

1 high quality draft sequences will be produced for both in the near future. It is
2 hoped that the *E. dispar* sequence will prove useful in identifying genomic
3 differences linked to disease causation while that of *E. invadens* will be used to
4 study patterns of gene expression during encystation. Small-scale genome
5 surveys have been performed for two other species: *E. moshkovskii*, which is
6 primarily a free-living species although it occasionally infects humans, and *E.*
7 *terrapinae*, a reptilian commensal species
8 (http://www.sanger.ac.uk/Projects/Comp_Entamoeba/).

9

10 **2. GENOME STRUCTURE**

11 **2.1 The *E. histolytica* Genome Sequencing, Assembly and** 12 **Annotation Process**

13 The first choice to be made in the genome project was perhaps the easiest - the
14 identity of the strain to be used for sequencing. A significant majority of the
15 existing sequence data prior to the genome project was derived from one strain:
16 HM-1:IMSS. This culture was established in 1967 from a rectal biopsy of a
17 Mexican man with amoebic dysentery and axenised shortly thereafter. It has
18 been used widely for virulence, immunology, cell biology and biochemistry in
19 addition to genetic studies. In an attempt to minimise the effects of long-term
20 culture cryopreserved cells that had been frozen in the early 1970s were revived
21 and this uncloned culture used to generate the DNA for sequencing.

22

23 Before undertaking a genome scale analysis it is important to understand the
24 quality and provenance of the underlying data. The *E. histolytica* genome was
25 sequenced by whole genome shotgun approach with each center generating
26 roughly half of the reads. Several different DNA libraries containing inserts of
27 different sizes were produced using DNA that had been randomly sheared and
28 sequences were obtained from both ends of each cloned fragment. The Phusion
29 assembler (Mullikin and Ning, 2003) was used to assemble the 450,000 short
30 reads into larger contigs (contiguous sequences), resulting in 1819 genome
31 fragments that were approximately 12X deep, which means that each base has

1 been sequenced 12 times, on average. While the genome shotgun sequence
2 provides high coverage of each base it is inevitable that there will be
3 misassemblies and sequencing errors in the final consensus particularly towards
4 each end of the contigs. Another problem with draft sequence is that it contains
5 gaps, and while most of these will be small and will mostly contain repetitive
6 non-coding “junk” sequence, some of the gaps will probably contain genes.
7 This makes it impossible to be absolutely certain of the absence of particular
8 genes in *E. histolytica* and, in some cases, the presence or absence of particular
9 biological pathways. Due to the high repeat content and low GC content
10 (24.1%) of the *E. histolytica* genome, closure of the remaining gaps is likely to
11 be a lengthy process. Therefore it was decided to undertake and publish an
12 analysis of the genome draft following assembly of the shotgun reads.

13

14 Annotation of the protein coding regions of the genome was initially carried out
15 using two genefinders (GlimmerHMM (Majoros *et al.*, 2004) and Phat (Cawley
16 *et al.*, 2001)) previously used successfully on another low G+C genome, that of
17 *P. falciparum*. The software was re-trained specifically for analysis of the *E.*
18 *histolytica* genome. The training process involved preparing a set of 600
19 manually edited genes to be used as models with the subsequent genefinding
20 then being carried out on all of the assembled contigs to generate a 'complete'
21 gene set. Predicted gene functions were generated automatically by homology
22 searches using public protein and protein-domain databases, with subsequent
23 refinement of identifications being carried out by manual inspection. For
24 particular genes and gene families of special interest, members of the
25 *Entamoeba* scientific community were involved throughout this process as
26 expert curators with each individual assisting in the analysis and annotation of
27 their genes of interest. Therefore although the manual curation of the genome
28 has not been systematic, those areas of biology that are of primary interest to
29 the *Entamoeba* community have been annotated most thoroughly. The
30 publication of the genome by Loftus *et al* therefore represents a “first draft” of
31 the complete genome sequence and the level of annotation is similar to the

1 initial publications of other genomes such as *Drosophila* (Adams *et al.*, 2000;
2 Myers *et al.*, 2000) and human (Lander *et al.*, 2001).

3

4 **2.2 Karyotype and Chromosome Structure.**

5 The current *E. histolytica* genome assembly is approximately 23.7 million
6 base pairs (Mbp) in size (Table 1). This figure is not likely to be a very accurate
7 measure. In part this is due to misassembly of repetitive regions, which will
8 cause the genome to appear smaller, and in part because of the possibility of
9 aneuploidy in some regions of the genome, which would cause them to appear
10 more than once in the assembly. Overall, however, this size is not inconsistent
11 with data from pulse-field gels (Willhoeft and Tannich, 1999) and kinetic
12 experiments (Gelderman *et al.*, 1971a,b) making the *E. histolytica* genome
13 comparable in size (24 Mbp) to that of *Plasmodium falciparum* (23 Mbp)
14 (Gardner *et al.*, 2002), *Trypanosoma brucei* (26 Mbp) (Berriman *et al.*, 2005),
15 and the free living amoeba *Dictyostelium discoideum* (34 Mbp) (Eichinger *et*
16 *al.*, 2005).

