

EAL50997	XP_656380.1	656	<i>Bacteroides thetaiotaomicron</i>	718	53	<i>Cryptococcus neoformans</i>	770	32	0	2.00E-69	0.00E+00
EAL51149	XP_656535.1	343	<i>Bacteroides fragilis</i>	359	43	<i>Pichia ozymaensis</i>	378	34	8.00E-84	1.00E-53	8.00E-31
			<i>Symbiobacterium thermophilum</i>	274	45	<i>Oryza sativa</i>	315	21	3.00E-51	0.003	1.00E-48
EAL51348	XP_656749.1	171	<i>Methanopyrus kandleri</i>	204	37	<i>Tetrahymena thermophila</i>	2872	22	3.00E-21	0.1	3.00E-20
EAL51525	XP_656903.1	316	<i>Bacteroides thetaiotaomicron</i>	300	29	<i>Candida boidinii</i>	314	32	8.00E-27	0.0007	1.14E-23
EAL51565	XP_656946.1	415	<i>Clostridium perfringens</i>	900	43	<i>Trichomonas vaginalis</i>	897	40	1.00E-89	5.00E-81	2.00E-09
EAL51925	XP_657304.1	448	<i>T. tengcongensis</i>	481	43	<i>Giardia lamblia</i>	937	33	3.00E-96	2.00E-60	1.50E-36
EAL52001	XP_657387.1	303	<i>Oceanobacillus theyensis</i>	306	27				2.00E-15	0.00E+00	
27 LGT cases that are more weakly supported than before according to our criteria											
EAL45152	XP_650539.1	122	<i>Shewanella oneidensis</i>	132	34	<i>Trypanosoma brucei</i>	385	24	5.00E-10	6.6	7.58E-11
EAL43347	XP_648734.1	848	<i>Burkholderia pseudomallei</i>	779	38	<i>Plasmodium falciparum</i>	2463	32	1.00E-136	4.00E-44	2.50E-93
EAL44257	XP_649643.1	407	<i>Clostridium acetobutylicum</i>	406	25	<i>Homo sapiens</i>	468	24	6.00E-23	1.00E-14	6.00E-09
EAL45586	XP_650972.1	460	<i>Clostridium tetani</i>	476	47	<i>Xenopus laevis</i>	513	38	1.00E-116	5.00E-84	2.00E-33
EAL46313	XP_651699.1	118	<i>Prochlorococcus marinus</i>	163	42	<i>Hordeum vulgare</i>	223	22	2.00E-21	1.4	1.43E-21
EAL46399	XP_651785.1	218	<i>Clostridium perfringens</i>	235	65	<i>Trypanosoma brucei</i>	295	52	3.00E-73	9.00E-54	3.33E-20
EAL46421	XP_651808.1	205	<i>Clostridium acetobutylicum</i>	230	40	<i>Arabidopsis thaliana</i>	241	33	7.00E-34	6.00E-12	1.17E-22
EAL46701	XP_652087.1	294	<i>Bacteroides fragilis</i>	308	45	<i>Thalassiosira pseudonana</i>	348	27	4.00E-63	1.00E-14	4.00E-49
EAL46757	XP_652143.1	95	<i>Lactococcus lactis</i>	103	31	<i>Tetrahymena thermophila</i>	112	32	3.00E-09	1.00E-07	3.00E-02
EAL46858	XP_652245.1	192	<i>Pseudomonas aeruginosa</i>	195	41	<i>Caenorhabditis briggsae</i>	229	40	6.00E-36	2.00E-17	3.00E-19
EAL47026	XP_652397.1	164	<i>Bacillus subtilis</i>	181	30	<i>Trichomonas vaginalis</i>	182	26	3.00E-10	2.00E-08	1.50E-02
EAL47464	XP_652839.1	504	<i>Treponema denticola</i>	509	39	<i>Piromyces sp.</i>	555	27	5.00E-88	2.00E-30	2.50E-58
EAL47648	XP_653034.1	259	<i>Methanosarcina mazi</i>	272	36	<i>Arabidopsis thaliana</i>	345	25	2.00E-39	3.00E-11	6.67E-29
EAL47787	XP_653173.1	546	<i>Spirochaeta thermophila</i>	571	56	<i>Solanum tuberosum</i>	552	46	1.00E-175	1.00E-135	1.00E-40
EAL48186	XP_653572.1	232	<i>Bacillus cereus</i>	279	34	<i>Thalassiosira pseudonana</i>	271	32	2.00E-10	2.00E-08	1.00E-02
EAL49309	XP_654698.1	358	<i>Methanosarcina mazi</i>	379	42	<i>Leishmania major</i>	373	31	5.00E-77	9.00E-44	5.56E-34
EAL48568	XP_653954.1	113	<i>Chlamydia pneumoniae</i>	271	38	<i>Debaryomyces hansenii</i>	699	38	5.00E-14	7.00E-16	7.14E+01
EAL48767	XP_654156.1	165	<i>Bacteroides fragilis</i>	177	40	<i>Trichomonas vaginalis</i>	189	28	7.00E-28	2.00E-05	3.50E-23
EAL48783	XP_654172.1	217	<i>Pseudomonas putida</i>	225	46	<i>Giardia lamblia</i>	239	35	2.00E-43	7.00E-24	2.86E-20
EAL49703	XP_655090.1	396	<i>Clostridium acetobutylicum</i>	398	34	<i>Tetrahymena thermophila</i>	445	29	4.00E-64	3.00E-44	1.33E-20
EAL49996	XP_655383.1	358	<i>Bacteroides thetaiotaomicron</i>	368	60	<i>Brachydanio rerio</i>	367	43	1.00E-121	5.00E-76	2.00E-46
EAL50325	XP_655711.1	447	<i>Clostridium tetani</i>	448	30	<i>Trichomonas vaginalis</i>	871	29	4.00E-46	1.00E-37	4.00E-09
EAL50521	XP_655905.1	285	<i>Streptococcus agalactiae</i>	323	29	<i>Leishmania major</i>	452	24	2.00E-22	3.00E-06	6.67E-17
EAL50620	XP_656005.1	261	<i>Wolinella succinogenes</i>	655	27	<i>Trichomonas vaginalis</i>	261	28	6.00E-21	1.00E-06	6.00E-15
EAL50838	XP_656225.1	299	<i>Anabaena sp.</i>	287	27	<i>Trichomonas vaginalis</i>	336	29	4.00E-15	0.0009	4.44E-12
EAL50986	XP_656369.1	219	<i>Bacteroides thetaiotaomicron</i>	240	31	<i>Xenopus laevis</i>	309	29	2.00E-20	1.00E-12	2.00E-08
EAL52121	XP_657511.1	220	<i>T. tengcongensis</i>	222	36	<i>Caenorhabditis elegans</i>	255	26	1.00E-30	1.00E-07	1.00E-23
14 cases where increased sampling has weakened that case for LGT											
EAL42539	XP_647925.1	213	<i>Bacteroides thetaiotaomicron</i>	319	47	<i>Entodinium caudatum</i>	411	43	3.00E-53	1.00E-32	3.00E-21
EAL42738	XP_648124.1	313	<i>Campylobacter jejuni</i>	324	40	<i>Trichomonas vaginalis</i>	313	36	1.00E-63	4.00E-42	2.50E-22
EAL44270	XP_649657.1	179	<i>Methanococcus maripaludis</i>	193	37	<i>Anopheles gambiae</i>	186	21	2.00E-27	2.00E-09	1.00E-18

