**). These results
201 indicate that Trp107f is necessary for pyruvate activity in apicomplexan LDHs, and that it has
202 become the new “specificity residue” despite the fact that Trp107f does not align in sequence
203 with the canonical specificity residue at position 102 (**Figure 2-figure supplement 1**).

204 ***The loop insert fails to swap specificity in modern LDH and MDH***

205 During evolution, the six-residue insertion displaced the canonical specificity residue at
206 position 102 and apparently switched substrate preference in apicomplexan LDHs. If this
207 insertion is sufficient for pyruvate recognition, then adding the insertion to a modern
208 apicomplexan MDH should convert the enzyme to an LDH. To test this hypothesis, we
209 incorporated the six-residue insertion from *Pf*LDH into the catalytic loop of *Pf*MDH (*Pf*MDH-
210 INS) and the *Cryptosporidium parvum* (*Cp*) MDH (*Cp*MDH-INS). The chimeric proteins
211 showed a >100-fold reduction in oxaloacetate activity with no significant gain in pyruvate
212 activity (**Figure 4**). Like other MDHs, the apicomplexan MDHs have an Arg at position 102
213 that is important for oxaloacetate recognition; in the modern apicomplexan LDHs position 102
214 is a Lys. The Arg102Lys mutation may be necessary to eliminate oxaloacetate activity and
215 increase pyruvate activity. Therefore, we also mutated Arg102 to Lys in the *Pf*MDH chimera
216 (*Pf*MDH-R102K-INS). However, this mutation reduced activity towards oxaloacetate by
217 another 100-fold, with no increase in pyruvate activity (**Figure 4**).

218 Alternatively, it may be possible to revert a modern apicomplexan LDH to MDH-like
219 specificity by deleting its six-residue loop insertion. To test this hypothesis we removed the
220 insertion from *Pf*LDH and from the *Toxoplasma gondii* (*Tg*) LDH2 (constructs *Pf*LDH-DEL and
221 *Tg*LDH2-DEL). However, deleting the insertion from the modern LDHs abolishes pyruvate
222 activity with no significant gain of oxaloacetate activity (**Figure 4**). Both of these deletion
223 mutants retain a Lys at position 102, but a specific MDH likely requires an Arg at position 102.
224 Mutating Lys102 to Arg in *Pf*LDH-DEL results in a two order-of-magnitude gain in oxaloacetate
225 activity (**Figure 4**). However, this mutant fails to recapitulate the level of oxaloacetate activity
226 seen in modern apicomplexan MDHs. In the modern enzymes, substrate specificity cannot be
227 switched with mutations involving the loop insert and position 102, indicating that additional
228 residues govern substrate preference.

229 ***The ancestral MDH and LDH enzymes are specific and highly active***

230 The apicomplexan LDH and MDH phylogeny strongly suggests that after (or coincident
231 with) the crucial gene duplication event, the nascent LDH branch gained pyruvate activity due to
232 the six-residue insertion in the specificity loop. This presents a conundrum, as our mutation
233 trials in the modern enzymes failed to recapitulate the historical swap in specificity. However,
234 the modern apicomplexan LDH and MDH enzymes differ by over 200 residues in addition to
235 the loop insert and Arg102Lys, differences that have accumulated in the descendants of the
236 ancestral MDH and LDH. Any of these differences may detrimentally affect the ability to switch
237 substrate specificity with the insertion in the modern enzymes. We therefore reasoned that the
238 ancestral background may be necessary for swapping specificity with the loop insertion. To test
239 this, we reconstructed and characterized four key ancestral enzymes: (1) AncMDH1, the
240 ancestral protein that was transferred from α -proteobacteria to the archaic Apicomplexa, (2)
241 AncMDH2, the last common ancestor of all apicomplexan MDHs and LDHs, found at the critical

242 duplication event, (3) AncMDH3, the last common ancestor of all modern apicomplexan MDHs,
243 and (4) AncLDH, the last common ancestor of modern apicomplexan LDHs (**Figure 3B**).

244 All four ancestral proteins are highly active in steady state kinetic assays, with substrate
245 preferences and catalytic efficiencies that are similar to their modern apicomplexan descendants
246 (**Figure 5**), despite sharing only 49-71% sequence identity with the modern apicomplexan
247 proteins (**Figure 5-figure supplement 1**). AncMDH1, AncMDH2, and AncMDH3 are highly
248 specific MDHs with negligible pyruvate activity, having even greater activity towards
249 oxaloacetate than modern *Plasmodium* and *Cryptosporidium* MDHs (**Figure 5**). AncLDH is a
250 highly active and specific LDH, with very low activity towards oxaloacetate (**Figure 5**).

251 ***The loop insert successfully swaps specificity in both ancestral LDH and MDH***

252 AncLDH differs from AncMDH2 by 66 residues, including the six-residue insertion and
253 Arg102Lys. We investigated the evolutionary trajectory from AncMDH2 to AncLDH by
254 characterizing three different mutations in the AncMDH2 background: (1) the addition of
255 AncLDH's six-residue insertion to the AncMDH2 specificity loop, (2) Arg102Lys, which assesses
256 the effect of changing the canonical specificity residue, and (3) the remaining 59 residues that
257 separate AncLDH from AncMDH2, simultaneously changed to their AncLDH identities.

258 Incorporating the loop insertion into AncMDH2 confers significant pyruvate activity
259 with minimal effect on oxaloacetate activity, resulting in a highly active, bifunctional enzyme
260 (AncMDH2-INS, **Figure 6**). In contrast, the Arg102Lys mutation in the AncMDH2 background
261 (AncMDH2-R102K, **Figure 6**) reduces oxaloacetate activity by more than a 100-fold, with no
262 increase in pyruvate activity. The 59 mutations in the AncMDH2 background have a minimal
263 effect on the activity towards both substrates (AncMDH2-59Mut, **Figure 6**). Note that the
264 AncMDH2-59Mut construct is equivalent to a modified AncLDH construct with the Lys102Arg
265 mutation and the insertion deleted from the loop. Therefore, only two changes — Lys102Arg
266 and the loop deletion — are sufficient to convert the AncLDH construct to a highly active and
267 specific MDH.

268 Combinations of these mutations confirm that the insertion is primarily responsible for
269 the evolution of pyruvate activity. Adding the 59 mutations to AncMDH2-INS (resulting in a
270 construct that differs from AncLDH by only one residue) has little additional effect (AncMDH2-
271 INS-59Mut, **Figure 6**). Surprisingly, the combination of Arg102Lys and the 59 mutations, a
272 construct that differs from AncLDH by just the six-residue insertion, yields a crippled MDH
273 enzyme with 1,000-fold less oxaloacetate activity than AncMDH2 (AncMDH2-R102K-59Mut,
274 **Figure 6**). However, the combination of Arg102Lys and the loop insertion in the AncMDH2
275 background is sufficient to confer pyruvate activity and specificity comparable to AncLDH
276 (AncMDH2-INS-R102K, **Figure 6**).

277 ***Ancestral kinetics are robust to reconstruction uncertainty***

278 Ancestral sequence reconstruction is a difficult statistical problem that strongly relies on
279 evolutionary assumptions, which may be unrealistic, and on available sequence data, which is
280 inherently incomplete. The likelihood and Bayesian ancestral reconstruction methodology that
281 we use produces the most probable ancestral sequence given certain evolutionary model
282 assumptions, along with a posterior probability for alternative amino acids at each position
283 (**Figure 6-figure supplement 1-6**). Ambiguous residues are generally associated with
284 positions of low conservation and presumably less functional importance. The reconstructed
285 AncMDH2 and AncLDH sequences have 31 and 48 ambiguous positions, respectively, all of
286 which are located outside of the “first active site shell” (defined as within 6 Å of the substrate).
287 In order to verify that these sequence ambiguities do not affect our kinetic results, alternative
288 ancestral sequences were reconstructed and assayed. We tested the robustness of our ancestral
289 proteins by constructing alternative ancestors based on perturbed sequence data, evolutionary
290 assumptions, and phylogenetic methodology. Both phylogenies give very similar relationships,
291 and Figure 2B summarizes both equally well. The alternative AncMDH2 (AncMDH2*) differs
292 from AncMDH2 by 27 residues; the alternative AncLDH (AncLDH*) differs from AncLDH by 19
293 residues.

294 The alternative ancestral reconstructions behave very similar to the prior reconstructions.
295 AncMDH2* is a strict MDH, and AncLDH* is a strict LDH (**Figure 6/****Figure 6-figure**
296 **supplement 7**). Addition of the six-residue insertion from AncLDH* to AncMDH2* confers
297 pyruvate specificity without adversely affecting oxaloacetate activity (AncMDH2*-INS, **Figure**
298 **6-figure supplement 7**). In the AncMDH2* background, mutating Arg102 to Lys together
299 with the 58 mutations from AncLDH* yields a poor enzyme with little pyruvate activity
300 (AncMDH2*-R102K-58Mut). The kinetic behavior of these AncMDH2* constructs closely
301 matches those seen with the corresponding AncMDH2 constructs (AncMDH2, AncMDH2-INS,
302 and AncMDH2-R102K-59Mut, **Figure 6**).

303 ***Crystal structures of ancestral MDH, LDH, and an evolutionary intermediate***

304 In order to understand the structural changes during evolution that shifted the
305 enzymatic substrate specificity of the apicomplexan dehydrogenases, we determined the high-
306 resolution crystal structures of three ancestral proteins bracketing the key duplication event:
307 AncMDH2 (1.9 Å), AncLDH* (1.3 Å), and AncMDH2-INS (1.8 Å). All three ancestral proteins
308 adopt the same overall fold and conformation as the modern, descendant enzymes. In
309 particular, the ancestral active sites and specificity loops are highly similar to their modern
310 counterparts.

311 **Ancestral malate dehydrogenase: AncMDH2**

312 The AncMDH2 structure superposes closely with the modern CpMDH structure (47)
313 (0.56 Å RMSD), although differing at ~119 residue positions (62% sequence identity, **Figure**
314 **7A**). In the modern and ancestral MDHs, all residues within the first shell of the active sites
315 (within 6 Å of the substrate) are identical, and the active site conformations are correspondingly
316 highly similar (**Figure 7B**). The first shell active site residues comprise Arg102, Arg109, Leu112,
317 Asn140, Leu167, Asp168, Arg171, His195, Met199, Gly236, Gly237, Ile239, Val240, Ser245,
318 Ala246, and Pro250.

319 Compared to the modern MDH, only slight differences are seen in the substrate loop
320 backbone and the positioning of the Arg102 and Arg109 sidechains, which are the only residues
321 from the specificity loop that directly interact with the substrate. However, these modest
322 conformational differences are largely within coordinate error, as the loop residues have some of
323 the highest B-factors in the structures. Furthermore, AncMDH2 was crystallized with lactate
324 and NADH while *Cp*MDH was crystallized with citrate and ADPR (an NADH analog lacking the
325 nicotinamide ring). Citrate is roughly three times larger than lactate and has likely affected the
326 position of substrate loop in the *Cp*MDH structure.

327 **Ancestral lactate dehydrogenase: AncLDH***

328 The ancestral AncLDH* and modern apicomplexan LDH structures are likewise highly
329 similar (6, 26, 28) (RMSD ~0.8 Å, **Figure 7C**), while sharing only 63-71% sequence identity.
330 The first shell active site residues are identical in the AncLDH* and modern *Toxoplasma* LDHs,
331 comprising the same residues as the apicomplexan MDH active site with the sole exception of
332 position 102, which is replaced by Trp107f in the LDHs. The modern *Plasmodium* LDHs have
333 two different residues in the active site first shell: Pro246 and Ala236, rather than Ala246 and
334 Gly236 as found in *Tg*LDH1, *Tg*LDH2, and AncLDH*. The conformations of the ancestral and
335 modern active sites are nearly indistinguishable, with only small differences in the specificity
336 loop conformation (**Figure 7D**).

337 In both the ancestral and modern LDH structures, Trp107f and Arg109 are the only
338 residues from the specificity loop that interact with the substrate. As in the modern LDH
339 structures, ancestral Lys102 does not interact with the substrate but points away from the active
340 site into solution. In contrast, Trp107f is buried within the active site, with the edge of the
341 indole ring interacting with the pyruvate C3 methyl, which is the very chemical moiety that
342 distinguishes pyruvate from oxaloacetate.

343 The largest differences between the modern and ancestral proteins are confined to two
344 regions: (1) a small shift of the entire C-terminal helix, and (2) a loop opposite the active site

345 specificity loop (residues 242-244, hereafter called the “opposing loop”). The modern
346 *Plasmodium* LDHs have a two-residue deletion within this opposing loop (highlighted in cyan in
347 **Figure 7C**), while the opposing loop is shared with AncLDH* and the *Toxoplasma* LDHs. The
348 ancestral LDHs also share very modest oxaloacetate activity with the modern *Toxoplasma* LDHs,
349 while the *Plasmodium* LDHs lack oxaloacetate activity (**Figure 5**). This correlation indicates
350 the opposing loop deletion (and perhaps Ala236 and Pro246) may be responsible for the
351 unusually strict substrate specificity of the modern *Plasmodium* LDHs.

