

What Does the Microsporidian *E. cuculii* Tell Us About the Origin of the Eukaryotic Cell?

Alexei Fedorov,¹ Hyman Hartman²

¹ Department of Medicine, Medical College of Ohio, Toledo OH 43614, USA

² Center for Biomedical Engineering NE47-377, MIT, Cambridge, MA 02139, USA

Received: 19 August 2003 / Accepted: 1 July 2004 [Reviewing Editor: Dr. Manyuan Long]

Abstract. The relationship among the three cellular domains Archaea, Bacteria, and Eukarya has become a central problem in unraveling the tree of life. This relationship can now be studied as the completely sequenced genomes of representatives of these cellular domains become available. We performed a bioinformatic investigation of the *Encephalitozoon cuculii* proteome. *E. cuculii* has the smallest sequenced eukaryotic genome, 2.9 megabases coding for 1997 proteins. The proteins of *E. cuculii* were compared with a previously characterized set of eukaryotic signature proteins (ESPs). ESPs are found in a eukaryotic cell, whether from an animal, a plant, a fungus, or a protozoan, but are not found in the Archaea and the Bacteria. We demonstrated that 85% of the ESPs have significant sequence similarity to proteins in *E. cuculii*. Hence, *E. cuculii*, a minimal eukaryotic cell that has removed all inessential proteins, still preserves most of the ESPs that make it a member of the Eukarya. The locations and functions of these ESPs point to the earliest history of eukaryotes.

Key words: Eukaryote — *Giardia* — *Encephalitozoon cuculii* — Microsporidia — Minimal cell — Nucleus — Endosymbiosis

Introduction

The microsporidia were once considered to be the deepest-branching eukaryote taxon. Carl Woese and his team found the microsporidia to be very close to the root of the Eukarya based on their ribosomal RNA phylogeny (Vossbrinck et al. 1987). However, this deep divergence was due to a long-branch attraction artifact resulting from rapid evolution of the microsporidia (Gribaldo and Phillippe 2002). The recent analyses of combined protein data resulted in the identification of the microsporidian *Encephalitozoon cuculii* as a member of the fungal kingdom (Baldauf et al. 2000). In fact, *E. cuculii* is an intracellular fungal parasite. This lifestyle has caused the shrinking of its genome. *E. cuculii* is thus a candidate for the “minimal” eukaryotic cell. It has the added virtue to have had its genome fully sequenced.

Let us examine the concept of a minimal cell. The search for the minimal cell began in the cellular domain of the Bacteria with the investigation of the *Mycoplasma*, a group of small parasitic bacteria (Razin 1997). *Mycoplasma genitalium*, an intracellular parasite, has one of the smallest known bacterial genomes. Because of its small size, it was among the first bacterial genomes to be sequenced. The genome of *M. genitalium* has only 580,000 base pairs (bp), which codes for 468 proteins (Fraser et al. 1995). The characterization of the minimal cell began by comparing the proteins found in *M. genitalium* to those found in *Haemophilus influenzae*, an early sequenced gram-negative bacteria causing ear infections (Mushegian and Koonin 1996). Computer analysis

produced a set of 240 proteins of *M. genitalium* that have orthologs in *H. influenzae*. This set was expanded by the addition of 16 other “gene displacement” proteins. Thus the genome of a minimal cell would code for the 256 proteins that would carry out the bare essential cellular functions such as DNA replication, protein synthesis (translation and transcription), metabolism (glycolysis), and various membrane related functions. It was assumed that the minimal bacterial cell would be a minimal cell for the other cellular domains. However, when Mushegian and Koonin compared these 256 proteins with those from other domains, a lack of correspondence unexpectedly appeared. For example, they found that different sets of proteins were involved in Eukarya and Archaea DNA replication. The authors concluded that the last common ancestor of the three cellular domains “had an RNA genome” and that DNA replication had evolved twice: once in Bacteria and independently in Archaea and Eukarya. (Mushegian and Koonin 1996). Hence, a minimal eukaryotic cell or a minimal archaeal cell would differ from a minimal bacterial cell.

To characterize a “minimal” eukaryotic cell we began our study with *E. cucuruli*. The genome of microsporidian *E. cucuruli* has been sequenced (Katinka et al. 2001). Its length is 2.9 million bp, close to the median size of sequenced prokaryotic genomes (2.6 million bp). *E. cucuruli* has 1997 predicted genes, 44% of which have known functions (Katinka et al. 2001). These data imply that *E. cucuruli* is approaching the size of a minimal eukaryotic cell. To understand the differences between minimal cells of Eukarya and Bacteria, we need to characterize the set of proteins in *E. cucuruli* that is unique to the eukaryotes but absent from other cellular domains. Previously, we collected a set of proteins that were found in all sequenced eukaryotic cells and absent in Bacteria or Archaea (Hartman and Fedorov 2002). We called this set eukaryotic signature proteins (ESPs). Each protein from the ESP set has homologs in all main eukaryotic branches—animals (*Drosophila melanogaster* and *Caenorhabditis elegans*), plants (*Arabidopsis thaliana*), fungi (*Saccharomyces cerevisiae*), and protists (*Giardia lamblia*)—but does not have homologs in the Archaea and Bacteria. The extracellular eukaryotic parasite *G. lamblia* was specifically chosen for the characterization of these ESPs, because it is still considered to be one of the deepest-branching taxon of the eukaryotes (i.e., the diplomonads). The numbers and composition of the proteins from the ESP set of *Giardia* led us to conclude that the origin of eukaryotes involved the formation of the nucleus from prokaryotic endosymbionts in an RNA-based cell (Hartman and Fedorov 2002). Here we compare the ESPs with proteins of a minimal and highly diverged eukaryotic

cell of *E. cucuruli*. The surprising result is the overwhelming agreement (85%) between the eukaryote-specific proteins of *E. cucuruli* and the previously characterized set of *Giardia* ESPs. What can the genome of the intracellular microsporidian parasite *E. cucuruli* tell us about the origin of the eukaryotic cell? The answer is a great deal as we investigate a minimal eukaryotic cell, which has eliminated all but the most essential functions.

Materials and Methods

Protein sequence databases of *S. cerevisiae*, *D. melanogaster*, *A. thaliana*, *G. lamblia*, and 44 bacteria and archaea were downloaded as previously described (Hartman and Fedorov 2002). Protein sequences of *E. cucuruli* (1997 entries) were downloaded from GenBank (Benson et al. 1999).

ESPs of *Giardia*. In our previous paper (Hartman and Fedorov 2002) we performed consecutive BLAST2.0 alignments of *S. cerevisiae* proteins with proteins of *D. melanogaster*, *C. elegans*, *A. thaliana*, and *G. lamblia* and those of 44 bacteria and archaea species. We used a blast score of 55 bits. This score was based on our consultation with experts in bioinformatics, as the *Giardia* database was in contigs only and was not assembled or annotated. As a result, 347 yeast proteins that have significant sequence similarity to proteins of all studied eukaryotes, but not to any bacteria or archaea (identity threshold of 55 blast score bits), were selected. We call this set ESPs of *G. lamblia*, since *Giardia* represents the most divergent eukaryotic proteins in the studied group.

