## Phylogenomic Analyses to Detect Prokaryote to Eukaryote Horizontal Gene Transfer

T. Martin Embley

Newcastle University

## Overview

- The approach used - detection of 'recent' HGT based upon phylogenetic trees
- The computer programmes used and some of their features
- Some results from analysis of the genome of *Entamoeba histolytica*
- Future directions and challenges

## Why study HGT?

- What is the role of HGT, outside of endosymbiosis, in the evolution of eukaryotes?
- How much HGT has occurred? What types of genes? Where from?
- How has HGT affected the evolution of parasitic eukaryotes? Does HGT provide drug targets?

## How have we chosen to look for HGT?

- A semi-automated phylogenomics approach based upon a published computer programme called PyPhy
- Use trees to detect HGT in a two-stage process:
  - Simple parsimony or p-distance trees from edited alignments
  - Bayesian trees with more sophisticated phylogenetic models as a second stage screen
  - Distance-based bootstrapping to gain a measure of confidence in groups on trees

**A phylogenomic approach to microbial evolution**
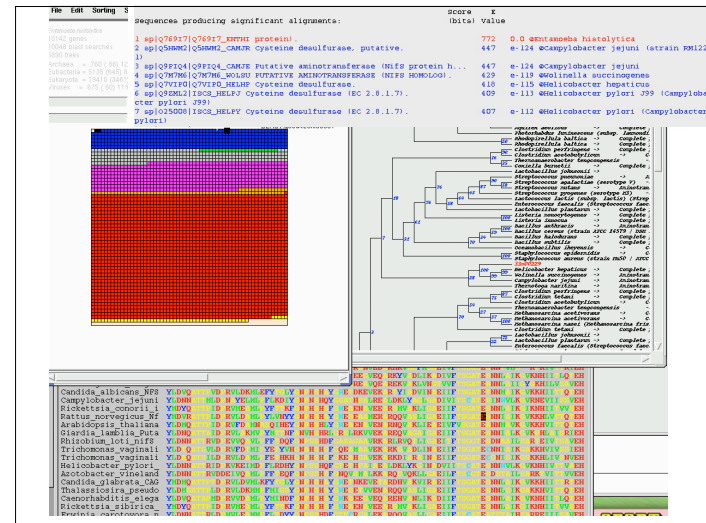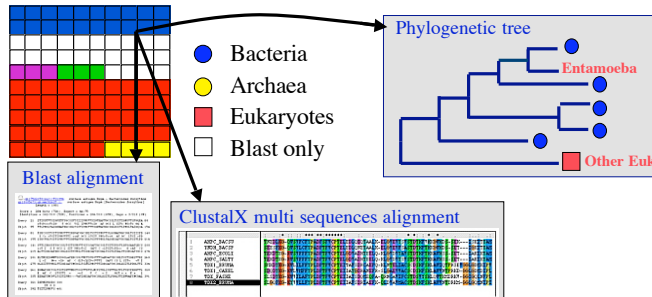
Thomas Sicheritz-Pontén and Siv G. E. Andersson*

Department of Molecular Evolution, Evolutionary Biology Center, Uppsala University, 752 36 Uppsala, Sweden
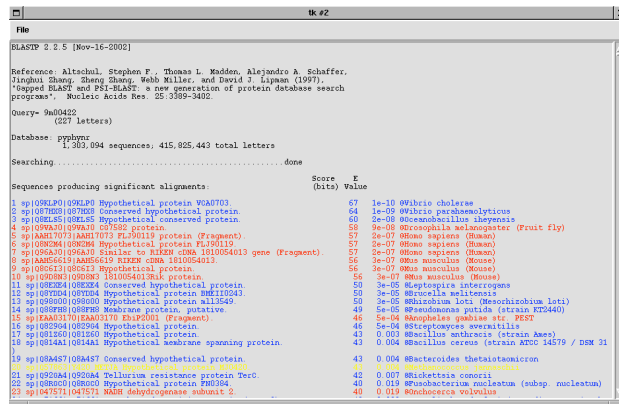
---