EAL44593	XP_649979.1	220	<i>Vibrio vulnificus</i>	244	24	<i>Trichomonas vaginalis</i>	238	21	0.0002	2.6	7.69E-05
EAL45320	XP_650707.1	154	<i>Geobacillus kaustophilus</i>	183	53	<i>Thalassiosira pseudonana</i>	182	43	8.00E-38	2.00E-32	4.00E-06
EAL45332	XP_650718.1	392	<i>Methanosarcina acetivorans</i>	420	48	<i>Trichomonas vaginalis</i>	396	47	8.00E-99	2.00E-93	4.00E-06
EAL45528	XP_650913.1	349	<i>Sulfolobus acidocaldarius</i>	343	28	<i>Cyanophora paradoxa</i>	313	27	1.00E-24	5.00E-17	2.00E-08
EAL45907	XP_651293.1	380	<i>Streptomyces coelicolor</i>	603	32	<i>Dictyostelium discoideum</i>	457	30	2.00E-39	2.00E-35	1.00E-04
EAL46026	XP_651412.1	176	<i>Bacteroides fragilis</i>	184	51	<i>Tetralymena thermophila</i>	323	32	2.00E-44	8.00E-08	2.50E-37
EAL46116	XP_651488.1	662	<i>Bacillus clausii</i>	684	48	<i>Solanum tuberosum</i>	761	48	0	1.00E-172	0.00E+00
EAL46656	XP_652044.1	419	<i>Dictyoglomis thermophilum</i>	579	30	<i>S. pombe</i>	493	41	2.00E-35	2.00E-19	1.00E-16
EAL50605	XP_655990.1	392	<i>Thermotoga maritima</i>	417	38	<i>Cryptococcus neoformans</i>	445	30	2.00E-69	1.00E-33	2.00E-36
EAL51270	XP_656656.1	251	<i>Porphyromonas gingivalis</i>	261	50	<i>Anopheles gambiae</i>	272	39	6.00E-53	1.00E-35	6.00E-18
EAL52102	XP_657492.1	345	<i>Bacteroides thetaiotaomicron</i>	358	54	<i>Thalassiosira pseudonana</i>	354	47	1.00E-105	7.00E-86	1.43E-20
Nine cases where <i>Entamoeba</i> is now recovered with a recently sequenced gene from another microbial eukaryote											
EAL44213	XP_649600.1	710	<i>Bdellovibrio bacteriovorus</i>	698	37	<i>Trichomonas vaginalis</i>	713	35	1.00E-127	1.00E-127	1.00E+00
EAL44435	XP_649823.1	250	<i>Bacteroides fragilis</i>	395	40	<i>Trichomonas vaginalis</i>	395	33	1.00E-43	3.00E-35	3.33E-09
EAL44766	XP_650152.1	401	<i>Porphyromonas gingivalis</i>	419	36	<i>Trichomonas vaginalis</i>	445	32	3.00E-65	1.00E-51	3.00E-14
EAL47785	XP_653171.1	234	<i>Bacillus anthracis</i>	242	32	<i>Trichomonas vaginalis</i>	256	39	2.00E-30	3.00E-33	6.67E+02
EAL47859	XP_653246.1	337	<i>Clostridium acetobutylicum</i>	322	50	<i>C. reinhardtii</i>	352	44	9.00E-74	0	N/A
EAL49158	XP_654544.1	397	<i>T. tengcongensis</i>	412	49	<i>Trichomonas vaginalis</i>	416	46	1.00E-100	4.00E-99	2.50E-02
EAL49488	XP_654874.1	320	<i>Geobacter sulfurreducens</i>	336	34	<i>Leishmania major</i>	357	31	1.00E-38	4.00E-30	2.50E-09
EAL49791	XP_655177.1	164	<i>Oceanobacillus iheyensis</i>	177	42	<i>Thalassiosira pseudonana</i>	96	38	8.00E-30	6.00E-09	1.33E-21
EAL50404	XP_655790.1	718	<i>T. tengcongensis</i>	717	37	<i>Trichomonas vaginalis</i>	721	34	1.00E-139	1.00E-118	1.00E-21
Five cases where vertical inheritance is now the simplest explanation for the new tree											
EAL44346	XP_649732.1	314	<i>Oceanobacillus iheyensis</i>	239	47	<i>Dictyostelium discoideum</i>	278	65	1.00E-52	3.00E-95	3.33E+42
EAL45466	XP_650849.1	209	<i>Agrobacterium tumefaciens</i>	254	31	<i>Thalassiosira pseudonana</i>	227	35	3.00E-23	1.00E-27	3.00E+04
EAL45548	XP_650934.1	259	<i>Bacillus cereus (strain ZK)</i>	233	29	<i>Candida glabrata</i>	270	30	7.00E-06	5.00E-05	1.40E-01
EAL45595	XP_650981.1	284	<i>Pyrobaculum aerophilum</i>	293	27	<i>Ashbya gossypii</i>	343	27	1.00E-23	7.00E-16	1.43E-08
EAL50185	XP_655571.1	186	<i>Aeropyrum pernix</i>	192	31	<i>Thalassiosira pseudonana</i>	149	30	4.00E-13	5.00E-06	8.00E-08