352 **Ancestral malate dehydrogenase with loop insertion: AncMDH2-INS**

353 We also crystallized AncMDH2-INS, a bifunctional AncMDH2 construct with the six-
354 residue specificity loop insertion. This AncMDH2-INS construct represents a possible
355 intermediate along the evolutionary trajectory between the MDH duplication event and the
356 ancestral apicomplexan LDH. AncMDH2-INS was successfully co-crystallized with both
357 oxamate/NADH and lactate/NADH. In the closed form, the specificity loop adopts an LDH-like
358 conformation with Trp107f occupying the specificity position and Arg102 oriented into solution,
359 similar to how Lys102 is positioned in the modern and ancestral LDH structures (**Figure 7F**).
360 The lactate and oxamate structures are highly similar (RMSD 0.19 Å), and the active site
361 architectures are nearly indistinguishable.

362 The three ancestral proteins, AncMDH2, AncLDH*, and AncMDH2-INS, are all highly
363 similar (RMSD 1.20 Å) with the main structural differences found in the conformation of the
364 specificity loop (**Figure 7E**, RMSD 0.86 Å excluding residues in the specificity loop).
365 Otherwise the first shell active site residues are identical between AncMDH2-INS and AncLDH*,
366 and the conformations of the active sites are correspondingly similar (**Figure 7F**).

367 **Convergent pathways available to the ancestral MDH**

368 Given the known importance of position 102, the “specificity residue”, in substrate
369 recognition, we wondered whether different residues at position 102 could confer pyruvate

370 activity. Position 102 in fact differs in the four convergent LDH families: Gln in canonical LDHs
371 (19), Lys in the apicomplexan LDHs, Gly in *Cryptosporidium* LDHs (39), and Leu in
372 trichomonad LDHs (40). Could the ancestral apicomplexan MDH have evolved pyruvate
373 specificity by any of these alternative routes? To answer this question, we evaluated the
374 potential of these different amino acids at the 102 position to confer pyruvate specificity in the
375 AncMDH2 background. Each mutation increases pyruvate activity, but none result in a highly
376 specific LDH. The canonical mutation (Arg102Gln) results in the largest gain in pyruvate
377 activity (2,800-fold) and the smallest loss of oxaloacetate activity (2,500-fold) (**Figure 8**).
378 Additionally, we tested whether the full six amino acid insertion was required to confer pyruvate
379 specificity in AncMDH2 or if simply mutating Arg102 to Trp was sufficient. The Arg102Trp
380 mutation all but abolishes activity towards both substrates, indicating that the loop insertion
381 was necessary to switch the specificity residue (**Figure 8**).

382 **Discussion**

383 ***An alternate mechanism of specificity in the convergent apicomplexan LDH***

384 Substrate recognition in the canonical MDHs and LDHs is thought to be determined by a
385 “specificity residue” in the active site loop at position 102. All known MDHs have Arg at
386 position 102, while canonical LDHs have Gln (21). In the classic explanation of the molecular
387 mechanism of substrate specificity, residue 102 discriminates between pyruvate and
388 oxaloacetate primarily via charge conservation (19). In MDHs, the positively charged Arg
389 interacts with and balances the negatively charged carboxylate of oxaloacetate. If pyruvate were
390 to bind in the active site, loop closure would result in a buried and unbalanced positive charge,
391 which is unfavorable. In canonical LDHs, the neutral Gln interacts with the neutral pyruvate
392 methyl group. Oxaloacetate binding would similarly result in the unfavorable burial of an
393 unbalanced negative charge.

394 In the apicomplexan LDHs, evolution has converged on pyruvate specificity using an
395 alternative molecular mechanism. Residue 102 is not a Gln but a positively charged Lys, similar
396 to Arg102 of MDHs, leading many researchers to wonder why apicomplexan LDHs lack activity
397 towards oxaloacetate (8, 24-28). However, during the evolution of the apicomplexan LDH from
398 the ancestral MDH, the six-residue insertion in the active site loop shifted both the position and
399 identity of the “specificity residue” from Arg102 to Trp107f. Due to the insertion, residue 102
400 no longer interacts with the substrate and is extruded from the active site. In contrast, the
401 hydrophobic Trp107f packs against the C3 methyl of the pyruvate substrate. Similar to the
402 canonical LDH, oxaloacetate binding would result in an unbalanced and buried negative charge.
403 As a large bulky residue, Trp107f can also occlude binding of the larger oxaloacetate, in which a
404 methylene carboxylate replaces the pyruvate methyl.

405 However, as discussed in detail below, this simplistic explanation is complicated by
406 long-range epistatic interactions. When the six-residue insertion is introduced into the modern

407 apicomplexan MDH, specificity is not switched; both specificity and activity are lost. Similarly,
408 removal of the insertion from the modern apicomplexan LDHs fails to swap specificity and kills
409 the enzymes. Therefore, while Trp107f is necessary for substrate specificity in the apicomplexan
410 enzymes (as indicated by the alanine scan mutations), it is insufficient to confer specificity.

411 The bifunctionality of AncMDH2-INS and AncMDH2-INS-59Mut also presents a
412 conundrum. Why do these constructs have high activity towards both pyruvate and oxaloacetate
413 substrates? The crystal structure of AncMDH2-INS offers few clues, since the loop insertion,
414 including Trp107f, adopts the same conformation as seen in AncLDH and the modern
415 apicomplexan enzymes. Both the AncMDH2-INS and AncMDH2-INS-59Mut constructs have
416 an Arg at position 102, like the MDHs. In fact, the bifunctional AncMDH2-INS-59Mut enzyme
417 differs from the strict AncLDH by only a R102K mutation, suggesting that Arg102 is responsible
418 for the oxaloacetate activity of AncMDH2-INS and AncMDH2-INS-59Mut. We speculate that
419 perhaps the enzymes change conformation depending upon the substrate. When using pyruvate,
420 these bifunctional enzymes may adopt an LDH-like conformation in which Trp107f interacts
421 with the substrate (as seen in the crystal structure). On the other hand, when presented with
422 oxaloacetate, perhaps Trp107f flips out of the active site, and Arg102 flips in to interact with
423 substrate in a manner similar to the canonical MDHs. We are currently testing this hypothesis.

424 ***Apicomplexan LDH evolved by classical neofunctionalization***

425 Our data show that apicomplexan LDHs evolved from a horizontally transferred
426 proteobacterial MDH by a classic neofunctionalization mechanism of gene duplication. Because
427 debasement to a pseudogene is much more likely to occur prior to the evolution of a novel
428 function, neofunctionalization has fallen out of favor as a mechanism for the evolution of novel
429 functions. A variety of alternative specialization models have been proposed that feature a
430 reduced risk of non-functionalization. Though differing in details, all specialization models
431 feature a promiscuous common ancestor of the duplicated proteins.

432 The reconstructed AncMDH2, which represents the last common ancestor of the
433 apicomplexan MDH and LDHs, is a highly active and specific MDH, preferring oxaloacetate over
434 pyruvate by seven orders of magnitude (**Figure 6**). The activity of AncMDH2 towards pyruvate
435 is barely detectable, requiring a high enzyme concentration to quantify. AncMDH2's k_{cat} for
436 pyruvate is 0.07 s^{-1} , with a K_m of 20 mM, while the physiological concentration of pyruvate is
437 estimated to be about three orders-of-magnitude lower (e.g., $\sim 50 \mu\text{M}$ in human erythrocytes
438 (48), the *Plasmodium* host during its blood stage). Based on these kinetic parameters, each
439 AncMDH2 reduces one pyruvate molecule per hour. While the enzyme can be forced to reduce
440 pyruvate *in vitro*, this negligible activity is unlikely to have been subjected to selection *in vivo*.
441 Therefore, the various specialization hypotheses, which require a promiscuous ancestor, are
442 poor models for apicomplexan LDH evolution. Activity towards pyruvate increased by over
443 seven orders of magnitude on the evolutionary lineage between AncMDH2 and AncLDH,
444 indicating neofunctionalization.

445 ***A highly active, promiscuous intermediate***

446 One of the most favored specialization models is “escape from adaptive conflict” (EAC)
447 (49). EAC holds that functional specialization is driven by an inability to simultaneously
448 optimize multiple functions on a single protein scaffold. Gene duplication relieves this
449 constraint and allows for the independent optimization of conflicting functions. Although the
450 apicomplexan AncMDH2 is highly specific, promiscuous intermediates did play a role in the
451 functional transition between AncMDH2 and AncLDH. AncMDH2-INS and AncMDH2-INS-
452 59Mut have high levels of MDH and LDH activity in a single protein scaffold (**Figure 6**). Both
453 the presence of bifunctional intermediates and the high specificity of AncMDH2 conflict with
454 fundamental predictions of the EAC specialization model.

455 **Convergent evolution of apicomplexan LDH involved long-range epistasis**

456 The evolution of apicomplexan LDHs involved strong epistasis that has profoundly
457 influenced the convergent evolution of pyruvate activity. Epistasis refers to interactions
458 between residues that potentiate the effects of a mutation depending on the presence or absence
459 of other residues (50). Epistasis can constrain the order of mutations and the pathways
460 accessible to evolution, and hence it is of great importance in understanding the evolution of
461 novel functions. In the apicomplexan dehydrogenases, the evolutionary mutations that switched
462 specificity from oxaloacetate to pyruvate (the six-residue insertion and Arg102Lys) are
463 insufficient to confer pyruvate activity in modern apicomplexan MDHs (*PfMDH-R102K*,
464 *PfMDH-INS*, *CpMDH-INS*, *PfMDH-R102K-INS*, **Figure 4**). However, these mutations are
465 sufficient to confer pyruvate function and specificity in the *AncMDH2* background (*AncMDH2-*
466 *INS*, *AncMDH2-INS-R102K*, **Figure 6**). Similarly, removal of the insert from the modern
467 LDHs (*PfLDH-DEL* and *TgLDH2-DEL*, **Figure 4**) kills the enzymes, while removal of the insert
468 from the ancestral LDH (*AncMDH2-R102K-59Mut*, **Figure 6**) results in a weak MDH. The
469 different effects of these mutations, depending upon the sequence of the rest of the protein,
470 provide direct evidence of epistatic interactions.

471 Why do these historical mutations “work” in the ancestral enzymes, but not in the
472 modern ones? Epistatic interactions are often mediated by direct physical contact. For example,
473 the active site of the ancestral MDH could have certain residues that the modern MDH lacks,
474 residues that interact with the insertion and allow it to preferentially bind pyruvate. However,
475 the active sites of the ancestral and modern MDHs are identical in sequence and virtually
476 indistinguishable in structure (**Figure 7B**), as are the active sites of the ancestral and modern
477 LDHs (**Figure 7D**) and the *AncMDH2-INS* intermediate (**Figure 7F**). In fact, the active sites
478 of the MDHs and the LDHs are also identical in sequence except for the 102 position, and they
479 are otherwise highly structurally similar. Therefore, residues remote from the active sites
480 necessarily affect the substrate specificity of the enzymes.

481 In principle, these long-range epistatic residue interactions could differentially modify
482 the structure of the active site. Certain residues found in the ancestral MDH, but not in the
483 modern MDH, could position the active site residues so that they allowed the insertion to confer
484 pyruvate activity. In this scenario the active site residues of the ancestral and modern MDHs
485 would be identical, but their conformations would differ due to interactions with residues in
486 other parts of the protein. However, the crystal structures reveal ancestral, intermediate, and
487 modern active sites that are nearly indistinguishable, suggesting that epistasis has modified the
488 protein dynamics or shifted the energy landscape, effects that are largely invisible to static
489 crystal structures.

490 ***Epistasis prevents mechanistic convergence in the LDH/MDH superfamily***

491 Interestingly, *Bacillus subtilis* (*Bs*) LDH reverts to an MDH with only a single mutation,
492 Gln102Arg, indicating a lack of complicating epistatic effects (19). The kinetics of wild-type
493 *Bs*LDH with pyruvate are comparable to those for the Gln102Arg mutant with oxaloacetate (e.g.,
494 *Bs*LDH has a $k_{\text{cat}}/K_{\text{M}}$ for pyruvate of $4.2 \times 10^6 \text{ M}^{-1} \text{ s}^{-1}$, and the *Bs*LDH-Q102R mutant has the
495 same $k_{\text{cat}}/K_{\text{M}}$ for oxaloacetate). However, *Bs*LDH likely is an exception in the LDH/MDH
496 superfamily, since the reverse mutation (Arg102Gln) fails to switch specificity in MDHs from
497 two other species (22, 23). In *Haloarcula marismortui* (*Hm*) MDH, the Arg102Gln mutation
498 switches specificity, but the mutant's $k_{\text{cat}}/K_{\text{M}}$ for pyruvate is 200-fold less than the wild-type's
499 $k_{\text{cat}}/K_{\text{M}}$ for oxaloacetate. The Arg102Gln mutation in *Escherichia coli* (*Ec*) MDH is even less
500 effective, as it converts a highly active MDH to an enzyme with low activity on both substrates
501 (10,000-fold lower $k_{\text{cat}}/K_{\text{M}}$). Hence, the strong epistasis observed in apicomplexan LDH and
502 MDHs is likely a general phenomenon within the superfamily.

