

# Classification and functional annotation of eukaryotic protein kinases

Diego Miranda-Saavedra and Geoffrey J. Barton\*

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

## ABSTRACT

Reversible protein phosphorylation by protein kinases and phosphatases is a ubiquitous signaling mechanism in all eukaryotic cells. A multilevel hidden Markov model library is presented which is able to classify protein kinases into one of 12 families, with a misclassification rate of zero on the characterized kinomes of *H. sapiens*, *M. musculus*, *D. melanogaster*, *C. elegans*, *S. cerevisiae*, *D. discoideum*, and *P. falciparum*. The Library is shown to outperform BLASTP and a general Pfam hidden Markov model of the kinase catalytic domain in the retrieval and family-level classification of protein kinases. The application of the Library to the 38 unclassified kinases of yeast enriches the yeast kinome in protein kinases of the families AGC (5), CAMK (17), CMGC (4), and STE (1), thereby raising the family-level classification of yeast conventional protein kinases from 66.96 to 90.43%. The application of the Library to 21 eukaryotic genomes shows seven families (AGC, CAMK, CK1, CMGC, STE, PIKK, and RIO) to be present in all genomes analyzed, and so is likely to be essential to eukaryotes. Putative tyrosine kinases (TKs) are found in the plants *A. thaliana* (2), *O. sativa* ssp. *Indica* (6), and *O. sativa* ssp. *Japonica* (7), and in the amoeba *E. histolytica* (7). To our knowledge, TKs have not been predicted in plants before. This also suggests that a primitive set of TKs might have predated the radiation of eukaryotes. Putative tyrosine kinase-like kinases (TKLs) are found in the fungi *C. neoformans* (2), *P. chrysosporium* (4), in the Apicomplexans *C. hominis* (4), *P. yoelii* (4), and *P. falciparum* (6), the amoeba *E. histolytica* (109), and the alga *T. pseudonana* (6). TKLs are found to be abundant in plants (776 in *A. thaliana*, 1010 in *O. sativa* ssp. *Indica*, and 969 in *O. sativa* ssp. *Japonica*). TKLs might have predated the radiation of eukaryotes too and have been lost secondarily from some fungi. The application of the Library facilitates the annotation of kinomes and has provided novel insights on the early evolution and subsequent adaptations of the various protein kinase families in eukaryotes.

Proteins 2007; 68:893–914.  
© 2007 Wiley-Liss, Inc.

**Key words:** database searching; signal transduction; phosphorylation; hidden Markov model; evolution.

## INTRODUCTION

Reversible protein phosphorylation by protein kinases and phosphatases is thought to regulate virtually every cellular activity.<sup>1</sup> Since abnormal levels of phosphorylation are known to be responsible for severe diseases,<sup>2</sup> there is considerable therapeutic promise in obtaining a detailed understanding of phosphorylation events, both within specific cell types and in an evolutionary context. A thorough understanding of the evolution of protein kinases will help decipher how signaling events have shaped the development, pathology, and biochemistry of eukaryotes and so lead to a better description of the biochemical circuitry of cells which may guide the development of more effective drugs.<sup>3</sup>

The KinBase resource (<http://www.kinase.com/kinbase/>),<sup>4</sup> reflects the currently accepted classification of eukaryotic protein kinases, which are split into two broad groups: “conventional” protein kinases (ePKs) and “atypical” protein kinases (aPKs). ePKs are the largest group, and have been subclassified into eight families by examining sequence similarity between catalytic domains, the presence of accessory domains, and by considering any known modes of regulation.<sup>5,6</sup> The eight ePK families defined in KinBase are: the AGC family (including cyclic-nucleotide and calcium-phospholipid-dependent kinases, ribosomal S6-phosphorylating kinases, G protein-coupled kinases, and all close relatives of these groups); the CAMKs (calmodulin-regulated kinases); the CK1 family (casein kinase 1, and close relatives); the CMGC family (including cyclin-dependent kinases, mitogen-activated protein kinases, CDK-like kinases, and glycogen synthase kinase); the RGC family (receptor guanylate cyclase kinases, which are similar in domain sequence to tyrosine kinases, TKs); the STE family (including many kinases functioning in MAP kinase cascades); the TK

The Supplementary Material referred to in this article can be found online at <http://www.interscience.wiley.com/jpages/0887-3585/suppmat/>  
Grant sponsor: Wellcome Trust.

\*Correspondence to: Geoffrey Barton, School of Life Sciences, University of Dundee, Dow St, Dundee DD1 5EH, Scotland, UK.

E-mail: [geoff@compbio.dundee.ac.uk](mailto:geoff@compbio.dundee.ac.uk)

Received 12 May 2006; Revised 8 November 2006; Accepted 24 January 2007

Published online 7 June 2007 in Wiley InterScience (www.interscience.wiley.com).

DOI: 10.1002/prot.21444

family (tyrosine kinases); and the TKL family (tyrosine kinase-like kinases (TKLs), a diverse group resembling TK but which are in fact serine-threonine kinases). A ninth group, called the “Other” group, consists of a mixed collection of kinases that could not be classified easily into the previous families.

The aPKs are a small set of protein kinases that do not share clear sequence similarity with ePKs, but have been shown experimentally to have protein kinase activity<sup>6</sup> and comprise the following families: Alpha (exemplified by myosin heavy chain kinase of *Dictyostelium discoideum*); PIKK (phosphatidylinositol 3′ kinase-related kinases); PHDK (pyruvate dehydrogenase kinases); RIO (“right open reading frame” as it was one of two adjacent genes that were found to be transcribed divergently from the same intergenic region<sup>7</sup>; BRD (bromodomain-containing kinases); ABC1 (ABC1 domain-containing kinases); and TIF1 (transcriptional intermediary factor 1). The aPKs also include H11 (a homolog of gene ICP10 of Herpes simplex virus), FASTK (Fas-activated serine-threonine kinase), G11 (reported kinase activity against alpha casein and histones), BCR (fused with Abl in chronic myelogenous leukemia), TAF (TATA binding factor-associated factor), and A6 (two genes in human, A6 and A6r). However, only the Alpha, PIKK, PHDK, and RIO families have strong experimental evidence for kinase activity. For the RIO and Alpha families, X-ray crystallography has revealed clear similarities to the ePK kinase fold.<sup>8,9</sup>

Entries in KinBase are filtered by stringent criteria, including verification by cDNA cloning, in order to reduce the possibility of incorrect classification. As a consequence, KinBase is favored by experimentalists working on kinases and signal transduction pathways. A disadvantage of KinBase is that it is only cross-referenced to a variety of genome-specific databases such as the *Saccharomyces* Genome Database rather than to universal databases such as UniProt.<sup>10</sup>

To obtain a clearer picture of the spectrum of kinase genes and so understand the events that have shaped the evolution of protein kinase function in various lineages, it is necessary to identify the protein kinase complements (the “kinomes”) of as many organisms as possible. The kinomes of several organisms have previously been determined by a combination of *in silico* and wet-lab studies, including those for *C. elegans*,<sup>11</sup> *S. cerevisiae*,<sup>12</sup> *D. melanogaster*,<sup>4,13</sup> *H. sapiens*,<sup>6</sup> *M. musculus*,<sup>14</sup> *P. falciparum*,<sup>15</sup> and *D. discoideum*.<sup>16</sup> Identification of protein kinase sequences in these organisms has typically relied on BLAST searches<sup>17</sup> and hidden Markov model (HMM) profiles designed to identify kinase catalytic domains.<sup>18</sup> Finer-grained classification into the 12 subfamilies has been achieved by the application of consensus sequences for the catalytic and activation loops (subdomains VIB and VIII,<sup>19</sup>), and by clustering on sequence similarity to previously classified protein kinases.<sup>6,15</sup> However, classifying all protein kinases in a newly sequenced genome by

these methods has proved very time-consuming. The growth in eukaryotic genome sequencing has highlighted the need for more efficient and reliable methods of kinase classification.

Profile HMMs are statistical descriptions of sequence conservation from multiple sequence alignments,<sup>20</sup> and have been shown to outperform standard pairwise sequence comparison methods, both in terms of sensitivity and specificity.<sup>21</sup> HMMs form the basis of protein family and domain description libraries such as SUPERFAMILY<sup>22</sup> and Pfam,<sup>18</sup> which have been used successfully in the automatic annotation of many genomes.<sup>23</sup> Brown et al. have described formally the division of protein families into subfamilies for the purpose of achieving finer classification with HMMs.<sup>24</sup> Brown et al. found that a better separation between subfamily members is possible by using subfamily HMMs rather than a single HMM for the entire family, and that homologues are recognized with stronger scores, and nonhomologues are rejected with larger E-values.

In this article, the development of a multilevel HMM library from functional subfamilies of eukaryotic protein kinases is described and applied to the classification of protein kinases into the 12 subfamilies for 21 eukaryotic genomes.

## METHODS

### Source of kinase sequences

KinBase (<http://www.kinase.com/kinbase/>)<sup>4</sup> was the source for the annotated protein kinase complement of *H. sapiens*, *M. musculus*, *D. melanogaster*, *C. elegans*, *S. cerevisiae*, and *D. discoideum* (Table I). The protein kinase complement of *P. falciparum* was kindly made available by Professor Christian Doerig (Wellcome Trust Centre for Molecular Parasitology, Glasgow, UK).