1 All 96 trees reanalysed here can be downloaded (in pdf format) from the following web site: http://www.ncl.ac.uk/microbial_eukaryotes/

2 ^aGenBank accession numbers and RefSeq accession numbers, respectively, for the 96 original candidates LGT identified by phylogenetic analysis (Lofus *et al.*, 2005)

3 ^bEhL, the length of the *E. histolytica* protein

4 ^cPL/EL, the protein length of the prokaryotic or eukaryotic top BlastP hit, respectively

5 ^d%ID, the percent identity between the *E. histolytica* protein and the top prokaryotic or eukaryotic protein in BlastP alignments (in respective columns)

- 1 ^ePE-score, the e-score of the top prokaryotic hit
- 2 ^fEE-score, the e-score of the top eukaryotic hit
- 3 ^sP/E Ratio, the e-score ratio between the top prokaryotic hit and top eukaryotic hit
- 4 Abbreviated taxon names (to fit the columns):
- 5 *Chlamydomonas reinhardtii*; *C. reinhardtii*; *Paracoccidioides brasiliensis*; *P. brasiliensis*; *Schizosaccharomyces pombe*; *S. pombe*;
- 6 *Thermoanaerobacter tengcongensis*; *T. tengcongensis*
- 7

Table 9. Examples of microarray-detected transcriptional changes in some gene families and the conditions tested

Gene family	Total number of genes in gene family	Number of genes transcriptionally regulated under condition tested	
		Heat shock ^a (1,131 genes on array)	Host colonisation and invasion ^b (9,435 genes on array)
Cysteine proteinases	29 ^c	2 upregulated (CPs 6, 4); 7 down-regulated (CPs 1, 2, 3, 8, 13, 17,	21 genes on array; 4 up-regulated (CPs 1, 9, 4, 6); 1 down-regulated (CP8)
Lectin (Heavy, Light, and Intermediate subunits)	12	1 up-regulated (Hgl-2); 5 down-regulated (Lgl-1 and 3, Igl 1 and 2, Hgl-3)	No change in heavy or intermediate subunits; Light subunit Igl2 and Igl3 down-regulated
Amoebapore	3	1 down-regulated (amoebapore C)	No substantial changes
Transmembrane receptor kinases	>80	NA	6 up-regulated (TMKs 69, 53, 95, 105, 63, 56)
AIG-1 (similar to plant antibacterial proteins)	15	NA	2 down-regulated (TMKs 03 and 17) 5 up-regulated at day 1; 6 down-regulated at day 29 (all non-overlapping)

^a Adapted from Weber *et al.* (2006); ^b Adapted from Gilchrist *et al.* (2006);

^cNumber of cysteine proteinase gene families in genome annotation at time studies were performed

1 Figure Legends

2

3 *Figure 1.*

4 **Positions of introns in the vacuolar ATPase subunit D gene in *P. falciparum*, *D.***
5 ***discoideum*, and *E. histolytica***

6

7 *Figure 2. Comparison of protein sizes in *E. histolytica* and *D. discoideum*.*

8 a: The graph shows the distribution of predicted amino acid length across sequenced
9 genomes from single celled eukaryotes: *D. discoideum* (DD) *Encephalitozoon*
10 *cuniculi* (EC), *P. falciparum* (PF), *E. histolytica* (EH), and *S. cerevisiae* (SC). *E.*
11 *histolytica* and *E. cuniculi* have a distribution that is skewed toward smaller proteins
12 relative to the other species.

13 b: The histogram displays the degree of size change of genes in *E. histolytica* relative
14 to *D. discoideum* when comparing orthologous genes identified by reciprocal best
15 blast hits. The black bars show genes that are smaller in *E. histolytica* where as the
16 grey bars are smaller in *D. discoideum*.

17

18 *Figure 3. Domain diagram of the Hgl subunit of the Gal/GalNAc lectin. CW-*
19 *Cysteine-Tryptophan region; CF- Cysteine free region; C-Rich- Cysteine rich region.*
20 The black vertical box near the carboxy-terminus of the protein represents the single
21 transmembrane domain. The horizontal black bars above the diagram indicate the
22 location of a carbohydrate recognition domain (CRD), the region with similarity to
23 the hepatic growth factor receptor, **c-Met**, and the region that has similarity to the
24 **CD59**, the membrane inhibitor of the complement membrane attack complex. The
25 numbers in parentheses indicate the location of these regions in the Hgl1 isoform
26 (Mann *et al.*, 1991), where the methionine of the immature protein is residue 1.

27

28 *Figure 4. Structural domains of the 3 different types of family C1-like cysteine*
29 **endopeptidases EhCP-A, EhCP-B and EhCP-C.** Shown are the location and length
30 of domains specific for each the 3 types as well as the conserved active site and
31 cysteine residue

32

33 *Figure 5. Predicted antioxidant system of *Entamoeba histolytica*.* A. Superoxide
34 radical anions are detoxified by an iron-containing superoxide dismutase (FeSOD).