ESPs of *E. cucuruli*. The comparison of *G. lamblia* and *E. cucuruli* proteins began with collecting ESPs of *E. cucuruli* which followed our previous approach for gathering ESPs of *Giardia*. We started from 6271 *S. cerevisiae* proteins and compared them with proteins of *D. melanogaster*, *C. elegans*, *A. thaliana*, *E. cucuruli*, and 44 sequenced bacteria and archaea using BLAST 2.0 alignment program (Altschul et al. 1997). This resulted in 401 *S. cerevisiae* sequences with significant similarity to all sequenced eukaryotes but not to any bacteria or archaea. The same similarity threshold of 55 blast score bits (used in our previous paper on *Giardia* ESPs [Hartman and Fedorov 2002]) has been employed in the present study which approximately corresponds to a P-value of 10^{-6} . This threshold is sufficiently stringent and, thus, allows us to assume that matched proteins share a common origin. We were also obliged to use the same 55 blast score as we were about to compare the *E. cucuruli* ESP results with those of *G. lamblia*. We then compared the 347 ESPs of *Giardia* with the 401 ESPs of *E. cucuruli*. It was found that 238 proteins are common to the ESPs of *G. lamblia* and *E. cucuruli*, while 109 ESPs are unique to *G. lamblia* and 163 are unique to ESP of *E. cucuruli*. These unique sets were studied further using PSI-BLAST programs.

PSI-BLAST. For each unique protein from the ESP set of *G. lamblia* that does not match an *E. cucuruli* ESP, we found the best-matched protein from *D. melanogaster* and *A. thaliana* proteome using the BLAST 2.0 program. These three protein sequences were aligned with each other by CLUSTALW 1.8 (Higgins and Sharp 1988) and the multiple alignments obtained was used as input for the PSI-BLAST program (Altschul et al. 1997) in the search of the *E. cucuruli* protein database. The results of the round 2 PSI-BLAST output were analyzed automatically by our PERL program (prog_PSI_401_2). When a protein from the *G. lamblia* ESP set had a PSI-BLAST alignment score with the *E. cucuruli*

Table 1. Comparison of ESP sets of *G. lamblia* and *E. cuculici*

| Protein set | BLAST 2.0 | | PSI-BLAST (round 2) | |
|------------------------------|-----------------------------|------------------------------|------------------------------------|--------------------------|
| | No. common (≥55 score bits) | No. unique (< 55 score bits) | No. weak-homolog (> 88 score bits) | unique (< 40 score bits) |
| 347 ESPs, <i>G. lamblia</i> | 238 | 109 | 57 | 52 |
| 401 ESPs, <i>E. cuculici</i> | 238 | 163 | | |

database higher than 88 bits, (in 98% of the cases it was > 100 bits), the protein was called a “putative homolog” and is shown in green in Fig. 1. Otherwise, the protein was called “unique” and is shown in red on in Fig. 1. At round 2 of PSI-BLAST all “unique” proteins had a psi-blast score of < 50 bits. We did not perform the reciprocal procedure for the PSI-BLAST comparison of unique *E. cuculici* ESPs with the *Giardia* database because the *G. lamblia* database consists of multiple short nucleotide sequences translated in six possible reading frames. This would lead to false negative results. All computational procedures were performed automatically by the PERL script prog_PSI_401_2. The program and all described protein sets are available at our web page, www.mco.edu/medicine/fedorov/E_cuculici.

Protein Groups. All 401 proteins were compared with each other by BLAST 2.0 binaries (Altschul et al. 1997). Next we performed the simplest grouping procedure: (1) two proteins were considered similar and put in the same group if they had a similarity score > 55 bits, and (2) groups were pooled together if any member of one had sequence similarity (≥55 bits) to any protein of another group. This procedure yielded 214 different protein groups.

Results

Since *Giardia* proteins (open reading frames) are still not assembled or annotated, we cannot compare eukaryote-specific proteins of *G. lamblia* with those of *E. cuculici* directly. In our previous paper, 914 proteins of *S. cerevisiae* which have sequence similarity to *D. melanogaster*, *C. elegans*, and *A. thaliana* (blast score, 55 bits) and no sequence similarity to Bacteria and Archaea were compared with the fragmented *Giardia* database. That comparison resulted in a set of 347 proteins comprised of 180 unrelated protein groups. We called these proteins ESPs of *G. lamblia*.

Here we used the same set of 914 eukaryote-specific proteins of *S. cerevisiae* and aligned it with the entire set of 1997 *E. cuculici* proteins. This analysis of 401 budding yeast proteins showed significant sequence similarity to the microsporidian proteome (similarity level, > 55 blast score bits). We called the 401 proteins that were obtained *E. cuculici* ESPs. The 401 ESPs of *E. cuculici* were compared to the 347 ESPs of *G. lamblia*. The main results of this comparison are summarized in Table 1 and the entire comparison is presented in Fig. 1. Two hundred thirty-eight proteins are common in the two sets, while 109 are unique to *G. lamblia* ESPs and 163 to ESPs of *E. cuculici* (see Table 1). However, when the 109 unique ESP *G. lamblia* proteins were compared

with the entire set of *E. cuculici* proteins using PSI-BLAST alignment, 57 of these proteins showed weak, yet significant, similarity to the microsporidian (shown in green in Fig. 1). The other 52 proteins remained specific to *Giardia* (shown in red in Fig. 1). Reciprocal psi-blast of *E. cuculici* unique ESPs was not performed for reasons outlined under Materials and Methods.

The ESP sets have some redundancy because of recent evolutionarily duplication of a substantial part of the *S. cerevisiae* genome. In order to take this redundancy into account, we divided the ESPs into unique protein groups. The similarity threshold of a blast score equaling 55 bits was used for grouping procedure. As a result, 401 ESPs of *E. cuculici* were divided into 214 unique groups, while 347 ESPs of *G. lamblia* were divided into 180 unique groups. Comparison of these groups demonstrated that 108 of the 180 *Giardia* ESP groups are common to *E. cuculici* ESPs (55-bit threshold of BLAST 2.0). When “putative homologs” revealed by psi-blast were taken into account, the number of common groups increased to 142.

Discussion

There are three groups of proteins in the ESPs of *E. cuculici*: (1) those that can be matched to the ESPs of *G. lamblia* by means of a blast score of 55 bits, (2) those that can be matched to the ESPs of *G. lamblia* by means of psi-blast, and (3) those that have no sequence similarity to the ESPs of *Giardia*. We discuss the relevance of these three groups to the structures of the eukaryotic cell.