## PyPhy  - the Primary Screen

- **The database**
  - **Translated (Tr)EMBL**
  - **Swissprot entries (for annotation)**
- **BlastP searches using query genome to recover homologues**
  - **user defined cut-off values (e.g. at least 40% similarity over >70<200% of the length of query sequence)**
- **Align sequences using ClustalX and trim the alignment using GBLOCKS**
- **Construct bootstrapped P-distance or parsimony trees**
- **Display results as a grid with coloured squares indicating neighbor relationships**

---

## PyPhy graphical interface



Bacteria
Archaea
Eukaryotes
Blast only

Phylogenetic tree
Entamoeba
Other Euk

Blast alignment

ClustalX multi sequences alignment

---

## BlastP search



## G-blocked Alignment



## Bootstrapped P-distance trees



## Proteins where the top BlastP hit is a bacterium

The analysis of 100 genes supports the grouping of three highly divergent amoebae: *Dictyostelium*, *Entamoeba*, and *Mastigamoeba*

Eric Bapteste*, Henner Brinkmann†, Jennifer A. Lee‡, Dorothy V. Moore‡, Christoph W. Sensen§, Paul Gordon¶, Laure Duruflé*, Terry Gaasterland‡, Philippe Lopez*, Miklós Müller‡, and Hervé Philippe*∥

**(Entamoeba + Dictyostelium) connections**

## PyPhy - the secondary screen I

- Search EST databases and hits added to the alignments
- Feed G-blocked alignment into MRBAYES (Huelsenbeck & Ronquist 2001)
  - WAG matrix, gamma correction for site rate variation and a proportion (pinvar) of invariant sites.
- Make Bayesian consensus trees with posterior probabilities as support values

## PyPhy - the secondary screen II

- Bootstrapping to provide an additional indication of support for relationships
- Each data set is bootstrapped (100 replicates) and used to make distance matrices under the same evolutionary model as in the Bayesian analysis, using custom (P4) software
- Trees are made from the distance matrices using FastME (Desper & Gascuel 2004) and a bootstrap consensus tree made using P4
- Transfer BS support values onto the Bayesian consensus tree



tRNA (Uracil-5-) - methyltransferase
Top Prok blast/Top Euk blast = E-59

## The genome of the protist parasite *Entamoeba histolytica*

Brendan Loftus[1], Iain Anderson[1], Rob Davies[2], U. Cecilia M. Alsmark[3], John Samuelson[4], Paolo Amedeo[1], Paola Roncaglia[1], Matt Berriman[2], Robert P. Hirt[3], Barbara J. Mann[5], Tomo Nozaki[6], Bernard Suh[1], Mihai Pop[1], Michael Duchene[7], John Ackers[8], Egbert Tannich[9], Matthias Leippe[10], Margit Hofer[7], Iris Bruchhaus[9], Ute Willhoeft[9], Alok Bhattacharya[11], Tracey Chillingworth[2], Carol Churcher[2], Zahra Hance[2], Barbara Harris[2], David Harris[2], Kay Jagels[2], Sharon Moule[2], Karen Mungall[2], Doug Ormond[2], Rob Squares[2], Sally Whitehead[2], Michael A. Quail[2], Ester Rabbinowitsch[2], Halina Norbertczak[2], Claire Price[2], Zheng Wang[1], Nancy Guillén[12], Carol Gilchrist[5], Suzanne E. Stroup[5], Sudha Bhattacharya[11], Anuradha Lohia[13], Peter G. Foster[14], Thomas Sicheritz-Ponten[15], Christian Weber[12], Upinder Singh[16], Chandrama Mukherjee[13], Najib M. El-Sayed[1], William A. Petri Jr[5], C. Graham Clark[8], T. Martin Embley[3], Bart Barrell[2], Claire M. Fraser[1] & Neil Hall[2]*

## E. histolytica trophozoites in situ
### (D. Mirelman 1996)



D. Mirelman © 1996



Life cycle of *Entamoeba histolytica* and the clinical manifestations of infection in humans