1 Molecular oxygen is reduced to hydrogen peroxide by a NADPH:flavin
 2 oxidoreductase (thioredoxin reductase, p34). Hydrogen peroxide is converted to water
 3 by rubrerythrin (Rbr). The nature of its redox partner is unknown. Hydrogen peroxide
 4 can also be converted to water via a classical thioredoxin redox system consisting of
 5 thioredoxin reductase (TrxR, p34), thioredoxin (Trx) and peroxiredoxin (Prx). B.
 6 Nitric oxide is reduced by an A-type flavoprotein (FprA) to nitrous oxide and water.
 7 For this reaction FprA receives electrons from NADH oxidase (Far).

8

9 **Figure 6. A phylogenetic tree of Rab proteins from *Entamoeba histolytica*,**
 10 **human, and yeast.** The number on the nodes represent the bootstrap proportions (%)
 11 of 1000 pseudo samples; only bootstrap proportions >30% are shown. *E. histolytica*
 12 Rab proteins are indicated in bold. Tentative subfamilies that revealed significant
 13 similarity (>40% identity) to their human or yeast counterpart are shaded dark, while
 14 *Entamoeba*-specific subfamilies have light shading. The scale bar indicates 0.1
 15 substitutions at each amino acid position. *: *EhRab* proteins that lack the conserved
 16 effector region, switch regions, or GTP-binding boxes. **: *EhRab* proteins that
 17 possess a non-conventional carboxyl-terminus or lack carboxyl-terminal cysteines.
 18 ***: Rab proteins that were not classified as isotypes based on <40% identity to other
 19 members of the subfamily. References on tree: (1), Temesvari *et al.* (1999); (2),
 20 Rodríguez *et al.* (2000); (3), Saito-Nakano *et al.* (2001); (4), Juárez *et al.* (2001); (5),
 21 Saito-Nakano *et al.* (2004); and (6), Okada *et al.* (2005).

22

23 **Figure 7. Synthesis of N-glycan precursors by *S. cerevisiae* (A) and *E. histolytica***
 24 **(B).** The N-glycan precursor of *S. cerevisiae* contains 14 sugars (Glc₃Man₉GlcNAc₂),
 25 each of which is added by a specific enzyme. The *E. histolytica* N-glycan precursor
 26 contains just seven sugars (Man₅GlcNAc₂), as the protist is missing enzymes that add
 27 mannose and glucose in the lumen of the ER. The figure is redrawn from Figure 1 of
 28 Samuelson *et al.* (2005). Glc = Glucose; GlcNAc = N-acetyl glucosamine; Man =
 29 Mannose.

30

31 **Figure 8. Selected N-glycans of mammals (A-E) and *Entamoeba* (F-H).** Precursors
 32 transferred to nascent peptide (A and F). Glycosylated products involved in N-
 33 glycan-associated QC of protein folding (B and G). Mannosidase product involved in

1 N-glycan-associated protein degradation (mammals only) (C). Trimmed product that
 2 is building block for complex N-glycans (mammals and *Entamoeba*) (D). Complex
 3 N-glycans made in the Golgi (E and H). Glc = Glucose; GlcNAc = N-acetyl
 4 glucosamine; Man = Mannose; Gal = Galactose; Fuc = Fucose.

5

6 **Figure 9. Model of quality control of protein folding in *Entamoeba*.** 1. N-glycan-
 7 dependent QC of protein folding. 2. N-glycan-independent QC of protein folding. 3.
 8 N-glycan-independent ERAD. 4. Ire1 and unfolded protein response (see text for
 9 details).

10

11 **Figure 10. Structure of cysteine-rich plasma membrane proteins of *E. histolytica*.**

12 These proteins include the various subunits of the Gal/GalNAc lectin, a cysteine
 13 protease, and numerous receptor kinases. Ire1, which is involved in the unfolded
 14 protein response, is also a receptor kinase but has no Cys-rich domain.

15

16 **Figure 11. Model for the *Entamoeba* cyst wall derived primarily from**

17 **experiments with *E. invadens*.** A. The cyst wall consists of chitosan fibrils, which
 18 are made by chitin synthase and chitin deacetylase. Wall proteins include Jacob
 19 lectins with tandem arrays of 6-Cys chitin-binding domains (CBDs), as well as
 20 chitinase and Jessie lectins that have a single 8-Cys CBD. The Gal/GalNAc lectin in
 21 the plasma membrane binds sugars on the Jacob and Jessie lectins. B. Structures of
 22 representative lectins illustrated in A.

23

24 **Figure 12. Phylogenetic relationships of *E. histolytica* glutamine synthase.** The
 25 gene encoding glutamine synthase (EC 6.3.1.2) is now shared by *E. histolytica* and the
 26 diatom *Thalassiosira*. This gene is mainly restricted to prokaryotic genomes
 27 (eukaryotes are highlighted by arrows). *T. vaginalis* also contains a homologue but in
 28 this case it clusters weakly with *Fusobacterium*. The scale bar represents 10% of
 29 inferred sequence divergence. Both the GenBank and RefSeq accession numbers are
 30 given for the *E. histolytica* entry.

31

32 **Figure 13. Phylogenetic relationships of *E. histolytica* tryptophanase.** This tree
 33 suggests that the *E. histolytica* gene encoding a tryptophanase was acquired by LGT
 34 from a relative of the anaerobic bacterium *Fusobacterium*. In contrast, the *T.*

1 *vaginalis* gene appears to have a separate origin with a LGT from a relative of the
2 anaerobic *Bacteroides* group. The scale bar represents 10% of inferred sequence
3 divergence. Both the GenBank and RefSeq accession numbers are given for the *E.*
4 *histolytica* entry. The EC number is also shown.