17

18 The current assembly does not represent complete chromosomes. Analysis of
19 pulse-field gels predicts 14 chromosomes ranging in sizes from 0.3 to 2.2 Mb
20 and possibly a ploidy of four (Willhoeft and Tannich, 1999). There is no current
21 information regarding the size and nature of the centromeres and there are no
22 contigs that appear to contain likely centromeric regions based on comparisons
23 with other organisms. A search for signature telomeric repeats within the data
24 indicates that these are either not present in the genome, not present in our
25 contigs, or are diverged enough to be unidentifiable. However, there is
26 circumstantial evidence that the chromosome ends may contain arrays of tRNA
27 genes (see 2.4 below).

28

29 **2.3 Ribosomal RNA Genes**

30 The organisation of the structural RNA genes in *E. histolytica* is unusual with
31 the rRNA genes carried exclusively on 24 kb circular episomes (Bhattacharya

1 *et al.*, 1998) that have two transcription units in an inverted repeat. These
2 episomes are believed to make up about 20% of the total cellular DNA; indeed,
3 roughly 15% of all of the sequencing reads generated in the genome project
4 were derived from this molecule with the exception of certain libraries where
5 attempts were made to exclude it. There are thought to be numerous other
6 circular DNA molecules of varying sizes present with unknown functions (Dhar
7 *et al.*, 1995; Lioutas *et al.*, 1995) but unfortunately they have not yet been
8 identified in the genome shotgun sequence data. The exact reasons for this are
9 unknown but the small size of the DNA may have prevented proper shearing
10 during the library construction process. These molecules represent an intriguing
11 unsolved aspect of the *E. histolytica* genome.

12

13 **2.4 tRNA Genes**

14 Perhaps the most unusual structural feature identified in the *E. histolytica*
15 genome is the unprecedented number and organisation of its tRNA genes (Clark
16 *et al.*, 2006a). Over 10% of the sequence reads contained tRNA genes and these
17 are (with a few exceptions) organised in linear arrays. The array organisation of
18 the tRNAs was immediately obvious in some cases from the presence of more
19 than one repeat unit in individual sequence reads and in other cases from their
20 presence in both reads from the two ends of the same clone. However because
21 of the near complete identity of the array units they were impossible to
22 assemble by the software used and therefore the size of the arrays cannot be
23 estimated accurately.

24

25 By manual assembly of tRNA gene-containing reads, 25 distinct arrays with
26 unit sizes ranging from under 500 bp to over 1750 bp were identified (Clark *et*
27 *al.*, 2006a). The arrayed genes are predicted to be functional because of the 42
28 acceptor types found in arrays none has been found elsewhere in the genome.
29 These array units encoded between one and five tRNAs and a few tRNA genes
30 are found in more than one unit. Three arrays also encode the 5S RNA and one
31 encodes what is thought to be a small nuclear RNA. Experimental quantitative

1 hybridisations suggest a copy number of between about 70 and 250 for various
2 array units. In total it is estimated that there are about 4500 tRNA genes in the
3 genome. The frequency of a particular tRNA isoacceptor appears to be
4 independent of the codon usage in *E. histolytica* protein-coding genes.

5

6 Between the genes in the array units are complex, non-coding, short tandem
7 repeats ranging in size from 5 to over 36 bp. Some variation in short tandem
8 repeat number is observed between copies of the same array unit but this
9 variation is usually minor and not visible when inter-tRNA PCR amplification
10 is performed. However, these regions often exhibit substantial variation when
11 different isolates of *E. histolytica* are compared and this is the basis of a
12 recently described genotyping method for this organism (Ali *et al.*, 2005).

13

14 There is indirect evidence to suggest that the tRNA arrays are present at the
15 ends of chromosomes. Although allelic *E. histolytica* chromosomes often differ
16 substantially in size in pulse-field gels, a central protein-encoding region
17 appears to be conserved as DNA digested with rare cutting enzymes gives only
18 a single band in Southern blots when most protein-coding genes are used as
19 probes. In contrast, when some tRNA arrays are used as probes on such blots,
20 the same number of bands is seen in digested and undigested DNA. It is
21 therefore tempting to conclude that the tRNA genes are at the ends of the
22 chromosomes and to speculate that these repeat units may perform a structural
23 role. In *D. discoideum* it is thought that rDNA may function as a telomere in
24 some cases (Eichinger *et al.*, 2005) and the tRNA arrays in *E. histolytica* may
25 perform a similar role.

26

27 The chromosomal regions flanking the tRNA arrays are generally devoid of
28 protein coding genes but often contain incomplete transposable elements (see
29 next section) and other repetitive sequences (Clark *et al.*, 2006a). This is also
30 consistent with a telomeric location.

31

1 2.5 LINES

2 The *E. histolytica* genome is littered with transposable elements. There are two
3 major types autonomous LINES (Long Interspersed Elements) of which there
4 are three subtypes (EhLINE 1, 2 and 3) and there are two types of SINEs (Short
5 Interspersed Elements) (Eh SINE1 and 2) (Table 2a). The classification of these
6 elements and their organisation has been reviewed recently (Bakre *et al.*, 2005).
7 Phylogenetic analysis of the EhLINES places them in the R4 clade of non- Long
8 Terminal Repeat (LTR) elements, a mixed clade of elements that includes
9 members from nematodes, insects, and vertebrates (Van Dellen *et al.*, 2002a).
10 Analysis of the *E. histolytica* genome shows no evidence for the presence of
11 LTR retrotransposons and very few DNA transposons (of the *Mutator* family)
12 (Pritham *et al.*, 2005).