Expert Reviews in Molecular Medicine ©1999 Cambridge University Press

## The *E. histolytica* genome

- Sequenced by TIGR and the Sanger institute (PIs Brendan Loftus and Neil Hall)
- ~ 23.7 Mbp
- 9938 predicted genes
- 1/3rd of genes have no detectable homologue in public data bases
- Large gene families, duplications are common

## Aim of our HGT screen

- To identify the strongest candidate 'recent' HGT from prokaryotes to *Entamoeba*
    - The tip of the HGT iceberg?
    - Help to validate HGT as a plausible explanation for topological incongruence deeper in tree
- These cases shoudl be the easiest to defend - paralogy or poor tree-building is unlikely to be an equally plausible explanation to HGT
- This screen is not seeking to detect more 'ancient' transfers, for example - at the base of major groups or from the alpha-proteobacterial mitochondrial endosymbiont

## Primary screen output for *E. histolytica*

- 9938 genes

- 5740 trees

- 819 bacterial and 129 archaeal connections in the trees



## Superimposed top blast hit = prokaryote

- 912 genes have a bacteria and 204 an archaea as the top BLAST hit, 107 of these genes show eukaryotic connections in the primary tree.



## Criteria for HGT 2° analysis

### (vary in degree of ad hoc-ness)

- Strongly supported tree placing *Entamoeba* within a conventional bacterial group
- No other eukaryotic sequence in tree
  - sampling is very poor, surprises occur!
- Two or more strongly supported nodes separating *Entamoeba* from any other eukaryotic sequence
- Eukaryotes never appear together in bootstrap partitions
- High ratio of prokaryote to eukaryote BlastP score

## Results of the 2nd generation analysis

- 96 Bayesian trees are most simply explained by prokaryote to eukaryote HGT

Zinc alcohol dehydrogenase (subtree)



Mapping HGT genes onto scaffolds



Mannose phosphate guanyltransferase
Top Prok Blast/Top Euk = E-88



Thiamine phosphokinase
Blast ratio Top Prok/Top Euk = E-25

## Top-left panel

**tRNA (Uracil-5-) - methyltransferase**
Top Prok blast/Top Euk blast = E-59



Tree tips:
- *Clostridium_tetani_CTC01941*
- 209m00106
- *Thermoanaerobacter_tengcongensis_TRMA_OR_TTE1797*
- *Arabidopsis_thaliana_AT3G21300 (top euk blast)*
- *Arabidopsis_thaliana_*
- *Chlorobium_tepidum_CT0009*
- *Bacteroides_thetaiotaomicron_BT0643*
- *Thermotoga_maritima_TM1094*
- *Leptospira_interrogans_LA0292*
- *Staphylococcus_epidermidis_SE1582*
- *Listeria_monocytogenes_LMO1703*
- *Listeria_innocua_LIN1815*
- *Enterococcus_faecalis_EF0728*
- *Enterococcus_faecalis_*
- *Lactobacillus_plantarum_LP_1151*
- *Streptococcus_pneumoniae_SPR0932*
- *Streptococcus_pneumoniae_SP1029*
- *Oceanobacillus_iheyensis_OB0768*
- *Bacillus_cereus_BC0364*
- *Bacillus_anthracis_BA0333*
- *Thermoanaerobacter_tengcongensis_TRMA2_OR_TTE1812*
- *Clostridium_perfringens_CPE0367*
- *Clostridium_tetani_CTC02481*
- *Clostridium_acetobutylicum_CAC0523*
- *Clostridium_perfringens_CPE2114*
- *Clostridium_acetobutylicum_CAC1435*

## Top-right panel

### Transferred genes are involved in metabolism

| | | |
|---|---|---|
| 121m00124 | Tartrate dehydrogenase | METABOLISM |
| 133m00129 | Fructokinase | METABOLISM |
| 400m00033 | Arginine decarboxylase | METABOLISM |
| 78m00151 | Lysophospholipase L2 | METABOLISM |
| 9m00390 | HesA/MoeB/ThiF family protein | METABOLISM |
| 24m00307 | NADH dehydrogenase, similar to nitrite reductase | ENERGY METABOLISM |
| 61m00186 | 5-nitroimidazole antibiotic resistance protein | CELL RESCUE, DEFENCE AND VIRULENCE |
| 30m00272 | ABC transporter | TRANSPORT |
| 289m00068 | Conserved hypothetical protein | UNKNOWN FUNCTION |