5

6 **Figure 14. Pie chart of functional categories for the 68 strongest LGT cases.** The
7 cases are those discussed in the text and listed in Table 8. Most entries encode
8 metabolic enzymes (KEGG annotation).

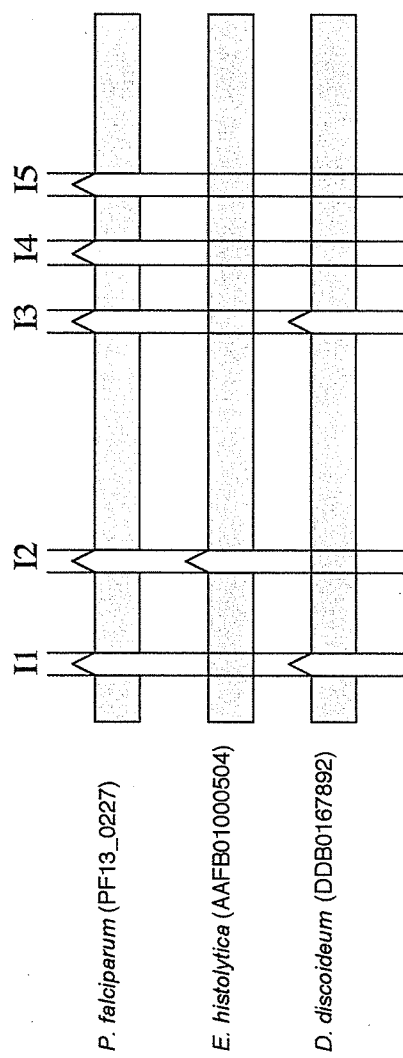


Fig. 1

1

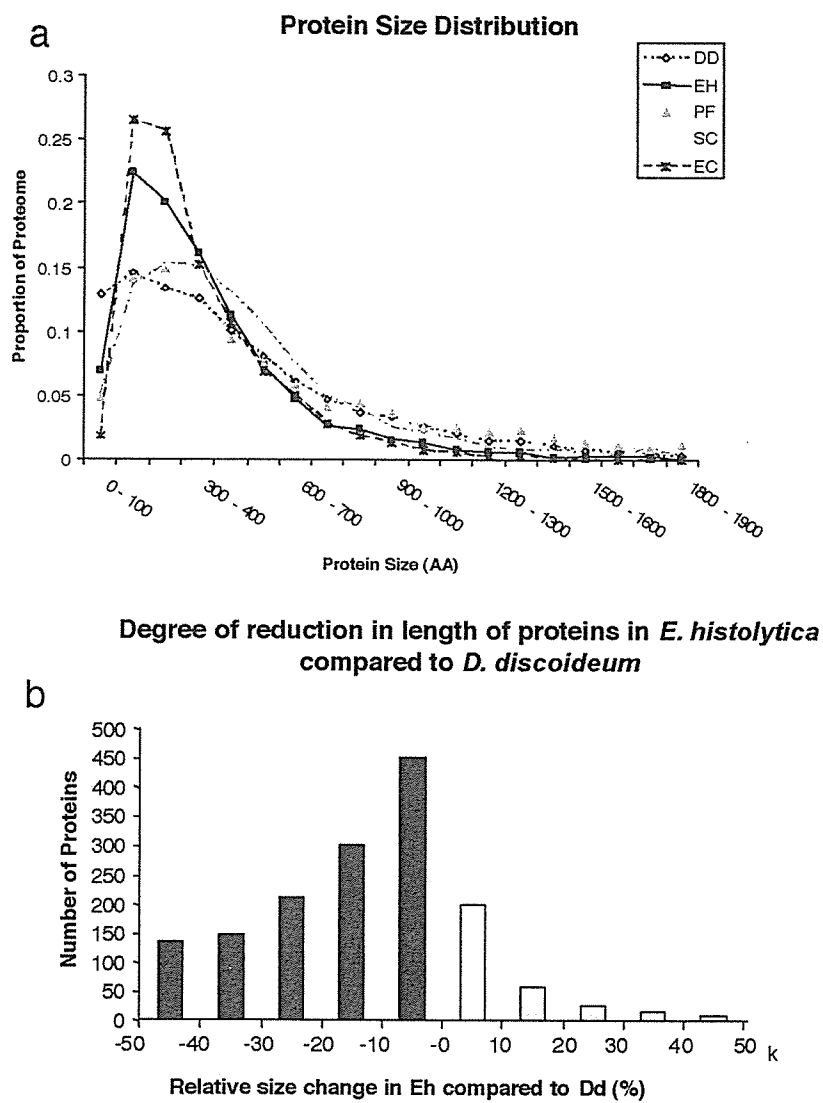


Fig. 2

1

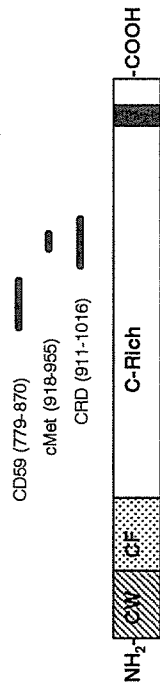


Fig. 3

2

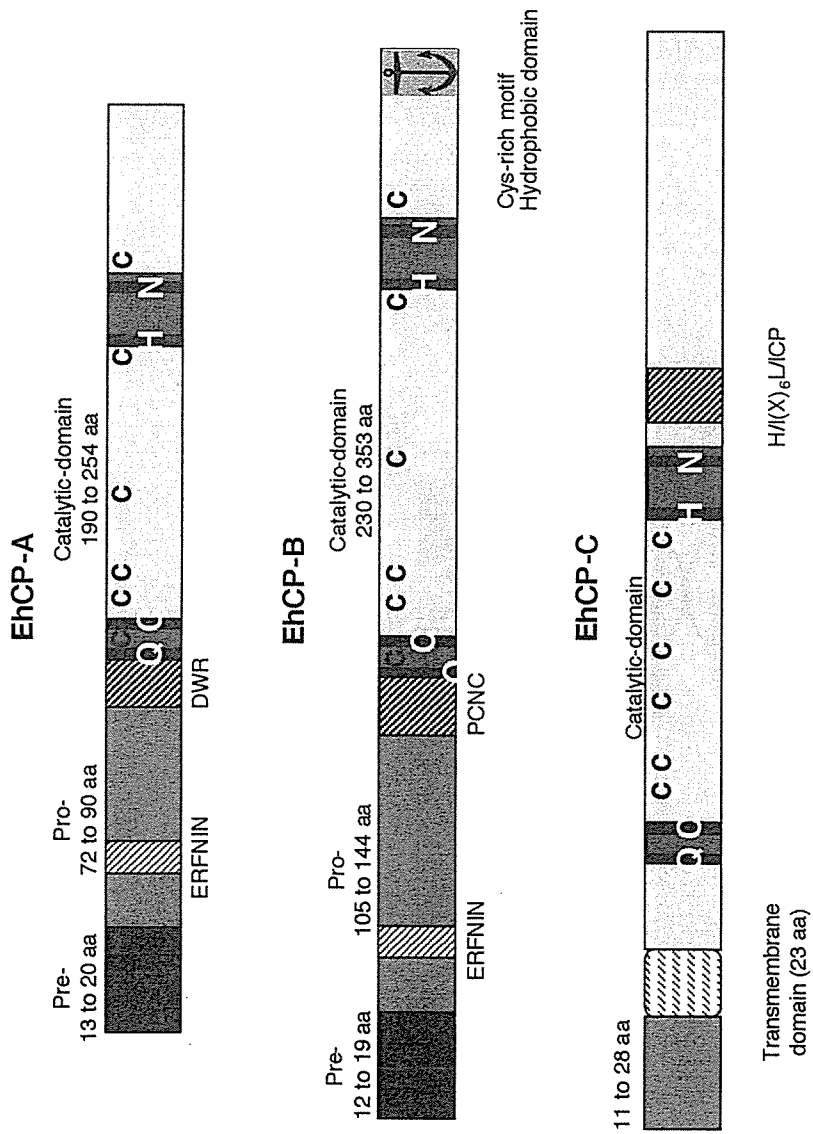