13

14 All copies of EhLINES examined encode non-conservative amino acid changes,
15 frame shifts, and/or stop codons and no copy with a continuous open reading
16 frame (ORF) has yet been found. This suggests that the majority of these
17 elements are inactive. However, a large number of EhLINE1 copies do contain
18 long ORFs without mutations in the conserved protein motifs of the RT and EN
19 domains, suggesting that inactivity is quite recent. ESTs corresponding to
20 EhLINES have been found suggesting that transcription of these elements still
21 occurs. Although most R4 elements insert in a site-specific manner, EhLINES
22 do not show strict site-specificity and are widely dispersed in the genome. They
23 are quite frequently found close to protein-coding genes and inserted near T-
24 rich stretches (Bakre *et al.*, 2005).

25

26 All three EhLINE subtypes are of approximately equal size ranging from 4715
27 to 4811 bp in length. Individual members within an EhLINE family typically
28 share >85% identity, while between families they are <60% identical. By
29 aligning the available sequences, each EhLINE can be interpreted to encode a
30 single predicted ORF that spans almost the entire element (EhLINE1, 1589 aa;
31 EhLINE2, 1567 aa; EhLINE3, 1587 aa). However, a precise 5bp duplication at

1 nucleotide position 1442 in about 80% of the copies of EhLINE1 creates a stop
2 codon, dividing the single ORF in two. Similarly in 92% of EhLINE2 copies,
3 the single ORF contains a precise deletion of two nucleotides at position 1272,
4 resulting in two ORFs. Very few intact copies of EhLINE3 are found. The
5 location of the stop codon leading to two ORFs appears to be conserved since
6 in both EhLINE1 and EhLINE2 the size of ORF1 is about half that of ORF2
7 (Bakre *et al.*, 2005). Among the identifiable domains in the predicted proteins
8 are reverse transcriptase (RT) and a restriction enzyme-like endonuclease (EN).
9 The putative 5' and 3' untranslated regions are very short (3-44 bp).

10

11 EhLINEs 1 and 2 appear to be capable of mobilising partner SINEs (see next
12 section) for which abundant transcripts have been detected in *E. histolytica*.
13 Putative LINE/SINE partners can be assigned on the basis of conserved
14 sequences at the 3' -ends of certain pairs, which otherwise showed no sequence
15 similarity. The relevance of this assignment for the EhLINE1/SINE1 pair has
16 recently been demonstrated (Mandal *et al.*, 2004).

17

18 **2.6 SINEs**

19 The two EhSINEs are clearly related to the EhLINEs as they have a conserved
20 3' sequence. They are nonautonomous, non-LTR retrotransposons
21 (nonautonomous SINEs). The genetic elements encoding the abundant
22 polyadenylated but untranslatable transcripts found in *E. histolytica* cDNA
23 libraries (initially designated IE elements (Cruz-Reyes and Ackers, 1992;
24 Cruz-Reyes *et al.*, 1995) or *ehapt2* (Willhoeft *et al.*, 2002)) have now been
25 designated EhSINE1 (VanDellen *et al.*, 2002a; Willhoeft *et al.*, 2002). BLAST
26 searching with representative examples of the first 44 EhSINE1s detected has
27 identified 90 full-length (= 99% complete) copies and at least a further 120
28 partial (= 50% of full length) copies in the genome. Length variation is
29 observed among EhSINE1s and is largely due to variable numbers of internal
30 26-27 bp repeats (Ackers, unpublished). The majority contain two internal
31 repeats and cluster closely around 546 bp in length.

1
2 A second *E. histolytica* SINE (EhSINE2) has recently been described (Van
3 Dellen *et al.*, 2002a; Willhoeft *et al.*, 2002). Examination of the four published
4 sequences again suggests the presence of variable numbers of short (20 bp)
5 imperfect repeats. BLAST searching identified a total of 47 full-length (= 99%)
6 and at least 60 partial copies in the genome. The 3'-end of EhSINE2 shows high
7 similarity (76%) to the 3' end of EhLINE2.

8
9 A polyadenylated transcript designated UEE1 found commonly in cDNA
10 libraries from *E. dispar* (Sharma *et al.*, 1999) is also a non-LTR
11 retrotransposon. A single copy of a UEE1-like element has been identified in
12 the *E. histolytica* genome and is here designated EhSINE3. There is no
13 significant sequence identity between EhSINE3 and EhLINE3 but the 3' end of
14 EhSINE3 is very similar to that of EhLINE1.

15
16 Analysis of an *E. histolytica* EST library identified over 500 significant hits to
17 both EhSINE1 and EhSINE2. No convincing transcript from EhSINE3 could
18 be identified although the nearly identical *E. dispar* UEE elements (EdSINE1;
19 Shire and Ackers, submitted) are abundantly transcribed.

20
21 A very abundant polyadenylated transcript, *ehapt1*, was described by Willhoeft
22 *et al.* (1999) in a cDNA library. However, only a small number of partial
23 matches could be found in the current *E. histolytica* assembly and only 10-20
24 strong hits in the much larger *E. histolytica* EST library now available. *ehapt1*
25 does not appear to be a SINE element and its nature is currently unclear. The
26 lack of matches in the genome suggests either that it is encoded in regions
27 missing from the current assembly or that it contains numerous introns.