503 LDH evolved convergently from MDH four separate times in the superfamily, but did the
504 activity evolve by the same mechanism each time? Each event has resulted in a different change
505 at the specificity residue (position 102) within the catalytic loop. However, the epistatic effects
506 seen in the apicomplexan, *H. marismortui*, and *E. coli* dehydrogenases indicate that in general

507 position 102 is not solely responsible for the transition from MDH to LDH. In order for the
508 historical LDH mutations to confer pyruvate specificity, additional residues must be present to
509 provide a permissive background (**Figure 8**). Due to the presence of different sets of
510 permissive mutations, LDH activity has evolved from an MDH under epistatic constraints by a
511 different mechanism four separate times.

512 ***Large effect, gain-of-function mutation***

513 The evolution of AncLDH from AncMDH2 involves a shift in substrate specificity by
514 twelve orders-of-magnitude. Through the characterization of possible evolutionary
515 intermediates, we have found that just two mutations are responsible for the great majority of
516 this switch: the six-residue insertion and the Arg102Lys point mutation. Mutagenesis within the
517 insertion indicates that only a single position, Trp107f, contributes strongly to pyruvate activity
518 and specificity. Both the insertion and Arg102Lys have a large effect on preference for pyruvate
519 vs oxaloacetate, although by differentially affecting activity towards each substrate.
520 Incorporating the six-residue insertion into AncMDH2's substrate loop results in a 12,000-fold
521 gain in pyruvate activity with little effect on oxaloacetate activity (**Figure 6**). Conversely,
522 mutating Arg102 to Lys reduces oxaloacetate activity by more than 2,500-fold, with minimal
523 effect on pyruvate activity (**Figure 6**).

524 The apicomplexan LDH six-residue insertion is an exceptionally large gain-of-function
525 mutation: it enhances pyruvate activity by more than four orders of magnitude while barely
526 affecting oxaloacetate activity. In contrast, other well-studied mutations of large effect are often
527 predominantly deleterious towards one function while modestly enhancing another. The
528 textbook example of a gain-of-function mutation is Gln102Arg in *BsLDH*, which causes a 10⁷-
529 fold change in the enzyme's specificity (19). The Gln102Arg mutation reduces pyruvate activity
530 by more than 8,000-fold, while enhancing activity towards oxaloacetate by only 1,000-fold.
531 Another example is given by *E. coli* isocitrate dehydrogenase (IDH), where seven mutations are
532 necessary to switch the cofactor specificity from a 7,000-fold preference for NADP to a 200-fold

533 preference for NAD (51). Within this set of mutations, two reduce specificity for NADP by
534 6,000-fold, whereas the rest enhance NAD usage 200-fold. Thus, while mutations can have
535 both deleterious and beneficial effects on different functions, the deleterious effects typically
536 appear greater than enhancement.

537 In previous ancestral sequence reconstruction studies, mutations of large effect are in
538 fact usually loss-of-function rather than gain-of-function (e.g., RNaseA (52), chymase (53), and
539 glucocorticoid receptors (54-57)). In these studies, the modern proteins are generally specific
540 for one substrate, whereas the ancestral proteins are promiscuous. Furthermore, the activity of
541 the ancestral protein is comparable to the modern descendants. Therefore, these proteins
542 specialized by accumulating deleterious mutations, with the modern, specialized activity being
543 the “last function standing”. For example, the ancestral glucocorticoid receptor binds three
544 steroid hormones tightly ($EC_{50} < 10$ nM for aldosterone, deoxycorticosterone, and cortisol),
545 while the modern receptors bind only cortisol ($EC_{50} \sim 100$ nM) (57). Seven historical mutations
546 produced the modern cortisol preference by completely eliminating aldosterone and
547 deoxycorticosterone sensitivity yet reducing cortisol sensitivity only 50-fold. In other ancestral
548 reconstruction studies, function-enhancing mutations have relatively minor effects, all less than
549 a 50-fold gain in k_{cat}/K_M (37, 38, 58).

550 During the evolution of the malate and lactate dehydrogenase superfamily,
551 pyruvate activity has converged multiple times despite strong constraints due to epistasis.
552 While epistasis may constrain evolutionary options locally, there are nevertheless multiple ways
553 to “skin the cat” in more distant regions of protein sequence space. The apicomplexan enzymes
554 provide a clear example of neofunctionalization in protein evolution and thereby validate the
555 plausibility of this particular mechanism of gene duplication. Specialization mechanisms may
556 be more common, but the evolution of novel function does not require a promiscuous genesis.

557 **Acknowledgements**

558 This work was supported by the National Institutes of Health, NIH grants R01GM096053 and
559 R01GM094468. The crystallographic data collection was conducted at the SIBYLS beamline at
560 the Advanced Light Source (ALS), a national user facility operated by Lawrence Berkeley
561 National Laboratory on behalf of the Department of Energy, Office of Basic Energy Sciences,
562 through the Integrated Diffraction Analysis Technologies (IDAT) program, supported by DOE
563 Office of Biological and Environmental Research. Additional support comes from the National
564 Institute of Health project MINOS (R01GM105404). We would also like to thank Chris Miller,
565 Phillip Steindel, and Catherine Theobald for critical commentary on the manuscript.

566

567 **Materials and Methods**

568 ***Modern sequences***

569 Protein sequences used in the phylogenetic analyses were identified through searches of
570 the non-redundant database (59) with the BLASTP algorithm (60) using selected query
571 sequences. All sequences from these searches that returned BLASTP E-values $<10^{-7}$ were
572 downloaded from NCBI (www.ncbi.nlm.nih.gov). Multiple complete apicomplexan genomes
573 (41-43) were also searched for LDH and MDH homologs in order to fill out the apicomplexan
574 portion of the tree (using a more lenient significance cutoff of E-values $<10^{-4}$). Redundant
575 sequences, synthetic constructs, and sequences from PDB files were removed. To reduce
576 phylogenetic complexity, sequences were curated based on character length and pairwise
577 sequence identity within each dataset (as described below).

578 The dataset used for the construction of the non-redundant phylogeny (**Figure 3A**) was
579 generated using four query sequences, UniProt IDs (61): MDHC_HUMAN, LDH_THEP1,
580 MDHP_YEAST, and LDH6A_HUMAN. Multiple sequences were necessary to generate full
581 coverage, due to the low sequence identity across the superfamily, which can be less than 20%
582 between members. Sequences were removed if their character length was less than 280 or
583 greater than 340. Limits were chosen to remove truncated/partial sequences and those
584 featuring large insertions or terminal extensions. Sequences greater than 97% identical,
585 determined by pairwise alignment within the dataset, were also removed. This level of identity
586 provides a high level of detail within the tree while accelerating computational time by removing
587 redundant taxa. The final dataset contains 1844 taxa.

588 Residue numbering in the text is based on the dogfish LDH convention (20) for
589 consistency with previous work.

590 ***Primary Phylogeny Construction***

591 A multiple sequence alignment of this dataset was generated using the program
592 MUSCLE (62). A maximum likelihood (ML) phylogenetic tree was inferred with PhyML 3.0
593 (63) using the LG substitution matrix (64) and estimating the gamma parameter (12 categories)
594 and empirical amino acid frequencies. The starting tree was generated by Neighbor-Joining
595 (BIONJ) and searched by Nearest Neighbor Interchange (NNI); tree topology, branch lengths,
596 and rate parameters were optimized. Branch supports were estimated with the approximate
597 likelihood ratio test (aLRT), as implemented in PhyML, represented as either the raw aLRT
598 statistic (roughly > 8 is considered highly significant) or the confidence level that the clade is
599 correct (65).

600 ***Phylogeny Rooting***

601 The outgroup for rooting the L/MDH phylogeny was identified through a profile analysis
602 of the Rossmann fold (66), based on a method used for OB folds and SH3 domains (67). All
603 structurally characterized Rossmann folds with 40% or less sequence identity were identified
604 from ASTRAL SCOP 1.73 protein domain sequence database (68). Each of the 193 domains
605 identified was searched against the SwissProt database (69) using BLASTP. A multiple
606 sequence alignment for each query and SwissProt sequences with BLASTP E-values $< 10^{-10}$ was
607 created using MUSCLE. Each alignment was cropped to the limits of the original query.
608 COMPASS (70) was then used to generate an all-against-all scoring matrix for the 193 multiple
609 sequence alignments. The E-values generated by COMPASS were converted to evolutionary
610 distances as described in Theobald & Wuttke (67). A weighted least-squares phylogenetic
611 analysis of the distance matrix was performed using PAUP (71). First order taxon jackknifing
612 (72, 73) was used to determine the robustness of tree topology, with a consensus tree calculated
613 from all analyses.

614 Rossmann fold domains from α - and β -glucosidases and aspartate dehydrogenases
615 (AspDH) were identified from the profile-profile analysis as grouping with the Rossmann fold
616 domain from L/MDHs. An L/MDH dataset was constructed for use with the outgroup to create
617 a rooted phylogeny. This dataset was generated by querying four sequences, UniProt IDs:
618 MDHC_HUMAN, LDH_THEP1, MDHP_YEAST, and LDH6A_HUMAN, against the SwissProt
619 database using BLASTP. All sequences from these searches that returned BLASTP E-values $<10^{-7}$
620 were downloaded from NCBI (www.ncbi.nlm.nih.gov). Redundant sequences, synthetic
621 constructs, and sequences from PDB files were removed. Also, four taxa identified as ubiquitin-
622 conjugating enzymes were removed due to sequence length. This SwissProt L/MDH dataset
623 contained 595 taxa.

624 An outgroup dataset was constructed by querying three sequences, UniProt IDs:
625 LICH_BACSU, AGAL_THEMA, and ASPD_THEMA, against the SwissProt database using
626 BLASTP. All sequences from these searches that returned BLASTP E-values $<10^{-7}$ were
627 downloaded from NCBI (www.ncbi.nlm.nih.gov). Redundant sequences, synthetic constructs,
628 and sequences from PDB files were removed. The outgroup dataset contained 62 taxa. The
629 SwissProt LDH, MDH, AspDH, and glucosidase datasets were combined and a multiple
630 sequence alignment was generated using the program MUSCLE. The C-terminal domain of the
631 glucosidases and AspDHs were removed from the MUSCLE alignment. A ML phylogenetic tree
632 was inferred from the alignment with PhyML using the LG substitution matrix (74) with the
633 gamma parameter estimated over 10 categories, no invariant sites, and estimating empirical
634 amino acid frequencies. The initial tree was obtained by BIONJ and searched by NNI; tree
635 topology, branch lengths, and rate parameters were optimized. Robustness of root positioning
636 was evaluated with two truncated alignments, one with the LDH and “LDH-like” MDH
637 sequences removed and the other with the cytosolic and mitochondrial MDH sequences
638 removed. Truncated alignments were input to PhyML for phylogenetic analysis using the
639 parameters described above.

640 ***Alternative phylogeny construction***

641 The dataset for the alternative phylogeny (used in reconstructing alternative ancestors)
642 is smaller and focused on apicomplexan taxa. It was generated by BLASTP searches with four
643 query sequences, UniProt IDs: MDHC_PIG, Q76NM3_PLAF7, C6KT25_PLAF, and
644 MDH_WOLPM for full coverage of the superfamily. Sequences were removed if their length
645 was less than 290 or greater than 340. The dataset was culled to 60% identity, but the
646 apicomplexan clade was filled back to 97% identity to gain resolution within the clade of interest.
647 The final dataset contained 277 taxa. A multiple sequence alignment of this dataset was
648 produced using the program MUSCLE. The ML tree was inferred with PhyML 3.0 using the LG
649 substitution matrix and estimating the gamma parameter (12 categories) and empirical amino
650 acid frequencies. The starting tree was generated by Neighbor-Joining (BIONJ) and searched
651 by Nearest Neighbor Interchange (NNI); tree topology, branch lengths, and rate parameters
652 were optimized.

653 ***Ancestral Sequence Reconstruction***

654 Sequences at internal nodes in phylogenies were inferred using the *codeml* program
655 from the PAML software package (75). Posterior amino acid probabilities at each site were
656 calculated using the LG substitution matrix, given the ML tree generated by PhyML. The initial
657 ancestral reconstruction assumed the background amino acid frequencies implicit in the LG
658 matrix, while the alternative reconstruction estimated background frequencies from the
659 sequence alignment of the alternative dataset. N-/C-termini of ancestral sequences were
660 modified manually to match those of the closest modern sequence (determined by branch
661 length).

662 ***Plasmid Construction and Mutation***

663 *Escherichia coli* codon-optimized coding sequences were constructed for the
664 *Plasmodium falciparum* MDH (gi#: 86171227), *Cryptosporidium parvum* MDH (gi#:

665 32765705), *Toxoplasma gondii* LDH1 (gi#: 237837615), *Toxoplasma gondii* LDH2 (gi#: 2497625), *Rickettsia bellii* MDH (gi#: 91205459), and ancestrally inferred protein sequences. 666 These coding sequences were synthesized and subcloned into pET-24a, bypassing the N- 667 terminal T7-tag but using the C-terminal 6xHis-tag. *Pf*LDH (gi#: 124513266) with six His 668 residues added to the C-terminus was synthesized and subcloned into pET-11b. All gene 669 synthesis and subcloning was performed by Genscript (Piscataway, NJ). All point mutations 670 were made using the QuikChange Lightning kit from Agilent (Santa Clara, CA) and synthesized 671 primers from IDT (Coralville, IA). 672

673 ***Protein Expression and Purification***

674 Plasmids were transformed in BL21 DE3 (pLysS) *E. coli* cells (Invitrogen, Grand Island, 675 NY) for expression. Cells were grown at 37°C with 225 rpm agitation in 2xYT media 676 supplemented with 30 mM potassium phosphate, pH 7.8 and 0.1% (w/v) glucose. Once cultures 677 reached an OD₆₀₀ between 0.5-0.8, cells were induced with 0.5 mM IPTG for 4 hours. Cells 678 were collected by centrifugation at 10,000xg for 15 minutes and stored at -80°C.