The Plasma Membrane and the Cytoskeleton. As we discussed in our previous paper, one of the deepest distinctions between prokaryotes and eukaryotes is found in the plasma membrane (Hartman and Fedorov 2002). There is a lack of clathrin and associated proteins in the ESPs of *E. cuculici* (see Fig. 1). Since *E. cuculici* is an intracellular parasite, it possibly has a diminished need for clathrin-based endocytosis. This conjecture points out a difference between being an intracellular parasite (*E. cuculici*) and being an extracellular parasite (*Giardia*).

| Giardia | E. cuniculi | Giardia | E. cuniculi |
|--|--|--|---|
| PLASMA MEMBRANE AND ENDOCYTOSIS | | Vacuole | |
| Clathrin (<i>Chc1</i>) | --- | vacuolar protein (<i>Pep8</i>) | --- |
| clathrin-associated proteins (<i>Apl2, Apm1, Aps1, Aps2, Aps3, Apl1, Apl3, Apl4, Apl5, Apm2, Apm4</i>) | clathrin-associated proteins (<i>Apl2, Apl6, Aps1, Aps2, Aps3, Apl1, Apm4, Apl4, Apm1</i>) | retromer complex component (<i>Vps35</i>) | --- |
| dynamitin (<i>Dnm1, Mgm1, Vps1</i>) | dynamitin (<i>Dnm1, Mgm1, Vps1</i>) | vacuolar ATPase V0 domain subunit c (<i>Cup5</i>) | vacuolar ATPase V0 domain subunit c (<i>Cup5</i>) |
| CYTOSKELETON | | vacuolar ATPase V0 domain c" (<i>Ppa1, Tfp3</i>) | vacuolar ATPase V0 domain c" (<i>Ppa1, Tfp3</i>) |
| <i>Tubulin</i> | | --- | Vacuolar membrane protein (<i>Vps13</i>) |
| alpha-tubulin (<i>Tub1, Tub3</i>) | alpha-tubulin (<i>Tub1, Tub3</i>) | --- | vacuolar ATPase V1 domain su F (<i>Vma7</i>) |
| beta-tubulin (<i>Tub2</i>) | beta-tubulin (<i>Tub2</i>) | --- | vacuolar ATPase V1 domain su H (<i>Vma13</i>) |
| gamma tubulin-like protein (<i>Tub4</i>) | gamma tubulin-like protein (<i>Tub4</i>) | SIGNALING CASCADE | |
| <i>Tubulin-associated proteins</i> | | <i>Calmodulin</i> | |
| kinesin-related protein (<i>Kip2, Kar3</i>) | kinesin-related protein (<i>Kip2, Kar3</i>) | (<i>Cmd1</i>) | (<i>Cmd1</i>) |
| kinesin-related protein involved in mitosis (<i>Kip3</i>) | kinesin-related protein involved in mitosis (<i>Kip3</i>) | Ca-binding protein (<i>Cdc31</i>) | Ca-binding protein (<i>Cdc31</i>) |
| kinesin heavy chain homolog (<i>Smy1</i>) | kinesin heavy chain homolog (<i>Smy1</i>) | <i>Phosphatidylinositol</i> | |
| microtubule-binding protein (<i>Bim1</i>) | microtubule-binding protein (<i>Bim1</i>) | phosphatidylinositol kinases (<i>Vps34, Pik1, Sit4, Mss4, Tel1, Tor2, Tor1, Mec1</i>) | phosphatidylinositol kinases (<i>Vps34, Pik1, Sit4, Mss4, Tel1, Tor2, Tor1, Mec1, Tra1</i>) |
| putative light chain of dynein (<i>Dyn2</i>) | putative light chain of dynein (<i>Dyn2</i>) | phosphatidylinositol phosphatases (<i>Inp51, Inp52, Inp53</i>) | phosphatidylinositol phosphatases (<i>Inp51, Inp52, Inp53</i>) |
| <i>Actin</i> | | <i>Ubiquitin</i> | |
| (<i>Act1</i>) | (<i>Act1</i>) | ubiquitin (<i>Ubi4</i>) | ubiquitin (<i>Ubi4</i>) |
| <i>Actin-related proteins</i> | | ubiquitin-like protein (<i>Sml3</i>) | ubiquitin-like protein (<i>Sml3</i>) |
| (<i>Arp1, Arp2, Arp3, Arp4, Arp5, Arp6, Arp7</i>) | (<i>Arp1, Arp2, Arp3, Arp4, Arp5, Arp6, Arp7</i>) | ubiquitin-like protein (<i>Rub1</i>) | ubiquitin-like protein (<i>Rub1</i>) |
| <i>Actin-associated proteins</i> | | <i>Ubiquitin conjugation enzymes</i> | |
| light chain for myosin (<i>Mic1</i>) | light chain for myosin (<i>Mic1</i>) | (<i>Cdc34, Ubc1, Ubc4, Ubc5, Ubc6, Ubc8, Ubc9, Ubc11, Ubc12, Ubc13, Pex4, Qn8, Rad6</i>) | (<i>Cdc34, Ubc1, Ubc4, Ubc5, Ubc6, Ubc8, Ubc9, Ubc11, Ubc12, Ubc13, Pex4, Qn8, Rad6</i>) |
| --- | myosin (<i>Myo2, Myo3, Myo4, Myo5</i>) | <i>Ubiquitin proteases</i> | |
| PROTEIN SYNTHESIS AND BREAKDOWN | | (<i>Ubp5, Ubp14, Ubp8, Ubp15, Ubp12, Doa4, Ubp6, Ubp10</i>) | (<i>Ubp5, Ubp8, Ubp12, Ubp14, Ubp15, Doa4, Ubp7, Ubp9, Ubp11, Ubp13, Ubp6, Ubp10</i>) |
| <i>Small ribosomal proteins</i> | | --- | ubiquitin ligases (<i>Hul5</i>) |
| ribosomal protein S7 [<i>Rps30</i>] (<i>Rps7a, Rps7b</i>) | --- | --- | miscellaneous (<i>Ufd2, Ufd4, Ulp1, Apc11</i>) |
| ribosomal protein S21 (<i>Rps21a, Rps21b</i>) | --- | <i>GTP-binding proteins</i> | |
| ribosomal protein S24 (<i>Rps24a, Rps24b</i>) | ribosomal protein S24 (<i>Rps24a, Rps24b</i>) | Ras (<i>Ras1, Ras2, Rsr1, Rsg1, Tem1</i>) | Ras (<i>Ras1, Ras2, Rsr1, Rsg1, Tem1</i>) |
| ribosomal protein S26A (<i>Rps26a, Rps26b</i>) | ribosomal protein S26A (<i>Rps26a, Rps26b</i>) | Rho (<i>Rho1, Rho2, Rho3, Rho4, Rho5, Cdc42, Rdi1</i>) | Rho (<i>Rho1, Rho2, Rho3, Rho4, Rho5, Cdc42</i>) |