28

29 **2.7 Other Repeats**

30 The *E. histolytica* genome contains a number of other repetitive elements
31 whose functions are not always clear. There are over 75 genes encoding

1 leucine-rich tandem repeats (LRR) of the type found in BspA-like proteins of
2 the *Treponema pallidum* LRR (TpLRR) subfamily, which has a consensus
3 sequence of LxxIxIxxVxxIgxxAfxxCxx (Davis *et al.*, 2006). These proteins
4 generally have a surface location and may be involved in cell-cell interaction.
5 Genes encoding such proteins are found mainly in Bacteria and some Archaea;
6 so far they have been identified in only one other eukaryote, *Trichomonas*
7 *vaginalis* (Hirt *et al.*, 2002). An extensive description of the BspA-like proteins
8 of *E. histolytica* has recently been published (Davis *et al.*, 2006) and one of
9 them has been shown to be surface exposed (Davis *et al.*, 2006).

10

11 *E. histolytica* stress sensitive protein (Ehssp) 1 is a dispersed, polymorphic and
12 multicopy gene family (Satish *et al.*, 2003) and is present in ca. 300 copies per
13 haploid genome as determined by hybridisation (Table 2a). The average Ehssp1
14 ORF is 1 kb in length with a centrally-located acidic-basic region (ABR) that is
15 highly polymorphic. Unlike other such domains no clear repetitive motifs are
16 present. The protein has, on average, 21% acidic (aspartate and glutamate) and
17 17% basic (arginine and lysine) amino acids, most of which are located in the
18 ABR. The ABR varies in size from 5 to 104 amino acids among the various
19 copies. No size polymorphism is seen outside the central ABR domain. The
20 genes have an unusually long 5' untranslated region (UTR; 280 nucleotides).
21 Only one or a few copies of the gene are transcribed during normal growth, but
22 many are turned on under stress conditions. Homologues of this gene are
23 present in *E. dispar*, but there is very little size polymorphism in the *E. dispar*
24 gene family.

25

26 Eukaryotic genomes usually contain numerous microsatellite loci with repeat
27 sizes of 2-3 basepairs. With the exception of di- and tri-nucleotides made up
28 entirely of A+T such sequences are rare in the *E. histolytica* genome. In
29 contrast, two dispersed repeated sequences of unknown function occur far more
30 frequently than would be expected at random. Family 16 has a 42 base
31 consensus sequence and occurs approximately 38 times in the genome while

1 family 17 has a 27 base consensus sequence and occurs 35 time in the genome
2 (Table 2b). The significance of these sequences remains to be determined.

3

4 **2.8 Gene Number**

5 The current assembly predicts that the genome contains around 10,000 genes,
6 almost twice as many as seen in *P. falciparum* (Gardner *et al.*, 2002) or
7 *Saccharomyces cerevisiae* (Goffeau *et al.*, 1996) but closer to that of the free
8 living protist *Dictyostelium discoideum* (ca. 12,500; Eichinger *et al.*, 2005). It
9 should be remembered that this number will change as the assembly improves,
10 and is likely to decrease somewhat. Nevertheless, the comparatively large gene
11 number when compared to some other parasitic organisms reflects both the
12 relative complexity of *E. histolytica* and the presence of large gene families,
13 despite the loss of certain genes as a consequence of parasitism. Gene loss and
14 gain can both represent an adaptive response to life in the human host. Gene
15 loss is most evident in the reconstruction of metabolic pathways of *E.*
16 *histolytica*, which show a consistent pattern of loss of synthetic capacity as a
17 consequence of life in an environment rich in complex nutrient sources.
18 Similarly, analyses of expanded gene families with identifiable functions
19 indicate that many are directly associated with the ability to sense and adapt to
20 the environment within the human host and the ability to ingest and assimilate
21 the nutrients present. One consequence of these gene family expansions being
22 linked to phagocytosis of bacteria and other cells may be an association
23 between many of these gene families and pathogenicity.

24

25 **2.9 Gene Structure**

26 Most *E. histolytica* genes comprise only a single exon; however as many as
27 25% may be spliced and 6% contain two or more introns. Therefore mRNA
28 splicing is far less common than in the related protist *D. discoideum* or the
29 malaria parasite *P. falciparum*. The genome contains all of the essential
30 machinery for splicing (section 2.14) and a comparison of intron positions
31 suggests that *D. discoideum* and *E. histolytica* have both lost introns since their

1 shared common ancestor with *P. falciparum*, although many more have been
2 lost in the *E. histolytica* lineage. A good example of this intron loss is the
3 vacuolar ATP synthase subunit D gene (Figure 1). This protein is highly
4 conserved but the number of introns in each gene varies. *P. falciparum* has 5
5 introns, *D. discoideum* has two and *E. histolytica* has one. The positions of
6 three of the five *P. falciparum* introns are conserved in one of the other species
7 which suggests that these three (at least) were present in the common ancestor
8 and that intron loss has led to the lower number seen in *E. histolytica* today.
9 This loss is consistent with reverse transcriptase mediated 3' intron loss (Roy
10 and Gilbert, 2005) as the 5' -most introns are retained. It would appear that this
11 process has been more active in the *E. histolytica* and *D. discoideum* lineages
12 than in *P. falciparum*, possibly because *Plasmodium* lacks a reverse
13 transcriptase.