679 Cell pellets were thawed on ice, releasing lysozyme produced by the pLysS plasmid from 680 within the cells, and resuspended in 15 mL lysis buffer (50 mM NaH₂PO₄, pH 8.0, 300 mM NaCl, 681 10 mM Imidazole) with 375 units of Pierce Universal Nuclease (Thermo Scientific, Rockford, IL) 682 per 1.5 L of culture. Once homogeneously resuspended, lysate was sonicated on ice at 35% 683 amplitude (30 sec ON, 20 sec OFF, 2 min total). Insoluble cell debris was separated by 684 centrifugation at 18,000xg for 20 min.

685 Proteins were purified by nickel affinity chromatography. Clarified lysate was applied to 686 a 5 mL HisTrap FF column (GE Healthcare, Piscataway, NJ) and eluted via an imidazole 687 gradient from 10 mM to 500 mM on an AKTA Prime (GE Healthcare, Piscataway, NJ). 688 Fractions were analyzed by SDS-PAGE, pooled, and concentrated using Amicon Ultracel-10K 689 centrifugal filters (Millipore, Billerica, MA). Finally, proteins were desalted into 50 mM Tris, pH

690 7.4, 100 mM NaCl, 0.1 mM EDTA and 0.01% azide by either PD10 column (GE Healthcare,
691 Piscataway, NJ) or gel filtration over a HiPrep 16/60 Sephacryl S-200 HR column (GE
692 Healthcare, Piscataway, NJ) on an AKTA Purifier (GE Healthcare, Piscataway, NJ). Enzyme
693 concentrations were determined by absorbance at 280 nm, using extinction coefficients and
694 molecular weights calculated by ExPASy's ProtParam tool (<http://web.expasy.org/protparam/>).

695 ***Steady-state Kinetic Assays***

696 Enzymatic reduction of pyruvate and oxaloacetate was monitored at 25°C by following
697 the decrease in absorbance at 340 nm due to NADH oxidation on a Cary 100 Bio (Agilent, Santa
698 Clara, CA) in 50 mM Tris, pH 7.5, 50 mM KCl. All substrates were purchased from Sigma-
699 Aldrich (St. Louis, MO). NADH concentration was held constant at 200 μM while
700 pyruvate/oxaloacetate concentrations were titrated. Enzyme concentrations ranged from 0.28
701 nM to 2.8 μM, depending on enzyme activity for a particular substrate. All experiments used 1-
702 cm path-length quartz cuvettes with 500 μL final volume of reaction mixture.

703 Kinetic parameters were estimated by chi-squared fitting to either the Michaelis-Menton
704 equation ($v/[E]_t = k_{cat} [S]/(K_M + [S])$) or a substrate inhibition equation ($v/[E]_t = k_{cat} [S]/(K_M +$
705 $[S] + [S]^2/K_i)$) using the KaleidaGraph software. Three datasets were fit using a modified
706 substrate inhibition equation with $K_M = K_i$ for identifiability and to prevent the K_i being less than
707 K_M . These datasets were: AncMDH2-INS-59Mut oxaloacetate and AncMDH2-R102Q for both
708 oxaloacetate and pyruvate.

709 Aqueous oxaloacetate spontaneously decarboxylates to pyruvate at 25 °C and neutral pH
710 at a rate of $\sim 3 \times 10^{-5} \text{ s}^{-1}$ (approximately 10% per hour) (76). As a result, oxaloacetate
711 preparations contain appreciable pyruvate contamination (approximately 1-3% from Sigma-
712 Aldrich, depending on batch) and must be handled with care. All oxaloacetate stock solutions
713 were made fresh before each assay and kept on ice to keep decarboxylation to a minimum. For
714 enzymes with low pyruvate activity, the oxaloacetate decarboxylation has a negligible affect on

715 measured rates. However, enzymes with appreciable pyruvate activity can display an apparent,
716 artifactual oxaloacetate activity that is due to pyruvate contamination (19, 27, 77). In this work,
717 seven such proteins are *PfLDH*, *PfLDH-K102R*, *TgLDH1*, *TgLDH2*, *AncLDH*, *AncLDH**, and
718 *AncMDH2-INS-R102K*. For these proteins, oxaloacetate activity was assayed at high enzyme
719 concentration (600 nM to 1 μ M), resulting in a biphasic ΔA_{340} trace with an initial burst in which
720 pyruvate is rapidly consumed followed by a slower linear phase representing oxaloacetate
721 reduction. The post-burst (slow) phase of the ΔA_{340} trace was used to quantify the oxaloacetate
722 catalytic rate (19, 77). This procedure controls for the standing pyruvate contamination but does
723 not account for the relatively slow spontaneous decarboxylation during the assay. Hence, the
724 oxaloacetate k_{cat}/K_m values for the seven enzymes with high pyruvate activity should be
725 considered upper limits on the true oxaloacetate activity. The low or negligible oxaloacetate
726 activities of these seven enzymes were further confirmed by (1) undetectable malate/NAD⁺
727 reactions in spectroscopic steady state enzyme assays, and (2) the absence of malate product as
728 determined from 1D proton NMR (3 μ M enzyme, 5 mM oxaloacetate, 5 mM NADH in
729 NaCl/P_i/D₂O pH 7.5 over four hour reaction) (27).

730 ***Protein Crystallization***

731 Crystallization trials were conducted by hanging-drop vapor-diffusion at room
732 temperature using Crystal Screen™ and Crystal Screen 2™ from Hampton Research (Aliso Viejo,
733 CA). Drops consisting of 2 μ L reservoir solution and 2 μ L protein stock were equilibrated
734 against 1 mL of reservoir solution. Crystals of the ancestral proteins were identified from
735 condition #43 of Crystal Screen™ (30% (w/v) polyethylene glycol 1,500) and further refined by
736 adding 0.1 M sodium HEPES.

737 Crystals of the ternary complexes were grown at room temperature by hanging-drop
738 vapor-diffusion with 4 μ L drops of 1:1 precipitating buffer:protein. Ancestral malate
739 dehydrogenase (*AncMDH2*, 25 mg/mL) was co-crystallized with 2 mM oxamate/NADH in 35%

740 (w/v) PEG-1500, 0.1 M sodium HEPES, pH 7.5 and with 2 mM L-lactate/NADH in 30% (w/v)
741 PEG-1500, 0.1 M sodium HEPES, pH 7.3. AncMDH2 with insertion (AncMDH2-INS, 18
742 mg/mL) was co-crystallized with 1 mM oxamate/NADH and 1 mM L-lactate/NADH in 25%
743 (w/v) PEG-1500, 0.1 M sodium HEPES, pH 8.1. Alternative ancestral lactate dehydrogenase
744 (AncLDH*, 20 mg/mL) was co-crystallized with 2 mM oxamate/NADH and 2 mM L-
745 lactate/NADH in 20% (w/v) PEG-1500, 0.1 M sodium HEPES, pH 7.5. *Pf*LDH_W107fA
746 (20mg/mL) was co-crystallized with 1.2 mM oxamate/2 mM NADH in 22% (w/v) PEG-1000.

747 All crystals were cryoprotected with a 30% (w/v) dextrose solution (15 mg dextrose
748 dissolved in 50 μ L reservoir solution). Crystals were harvested from the drop, soaked in 15%
749 (w/v) dextrose solution for 3 minutes, transferred to the 30% solution, and flash-frozen
750 immediately in liquid N₂.

751 **Structure Determination**

752 Diffraction datasets were collected at the SIBYLS beamline (12.3.1, Lawrence Berkeley
753 National Laboratory, Berkeley, CA). All datasets were indexed, integrated, and scaled with
754 XDS/XSCALE (78). Structures were solved by molecular replacement using AutoMR in
755 PHENIX (79). Homology models for the AncMDH2 and AncLDH* datasets were generated by
756 the Phyre2 server (80). The AncMDH2 homology model was based on the structure for
757 *Cryptosporidium parvum* MDH (PDB entry: 2hjr, 62% sequence identity, (47)), while the
758 model for AncLDH* was based on the *Toxoplasma gondii* LDH1 structure (PDB entry: 1pzf, 65%
759 sequence identity, (26)). AncMDH2-INS datasets were solved using the AncMDH2 structure as
760 the model. *Pf*LDH_W107fA dataset was solved using the *P. falciparum* LDH structure (pdb id:
761 1t2d) structure as a model. All models were improved by rounds of manual building in Coot (81)
762 and refinement by phenix.refine in PHENIX. Model quality of all structures was validated with
763 MolProbity (82, 83) in PHENIX. All structural alignments were generated using THESEUS (84).
764 Structure images were rendered with PyMOL.

765 **Figure legends**

766 **Figure 1. Schematic of M/LDH superfamily active site and catalytic mechanism.**

767 MDH reduces oxaloacetate to malate, in which the R-group is a methylene carboxylate. LDH
768 reduces pyruvate to lactate, in which the R-group is a methyl. Key conserved active site residues
769 are shown in black; substrate is shown in blue. The oxidized 2-ketoacid form of the substrate is
770 at left; the reduced 2-hydroxy acid form is shown at right. The R-group of the substrate
771 interacts with Arg102 in MDHs and Gln102 in LDHs. Both Arg109 and position 102 are found
772 in the “specificity loop” that closes over the active site.

773

774 **Figure 2. Apicomplexan M/LDH active sites.** Structures of *CpMDH* (blue, PDB ID: 2hjr)

775 and *PfLDH* (vermilion, PDB ID: 1t2d) superposed using THESEUS. The ligands (oxalate and
776 NAD⁺) are from 1t2d and colored WHITE. Side chains of important residues are shown as sticks
777 and the six-residue insert of *PfLDH* is highlighted in YELLOW. Note how the *PfLDH* Trp107f
778 overlays Arg102 from *CpMDH*. Residues in the insertion are labeled using numbers and letters
779 to maintain consistency with homologous positions in the dogfish LDH.

780

781 **Figure 3. Phylogeny of M/LDH superfamily. A.** 1844 taxa. The tree is colored according

782 to function (LDH – vermilion; MDH – blue; HicDH – moss). The N-terminal Rossmann-fold of
783 glucosidases and aspartate dehydrogenases (AspDHs) was used to root the phylogeny. Numbers
784 highlight convergent events of LDH evolution from MDHs: 1 - Canonical LDHs, 2 -

785 Trichomonad LDHs, and 3,4 - apicomplexan LDHs. The shaded clades have aLRT supports of
786 42, 57, 75 and 117, respectively (roughly, an aLRT > 8 is considered highly significant (65)). **B.**

787 **Apicomplexan M/LDH Clade.** A close-up of the apicomplexan portion of the phylogeny in **A**,

788 similarly colored by function. aLRT supports for each group: α -proteobacteria MDHs, 15;

789 apicomplexan LDHs, 11; *Plasmodium* LDHs, 333; *Cryptosporidium* MDHs, 54;

790 *Cryptosporidium* LDHs, 202. Ancestral reconstructed proteins are labeled at internal nodes

791 (AncMDH1, AncMDH2, AncMDH3, AncLDH). The focus of the present work is the gene
792 duplication at node 3.

793

794 **Figure 4. Specificity switching in apicomplexan M/LDHs.** Blue horizontal bars (left)
795 quantify activity towards oxaloacetate; red horizontal bars (right) quantify activity towards
796 pyruvate. Error bars are shown as small black brackets and represent 1 SD from triplicate
797 measurements. INS refers to the presence of the six-residue insertion from *Pf*LDH, DEL refers
798 to the removal of the six-residue insertion. Relative specificity (RS) is the ratio of k_{cat}/K_M for
799 pyruvate vs oxaloacetate, with positive $\log_{10}(RS)$ representing a preference for pyruvate and
800 negative $\log_{10}(RS)$ representing a preference for oxaloacetate. All logarithms are base 10.

801

802 **Figure 5. Evolution of novel LDHs in Apicomplexa.** The y-axis of the bar graphs is
803 $\log(k_{cat}/K_M)$, with oxaloacetate in blue and pyruvate in vermilion. Blue vertical bars (left)
804 quantify activity of the given enzyme towards oxaloacetate; red vertical bars (right) quantify
805 activity towards pyruvate. *Rb*MDH is a representative α -proteobacterial MDH from *Rickettsia*
806 *bellii*. *T. gondii* has two LDH proteins (TgLDH1 and TgLDH2), each expressed at different
807 stages of the life cycle (25).

808

809 **Figure 6. Specificity switching in ancestral MDH2.** INS refers to the reconstructed six-
810 residue insertion from AncLDH. 59Mut is described in the text. Relative specificity (RS) is
811 described in legend of **Figure 4**.

