

The genome of the social amoeba *Dictyostelium discoideum*

L. Eichinger^{1*}, J. A. Pachebat^{1,2*}, G. Glöckner^{3*}, M.-A. Rajandream^{4*}, R. Sugang^{5*}, M. Berriman⁴, J. Song⁵, R. Olsen⁹, K. Szafranski³, Q. Xu^{6,7}, B. Tunggal¹, S. Kummerfeld², M. Madera², B. A. Konfortov², F. Rivero¹, A. T. Bankier², R. Lehmann³, N. Hamlin⁴, R. Davies⁴, P. Gaudet¹⁰, P. Fey¹⁰, K. Pilcher¹⁰, G. Chen⁵, D. Saunders⁴, E. Sodergren^{6,8}, P. Davis⁴, A. Kerhornou⁴, X. Nie⁵, N. Hall^{4,†}, C. Anjard⁹, L. Hemphill⁵, N. Bason⁴, P. Farbrother¹, B. Desany⁵, E. Just¹⁰, T. Morio¹¹, R. Rost¹², C. Churcher⁴, J. Cooper⁴, S. Haydock¹³, N. van Driessche⁶, A. Cronin⁴, I. Goodhead⁴, D. Muzny⁸, T. Mourier⁴, A. Pain⁴, M. Lu⁵, D. Harper⁴, R. Lindsay⁵, H. Hauser⁴, K. James⁴, M. Quiles⁸, M. Madan Babu², T. Saito¹⁴, C. Buchrieser¹⁵, A. Wardroper^{2,16}, M. Felder³, M. Thangavelu¹⁷, D. Johnson⁴, A. Knights⁴, H. Louseged⁸, K. Mungall⁴, K. Oliver⁴, C. Price⁴, M. A. Quail⁴, H. Urushihara¹¹, J. Hernandez⁸, E. Rabinowitsch⁴, D. Steffen⁸, M. Sanders⁴, J. Ma⁵, Y. Kohara¹⁸, S. Sharp⁴, M. Simmonds⁴, S. Spiegler⁴, A. Tivey⁴, S. Sugano¹⁹, B. White⁴, D. Walker⁴, J. Woodward⁴, T. Winckler²⁰, Y. Tanaka¹¹, G. Shaulsky^{6,7}, M. Schleicher¹², G. Weinstock^{6,8}, A. Rosenthal³, E. C. Cox²¹, R. L. Chisholm¹⁰, R. Gibbs^{6,8}, W. F. Loomis⁹, M. Platzer³, R. R. Kay², J. Williams²², P. H. Dear², A. A. Noegel¹, B. Barrell⁴ & A. Kuspa^{5,6}

¹Center for Biochemistry and Center for Molecular Medicine Cologne, University of Cologne, Joseph-Stelzmann-Str. 52, 50931 Cologne, Germany

²Laboratory of Molecular Biology, MRC Centre, Cambridge CB2 2QH, UK

³Genome Analysis, Institute for Molecular Biotechnology, Beutenbergstr. 11, D-07745 Jena, Germany

⁴The Wellcome Trust Sanger Institute, Wellcome Trust Genome Campus, Hinxton, Cambridgeshire CB10 1SA, UK

⁵Verna and Marrs McLean Department of Biochemistry and Molecular Biology, ⁶Department of Molecular and Human Genetics, ⁷Graduate Program in Structural and Computational Biology and Molecular Biophysics, and ⁸Human Genome Sequencing Center, Baylor College of Medicine, Houston, Texas 77030, USA

⁹Section of Cell and Developmental Biology, Division of Biology, University of California, San Diego, La Jolla, California 92093, USA

¹⁰dictyBase, Center for Genetic Medicine, Northwestern University, 303 E Chicago Ave, Chicago, Illinois 60611, USA

¹¹Graduate School of Life and Environmental Sciences, University of Tsukuba, Tsukuba, Ibaraki 305-8572, Japan

¹²Adolf-Butenandt-Institute/Cell Biology, Ludwig-Maximilians-University, 80336 Munich, Germany

¹³Biochemistry Department, University of Cambridge, Cambridge CB2 1QW, UK

¹⁴Division of Biological Sciences, Graduate School of Science, Hokkaido University, Sapporo 060-0810, Japan

¹⁵Unité de Genomique des Microorganismes Pathogènes, Institut Pasteur, 28 rue du Dr Roux, 75724 Paris Cedex 15, France

¹⁶Department of Biology, University of York, York YO10 5YW, UK

¹⁷MRC Cancer Cell Unit, Hutchison/MRC Research Centre, Hills Road, Cambridge CB2 2XZ, UK

¹⁸Centre for Genetic Resource Information, National Institute of Genetics, Mishima, Shizuoka 411-8540, Japan

¹⁹Department of Medical Genome Sciences, Graduate School of Frontier Sciences, The University of Tokyo, Minato, Tokyo 108-8639, Japan

²⁰Institut für Pharmazeutische Biologie, Universität Frankfurt (Biozentrum), Frankfurt am Main 60439, Germany

²¹Department of Molecular Biology, Princeton University, Princeton, New Jersey 08544-1003, USA

²²School of Life Sciences, University of Dundee, Dow Street, Dundee DD1 5EH, UK

* These authors contributed equally to this work

† Present address: The Institute for Genomic Research, 9712 Medical Center Drive, Rockville, Maryland 20850, USA

The social amoebae are exceptional in their ability to alternate between unicellular and multicellular forms. Here we describe the genome of the best-studied member of this group, *Dictyostelium discoideum*. The gene-dense chromosomes of this organism encode approximately 12,500 predicted proteins, a high proportion of which have long, repetitive amino acid tracts. There are many genes for polyketide synthases and ABC transporters, suggesting an extensive secondary metabolism for producing and exporting small molecules. The genome is rich in complex repeats, one class of which is clustered and may serve as centromeres. Partial copies of the extrachromosomal ribosomal DNA (rDNA) element are found at the ends of each chromosome, suggesting a novel telomere structure and the use of a common mechanism to maintain both the rDNA and chromosomal termini. A proteome-based phylogeny shows that the amoebozoa diverged from the animal–fungal lineage after the plant–animal split, but *Dictyostelium* seems to have retained more of the diversity of the ancestral genome than have plants, animals or fungi.

The amoebozoa are a richly diverse group of organisms whose genomes remain largely unexplored. The soil-dwelling social amoeba *Dictyostelium discoideum* has been actively studied for the past 50 years and has contributed greatly to our understanding of cellular motility, signalling and interaction¹. For example, studies in *Dictyostelium* provided the first descriptions of a eukaryotic cell chemoattractant and a cell–cell adhesion protein^{2,3}.

Dictyostelium amoebae inhabit forest soil and consume bacteria and yeast, which they track by chemotaxis. Starvation, however, prompts the solitary cells to aggregate and develop as a true multicellular organism, producing a fruiting body comprised of a cellular, cellulosic stalk supporting a bolus of spores. Thus, *Dictyostelium* has evolved mechanisms that direct the differentiation of a homogeneous population of cells into distinct cell types, regulate

the proportions between tissues and orchestrate the construction of an effective structure for the dispersal of spores⁴. Many of the genes necessary for these processes in *Dictyostelium* were also inherited by Metazoa and fashioned through evolution for use within many different modes of development.

The amoebozoa are also noteworthy as representing one of the earliest branches from the last common ancestor of all eukaryotes. Each of the surviving branches of the crown group of eukaryotes provides an example of the ways in which the ancestral genome has been sculpted and adapted by lineage-specific gene duplication, divergence and deletion. Comparison between representatives of these branches promises to shed light not only on the nature and content of the ancestral eukaryotic genome, but on the diversity of ways in which its components have been adapted to meet the needs

of complex organisms. The genome of *Dictyostelium*, as the first free-living protozoan to be fully sequenced, should be particularly informative for these analyses.

Mapping, sequencing and assembly

An international initiative to sequence the genome of *Dictyostelium discoideum* AX4 (refs 5, 6) was launched in 1998. The high repeat content and (A+T)-richness of the genome (the latter rendering large-insert bacterial clones unstable) posed severe challenges for sequencing and assembly. The response to these challenges was to use a whole-chromosome shotgun (WCS) strategy, partially purifying each chromosome electrophoretically and treating it as a separate project. This approach was supported by novel statistical tools to recover chromosome specificity from the impure WCS libraries, and by highly detailed HAPPY maps that provided a framework for sequence assembly. These approaches have enabled the completion of this difficult genome to a high standard, and are likely to be valuable in tackling the many other genomes that present challenges of composition and complexity.

Genome mapping

To support sequence assembly, we made high-resolution maps of the chromosomes using HAPPY mapping^{7–9}, which relies on analysing the sequence content of single DNA molecules prepared by limiting dilution. A total of 3,902 markers selected mostly from the emerging shotgun data were mapped, and maps of all six chromosomes were assembled (see Methods and Table 1; see also Supplementary Fig. 1 and Supplementary Table 1).

Genome sequencing and assembly

Two strategies were used to recover chromosome-specific data from impure WCS libraries (see Methods). The first (for chromosomes 1, 2 and 3) used enrichment of the respective libraries as the main statistical indicator of the chromosomal assignment of contigs, and HAPPY maps were used to guide assembly. The second strategy (for chromosomes 4, 5 and most of 6) used mapping data to assign sequences to chromosomes initially, with detailed HAPPY maps being used to validate final assemblies. A 1,508-kilobase (kb) portion of chromosome 6 was sequenced as a pilot project using a combination of approaches (see Methods).

Repetitive tracts complicated assembly. For chromosomes 1, 2 and 3, inspection of polymorphisms, combined with HAPPY maps, allowed unambiguous assembly in many cases. For chromosomes 4,

5 and 6, low-coverage sequencing of AX4-derived yeast artificial chromosomes (YACs) alleviated the problems by providing a local data set within which the troublesome repeat element was present as a single copy. Nevertheless, some repeat tracts proved intractable and remain as gaps. Thirty-four unlinked (floating) contigs of >1 kb, totalling 225,339 base pairs (bp), remain unpositioned in the genome, but can be provisionally assigned to specific chromosomes based on their content of reads from the WCS libraries. Most or all of these floating contigs are bounded by repetitive regions. The chromosome 2 sequence in the current assembly supersedes that previously published⁹, having benefited from further HAPPY mapping and manual sequence finishing.

The six chromosomal assemblies span 33,817 kb (Table 1), including ~156 kb in the form of clone-, sequence- and repeat gaps. Assuming that most of the floating contigs lie beyond the termini of the assemblies, the total genome size is estimated at 34,042,810 bp. In estimating the completeness of the sequence, we note that of 967 well-characterized *D. discoideum* genes, 957 (99%) were found initially in the assemblies. Of the remaining ten, seven (*cupE*, *trxA*, *trxB*, *trxC*, *staA*, *staB* and *cinB*) have close matches, suggesting that their GenBank entries may contain errors or represent alternative alleles. Only three (*fcpA*, *wasA* and *roco5*) had no matches in the initial assemblies, although the first two of these were recovered by searches of unincorporated sequence followed by local reassembly. Of 133,168 'qualified' *D. discoideum* AX4 expressed sequence tags (ESTs of >200 bp and >20% G+C, and not matching mitochondrial sequence; ref. 10 and H. Urushihara *et al.*, unpublished data), 128,207 (96.3%) are found in the assemblies (the higher proportion of missing sequences among the ESTs probably reflects the higher error rate inherent in EST data).

We conclude that the current assembly represents >>95% of the chromosomal sequence (less than 1% of which is in floating contigs) and >99% of genes, with most of the missing sequence comprising complex or simple repeats. The most stringent test of the medium-to long-range accuracy of the assembly comes from comparison with the HAPPY maps. This is particularly true for chromosomes 4, 5 and 6, where HAPPY markers were used to nucleate contigs but not to guide their assembly or ordering, specifically to allow such a comparison to be made without circularity of argument. As can be seen, good agreement between map and sequence confirms the accuracy of the assembly (Fig. 1).

Table 1 Sequence assembly details

Feature	Chromosome						All
	1	2	3	4	5	6	
Chromosomal assemblies							
Assembly span (bp)*	4,919,822	8,467,571	6,358,352	5,430,575	5,062,323	3,578,828	33,817,471
Assembly sequence (bp)†	4,911,622	8,437,971	6,334,852	5,397,875	5,032,273	3,547,128	33,661,721
Total contigs	11	40	32	65	107	44	309
Mean contig size (bp)	446,511	210,949	197,964	83,044	47,031	80,617	108,938
Number of sequence gaps	4	12	10	34	81	14	155
Number of repeat gaps	8	29	23	9	4	11	84
Number of clone gaps	0	0	0	22	22	20	64
Total estimated gap size (bp)‡	8,200	29,600	23,500	32,700	30,050	31,700	155,750
Number of HAPPY markers (mean spacing in kb)	749 (6.6)	615 (12.5)§	684 (9.3)	628 (8.6)	628 (8.1)	598 (6.0)	3,902 (8.7)
Floating contigs							
Number of floating contigs	0	22	3	9¶		0	34
Total size of floating contigs (bp)	0	171,670	16,360	37,309		0	225,339
Combined (assemblies plus floating contigs)							
Total sequence (bp)	4,911,622	8,609,641	6,351,212	5,416,529¶	5,050,928¶	3,547,128	33,887,060
Mean coverage (fold)	9.1	6.5	6.7	9.6	9.9	10.3	8.3

*Total end-to-end length of the chromosomal assembly, including any gaps.

†Sequenced bases covered by chromosomal assembly, not counting gaps.

‡Sequence, repeat and clone gaps are taken to have average sizes of 50 bp, 1,000 bp and 1,000 bp, respectively.

§Does not include the second copy of the 755-kb inverted duplication.

||Includes only those contigs that can be assigned to specific chromosomes.

¶Floating contigs from chromosomes 4 and 5 cannot be distinguished. In calculating total chromosomal sequence, we assume that half of these floating contigs are from each of chromosomes 4 and 5.