Since KinBase does not map its entries onto Swiss-Prot/TrEMBL, this mapping was performed for sequences from *H. sapiens*, *D. melanogaster*, *S. cerevisiae*, and *C. elegans* (Swiss-Prot/TrEMBL Release 47, July 2004). The total number of protein kinases for *D. melanogaster*, *C. elegans*, *H. sapiens*, and *S. cerevisiae* in KinBase was 1487, which included 107 pseudogenes. Although it has been noted that many kinase pseudogenes are transcribed and could have a residual or scaffolding function,<sup>6</sup> kinase pseudogenes were not mapped onto Swiss-Prot/TrEMBL since many of them are partial transcripts or have non-amino acid characters in their sequences that indicate STOP codons.

The mapping of KinBase sequences onto Swiss-Prot/TrEMBL was performed by BLASTP-searching each KinBase sequence against Swiss-Prot/TrEMBL (Release 47, July 2004). Each of the 1380 result files was parsed to retrieve the top hit belonging to the same species as the query sequence. Three human sequences (SgK110, SgK424, and Slob) and two worm sequences (D1073.1 and aSWk029) were found to have no corresponding sequence

**Table I**

The ePK Complement of *H. sapiens*, *C. elegans*, *D. melanogaster*, *M. musculus*, *S. cerevisiae*, and *D. discoideum* as Described in KinBase (<http://www.kinase.com/kinbase/>)

Family	<i>H. sapiens</i>	<i>C. elegans</i>	<i>D. melanogaster</i>	<i>S. cerevisiae</i>	<i>M. musculus</i>	<i>D. discoideum</i>	Total
AGC	63	30	30	17	60	22	222
CAMK	74	40	32	21	96	21	284
CK1	12	85	10	4	11	3	125
CMGC	61	49	33	21	60	30	254
RGC	5	27	6	N/A	7	N/A	45
STE	47	25	18	14	47	45	196
TK	90	88	31	N/A	90	N/A	299
TKL	43	15	17	N/A	43	68	186
Total	395	359	177	77	414	189	1611

Protein kinase families: AGC (including cyclic-nucleotide and calcium-phospholipid-dependent kinases, ribosomal S6-phosphorylating kinases, G protein-coupled kinases, and all close relatives of these groups); the CAMKs (calmodulin-regulated kinases); the CK1 family (casein kinase 1, and close relatives); the CMGC family (including cyclin-dependent kinases, mitogen-activated protein kinases, CDK-like kinases, and glycogen synthase kinase); the RGC family (receptor guanylate cyclase kinases, which are similar in domain sequence to tyrosine kinases); the STE family (including many kinases functioning in MAP kinase cascades); the TK family (tyrosine kinases); and the TKL family (tyrosine kinase-like kinases, a diverse group resembling TK but which are in fact serine-threonine kinases). This classification excludes the "Other" group, which consists of kinases apparently not easy to classify into any of the groups below. *S. cerevisiae* lacks kinases from the groups: receptor guanylate cyclase kinase (RGC), tyrosine kinase (TK), and tyrosine kinase-like (TKL).

in Swiss-Prot/TrEMBL. 993 (75.4%) of KinBase ePK sequences could be mapped directly to identical sequences in Swiss-Prot/TrEMBL, whereas 162 (12.3%) of ePKs were mapped to sequences of nonidentical length but whose alignment contained no gaps. A further 161 (12.2%) of KinBase ePK sequences matched to Swiss-Prot/TrEMBL polypeptides of nonidentical length but the alignment contained gaps. This group encompasses pairs of alternatively spliced forms of kinase sequences, or pairs of very similar sequences with major DNA sequencing errors. All aPK sequences from KinBase were found to map the identical sequences in Swiss-Prot/TrEMBL.

KinBase contains full-length kinase sequences, but does not annotate the location of the catalytic domain in all sequences. To retrieve the catalytic domains of ePKs from *H. sapiens*, *D. melanogaster*, *C. elegans*, and *S. cerevisiae*, HMMs were generated from known human kinase catalytic domains to represent each family of kinases, one HMM per family. These HMMs were used to scan the full-length protein kinase sequences from each species. All kinase catalytic domains retrieved in this way were inspected and found to have the sequence features typical of the kinase catalytic domain.<sup>5</sup>

The catalytic domains of the aPKs, PIKK, Alpha, and RIO were retrieved by scanning the full-length protein sequences with the Pfam HMMs PF00454, PF02816, and PF01163, respectively. The catalytic domains of PDHKs were manually extracted from a multiple alignment by following the description given in Ref. 25.

Table II summarizes the aPK complement of *H. sapiens*, *M. musculus*, *S. cerevisiae*, *D. melanogaster*, and *C. elegans*, and the mapping of aPKs to the Swiss-Prot database.

### Library 1: Single HMM per family, and validation

A set of 12 profile HMMs<sup>20</sup> was built, one for each of the families of eight ePK and four aPK kinases from an or-

ganism (the "training set"), and tested to see if it could classify correctly the kinases of a different organism (the "test set"). Models from *H. sapiens*, *C. elegans*, and *D. melanogaster* were used as both training and test sets; however, since *S. cerevisiae* does not possess kinases of the RGC, TK, or TKL families, the ePKs of *S. cerevisiae* were only used as a test set. For the aPKs, only aPK catalytic domains from *H. sapiens* and *C. elegans* were used to build HMMs, since these organisms contain all four families of aPKs (RIO, Alpha, PDHK, and PIKK), whereas *D. melanogaster* and *S. cerevisiae* lack alpha kinases.<sup>4</sup>

The HMMs representing the families of ePKs and aPKs were created from multiple sequence alignments of all the protein kinase catalytic domain sequences in a family, generated by the AMPS suite of programs.<sup>26</sup> Since an HMM profile is only as reliable as the underlying alignment from which it is derived the quality of the alignment was assessed by Z-score analysis. In benchmarks, alignments that gave Z-scores above 5.0 S.D. had better than 70% accuracy within the core secondary structural regions of the proteins. All 32 ePK and aPK alignments gave Z-scores of >19 (data not shown), suggesting that the alignments are likely to be of high quality and so suitable for the derivation of HMMs. As a further test of quality, each alignment was inspected by eye to verify conservation of the core kinase catalytic domain motifs.<sup>5</sup> HMMs were built with the HMMER suite of programs (<http://hmmer.wustl.edu>),<sup>27</sup> HMMER version 2.1.1). The hmmbuild program was set to build HMMs that could identify one or more non-overlapping alignments to the complete model (multiple global alignments with respect to the model, and local with respect to the sequence) and the HMMs were calibrated by the hmmcalibrate program.

The classification performance of the family-specific models was tested by searching against the full-length protein kinase sequences of the characterized kinomes