14

15 **2.10 Gene Size**

16 Genes in *E. histolytica* are surprisingly short, not only due to the loss of introns
17 but also in the predicted lengths of the proteins they code for. On average the
18 predicted length of a protein in *E. histolytica* is 389 amino acids (aa) which is
19 129 aa and 372 aa shorter than in *D. discoideum* and *P. falciparum* respectively.
20 In fact the protein length distribution is most similar to that of the
21 microsporidian *Encephalitozoon cuniculi* (Figure 2) which has a very compact
22 genome of 3Mb and less than 2000 genes. Direct comparison of orthologous
23 genes between *E. histolytica* and its closest sequenced relative *D. discoideum*
24 demonstrates this phenomenon quite well, with the majority of *E. histolytica*
25 proteins being shorter than the *D. discoideum* counterpart (Hall, unpublished).
26 Protein length is normally very well conserved among eukaryotes so the reason
27 for protein shortening is unclear. It has been postulated that in bacteria reduced
28 protein lengths reflects a reduced capacity for signaling (Zhang, 2000). This
29 would not seem to be the case here as the number of genes identified as having
30 a role in signaling suggests quite the opposite. An alternative theory is that as *E.*

1 *histolytica* has reduced organelles it is possible that its proteins contain fewer or
2 simpler targeting signals.

3

4 **2.11 Protein Domain Content**

5 The most common protein family (Pfam) domains of *E. histolytica* are shown in
6 Table 3. The domains that are unusually common in *E. histolytica* reflect some
7 of the more unusual aspects of the biology of this protist. For example, the Rab
8 and Rho families that are involved in signaling and vesicle trafficking are
9 among the most common domains in *E. histolytica* while in other species they
10 are not often among the top 50 families. This could well be due to the fact that
11 *E. histolytica* has a 'predatory' life style and these domains are intimately
12 involved in environmental sensing, endocytosis and delivery of lysosomes to
13 the phagosome. There are also a number of domains involved in actin
14 dynamics and cytoskeletal rearrangement that are not common in non-
15 phagocytic species, such as the gelsolin and SH3 domains. Myb domains are
16 the most common transcription regulatory domains in *E. histolytica*; this
17 domain is also common in plants where the proteins regulate many plant-
18 specific pathways (Ito, 2005). An important finding from an initial analysis
19 was the presence of unusual multidomain proteins, including five proteins
20 containing both RhoGEF and Arf-GAP domains, suggesting a mechanism for
21 direct communication between the regulators of vesicle budding and
22 cytoskeletal rearrangement. Over 80 receptor kinases were identified (section
23 7.2.2), each containing a kinase domain and a C rich extracellular domain.
24 These kinases fall into distinct classes depending on the presence of CXC or
25 CXXC repeats. There are also domains that are common in most other
26 sequenced genomes but rare or missing from *E. histolytica*. For example, most
27 mitochondrial carrier domain proteins are not needed in *E. histolytica* as it lacks
28 a normal mitochondrion (section 8).

29

30 **2.12 Translation-Related Proteins**

1 Two of the predicted tRNAs (Ile^{TAT} and Tyr) need to be spliced due to the
2 presence of an intron. tRNA introns are distinct in structure from those in
3 protein-coding genes and require a distinct splicing machinery. The expected
4 enzymes required for this splicing are present as are a number of tRNA
5 modification enzymes (including those for synthesising queuine and
6 pseudouridine) and rRNA methylases that act on specific bases in their
7 respective RNA molecules. The expected panel of tRNA synthetases necessary
8 for aminoacylating the tRNAs is also present, with one or two gene copies for
9 each type.

10

11 The majority of ribosomal protein genes are well-conserved in *E. histolytica*
12 and only the gene for large subunit protein L41 could not be identified. The
13 missing protein is only 25 amino acids in length, 17 of which are arginines or
14 lysines, which would make it difficult to identify in this A+T-rich genome, but
15 it is highly conserved, having been reported from Archaea to mammals.
16 However, it also appears to be dispensable, as *S. cerevisiae* can grow relatively
17 normally after deletion of both its copies (Yu and Warner, 2001). Nevertheless,
18 deletion of L41 in *S. cerevisiae* reduces the level of 80S ribosomes, suggesting
19 that it is involved in ribosomal subunit association, reduces peptidyl transferase
20 activity, and increases translocation (Dresios *et al.*, 2003). In addition, L41 has
21 been shown to interact with the beta subunit of protein kinase CKII and to
22 stimulate phosphorylation of DNA topoisomerase II alpha by CKII (Lee *et al.*,
23 1997b). If this gene is truly absent from *E. histolytica* it may have important
24 consequences for the cell.

25

26 No genes for mitochondrial ribosomal proteins were found. Their absence is not
27 surprising since *E. histolytica* lacks typical mitochondria (see section 8 below).