Sequence characteristics of the genome

The genome is (A+T)-rich (77.57%) and has a broadly uniform composition, apart from the more (G+C)-rich repeat-dense regions (Fig. 2). On a finer scale, nucleotide composition tracks the distribution of exons (see below). Among dinucleotides, CpG is under-represented, not just in absolute terms but also relative to its isomer GpC (the former occurring only 62% as often as the latter). This bias normally reflects cytosine methylation at CpG sequences, promoting their mutation to TpG (which is over-represented relative to GpT by 38%). Hence, these observations suggest that cytosine methylation may occur in *Dictyostelium*, contrary to earlier findings¹¹.

Simple sequence repeats are abundant and unusual

Simple sequence repeats (SSRs) are more abundant in *Dictyostelium* than in any other genome sequenced so far, comprising >11% of bases (Supplementary Fig. 2). In non-coding sequence, tracts of dinucleotides or longer motifs occur every 392 bp on average and comprise 6.4% of the bases. There is a bias towards repeat units of 3–6 bases, whereas dinucleotide tracts predominate in most other genomes. Homopolymer tracts are also abundant, comprising a further 16% of non-coding sequence. The base composition of non-coding SSRs and homopolymer tracts (99.2% A+T content) is even more biased than that of the surrounding sequence, suggesting that either selection or the mechanism of repeat expansion favours (A+T)-rich repeats.

Notably, SSRs are also abundant in protein-coding sequence, occurring on average every 724 bp within exons. We consider these coding SSRs in further detail below, in the context of proteins.

Transposable elements are clustered

The genome is rich in transposable elements^{9,12}. Completion of the sequence confirms the earlier observation that transposable elements of the same type are clustered, suggesting their preferential insertion within similar resident elements. However, none of the elements appears to use a specific sequence as a target for insertion: they insert at random within other elements of the same type. Non-long terminal repeat (LTR) retrotransposons are known to insert next to transfer RNA genes; we find many such instances

(Fig. 2), but again no specific sequences were identified as insertion targets.

tRNAs are numerous and paired by specificity

The sequenced genome encodes 390 tRNAs, a number at the upper end of the eukaryotic spectrum (for example, *Plasmodium falciparum* = 43, *Drosophila melanogaster* = 284, *Homo sapiens* = 496). Allowing for the normal wobble rules in codon–anticodon pairing^{13,14}, every sense codon can be decoded, apart from the rare alanine codon GCG; we infer that the missing tRNA(s) lie in one or more gaps in the sequence. We also find a possible selenocysteine tRNA in the genome, as well as corresponding selenocysteine insertion targets in two predicted proteins (see Supplementary Fig. 3).

Dictyostelium, in common only with *Acanthamoeba castellanii*¹⁵, has been shown to lack certain apparently essential tRNAs in its mitochondrial genome¹⁶. It therefore seems likely that at least some chromosomally encoded tRNAs (those for valine, threonine, asparagine and glycine, as well as one arginine and two serine tRNAs) are imported into mitochondria.

Although the gross distribution of tRNAs is uniform, organization of tRNAs on a finer scale is striking: about 20% occur as pairs or triplets with identical anticodons (and usually 100% sequence identity), separated by <20 kb and often by <5 kb (Fig. 2). There are 41 such groups in the genome; a random distribution would produce few, if any. This pattern is unique among sequenced genomes, and suggests a wave of recent duplications. However, tRNA pairs are found in tandem, converging and diverging orientations with comparable frequencies, suggesting no straightforward duplication mechanism; nor is there usually duplication of extensive flanking sequences. Whether the preference of TRE elements for inserting adjacent to tRNAs is related to the large number and unusual distribution of tRNAs is unclear.

A chromosomal master copy of the extrachromosomal rDNA element

In *Dictyostelium*, ribosomal RNA genes lie on an 88-kb palindromic extrachromosomal element¹⁷, present at ~100 copies per nucleus (Fig. 2). Evidence also exists of chromosomal copies: at least the central 3.2 kb of the element is located¹⁷ on chromosome 4, whereas

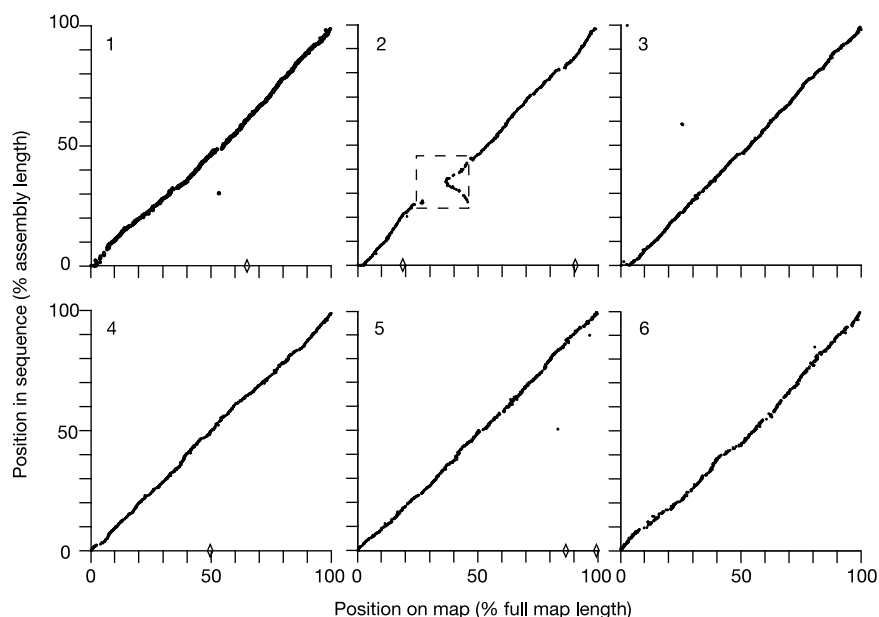


Figure 1 Chromosomal assemblies compared against HAPPY map data. The locations of markers as found in the sequence (y axis) are plotted against their location in HAPPY maps (x axis) for chromosomes 1–6. Markers mapped to one chromosome but found in the

assembled sequence of another are indicated by diamonds on the x axis. The dashed box indicates a large inverted duplication on chromosome 2: markers in this region are shown at one of their two possible map locations but are found at two points in the sequence.

chromosome 2 carries both a partial rDNA sequence and a 5S rRNA pseudogene^{9,18}.

In this study, two unanchored contigs assigned to chromosomes 4 and 5 contained junctions between rDNA sequences and complex repeats—attempts to extend the sequence and integrate these contigs into the assemblies failed owing to the highly repetitive nature of the adjoining sequences. We postulate that these contigs represent the junctions between a 'master copy' of the rDNA and the remainder of chromosome 4 (Fig. 2). One contig contains sequence matching a region of (G+C)-rich repeats near the centre of the palindrome, whereas the other matches sequence near the tip of the palindrome arm, adjacent to the one unclosed gap in the rDNA element sequence¹⁷. This gap is believed to represent a tandem array of short repeats, probably added post-synthetically to the extrachromosomal elements.

The structure of this master copy suggests a mechanism for generating the extrachromosomal copies by a process of transcription, hairpin formation and second-strand synthesis (Fig. 2). This process would account for the complete absence of sequence variation between the two arms of the palindrome.

Centromeres, telomeres and rearrangements

Repeat clusters may serve as centromeres

Centromeres mobilize eukaryotic chromosomes during cell division but vary widely in their structure and organization¹⁹, making them difficult to identify. Each *Dictyostelium* chromosome carries a single cluster of repeats rich in DIRS (*Dictyostelium* intermediate repeat sequence) elements^{20,21} near one end²², and this sole but striking structural consistency suggests that these clusters may serve as centromeres. Although the repetitive nature of the chromosomal termini impeded their assembly, most of the cluster on chromosome 1 was assembled (Fig. 3) and shows a complex pattern of DIRS and related Skipper elements, each preferentially associated with others of the same type. Frequent insertions and partial deletions have created a mosaic with little long-range order.

In *Dictyostelium* cells demonstrating condensed chromosomes characteristic of mitosis, DIRS-element probes hybridize to one end of each chromosome (Supplementary Fig. 4), consistent with the mapping data. DIRS-like elements in other species are more uniformly scattered along the chromosomes²³, suggesting that their restricted distribution in *Dictyostelium* chromosomes is functionally important. Furthermore, the DIRS-containing ends of the chromosomes cluster not only during mitosis, but also during interphase (Supplementary Fig. 4), as has been observed for centromeres in *Schizosaccharomyces pombe*²⁴.

rDNA sequences seem to act as telomeres

No (G+T)-rich telomere-like motifs were identified in the sequence; however, earlier findings²² suggested that the chromosomes terminate in the same (G+A)-rich repeat motif that caps the extrachromosomal rDNA element. We therefore surveyed all shotgun sequence to identify reads containing a junction between complex repetitive elements and rDNA-like sequence. Only 556 such reads were identified, of which 221 could be built into 13 contigs, which we refer to as C/R (complex-repeat/rDNA) junctions.

Of the 13 junctions, two represent known regions lying internally within the chromosomal assemblies. Of the remaining 11, one had twice the sequence coverage of the others, suggesting that it represents two distinct but identical portions of the genome (a possibility supported by the fact that another two of the junctions differed from each other by only two bases). Hence, we infer that the 11 remaining contigs represent 12 distinct junctions between repetitive elements and rDNA-like sequences—potentially one for every chromosomal end.

On the basis of their content of sequence reads from each of the whole-chromosome libraries, we assigned two of the C/R junctions

to each of the chromosomes. Chromosomes 4 and 5 cannot be distinguished in this way, but three junctions, including the one believed to be present as two copies, are assigned to this chromosome pair. The point in the rDNA palindrome that is represented differs from one junction to the next (Supplementary Fig. 5), but several junctions fall at common parts of the palindrome. This may reflect a preference in the mechanism that forms or maintains the junctions, or may result from a homogenizing recombination between them or with other rDNA sequences. Certainly the low frequency of differences between the rDNA components of the junction fragments and the extrachromosomal rDNA element argues for some process that limits or rectifies mutation. At each junction, we see only the rDNA sequence that immediately adjoins the complex repeat, as further assembly is precluded by the multi-copy nature of rDNA. Therefore we cannot tell whether each junctional rDNA sequence extends to the telomere-repeat-carrying tip of the rDNA palindrome sequence, nor whether other sequences lie beyond the rDNA components.

HAPPY mapping of markers derived from six of these C/R junctions confirmed not only the chromosomal assignments that had been made based on the origins of their component sequences, but also their locations at the termini of the mapped regions of the chromosomes. For the other junctions, the absence of unique sequence features precluded such mapping. Taken as a whole, this evidence strongly suggests that rDNA-like elements form part of the telomere structure in *D. discoideum*, and that common mechanisms stabilize both the extrachromosomal rDNA element and the chromosomal termini.

Chromosome 2 duplication

Chromosome 2 of *D. discoideum* AX4 carries a perfect inverted 1.51-

Figure 2 The genome of *Dictyostelium discoideum*. On each of the chromosomal assemblies (numbered 1–6) the diameter of the tube represents coding density (proportion of coding bases summed over both strands; centre-weighted sliding window of 100 kb; scale on right). The coloured bands on the chromosomes represent tRNAs (red), complex repeats (blue), gaps (black) and ribosomal DNA sequences (yellow). G+C content is plotted above each chromosome (centre-weighted sliding window of 100 kb; scale on left). The locations of HAPPY markers are indicated by short green ticks immediately below the distance scale. Immediately beneath each chromosome, the locations (short vertical ticks) of genes known to be upregulated (red), downregulated (blue) or whose level of expression does not change significantly (grey) in the transition from solitary to aggregative existence (expression data from ref. 91) are indicated; coloured horizontal bars below this indicate significant clusters of genes that are preferentially expressed in germinating spores (red), de-differentiating cells (green), pre-spore cells (blue) or in pre-stalk cells (yellow). The translucent 'hourglass' shape on chromosome 2 is centred on a large inverted duplication. The translucent cylinder on chromosome 3 indicates a typical 300-kb region, which is shown in expanded form in inset **a** to illustrate the clustering of identical tRNA genes (red arrows indicate polarity of tRNA genes); a 50-kb section of this region is expanded further in inset panel **b**, revealing the close association of TRE elements (specific family named above) with tRNAs. The translucent yellow disc on chromosome 4 indicates the location of the presumed chromosomal master copy of the rDNA element. In inset panel **c**, the structure of the palindromic extrachromosomal element is shown schematically. (I) Magenta bands indicate rDNA genes; green bands indicate G+C-rich regions; red end caps indicate short repetitive telomere structures; the translucent hoop indicates the central region of asymmetry. (II) Two chromosomal sequence contigs, each carrying an rDNA-like sequence (green or yellow; dotted lines indicate corresponding part of element) flanked by complex repeats (blue). From these contigs, we infer the probable structure (III) of the genomic master copy (grey indicates flanking sequence on chromosome 4). This structure suggests a mechanism for regenerating the extrachromosomal copies by transcription of a single strand (IV), hairpin formation and strand extension (V; broken line indicates synthesis of complementary strand), unfolding of the hairpin and synthesis of a fully complementary strand (VI; broken line indicates synthesis of second strand; telomeric caps added post-synthetically).