**Table II***The Revised Classification of Eukaryotic Atypical Protein Kinases (aPKs)*

Organism/family	Protein length	KinBase accession number	Swiss-Prot accession number	Notes
<i>S. cerevisiae</i> /PIKK	2368	6319612_MEC1	P38111	ESR1 protein
<i>S. cerevisiae</i> /PIKK	2787	TEL1	P38110	Telomere length regulation protein TEL1
<i>S. cerevisiae</i> /PIKK	2470	TOR1	P35169	PI3K-related kinase TOR1
<i>S. cerevisiae</i> /PIKK	2473	TOR2	P32600	PI3K-related kinase TOR2
<i>S. cerevisiae</i> /PIKK	3744	TRA1	P38811	Transcription-associated protein 1
<i>D. melanogaster</i> /PIKK	3741	CG2905	Q818U7	DTRA1 (Transcription-associated protein 1). CG2905 and Q818U7 are isoforms
<i>D. melanogaster</i> /PIKK	3218	CG4549	Q9W3V6	Smg-1
<i>D. melanogaster</i> /PIKK	2429	CG6535	Q9VFB1	CG6535-PA
<i>D. melanogaster</i> /PIKK	2470	Tor	Q9VK45	Tor
<i>D. melanogaster</i> /PIKK	2354	mei-41	Q9VXG8	mei-41
<i>C. elegans</i> /PIKK	2697	B0261.2	Q95Q95	Target of rapamycin homolog (CeTOR), gene let-363
<i>C. elegans</i> /PIKK	3944	C47D12.1	Q6A4L2	Hypothetical protein C47D12.1b, gene trr-1
<i>C. elegans</i> /PIKK	2531	atl-1	Q22258	Hypothetical protein T06E4.3a
<i>C. elegans</i> /PIKK	649	atm-1	Q9N3Q4	ATM family protein 1
<i>C. elegans</i> /PIKK	2327	Smg-1	O01510	Smg-1 (suppressor with morphological effect on genitalia protein 1). KinBase smg-1 and O01510 are isoforms
<i>M. musculus</i> /PIKK	2635	ATR	Q9JKK8	No direct correspondence. Q9JKK8 is a fragment of KinBase ATR
<i>M. musculus</i> /PIKK	4128	DNAPK	P97313	DNA-dependent protein kinase
<i>M. musculus</i> /PIKK	2549	FRAP	Q9JLN9	FKBP12-rapamycin complex-associated protein
<i>M. musculus</i> /PIKK	3658	SMG1	Q8BLU4	smg-1
<i>M. musculus</i> /PIKK	3877	TRRAP	Q80YV3	Transformation/transcription-associated protein. TRRAP and Q80YV3 are isoforms
<i>H. sapiens</i> /PIKK	3056	ATM	Q13315	ATM ser/thr protein kinase
<i>H. sapiens</i> /PIKK	2644	ATR	Q13535	ATR
<i>H. sapiens</i> /PIKK	4128	DNAPK	P78527	DNA-dependent protein kinase
<i>H. sapiens</i> /PIKK	2549	FRAP	P42345	FKBP12-rapamycin complex-associated protein
<i>H. sapiens</i> /PIKK	3657	SMG-1	Q96Q15	PI3K-related protein kinase smg-1
<i>H. sapiens</i> /PIKK	3877	TRRAP	Q9Y4A5	Transformation/Transcription domain-associated protein
<i>S. cerevisiae</i> /RIO	484	6324693_RIO1	Q12196	Ser/Thr protein kinase RIO1
<i>S. cerevisiae</i> /RIO	425	6324122_RIO2	P40160	Ser/Thr protein kinase RIO2
<i>D. melanogaster</i> /RIO	585	CG11660	Q9VTL5	CG11660-PA, isoform A
<i>D. melanogaster</i> /RIO	538	CG11859	Q9UBU2	CG11859-PA
<i>D. melanogaster</i> /RIO	603	CG3008	Q9VR42	CG3008-PA
<i>C. elegans</i> /RIO	548	M01B12.5	O44959	Hypothetical protein M01B12.5
<i>C. elegans</i> /RIO	510	ZK632.3	P34649	Putative RIO-type ser/thr protein kinase
<i>C. elegans</i> /RIO	529	aSWK440	Q95Q34	Hypothetical protein Y105E8B.3
<i>M. musculus</i> /RIO	567	RIOK1	Q9CU84	RIOK1 ser/thr protein kinase
<i>M. musculus</i> /RIO	547	RIOK2	Q9CQS5	RIOK2 ser/thr protein kinase
<i>M. musculus</i> /RIO	519	RIOK3	Q9DBU3	RIOK3 ser/thr protein kinase
<i>H. sapiens</i> /RIO	568	RIOK1	Q9BRS2	RIO1 ser/thr protein kinase
<i>H. sapiens</i> /RIO	552	RIOK2	Q9BVS4	RIO2 ser/thr protein kinase
<i>H. sapiens</i> /RIO	519	RIOK3	Q14730	RIO3 ser/thr protein kinase
<i>S. cerevisiae</i> /PDHK	445	6321379_YGL059W	P53170	Hypothetical 51.9 kDa protein in PYC-UBC2 intergenic region
<i>S. cerevisiae</i> /PDHK	394	6322147_YIL042C	P40530	Hypothetical 45.4 kDa protein in CBR5-NOT3 intergenic region
<i>D. melanogaster</i> /PDHK	413	PDK	P91622	Pyruvate dehydrogenase kinase
<i>C. elegans</i> /PDHK	401	ZK370.5	Q02332	Probable PDHK, mitochondrial precursor
<i>M. musculus</i> /PDHK	412	BCKDK	Q55028	BCKDK, mitochondrial precursor
<i>M. musculus</i> /PDHK	434	PDHK1	Q8BFP9	PDHK, isozyme 1, mitochondrial precursor
<i>M. musculus</i> /PDHK	407	PDHK2	Q9JK42	PDHK2, mitochondrial precursor
<i>M. musculus</i> /PDHK	415	PDHK3	Q922H2	PDHK3, mitochondrial precursor
<i>M. musculus</i> /PDHK	412	PDHK4	O70571	PDHK4, mitochondrial precursor
<i>H. sapiens</i> /PDHK	412	BCKDK	O14874	BCKDK, mitochondrial precursor
<i>H. sapiens</i> /PDHK	436	PDHK1	Q15118	PDHK, isozyme 1, mitochondrial precursor
<i>H. sapiens</i> /PDHK	407	PDHK2	Q15119	PDHK, isozyme 2, mitochondrial precursor
<i>H. sapiens</i> /PDHK	406	PDHK3	Q15120	PDHK, isozyme 3, mitochondrial precursor

*Continued*



**Table II**  
Continued

Organism/family	Protein length	KinBase accession number	Swiss-Prot accession number	Notes
<i>H. sapiens</i> /PDHK	411	PDHK4	Q16654	PDHK, isozyme 4, mitochondrial precursor
<i>C. elegans</i> /Alpha	760	efk-1	O01991	Elongation factor-2 kinase
<i>M. musculus</i> /Alpha	1862	ChaK1	Q923J1	TRMP7
<i>M. musculus</i> /Alpha	2028	ChaK2	Q8CIR4	TRMP6
<i>M. musculus</i> /Alpha	724	eEF2K	O08796	Elongation factor-2 kinase
<i>H. sapiens</i> /Alpha	1907	AlphaK1	Q96L96	Muscle Alpha-kinase
<i>H. sapiens</i> /Alpha	1531	AlphaK2	Q96L95	Heart Alpha-kinase
<i>H. sapiens</i> /Alpha	1244	AlphaK3	Q96QP1	Lymphocyte Alpha-kinase
<i>H. sapiens</i> /Alpha	1865	ChaK1	Q96QT4	TRPM7 human
<i>H. sapiens</i> /Alpha	2012	ChaK2	Q9BX84-2	TRPM6
<i>H. sapiens</i> /Alpha	725	eEF2K	O00418	Elongation factor-2 kinase

Four families—PIKK, RIO, PDHK, and Alpha—are present in *H. sapiens* and *C. elegans*, whereas *S. cerevisiae* and *D. melanogaster* only possess members of the RIO, PIKK, and PDHK families.

from KinBase. All kinases gave *E*-values better than  $1e-05$ . The classification was regarded as correct whenever the model that aligned with the best *E*-value to a given query sequence belonged to the same family as the sequence being classified. Table III summarizes the classification performance of Library 1 built to represent the families of protein kinase catalytic domains from *H. sapiens*. Library 1 correctly classified 394/395 of human ePK domains. The misclassified human kinase turned out to be RSK2, an AGC kinase. RSK kinases harbor two functional kinase catalytic domains, and whereas the first domain was correctly classified as “AGC”, the second domain would rather be classified as “CAMK”. This appears to be a typical case of neo-functionalization following

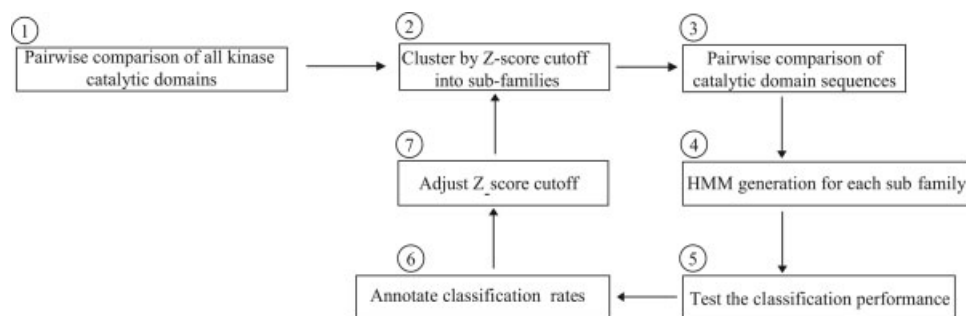
domain duplication. Since the classification of RSK kinases is based on the identity and behavior of the first kinase catalytic domain, this misclassified kinase may be ignored.

Similarly good performance was obtained for the classification of *S. cerevisiae* (75/77) and *D. melanogaster* (175/177) ePKs. However, the HMM library performed less well on the ePKs of *C. elegans*, where only 338/359 ePKs (94.15%) were classified correctly. The misclassification rate was particularly marked in the TKL family of kinases from *C. elegans* where only 6/15 (40%) were classified correctly. In contrast, the HMM library derived from human aPKs correctly classified all atypical kinases of *C. elegans*, *H. sapiens*, *D. melanogaster*, and *S. cerevi-*

**Table III**  
Classification Performance of Library 1 from *H. sapiens*

<i>H. sapiens</i> -based models	<i>C. elegans</i>			<i>H. sapiens</i>			<i>D. melanogaster</i>			<i>S. cerevisiae</i>		
Family	No. of K	No. of CC	% of CC	No. of K	No. of CC	% of CC	No. of K	No. of CC	% of CC	No. of K	No. of CC	% of CC
AGC	30	27	90.00	63	62	98.41	30	30	100	17	16	94.12
CAMK	40	36	90.00	74	74	100	32	31	96.87	21	20	95.24
CK1	85	83	97.65	12	12	100	10	10	100	4	4	100
CMGC	49	49	100	61	61	100	33	33	100	21	21	100
RGC	27	27	100	5	5	100	6	6	100	N/A	N/A	N/A
STE	25	22	88.00	47	47	100	18	18	100	14	14	100
TK	88	88	100	90	90	100	31	31	100	N/A	N/A	N/A
TKL	15	6	40.00	43	43	100	17	16	94.12	N/A	N/A	N/A
Total (ePKs)	359	338	94.15	395	394	99.77	177	175	98.87	77	75	97.40
Alpha	1	1	100	6	6	100	N/A	N/A	N/A	N/A	N/A	N/A
PIKK	5	5	100	6	6	100	5	5	100	5	5	100
RIO	1	1	100	5	5	100	1	1	100	2	2	100
PDHK	3	3	100	3	3	100	3	3	100	2	2	100
Total (aPKs)	10	10	100	20	20	100	9	9	100	9	9	100

The classification was regarded as correct whenever the model aligning with the best *E*-value to a given query sequence belonged to the same family as the sequence being classified. No. of K, number of protein kinases; No. of CC, number of correctly classified protein kinases; % of CC, percentage of correctly classified protein kinases.