## Bottom-left panel

## Functional distribution of HGT



Legend:
- Metabolism amino acids — 18%
- Metabolism carbohydrate — 18%
- Metabolism cofactors — 7%
- Metabolism energy — 7%
- Metabolism nucleotides — 4%
- Metabolism lipids — 1%
- Environmental information processing — 3%
- Unclassified — 21%
- Hypothetical — 21%

## Bottom-right panel



9

## We see a broad collection of donors

| Name | NN MrBayes | Best Prokaryote blast hit |
|---|---|---|
| 100m00122 | DT Deinococcus | Clostridium perfringens |
| 189m00102 | Archeal like | Pyrobaculum aerophilum |
| 10m00331 | Proteo /Xanthobact | Methanobacterium thermoautotrophicum |
| 119m00142 | Proteo /Vibrio | Bacillus subtilis |
| 11m00315 | Bact/Chlorob | Bacteroides thetaiotaomicron |
| 209m00106 | Firmicutes | Clostridium perfringen |
| 126m00107 | Bacteria like | Pseudomonas aeruginosa |
| 130m00124 | Firmicutes | Lactococcus lactis (subsp. lactis) (Streptococcus lactis) |
| 133m00129 | Bacteria like | Bacteroides thetaiotaomicron |
| 133m00136 | Bact/Chlorob | Methanosarcina mazei (Methanosarcina frisia) |