| ribosomal protein S27 (<i>Rps27a, Rps27b</i>) | ribosomal protein S27 (<i>Rps27a, Rps27b</i>) | Arf (<i>Arf1, Arf2, Arf3, Sar1, Arl1</i>) | Arf (<i>Arf1, Arf2, Arf3, Sar1, Arl1</i>) |
| ribosomal protein S31 [ubiquitin related] (<i>Rps31</i>) | ribosomal protein S31 [ubiquitin related] (<i>Rps31</i>) | Arf gap (<i>Age2, Gcs1, Age1</i>) | Arf gap (<i>Age2, Gcs1, Age1</i>) |
| <i>Large ribosomal proteins</i> | | Ran (<i>Gsp1, Gsp2, Yrb1</i>) | Ran (<i>Gsp1, Gsp2, Yrb1</i>) |
| --- | ribosomal protein L6 (<i>Rpl6b, Rpl6a</i>) | Rab (<i>Ypt1, Ypt6, Ypt7, Ypt10, Ypt13, Ypt32, Ypt52, Ypt53, Vps21</i>) | Rab (<i>Ypt1, Ypt6, Ypt7, Ypt10, Ypt13, Ypt32, Ypt52, Ypt53, Vps21</i>) |
| ribosomal protein L13 (<i>Rpl13b, Rpl13a</i>) | ribosomal protein L13 (<i>Rpl13b, Rpl13a</i>) | Rab gap (<i>Mdr1, Msb3</i>) | Rab gap (<i>Mdr1, Msb3</i>) |
| ribosomal protein L14 (<i>Rpl14a, Rpl14b</i>) | --- | GTP-binding related (<i>Cin4, Ypt11, Arl3</i>) | GTP-binding related (<i>Cin4, Ypt11, Arl3, Rna1, Srm1, Gyp1, Gyp7</i>) |
| ribosomal protein L18 (<i>Rpl18a, Rpl18b</i>) | ribosomal protein L18 (<i>Rpl18a, Rpl18b</i>) | <i>Cyclin</i> | |
| ribosomal protein L20 (<i>Rpl20b, Rpl20a</i>) | ribosomal protein L20 (<i>Rpl20b, Rpl20a</i>) | B-type cyclin (<i>Cib1, Cib2, Cib3, Cib4, Cib5, Cib6</i>) | B-type cyclin (<i>Cib1, Cib2, Cib3, Cib4, Cib5, Cib6</i>) |
| ribosomal protein L21 (<i>Rpl21b, Rpl21a</i>) | ribosomal protein L21 (<i>Rpl21b, Rpl21a</i>) | cell cycle checkpoint protein (<i>Bub3</i>) | cell cycle checkpoint protein (<i>Bub3</i>) |
| ribosomal protein L24 (<i>Rpl24a, Rpl24b</i>) | ribosomal protein L24 (<i>Rpl24a, Rpl24b</i>) | cyclin-dependent kinase-activating kinase (<i>Cak1</i>) | cyclin-dependent kinase-activating kinase (<i>Cak1</i>) |
| ribosomal protein L29 (<i>Rpl29</i>) | --- | <i>Kinases and phosphatases</i> | |
| ribosomal protein L33 (<i>Rpl33a</i>) | ribosomal protein L33 (<i>Rpl33a</i>) | serine/threonine protein kinase (<i>Cdc7, Sky1, Yki171w, Vps15, Iks1</i>) | serine/threonine protein kinase (<i>Cdc7, Sky1, Yki171w, Vps15, Iks1</i>) |
| ribosomal protein L35 (<i>Rpl35b, Rpl35a</i>) | ribosomal protein L35 (<i>Rpl35b, Rpl35a</i>) | involved in cell cycle (<i>Cdc50</i>) | involved in cell cycle (<i>Cdc50</i>) |
| ribosomal protein L36 (<i>Rpl36a, Rpl36b</i>) | ribosomal protein L36 (<i>Rpl36a, Rpl36b</i>) | subunit of the Cdc28 protein kinase (<i>Cks1</i>) | --- |
| ribosomal protein L40 [ubiquitin related] (<i>Rpl40a, Rpl40b</i>) | ribosomal protein L40 [ubiquitin related] (<i>Rpl40a, Rpl40b</i>) | LAMMER Protein Kinases (<i>Kns1</i>) | LAMMER Protein Kinases (<i>Kns1</i>) |
| <i>Translation factors</i> | | β subunit of casein kinase II (<i>Ckb1, Ckb2</i>) | β subunit of casein kinase II (<i>Ckb1, Ckb2</i>) |
| translation elongation factor EF-1beta (<i>Efb1</i>) | --- | dual-specificity tyrosine phosphatases (<i>Pps1, Yvh1, Tep1, Cdc14</i>) | --- |
| translation elongation factor EF-1gamma (<i>Cam1, Tef4</i>) | --- | protein phosphatase regulatory subunits (<i>Cdc55, Cnb1, Sds22, Rts1, Tpd3</i>) | protein phosphatase regulatory subunits (<i>Cdc55, Cnb1, Sds22, Rts1, Tpd3</i>) |
| <i>Proteasome associated proteins</i> | | protein phosphatase type 2C (<i>Ptc1, Ptc2, Ptc3, Ptc4</i>) | --- |
| subunits of proteasome regulatory particle (<i>Rpn1, Rpn8, Rpn11, Rpn10</i>) | subunits of proteasome regulatory particle (<i>Rpn1, Rpn8, Rpn11, Rpn2, Rpn3, Rpn5, Rpn6, Rpn7, Rpn10</i>) | myotubularin dual specific phosphatase (<i>Yjr110w</i>) | --- |
| <i>Signal peptidase</i> | | <i>14-3-3 proteins</i> | |
| (<i>Spe3</i>) | (<i>Spe3</i>) | (<i>Bmh1, Bmh2</i>) | (<i>Bmh1, Bmh2</i>) |
| MEMBRANE | | | |
| <i>ER and Golgi</i> | | | |
| transport protein particle TRAPP comp. (<i>Bet3</i>) | transport protein particle TRAPP comp. (<i>Bet3</i>) | | |
| HDEL receptor (<i>Erd2</i>) | HDEL receptor (<i>Erd2</i>) | | |
| integral membrane proteins (<i>Sac1, Fig4</i>) | integral membrane proteins (<i>Sac1, Fig4</i>) | | |
| --- | SNARE docking (<i>Sec1, Sly1, Vps45</i>) | | |
| subunit of coatomer (<i>Sec26</i>) | subunit of coatomer (<i>Sec26</i>) | | |
| vesicle coat component (<i>Sec24</i>) | vesicle coat component (<i>Sec24</i>) | | |
| miscellaneous (<i>Gpi8</i>) | miscellaneous (<i>Gpi8, Hrd1, Rer1, Yip1, Sso1, Sso2, Sec21, Syg1</i>) | | |

A
Fig. 1. Continued.