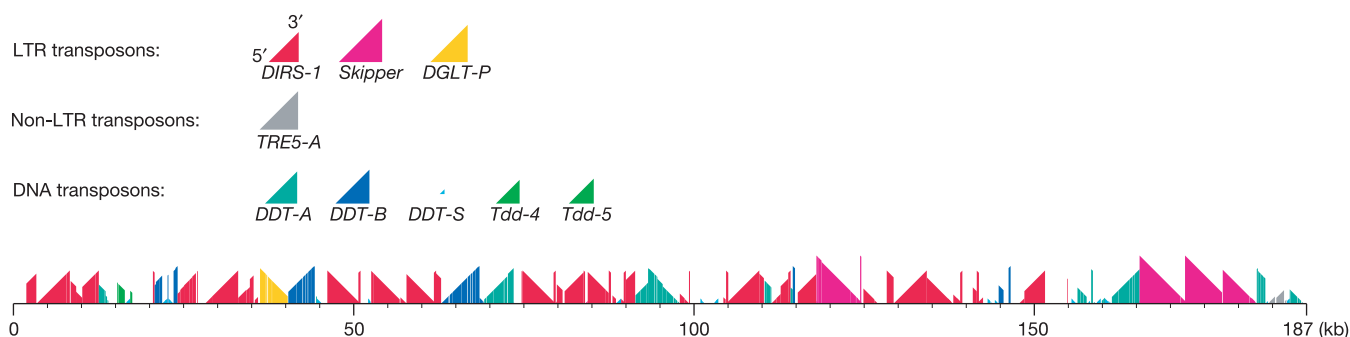


Figure 3 DIRS repeat region of chromosome 1. Complete complex repeat units are represented by coloured triangles whose size corresponds to the sequence length of the repeat unit (see key at top of figure). The bottom-left and top-right corners of each triangle both at about 3.74 Mb (Fig. 2)—have been implicated in centromeric and telomeric functions, respectively, elsewhere in the genome.

megabase (Mb) duplication (Fig. 2; see also refs 9, 25). This duplication, containing 608 genes, is known²⁵ to be absent from the wild-type isolate NC4 and from one of its direct descendents (AX2), but present in another (AX3); AX4 in turn is derived from AX3. The sequences adjoining the right-hand end of the duplication—a partial copy of a DIRS element (and a partial DDT-A element) and a region identical to part of the rDNA palindrome, both at about 3.74 Mb (Fig. 2)—have been implicated in centromeric and telomeric functions, respectively, elsewhere in the genome.

We propose that this duplication arose from a ‘breakage-fusion-bridge’ cycle as first described in maize²⁶ and since observed in many genomes. The nearby DIRS and rDNA components, in this view, represent abortive attempts to stabilize the halves of the broken chromosome by establishing new telomeres and centromeres, followed by re-fusion of the pieces to create a restored and enlarged chromosome (Supplementary Fig. 6).

Chromosome 2 (the largest of the chromosomes, even discounting the duplication in AX4) may be prone to breakage: in the Bonner isolate of NC4, maintained in vegetative growth for 50 years, chromosome 2 is represented by two smaller fragments²⁷. Comparison with more recent data²² indicates that the break point in NC4-Bonner lies in the same region as the duplication in AX4, suggesting that NC4-Bonner underwent the early stages of this process, but that the chromosome fragments were stabilized and maintained after the initial breakage. Preliminary results (data not shown) from HAPPY mapping also suggest that although wild-type isolates V12M2 and NC4 both lack the duplication seen in AX4, NC4 may carry a duplication of ~300 kb near the opposite end of chromosome 2.

partial repeat units within the first 187 kb of *D. discoideum* chromosome 1 is shown (bottom) by corresponding portions of the triangles; the orientation of the triangles indicates the direction in which each repeat unit lies. The vertical scale (sizes of repeat units) is the same as the horizontal scale (chromosomal distances).

Content and organization of the proteome

Prediction of protein-coding genes (see Methods) was performed on the complete set of chromosomes and floating contigs (Table 2). In assessing the completeness and accuracy of the predictions, we find that of the 957 well-characterized *D. discoideum* genes that are present in the current sequence, 823 (86%) are predicted as transcripts with structures matching the experimentally determined ones. For a further 123 (13%), the predicted transcript differs from the experimentally determined one, about one-half of these differing only in their 5' boundary; the remaining 11 (1%), although present in the sequence, were not predicted as transcripts. Similarly, of the 128,207 qualified ESTs present in the current sequence, 127,097 (99.1%) fall within predicted transcripts. Combining our estimate of sequence coverage (above) with these estimates of the success of gene prediction, we infer that approximately 98% of all *D. discoideum* genes are present in the predicted set.

The level of overprediction, conversely, is harder to estimate: prediction was performed generously to ensure that most true genes were represented. Of the 13,541 predicted proteins, 47.5% are represented by qualified ESTs, reflecting the inevitable bias in EST sampling. Among the shortest predicted proteins, fewer are represented by ESTs (for example, 21% of those of <60 amino acids); this is at least partly due to a higher level of overprediction. On the basis of the simplifying assumption that 50% of all genes coding for proteins of <100 amino acids are mis-predictions, we estimate the true number of genes at roughly 12,500. This number is closer to that seen in multicellular organisms rather than in most unicellular eukaryotes (Table 2). The same relative complexity is seen in the total number of amino acids encoded by the respective genomes; this measure of complexity is less affected by the inclusion

Table 2 Comparison between the predicted protein-coding gene set of *D. discoideum* and those of other organisms

Feature	<i>D. discoideum</i>	<i>P. falciparum</i>	<i>S. cerevisiae</i>	<i>A. thaliana</i>	<i>D. melanogaster</i>	<i>C. elegans</i>	Human
Genome size (Mb)	34	23	13	125	180	103	2,851
Number of genes	12,500*	5,268	5,538	25,498	13,676	19,893	22,287
Gene spacing (kb per gene)	2.5	4.3	2.2	4.9	13.2	5.0	127.9
Mean gene length (bp)	1,756	2,534	1,428	2,036	1,997	2,991	27,000
Mean coding size (amino acids)	518	761	475	437	538	435	509
Genes with introns (%)	69	54	5	79	38	5	85
Mean intron size (bp)	146	179	ND	170	ND	270	3,365
Mean no. of introns (in spliced genes)	1.9	2.6	1.0	5.4	4.0	5.0	8.1
Total amino acids encoded (thousands)	7,021	4,009	2,471	11,143	7,358	9,038	11,333
Codon A + T bias†	86	83	62	57	50	64	41
Mean A + T percentage (exons)	73	76	72	72	45	58	55
Mean A + T percentage (introns)	88	87	51	55	38	71	62
Mean A + T percentage (intergenic)	85	86	51	56	38	72	62

ND, not determined.

* See text. The estimated number of true transcripts for *D. discoideum* is given here for comparability with other species; however, the total predicted gene number of 13,541 is used in calculating the figures below.

† Percentage of all codons used which have A/T at their third base.

of shorter (and hence more dubious) gene predictions. Introns in *Dictyostelium* are few and short, and intergenic regions are small, producing a compact genome of which 62% encodes protein.

Genes are distributed approximately uniformly across the genome (Fig. 2). Although we do not see widespread clustering of genes with coordinated expression patterns (see Methods), we do find statistically significant ($P < 0.01$) clusters of genes expressed predominantly at some developmental stages or in specific cell types (Fig. 2).

(A+T)-richness influences protein composition and codon usage

Codon usage in *Dictyostelium* favours codons of the form NNT or NNA over their NNG or NNC synonyms, the bias being even greater than for the (A+T)-rich *Plasmodium* genome. Comparison of tRNA and codon frequencies (Supplementary Table 2) reveals a similar picture to that in human²⁸ and other eukaryotes, suggesting that the same use is made of 'wobble' and of base modifications (for example, of adenine to inosine in some tRNAs) to expand the effective repertoire of tRNAs.

As in *Plasmodium*²⁹, the extreme (A+T)-richness is reflected not just in the choice of synonymous codons, but also in the amino acid composition of the proteins. Amino acids encoded solely by codons of the form WVN (where W indicates A or T and N indicates any base; these are Asn, Lys, Ile, Tyr and Phe) are much commoner in *Dictyostelium* proteins than in human ones; the reverse is true for those encoded solely by SSN codons (where S indicates C or G; these are Pro, Arg, Ala and Gly).

Geometry reflects phylogeny—duplications in the genome

The predicted gene set of *Dictyostelium* is rich in relatively recently duplicated genes. Of the 13,498 predicted proteins analysed, 3,663 fall into 889 families clustered by BLASTP similarities of $e < 10^{-40}$. Most (538) families contain only two members, but 351 families contain between three and 81 proteins (Supplementary Table 3). Hence, 2,774 (20%) of all predicted proteins have arisen by relatively recent duplication, potentially accounting for much of *Dictyostelium*'s excess gene number compared with typical unicellular eukaryotes.

We tried to infer the mechanisms by which such duplications

arise and propagate in the genome. Where members of a family are clustered on one chromosome, the physical distance between family members often (23 out of 86 families examined) correlates strongly with their evolutionary divergence (see Methods). Where a family is split between different chromosomes, members on the same chromosome are often (23 out of 50 families examined) more related to each other than to members on different chromosomes; the reverse is never observed.

These findings suggest that three processes combine to account for most of the duplications in *Dictyostelium*: tandem duplication, local inversion and interchromosomal exchange. In this model, gene families expand by tandem duplication of either single genes or blocks containing several consecutive genes, as in an earlier model³⁰; inversions within these expanding clusters may reverse local gene order. An elegant illustration of these two processes is provided by a cluster of acetyl-coA synthetases on chromosome 2 (Fig. 4). The third process (exchange of segments between chromosomes) may fragment these clusters at any stage. If such an interchromosomal exchange splits a gene family early in its expansion, then each of the two resulting subfamilies has a long subsequent period of evolution independent of the other, so similarities will be greatest between genes on the same chromosome. If, conversely, the split occurs later, then all family members, whether on the same chromosome or on different chromosomes, will tend to resemble each other equally closely. We cannot exclude the possibility of duplication occasionally creating a second copy of a gene, or group of genes, directly on a different chromosome from the first. However, all instances that we have examined can be accounted for without such intermolecular duplication.

Amino acid repeats

Tandem repeats of trinucleotides (and of motifs of 6, 9, 12, and so on, bases) are unusually abundant in *Dictyostelium* exons and naturally correspond to repeated sequences of amino acids. However, at the protein level the situation is even more extreme: there are many further amino acid repeats that use different synonymous codons, and so do not arise from perfect nucleotide repeats. Among the predicted proteins, there are 9,582 SSRs of amino acids (homopolymers of length ≥ 10 , or ≥ 5 consecutive repeats of a motif of two

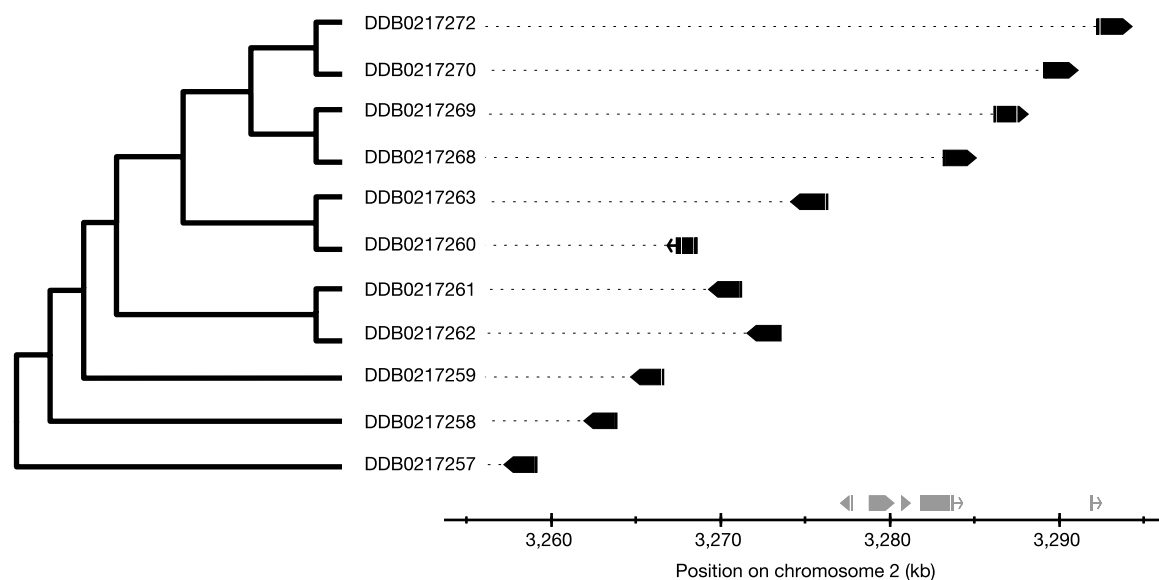


Figure 4 Phylogeny of gene family members compared to their physical order. The optimally parsimonious phylogenetic tree of 11 acetyl-CoA synthase genes, computed using the PHYLIP module 'Protpars' (<http://evolution.gs.washington.edu/phylip/doc/protpars.html>), is shown to the left; dictyBase identification numbers are shown at the end of each branch. The graph (right) indicates the arrangement on chromosome 2 of these

genes (solid black boxes; gaps indicate introns; pointed ends indicate direction of transcription). Chromosomal distance scale is given along the bottom and other unrelated genes in the same region are indicated in grey above the x axis. The correspondence between phylogeny and physical order implies that the cluster has arisen by a series of segmental tandem duplications and local inversions in parallel with sequence divergence.

or more amino acids). Of these, the most striking are polyasparagine and polyglutamine tracts of ≥ 20 residues, present in 2,091 of the predicted proteins. Also abundant are low-complexity regions such as QLQLQQQQQLQLQQ: there are 2,379 tracts of ≥ 15 residues composed of only two different amino acids. In total, repeats or simple-sequence tracts of amino acids (even by these conservative definitions) occur in 34% of predicted proteins and encode 3.3% of all amino acids.

It seems likely that these repeats have arisen through nucleotide expansion, but have been selected at the protein level. Evidence for selection at the protein level is that any given trinucleotide repeat occurs predominantly in only one of the three reading frames. For example, the repeat ...ACAACAACAACA... is usually translated as polyglutamine ([CAA] n) rather than polythreonine ([ACA] n) or polyasparagine ([AAC] n). Further evidence comes from the many trinucleotide repeats that have apparently mutated to produce only synonymous codons (for example, ...GATGACGATGATGAC..., translated as polyaspartate). Moreover, the distribution of repeats and simple-sequence tracts is nonrandom: most proteins either have no such features (66% of proteins) or have two or more (18% of proteins), suggesting that they are tolerated only in certain types of protein. The polyasparagine- and polyglutamine-containing proteins appear to be over-represented in protein kinases, lipid kinases, transcription factors, RNA helicases and messenger RNA binding proteins such as spliceosome components (Supplementary Fig. 9). Protein kinases and transcription factors are also over-represented in the polyasparagine- and polyglutamine-containing proteins of *Saccharomyces cerevisiae*, so it is possible that these homopolymers serve some functional role in these protein classes. A more detailed analysis of amino acid homopolymers is given in Supplementary Tables 4–6 and Supplementary Figs 7–10.