## Summary BlastP Statistics
### (provided as supplemental info to Loftus et al 2005)

Summary BlastP statistics for 96 candidate LGT

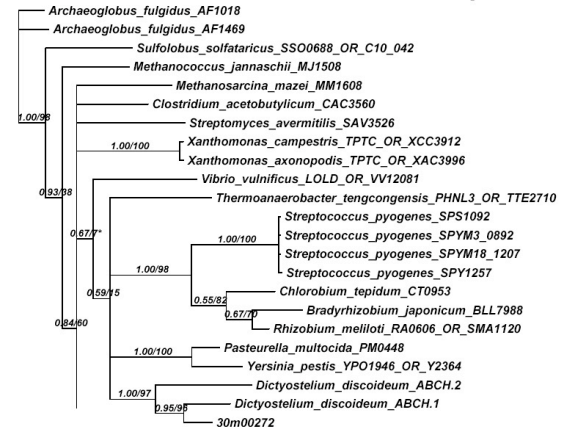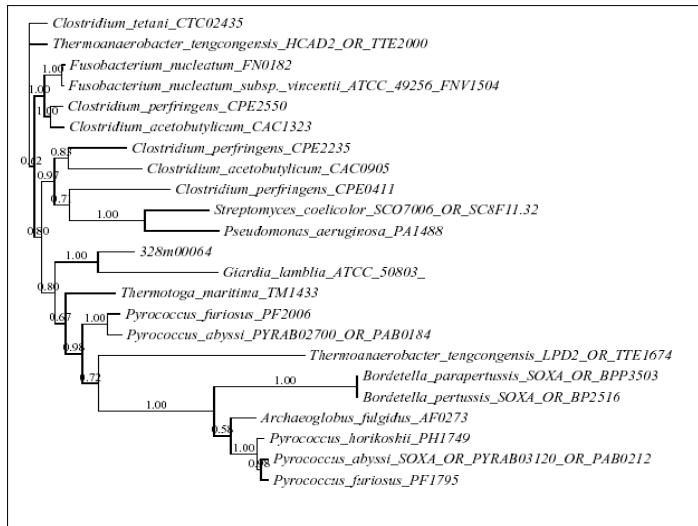| Acc. No. | Length | Top Prokaryotic blast hit | PL | ID level | %ID | Top Eukaryotic hit | EL | ID level | %ID | P E-score | E E-score | Ratio (P/E) |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| EAL47525 | 380 | Clostridium perfringens | 284 | 77/324 | 23 | none | | none | | 1.00E-13 | none | na |
| EAL47464 | 504 | Bacteroides thetaiotaomicron | 482 | 176/495 | 38 | Piromyces sp. E2 | 555 | 110/400 | 27 | 1.00E-86 | 9.00E-31 | 1.11E-56 |
| EAL51236 | 259 | Methanobacterium thermoautotrophicum | 258 | 101/256 | 39 | Arabidopsis thaliana | 286 | 60/276 | 21 | 1.00E-47 | 0.011 | 9.09E-46 |
| EAL47026 | 164 | Bacillus subtilis | 171 | 42/138 | 30 | Homo sapiens | 170 | 20/67 | 29 | 1.00E-10 | 0.002 | 5.00E-08 |
| EAL51149 | 343 | Bacteroides thetaiotaomicron | 350 | 147/344 | 42 | Phytomonas sp. | 361 | 115/352 | 32 | 4.00E-83 | 2.00E-45 | 2.00E-38 |
| EAL46975 | 370 | Bordetella pertussis | 355 | 160/346 | 46 | Anopheles gambiae | 351 | 124/360 | 34 | 5.00E-83 | 8.00E-49 | 6.25E-35 |
| EAL46858 | 192 | Pseudomonas aeruginosa | 184 | 68/162 | 41 | Caenorhabditis elegans | 212 | 54/156 | 34 | 4.00E-36 | 3.00E-17 | 1.33E-19 |
| EAL46757 | 95 | Lactococcus lactis | 94 | 28/90 | 31 | Mus musculus | 352 | 17/74 | 22 | 1.00E-09 | 1.4 | 7.14E-10 |
| EAL46701 | 294 | Bacteroides thetaiotaomicron | 295 | 130/282 | 46 | Arabidopsis thaliana | 325 | 76/286 | 26 | 3.00E-62 | 8.00E-14 | 3.75E-49 |
| EAL46679 | 218 | Methanosarcina mazei | 219 | 73/196 | 37 | Oryza sativa | 4572 | 28/122 | 22 | 5.00E-31 | 3.7 | 1.35E-31 |
| EAL46656 | 419 | Dictyoglomus thermophilum | 562 | 122/403 | 30 | Schizosaccharomyces pombe | 478 | 46/111 | 41 | 2.00E-35 | 1.00E-19 | 2.00E-16 |
| EAL50986 | 219 | Bacteroides thetaiotaomicron | 229 | 59/186 | 31 | Anopheles gambiae | 312 | 49/166 | 29 | 1.00E-20 | 4.00E-12 | 2.50E-09 |
| EAL50992 | 140 | Archaeoglobus fulgidus | 174 | 68/170 | 40 | Giardia lamblia | 264 | 24/80 | 30 | 5.00E-28 | 3.4 | 1.47E-28 |
| EAL50997 | 656 | Bacteroides thetaiotaomicron | 699 | 386/685 | 53 | Schizosaccharomyces pombe | 683 | 155/500 | 31 | 0 | 3.00E-67 | 0.00E+00 |
| EAL46421 | 205 | Clostridium acetobutylicum | 219 | 84/209 | 40 | Arabidopsis thaliana | 230 | 52/157 | 33 | 5.00E-34 | 4.00E-12 | 1.25E-22 |
| EAL46399 | 218 | Clostridium perfringens | 224 | 143/217 | 65 | Giardia lamblia | 258 | 101/217 | 46 | 2.00E-73 | 2.00E-43 | 1.00E-30 |
| EAL46311 | 248 | Synechococcus elongatus | 258 | 94/257 | 36 | Plasmodium yoelii | 407 | 33/177 | 24 | 9.00E-31 | 2.1 | 4.29E-31 |
| EAL46313 | 118 | Prochlorococcus marinus | 153 | 47/111 | 42 | Hordeum vulgare | 212 | 18/79 | 22 | 1.00E-21 | 0.58 | 1.72E-21 |
| EAL46110 | 157 | Bacteroides thetaiotaomicron | 155 | 75/153 | 49 | Anopheles gambiae | 446 | 27/139 | 19 | 1.00E-34 | 0.59 | 1.69E-34 |
| EAL46116 | 661 | Bacillus halodurans | 666 | 323/662 | 48 | Craterostigma plantagineum | 679 | 331/659 | 50 | 0 | 1.00E-172 | 0.00E+00 |
| EAL46026 | 176 | Bacteroides thetaiotaomicron | 174 | 85/170 | 50 | Arabidopsis thaliana | 147 | 30/64 | 46 | 1.00E-40 | 3.00E-07 | 3.33E-34 |
| EAL50801 | 499 | Bacteroides thetaiotaomicron | 497 | 257/494 | 52 | Piromyces sp. E2 | 494 | 118/414 | 28 | 1.00E-145 | 8.00E-37 | 1.25E-109 |
| EAL50838 | 299 | Anabaena sp. (strain PCC 7120) | 275 | 77/280 | 27 | Caenorhabditis elegans | 446 | 40/203 | 19 | 2.00E-15 | 0.002 | 1.00E-12 |
| EAL45907 | 300 | Streptomyces coelicolor | 586 | 102/318 | 32 | Dictyostelium discoideum | 442 | 110/365 | 30 | 1.00E-39 | 1.00E-35 | 0.0001 |
| EAL45618 | 159 | Bacteroides thetaiotaomicron | 147 | 67/144 | 46 | none | | none | | 2.00E-28 | none | na |
| EAL45586 | 460 | Clostridium tetani | 461 | 213/450 | 47 | Dictyostelium discoideum | 503 | 181/458 | 39 | 1.00E-116 | 3.00E-64 | 3.33E-53 |
| EAL45595 | 284 | Pyrobaculum aerophilum | 281 | 78/286 | 27 | Arabidopsis thaliana | 291 | 63/230 | 27 | 8.00E-24 | 0.011 | 8.00E-10 |