| Giardia | E. cucuruli | Giardia | E. cucuruli |
|---|---|---|--|
| Lipid attachments | | OTHERS | |
| geranylgeranyltransferase type II β (Bet2) | geranylgeranyltransferase type II β (Bet2) | riboflavin kinase (<i>Fmn1</i>) | Enzymes |
| geranylgeranyltransferase type II α (Bet4) | geranylgeranyltransferase type II α (Bet4) | FAD synthetase (<i>Fad1</i>) | ---- |
| geranylgeranyltransferase type I (Cdc43) | geranylgeranyltransferase type I (Cdc43) | protein carboxyl methylase (Ycr047c) | protein carboxyl methylase (Ycr047c) |
| farnesyltransferase beta subunit (Ram1) | farnesyltransferase beta subunit (Ram1) | N-terminal acetyltransferase (Nat3) | N-terminal acetyltransferase (Nat3) |
| CAAX farnesyltransferase α (Ram2) | CAAX farnesyltransferase α (Ram2) | acetyltransferase (SAS gene family) (Esa1) | acetyltransferase (SAS gene family) (Esa1) |
| farnesyl cysteine-carboxyl methyltransferase (Ste14) | farnesyl cysteine-carboxyl methyltransferase (Ste14) | glucosamine-phosphate N-acetyltransferase (Gna1) | glucosamine-phosphate N-acetyltransferase (Gna1) |
| N-myristoyl transferase (Nmt1) | N-myristoyl transferase (Nmt1) | UDP Glucose pyrophosphorylase (Yhi012w) | UDP Glucose pyrophosphorylase (Yhi012w) |
| Rab geranyltransferase regulatory (Mrs6) | Rab geranyltransferase regulatory (Mrs6) | phosphoryltransferase (Gpi13) | phosphoryltransferase (Gpi13) |
| NUCLEUS | | ---- | Similar to NADH dehydrogenase (<i>Yif1</i>) |
| Histones | | ---- | Phosphoacetylglucosamine mutase (<i>Pcm1</i>) |
| histone H2A (Hta1; Hta2) | histone H2A (Hta1; Hta2) | ---- | Desaturase/hydroxylase enzyme (<i>Scs7</i>) |
| histone H2B (Htb1; Htb2) | histone H2B (Htb1; Htb2) | ---- | N-AGPinositol de-N-acetylase (<i>Gpi12</i>) |
| histone H3 (Hht2; Hht1) | histone H3 (Hht2; Hht1) | ---- | Polyphosphate synthetase (<i>Vlc2; Vlc3; Vlc4</i>) |
| histone H4 (Hhf1; Hhf2) | histone H4 (Hhf1; Hhf2) | ---- | Mannosyltransferase (<i>Pmi2; Pmi3; Pmi5</i>) |
| Histone-associated proteins | | ---- | Lipid phosphate phosphatase (<i>Dpp1; Lpp1</i>) |
| histone acetyltransferase (Gcn5; Hal2) | histone acetyltransferase (Gcn5; Hal2) | ---- | Phosphorylcholine transferase (<i>Muq1; Pct1</i>) |
| Chromatin binding proteins (Cse4) | Chromatin binding proteins (Cse4; <i>Esp1; Rsc4; Pho23; Yng1</i>) | ---- | Sterol-ester synthetase (<i>Are1; Are2</i>) |
| Topoisomerase I | | Clusters of unknown proteins | |
| (Trf5; Trf4) | (Trf5; Trf4) | (Ydr126w; Erf2; Ydr459c; Ynl326c; Yolo03c) | (Ydr126w; Erf2; Ydr459c; Ynl326c; Yolo03c) |
| Transcriptional factors | | (Psr2; Ypl063w; Psr1; Nem1) | (Psr2; Ypl063w; Psr1; <i>Nem1</i>) |
| (Hap3; Set2; Sps18; Ssl1; Gts1; Htz1; <i>Mob1; Mob2; Sip2</i>) | (Hap3; Set2; Sps18; Ssl1; <i>Hap2; Hap5; Fcp1; Bdf1; Bdf2; Spl5; Spl7; Taf60; Tfa1; Taf19; Toa1; Asf1; Fhl1; Tsm1; Cdc39; Bas1; Ada2; Bur6; Arg80; Cor4; Sin3; Azf1; Pzf1; Fzf1; Crz1; Gts1; Mob2</i>) | (Ygl014w; Mpt5; Yli013c) | (Ygl014w; Mpt5; Yli013c) |
| Zinc fingers | | (Tom1; Rsp5; Hul4) | (Tom1; Rsp5; Hul4) |
| (Mot2; Cth1; Sas2; Glo3; Tis11; Ybr267w) | (Mot2; Cth1; Sas2; Glo3; Tis11; <i>Ssu72; Sfb2; Ybr267w</i>) | (Yli088c; Ybl089w; Ynl101w) | (Yli088c; Ybl089w; Ynl101w) |
| DNA-repair enzymes | | (Ypl249c; Msb4) | (Ypl249c; Msb4) |
| --- | (Pol12; Pob3; Hys2; Soh1; Mus81; Rfa1; <i>Ixr1; Mre11</i>) | (Msi1; Rsa2) | (Msi1; Rsa2) |
| RNA-polymerase | | (Imp4; Rpf1) | (Imp4; Rpf1) |
| subunits of RNA polymerases (Rpc19; <i>Rpb8</i>) | subunits of RNA polymerases (<i>Rpc19</i>) | (Mrd1; Ynl110c) | (Mrd1; Ynl110c) |
| Spliceosomal proteins | | (Gdi1) | (Gdi1) |
| (Smd3; Prp8; Prp9; Lsm2) | (Smd3; Prp8; Prp9; <i>Smd2; Smx3; Cef1; Hsh49; Cus1; Cbc2; Luc7; Lsm2</i>) | (Ydi060w; Bms1) | (Ydi060w; Bms1) |
| RNA enzymes | | (Yfr021w; Ypl100w; Ygr223c) | (Yfr021w; Ypl100w; Ygr223c) |
| RNA exonuclease (Rex3) | RNA exonuclease (Rex3; <i>Rex4</i>) | (Ykl121w; Ymr102c) | (Ykl121w; Ymr102c) |
| ribonuclease H (Rnh70) | ribonuclease H (Rnh70) | (Ssf1; Ssf2) | (Ssf1; Ssf2) |
| mRNA guanylyltransferase [capping] (Ceg1) | mRNA guanylyltransferase [capping] (Ceg1) | (Vps24; F51; Ykl002w) | --- |
| RNA (guanine-7-methyltransferase (Abd1)) | RNA (guanine-7-methyltransferase (Abd1)) | (Ylr328w; Ygr010w) | --- |
| poly(A) polymerase (Pap1) | poly(A) polymerase (Pap1) | --- | (Ymr285c; Yoi042w; Yor353c) |
| --- | miscellaneous (<i>Rna14; Ctf2</i>) | --- | (Ymi072c; Yrl087w; Yor086c) |
| Nucleolus | | Unique unknown proteins | |
| nucleolar protein (<i>Ebp2</i>) | nucleolar protein (<i>Nop13</i>) | Sas3; Ydi216c; Sfb3; Las21; Ynr048w; | Sas3; Ydi216c; Sfb3; Las21; Ynr048w; |
| small nucleolar RNP proteins (Gar1) | small nucleolar RNP proteins (Gar1) | Rlp7; Yhr122w; Ylr409c; Nip7; Ypr031w; | Rlp7; Yhr122w; Ylr409c; Nip7; Ypr031w; |
| protein required for biogenesis of the 60S ribosomal subunit (Brx1) | protein required for biogenesis of the 60S ribosomal subunit (Brx1) | Hrt1; Nmd3; Yer082c; Ykl099c; Pri2; Yth1; | Hrt1; Nmd3; Yer082c; Ykl099c; Pri2; Yth1; |
| Nuclear Pore and Transport | | Yer126c; Nud1; Ypl247c; Ypl236c; | Yer126c; Nud1; Ypl247c; Ypl236c; |
| nuclear pore protein (Gle2; Nsp1; <i>Ntf2</i>) | nuclear pore protein (Gle2; <i>Nsp1</i>) | Pip1; Ufd1; Ygr145w; Crm1; Ptk1; Ykt6; | Pip1; Ufd1; Ygr145w; Crm1; Ptk1; Ykt6; |
| karyopherin alpha (Srp1) | karyopherin alpha (Srp1; <i>Kap95</i>) | Ydr266c; Ydr339c; Ydr411c; Ydr083w; | Ydr266c; Ydr339c; Ydr411c; Ydr083w; |
| putative nuclear protein (Mak16) | putative nuclear protein (Mak16) | Yli005w; Yhr186c; Ydr365c; Abp140; | Yli005w; Yhr186c; Ydr365c; Abp140; |
| --- | nuclear pore complex (<i>Nup100; Nup116</i>) | Bph1; Ymr068w; Yjl109c; Sgl1; | Bph1; Ymr068w; Yjl109c; Sgl1; |
| RNA binding and export | | <i>Vip1; Ybr228w; Yol093w; Enp1; Yil113w; Ykl100c; Ygl047w; Ykl056c; Pac10; Ent3; Hym1; Yor289w</i> | <i>Vip1; Ybr228w; Yol093w; Enp1; Yil113w; Ykl100c; Ygl047w; Ykl056c; Pac10; Ent3; Hym1; Yor289w</i> |
| --- | (Mip6; Lhp1; Gbp2; Pes4) | --- | <i>Ama1; Caf40; Cdc20; Cdc33; Dcp2; Dst1; Erv1; Fkh2; Fra1; Grd19; Hul1; Ybi071w-a; Lem3; Lsm1; Per1; Sac6; Smp2; Yip2; Yaf9; Ybr042c; Ycr063w; Ycs4; Ydr018c; Ydr128w; Ydr255c; Ygl131c; Ygr054w; Yjr070c; Yju2; Ykl059c; Ykl146w; Ylr002c; Ylr323c; Ymr073c; Ymr093w; Ymr288w; Yor091w; Yor175c; Ypl110c; Ypr037c; Tos8; Yng2; Ypr133c; Zap1; Jsn1</i> |