28

29 In eukaryotic translation, elongation factor EF-1 is activated upon GTP binding
30 and forms a ternary complex with aminoacyl tRNAs and ribosomes. EF-1 beta
31 and delta subunits work as GDP-GTP exchange factors to cycle EF-1 alpha

1 between two forms while EF-1 gamma provides structural support for the
2 formation of this multimeric complex. EF2 assists in the translocation of tRNAs
3 on the mRNA by exactly one codon. *E. histolytica* has most of the expected
4 factors except for EF-1 delta, a protein involved in exchanging GDP with GTP.
5 This is also absent from *S. cerevisiae* and *P. falciparum*. It is likely that EF-1
6 beta carries out this activity. It is thought that the EF-1 complex can exist in
7 two forms, EF-1-alpha/beta/gamma and EF-1-alpha/delta/gamma. In *E.*
8 *histolytica*, probably only the former complex exists.

9
10 Eukaryotes typically have two polypeptide release factors, eRF1 and eRF3.
11 Both of these factors have been found in *E. histolytica*.

12 13 **2.13 Analysis of Cell Cycle Genes**

14 Alternation of DNA duplication and chromosome segregation is a hallmark in
15 the cell cycle of most eukaryotes. Carefully orchestrated processes coordinate
16 an ensemble of cell cycle regulating 'checkpoint' proteins ensure that progeny
17 cells receive an exact copy of the parental genetic material (Hartwell and
18 Weinert, 1989). Unlike most eukaryotes, *Entamoeba histolytica* cells can
19 reduplicate their genome several times before cell division occurs
20 (Gangopadhyay *et al.*, 1997). Approximately 5-20% of the trophozoites
21 (depending on the growth phase) of axenic culture are multi-nucleated.
22 Additionally, DNA reduplication may occur without nuclear division so that
23 single nuclei contain 1X -6X or more genome contents (Das and Lohia, 2002).
24 Thus axenically cultured *E. histolytica* trophozoites display heterogeneity in
25 their genome content suggesting that eukaryotic cell cycle checkpoints are
26 either absent or altered in this organism. Around 200 genes have been identified
27 in yeast that play a direct role in cell cycle progression.

28 29 *2.13.1 DNA replication initiation and DNA duplication*

30 The DNA replication licensing system is one of the crucial mechanisms that
31 ensures the alternation of S-phase with mitosis in most cells (Tye, 1999).

1 Initiation of DNA replication involves binding of the replicative helicases to
2 DNA replication origins in late mitosis. Loading of the replicative helicase
3 Mcm2-7 proteins is preceded by formation of the pre-replicative complex (pre-
4 RC) and its subsequent activation. Formation of pre-RC requires the ordered
5 assembly of the origin recognition complex (ORC), Cdc6, Cdt1 and the Mcm2-
6 7 proteins. The pre-RC is activated by the protein kinase Cdc7p and its
7 regulatory subunit Dbf4 (Masai and Arai, 2002). Other factors that regulate the
8 transition from pre-RC to replication initiation are Mcm10p, Cdc45p, TopBP1,
9 RecQL4 and the GINS complex (Gregan *et al.*, 2003; Machida *et al.*, 2005;
10 Merchant *et al.*, 1997; Wohlschlegel *et al.*, 2002). Two other Mcm proteins –
11 Mcm8 and Mcm9 - have been identified in metazoan systems and are believed
12 to be part of the replicative helicase (Maiorano *et al.*, 2006). Replication origin
13 licensing is inactivated during S-phase but Mcm2-9p may function as a
14 helicase that unwinds DNA ahead of the replication fork during S-phase
15 (Maiorano *et al.*, 2006). Once S-phase has begun, the formation of new pre-RC
16 is kept in check by high CDK activity and by the activity of the protein geminin
17 Bell and Dutta, 2002).

18

19 A detailed analysis of the *E. histolytica* genome shows that homologues of
20 several proteins required for DNA replication initiation are absent. These
21 include ORC (Origin Recognition Complex) 2-6, Cdt1, geminin, Cdc7/Dbf4
22 and Mcm10. A single gene encoding a homologue of the archaeal and human
23 Cdc6/Orc1p (Capaldi and Berger, 2004) was identified. This suggests that DNA
24 replication initiation in *E. histolytica* is likely similar to archaeal replication
25 initiation where a single Cdc6p/ORC1p replaces the hetero-hexameric ORC
26 complex (Kelman and Kelman, 2004). Several proteins described from
27 metazoa, such as Cdt1, geminin, Mcm8 and Mcm9, have not been found in
28 yeast. Surprisingly, Mcm8 and Mcm9 were identified in the *E. histolytica*
29 genome.

30

1 Of the four known checkpoint genes that regulate DNA replication in *S.*
2 *cerevisiae* only Mec1 and Mrc1 have homologues in *E. histolytica*. *E.*
3 *histolytica* homologues of several proteins involved in G1-S transitions are
4 absent, such as Sic1, Chk1. The S-phase checkpoint genes p21, p27, p53 and
5 retinoblastoma (RB) required for transition from G1 to S-phase in humans were
6 absent in *E. histolytica*. Chk1 and Chk2 genes encode kinases that act
7 downstream from the ATM and ATR kinases (intra-S phase checkpoint genes).
8 The Chk1 homologue is absent but a Chk2 homologue has been identified in *E.*
9 *histolytica* and partially characterised (Iwashita *et al.*, 2005).