812

813 **Figure 7. Ancestral and modern dehydrogenase structures. A. Superposition of**
814 **CpMDH and AncMDH2.** Superposition of AncMDH2 structure (blue) and CpMDH
815 (aquamarine, PDB ID: 2hjr). Ligands from AncMDH2 are shown in gray; ligands from CpMDH
816 are in white. **B. Active site detail of Cp MDH and AncMDH2.** Side chains of catalytic

817 residues highlighted as sticks. **C. Superposition of apicomplexan LDHs and AncLDH*.**
818 Superposition of AncLDH* structure (vermilion) and four apicomplexan LDHs (deep olive,
819 *Pf*LDH (PDB ID: 1t2d), *Plasmodium berghei* (*Pb*) LDH (PDB ID: 1oc4), *Tg*LDH1 (PDB ID: 1pzh),
820 *Tg*LDH2 (PDB ID: 1sow)). Ligands from AncLDH* are shown in gray, ligands from
821 apicomplexan LDHs are in white. The “opposing loop” and residues 236 and 246 (discussed in
822 text) are highlighted in cyan. **D. Active site detail of apicomplexan LDHs and**
823 **AncLDH*.** Side chains of catalytic residues highlighted as sticks. **E. Superposition of**
824 **ancestral dehydrogenases.** Superposition of AncMDH2 (blue), AncLDH* (vermilion), and
825 AncMDH2-INS (magenta). Ligands are shown in gray. **F. Active site detail of ancestral**
826 **dehydrogenases.** Side chains of catalytic residues highlighted as sticks.

827

828 **Figure 8. Alternative LDH mutations in AncMDH2.** Relative specificity (RS) is
829 described in legend of **Figure 4.**

830

831 **Figure 1-figure supplement 1. Fold architecture in the LDH and MDH superfamily.**

832 At left is *Cp*MDH (blue, PDB ID: 2hjr), at right is *Pf*LDH (vermilion and olive, PDB ID: 1t2d).
833 The Rossmann fold domain, which binds the NADH cofactor, is show as light blue in *Cp*MDH
834 and vermilion in *Pf*LDH. The active site is found at the interface of the two domains. In
835 *Cp*MDH, the “opposing loop” is highlighted in yellow (see text). In *Pf*LDH, the six-residue
836 insertion is highlighted in yellow.

837

838 **Figure 2-figure supplement 1. Sequence alignment of the specificity loop from**
839 **apicomplexan M/LDHs with ancestral sequences.**

840

841 **Figure 2-figure supplement 2. Alanine scanning of *Pf*LDH specificity loop.**
842 Logarithm of pyruvate k_{cat}/K_M of *Pf*LDH and each mutant. Labels on x-axis describe the
843 mutation tested in the WT *Pf*LDH background.
844
845 **Figure 2-figure supplement 3. Crystal structure of *Pf*LDH-W107fA mutant.** A.
846 Crystal lattice of the W107fA mutant (left) compared to the WT *Pf*LDH (right). B. Superposition
847 of the WT *Pf*LDH (olive) and the W107fA mutant (vermilion). The structures are highly similar
848 throughout, expect for the active site loop (at top), which is closed in the WT and partially
849 disordered and open in the mutant.
850
851 **Figure 3-figure supplement 1. Phylogeny of M/LDH superfamily.** Same phylogeny as
852 **Figure 3A** with select branch supports shown (aLRT supports).
853
854 **Figure 3-figure supplement 2. Phylogeny of M/LDH superfamily.** Same phylogeny as
855 **Figure 3A**. The tree is colored by domain of life (eubacterial – vermilion; eukaryotic – blue;
856 archeal – magenta).
857
858 **Figure 3-figure supplement 3. Apicomplexan M/LDH Clade.** Same phylogeny as
859 **Figure 3B** with aLRT branch supports and clades shown in full detail.
860
861 **Figure 5-figure supplement1. Sequence identity of ancestral and modern proteins.**
862
863 **Figure 6-figure supplement 1-6. Histograms of ancestral reconstructions.**
864 Reconstructed residues binned according to posterior probability (PP) of the predicted residue.
865

866 **Figure 6-figure supplement 7. Alternative ancestral enzymes.** INS refers to the
867 reconstructed six amino acid insertion from AncLDH*. 58Mut refers to remaining residue
868 differences between AncMDH* and AncLDH* that are not R102K or INS. Relative specificity
869 (RS) is described in legend of **Figure 4**.

870

871 **Figure 7-figure supplement 1. Statistics table for AncMDH2 structures.** Statistics for
872 highest resolution shell are shown in parentheses.

873

874 **Figure 7-figure supplement 2. Statistics table for AncLDH* structures.** Statistics for
875 highest resolution shell are shown in parentheses.

876

877 **Figure 7-figure supplement 3. Statistics table for AncMDH2-INS structures.**
878 Statistics for highest resolution shell are shown in parentheses.

879

880 **Figure 2 -source data 1. Source data for figure supplement 2. Kinetic parameters**
881 **for PflDH alanine-scan.**

882

883 **Figure 4-source data 1. Kinetic parameters for modern constructs.**

884

885 **Figure 5-source data 1. Kinetic parameters for ancestral/modern phylogeny.**

886

887 **Figure 6-source data 1. Kinetic parameters for ancestral specificity switch mutants.**

888

889 **Figure 6 -source data 2. Source data for figure supplement 7. Kinetic parameters**
890 **for alternative ancestral proteins.**

891

892 **Figure 8-source data 1. Kinetic parameters for specificity residue mutants.**

893

894 **Supplementary File 1. Sequences, alignments, and trees.** Alignments and tree files for
895 both the original (**Figure 3**) and the alternative phylogeny. Alignment for **Figure 5-figure**
896 **supplement 1.** Ancestral FASTA files and posterior probabilities for each ancestral sequence
897 (parsed in **Figure 6-figure supplement 1-6.**

898

899 **Supplementary File 2. Molecular weights and extinction coefficients.** ExPASy
900 calculated molecular weights and extinction coefficients for all proteins used within this study.

901

902

References

- 904 1. Douzery EJP, Snell EA, Bapteste E, Delsuc F, Philippe H. The timing of eukaryotic
905 evolution: Does a relaxed molecular clock reconcile proteins and fossils? *Proceedings of the*
906 *National Academy of Sciences of the United States of America*. 2004 Oct 26;101(43):15386-91.
- 907 2. Madern D. Molecular Evolution Within the L-Malate and L-Lactate Dehydrogenase
908 Super-Family. *Journal of Molecular Evolution*. 2002;54(6):825-40.
- 909 3. Zhu G, Keithly JS. Alpha-proteobacterial relationship of apicomplexan lactate and
910 malate dehydrogenases. *The Journal of eukaryotic microbiology*. 2002 May-Jun;49(3):255-61.
- 911 4. Golding GB, Dean AM. The structural basis of molecular adaptation. *Molecular biology*
912 *and evolution*. 1998 Apr;15(4):355-69.
- 913 5. Royer RE, Deck LM, Campos NM, Hunsaker LA, Vanderjagt DL. Biologically-Active
914 Derivatives of Gossypol - Synthesis and Antimalarial Activities of Peri-Acylated Gossylic Nitriles.
915 *Journal of Medicinal Chemistry*. 1986 Sep;29(9):1799-801.
- 916 6. Cameron A, Read J, Tranter R, Winter VJ, Sessions RB, Brady RL, et al. Identification
917 and activity of a series of azole-based compounds with lactate dehydrogenase-directed anti-
918 malarial activity. *Journal of Biological Chemistry*. 2004 Jul 23;279(30):31429-39.
- 919 7. Conners R, Schambach F, Read J, Cameron A, Sessions RB, Vivas L, et al. Mapping the
920 binding site for gossypol-like inhibitors of *Plasmodium falciparum* lactate dehydrogenase.
921 *Molecular and Biochemical Parasitology*. 2005 Aug;142(2):137-48.
- 922 8. Gomez MS, Piper RC, Hunsaker LA, Royer RE, Deck LM, Makler MT, et al. Substrate
923 and cofactor specificity and selective inhibition of lactate dehydrogenase from the malarial
924 parasite *P-falciparum*. *Molecular and Biochemical Parasitology*. 1997 Dec 1;90(1):235-46.
- 925 9. Read JA, Wilkinson KW, Tranter R, Sessions RB, Brady RL. Chloroquine binds in the
926 cofactor binding site of *Plasmodium falciparum* lactate dehydrogenase. *Journal of Biological*
927 *Chemistry*. 1999 Apr 9;274(15):10213-8.
- 928 10. Rossmann MG, Liljas, A., Branden, C-I., Banaszak, L. J. Evolutionary and Structural
929 Relationships among Dehydrogenases. In: Boyer PD, editor. *The Enzymes*. 3rd ed: Academic
930 Press; 1975. p. 61-102.
- 931 11. Birkoft JJ, Banaszak LJ. The presence of a histidine-aspartic acid pair in the active site
932 of 2-hydroxyacid dehydrogenases. X-ray refinement of cytoplasmic malate dehydrogenase. *J*
933 *Biol Chem*. 1983 Jan 10;258(1):472-82.
- 934 12. Clarke AR, Wigley DB, Chia WN, Barstow D, Atkinson T, Holbrook JJ. Site-directed
935 mutagenesis reveals role of mobile arginine residue in lactate dehydrogenase catalysis. *Nature*.
936 1986 Dec 18-31;324(6098):699-702.
- 937 13. Clarke AR, Wilks HM, Barstow DA, Atkinson T, Chia WN, Holbrook JJ. An investigation
938 of the contribution made by the carboxylate group of an active site histidine-aspartate couple to
939 binding and catalysis in lactate dehydrogenase. *Biochemistry*. 1988 Mar 8;27(5):1617-22.
- 940 14. Hart KW, Clarke AR, Wigley DB, Chia WN, Barstow DA, Atkinson T, et al. The
941 importance of arginine 171 in substrate binding by *Bacillus stearothermophilus* lactate
942 dehydrogenase. *Biochem Biophys Res Commun*. 1987 Jul 15;146(1):346-53.
- 943 15. Hart KW, Clarke AR, Wigley DB, Waldman AD, Chia WN, Barstow DA, et al. A strong
944 carboxylate-arginine interaction is important in substrate orientation and recognition in lactate
945 dehydrogenase. *Biochim Biophys Acta*. 1987 Aug 21;914(3):294-8.
- 946 16. Waldman ADB, Hart KW, Clarke AR, Wigley DB, Barstow DA, Atkinson T, et al. The Use
947 of a Genetically Engineered Tryptophan to Identify the Movement of a Domain of B-
948 *Stearothermophilus* Lactate-Dehydrogenase with the Process Which Limits the Steady-State
949 Turnover of the Enzyme. *Biochemical and Biophysical Research Communications*. 1988 Jan
950 29;150(2):752-9.

- 951 17. Bzik DJ, Fox BA, Gonyer K. Expression of Plasmodium-Falciparum Lactate-
952 Dehydrogenase in Escherichia-Coli. *Molecular and Biochemical Parasitology*. 1993
953 May;59(1):155-66.
- 954 18. Dunn CR, Banfield MJ, Barker JJ, Higham CW, Moreton KM, Turgut-Balik D, et al. The
955 structure of lactate dehydrogenase from Plasmodium falciparum reveals a new target for anti-
956 malarial design. *Nature structural biology*. 1996 Nov;3(11):912-5.
- 957 19. Wilks HM, Hart KW, Feeney R, Dunn CR, Muirhead H, Chia WN, et al. A specific, highly
958 active malate dehydrogenase by redesign of a lactate dehydrogenase framework. *Science*. 1988
959 Dec 16;242(4885):1541-4.
- 960 20. Eventoff W, Rossmann MG, Taylor SS, Torff HJ, Meyer H, Keil W, et al. Structural
961 Adaptations of Lactate-Dehydrogenase Isozymes. *Proceedings of the National Academy of
962 Sciences of the United States of America*. 1977;74(7):2677-81.
- 963 21. Chapman AD, Cortes A, Dafforn TR, Clarke AR, Brady RL. Structural basis of substrate
964 specificity in malate dehydrogenases: crystal structure of a ternary complex of porcine
965 cytoplasmic malate dehydrogenase, alpha-ketomalonate and tetrahydroNAD. *J Mol Biol*. 1999
966 Jan 15;285(2):703-12.
- 967 22. Cendrin F, Chroboczek J, Zaccai G, Eisenberg H, Mevarech M. Cloning, sequencing, and
968 expression in Escherichia coli of the gene coding for malate dehydrogenase of the extremely
969 halophilic archaeobacterium Haloarcula marismortui. *Biochemistry*. 1993 Apr 27;32(16):4308-13.
- 970 23. Nicholls DJ, Miller J, Scawen MD, Clarke AR, Holbrook JJ, Atkinson T, et al. The
971 importance of arginine 102 for the substrate specificity of Escherichia coli malate
972 dehydrogenase. *Biochem Biophys Res Commun*. 1992 Dec 15;189(2):1057-62.
- 973 24. Brown WM, Yowell CA, Hoard A, Jagt TAV, Hunsaker LA, Deck LM, et al. Comparative
974 structural analysis and kinetic properties of lactate dehydrogenases from the four species of
975 human malarial parasites. *Biochemistry*. 2004 May 25;43(20):6219-29.
- 976 25. Dando C, Schroeder ER, Hunsaker LA, Deck LM, Royer RE, Zhou XL, et al. The kinetic
977 properties and sensitivities to inhibitors of lactate dehydrogenases (LDH1 and LDH2) from
978 Toxoplasma gondii: comparisons with pLDH from Plasmodium falciparum. *Molecular and
979 Biochemical Parasitology*. 2001 Nov;118(1):23-32.
- 980 26. Kavanagh KL, Elling RA, Wilson DK. Structure of Toxoplasma gondii LDH1: Active-site
981 differences from human lactate dehydrogenases and the structural basis for efficient APAD(+)
982 use. *Biochemistry*. 2004 Feb 3;43(4):879-89.
- 983 27. Shoemark DK, Cliff MJ, Sessions RB, Clarke AR. Enzymatic properties of the lactate
984 dehydrogenase enzyme from Plasmodium falciparum. *Febs Journal*. 2007 Jun;274(11):2738-48.
- 985 28. Winter VJ, Cameron A, Tranter R, Sessions RB, Brady RL. Crystal structure of
986 Plasmodium berghei lactate dehydrogenase indicates the unique structural differences of these
987 enzymes are shared across the Plasmodium genus. *Molecular and Biochemical Parasitology*.
988 2003 Sep;131(1):1-10.
- 989 29. Innan H, Kondrashov F. The evolution of gene duplications: classifying and
990 distinguishing between models. *Nature Reviews Genetics*. 2010 Feb;11(2):97-108.
- 991 30. Ohno S. *Evolution by Gene Duplication*. New York: Springer; 1970.
- 992 31. Lynch M, Conery JS. The evolutionary fate and consequences of duplicate genes. *Science*.
993 2000 Nov 10;290(5494):1151-5.
- 994 32. Walsh JB. How Often Do Duplicated Genes Evolve New Functions. *Genetics*. 1995
995 Jan;139(1):421-8.
- 996 33. Conant GC, Wolfe KH. Turning a hobby into a job: How duplicated genes find new
997 functions. *Nature Reviews Genetics*. 2008 Dec;9(12):938-50.
- 998 34. Force A, Lynch M, Pickett FB, Amores A, Yan YL, Postlethwait J. Preservation of
999 duplicate genes by complementary, degenerative mutations. *Genetics*. 1999 Apr;151(4):1531-45.
- 1000 35. Soskine M, Tawfik DS. Mutational effects and the evolution of new protein functions.
1001 *Nat Rev Genet*. 2010 Aug;11(8):572-82.