Phylogeny, evolution and comparative proteomics

The organisms that diverged from the last common ancestor of all eukaryotes followed different evolutionary paths, but all retained the basic properties of eukaryotic cells. Their genomes have been sculpted by chromosomal deletions and duplications that led to lineage-specific gene family expansions, reductions and losses, as

well as genes with new functions^{31,32}. Our analysis of *Dictyostelium*'s proteome shows that similar mechanisms have shaped its genome, augmented by horizontal gene transfer from bacterial species.

Phylogeny of eukaryotes based on complete proteomes

Using morphological criteria, early workers were unsure whether to classify Dictyostelids as fungi or protozoa³³. Molecular methods indicated that they were amoebozoia and also suggested that *Dictyostelium* diverged from the line leading to animals at about the same time as plants^{34,35}. A study of more than 100 proteins suggested that *Dictyostelium* diverged after the plant–animal split, but before the divergence of the fungi³⁶. The recent finding of a gene fusion encoding three pyrimidine biosynthetic enzymes, shared only by *Dictyostelium*, fungi and Metazoa, indicates that the amoebozoia are a true sister group of the fungi and Metazoa³⁷.

To examine the phylogeny of *Dictyostelium* on a genomic scale, we applied an improved method for predicting orthologous protein clusters to complete eukaryotic proteomes³⁸ (for details, see Supplementary Information). The data were used to construct a phylogenetic tree that confirms the divergence of *Dictyostelium* along the branch leading to the Metazoa soon after the plant–animal split (Fig. 5). Despite the earlier divergence of *Dictyostelium*, many of its proteins are more similar to human orthologues than are those of *S. cerevisiae*, probably due to higher rates of evolutionary change along the fungal lineage. Whether the greater similarity between amoebozoia and Metazoa proteins translates into a generally higher degree of functional conservation between them compared to the fungi remains to be seen.

Proteins shared by *Dictyostelium* and major organism groups

To examine shared functions, we identified eukaryote-specific Superfamily and Pfam protein domains, and sorted them according to their presence or absence within 12 completely sequenced genomes to arrive at their distribution among the major organismal groups (see Supplementary Tables 7–10 and Supplementary Fig. 11). Plants, Metazoa, fungi and *Dictyostelium* all share 32% of the eukaryotic Pfam domains (Fig. 6). The protein domains present

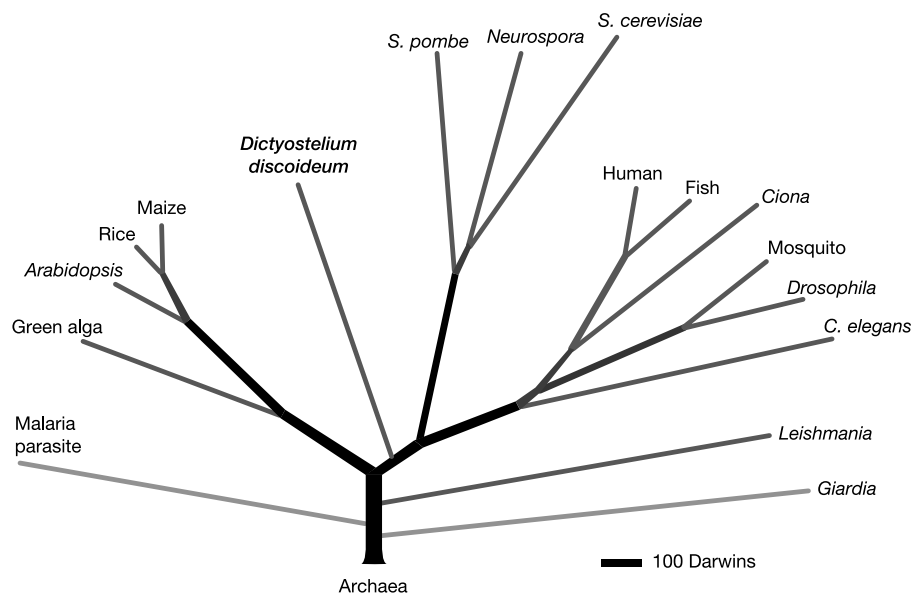


Figure 5 Proteome-based eukaryotic phylogeny. The phylogenetic tree was reconstructed from a database of 5,279 orthologous protein clusters drawn from the proteomes of the 17 eukaryotes shown, and was rooted on 159 protein clusters that had representatives from six archaeobacterial proteomes. Tree construction, the database of protein clusters and a model of protein divergence used for maximum likelihood

estimation are described in Supplementary Information. The relative lengths of the branches are given as Darwins (where 1 Darwin = 1/2,000 of the divergence between *S. cerevisiae* and humans). Species that are not specified are *Plasmodium falciparum* (malaria parasite), *Chlamydomonas reinhardtii* (green alga), *Oryza sativa* (rice), *Zea mays* (maize), *Takifugu rubripes* (fish) and *Anopheles gambiae* (mosquito).

in *Dictyostelium*, Metazoa and fungi, but absent in plants, are interesting because they probably arose soon after plants diverged and before *Dictyostelium* diverged from the line leading to animals. The major classes of domains in this group of proteins include those involved in small and large G-protein signalling (for example, RGS proteins), cell cycle control and other domains involved in signalling (Supplementary Tables 8 and 9). It also appears that glycogen storage and usage arose as a metabolic strategy soon after the plant–animal divergence, because glycogen synthetase seems to have appeared in this evolutionary interval.

Particularly notable are the cases where otherwise ubiquitous domains appear to be completely absent in one group or another. For instance, *Dictyostelium* seems to have lost the genes that encode collagen domains, the circadian rhythm control protein timeless and basic helix–loop–helix transcription factors (Supplementary Table 7). Metazoa, on the other hand, appear to have lost receptor histidine kinases that are common in bacteria, plants and fungi, whereas *Dictyostelium* has retained and expanded its complement to 14 members³⁹.

Orthologues of human disease genes

An important motivation for sequencing the *Dictyostelium* genome was to aid the discovery of proteins that would facilitate studies of orthologues in human, with possible implications for human health. Although orthologues of human genes implicated in disease are of course present in many species, *Dictyostelium* provides a potentially valuable vehicle for studying their functions in a system that is experimentally tractable and intermediate in complexity between the yeasts and the higher multicellular eukaryotes. To assess the usefulness of *Dictyostelium* for investigating the functions of genes related to human disease we used the protein sequences of 287 confirmed human disease genes as queries and carried out a systematic search for putative orthologues in the *Dictyostelium* proteome⁴⁰. At a stringent threshold value of $e \leq 10^{-40}$, we identified 64 such proteins. Of these, 33 were similar in length to the human protein and had similarity extending over >70% of the two proteins (Table 3). The number of *Dictyostelium* orthologues of human disease genes is lower than in *D. melanogaster* or *Caenorhabditis elegans* but higher than in *S. cerevisiae* or *S. pombe*. Of the 33 putative orthologues of confirmed human disease genes in *Dictyostelium*, five are absent in both *S. cerevisiae* and *S. pombe* (e -value $\leq 10^{-30}$), a further four are absent from *S. cerevisiae* and two are not found in *S. pombe*.

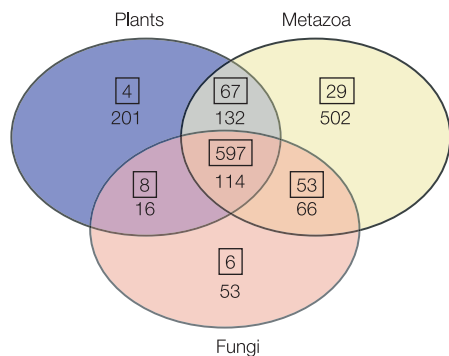


Figure 6 Distribution of Pfam domains among eukaryotes. The number of eukaryote-specific Pfam domains present in each group of eukaryotic organisms is shown. The boxed numbers are the domains that are present in *Dictyostelium* and the other numbers are those domains that are absent from *Dictyostelium*. The animals are *H. sapiens*, *T. rubripes*, *C. elegans*, *D. melanogaster*; the fungi are *N. crassa*, *Aspergillus nidulans*, *S. pombe* and *S. cerevisiae*; and the plants are *Arabidopsis thaliana*, *O. sativa* and *C. reinhardtii*. A complete listing of the domains can be found in the Supplementary Information.

Horizontal gene transfer

The acquisition of genes by horizontal transfer from one species to another (HGT) has become increasingly recognized as a mechanism of genome evolution^{41–43}. We identified 18 potential instances of HGTs, by screening *Dictyostelium* protein domains that are similar to bacteria-specific Pfam domains and have phyletic relationships consistent with HGT (see Supplementary Information). The transferred domains appear to have replaced functions, added new functions or evolved into new functions (Table 4). The *thy1* gene, which encodes an alternative form of thymidylate synthase (ThyX), appears to have replaced the endogenous gene, as the conventional thymidylate synthase (ThyA) is not present⁴⁴. Other HGT domains also have established functions, which are presumably retained and give *Dictyostelium* the ability to degrade bacterial cell walls (dipeptidase), scavenge iron (siderophore), or resist the toxic effects of tellurite in the soil (terD). Still other horizontally transferred domains have become embedded within *Dictyostelium* genes that encode larger proteins. An example of this is the Cna B domain that is found within four large predicted proteins, one of which, colossin A, is predicted to be 1.2 MDa (Supplementary Fig. 12).

Dictyostelium ecology

Dictyostelium faces many complex ecological challenges in the soil. Amoebae, fungi and bacteria compete for limited resources in the soil while defending themselves against predation and toxins. For instance, the nematode *C. elegans* is a competitor for bacterial food

Table 3 *Dictyostelium* genes related to human disease genes

Disease category*	SwissProt†	dictyBase ID‡
Cancer		
Colon cancer (MSH2)	MSH2_HUMAN	DDB0202539
Colon cancer (MLH1)	MLH1_HUMAN	DDB0187465
Colon cancer (MSH3)	MSH3_HUMAN	DDB0204604
Colon cancer (PMS2)	PMS2_HUMAN	DDB0185791
Xeroderma pigmentosum (ERCC3)	XPB_HUMAN	DDB0206281
Xeroderma pigmentosum (XPD)	XPD_HUMAN	DDB0189539
Oncogene (AKT2)	AKT2_HUMAN	DDB0189970
Oncogene (RAS)	RAS_HUMAN	DDB0191937
Cyclin-dependent kinase 4 (CDK4)	CDK4_HUMAN	DDB0188077
Neurological		
Lowe oculocerebrorenal (OCRL)	OCRL_HUMAN	DDB0189888
Miller–Dieker lissencephaly (PAF)	LIS1_HUMAN	DDB0219335
Adrenoleukodystrophy (ABCD1)	ABCD1_HUMAN (P)	DDB0219834
Angelman (UBE3A)	UBE3A_HUMAN	DDB0188760
Ceroid lipofuscinosis (CLN2)	TPP1_HUMAN (C, P)	DDB0190668
Tay–Sachs (HEXA)	HEXA_HUMAN (C, P)	DDB0187255
Ceroid lipofuscinosis (PPT)	PPT1_HUMAN (C)	DDB0186550
Thomsen myotonia congenita (CLCN1)	CLCN1_HUMAN	DDB0191805
Choroideremia (CHM)	RAE1_HUMAN	DDB0206402
Amyotrophic lateral sclerosis (SOD1)	SODC_HUMAN	DDB0188850
Parkinson's (UCHL1)	UCHL1_HUMAN (C, P)	DDB0205083
Cardiovascular		
Hypertrophic cardiomyopathy	MYH7_HUMAN	DDB0186963
Renal		
Renal tubular acidosis (ATP6B1)	VAB1_HUMAN	DDB0169211
Hyperoxaluria (AGXT)	SPYA_HUMAN (C, P)	DDB0188646
Metabolic/endocrine		
Niemann–Pick type C (NPC1)	NPC1_HUMAN (P)	DDB0191057
Hyperinsulinism (ABCC8)	ACC8_HUMAN	DDB0187670
McCune–Albright (GNAS1)	GBAS_HUMAN	DDB0185461
Pendred (PDS)	PEND_HUMAN (C)	DDB0202939
Haematological/immune		
G6PD deficiency (G6PD)	G6PD_HUMAN	DDB0168147
Chronic granulomatous (CYBB)	C24B_HUMAN (C, P)	DDB0188527
Malformation		
Diastrophic dysplasia (SLC26A2)	DTD_HUMAN (C)	DDB0202939
Other		
Cystic fibrosis (ABCC7)	CFTR_HUMAN	DDB0186232
Darier–White (SERCA)	ATA2_HUMAN	DDB0169159
Congenital chloride diarrhoea (DRA)	DRA_HUMAN (C)	DDB0202939

*From a list of 287 confirmed human disease protein sequences⁴⁰. Those listed match a predicted *Dictyostelium* protein with a BLASTP probability of $e \leq 10^{-40}$, are similar in length ($\pm 25\%$ in comparison to the *Dictyostelium* protein) and both proteins align over more than 70% of their respective lengths.

†SwissProt identifiers for the human proteins. Letters in brackets indicate that the protein has no homologue (BLASTP probability of $e \leq 1.0 \times 10^{-30}$) in *S. cerevisiae* (C) or *S. pombe* (P).

‡The best match to the human gene is listed by its dictyBase identification number. Matches with a BLASTP probability of $e \leq 10^{-100}$ are indicated in bold.

and a predator of *Dictyostelium* amoebae, but also a potential dispersal agent for *Dictyostelium* spores⁴⁵. *Dictyostelium* has expanded its repertoire of several protein classes that are probably crucial for such interspecies interactions and for survival and motility in this complex ecosystem.