10

11 2.13.2 Chromosome segregation and cell division

12 A large number of genes are known to regulate different events during the
13 transition from G2-Mitosis - spindle formation checkpoint, chromosome
14 segregation, mitosis, exit from mitosis, and cytokinesis - in *S.cerevisiae*. Many
15 of the proteins required by yeast for kinetochore formation have no obvious
16 homologues in *E. histolytica* suggesting that amoeba kinetochores may have an
17 altered composition and structure. Proteins of the Anaphase Promoting
18 Complex (APC) regulate transition from metaphase to anaphase. With the
19 exception of APC11, none of the APC proteins could be identified in
20 *E.histolytica*. In contrast two genes encoding CDC20 homologues, which are
21 known to activate the APC complex, were identified in *E.histolytica* along with
22 ubiquitin and related proteins (Wöstmann *et al.*, 1992), indicating that although
23 most APC subunit homologues were absent the pathway of proteasomal
24 degradation for regulation of cell cycle proteins may still be functional in
25 *E.histolytica*. Effectors of the apoptotic pathway and meiosis were also largely
26 absent.

27

28 2.13.3 CDKs and cyclins

29 The CDC28 gene encodes the single cyclin dependant kinase (CDK) in *S.*
30 *cerevisiae* and regulates cell cycle progression by binding to different cyclins at
31 the G1/S or G2/M boundaries (Reed, 1992; Surana *et al.*, 1991; Wittenberg *et*

1 *al.*, 1990). Similarly, *Schizosaccharomyces pombe* also encodes a single CDK
2 (*cdc2*) (Simanis and Nurse, 1986). Mammals and plants can encode multiple
3 CDKs and an equally large number of cyclins (Morgan, 1995; Vandepoele *et*
4 *al.*, 2002). Association of different CDKs with specific cyclins regulates the
5 cell cycle in different developmental stages as well as in specific tissues. CDKs
6 belong to the serine/threonine family of kinases with a conserved PSTAIRE
7 domain where cyclins are believed to bind (Jeffrey *et al.*, 1995; Morgan, 1996)
8 although some mammalian and plant CDKs have been shown to have divergent
9 PSTAIRE motifs. This heterogeneity may or may not affect cyclin binding
10 (Poon *et al.*, 1997). The *E. histolytica* genome encodes at least 9 different
11 CDKs among which not even one has the conserved PSTAIRE motif. The
12 closest homologue of the CDC28/*cdc2* gene, which shows only conservative
13 substitutions in the PSTAIRE motif (PVSTVRE), was cloned previously (Lohia
14 and Samuelson, 1993). The remaining 8 CDK homologues exhibit even greater
15 divergence in this motif. Eleven putative cyclin homologues with a high degree
16 of divergence have been found. Identifying their CDK/cyclin partner along with
17 their roles in the cell cycle is a major task that lies ahead. Some of the CDKs
18 may not function by associating with their functional cyclin partners but may
19 play a role in regulating global gene expression, either by activation from non-
20 cyclin proteins or by other mechanisms (Nebreda, 2006).

21

22 *E. histolytica* presents a novel situation where the eukaryotic paradigm of a
23 strictly alternating S-phase and mitosis is absent. Discrete G1, S and G2
24 populations of cells are not routinely found in axenic cultures. Instead cells in
25 S-phase show greater than 2x genome contents, suggesting that the G2 phase is
26 extremely short and irregular. This observation together with the absence of a
27 large number of checkpoint genes suggests that regulation of genome
28 partitioning and cell division in *E. histolytica* may be additionally dependant on
29 extracellular signals. *E. histolytica* must however contain regulatory
30 mechanisms to ensure that its genome is maintained and transmitted with
31 precision even in the absence of the expected checkpoint controls. The

1 discovery of these mechanisms will be crucial to our understanding of how the
2 *E. histolytica* cell divides.

3

4 **2.14 Transcription**

5 RNA polymerase II transcription in *E. histolytica* is known to be α -amanitin-
6 resistant (Lioutas and Tannich, 1995). The F homology block of the RNA
7 polymerase II large subunit has been identified as the putative α -amanitin
8 binding site. This block is highly divergent in the α -amanitin resistant
9 *Trichomonas vaginalis* RNA polymerase II (Quon *et al.*, 1996). The *E.*
10 *histolytica* RPB1 homologue also diverges from the consensus in this region
11 but, interestingly, it is also quite dissimilar to the *T. vaginalis* sequence.

12

13 The heptapeptide repeat (TSPTSPS) common to other eukaryotic RNA
14 polymerase II large subunit C terminal domains (CTD) is not present in the *E.*
15 *histolytica* protein. Indeed, the *E. histolytica* CTD is not similar to any other
16 RNA polymerase II domain in the current database. However, the CTD of the
17 *E. histolytica* enzyme does remain proline/serine-rich (these amino acids
18 constitute 40% of the CTD sequence). The *E. histolytica* CTD also retains the
19 potential to be highly phosphorylated: of the 24 serines, 6 threonines and 3
20 tyrosines within the CTD, 9 serines, 3 threonines and 1 tyrosine are predicted to
21 be within potential phosphorylation sites. It is therefore possible that, despite
22 its divergence, modification of the CTD by kinases and phosphatases could
23 modulate protein-protein interactions as is postulated to occur in other RNA
24 polymerases (Yeo *et al.*, 2003). In yeast, phosphorylation of the CTD regulates
25 association with the mediator protein (Davis *et al.*, 2002; Kang *et al.*, 2001;
26 Komberg, 2001). The yeast mediator protein complex consists of 20 subunits.
27 However, perhaps due to the divergence of the CTD, only two of these proteins
28 have been identified in *E. histolytica* (Med7 and Med10). Homologues of the
29 Spt4 and Spt5 elongation factors, also thought to interact with the CTD, have
30 been identified.