Fig. 4

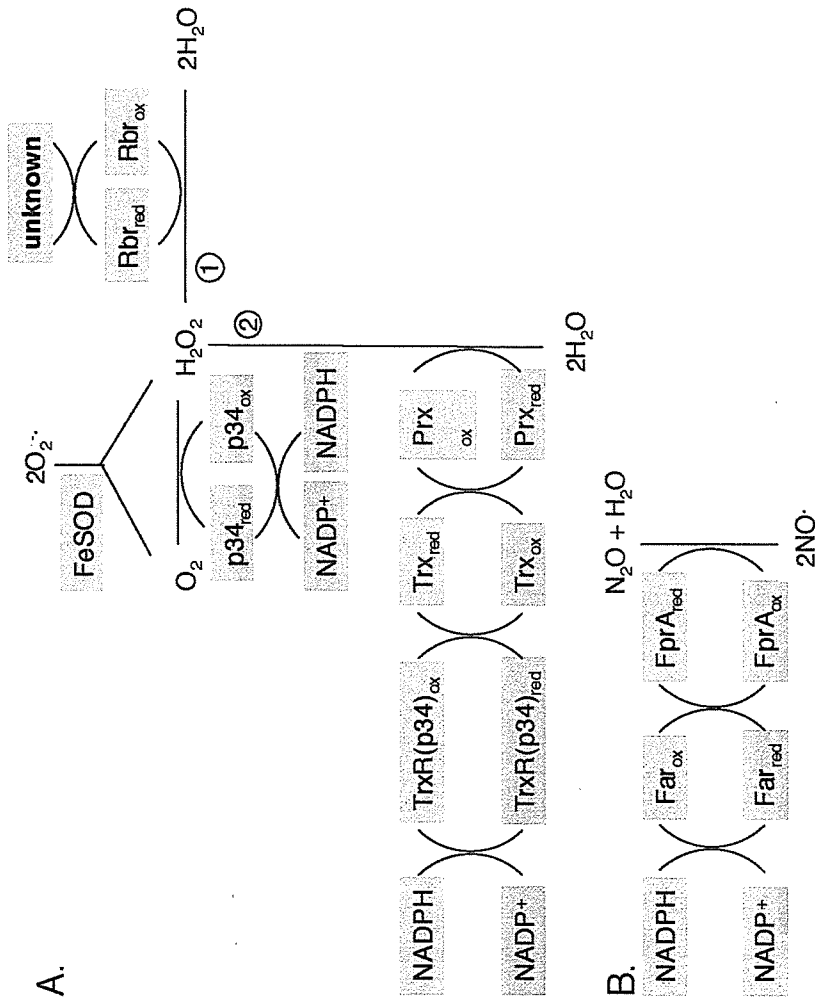


Fig. 5

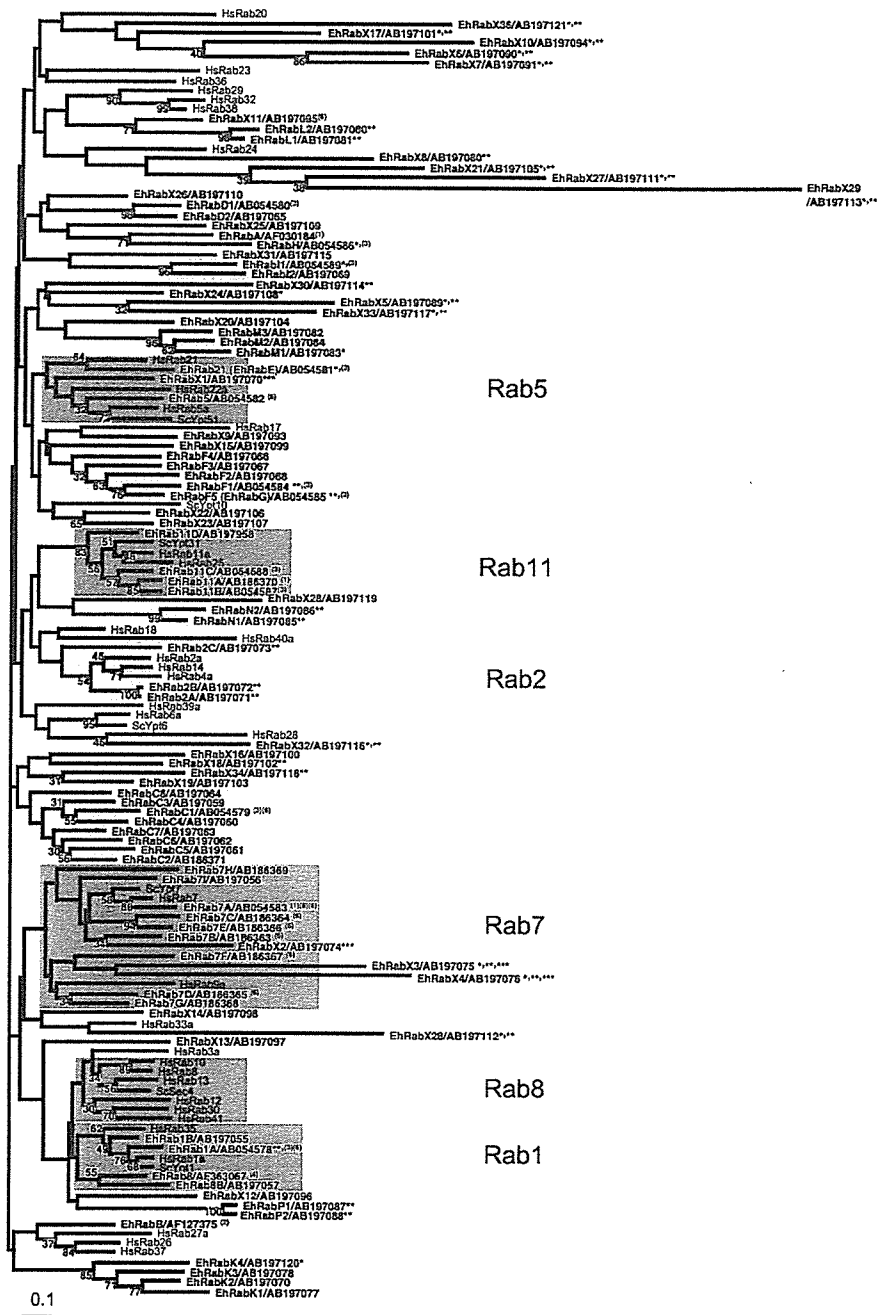
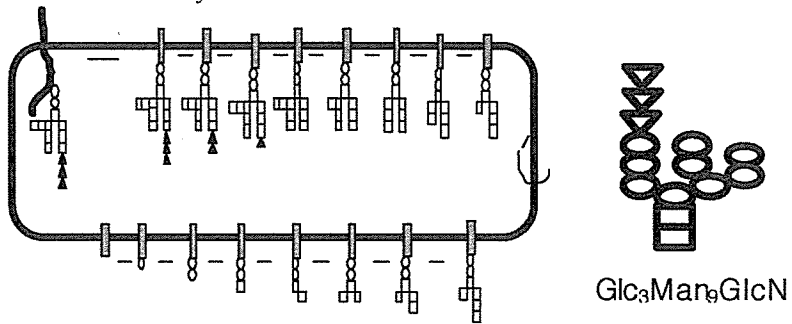
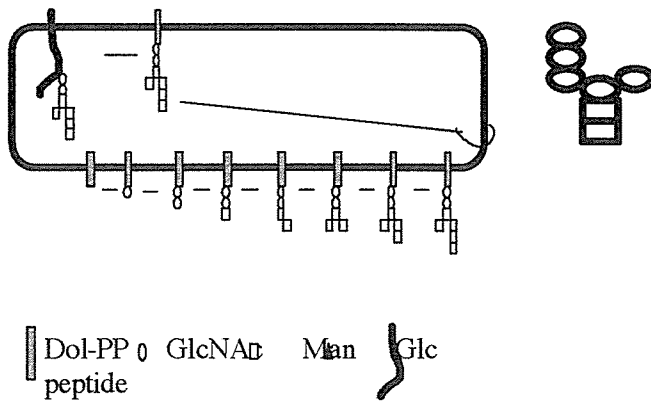