Polyketide synthases

A small number of natural products have already been identified from *Dictyostelium*, but the gene content suggests that it is a prolific producer of such molecules. Some of them may act as signals during development, such as the dichlorohexanophenone DIF-1, but others are likely to mediate currently unknown ecological interactions⁴⁶. Many antibiotics and secondary metabolites destined for export are produced by polyketide synthases, modular proteins of around 3,000 amino acids⁴⁷. We identified 43 putative polyketide synthases in *Dictyostelium* (see Supplementary Information). By contrast, *S. cerevisiae* completely lacks polyketide synthases and *Neurospora crassa* has only seven. Furthermore, two of the *Dictyostelium* proteins have an additional chalcone synthase domain, representing a type of polyketide synthase most typical of higher plants and found to be exclusively shared by *Dictyostelium*, fungi and plants. In addition to polyketide synthases, the predicted proteome has chlorinating and dechlorinating enzymes as well as O-methyl transferases, which could increase the diversity of natural products made. Thus, *Dictyostelium* appears to have a large secondary metabolism, which warrants further investigation.

ABC transporters

ATP-binding cassette (ABC) transporters are prevalent in the proteomes of soil microorganisms and are thought to provide resistance to xenobiotics through their ability to translocate small-molecule substrates across membranes against a substantial concentration gradient^{48–51}. There are 66 ABC transporters encoded by the genome, which can be classified according to the subfamilies defined in humans (ABCA, ABCB, ABCC, ABCD, ABCE, ABCF and ABCG) based on domain arrangement and signature sequences⁵². At least 20 of them are expressed during growth and are probably involved in detoxification and the export of endogenous secondary metabolites.

Cellulose degradation

Many of the predicted cellulose-degrading enzymes in the proteome

(see Supplementary Information) that have secretion signals are expressed in growing cells that do not produce cellulose⁵³. The proteome also contains one xylanase enzyme that can degrade the xylan polymers that are often found associated with the cellulose of higher plants. Perhaps *Dictyostelium* uses these enzymes to degrade plant tissue into particles that are then taken up by cells. These enzymes may also aid in the breakdown of cellulose-containing microorganisms upon which *Dictyostelium* feeds. Alternatively, these enzymes may promote the growth of bacteria that can serve as food, because *Dictyostelium*'s habitat also contains cellulose-degrading bacteria.

Specializations for cell motility

During both growth and development, *Dictyostelium* amoebae display motility that is characteristic of human leukocytes⁵⁴. As a consequence, studies of *Dictyostelium* have contributed significantly to cytoskeleton research⁵⁵. *Dictyostelium*'s survival depends on an ability to efficiently sense, track and consume soil bacteria using sophisticated systems for chemotaxis and phagocytosis. Its multicellular development depends on chemotactic aggregation of individual amoebae and the coordinated movement of thousands of cells during fruiting body morphogenesis. The proteome reveals an astonishing assortment of proteins that are used for robust, dynamic control of the cytoskeleton during these processes. As suggested by functional parallels to human cells, these proteins are most similar to metazoan proteins in their variety and domain arrangements (Fig. 7; see also Supplementary Table 11). Surprisingly, although the actin cytoskeleton has been studied for over 25 years, 71 putative actin-binding proteins apparently escaped classical methods of discovery. For example, actobindins had not been previously recognized in *Dictyostelium*. Curiously, the actin depolymerization factor (ADF) and calponin homology (CH) domain proteins appear to have diversified by domain shuffling, a substantial fraction having domain combinations unique to *Dictyostelium* (Supplementary Table 12 and Supplementary Fig. 13). In addition to 30 actin genes, there are also orthologues of all actin-related protein (ARP) classes present in mammals, as well as three founding members of a new class (Supplementary Fig. 14).

Cytoskeletal remodelling during chemotaxis and phagocytosis is regulated by a considerable number of upstream signalling components. Of the 18 Rho family GTPases in *Dictyostelium*, some are

Table 4 Candidate horizontal gene transfers from bacteria

Function*	Pfam†	Number of proteins‡	DictyBase ID§	Length (aa)	Region matched¶	e-value#
Aromatic amino acid lyase	Beta_elim_lyase	2	DDB0204031	170	4–170	3.2×10^{-65}
Biotin metabolism	BioY	1	DDB0184375	338	145–299	5.8×10^{-20}
Unknown	Cna_B	4	DDB0184530	11,103	Multiple††	1.1×10^{-10}
Peroxidase	Dyp_peroxidase	1	DDB0168077	306	3–303	1.4×10^{-82}
Insecticide	Endotoxin_N	2	DDB0188332	628	38–210	1.2×10^{-32}
Isopentenyl transferase	IPT	1	DDB0169077	283	1–63	5.1×10^{-12}
Siderophore	lucA_lucC	2	DDB0219918	739	183–350	2.3×10^{-18}
Osmoregulation	OsmC	2	DDB0190102	156	16–156	9.8×10^{-22}
Dipeptidase/β-lactamase	Peptidase M15	1	DDB0205124	897	68–406; 711–879	3.4×10^{-16}
Dipeptidase/β-lactamase	Peptidase S13	1	DDB0168572	522	337–495	4.2×10^{-25}
Polyphosphate synthesis	PP_kinase	1	DDB0192001	1,053	372–1045	1.6×10^{-234}
Tellurite resistance	TerD	2	DDB0169240	287	152–279	2.1×10^{-67}
Thymidylate synthesis	Thy1	1	DDB0214905	303	38–254	9.9×10^{-117}
Unknown	DUF84	1	DDB0203145	179	5–175	1.6×10^{-20}
Unknown	DUF885	2	DDB0205394	689	318–685	1.5×10^{-124}
(Prespore protein 3B)	DUF1121	3	DDB0169184	226	1–226	8.7×10^{-134}
Unknown	DUF1289	1	DDB0204782	88	29–85	3.3×10^{-15}
Unknown	DUF1294	1	DDB0186703	155	2–73	8.9×10^{-18}

* Confirmed or proposed function of the prokaryotic orthologue is given. For domains without function information, information on any *Dictyostelium* protein in the set is given in parentheses.

† The Pfam domain designation (<http://www.sanger.ac.uk/Software/Pfam/>).

‡ The number of gene models in which the domain appears. Bold numbers indicate gene sets where there are pairs of genes that map within 10 kb of each other.

§ The gene identification number for the example given in the rest of the table (release v2.0 at <http://www.dictybase.org/>).

|| Number of amino acid (aa) residues in the predicted *Dictyostelium* protein containing the domain.

¶ The region of the *Dictyostelium* protein that matched the prokaryotic domain. The amino acid sequence identity between this region and the most highly related prokaryotic protein was 21–52%.

The e-value for the domain against the Pfam model library used to identify it (see Supplementary Information).

†† The protein colossin A consists of an array of 91 partial Cna_B domains within 18 larger repeats, and the e-value corresponds to one domain.

Class	Protein	No.	Module	Occurrence			
				D	M	F	P
G-actin binding	Profilin	3	PRO				
	Actobindin-like	3	WH2				
	CAP	1	WH2				
	WH2-containing	5	WH2				
	Twinfilin-like	1	ADP				
Capping and/or severing	Cap32/34 (Aginactin)	2	CAP				
	Cofilin	6	ADF				
	Severin	1	GEL				
	GRP125	1	GEL				
	Gelsolin-related	2	GEL				
Actin capping and nucleation	Arp2/3 complex	7	ACT				
	Scar	1	WH2				
	WASP	1	WH2				
	WASP-related	2	WH2				
	VASP	1	EVH				
Actin cross-linking	Formin	10	FH2				
	ABP34	1					
	eEF1- α (ABP50)	2	EF1- α				
	eEF1- β	3					
	Dynactin	1					
	Fimbrin	1	CH				
	Fimbrin type ABD-containing	5	CH				
	Filamin (gelation factor)	1	CH				
	α -actinin	1	CH				
	Cortexillin	2	CH				
	α -actinin type ABD-containing	3	CH				
	Protovillin (Cap100)	1	GEL				
	Villin-related	1	GEL				
	Flightless/Villin-related	1	GEL				
	Villin	1	GEL				
Lateral actin binding	Kelch-related	1	KELCH				
	Smoothelin-related	1	CH				
	GAS2-related	1	CH				
	CH-containing	19	CH				
	VHP-containing	3	VHP				
	Coronin	1					
	Coronin-like	1					
	Aip	1					
	Coactosin	1	ADF				
	Coactosin-related	3	ADF				
Membrane-associated	Abp1	1	ADF				
	Glia maturation factor-related	1	ADF				
	UIM domain-containing	3					
	Interaptin	1	CH				
	Ponticulin	2					
	Ponticulin-related	2					
	Comitin	1					
	Comitin-related	1					
	Hisactophilin	3	TRE				
	Talin A (filopodin)	1	TAL				
Motors	Talin B	1	TAL				
	SLA-2-like	1	TAL				
	Annexin	2					
	Vinculin/ α -catenin-related	2					
	Conventional myosin	1	MYO				
	Unconventional myosins	12	MYO				
Total 138 (71 new)							

Figure 7 Microfilament system proteins. Proteins with probable interactions with the actin cytoskeleton are tabulated by their documented or predicted functions. Coloured boxes indicate the presence of a protein related to the *Dictyostelium* (D) protein in Metazoa (M), fungi (F) or plants (P). *Dictyostelium*-specific proteins have no recognizable relatives or differ from relatives due to extensions or unusual domain compositions. For details see Supplementary Information. Actin-binding modules: ACT, actin fold; ADF, actin depolymerization factor/cofilin-like domain; CAP, capping protein fold; CH, calponin homology domain; EVH, Ena/VASP homology domain 2; FH2, formin homology 2 domain; GEL, gelsolin repeat domain; KELCH, Kelch repeat domain; MYO, myosin motor domain; PRO, profilin fold; TAL, the I/LWEQ actin-binding domain of talin and related proteins; TRE, trefoil domain; VHP, villin head piece; WH2, Wiskott Aldrich syndrome homology region 2.

clear Rac orthologues and one belongs to the RhoBTB subfamily⁵⁶. However, the Cdc42 and Rho subfamilies characteristic of Metazoa and fungi are absent, as are the Rho subfamily effector proteins. The activities of these GTPases are regulated by two members of the RhoGDI family, by components of ELMO1–DOCK180 complexes and by a large number of proteins carrying RhoGEF and RhoGAP domains (>40 of each), most of which show domain compositions not found in other organisms. Remarkably, *Dictyostelium* appears to be the only lower eukaryote that possesses class I phosphatidylinositol-3-OH kinases, which are at the crossroad of several critical signalling pathways (for details of the regulators and their effectors, see Supplementary Table 13)⁵⁷. The diverse array of these regulators and the discovery of many additional actin-binding proteins suggest that there are many aspects of cytoskeletal regulation that have yet to be explored.

Multicellularity and development

The evolution of multicellularity was arguably as significant as the origin of the eukaryotic cell in enabling the diversification of life. The common unicellular ancestor of the crown group of organisms must have possessed the basic machinery to regulate nutrient uptake, metabolism, cellular defence and reproduction, and it is likely that these mechanisms were adapted to integrate the functions of cells in multicellular organisms. *Dictyostelium* achieved multicellularity through a different evolutionary route compared with plants and animals, yet the ancestors of these respective groups probably started with the same endowment of genes and faced the same problem of achieving cell specialization and tissue organization.

When starved, *Dictyostelium* develops as a true multicellular organism, organizing distinct tissues within a motile slug and producing a fruiting body comprised of a cellular, cellulosic stalk supporting a bolus of spores⁴. Thus, *Dictyostelium* has evolved differentiated cell types and the ability to regulate their proportions and morphogenesis. A broad survey of proteins required for multicellular development shows that *Dictyostelium* has retained cell adhesion and signalling modules normally associated exclusively with animals, whereas the structural elements of the fruiting body and terminally differentiated cells clearly derive from the control of cellulose deposition and metabolism now associated with plants. The *Dictyostelium* genome offers a first glimpse of how multicellularity evolved in the amoebozoan lineage. In the following sections, we consider some of the systems that are particularly relevant to cellular differentiation and integration in a multicellular organism.

Signal transduction through G-protein-coupled receptors

The needs of multicellular development add greatly to those of chemotaxis in demanding dynamically controlled and highly selective signalling systems. G-protein-coupled cell surface receptors (GPCRs) form the basis of such systems in many species, allowing the detection of a variety of environmental and intra-organismal signals such as light, Ca²⁺, odorants, nucleotides and peptides. They are subdivided into six families, which, despite their conserved secondary domain structure, do not share significant sequence similarity⁵⁸. Until recently, in *Dictyostelium* only the seven CAR/CRL (cAMP receptor/ cAMP receptor-like) family GPCRs had been examined in detail^{59,60}. Surprisingly, a detailed search uncovered 48 additional putative GPCRs of which 43 can be grouped into the secretin (family 2), metabotropic glutamate/GABA_B (family 3) and the frizzled/smoothed (family 5) families of receptors (Fig. 8; see also Supplementary Information). The presence of family 2, 3 and 5 receptors in *Dictyostelium* was surprising because they had been thought to be specific to animals. Their occurrence in *Dictyostelium* suggests that they arose before the divergence of the animals and fungi and were later lost in fungi, and that the radiation of GPCRs pre-dates the divergence of the animals and fungi. The secretin

family is particularly interesting because these proteins were thought to be of relatively recent origin, appearing closer to the time of the divergence of animals⁶¹. The putative *Dictyostelium* secretin GPCR does not contain the characteristic GPCR proteolytic site, but its transmembrane domains are clearly more closely related to secretin GPCRs than to other families (Fig. 8). Many downstream signalling components that transduce GPCR signals could also be recognized in the proteome, including heterotrimeric G-protein subunits (fourteen G α , two G β and one G γ proteins) and seven regulators of G-protein signalling (RGS) that share highest similarity with the R4 subfamily of mammalian RGS proteins.