31

1 The RNA polymerase core is composed of 12 putative subunits in *S. cerevisiae*
2 (Young, 1991), while *S. pombe* contains a subset of 10 of these proteins,
3 lacking the equivalents of subunits 4 and 9 (Yasui *et al.*, 1998). In *E.*
4 *histolytica* only 10 of the RNA polymerase subunits have been identified,
5 identifiable homologues of subunits 4 and 12 being absent. While the
6 homologue to subunit 9 was present it lacks the first of the two characteristic
7 zinc binding motifs of this protein and the DPTLPR motif in the C terminal
8 region. A similar sequence, DPTYPK, is however present and a homologue of
9 the TFIIE large subunit Tfa1, which is proposed to interact with this region of
10 the protein, has been identified (Hemming and Edwards, 2000; Van Mullem *et*
11 *al.*, 2002). The conserved N terminal portion (residues 1-52) of Rpb9 is
12 thought to interact with both Rpb1 and Rpb2 in *S. cerevisiae* (Hemming and
13 Edwards, 2000) and homologues of these have been identified.

14

15 The core promoter of *E. histolytica* has an unusual tripartite structure consisting
16 of the three conserved elements TATA, GAAC and INR (Purdy *et al.*, 1996;
17 Singh and Rogers, 1998; Singh *et al.*, 2002; Singh *et al.*, 1997). Singh and
18 Rogers (1998) have speculated that the GAAC motif may be the binding site of
19 a second or alternative *E. histolytica* DNA binding protein in the preinitiation
20 complex. It is therefore of interest that, in addition to the *E. histolytica* TATA-
21 binding protein (TBP), two other proteins contain the TATA-binding motif
22 (Hernandez *et al.*, 1997). TBP is a subunit of the TFIID general transcription
23 factor (GTF) which in other organisms is required for the recognition of the
24 core promoter. In light of the variation in the core promoter previously
25 mentioned, and the divergence in proteins that bind to the core promoter in
26 other parasitic protists, it is not surprising that only six of the 14 evolutionary
27 conserved subunits of TFIID, TBP Associated Factors (TAFs) 1, 5, 6, 10, 12
28 and 13 were identified. Homologues of some of the global regulatory subunits
29 of the Ccr4/Not complex, which interacts with TBP and TAFs 1 and 13, have
30 also been identified.

31

1 TAFs 5, 6, 10 and 12 are also components of the histone acetyltransferase
2 (HAT) complexes in other organisms as is SPT6 and 16 (Carrozza *et al.*, 2003).
3 While all known components of the HAT complexes have by no means been
4 identified or the role of the previously unknown bromodomain containing
5 proteins encoded in the *E. histolytica* genome, histone acetylation complexes
6 are known to be active in *E. histolytica* (Ramakrishnan *et al.*, 2004). Other
7 potential members of chromatin remodeling complexes of *E. histolytica* include
8 the TBP interacting helicase (RVB1 & 2) and the SNF2 subunit of the
9 SWI/SNF complex.

10

11 Homologues of some of the other GTFs (TFII E, F and H) but not the large or
12 small subunits were identified. In contrast to the difficulty identifying some of
13 the GTFs, the *E. histolytica* spliceosomes components U1, U2, U4/6, U5 and
14 the Prp19 complex have all been identified. In fact homologues of ten of the
15 fourteen “core” snRNP proteins, two of the U1 specific snRNPs, seven of the
16 ten U2 specific snRNPs, five of the six U5 specific snRNPs, three of the U4/6
17 specific snRNPs, and four of the nine subunits of the Prp19 complex have been
18 found. In fact *E. histolytica* has homologues of approximately 80% of the *S.*
19 *cerevisiae* splicing machinery (Jurica and Moore, 2003).

20

21 Like *G. intestinalis*, *E. histolytica* has short 5' untranslated regions on its
22 mRNAs. However, unlike those of *G. intestinalis*, *E. histolytica* mRNA has
23 been shown to be capped (Ramos *et al.*, 1997; Vanacová *et al.*, 2003).
24 Identification of homologues of the Ceg1 RNA guanylyltransferase - an
25 enzyme which adds an unmethylated GpppRNA cap to new transcripts - and of
26 Abd1 - which methylates the cap to form m7GpppRNA - gives new insight into
27 the probable cap structure in *E. histolytica* (Hausmann *et al.*, 2001; Pillutla *et*
28 *al.*, 1998). It has been proposed that the capping enzymes interact with the
29 phosphorylated CTD of RNA polymerase (Schroeder *et al.*, 2000). The CTD of
30 *E. histolytica* large subunit is, as discussed earlier, not well conserved but
31 contains several probable phosphorylation sites.