B

Fig. 1. Comparison of proteins from the 347-ESP set of *G. lamblia* and 401-ESP set of *E. cucuruli*. The unique identifier symbols for the proteins are from the Saccharomyces Genome Database (<http://genome-www.stanford.edu/Saccharomyces>) and are shown in parentheses in different colors. The **black** color of a protein identifier from the ESP set of *G. lamblia* shows that this protein is also present among 401 ESPs of *E. cucuruli* or has a "putative homolog" in the *E. cucuruli* proteome revealed by PSI-BLAST. The **black** color of a protein identifier from the ESP set of *E. cucuruli* means that this protein is also present among the 347-ESP set of *G. lamblia*. The **green** color of a protein identifier shows

that this protein is missing from the 401-ESP set of *E. cucuruli*, however, its "putative homolog," detected by psi-blast, is present in the *E. cucuruli* proteome. The **red** color of a protein identifier from the 347-ESP set of *G. lamblia* shows that this protein does not have any sequential similarity to *E. cucuruli* proteins. The **red** color of a protein identifier from the 401-ESP set of *E. cucuruli* means that this protein does not show sequence similarity above 55 bits (BLAST 2.0 score) while screening the *G. lamblia* contig database. The sum of all black and red identifiers for *G. lamblia* is equal to 347, and that for *E. cucuruli* is equal to 401.

The ESPs most closely associated with the plasma membrane are cytoskeletal proteins, such as actin and associated proteins. *G. lamblia* and *E. cuciculi* have a full complement of actin and actin-related proteins (see Fig. 1). The three-dimensional structure of actin showed unexpected structural similarities to hexokinase and HSP 70 (Kabsch and Holmes 1995). These three proteins have a common nucleotide-binding motif, as they all bind ATP in the presence of calcium or magnesium ions. The “actin fold” was also found in the bacterial proteins FtsA and MreB (Kabsch and Holmes 1995). What is the evolutionary relationship between actin and bacterial proteins that have the “actin fold”? The role of MreB filaments in the cell shape of bacteria has suggested that the actin cytoskeleton evolved from MreB filaments (van den Ent et al. 2001). However, a careful sequence comparison of MreB and actin demonstrated that MreB and actin most likely had a common ancestor in the very distant past but that MreB was not an evolutionary precursor of actin (Doolittle and York 2002). Presumably, actin and all the related nucleotide-binding proteins evolved from a common ancestor which bound a nucleotide. This gene may have undergone gene doubling and various insertions, eventually evolving into MreB and, independently, actin proteins (Kabsch and Holmes 1995). Finally, in all such cases there is an alternative explanation of an independent origin of proteins with very little or without sequence similarity due to convergent evolution. Both *G. lamblia* and *E. cuciculi* have a full set of actin-related proteins (abbreviated *Arp* in Fig. 1). *Arp1* is involved in spindle alignment. *Arp2* and *Arp3* are involved in actin polymerization. *Arp4*, *Arp5*, and *Arp7* are involved in chromatin remodeling. *Arp6* is localized to heterochromatin (Goodson and Hawse 2002).

Both *G. lamblia* and *E. cuciculi* have tubulins among their ESPs. The three-dimensional structure of tubulin showed structural similarities to glyceraldehyde-3-phosphate dehydrogenase and bacterial protein FtsZ (Nogales et al. 1998). These three types of proteins have a common nucleotide-binding motif which is similar to the Rossmann fold, distinct from the actin fold (Doolittle 1995). As in the case with actin, we consider the existence of a deep ancestor to tubulin and FtsZ rather than bacterial FtsZ as the possible precursor to tubulin.