SH2 domain signalling

In animals, SH2 domains act as regulatory modules of proteins in intracellular signalling cascades, interacting with phosphotyrosine-containing peptides in a sequence-specific manner. *Dictyostelium* is the only organism, outside of the animal kingdom, where SH2 domain phosphotyrosine signalling has been shown to occur⁶². What has been lacking in *Dictyostelium* is evidence of the other components of such signalling pathways; that is, equivalents of the metazoan SH2-domain-containing receptors, adaptors and targeting proteins. Three newly predicted proteins are strong candidates for these roles (Supplementary Fig. 15). One of them, CblA, is highly related to the metazoan Cbl proto-oncogene product. This is entirely unexpected because it is the first time that a Cbl homologue has been observed outside the animal kingdom. The Cbl protein is a 'RING finger' ubiquitin-protein ligase that recognizes activated

receptor tyrosine kinases and various molecular adaptors⁶³. Remarkably, the Cbl SH2 domain went unrecognized in the protein sequence, but it was revealed when the crystal structure of the protein was determined⁶⁴. Thus, although SH2 domain proteins are less prevalent in *Dictyostelium*, there is the potential for the kind of complex interactions that typify metazoan SH2 signalling pathways.

ABC transporter signalling

Dictyostelium, like other organisms, has adapted ABC transporters to control various developmental signalling events. Several ABC transporters (TagA, TagB and TagC) are used for peptide-based signalling, similar to that previously observed for mating in *S. cerevisiae* and antigen presentation in human T cells^{65–67}. The novel domain arrangement of the Tag proteins—a serine protease domain fused to a single transporter domain—suggests that they have been selected for improved efficiency in signal production. Additional ABC transporters are needed for cell fate determination in *Dictyostelium*, suggesting that this ubiquitous protein family may be used in similar developmental contexts within many different species⁶⁸.

Kinases and transcription factors

Much cellular signal transduction involves the regulation of protein function through phosphorylation by protein kinases, often leading to the reprogramming of gene transcription in response to extracellular signals. The *Dictyostelium* proteome contains 295 predicted protein kinases, representing as wide a spectrum of kinase families

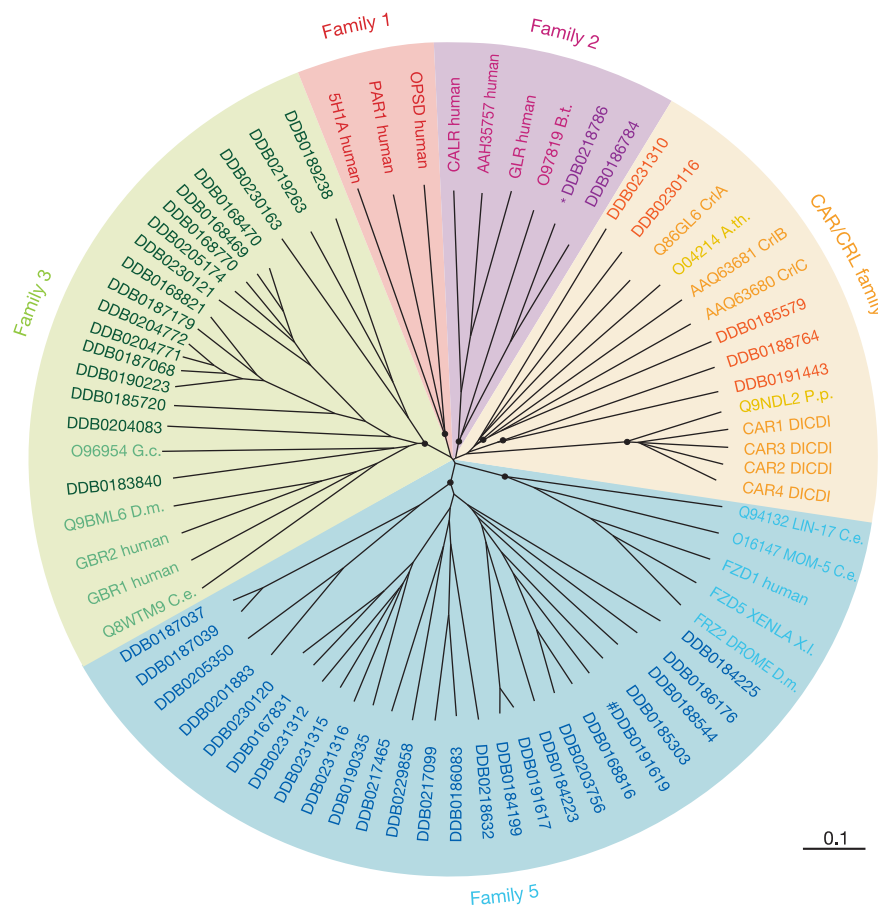


Figure 8 The G-protein-coupled receptors. A CLUSTALX alignment of the sequences encompassing the seven transmembrane domains of all *Dictyostelium* GPCRs, and selected GPCRs from other organisms, was used to create an unrooted dendrogram with the TreeView program. A black circle marks the innermost node of each branch supported by >60% bootstraps. The hash symbol indicates that this gene model has to be split, and

the asterisk indicates a putative pseudogene. DictyBase identifiers (DDB) were used for the newly discovered *Dictyostelium* receptors and SwissProt identifiers for all other receptors. A.th., *A. thaliana*; B.t., *Bos taurus*; CAR/CRL, CAMP receptor/CAMP receptor-like; C.e., *C. elegans*; DICD1, *D. discoideum*; D.m., *D. melanogaster*; G.c., *Geodia cydonium*; P.p., *Polysphondylium pallidum*; X.l., *Xenopus laevis*.

as that observed in Metazoa (Supplementary Tables 14–16 and Supplementary Fig. 16). Given the presence of SH2-domain-based signalling it was surprising that no receptor tyrosine kinases could be recognized in the genome. However, *Dictyostelium* has a number of other receptor kinases, such as the histidine kinases and a group of eight novel putative receptor serine/threonine kinases, which are involved in nutrient and starvation sensing⁶⁹. Most of the ubiquitous families of transcription factors are represented in *Dictyostelium*, with the notable exception of the otherwise ubiquitous basic helix–loop–helix proteins (Supplementary Table 17 and Supplementary Fig. 17). Compared with other eukaryotes, *Dictyostelium* appears to have fewer transcription factors relative to the total number of genes, suggesting that many transcription factors have yet to be defined, or that the activities of a smaller repertoire of factors are combined and controlled to achieve complex regulation (Supplementary Table 18 and Supplementary Fig. 18).

Cell adhesion

Throughout *Dictyostelium* development, cells must modulate their adhesiveness to the substrate, to the extracellular matrix and to other cells in order to create tissues and carry out morphogenesis. To accomplish this, *Dictyostelium* uses a surprising number of components that have been normally only associated with animals. For example, disintegrin proteins regulate cell adhesiveness and differentiation in a number of Metazoa, and at least one *Dictyostelium* disintegrin, AmpA, is needed throughout development for cell fate specification⁷⁰. We also identified distant relatives of vinculin and α -catenin—normally associated with adherens junctions—which support the idea that the epithelium-like sheet of cells that surrounds the stalk tube contains such junctions⁷¹. Consistent with this, the *Dictyostelium* genome encodes numerous proteins previously described as components of adherens junctions in Metazoa, such as β -catenin (Aardvark), α -actinin, formins, VASP and myosin VII.

In animals, tandem repeats of immunoglobulin, cadherin, fibronectin III or E-set domains are often present in cell adhesion proteins, although their common protein fold pre-dates the emergence of eukaryotes. EGF/laminin domains are also found in adhesion proteins but, before the analysis of the *Dictyostelium* genome, no non-metazoan was known to have more than two EGF repeats in a single predicted protein. *Dictyostelium* has 61 predicted proteins containing repeated E-set or EGF/laminin domains, and many of these contain additional domains that suggest they have roles in cell adhesion or cell recognition, such as mannose-6-phosphate receptor, fibronectin III, or growth factor receptor domains and transmembrane domains (Fig. 9). In support of this idea, four of these proteins (LagC, LagD, AmpA and ComC) have been shown to be required for cell adhesion and signalling during development^{70,72–74}.

Cellulose-based structures

During development, *Dictyostelium* cells produce a number of cellulose-based structural elements. *Dictyostelium* slugs synthesize an extracellular matrix, or sheath, around themselves that is comprised of proteins and cellulose. Several of the smaller sheath proteins bind cellulose and are believed to have a role in slug migration, whereas the larger, cysteine-rich EcmA protein is essential for full integrity of the sheath and for establishing correct slug shape^{75,76}. During terminal differentiation, cellulose is deposited in the stalk and in the cell walls of the stalk and spore cells^{77–79}. The first confirmed eukaryotic gene for cellulose synthase was discovered in *Dictyostelium* and this gene has since been recognized in many plants, *N. crassa* and the ascidian *Ciona intestinalis*⁸⁰. The fungal and urochordate enzymes are more closely related to the *Dictyostelium* homologue than to plant or bacterial cellulose synthases, indicating that the common ancestor of fungi and animals carried a gene for

cellulose synthase that was subsequently lost in most animals. The *Dictyostelium* genome encodes more than 40 additional proteins that are likely to be involved in cellulose synthesis or degradation, and are probably involved in the production and remodelling of cellulose fibres of the slug sheath, stalk tube and cell walls (see Supplementary Information).

The fundamental similarities in cellular cooperation found in *Dictyostelium* and in the Metazoa clearly resulted in a parallel positive selection for structural and regulatory genes required for cell motility, adhesion and signalling. *Dictyostelium* uses a set of signals and adhesion proteins that are distinct from those employed for similar purposes in Metazoa but, like the Metazoa, *Dictyostelium* has maintained a diversity of GPCRs, protein kinases and ABC transporters that enable it to respond to those signals. *Dictyostelium* has also retained and modified an organizational strategy perfected in plants, basing several structural elements on cellulose. At one

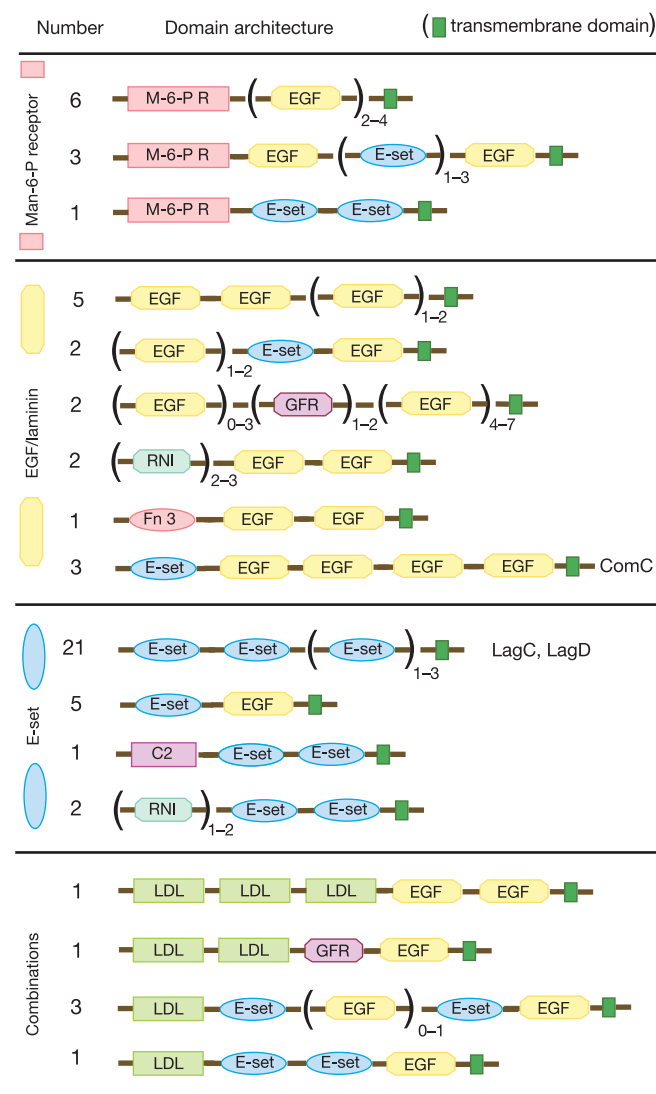


Figure 9 Putative adhesion/signalling proteins. Proteins containing repeated EGF/laminin and/or E-set SCOP Superfamily domains are classified into groups containing mannose-6-phosphate receptor, mainly EGF/laminin, mainly E-set, or combinations of domains. Most of these proteins have predicted transmembrane domains and so are expected to be cell surface proteins. ComC, LagC and LagD are proteins that have been characterized to have adhesion and/or signalling functions during multicellular development^{72–74}. C2, calcium-dependent lipid binding; Fn 3, fibronectin type III; GFR, growth factor receptor; LDL, L domain-like leucine-rich repeat; M-6-P R, mannose-6-phosphate receptor; RNI, RNI-like.

level *Dictyostelium* has achieved multicellularity by using strategies that are similar to plants and Metazoa, but the differences between them suggest convergent evolution, rather than lineal descent from an ancestor with overt or latent multicellular capacities.

Conclusion

The complete protein repertoire of *Dictyostelium* provides a new perspective for studying its cellular and developmental biology. At a systems level, *Dictyostelium* provides a level of complexity that is greater than the yeasts, but much simpler than plants or animals. Thus, high-resolution molecular analyses in this system may reveal control networks that are difficult to study in more complex systems, and may presage regulatory strategies used by higher organisms^{81–83}. At a practical level, the comparative genomics of *Dictyostelium* and related pathogens, such as *Entamoeba histolytica*, should aid in the functional definition of amoebozoan-specific genes that may open new avenues of research aimed at controlling amoebic diseases. *Dictyostelium*'s adeptness at hunting bacteria also renders it susceptible to infections by intracellular bacterial pathogens^{84,85}. *Dictyostelium* and human macrophages display fundamental similarities in their cell biology, which has spurred the use of *Dictyostelium* as a model host for bacterial pathogenesis. It is also an attractive model in which to study other disease processes: for a number of human disease-related proteins, it provides a test-bed for studying their functions in a model organism that has greater similarity to higher eukaryotes than do the yeasts, yet shares the latter's experimental tractability.