Phylogenetic tree (lower left panel):

```
Synechococcus_sp._CYSE
Synechocystis_sp._CYSE_OR_SLR1348
  Helicobacter_pylori_J99_CYSE_OR_JHP1133
  Helicobacter_pylori_CYSE_OR_HP1210
  Bacillus_subtilis_CYSE_OR_CYSA_OR_BSU00930
  Staphylococcus_xylosus_CYSE
  Pseudomonas_syringae_PSPTO5178
  Rhodopirellula_baltica_RB5098
  Bacteroides_thetaiotaomicron_BT3256
  E_moshkovski
  200m00078
  Entamoeba_histolytica_EHSAT1
  Entamoeba_dispar_EDSAT1
  Entamoeba_dispar_EDSAT2
  Nicotiana_tabacum_SAT1
  Arabidopsis_thaliana_
  Buchnera_aphidicola_CYSE_OR_BU054
  Buchnera_aphidicola_CYSE_OR_BUSG051
  Pasteurella_multocida_CYSE_OR_PM1430
  Haemophilus_influenzae_CYSE_OR_HI0606
  Yersinia_pestis_CYSE_OR_YPO0070
  Yersinia_pestis_CYSE_OR_Y0072
  Escherichia_coli
  Salmonella_typhimurium_CYSE_OR_STM3699
```