In close association with cytoskeleton composed of actin and tubulin, there are the motor proteins dynein, myosin, and kinesin. This association allows the eukaryotic cell to move, phagocytize, and divide. The control of this structure is modulated through interaction with calcium ions. Both *E. cuciculi* and *G. lamblia* have kinesins in their ESPs (Fig. 1). *E. cuciculi* has four myosins among its ESPs, whereas none exists in *G. lamblia*. It appears that myosin is a later

addition to the evolving eukaryotic cell, as fungi branched much later in the phylogenetic tree than the diplomonads. The absence of dynein in the ESPs of both *E. cuciculi* and *G. lamblia* is due to its relationship to AAA proteins which are widely dispersed among the bacteria (Vale 2000). The motor proteins myosin and kinesin have a common structural motif with the G proteins (Kull et al. 1998). The similarity of these proteins is due to a common nucleotide-binding motif, which differs from that of actin or tubulin.

The fact that many ESPs have structural similarity to the bacterial proteome is not surprising because, likely, all cells share a common origin. However, the lack of sequential similarity of the discussed conservative structural proteins with their prokaryotic counterparts gave us a foundation to conjecture that the cytoskeleton including the motor proteins of the eukaryotic cell evolved out of a set of nucleotide-binding proteins from a primitive RNA-based cell and not out of a bacterial or an archaeal cell (Hartman and Fedorov 2002).

The Endoplasmic Reticulum and Protein Synthesis. The cytoskeleton of the eukaryotic cell is closely associated with the endoplasmic reticulum (ER). *E. cuciculi* has a significantly larger number of ER and Golgi ESPs than found in *G. lamblia* (16 and 6 proteins, respectively). This difference implies that the ER and Golgi of the parasitic fungus *E. cuciculi* are more complex than those of *G. lamblia*. The evolutionary study of the ER and Golgi complex is still in its infancy (Beznoussenko and Mironov 2002). We need a much larger sample of ER and Golgi proteins from single-celled protozoans, which are not parasitic, to infer the origin and evolution of these membranes.

Prominent among the ESPs of *E. cuciculi* and *G. lamblia* are ubiquitin, ubiquitin ligases, and proteases (Fig. 1). Some proteins of Archaea and Bacteria have a 3D structural fold similar to that of ubiquitin and are involved in the biosynthesis of sulfur-containing coenzymes (Wang et al. 2001). These prokaryotic proteins and ubiquitin might have diverged from a common ancestor, yet they have evolved independently in prokaryotes and eukaryotes. We found several proteasome-associated proteins from the 19S proteasome regulatory particle in the ESP sets of both *G. lamblia* and *E. cuciculi* (Fig. 1). There are more proteins associated with the proteasome regulatory particle of *E. cuciculi* than that of *G. lamblia* (eight and four proteins, respectively). This might point to an evolution from the simpler regulatory system of *G. lamblia* to that of the fungi. We hypothesize that the eukaryotic protein degradation complex composed of the ubiquitin, ubiquitin ligases, ubiquitin proteinases, and 19S regulatory proteasome

particle did not originate from prokaryotes. At the same time, the ancient proteasome likely came from an archaeal endosymbiont (Bouzat et al. 2000).

The Nucleus. Using the relatively strong threshold for protein similarity of 55 blast score bits ($\sim 10^{-6}$ e-values), we revealed four histones (H2A, H2B, H3, and H4) in the nuclear ESPs of *G. lamblia* and only two histones (H3 and H4) in *E. cucurbiti*. However, if we search for sequence similarity using psi-blast, as opposed to blast, then we pick up highly diverged histones H2A and H2B of *E. cucurbiti*. The eukaryotic histones share the same 3D structure with the archaeal histone-like proteins of the Euryarchaeota (methanogens, etc.) (Arents and Moudrianakis 1995). Unlike actin, tubulin, ubiquitin, and the GTP-binding proteins whose 3D counterparts are found throughout the Archaea and Bacteria, the histone fold is only found in the Euryarchaeota, and not in the Crenarchaeota or the Bacteria. Presently, the simplest explanation for the evolution of histones is that a histone-like protein came from an ancient archaeal cell and subsequently evolved into the full eukaryotic complement of histones. As for other proteins connected to the DNA structure, there are two topoisomerase I ESPs (Trf4 and Trf5) found in both *E. cucurbiti* and *G. lamblia*. There are eight *E. cucurbiti* ESPs involved with DNA repair that are not found in *G. lamblia*.

The ESPs found in the nucleus are dominated by proteins involved in the synthesis, processing, and transport of RNAs out of the nucleus into the cytoplasm. There are two ESPs (Rpc19 and Rpb8) associated with the RNA polymerases of *G. lamblia*. We could only detect by psi-blast one diverged ESP (Rpc19) in *E. cucurbiti*. The ESPs representing transcription factors are much more diverse in *E. cucurbiti* (29 proteins) than in *G. lamblia* (9 proteins). They share only four transcription factors in common. These deep distinctions show that transcriptional factors play a major role in the evolution of eukaryotes. However, when we compare groups of zinc finger ESPs between these two cells, we find more than 70% identity. The groups of nucleolar proteins ESPs associated with the synthesis and transport of ribosomal RNA are very similar in both *E. cucurbiti* and *G. lamblia*. There are more ESP spliceosomal proteins in *E. cucurbiti* (10) than in *G. lamblia* (4). The ESPs involved in the transport proteins and the nuclear pore proteins are found in equal abundance in *E. cucurbiti* and *G. lamblia*.

The Cell Cycle and Cellular Coordination. The regulators of the eukaryotic cell cycle (cyclins, serine/threonine kinases, and ubiquitin proteins) are present among ESPs (Fig. 1). When we compare the cyclin ESPs of *E. cucurbiti* with those of *G. lamblia*, we see

complete agreement between these two cells. However, when we compare the ESPs representing kinases and phosphatases of these two groups, we get a different picture. There is a significantly higher number of kinases and phosphatases in *G. lamblia* than in *E. cucurbiti* (24 versus 10 proteins). This must be due to the difference in intracellular versus extracellular lifestyles.

The GTP-binding proteins *ras* (plasma membrane), *rho* (cytoskeleton), *rab* (endoplasmic reticulum), *arf* (Golgi), and *ran* (nucleus–cytoplasm) are very prominent among the ESPs of *E. cucurbiti* and *Giardia*. We hypothesize that these proteins *ras*, *rho*, *rab*, and *arf* have evolved from the membrane-protein synthesizing machinery and the cytoskeleton of the host cell. They are now localized on the cytoskeleton and membranes (the plasma, endoplasmic reticulum, and Golgi) of the eukaryotic cell (Hartman and Fedorov 2002).

Enzymes. There are 10 enzymes found in *E. cucurbiti* that are not found in the ESPs of *G. lamblia*. This may be due to the fact that *E. cucurbiti*, as a fungus, belongs to the crown of the eukaryotic tree (Katinka et al. 2001), while *G. lamblia*, as a diplomonad, branched earlier. Finally, there are 108 and 91 of ESPs of *E. cucurbiti* and *G. lamblia*, respectively, that have no assigned function. Among them, 57 proteins were found in both *E. cucurbiti* and *G. lamblia*.