**Figure 1**

Schematic representation of the process for arriving at a library of human-derived HMMs that would optimize the family-specific classification of protein kinases across a large evolutionary distance. In step 1, all human kinase catalytic domain sequences are compared pairwise, and the chosen Z-score cutoff splits the family into a number of subfamilies. The sequences in each subfamily are compared pairwise, followed by their multiple alignment and their translation into an HMM. At this stage a new HMM library has been created where each human kinase family is represented by a number of HMMs. The classification performance of the HMM library is tested on the four characterised kinomes of *H. sapiens* (self-classification), *D. melanogaster*, *C. elegans*, and *S. cerevisiae*. The Z-score cutoff is adjusted empirically in this iterative process until the best classification rate has been achieved on the four kinomes.

*siae*, reflecting the high degree of sequence conservation of aPKs across a large evolutionary distance.

The classification performance of the HMM libraries built from *D. melanogaster* and *C. elegans* kinase catalytic domains were assessed in a similar way. The *D. melanogaster*-derived library had a performance similar to the *H. sapiens*-derived library, yielding a correct self-classification rate of 176/177 (99.44%), and cross-species correct classification rates of 77/77 (100%), 389/395 (98.48%), and 347/359 (96.66%) for *S. cerevisiae*, *H. sapiens*, and *D. melanogaster*. Self-classification by the *C. elegans*-derived library yielded a correct classification rate of only 161/359 (44.84%), and similarly poor rates on *H. sapiens* (185/395, 53.17%), *D. melanogaster* (99/177, 55.93%), and *S. cerevisiae* (50/77, 64.93%) (data not shown).

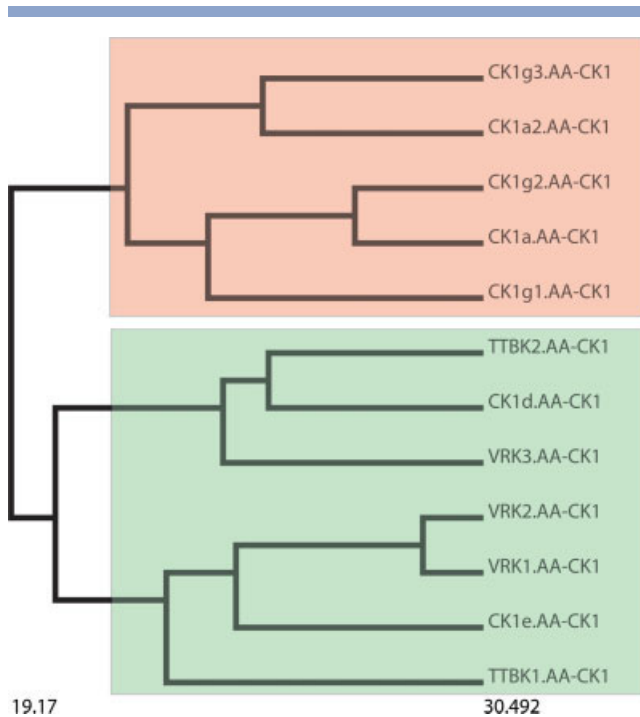
### Library 2: Multiple HMMs per family, and assessment of performance

With the exception of *C. elegans*, the classification performed by Library 1 HMMs appeared good. In an attempt to raise the accuracy of classification still further and in particular for *C. elegans*, the larger sequence families were subdivided to generate multiple HMMs for each family to create Library 2. Subdividing large families can raise recognition accuracy for HMMs by allowing the unique features of each subfamily to be captured more effectively. It can also give benefits by eliminating the need to align divergent sequences and so raise the accuracy of the alignments from which the HMMs are derived.<sup>23</sup>

Given the wide spectrum of protein kinases in *H. sapiens* and the high classification rates of the human kinase-derived library (Library 1, Table III), it was decided to base all further work on the available human kinases.

Attempts were only made to improve the library of HMMs representing human ePKs since the HMM library of aPKs of human origin already gave a misclassification rate of zero on the aPK complement of all four organisms.

Figure 1 outlines the iterative process devised to optimize the representation of each kinase family by multiple subfamily HMMs. Z-scores were calculated for each pairwise alignment of the ePK catalytic domain sequences in each family. Complete linkage clustering was then applied to the Z-scores and clusters selected according to a Z-score cutoff. This is illustrated for the CK1 family in Figure 2. An initial cutoff of 19.2 clustered the 12 kinase catalytic domain sequences of human CK1 protein kinases into a single group. A cutoff of 20 S.D. split the family into two groups of five and seven domains (Fig. 2). For each subfamily a multiple alignment and HMM was generated. As for Library 1, the classification performance of the resulting library of HMMs was tested on the full-length kinases of *H. sapiens*, *C. elegans*, *S. cerevisiae*, and *D. melanogaster*. It was found empirically that generating models for any family where each model consisted of only a handful of sequences resulted in a serious decrease in the rate of correct classification. The trade-off between specificity and number of models for each family led to the number of models summarized in Table IV. Seven models were needed to represent the largest family of kinases from *H. sapiens* (the CAMK family), while the small RGC family and all the atypical kinases continued to be represented by single models from Library 1. Figure 3 provides an overview of the development of Library 2 which has multiple HMMs for larger families from Library 1. The resulting library of subfamily HMMs was referred to as "Library 2".

**Figure 2**

The 12 human kinase catalytic domains of the CK1 family were compared pairwise, and a Z-score cutoff of 20.0 led to their division into two subfamilies. The sequences in each subfamily were then compared pairwise, followed by their multiple alignment and the generation of an HMM from each alignment.

Table V shows the classification rate for Library 2 on the ePKs from *H. sapiens*, *D. melanogaster*, *C. elegans*, and *S. cerevisiae*. With the exception of *C. elegans*, the library gave perfect classification, while for *C. elegans*, the accuracy rose by 14 kinases from the 338/359 achieved by Library 1 to 354/359. The largest improvement for *C. elegans* was seen for the TKL family, where the number of misclassified kinases dropped from 9 to 0. No change in classification accuracy was seen for the CK1 family (83/85) while for CMGC, a single error was introduced with 48/49 correctly classified kinase domains compared to 49/49 for Library 1. The misclassified CMGC kinase was C34G6.5 of *C. elegans*. The Library 2 classification suggests that this might be a kinase of the AGC family since the HMM classifying it with the lowest *E*-value was AGC\_sub4.hmm (*E*-value = 5.6e-15), followed by CMGC\_sub2.hmm with an *E*-value of 1.6e-12. Likewise, the AGC kinase that was not classified as an AGC (YL3D4A.6) was, according to Library 2 a CAMK (*E*-value of CAMK\_sub4.hmm 8.7e-40 versus *E*-value of AGC\_sub4.hmm 1.5e-29, the top-matching AGC model). This suggests that CAMK R166.5 might indeed be an AGC (*E*-value of AGC\_sub4.hmm 1.3e-57) and CK1s F26A1.3 and F39F10.2 might be CAMKs (*E*-values of CAMK\_sub7.hmm 3.2e-5 and 4.1e-11, respectively). The

division of each family of human kinase catalytic domains into a number of subfamilies, followed by their representation in the form of HMMs has been shown to increase the classification rate for KinBase sequences and suggest possible incorrect annotations for C34G6.5, YL3D4A.6, R166.5, F26A1.3, and F39F10.2 of *C. elegans*.

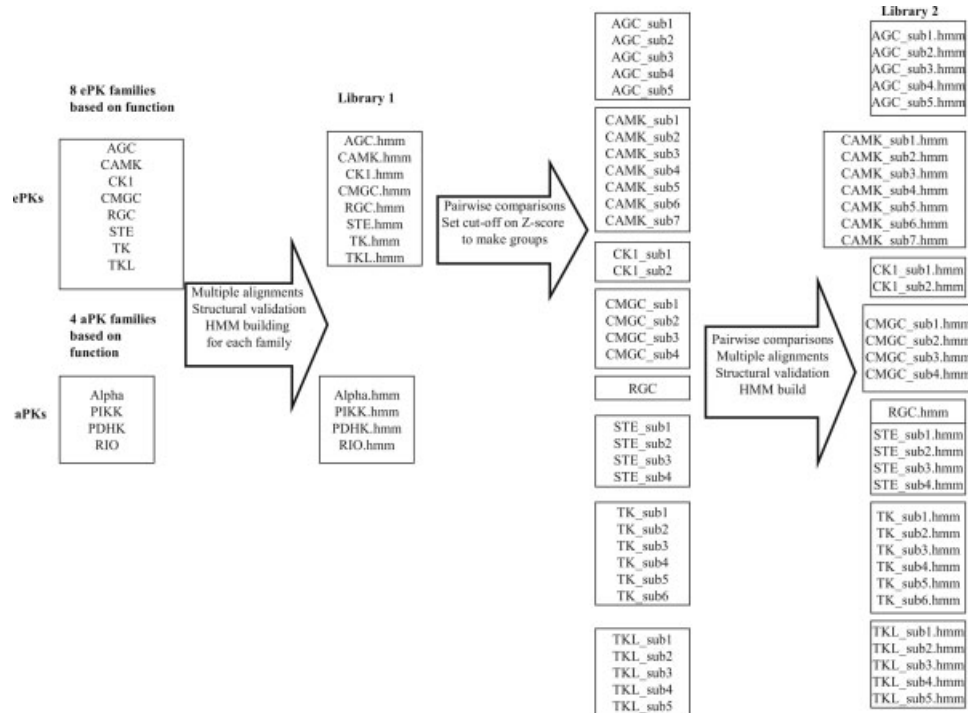
Table V shows best case results for Library 2 since they are for organisms against which the subfamilies that make up the Library were optimized. To obtain a more realistic estimate of the performance of Library 2, it was applied to the kinomes of *M. musculus*<sup>14</sup> and the phylogenetically distant human malaria parasite *P. falciparum*<sup>15</sup> and *D. discoideum*.<sup>16</sup> These organisms have had their kinomes annotated and so provide a standard against which to test the approach described here. The results for Library 2 are summarized in Table VI. The *M. musculus* kinome is very similar to that of *H. sapiens* and so the 100% correct classification shown is no surprise. Library 2 was able to classify correctly all protein kinases of *P. falciparum* that had been assigned to families by Ward et al.<sup>15</sup> In the case of the *Dictyostelium* kinome, Library 2 was able to classify correctly all the kinases of the AGC, CK1, CMGC, and TKL families. The protein kinase DDB0229351 had been classified by Goldberg et al. as a CAMK kinase.<sup>16</sup> However, Library 2 suggests that this sequence would rather belong to the AGC family (*E*-value for AGC\_sub4.hmm of 2.6e-72) than to the CAMK family (*E*-value for the best hit to a CAMK sub-HMM was 1.6e-67). When the phylogenetic tree provided by Goldberg et al. was inspected in detail, DDB0229351 was found to be classified as a CAMK because it groups loosely with a cluster of bonafide CAMK kinases (the “FHAK”).<sup>16</sup> However, Goldberg et al. could not provide