- 1002 36. Bridgham JT, Brown JE, Rodriguez-Mari A, Catchen JM, Thornton JW. Evolution of a
1003 new function by degenerative mutation in cephalochordate steroid receptors. *PLoS Genet.*
1004 2008;4(9):e1000191.
- 1005 37. Voordeckers K, Brown CA, Vanneste K, van der Zande E, Voet A, Maere S, et al.
1006 Reconstruction of Ancestral Metabolic Enzymes Reveals Molecular Mechanisms Underlying
1007 Evolutionary Innovation through Gene Duplication. *Plos Biology.* 2012 Dec;10(12).
- 1008 38. Zhang JZ, Rosenberg HF. Complementary advantageous substitutions in the evolution of
1009 an antiviral RNase of higher primates. *Proceedings of the National Academy of Sciences of the*
1010 *United States of America.* 2002 Apr 16;99(8):5486-91.
- 1011 39. Madern D, Cai XM, Abrahamsen MS, Zhu G. Evolution of *Cryptosporidium parvum*
1012 lactate dehydrogenase from malate dehydrogenase by a very recent event of gene duplication.
1013 *Molecular biology and evolution.* 2004 Mar;21(3):489-97.
- 1014 40. Wu G, Fiser A, ter Kuile B, Sali A, Muller M. Convergent evolution of *Trichomonas*
1015 *vaginalis* lactate dehydrogenase from malate dehydrogenase. *Proc Natl Acad Sci U S A.* 1999
1016 May 25;96(11):6285-90.
- 1017 41. Gajria B, Bahl A, Brestelli J, Dommer J, Fischer S, Gao X, et al. ToxoDB: an integrated
1018 *Toxoplasma gondii* database resource. *Nucleic acids research.* 2008 Jan;36(Database
1019 issue):D553-6.
- 1020 42. Aurrecochea C, Brestelli J, Brunk BP, Dommer J, Fischer S, Gajria B, et al. PlasmoDB: a
1021 functional genomic database for malaria parasites. *Nucleic acids research.* 2009
1022 Jan;37(Database issue):D539-43.
- 1023 43. Heiges M, Wang H, Robinson E, Aurrecochea C, Gao X, Kaluskar N, et al. CryptoDB: a
1024 *Cryptosporidium* bioinformatics resource update. *Nucleic acids research.* 2006 Jan
1025 1;34(Database issue):D419-22.
- 1026 44. Templeton TJ, Enomoto S, Chen WJ, Huang CG, Lancto CA, Abrahamsen MS, et al. A
1027 genome-sequence survey for *Ascogregarina taiwanensis* supports evolutionary affiliation but
1028 metabolic diversity between a Gregarine and *Cryptosporidium*. *Molecular biology and evolution.*
1029 2009 Feb;27(2):235-48.
- 1030 45. Feng ZP, Zhang X, Han P, Arora N, Anders RF, Norton RS. Abundance of intrinsically
1031 unstructured proteins in *P. falciparum* and other apicomplexan parasite proteomes. *Mol*
1032 *Biochem Parasitol.* 2006 Dec;150(2):256-67.
- 1033 46. Kissinger JC, DeBarry J. Genome cartography: charting the apicomplexan genome.
1034 *Trends Parasitol.* 2011 Aug;27(8):345-54.
- 1035 47. Vedadi M, Lew J, Artz J, Amani M, Zhao Y, Dong AP, et al. Genome-scale protein
1036 expression and structural biology of *Plasmodium falciparum* and related Apicomplexan
1037 organisms. *Molecular and Biochemical Parasitology.* 2007 Jan;151(1):100-10.
- 1038 48. Garrett R, Grisham, C.M. *Biochemistry.* 3rd ed. Belmont, CA: Thomson Brooks/Cole;
1039 2005.
- 1040 49. Des Marais DL, Rausher MD. Escape from adaptive conflict after duplication in an
1041 anthocyanin pathway gene. *Nature.* 2008 Aug 7;454(7205):762-U85.
- 1042 50. Harms MJ, Thornton JW. Analyzing protein structure and function using ancestral gene
1043 reconstruction. *Curr Opin Struct Biol.* 2010 Jun;20(3):360-6.
- 1044 51. Chen RD, Greer A, Dean AM. A Highly-Active Decarboxylating Dehydrogenase with
1045 Rationally Inverted Coenzyme Specificity. *Proceedings of the National Academy of Sciences of*
1046 *the United States of America.* 1995 Dec 5;92(25):11666-70.
- 1047 52. Jermann TM, Opitz JG, Stackhouse J, Benner SA. RECONSTRUCTING THE
1048 EVOLUTIONARY HISTORY OF THE ARTIODACTYL RIBONUCLEASE SUPERFAMILY.
1049 *Nature.* 1995 Mar;374(6517):57-9.
- 1050 53. Wouters MA, Liu K, Riek P, Husain A. A despecialization step underlying evolution of a
1051 family of serine proteases. *Molecular Cell.* 2003 Aug;12(2):343-54.

1052 54. Bridgham JT, Carroll SM, Thornton JW. Evolution of hormone-receptor complexity by
1053 molecular exploitation. *Science*. 2006 Apr 7;312(5770):97-101.

1054 55. Carroll SM, Bridgham JT, Thornton JW. Evolution of Hormone Signaling in
1055 Elasmobranchs by Exploitation of Promiscuous Receptors. *Molecular biology and evolution*.
1056 2008 Dec;25(12):2643-52.

1057 56. Carroll SM, Ortlund EA, Thornton JW. Mechanisms for the Evolution of a Derived
1058 Function in the Ancestral Glucocorticoid Receptor. *Plos Genetics*. 2011 Jun;7(6).

1059 57. Ortlund EA, Bridgham JT, Redinbo MR, Thornton JW. Crystal structure of an ancient
1060 protein: Evolution by conformational epistasis. *Science*. 2007 Sep 14;317(5844):1544-8.

1061 58. Risso VA, Gavira JA, Mejia-Carmona DF, Gaucher EA, Sanchez-Ruiz JM. Hyperstability
1062 and Substrate Promiscuity in Laboratory Resurrections of Precambrian beta-Lactamases.
1063 *Journal of the American Chemical Society*. 2013 Feb 27;135(8):2899-902.

1064 59. Pruitt KD, Tatusova T, Klimke W, Maglott DR. NCBI Reference Sequences: current
1065 status, policy and new initiatives. *Nucleic acids research*. 2009 Jan;37(Database issue):D32-6.

1066 60. Altschul SF, Gish W, Miller W, Myers EW, Lipman DJ. Basic local alignment search tool.
1067 *J Mol Biol*. 1990 Oct 5;215(3):403-10.

1068 61. Consortium TU. Update on activities at the Universal Protein Resource (UniProt) in
1069 2013. *Nucleic acids research*. 2013 Jan;41(Database issue):D43-7.

1070 62. Edgar RC. MUSCLE: multiple sequence alignment with high accuracy and high
1071 throughput. *Nucleic acids research*. 2004;32(5):1792-7.

1072 63. Guindon S, Dufayard JF, Lefort V, Anisimova M, Hordijk W, Gascuel O. New algorithms
1073 and methods to estimate maximum-likelihood phylogenies: assessing the performance of
1074 PhyML 3.0. *Systematic biology*. 2010 May;59(3):307-21.

1075 64. Le SQ, Gascuel O. An improved general amino acid replacement matrix. *Molecular
1076 biology and evolution*. 2008 Jul;25(7):1307-20.

1077 65. Anisimova M, Gascuel O. Approximate likelihood-ratio test for branches: A fast, accurate,
1078 and powerful alternative. *Systematic biology*. 2006 Aug;55(4):539-52.

1079 66. Rao ST, Rossmann MG. Comparison of super-secondary structures in proteins. *J Mol
1080 Biol*. 1973 May 15;76(2):241-56.

1081 67. Theobald DL, Wuttke DS. Divergent evolution within protein superfolds inferred from
1082 profile-based phylogenetics. *J Mol Biol*. 2005 Dec 2;354(3):722-37.

1083 68. Chandonia JM, Hon G, Walker NS, Lo Conte L, Koehl P, Levitt M, et al. The ASTRAL
1084 Compendium in 2004. *Nucleic acids research*. 2004 Jan 1;32(Database issue):D189-92.

1085 69. Boeckmann B, Bairoch A, Apweiler R, Blatter MC, Estreicher A, Gasteiger E, et al. The
1086 SWISS-PROT protein knowledgebase and its supplement TrEMBL in 2003. *Nucleic acids
1087 research*. 2003 Jan 1;31(1):365-70.

1088 70. Sadreyev RI, Baker D, Grishin NV. Profile-profile comparisons by COMPASS predict
1089 intricate homologies between protein families. *Protein Sci*. 2003 Oct;12(10):2262-72.

1090 71. Swofford DL. PAUP*. Phylogenetic Analysis Using Parsimony (*and Other Methods).
1091 Version 4 Sinauer Associatesm Sunderland, Massachusetts. 2003.

1092 72. Lanyon SM. Detecting Internal Inconsistencies in Distance Data. *Systematic Zoology*.
1093 1985 Dec;34(4):397-403.

1094 73. Siddall ME. Another monophyly index: Revisiting the jackknife. *Cladistics-the
1095 International Journal of the Willi Hennig Society*. 1995 Mar;11(1):33-56.

1096 74. Whelan S, Goldman N. A general empirical model of protein evolution derived from
1097 multiple protein families using a maximum-likelihood approach. *Molecular biology and
1098 evolution*. 2001 May;18(5):691-9.

1099 75. Yang Z. PAML 4: phylogenetic analysis by maximum likelihood. *Molecular biology and
1100 evolution*. 2007 Aug;24(8):1586-91.

1101 76. Wolfenden R, Lewis CA, Jr., Yuan Y. Kinetic challenges facing oxalate, malonate,
1102 acetoacetate, and oxaloacetate decarboxylases. *J Am Chem Soc*. Apr 20;133(15):5683-5.