Conclusions

We found 401 ESPs in *E. cucurbiti*. These ESPs represent 214 unique protein groups. Comparison of the ESPs of *G. lamblia* with the proteins of *E. cucurbiti* demonstrated that 85% of the 347 ESPs of *G. lamblia* have sequence similarity to ESPs found in *E. cucurbiti*. Proteins from the ESP sets fall into two main categories: (1) proteins related to the observable structures in the cytoplasm of the eukaryotic cell such as the plasma membrane (clathrin), the cytoskeleton (actin and arps, tubulin, and associated kinesins), the endoplasmic reticulum, the Golgi, and the nucleus (histones, etc.), and (2) proteins involved in the coordination of the eukaryotic cell such as GTP-binding proteins (i.e., *ras*, *rho*, *rab*, *arf*, and *ran*), calmodulin, ubiquitin, cyclin, serine–threonine kinases, and phosphatases, 14–3–3 proteins, and enzymes modulating PIP (phosphatidyl inositol phosphates). These cellular structures and their defining proteins are unique to the eukaryotic cell and so are the control proteins.

A significant number of ESP sets prominent in the cytoplasm, such as actin, tubulin, kinesins, ubiquitin, and GTP-binding proteins, all have counterparts in

prokaryotes with similar 3D folds but no significant sequential similarity. These proteins are also found among the most conserved proteins in the eukaryotic cell (Copley et al. 1999). The best that can be inferred from these facts is that prokaryotic proteins and their ESP counterparts had a common ancestor, and that ancestor was a more primitive RNA-based cell (Doolittle and York 2002; Nogales et al. 1998).

The large number and diversity of ESPs point to a very ancient origin of the eukaryotic cell. These proteins are congruent with the recent hypothesis which places eukaryotes at the root of the universal tree of life (Gribaldo and Phillippe 2002; Penny and Poole 1999). The data presented here are also consistent with our previous hypothesis that the eukaryotic cell had an RNA-based cell as one of its ancestors (chronocyte) and that the nucleus was formed by the engulfing of prokaryotic cells that became the nuclear endosymbiont (Hartman and Fedorov 2002).

Acknowledgments. This research was supported in part by NSF Grant DBI-0205512. Support for this work was provided by the Medical College of Ohio Foundation. We would also like to thank Ms. Lisa Johnston for her excellent secretarial assistance.

References

- Altschul SF, Madden TL, Schaffer AA, Zhang J, Zhang Z, Miller W, Lipman DJ (1997) Gapped BLAST and PSI-BLAST: A new generation of protein database search programs. *Nucleic Acids Res* 25:3389–3402
- Arents G, Moudrianakis EN (1995) The histone fold: A ubiquitous architectural motif utilized in DNA compaction and protein dimerization. *Proc Natl Acad Sci USA* 92:11170–11174
- Baldauf SL, Roger AJ, Wenk-Siefert I, Doolittle WF (2000) A kingdom-level phylogeny of eukaryotes based on combined protein data. *Science* 290:972–976
- Benson DA, Boguski MS, Lipman DJ, Ostell J, Ouellette BF, Rapp BA, Wheeler DL (1999) GenBank. *Nucleic Acids Res* 27:12–17
- Beznoussenko GV, Mironov AA (2002) Models of intracellular transport and evolution of the Golgi complex. *Anat Rec* 268:226–238
- Bouzat JL, McNeil LK, Robertson HM, Solter LF, Nixon JE, Beever JE, Gaskins HR, Olsen G, Subramaniam S, Sogin ML, Lewin HA (2000) Phylogenomic analysis of the alpha protease gene family from early-diverging eukaryotes. *J Mol Evol* 51:532–543
- Copley RR, Schultz J, Ponting CP, Bork P (1999) Protein families in multicellular organisms. *Curr Opin Struct Biol* 9:408–415
- Doolittle RF (1995) The origins and evolution of eukaryotic proteins. *Philos Trans R Soc Lond B Biol Sci* 349:235–240
- Doolittle RF, York AL (2002) Bacterial actins? An evolutionary perspective. *Bioessays* 24:293–296
- Fraser CM, Gocayne JD, White O, et al. (1995) The minimal gene complement of *Mycoplasma genitalium*. *Science* 270:397–403
- Goodson HV, Hawse WF (2002) Molecular evolution of the actin family. *J Cell Sci* 115:2619–2622
- Gribaldo S, Phillippe H (2002) Ancient phylogenetic relationships. *Theor Population Biol* 61:391–408
- Hartman H, Fedorov A (2002) The origin of the eukaryotic cell: A genomic investigation. *Proc Natl Acad Sci USA* 99:1420–1425
- Higgins DG, Sharp PM (1988) CLUSTAL: A package for performing multiple sequence alignment on a microcomputer. *Gene* 73:237–244
- Kabsch W, Holmes KC (1995) The actin fold. *FASEB J* 9:167–174
- Katinka MD, Duprat S, Cornillot E, Metenier G, Thomarat F, Prensier G, Barbe V, Peyretailade E, Brottier P, Wincker P, Delbac F, El Alaoui H, Peyret P, Saurin W, Gouy M, Weissenbach J, Vivares CP (2001) Genome sequence and gene compaction of the eukaryote parasite *Encephalitozoon cuniculi*. *Nature* 414:450–453
- Kull FJ, Vale RD, Fletterick RJ (1998) The case for a common ancestor: Kinesin and myosin motor proteins and G proteins. *J Muscle Res Cell Motil* 19:877–886
- Mushegian AR, Koonin EV (1996) A minimal gene set for cellular life derived by comparison of complete bacterial genomes. *Proc Natl Acad Sci USA* 93:10268–10273
- Nogales E, Downing KH, Amos LA, Lowe J (1998) Tubulin and FtsZ form a distinct family of GTPases. *Nat Struct Biol* 5:451–458
- Penny D, Poole A (1999) The nature of the last universal common ancestor. *Curr Opin Genet Dev* 9:672–677
- Razin S (1997) The minimal cellular genome of mycoplasma. *Indian J Biochem Biophys* 34:124–130
- Vale RD (2000) AAA proteins: Lords of the ring. *J Cell Biol* 150:F13–F19
- Van den Ent F, Amos L, Lowe J (2001) Prokaryotic origin of the actin cytoskeleton. *Nature* 413:39–44
- Vossbrinck CR, Maddox JV, Friedman S, Debrunner-Vossbrinck BA, Woese CR (1987) Ribosomal RNA sequence suggests microsporidia are extremely ancient eukaryotes. *Nature* 326:411–414
- Wang C, Xi J, Begley TP, Nicholson LK (2001) Solution structure of ThiS and implications for the evolutionary roots of ubiquitin. *Nat Struct Biol* 8:47–51