**Table IV**

Make-Up of HMM Library 2

Family	Number of catalytic domain sequences	Z-score cutoff	Number of groups generated for each family
AGC	78	19.5	5
CAMK	79	19.5	7
CK1	12	20	2
CMGC	61	19	4
RGC	5	N/A	1
STE	47	19	4
TK	102	19	6
TKL	43	19	5
<i>Total ePKs</i>	<i>407</i>		<i>34</i>
Alpha	6	N/A	1
PIKK	5	N/A	1
PDHK	5	N/A	1
RIO	3	N/A	1
<i>Total aPKs</i>	<i>20</i>		<i>4</i>

Each family used in the construction of Library 1 was split into subfamilies as shown in this table. For example, the AGC family was represented by 5 subfamily HMMs rather than a single HMM for the entire family.

**Figure 3**

Graphic summary of the process for arriving at Library 2 from Library 1.

sound bootstrap support for this association. Eight protein kinases that were classified as STE kinases<sup>16</sup> were not classified as such by Library 2. Examination of the phylogenetic tree provided by Goldberg et al.<sup>16</sup> showed that all these protein kinases were found to belong to a self-contained group that has some loose relationship to the main group of STE kinases and to which bootstrap support was not provided. The authors of the *Dictyostelium* kinome<sup>16</sup> produced a phylogenetic tree from a multiple alignment of

the entire kinome, including a number of representative sequences from other kinomes. A phylogenetic tree is only as good as the underlying alignment. Given the degree of sequence divergence of the protein kinase superfamily members, the generation of whole kinome multiple alignments is likely to produce trees have little or no bootstrap support for some of the branching patterns observed. Therefore, the generation of protein kinase family-specific multiple alignments is bound to produce more meaningful

**Table V**

Classification Performance of Library 2

Family	<i>H. sapiens</i>			<i>D. melanogaster</i>			<i>C. elegans</i>			<i>S. cerevisiae</i>		
	No. of K	No. of CC	% of CC	No. of K	No. of CC	% of CC	No. of K	No. of CC	% of CC	No. of K	No. of CC	% of CC
AGC	63	63	100	30	30	100	30	29	96.67	17	17	100
CAMK	74	74	100	32	32	100	40	39	97.50	21	21	100
CK1	12	12	100	10	10	100	85	83	97.65	4	4	100
CMGC	61	61	100	33	33	100	49	48	97.96	21	21	100
RGC	5	5	100	6	6	100	27	27	100	N/A	N/A	N/A
STE	47	47	100	18	18	100	25	25	100	14	14	100
TK	90	90	100	31	31	100	88	88	100	N/A	N/A	N/A
TKL	43	43	100	17	17	100	15	15	100	N/A	N/A	N/A
Total	395	395	100	177	177	100	359	354	98.61	63	63	100

The classification was regarded as correct whenever the model aligning with the best E-value to a given query sequence belonged to the same family as the sequence being classified. No. of K, number of protein kinases; No. of CC, number of correctly classified protein kinases; % of CC, percentage of correctly classified protein kinases.



**Table VI**Performance of Library 2 in the Classification of the Kinomes of *M. musculus*, the malaria parasite *P. falciparum*, and *D. discoideum*

Family	<i>M. musculus</i>			<i>P. falciparum</i>			<i>D. discoideum</i>		
	No. of K	No. of CC	% of CC	No. of K	No. of CC	% of CC	No. of K	No. of CC	% of CC
AGC	60	60	100	5	5	100	22	22	100
CAMK	97	97	100	13	13	100	21	20	95.24
CK1	11	11	100	1	1	100	3	3	100
CMGC	60	60	100	18	18	100	30	30	100
RGC	7	7	100	N/A	N/A	N/A	N/A	N/A	N/A
STE	47	47	100	N/A	N/A	N/A	45	37	82.22
TK	91	91	100	N/A	N/A	N/A	N/A	N/A	N/A
TKL	44	44	100	5	5	100	68	68	100
Total	417	417	100	42	42	100	189	180	95.24

No. of K, number of protein kinases; No. of CC, number of correctly classified protein kinases; % of CC, percentage of correctly classified protein kinases.

alignments from which more robust phylogenetic trees can be derived.

In conclusion, Library 2 was able to classify correctly all *bona fide* protein kinases of *M. musculus*, *P. falciparum*, and *D. discoideum*. *Dictyostelium* and the malaria parasite represent early branching points in eukaryote evolution,<sup>28–30</sup> and the fact that Library 2 can classify their kinomes correctly suggests that it should be generally useful for the classification of eukaryotic protein kinases.

## RESULTS AND DISCUSSION

A multilevel HMM library (Library 2) has been developed that is able to classify eukaryotic protein kinases into one of the 12 previously established families. The ePK families include the AGC, CAMK, CK1, CMGC, RGC, STE, TK, and TKL, and the *bona fide* aPK families Alpha, PIKK, PDHK, and RIO. As explained in Methods, the reliability of Library 2 has been tested by cross-validation between the kinomes of *H. sapiens*, *D. melanogaster*, *C. elegans*, *S. cerevisiae*, *M. musculus*, *D. discoideum*, and *P. falciparum*.

### Comparison of HMM library 2 to BLAST and a general HMM for kinase classification

HMMs are known to be more sensitive and selective than pairwise comparison methods such as BLAST.<sup>21</sup> However, since BLAST has been widely applied in kinome analysis, the performance of Library 2 and BLASTP were compared in database searches. For this, a database was assembled that consisted of Uniref100<sup>10</sup> plus the well-characterized protein kinase complements of *H. sapiens*, *D. melanogaster*, *C. elegans*, and *S. cerevisiae* from KinBase. The database was scanned with Library 2 and, for comparison, the domains making up each HMM were also used to search the database by BLASTP (version 2.2.10) with default parameters (matrix = BLOSUM62; cost to extend/open a gap = 0; word size = 3; but with the “Filter query sequence” parameter set to F).

For Library 2, the number of true positives (TP) was determined as the number of KinBase sequences returned with scores better than 1e-05. Library 2 was able to retrieve and correctly classify all ePKs of *H. sapiens*, *C. elegans*, *D. melanogaster*, and *S. cerevisiae*, respectively (Table VII). The Library 2 database search results illus-

**Table VII**

Result of Using the Library 2 to Search Uniref100 for Annotated Protein Kinases

Family	<i>H. sapiens</i>			<i>C. elegans</i>			<i>D. melanogaster</i>			<i>S. cerevisiae</i>		
	No. of K	TP	E-value	No. of K	TP	E-value	No. of K	TP	E-value	No. of K	TP	E-value
AGC	63	63	8.50e-014	30	30	2.7e-07	30	30	5.4e-09	17	17	1.6e-58
CAMK	74	74	2.2e-16	40	40	5.6e-24	32	32	9e-15	21	21	3.2e-14
CK1	12	12	5.8e-165	85	85	3.2e-05	10	10	1.2e-08	4	4	4.3e-126
CMGC	61	61	5.6e-108	49	49	6.7e-12	33	33	7.1e-39	21	21	1.2e-07
RGC	5	5	7.5e-168	27	27	4.8e-05	6	6	4.9e-59	0	0	N/A
STE	47	47	1.6e-98	25	25	1.4e-06	18	18	2.6e-13	14	14	8.2e-56
TK	90	90	7.3e-89	88	88	1.1e-09	31	31	7.7e-13	0	0	N/A
TKL	43	43	1.7e-60	15	15	1.7e-12	17	17	1.4e-14	0	0	N/A
Total	395	395/395 (100%)		359	359/359 (100%)		177	177/177 (100%)		77	77/77 (100%)	

Library 2 retrieves all of the KinBase kinases in a database search. No. of K: Number of KinBase sequences; TP indicates the number of KinBase sequences of that family retrieved by the library of models; E-value: Indicates the worst E-value obtained for a protein kinase member of the family.