The high frequency of repeated amino acid tracts in *Dictyostelium* proteins has long been known anecdotally, but we can now survey their precise nature and number, and find them to be more abundant than in any other sequenced genome. Many human diseases result from the expansion of triplet nucleotide repeats, some of which encode polyglutamine tracts that cause cell degeneration^{86,87}. Learning how *Dictyostelium* cells tolerate so many proteins with amino acid homopolymers will, we hope, help to elucidate the roles of these motifs in protein function and dysfunction.

Comparative genomic studies in eukaryotes are providing the raw material for global examinations of the evolution of cellular regulation and developmental mechanisms⁸⁸. Many genes have been lost in one species but retained in others, such that each new genome sequence adds to our understanding of the genetic complement of the eukaryotic progenitor. Thus, our understanding of eukaryotes will continue to be refined as more genome sequences become available from representatives of large groups of organisms whose genomes remain largely unexplored, such as the amoebozoan. The surprising molecular diversity of the *Dictyostelium* proteome, which includes protein assemblages usually associated with fungi, plants or animals, suggests that their last common ancestor had a greater number of genes than had been previously appreciated. □

Methods

Details on the availability of reagents can be found in the Supplementary Information. All analyses described here were performed on Version 2.0 of the genome sequence. Updates to the sequence and annotation are available at <http://www.dictybase.org> and <http://www.genedb.org/genedb/dicty/index.jsp>. Further details of analyses not explicitly described below can be found in the Supplementary Information.

HAPPY mapping

A short-range (~100-kb), high-resolution (±8.54-kb) mapping panel was prepared as described⁹. Briefly, 96 aliquots each containing ±0.52 haploid genome equivalents of sheared AX4 genomic DNA were pre-amplified by PEP (primer extension pre-amplification⁸⁹). A total of 4,913 STS markers (Supplementary Table 1) were typed by two-phase hemi-nested polymerase chain reaction (PCR; multiplexed for up to 1,200 markers in the first phase) on aliquots of the diluted PEP products. Maps were assembled from good-quality data essentially as described previously⁹. A second, longer-range (±150 kb) mapping panel was used to confirm some linkages on chromosomes 2 and 5. HAPPY map analysis and PCR primer design for HAPPY mapping was performed using various custom programs (P.H.D. and A.T.B., unpublished).

Chromosome purification

Genomic DNA from *D. discoideum* strain AX4 was prepared and separated by pulsed field gel electrophoresis essentially as described^{27,9}, except that gels were run in stacked pairs; one member of each pair was stained with ethidium bromide, and bands excised from its unstained counterpart by alignment.

WCS and YAC subclone libraries

For WCS libraries, gel slices (above) were disrupted by several passages through a 30-gauge syringe needle, digested with β-agarase (NEB) and phenol-extracted. DNA was concentrated by ethanol precipitation, sonicated, end-blunted using mung bean nuclease and size-fractionated on 0.8% low-melting-point agarose gels. Fractions of 1.4–2 kb and 2–4 kb were excised, DNA extracted as before and ligated into the *Sma*I site of pUC18 or pUC19. Clone propagation and template preparation followed standard protocols.

For YAC subclone libraries, AX4-derived YACs were identified (and their position and integrity confirmed) by screening the set described by ref. 22 using markers from the HAPPY map. Subclones were prepared from PFG-purified YACs essentially as for the WCS libraries; contaminating yeast-derived sequences were filtered out *in silico*.

Sequencing and assembly

Details of the sequencing and assembly methods can be found in Supplementary Information. Generally, mapped sequence features were used to nucleate sequence contigs assembled from the WCS data, and extended using read-pair information and iterative searches for overlapping sequences, followed by directed gap closure using a range of approaches.

Fluorescent *in situ* hybridization

In situ hybridization was performed as in ref. 17.

Gene prediction and identification of sequence features

Full details are provided in the Supplementary Information. Briefly, automated gene prediction was performed using a combination of programs that had been trained on well-characterized *D. discoideum* genes, and the results integrated with reference to *D. discoideum* complementary DNA sequences and homology to genes in other species. Other features in the predicted proteins, and other sequence features, were identified using a variety of software packages.

Analysis of functional gene clustering

Microarray targets (refs 53, 90, 91; and N. Van Driessche and G. Shaulsky, unpublished data) and gene models were mapped onto the genome sequence using BLAST⁹² and the modified LIS algorithm⁹³. To look for clustering of genes with correlated temporal expression profiles, pairwise correlation coefficients were calculated for genes with known expression profiles on each chromosome⁹¹. Blocks of ≥6 consecutive genes were sought, for which either (1) all pairwise correlation coefficients were positive and ≥70% were >0.2 (genes with similar developmental trajectories) or (2) each gene had a partner with an absolute correlation coefficient value of >0.6 (tightly co-regulated genes); no statistically significant clusters met these criteria.

To look for clustering of genes associated with specific developmental stages^{94,95} or cell types^{90,96}, the genome was scanned with various sized windows⁹⁷ for regions with significant ($P < 0.01$) over-representation of genes in any one of these groups.

Analysis of duplicated genes

Predicted protein sequences were clustered using TribeMCL⁹⁸, using a BLASTP expectation of $<10^{-40}$ as a cutoff. A χ^2 test invalidated the hypothesis that members of a family are randomly distributed in the genome. Within each family, protein divergences (similarity distances computed using the 'ProtDist' module of PHYLIP; <http://evolution.genetics.washington.edu/phylip.html>) and physical intergenic distances between all pairs of family members were tabulated, and the correlation coefficient between the former and latter values was calculated. Analysis was performed on the 86 gene families (representing 155 gene pairs) with at least 10 intrachromosomal distance pairings to provide robust statistical confidence.

Other sequence analyses and graphical representation

Other sequence analyses (nucleotide and dinucleotide composition; identification of simple-sequence repeats in nucleotide and protein sequence; coding density computation; tRNA cluster identification) were performed using a range of custom software (P.H.D. and A.T.B., unpublished). Graphical representation of chromosomes in Fig. 2 was done primarily using Cinema4D-8.5 (Maxon Computer GmbH) after pre-processing using custom software (P.H.D.).

Received 16 September 2004; accepted 17 February 2005; doi:10.1038/nature03481.

1. Kessin, R. H. *Dictyostelium—Evolution, Cell Biology, and the Development of Multicellularity*, xiv, 294 (Cambridge Univ. Press, Cambridge, 2001).
2. Konijn, T. M. *et al.* The acrasin activity of adenosine-3',5'-cyclic phosphate. *Proc. Natl Acad. Sci. USA* **58**, 1152–1154 (1967).
3. Müller, K. & Gerisch, G. A specific glycoprotein as the target site of adhesion blocking Fab in aggregating *Dictyostelium* cells. *Nature* **274**, 445–449 (1978).
4. Raper, K. B. Pseudopodium formation and organization in *Dictyostelium discoideum*. *J. Elisha Mitchell Sci. Soc.* **56**, 241–282 (1940).