## HGT at the base of the amoebozoa?
### ABC transporter

Phylogenetic tree (lower right panel):

```
Archaeoglobus_fulgidus_AF1018
Archaeoglobus_fulgidus_AF1469
Sulfolobus_solfataricus_SSO0688_OR_C10_042
Methanococcus_jannaschii_MJ1508
Methanosarcina_mazei_MM1608
Clostridium_acetobutylicum_CAC3560
Streptomyces_avermitilis_SAV3526
Xanthomonas_campestris_TPTC_OR_XCC3912
Xanthomonas_axonopodis_TPTC_OR_XAC3996
Vibrio_vulnificus_LOLD_OR_VV12081
Thermoanaerobacter_tengcongensis_PHNL3_OR_TTE2710
Streptococcus_pyogenes_SPS1092
Streptococcus_pyogenes_SPYM3_0892
Streptococcus_pyogenes_SPYM18_1207
Streptococcus_pyogenes_SPY1257
Chlorobium_tepidum_CT0953
Bradyrhizobium_japonicum_BLL7988
Rhizobium_meliloti_RA0606_OR_SMA1120
Pasteurella_multocida_PM0448
Yersinia_pestis_YPO1946_OR_Y2364
Dictyostelium_discoideum_ABCH.2
Dictyostelium_discoideum_ABCH.1
30m00272
```

## Panel 1 (top left — phylogenetic tree)



```
Clostridium_tetani_CTC02435
Thermoanaerobacter_tengcongensis_HCAD2_OR_TTE2000
Fusobacterium_nucleatum_FN0182
1.00
1.00  Fusobacterium_nucleatum_subsp._vincentii_ATCC_49256_FNV1504
1.00  Clostridium_perfringens_CPE2550
      Clostridium_acetobutylicum_CAC1323
0.88  Clostridium_perfringens_CPE2235
0.9       Clostridium_acetobutylicum_CAC0905
          Clostridium_perfringens_CPE0411
0.7   1.00  Streptomyces_coelicolor_SCO7006_OR_SC8F11.32
0.80        Pseudomonas_aeruginosa_PA1488
1.00  328m00064
0.80  Giardia_lamblia_ATCC_50803_
      Thermotoga_maritima_TM1433
0.67  1.00  Pyrococcus_furiosus_PF2006
              Pyrococcus_abyssi_PYRAB02700_OR_PAB0184
0.98  Thermoanaerobacter_tengcongensis_LPD2_OR_TTE1674
0.72  1.00  Bordetella_parapertussis_SOXA_OR_BPP3503
              Bordetella_pertussis_SOXA_OR_BP2516
1.00  Archaeoglobus_fulgidus_AF0273
0.58  Pyrococcus_horikoshii_PH1749
1.00  Pyrococcus_abyssi_SOXA_OR_PYRAB03120_OR_PAB0212
      Pyrococcus_furiosus_PF1795
```

## Panel 2 (top right — phylogenetic tree)

**Zinc alcohol dehydrogenase (subtree)**



```
Clostridium_beijerinckii_ADH
Thermoanaerobacter_tengcongensis_TDH_OR_TTE0695
Thermoanaerobacter_brockii_ADH
Thermoanaerobacter_ethanolicus_ADHB
Entamoeba_histolytica_ADH1
Trichomonas_vaginalis_
Mycoplasma_penetrans_MYPE4620
Rhizobium_meliloti_RA0626_OR_SMA1156
Phytomonas_sp._ADU_2003_IPDH
Chromobacterium_violaceum_CV2051
Alcaligenes_eutrophus_ADH
Pseudomonas_aeruginosa_PA2119
11m00315
Bacteroides_thetaiotaomicron_BT4512
Bacillus_subtilis_YJMD
Tropheryma_whipplei_ADH_OR_TWT326
Tropheryma_whipplei_ADH_OR_TW445
Streptomyces_coelicolor_SCO0179_OR_SCJ1.28C
Streptomyces_avermitilis_SAV1272
Ralstonia_solanacearum_RSC0194_OR_RS00626
Neisseria_meningitidis_NMB1395
Staphylococcus_epidermidis_SE2098
Streptococcus_pneumoniae_ADH_OR_SPR1866
Streptococcus_pneumoniae_SP2055
Streptococcus_pyogenes_SPYM18_1073
Streptococcus_pyogenes_SPYM3_0772_OR_SPS0972
Streptococcus_pyogenes_SPY1111
Streptococcus_agalactiae_ADH_OR_SAG1637
Streptococcus_agalactiae_GBS1684
```
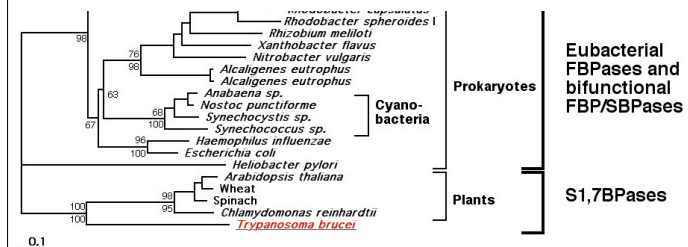
## Panel 3 (bottom left — slide)