**Table VIII**

Result of Using BLAST with Human Kinase Catalytic Domains as Query Sequences to Mine Uniref 100 KinBase Sequences

Kinase family	<i>H. sapiens</i>			<i>C. elegans</i>			<i>D. melanogaster</i>			<i>S. cerevisiae</i>		
	No. of K	TP/FP	E-value	No. of K	TP/FP	E-value	No. of K	TP/FP	E-value	No. of K	TP/FP	E-value
AGC	63	20.44/3.88	6e-15	30	6.47/0.68	5e-15	30	7.13/0.96	4e-15	17	4.72/0.26	2e-15
CAMK	74	20.03/1.32	2e-1	40	6.26/0.16	5e-16	32	6.70/0.16	1e-15	21	2.36/0.03	1e-17
CK1	12	11.08/0.00	2e-04	85	48.00/0.08	2e-04	10	8.83/0.00	7e-08	4	4.00/0.00	5e-08
CMGC	61	18.02/0.00	4e-07	49	7.89/0.00	3e-07	33	8.15/0.00	1e-08	21	3.84/0.00	1e-07
RGC	5	5.00/11.40	2e-20	27	21.80/2.20	1e-11	6	6.00/1.60	1e-23	0	0.00/0.00	N/A
STE	47	22.38/0.36	2e-16	25	6.40/0.04	5e-18	18	7.55/0.09	3e-16	14	3.83/0.04	5.00E-015
TK	90	19.34/2.19	4e-15	88	2.99/0.68	3e-15	31	5.15/0.69	2e-15	0	0.00/0.07	N/A
TKL	43	10.42/6.98	2e-09	15	2.98/2.16	2e-09	17	3.86/2.21	2e-09	0	0.00/0.23	N/A
Total	395			359			177			77		

BLAST was not found to be as good as Library 2, only returning a fraction of the KinBase sequences. No. of K, number of KinBase sequences; TP/FP, average number of true positives (TP), and false positives (FP). A FP is defined as a KinBase sequence that has been retrieved in the search but which belongs to a distinct ePK family as the query sequence; E-value: E-value of lowest-scoring match.

trate the *E*-value cutoffs that should be applied to each family of protein kinases in database searches. Over 99% of the proteins from Uniref100 retrieved alongside the KinBase sequences with *E*-values above the cutoff for each family had been annotated as protein kinases.

As expected, BLASTP searching with the kinase catalytic domains was found to be less sensitive and less specific than the search performed with Library 2. The detailed results of using KinBase human kinase catalytic domains as query sequences to search for KinBase sequences of the four organisms are presented in Table VIII. The typical BLAST behavior is illustrated by the search performed with the human AGC catalytic domains. Searching with AGC kinase catalytic domains returned an average of 20.44/63 (32.4%), AGC kinases of *H. sapiens* above the default BLAST cutoff *E*-value of 10.0 including kinases belonging to families other than AGC. The average performance of human AGC kinase domains as query sequences was similarly poor with respect to the AGC kinases of *C. elegans*, *D. melanogaster*, and *S. cerevisiae*, retrieving on average 6.47/30 (21.57%), 7.13/30 (23.77%), and 4.72/17 (27.76%), respectively, although the average number of false positives was lower in these organisms as smaller families tend to harbor less divergent sequences. The *E*-values returned by BLAST are higher than those returned by Library 2 for the majority of families, and since BLAST only retrieves a fraction of all existing protein kinases, multiple overlapping BLAST searches would be necessary to arrive at the same result that is provided by Library 2 in a single step. Furthermore, Library 2 produces automatically, a classification of protein kinases that is correct at the family level, whereas the hits returned by BLASTP searches retrieve a mixture of kinases from different families.

In kinome analysis papers, researchers have used a combination of BLAST and general HMMs of the kinase catalytic domain to mine databases for protein kinases. Here, we have compared the database scanning perform-

ance of a general, multispecies, HMM of the kinase catalytic domain from Pfam (PF00069, Pfam\_fs model) with Library 2. For a full comparison of the database retrieval capacities of the two approaches, we have also included the protein kinases belonging to the “Other” family of *S. cerevisiae*, *C. elegans*, *D. melanogaster*, and *H. sapiens*. Table IX shows the results of retrieving the kinomes of *H. sapiens*, *C. elegans*, *D. melanogaster*, and *S. cerevisiae* with the Pfam HMM PF00069. Whereas Library 2 could retrieve all protein kinases, PF00069 failed to identify a number of protein kinases, especially those belonging to the “Other” family. This is not surprising as protein kinases belonging to this family are typically the most divergent ones, making them more difficult to retrieve and classify. In these cases, researchers typically use BLAST for identifying select, remote, homologues.<sup>16</sup> Overall PF00069 failed to identify 2.93, 10.09, 3.60, and 2.61% of the kinases of human, worm, fly, and yeast origin. This indicates that Library 2 is superior not only to BLAST but also to a general HMM of the kinase catalytic domain for database searches. In addition, Library 2 is capable of doing automatic classification of protein kinases into families.

#### Classification of the “Other” kinases in *S. cerevisiae*

In any characterized kinome, a group called “Other” always exists which consists of protein kinase sequences known to belong to the ePKs but which annotators have not been able to classify because of lack of clear similarity to kinases in the eight ePK families. The current *H. sapiens* kinome annotation<sup>6</sup> contains 83 kinases belonging to the “Other” group, whereas 45 are found in the *D. melanogaster* kinome<sup>4</sup>, and 38 and 67 in the *S. cerevisiae*<sup>12</sup> and *C. elegans*<sup>11</sup> kinomes, respectively.

In an attempt to classify the “Other” ePKs of yeast, these were scanned through Library 2, and also their syn-

**Table IX**

Result of Using the General HMM of the Kinase Catalytic Domain PF00069 for Searching Uniref 100 KinBase Sequences

Kinase family	<i>H. sapiens</i>			<i>C. elegans</i>			<i>D. melanogaster</i>			<i>S. cerevisiae</i>		
	No. of K	No. missed	% missed	No. of K	No. missed	% missed	No. of K	No. missed	% missed	No. of K	No. missed	% missed
AGC	63	0	0	30	0	0	30	0	0	17	0	0
CAMK	74	0	0	40	1	2.50	32	0	0	21	0	0
CK1	12	1	8.33	85	13	15.29	10	1	10	4	0	0
CMGC	61	0	0	49	0	0	33	0	0	21	0	0
RGC	5	0	0	27	4	14.81	6	0	0	0	0	0
STE	47	0	0	25	0	0	18	0	0	14	0	0
TK	90	0	0	88	1	1.14	31	0	0	0	0	0
TKL	43	0	0	15	0	0	17	0	0	0	0	0
Other	83	13	15.66	67	24	35.82	45	7	15.55	38	3	7.89
Total	478	14	2.93	426	43	10.09	222	8	3.60	115	3	2.61

The HMM failed to identify between 2.61 and 10.09% of the ePKs from characterized kinomes, especially the more divergent “Other” ePKs. No. of K, number of KinBase sequences.

tenic homologues from the related fungus *Ashbya gossypii*. Although *S. cerevisiae* and *A. gossypii* shared a common ancestor over 100 million years ago, the order and orientation of >95% of *Ashbya*'s genes is conserved in the genome of *S. cerevisiae*.<sup>31</sup> Therefore, the automatic classification of syntenic homologues from two different species into the same family is an acceptable measure of family membership. Table X shows the classification of the “Other” kinases of yeast into the main ePK families, together with the classification of the homologous genes from *A. gossypii*. 27/38 yeast kinases of the “Other” group, together with their corresponding, syntenic, homologues in *A. gossypii*, were classified by Library 2 automatically into the same family, and above the family-specific *E*-value cut-off. These 27 kinases of yeast correspond to 23 loci in the ancestral fungal genome. This suggests that, following the whole genome duplication event in the *Saccharomyces* lineage, most duplicate kinase genes, like most duplicate *S. cerevisiae* genes, were lost secondarily.<sup>32</sup> However, those pairs of yeast kinases that originated from the same ancestral locus (the pairs are KKQ8 and HAL5 (*Ashbya* AEL118C); PAK1 and TOS3 (*Ashbya* ACL053C); NPR1 and PRR2 (*Ashbya* ABL143C); and PTK1 and PTK2 (*Ashbya* AFR372W), have retained the same family classification for both of the yeast kinases.

The yeast kinase of the “Other” group VPS15 corresponds with locus ADL316C in *A. gossypii*, and which in *Ashbya* codes for a degenerate, nonfunctional, form of a kinase catalytic domain. The remaining 10 yeast “Other” kinases could not be classified with strict confidence into the same ePK family together with their syntenic homologues from *A. gossypii*. Even though these two fungal species shared a common ancestor 100 million years ago, the two have radically different life styles. *A. gossypii* is a fungal pathogen of cotton plants, whereas *S. cerevisiae* has been used as a model organism for decades in the laboratory. The real power of comparative kinomics will

come from comparing the kinomes of all yeast species that have been, and are being sequenced. It might be possible that a number of yeast “Other” kinases constitute the seeds of fungi-specific kinase subfamilies that represent independent clusters themselves.

The application of Library 2 enriched the repertoire of yeast ePKs as follows: AGC(+5), CAMK(+17), CMGC(+4), STE(+1). This has raised the classification rate of yeast ePKs from 77/115 (66.96%) to 104/115 (90.43%). While it is not possible to prove computationally that the suggested new classification for these 27 kinases is correct, the consistency at syntenic loci for the two organisms lends support to the classification. The high degree of selectivity of Library 2 suggests that it might be useful for reassigning the “Other” kinases of a number of kinomes. Representative alignments of yeast kinases of the “Other” group, now reclassified into one of the main ePK families, are provided in Figure S1.