- 1103 77. Parker DM, Holbrook JJ. The oxaloacetate reductase activity of vertebrate lactate
1104 dehydrogenase. *Int J Biochem.* 1981;13(10):1101-5.
- 1105 78. Kabsch W. Xds. *Acta Crystallographica Section D-Biological Crystallography.* 2010
1106 Feb;66:125-32.
- 1107 79. Adams PD, Afonine PV, Bunkoczi G, Chen VB, Davis IW, Echols N, et al. PHENIX: a
1108 comprehensive Python-based system for macromolecular structure solution. *Acta*
1109 *Crystallographica Section D-Biological Crystallography.* 2010 Feb;66:213-21.
- 1110 80. Kelley LA, Sternberg MJE. Protein structure prediction on the Web: a case study using
1111 the Phyre server. *Nature Protocols.* 2009;4(3):363-71.
- 1112 81. Emsley P, Lohkamp B, Scott WG, Cowtan K. Features and development of Coot. *Acta*
1113 *Crystallographica Section D-Biological Crystallography.* 2010 Apr;66:486-501.
- 1114 82. Chen VB, Arendall WB, Headd JJ, Keedy DA, Immormino RM, Kapral GJ, et al.
1115 MolProbity: all-atom structure validation for macromolecular crystallography. *Acta*
1116 *Crystallographica Section D-Biological Crystallography.* 2010 Jan;66:12-21.
- 1117 83. Davis IW, Leaver-Fay A, Chen VB, Block JN, Kapral GJ, Wang X, et al. MolProbity: all-
1118 atom contacts and structure validation for proteins and nucleic acids. *Nucleic acids research.*
1119 2007 Jul;35:W375-W83.
- 1120 84. Theobald DL, Wuttke DS. Accurate structural correlations from maximum likelihood
1121 superpositions. *Plos Computational Biology.* 2008 Feb;4(2).
1122
1123

Figure 1

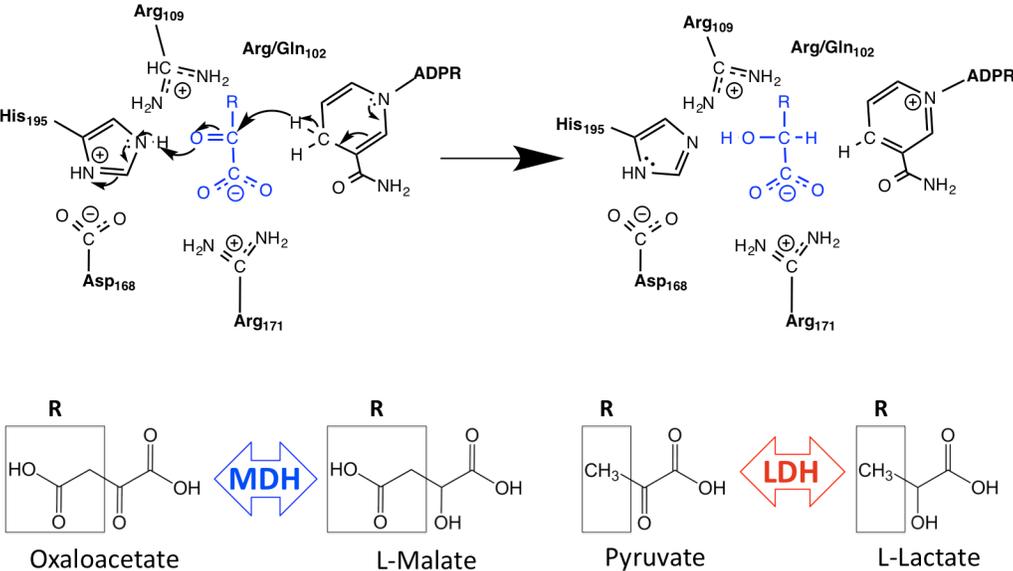


Figure 2

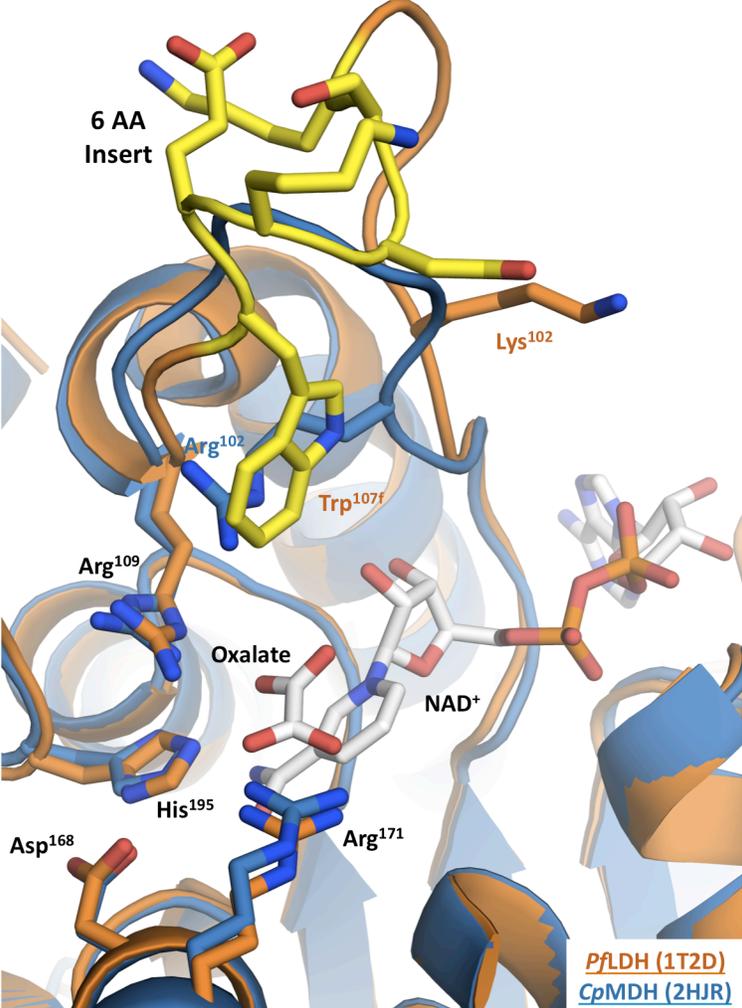


Figure 3

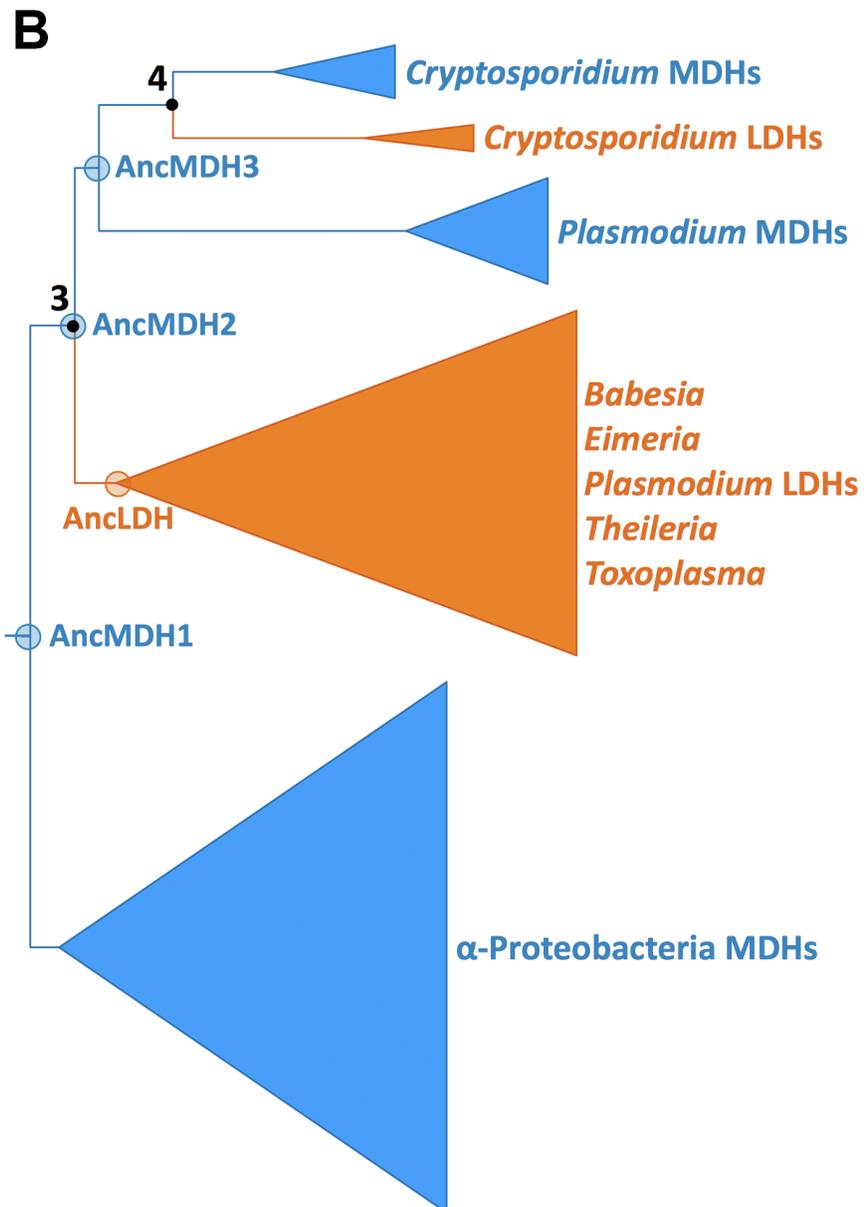
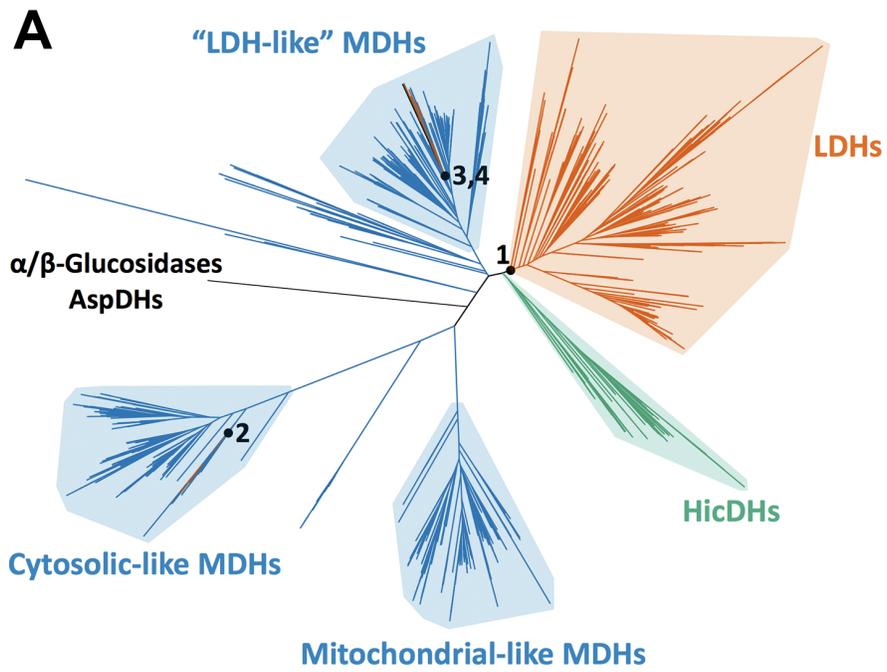