# Plant-like traits associated with metabolism of *Trypanosoma* parasites
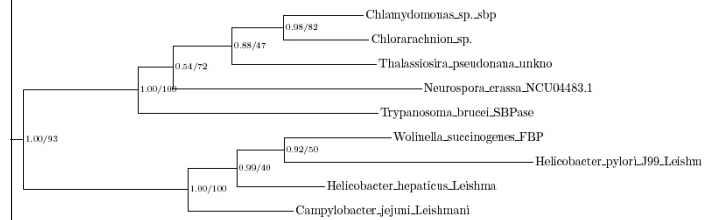
Véronique Hannaert*[†], Emma Saavedra*[†], Francis Duffieux*[‡], Jean-Pierre Szikora*, Daniel J. Rigden[§], Paul A. M. Michels*, and Fred R. Opperdoes*[¶]

- *Trypanosoma brucei* contains a complete open-reading frame encoding a homologue of sedoheptulose-1,7-bisphosphatase (SBPase).
- This enzyme is typical for the Calvin cycle of photosynthetic organisms and only encountered in the chloroplasts of green algae and plants
- This was taken as evidence by that *Trypanosoma* once had a plastid

## Panel 4 (bottom right — slide)

http://www.icp.ucl.ac.be/~opperd/Supplementary/sbpase_nbj_tree.html



```
                    Rhodobacter capsulatus
                    Rhodobacter spheroides I
                    Rhizobium meliloti
98      76          Xanthobacter flavus
        98          Nitrobacter vulgaris
              Alcaligenes eutrophus
        63    Alcaligenes eutrophus
              Anabaena sp.
        68    Nostoc punctiforme
        100   Synechocystis sp.
67            Synechococcus sp.
        96    Haemophilus influenzae
        100   Escherichia coli
                    Heliobacter pylori
                    Arabidopsis thaliana
              98    Wheat
100           95    Spinach
100                 Chlamydomonas reinhardtii
                    Trypanosoma brucei
0.1
```

Cyano-bacteria

Prokaryotes — Eubacterial FBPases and bifunctional FBP/SBPases

Plants — S1,7BPases

11

## Fungi contain SBPase too!



```
                                    Chlamydomonas_sp._sbp
                        0.98/82
                                    Chlorarachnion_sp.
                 0.88/47
                                    Thalassiosira_pseudonana_unkno
         0.54/72
                                    Neurospora_crassa_NCU04483.1
         1.00/100
                                    Trypanosoma_brucei_SBPase
                                    Wolinella_succinogenes_FBP
                        0.92/50
                                    Helicobacter_pylori_J99_Leishm
                 0.99/40
                                    Helicobacter_hepaticus_Leishma
         1.00/100
                                    Campylobacter_jejuni_Leishmani
1.00/93
```

## The people who do the work



Robert Hirt and
Thomas Sicheritz
Ponten

Cessie Alsmark

Peter Foster

Collaborators:
Brendan Loftus and Neil Hall (TIGR)
Matt Berriman (Sanger)

## Useful URLs

- **Pyphy** http://www.cbs.dtu.dk/staff/thomas/pyphy/
- **P4** http://www.nhm.ac.uk/zoology/external/p4.htm
- **ClustalX** http://www.hgmp.mrc.ac.uk/Registered/Option/clustalx.html
- **Gblocks** http://molevol.ibmb.csic.es/Gblocks/Gblocks.html
- **MrBayes** http://morphbank.ebc.uu.se/mrbayes/
- **FastME** http://www.lirmm.fr/~w3ifa/MAAS/FastME/FastME.html