### Classification of the protein kinase complement of 21 eukaryotic genomes

To allow a comparison of kinomes between divergent species, Library 2 was used to search for protein kinases in 21 completed and published eukaryotic genomes. These included two algal species, *Thalassiosira pseudonana*<sup>33</sup> and *Cyanidioschyzon merolae*,<sup>34</sup> the arthropod *Anopheles gambiae*,<sup>35</sup> the chordate *Ciona intestinalis*,<sup>36</sup> the fishes *Tetraodon nigroviridis*,<sup>37</sup> and *Takifugu rubripes*,<sup>38</sup> eight fungal species: *Candida glabrata*, *Kluyveromyces lactis*, *Debaryomyces hansenii*, and *Yarrowia lipolytica*,<sup>39</sup> *Cryptococcus neoformans*,<sup>40</sup> *Neurospora crassa*,<sup>41</sup> *Phanerochaete chrysosporium*,<sup>42</sup> and *Schizosaccharomyces pombe*,<sup>43</sup> the mammal *Rattus norvegicus*,<sup>44</sup> three plants: *Arabidopsis thaliana*,<sup>45</sup> *Oryza sativa* L. ssp. *Indica*,<sup>46</sup> and *Oryza sativa* L. ssp. *Japonica*,<sup>47</sup> and three parasitic protozoans: *Cryptosporidium hominis*,<sup>48</sup> *Entamoeba histolytica*,<sup>49</sup> and *Plasmodium yoelii*<sup>50</sup> (Table XI). These

**Table X**

Classification of the Yeast Protein Kinases of the "Other" Group into the ePK Families AGC, CAMK, CMGC, and STE

Yeast protein kinase of the "Other" group	Locus ID on the genome of <i>S. cerevisiae</i>	Syntenic locus on the genome of <i>A. gossypii</i>	New family (reclassification)
BUB1	YGR188C	AGR315C	AGC
CDC5	YMR001C	ACL006W	AGC
IRE1	YHR079C	ADR293C	AGC
IPL1	YPL209C	AFL101C	AGC
YKL171W	YKL171W	AEL120W	AGC
APG1	YGL180W	ACL054W	CAMK
KKQ8	YKL168C	AEL118C	CAMK
HAL5	YJL165C	AEL118C	CAMK
IKS1	YJL057C	AEL173W	CAMK
ISR1	YPR106W	AEL330C	CAMK
PAK1	YER129W	ACL053C	CAMK
TOS3	YGL179C	ACL053C	CAMK
NPR1	YNL183C	ABL143C	CAMK
PRR2	YDL214C	ABL143C	CAMK
KSP1	YHR082C	ADR300C	CAMK
PTK1	YKL198C	AFR372W	CAMK
PTK2	AFR372W	AFR372W	CAMK
SAT4	YCR008W	AER195C	CAMK
YDL025C	YDL025C	ADR313W	CAMK
YGR052W	YGR052W	AEL284C	CAMK
YOR267C	YOR267C	ADR174C	CAMK
YPL236C	YPL236C	AFL143C	CAMK
CDC7	YDL017W	AER216C	CMGC
CKA1	YIL035C	ADL102C	CMGC
CKA2	YOR061W	ADR204W	CMGC
MPS1	YDL028C	ADR317C	CMGC
SWE1	YJL187C	AEL149C	STE

genomes are representative of a diversity of phylogenetic lineages and display a large variation in their predicted gene numbers, ranging from 3994 for *Cryptosporidium hominis* (an Apicomplexan pathogen causing diarrhoea and acute gastroenteritis in humans) to 46,000–55,000 for *Oryza sativa ssp. indica* (rice).

The set of predicted peptides for each proteome was downloaded from the sources indicated in Table XI. Library 2 was then applied uniformly to the 21 genomes, with an *E*-value cutoff specific for each family (Table VII). The "Other" kinases of the characterized kinomes of *H. sapiens*, *D. melanogaster*, *C. elegans*, *S. cerevisiae*, *M. musculus*, *P. falciparum*, *S. cerevisiae*, and *A. gossypii* were also classified by Library 2 into the main ePK families. This was done in order to compare the number of kinases per family across organisms.

If we accept the classification as presented here (Table XII), the 21-genome kinase-family annotation showed that five ePK families (AGC, CAMK, CK1, CMGC, and STE) and two aPK families (RIO and PIKK) were present in all eukaryotic organisms whose genomes have been sequenced. As a consequence, these seven kinase families are likely to be indispensable to eukaryotic life. The other families (TK, TKL, Alpha, PDHK, and RGC) appear to

be restricted to specific lines of descent, as discussed below.

### TKs

In metazoans, tyrosine phosphorylation plays an essential role in intercellular communication (organ development and tissue homeostasis) as well as intracellular communication (transcriptional control, proliferative/differentiation decisions, cell shape, and cell motility). Abnormal level of tyrosine phosphorylation result in conditions such as cancer and immune diseases.<sup>51–54</sup> TKs were absent from fungi, and were mainly found in metazoan species, in agreement with their role. Interestingly, the three plant species seemed to encode putative TKs: 2 in *A. thaliana*, and 6 and 7 in the two rice species. When looked at in detail, the two putative TKs of *A. thaliana* lack the conserved lysine residue of subdomain II and the conserved aspartate of the activation loop in subdomain VII.<sup>5</sup> This suggests that these two enzymes are likely to be catalytically inactive. The two rice species were found to encode six and seven putative TKs. Out of the six predicted TKs of *O. sativa ssp. Indica*, only three of them were found to include all the residues known to be essential for catalytic activity. Out of the seven putative TKs of *O. sativa ssp. Japonica*, four were found to harbor all residues known to be required for catalytic activity. To our knowledge, TKs have never been predicted in plants before,<sup>55</sup> although the existence of tyrosine phosphorylation in plants has been previously documented. Tyrosine phosphorylation in plants is thought to be involved in stress-related responses, embryogenesis, and tissue differentiation,<sup>56</sup> and in response to movement.<sup>57</sup> The small number of putative TKs identified in the plant species does not necessarily mean that tyrosine phosphorylation in plants is a limited phenomenon, but that tyrosine phosphorylation as carried out by putative TKs is probably more limited than that carried out by dual-specificity kinases in plants. A recent *in silico* survey of the *Arabidopsis* proteome suggested that this plant lacks bona fide TKs, and that tyrosine phosphorylation in plants might be carried out by dual-specificity protein kinases.<sup>55</sup>

The human intestinal parasite *Entamoeba histolytica* encodes seven putative TKs, which had already been noted.<sup>49,58</sup> The authors identified 270 ePKs, including TKLs, and at least 90 putative receptor serine/threonine kinases, which are uncommon in protists and usually only found in plants and animals. Library 2 predicted a total of 295 ePKs for *E. histolytica*. This makes *E. histolytica* the single-celled eukaryote with the second largest kinome after *Trichomonas vaginalis*.<sup>59</sup>

The plant TKs, together with the seven putative TKs of *E. histolytica* identified here, and the absence of TKs among the fungal species surveyed in this study, suggest that a primitive set of TKs might have predated the radi-



**Table XI**

Source and Statistics on the 21 Eukaryotic Genomes Whose Protein Kinase Complement was Determined by Library 2, and Which Represent a Variety of Taxonomic Groups

Taxonomy and species <sup>a</sup>	Common name/description	Source database <sup>b</sup>	Release <sup>c</sup>	Haploid genome size (Mb)	Predicted number of genes
<b>Fungi</b>					
<i>C. glabrata</i>	Pathogen causing human candidiasis	Genolevures	N/A	12.3	5283
<i>C. neoformans</i>	Basidiomycetous yeast	TIGR	N/A	19	6572
<i>D. hansenii</i>	Halotolerant yeast	Genolevures	N/A	12.2	6906
<i>K. lactis</i>	Hemiascomycete fungus	Genolevures	N/A	10.6	5329
<i>N. crassa</i>	Multic filamentous fungus	BROAD-MIT	7	40	10,082
<i>P. chrysosporium</i>	White rot fungus	JGI	1.0	30	11,777
<i>S. pombe</i>	Fission yeast	Sanger	N/A	13.8	4940
<i>Y. lipolytica</i>	Alkane-using yeast	Genolevures	N/A	20.5	6703
<b>Animals</b>					
<i>A. gambiae</i>	Malaria mosquito	TIGR	8	278	14,000
<i>C. intestinalis</i>	Ascidian tadpole	JGI	1.0	160	15,852
<i>T. nigroviridis</i>	Freshwater pufferfish	BROAD-MIT	Data version 10/31/01	381	27,918
<i>T. rubripes</i>	Pufferfish	JGI	3.0	365	40,000
<i>R. norvegicus</i>	“Lab rat”	TIGR	10	2.75 Gb	~25,000
<b>Plants</b>					
<i>A. thaliana</i>	Flowering plant	Arabidopsis.org	Data version 31/01/03	125	25,426
<i>O. sativa indica</i>	Rice	NCBI	2.1	466	46,000–55,000
<i>O. sativa jap.</i>	Rice	NCBI	N/A	420	32,000–50,000
<b>Apicomplexa</b>					
<i>C. hominis</i>	Apicomplexan protozoan pathogen	VCU	N/A	9.2	3994
<i>P. yoelii</i>	Rodent malaria parasite	TIGR	5	23.1	5878
<b>Amoebozoa</b>					
<i>E. histolytica</i>	Intestinal protozoan parasite	TIGR	N/A	18–20	9938
<b>Red algae</b>					
<i>C. merolae</i>	Unicellular red alga	Tokyo Univ.	N/A	16	5331
<b>Diatoms</b>					
<i>T. pseudonana</i>	Unicellular alga (diatom)	JGI	1.0	34.5	11,242

<sup>a</sup>Taxonomy follows Baldauf et al. (2000).<sup>28</sup>

<sup>b</sup>The URLs for the source databases are as follows: JGI, <http://genome.jgi-psf.org/>; Tokyo University, <http://merolae.biol.s.u-tokyo.ac.jp/>; NCBI, <ftp://ftp.ncbi.nih.gov/genbank/genomes/>; BROAD-MIT, <ftp://ftp.broad.mit.edu/pub/annotation/>; Genolevures, <http://cbi.labri.fr/Genolevures/download.php#codingseq>; TIGR, <http://www.tigr.org/>; Sanger, [http://www.sanger.ac.uk/Projects/S\\_pombe/](http://www.sanger.ac.uk/Projects/S_pombe/); VCU, [www.parvum.mic.vcu.edu.](http://www.parvum.mic.vcu.edu/); <http://www.arabidopsis.org>. For full references to the genomes, see the text.