5. Raper, K. B. *Dictyostelium discoideum*, a new species of slime mold from decaying forest leaves. *J. Agr. Res.* **50**, 135–147 (1935).
6. Knecht, D. A., Cohen, S. M., Loomis, W. F. & Lodish, H. F. Developmental regulation of *Dictyostelium discoideum* actin gene fusions carried on low-copy and high-copy transformation vectors. *Mol. Cell. Biol.* **6**, 3973–3983 (1986).
7. Dear, P. H. & Cook, P. R. HAPPY mapping—linkage mapping using a physical analog of meiosis. *Nucleic Acids Res.* **21**, 13–20 (1993).
8. Konfortov, B. A., Cohen, H. M., Bankier, A. T. & Dear, P. H. A high-resolution HAPPY map of *Dictyostelium discoideum* chromosome 6. *Genome Res.* **10**, 1737–1742 (2000).
9. Glöckner, G. *et al.* Sequence and analysis of chromosome 2 of *Dictyostelium discoideum*. *Nature* **418**, 79–85 (2002).
10. Urushihara, H. *et al.* Analyses of cDNAs from growth and slug stages of *Dictyostelium discoideum*. *Nucleic Acids Res.* **32**, 1647–1653 (2004).
11. Smith, S. S. & Ratner, D. I. Lack of 5-methylcytosine in *Dictyostelium discoideum* DNA. *Biochem. J.* **277**, 273–275 (1991).
12. Glöckner, G. *et al.* The complex repeats of *Dictyostelium discoideum*. *Genome Res.* **11**, 585–594 (2001).
13. Crick, F. H. Codon-anticodon pairing: the wobble hypothesis. *J. Mol. Biol.* **19**, 548–555 (1966).
14. Soll, D. & Rajbhandary, U. (ed.) *tRNA: Structure, Biosynthesis and Function* (ASM, Washington DC, 1995).
15. Burger, G., Plante, I., Lonergan, K. M. & Gray, M. W. The mitochondrial DNA of the amoeboid protozoan, *Acanthamoeba castellanii*: complete sequence, gene content and genome organization. *J. Mol. Biol.* **3**, 522–537 (1995).
16. Ogawa, S. *et al.* The mitochondrial DNA of *Dictyostelium discoideum*: complete sequence, gene content and genome organization. *Mol. Gen. Genet.* **263**, 514–519 (2000).
17. Sugcan, R. *et al.* Sequence and structure of the extrachromosomal palindrome encoding the ribosomal RNA genes in *Dictyostelium*. *Nucleic Acids Res.* **31**, 2361–2368 (2003).
18. Szafranski, K., Dingermann, T., Glöckner, G. & Winckler, T. Template jumping by a LINE reverse transcriptase has created a SINE-like 5S rRNA retropseudogene in *Dictyostelium*. *Mol. Genet. Genom.* **271**, 98–102 (2004).
19. Eichler, E. E. & Sankoff, D. Structural dynamics of eukaryotic chromosome evolution. *Science* **301**, 793–797 (2003).
20. Cappello, J., Cohen, S. M. & Lodish, H. F. *Dictyostelium* transposable element DIRS-1 preferentially inserts into DIRS-1 sequences. *Mol. Cell. Biol.* **4**, 2207–2213 (1984).
21. Cappello, J., Handelsman, K. & Lodish, H. F. Sequence of *Dictyostelium* DIRS-1: an apparent retrotransposon with inverted terminal repeats and an internal circle junction sequence. *Cell* **43**, 105–115 (1985).
22. Loomis, W. F., Welker, D., Hughes, J., Maghakian, D. & Kuspa, A. Integrated maps of the chromosomes in *Dictyostelium discoideum*. *Genetics* **141**, 147–157 (1995).
23. Goodwin, T. J. & Poulter, R. T. Multiple LTR-retrotransposon families in the asexual yeast *Candida albicans*. *Genome Res.* **10**, 174–191 (2000).
24. Appeltren, H., Knöhl, B. & Ekwall, K. Distinct centromere domain structures with separate functions demonstrated in live fission yeast cells. *J. Cell Sci.* **116**, 4035–4042 (2003).
25. Kuspa, A., Maghakian, D., Bergesch, P. & Loomis, W. F. Physical mapping of genes to specific chromosomes in *Dictyostelium discoideum*. *Genomics* **13**, 49–61 (1992).
26. McClintock, B. The production of homozygous deficient tissues with mutant characteristics by means of the aberrant mitotic behaviour of ring-shaped chromosomes. *Genetics* **23**, 315–376 (1938).
27. Cox, E. C., Vocke, C. D., Walter, S., Gregg, K. Y. & Bain, E. S. Electrophoretic karyotype for *Dictyostelium discoideum*. *Proc. Natl Acad. Sci. USA* **87**, 8247–8251 (1990).
28. International Human Genome Sequencing Consortium. Initial sequencing and analysis of the human genome. *Nature* **409**, 860–941 (2001).
29. Gardner, M. J. *et al.* Genome sequence of the human malaria parasite *Plasmodium falciparum*. *Nature* **419**, 498–511 (2002).
30. Trusov, Y. A. & Dear, P. H. A molecular clock based on the expansion of gene families. *Nucleic Acids Res.* **24**, 995–999 (1996).
31. Kellis, M., Birren, B. W. & Lander, E. S. Proof and evolutionary analysis of ancient genome duplication in the yeast *Saccharomyces cerevisiae*. *Nature* **428**, 617–624 (2004).
32. Dujon, B. *et al.* Genome evolution in yeasts. *Nature* **430**, 35–44 (2004).
33. Raper, K. B. *The Dictyostelids* (Princeton Univ. Press, Princeton, New Jersey, 1984).
34. Loomis, W. F. & Smith, D. W. Consensus phylogeny of *Dictyostelium*. *Experientia* **51**, 1110–1115 (1995).
35. Baldauf, S. L. & Doolittle, W. F. Origin and evolution of the slime molds (Mycetozoa). *Proc. Natl Acad. Sci. USA* **94**, 12007–12012 (1997).
36. Baptiste, E. *et al.* The analysis of 100 genes supports the grouping of three highly divergent amoebae: *Dictyostelium*, *Entamoeba*, and *Mastigamoeba*. *Proc. Natl Acad. Sci. USA* **99**, 1414–1419 (2002).
37. Nara, T., Hshimoto, T. & Aoki, T. Evolutionary implications of the mosaic pyrimidine-biosynthetic pathway in eukaryotes. *Gene* **257**, 209–222 (2000).
38. Olsen, R. & Loomis, W. F. A model of orthologous protein sequence divergence. *J. Mol. Evol.* (in the press).
39. Thomason, P. & Kay, R. Eukaryotic signal transduction via histidine-aspartate phosphorelay. *J. Cell Sci.* **113**, 3141–3150 (2000).
40. Fortini, M. E., Skupski, M. P., Boguski, M. S. & Hariharan, I. K. A survey of human disease gene counterparts in the *Drosophila* genome. *J. Cell Biol.* **150**, F23–F30 (2000).
41. Jain, R., Rivera, M. C., Moore, J. E. & Lake, J. A. Horizontal gene transfer in microbial genome evolution. *Theor. Popul. Biol.* **61**, 489–495 (2002).
42. Richards, T. A., Hirt, R. P., Williams, B. A. & Embley, T. M. Horizontal gene transfer and the evolution of parasitic protozoa. *Protist* **154**, 17–32 (2003).
43. Iyer, L. M., Aravind, L., Coon, S. L., Klein, D. C. & Koonin, E. V. Evolution of cell-cell signaling in animals: did late horizontal gene transfer from bacteria have a role? *Trends Genet.* **20**, 292–299 (2004).
44. Myllykallio, H. *et al.* An alternative flavin-dependent mechanism for thymidylate synthesis. *Science* **297**, 105–107 (2002).
45. Kessin, R. H., Gundersen, G. G., Zaydfudim, V., Grimson, M. & Blanton, R. L. How cellular slime molds evade nematodes. *Proc. Natl Acad. Sci. USA* **93**, 4857–4861 (1996).
46. Morris, H. R., Taylor, G. W., Masento, M. S., Jermyn, K. A. & Kay, R. R. Chemical structure of the morphogen differentiation inducing factor from *Dictyostelium discoideum*. *Nature* **328**, 811–814 (1987).
47. Cane, D. E., Walsh, C. T. & Khosla, C. Harnessing the biosynthetic code: combinations, permutations, and mutations. *Science* **282**, 63–68 (1998).
48. Holland, I. B. & Blight, M. A. ABC-ATPases, adaptable energy generators fuelling transmembrane movement of a variety of molecules in organisms from bacteria to humans. *J. Mol. Biol.* **293**, 381–399 (1999).
49. Andrade, A. C., Van Nistelrooy, J. G., Peery, R. B., Skatrud, P. L. & De Waard, M. A. The role of ABC transporters from *Aspergillus nidulans* in protection against cytotoxic agents and in antibiotic production. *Mol. Gen. Genet.* **263**, 966–977 (2000).
50. Mendez, C. & Salas, J. A. The role of ABC transporters in antibiotic-producing organisms: drug secretion and resistance mechanisms. *Res. Microbiol.* **152**, 341–350 (2001).
51. Schoonbeek, H. J., Raaijmakers, J. M. & De Waard, M. A. Fungal ABC transporters and microbial interactions in natural environments. *Mol. Plant Microbe Interact.* **15**, 1165–1172 (2002).
52. Anjard, C. & Loomis, W. F. Evolutionary analyses of ABC transporters of *Dictyostelium discoideum*. *Eukaryot. Cell* **1**, 643–652 (2002).
53. Iranfar, N., Fuller, D. & Loomis, W. F. Genome-wide expression analyses of gene regulation during early development of *Dictyostelium discoideum*. *Eukaryot. Cell* **2**, 664–670 (2003).
54. Devreotes, P. N. & Zigmond, S. H. Chemotaxis in eukaryotic cells: A focus on leukocytes and *Dictyostelium*. *Annu. Rev. Cell Biol.* **4**, 649–686 (1988).
55. Noegel, A. A. & Schleicher, M. The actin cytoskeleton of *Dictyostelium*: a story told by mutants. *J. Cell Sci.* **113**, 759–766 (2000).
56. Rivero, F. & Somesh, B. P. Signal transduction pathways regulated by Rho GTPases in *Dictyostelium*. *J. Muscle Res. Cell Motil.* **23**, 737–749 (2002).
57. Merlot, S. & Firtel, R. A. Leading the way: directional sensing through phosphatidylinositol 3-kinase and other signalling pathways. *J. Cell Sci.* **116**, 3471–3478 (2003).
58. Bockaert, J. & Pin, J. P. Molecular tinkering of G protein-coupled receptors: an evolutionary success. *EMBO J.* **18**, 1723–1729 (1999).
59. Ginsburg, G. T. *et al.* The regulation of *Dictyostelium* development by transmembrane signalling. *J. Eukaryot. Microbiol.* **42**, 200–205 (1995).
60. Raisley, B., Zhang, M., Herold, D. & Hadwiger, J. A. A cAMP receptor-like G protein-coupled receptor with roles in growth regulation and development. *Dev. Biol.* **265**, 433–445 (2004).
61. King, N., Hittinger, C. T. & Carroll, S. B. Evolution of key cell signaling and adhesion protein families predates animal origins. *Science* **301**, 361–363 (2003).
62. Kawata, T. *et al.* SH2 signaling in a lower eukaryote: A STAT protein that regulates stalk cell differentiation in *Dictyostelium*. *Cell* **89**, 909–916 (1997).
63. Thien, C. B. & Langdon, W. Y. Cbl: many adaptations to regulate protein tyrosine kinases. *Nature Rev. Mol. Cell Biol.* **2**, 294–307 (2001).
64. Meng, W. *et al.* Structure of the amino-terminal domain of Cbl complexed to its binding site on ZAP-70 kinase. *Nature* **398**, 84–90 (1999).
65. Shaulsky, G., Kuspa, A. & Loomis, W. F. A multidrug resistance transporter serine protease gene is required for prestalk specialization in *Dictyostelium*. *Genes Dev.* **9**, 1111–1122 (1995).
66. Anjard, C., Zeng, C., Loomis, W. F. & Nellen, W. Signal transduction pathways leading to spore differentiation in *Dictyostelium discoideum*. *Dev. Biol.* **193**, 146–155 (1998).
67. Good, J. R. *et al.* TagA, a putative serine protease/ABC transporter of *Dictyostelium* that is required for cell fate determination at the onset of development. *Development* **130**, 2953–2965 (2003).
68. Good, J. R. & Kuspa, A. Evidence that a cell-type-specific efflux pump regulates cell differentiation in *Dictyostelium*. *Dev. Biol.* **220**, 53–61 (2000).
69. Chibalina, M. V., Anjard, C. & Insall, R. H. Gdt2 regulates the transition of *Dictyostelium* cells from growth to differentiation. *BMC Dev. Biol.* **4**, 8 (2004).
70. Blumberg, D. D., Ho, H. N., Petty, C. L., Varney, T. R. & Gandham, S. AmpA, a modular protein containing disintegrin and ornatin domains, has multiple effects on cell adhesion and cell fate specification. *J. Muscle Res. Cell Motil.* **23**, 817–828 (2002).
71. Grimson, M. J. *et al.* Adherens junctions and β -catenin-mediated cell signalling in a non-metazoan organism. *Nature* **408**, 727–731 (2000).
72. Dynes, J. L. *et al.* LagC is required for cell-cell interactions that are essential for cell-type differentiation in *Dictyostelium*. *Genes Dev.* **8**, 948–958 (1994).
73. Wang, J. *et al.* The membrane glycoprotein gp150 is encoded by the *lagC* gene and mediates cell-cell adhesion by heterophilic binding during *Dictyostelium* development. *Dev. Biol.* **227**, 734–745 (2000).
74. Kibler, K., Svetz, J., Nguyen, T. L., Shaw, C. & Shaulsky, G. A cell-adhesion pathway regulates intercellular communication during *Dictyostelium* development. *Dev. Biol.* **264**, 506–521 (2003).
75. Morrison, A. *et al.* Disruption of the gene encoding the EcmA, extracellular matrix protein of *Dictyostelium* alters slug morphology. *Dev. Biol.* **163**, 457–466 (1994).
76. Wang, Y. Z., Slade, M. B., Gooley, A. A., Atwell, B. J. & Williams, K. L. Cellulose-binding modules from extracellular matrix proteins of *Dictyostelium discoideum* stalk and sheath. *Eur. J. Biochem.* **268**, 4334–4345 (2001).
77. Freeze, H. & Loomis, W. F. Chemical analysis of stalk components of *Dictyostelium discoideum*. *Biochim. Biophys. Acta* **539**, 529–537 (1978).
78. Zhang, P., McGlynn, A., Loomis, W. F., Blanton, R. L. & West, C. M. Spore coat formation and timely sporulation depend on cellulose in *Dictyostelium*. *Differentiation* **67**, 72–79 (2001).
79. West, C. M., Zhang, P., McGlynn, A. C. & Kaplan, L. Outside-in signaling of cellulose synthesis by a spore coat protein in *Dictyostelium*. *Eukaryot. Cell* **1**, 281–292 (2002).
80. Blanton, R. L., Fuller, D., Iranfar, N., Grimson, M. J. & Loomis, W. F. The cellulose synthase gene of *Dictyostelium*. *Proc. Natl Acad. Sci. USA* **97**, 2391–2396 (2000).
81. Thomason, P. A. *et al.* An intersection of the cAMP/PKA and two-component signal transduction systems in *Dictyostelium*. *EMBO J.* **17**, 2838–2845 (1998).
82. Maeda, M. *et al.* Periodic signaling controlled by an oscillatory circuit that includes protein kinases ERK2 and PKA. *Science* **304**, 875–878 (2004).
83. Soler-Lopez, M. *et al.* Structure of an activated *Dictyostelium* STAT in its DNA-unbound form. *Mol. Cell* **13**, 791–804 (2004).
84. Solomon, J. M., Rupper, A., Cardelli, J. A. & Isberg, R. R. Intracellular growth of *Legionella pneumophila* in *Dictyostelium discoideum*, a system for genetic analysis of host-pathogen interactions. *Infect. Immun.* **68**, 2939–2947 (2000).
85. Skriwan, C. *et al.* Various bacterial pathogens and symbionts infect the amoeba *Dictyostelium discoideum*. *Int. J. Med. Microbiol.* **291**, 615–624 (2002).
86. Zoghbi, H. Y. & Orr, H. T. Glutamine repeats and neurodegeneration. *Annu. Rev. Neurosci.* **23**, 217–247 (2000).

87. Brown, L. Y. & Brown, S. A. Alanine tracts: the expanding story of human illness and trinucleotide repeats. *Trends Genet.* **20**, 51–58 (2004).
88. Rubin, G. M. *et al.* Comparative genomics of the eukaryotes. *Science* **287**, 2204–2215 (2000).
89. Zhang, L. *et al.* Whole genome amplification from a single cell: Implications for genetic analysis. *Proc. Natl Acad. Sci. USA* **89**, 5847–5851 (1992).
90. Iranfar, N. *et al.* Expression patterns of cell-type-specific genes in *Dictyostelium*. *Mol. Biol. Cell* **12**, 2590–2600 (2001).
91. Van Driessche, N. *et al.* A transcriptional profile of multicellular development in *Dictyostelium discoideum*. *Development* **129**, 1543–1552 (2002).
92. Altschul, S. F., Gish, W., Miller, W., Myers, E. W. & Lipman, D. J. Basic local alignment search tool. *J. Mol. Biol.* **215**, 403–410 (1990).
93. Zhang, H. Alignment of BLAST high-scoring segment pairs based on the longest increasing subsequence algorithm. *Bioinformatics* **19**, 1391–1396 (2003).
94. Katoh, M. *et al.* An orderly retreat: dedifferentiation is a regulated process. *Proc. Natl Acad. Sci. USA* **101**, 7005–7010 (2004).
95. Xu, Q. *et al.* Transcriptional transitions during *Dictyostelium* spore germination. *Eukaryot. Cell* **3**, 1101–1110 (2004).
96. Maeda, M. *et al.* Changing patterns of gene expression in *Dictyostelium* prestalk cell subtypes recognized by *in situ* hybridization with genes from microarray analyses. *Eukaryot. Cell* **2**, 627–637 (2003).
97. Cohen, B. A., Mitra, R. D., Hughes, J. D. & Church, G. M. A computational analysis of whole-genome expression data reveals chromosomal domains of gene expression. *Nature Genet.* **26**, 183–186 (2000).
98. Enright, A. J., Van Dongen, S. & Ouzounis, C. A. An efficient algorithm for large-scale detection of protein families. *Nucleic Acids Res.* **30**, 1575–1584 (2002).

Supplementary Information accompanies the paper on www.nature.com/nature.

Acknowledgements Sequencing and analysis of chromosomes 1, 2 and 3 were supported by grants from the DFG and by Köln Fortune, and that of chromosomes 4, 5 and 6 in the USA by grants from NICHD/NIH. Work in the UK/European Union (EU) was supported by a

programme grant from the MRC to J.W., R.R.K., B.B. and P.H.D., and by the EU. Analyses at dictyBase were supported by a grant from the NIGMS/NIH to R.L.C. The *Dictyostelium* cDNA project was supported by Research for the Future of JSPS and by Grants-in-Aid for Scientific Research on Priority Areas of MEXT of Japan. The German team wishes to thank S. Förste, N. Zeisse, S. Rothe, S. Landmann, R. Schultz, C. Neuhoff and R. Müller for technical assistance. The US team thanks S. Kaminsky, S. Klein and T. Hewitt for their scientific foresight in the early stages of this project and for their support of *Dictyostelium* as a model system, and H. Hosak, O. Delgado, L. Lewis, K. Hamilton, J. Hume, C. Kovar Smith, D. Neal, P. Havlak, K. J. Durbin and P. Burch of the HGSC at Baylor College of Medicine. R.S. thanks L. Cortez, E. Joyner and B. Hill for their assistance during the course of the project. C.B. thanks M. Veron and P. Glaser for their support and discussions. P.H.D. thanks H. O'Hare for early involvement in the project, A. Ivens for discussions, and the MRC Centre Visual Aids department, Cambridge, for advice on graphics. We also thank S. Bowman and D. Lawson for their contribution to the EUDICT region in the initial stages of the project, and D. Martin at the Wellcome Trust Biocentre, University of Dundee, for running the GOTcha search of our gene models. The Japanese cDNA project thanks N. Ogasawara and I. Takeuchi for comments and encouragement, and others who participated in earlier stages of the project.

Author contributions M. Platzer, R. R. Kay, J. Williams, P. H. Dear, A. A. Noegel, B. Barrell and A. Kuspa are co-senior authors.

Competing interests statement The authors declare that they have no competing financial interests.

Correspondence and requests for materials should be addressed to P.H.D. (phd@mrc-lmb.cam.ac.uk). Sequence data for the genome were deposited in the GenBank nucleotide database under the project accession number AAFI00000000. The six chromosomal assemblies have accession numbers CM000150–CM000155, and their component sequence contigs have accession numbers AAFI01000001–AAFI01000336.