Fig. 6

1
2

A. *Saccharomyces*:B. *Entamoeba*:

1

2 Fig. 7

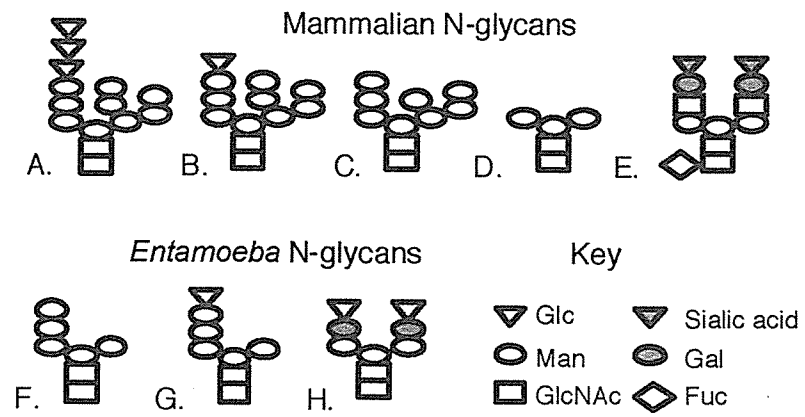


Fig. 8

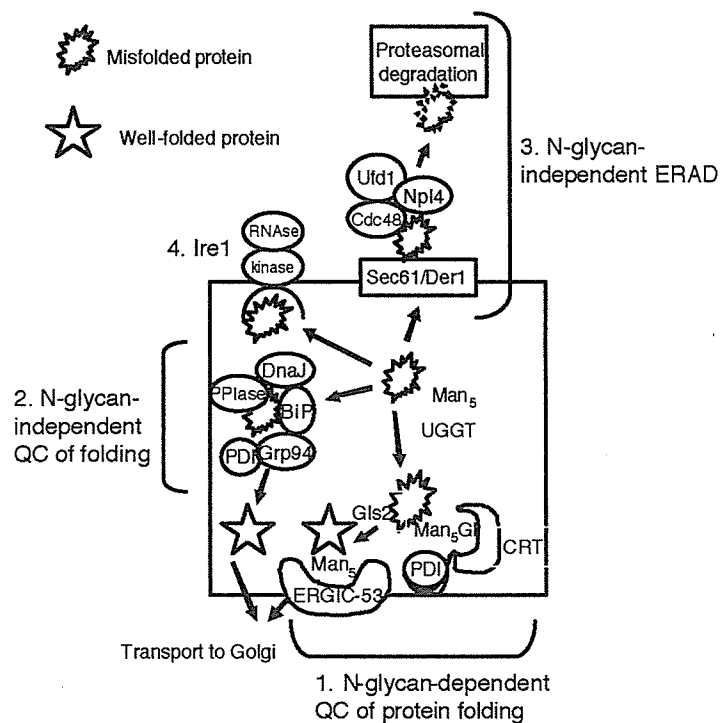


Fig. 9

1

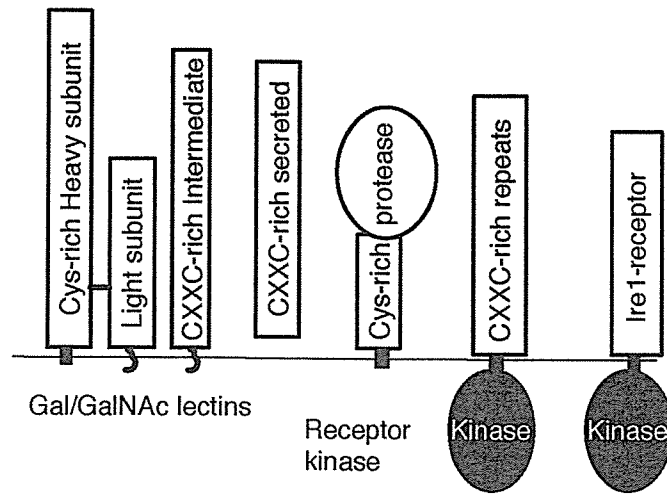
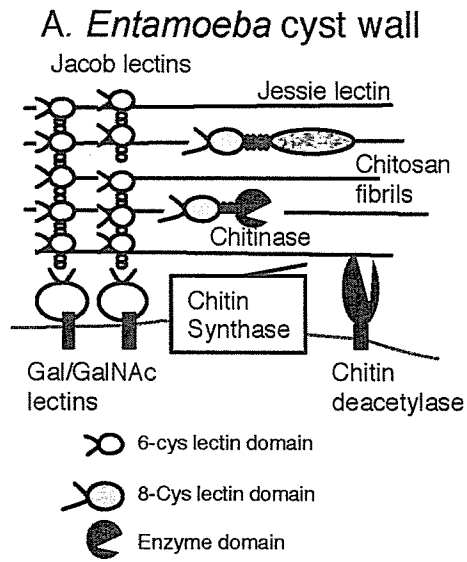


Fig. 10

1



B. *Entamoeba* cyst wall-associated lectins

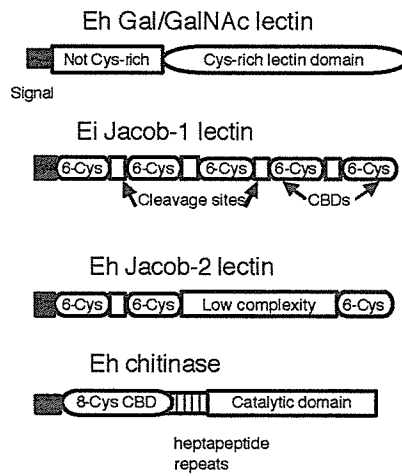


Fig. 11

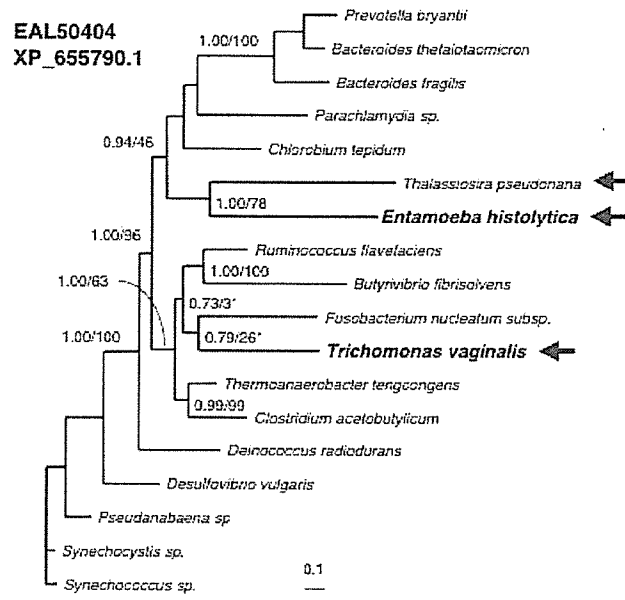


Fig. 12

1
2