Figure 4

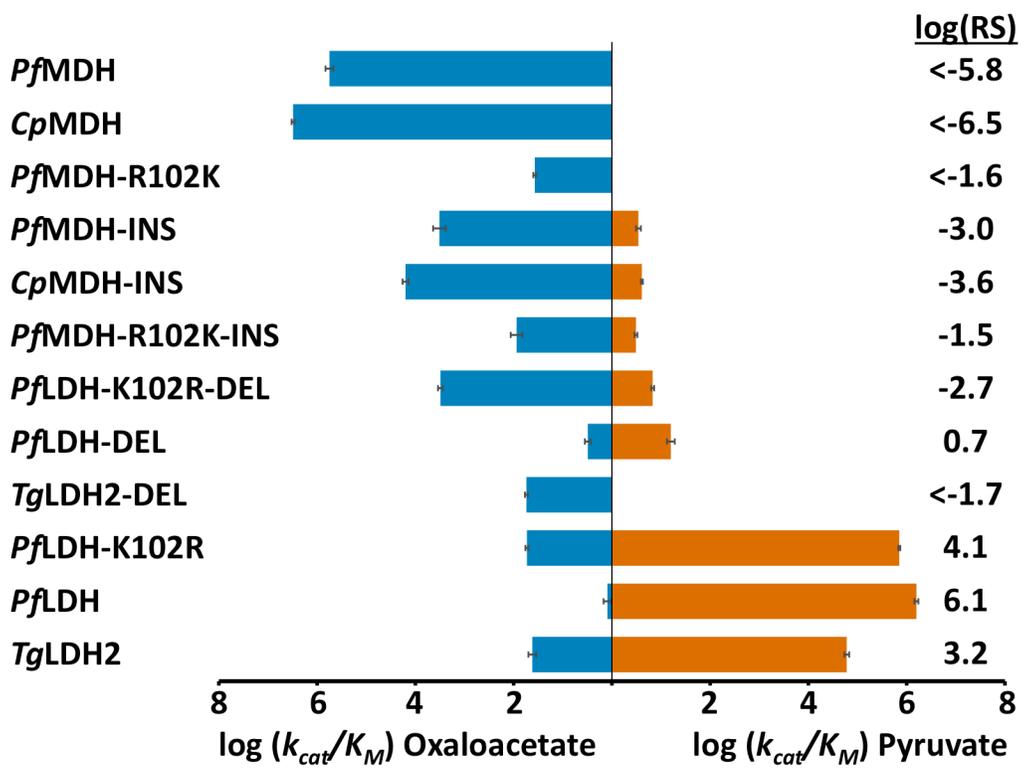


Figure 5

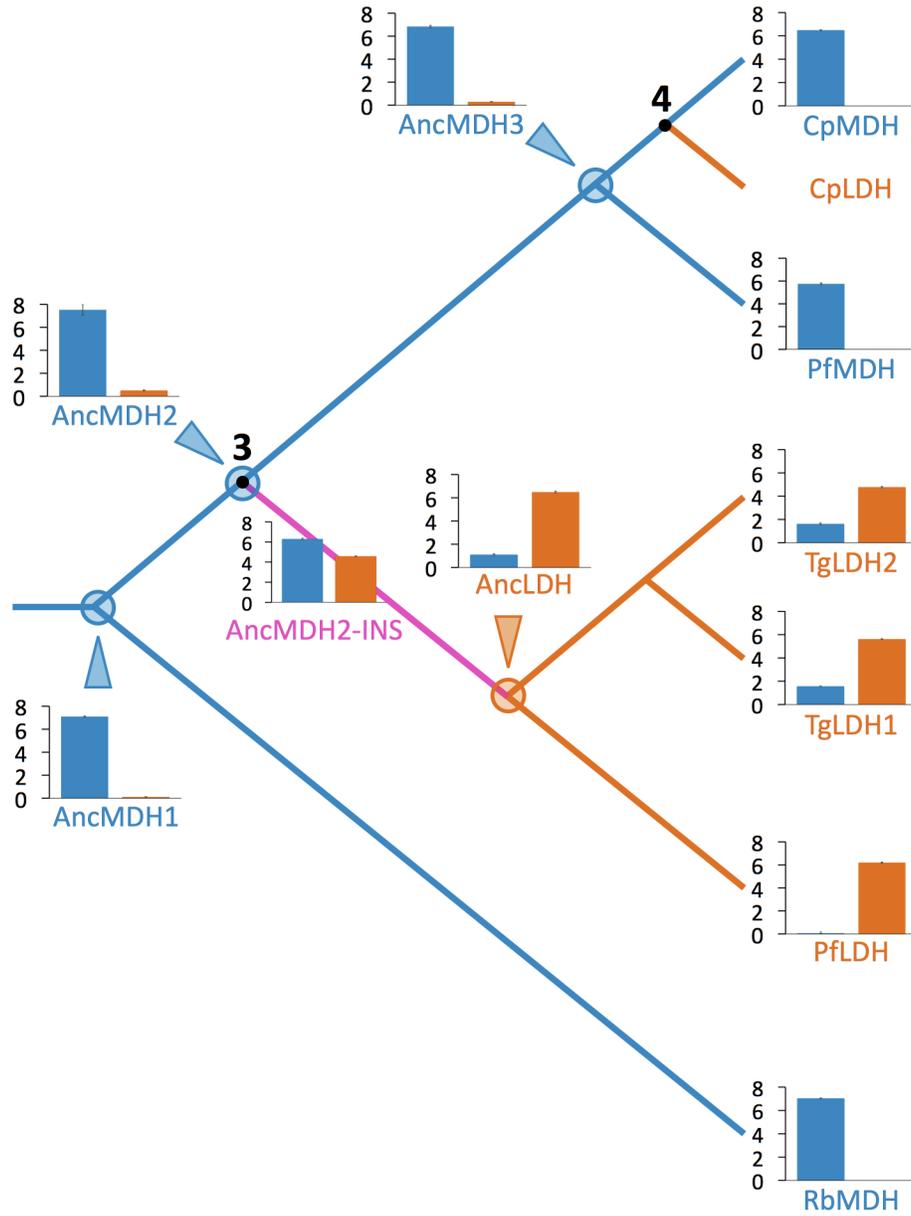


Figure 6

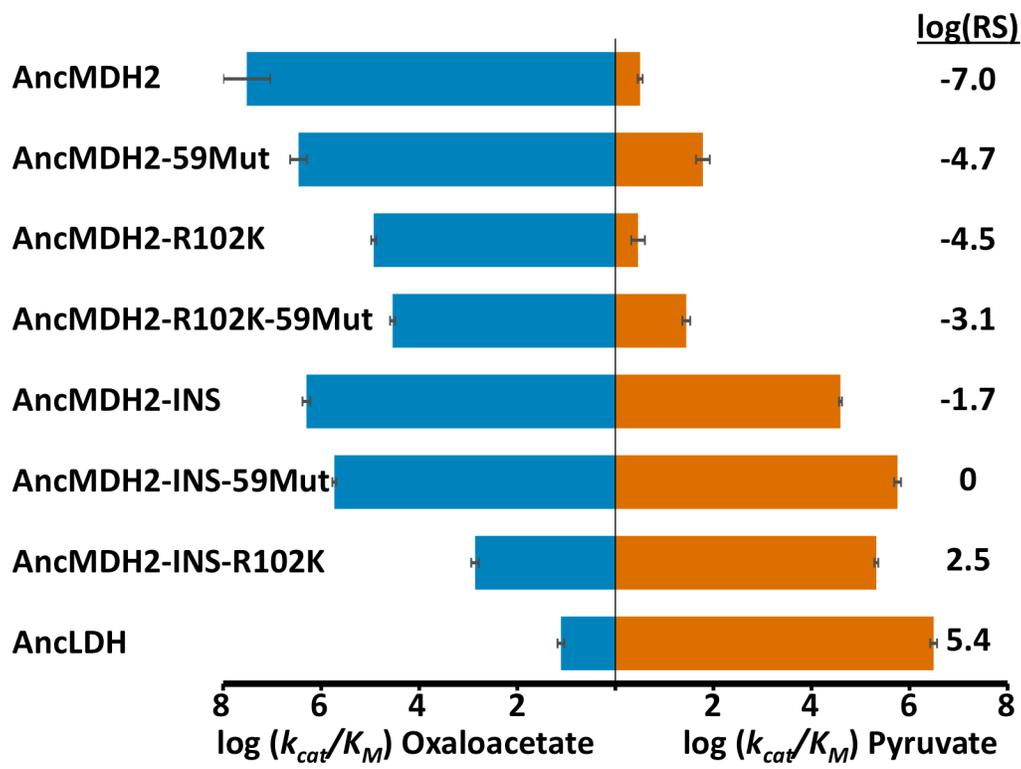
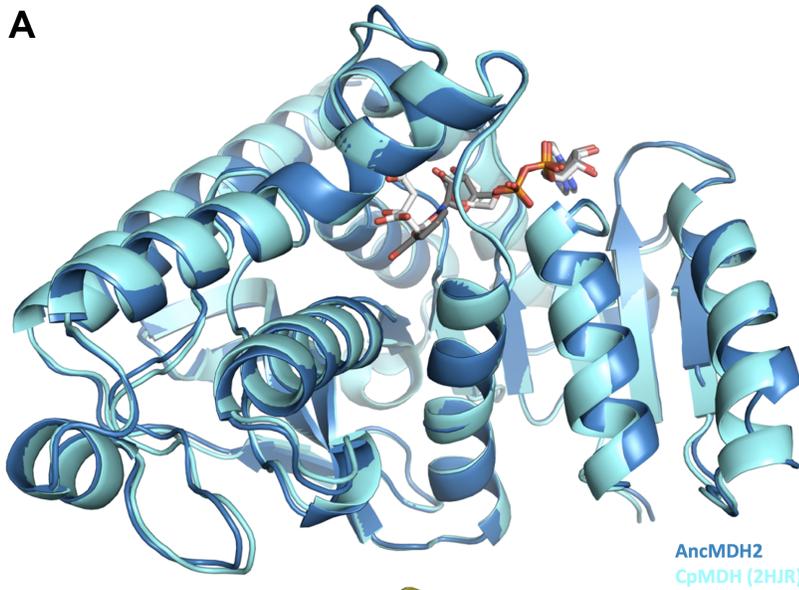
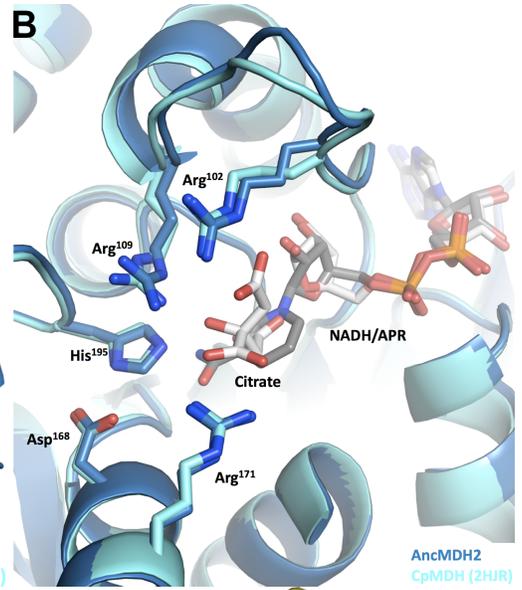


Figure 7

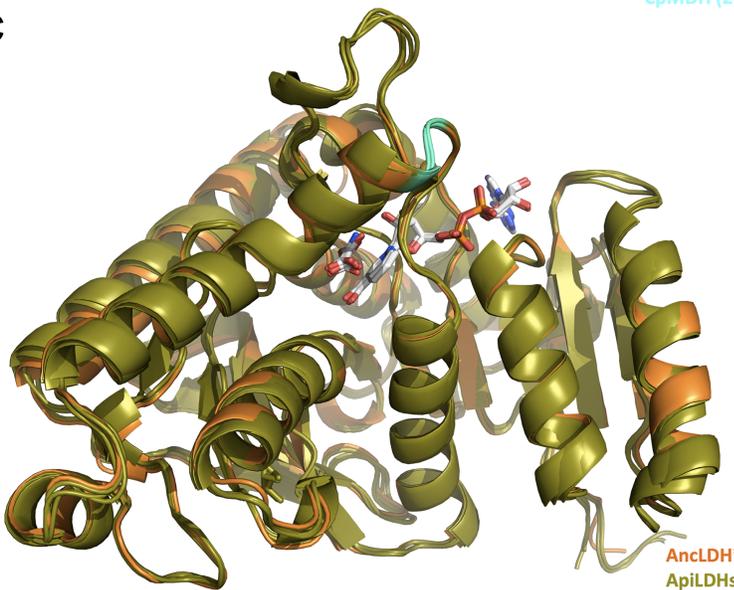
A



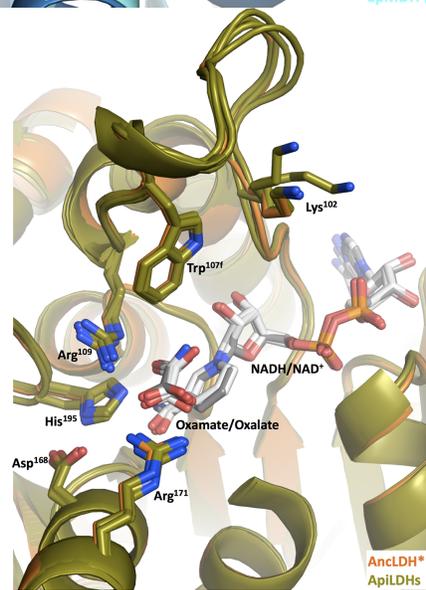
B



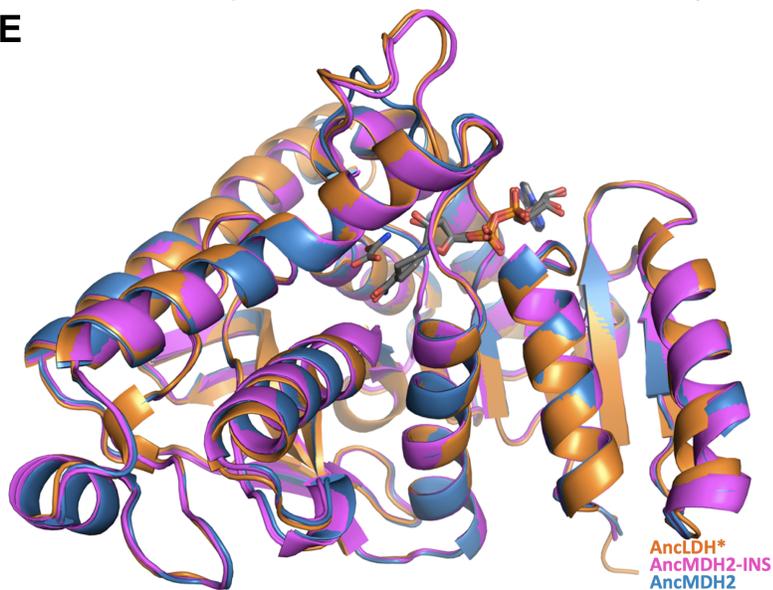
C



D



E



F

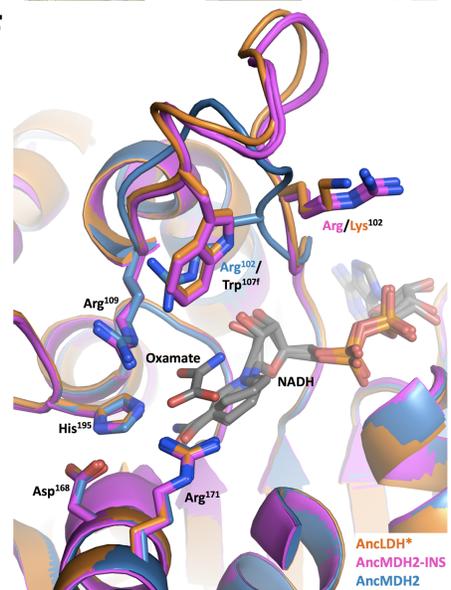


Figure 8