<sup>c</sup>N/A, version number not available.

ation of eukaryotic life forms.<sup>28</sup> The number of TKs per predicted gene complement is smaller in the two plant species and *E. histolytica* than in the arthropods *Drosophila* and *Anopheles*, which harbor 33 and 32 TKs, respectively. The fishes *T. nigroviridis* and *T. rubripes* have about 2.5 times as many putative TKs as does the chordate *Ciona intestinalis*. Whole genome duplication events are known to have occurred in the teleost fish lineage after its divergence from the lineage leading to mammals,<sup>37</sup> hence explaining why the two fish species have more TKs than humans do (despite arguably being less complex phenotypically). Receptor and non-receptor TKs have previously been found in choanoflagellates.<sup>60,61</sup> Choanoflagellates are unicellular protozoa related to metazoans and which have been studied in detail to try to explain the molecular innovations that gave rise to animal multicellularity from an ancestral protozoan. TKs have adopted a prominent role in cellular signaling in metazoans when compared with non-metazoans. The radiation of the tyrosine kinase family suggests that this

family has progressively become adapted to fulfilling sophisticated signaling roles in the more complex animals.

### Phosphotyrosine signalling in evolution

The presence of putative TKs in *Entamoeba* and the plant species prompted us to search for further evidence of tyrosine phosphorylation in these organisms. Strong support for tyrosine phosphorylation mechanisms may be provided by the presence of phosphotyrosine-binding domains, such as SH2 and PTB, and the existence of tyrosine phosphatases. Phosphotyrosines are known to bind target proteins via their SH2 or PTB domains to form multiprotein signaling complexes. SH2 domains have previously been described in plants (two in *A. thaliana*, one in *O. sativa*). Plant SH2 domains were found to be shorter than their animal counterparts, but seemed to fold into a functional phosphotyrosine binding domain.<sup>62</sup> A search with a Pfam HMM specific for the SH2 domain (PF00017)<sup>18</sup> identified the same SH2

**Table XII**

The Protein Kinase Complement of the Eukaryotic Organisms with Completed Genomes Divided into 8 ePK and 4 aPK Families

Taxonomy and species	Gene No.	AGC	CAMK	CK1	CMGC	RGC	STE	TK	TKL	Total ePKs	PIKK	PDHK	RIO	Alpha	Total aPKs
<b>Fungi</b>															
<i>C. glabrata</i>	5283	22	29	4	26	0	14	0	0	95	5	2	1	0	8
<i>C. neoformans</i>	6572	22	16	4	24	0	11	0	2	79	4	3	2	0	9
<i>D. hansenii</i>	6906	18	20	3	26	0	14	0	0	81	4	3	1	0	8
<i>K. lactis</i>	5329	20	25	3	25	0	12	0	0	85	4	3	1	0	8
<i>N. crassa</i>	10082	20	20	2	22	0	14	0	0	78	4	3	1	2	10
<i>P. chrysosporium</i>	11777	31	24	5	27	0	17	0	4	108	5	2	1	0	8
<i>S. pombe</i>	4940	20	29	5	26	0	17	0	0	97	6	1	1	0	8
<i>Y. lipolytica</i>	6703	21	19	2	22	0	13	0	0	77	4	3	1	0	8
<i>S. cerevisiae</i>	5885	22	44	4	29	0	16	0	0	115	5	2	2	0	9
<b>Animals</b>															
<i>C. elegans</i>	~19000	35	65	92	56	28	36	94	20	426	5	1	3	1	10
<i>D. melanogaster</i>	13600	42	45	10	39	6	24	33	23	222	5	1	3	0	9
<i>A. gambiae</i>	14000	42	41	9	33	6	24	32	18	205	6	1	3		
<i>C. intestinalis</i>	15852	47	58	6	36	5	27	49	29	257	6	1	3	2	12
<i>T. nigroviridis</i>	27918	105	109	14	90	14	62	120	56	570	5	5	3	6	19
<i>T. rubripes</i>	40000	125	106	17	110	14	76	135	50	633	3	8	3	6	20
<i>R. norvegicus</i>	~25000	96	135	26	84	6	73	106	52	578	5	6	4	6	21
<i>M. musculus</i>	~30000	81	125	11	67	7	59	95	54	499	6	5	3	6	20
<i>H. sapiens</i>	~30000	84	98	12	70	5	61	93	55	478	6	5	3	6	20
<b>Plants</b>															
<i>A. thaliana</i>	25426	67	158	28	149	0	104	2	776	1284	10	3	4	0	17
<i>O. sativa indica</i>	46000–55000	54	93	15	97	0	55	6	1010	1330	5	2	1	0	8
<i>O. sativa jap.</i>	32000–50000	52	88	15	91	0	54	7	969	1276	5	2	1	0	8
<b>Apicomplexa</b>															
<i>C. hominis</i>	3994	11	14	2	18	0	6	0	4	55	2	0	1	0	3
<i>P. yoelii</i>	5878	14	11	1	18	0	0	0	4	48	3	0	2	0	5
<i>P. falciparum</i>	5268	14	18	1	23	0	3	0	6	65	3	0	2	0	5
<b>Amoebozoa</b>															
<i>E. histolytica</i>	9938	38	49	9	48	0	35	7	109	295	6	0	3	4	13
<b>Red algae</b>															
<i>C. merolae</i>	5331	10	9	2	16	0	7	0	9	53	3	1	1	0	5
<b>Diatoms</b>															
<i>T. pseudonana</i>	11242	33	41	3	21	0	8	0	6	112	4	3	2	2	11

The peptide complement of each genome was scanned through Library 2. The assignment of kinase family was performed by retrieving the HMM that classified a peptide with the best *E*-value.

domains as described before for the plant species. Interestingly, *Entamoeba* was found to harbor four polypeptide sequences with SH2 domains, three of which were found to be protein kinases. This comes in contrast with the plant SH2-containing polypeptide chains, where the regions out with their SH2 domains seem to have no inferable function.<sup>62</sup>

PTB domains are also known to bind phosphotyrosines.<sup>63</sup> However, a search with a Pfam HMM specific for the PTB domain (PF08416) returned no plausible candidates.

A third type of phosphotyrosine-binding domain has recently been characterized, the C2 domain of human protein kinase C delta.<sup>64</sup> Crystallographic analysis showed that the phosphotyrosine binding ability of the C2 domain is dependent on a number of key residues. A search with an HMM specific for the C2 domain (PF00168) returned 25 proteins containing C2 domains in *Entamoeba* (2 of which have double C2 domains). *A.*

*thaliana*, *O. sativa ssp. Indica*, and *O. sativa ssp. Japonica* were found to possess 85, 103, and 106 C2 domain-containing proteins, respectively. However, whereas no *Arabidopsis* protein contains more than one C2 domain, 19 and 20 proteins in the *O. sativa Indica* and *Japonica* species harbor two, three, or four C2 domains. None of the C2 domains of the plant species were found in the predicted protein kinases, whereas one C2 domain of *Entamoeba* was found N-terminal to the kinase domain of the *Dictyostelium* homologue of myosin light chain kinase (MLCK). *Dictyostelium's* MLCK (UniProt Q54W26) lacks the C2 domain N-terminal to the kinase catalytic domain.

The phosphotyrosine binding ability of the C2 domain of human PKCdelta was found to be dependent on a number of key residues.<sup>64</sup> These residues are found in an internal extension of the C2 domain that appears to be particular to PKCs. This extension was not found in the plant or amoeba C2 domains retrieved. It remains to be determined

whereas these residues are critical and confer the phosphotyrosine binding ability to PKC $\delta$ s or whether the C2 domain has a general phosphotyrosine binding ability that can be carried out by alternative residues.

Dual-specificity phosphatases and, most importantly, tyrosine phosphatases are strong indicators of a phosphotyrosine signaling system. A search with an HMM specific for dual-specificity phosphatases (PF00782) retrieved 15 candidates in *Entamoeba* and six in each of the plant species which could align meaningfully with previously characterized dual-specificity phosphatases (Fig. 4). A search with an HMM specific for tyrosine phosphatases (PF00102) returned two candidates in *Entamoeba*, one in *A. thaliana* and two in each of the *O. sativa* species (Fig. 5).

The presence of putative TKs, phosphotyrosine-binding domains as well as tyrosine phosphatases in *Entamoeba* and the three plant species lends support to the concept that phosphotyrosine signaling in eukaryotes appeared early in evolution, and is not a phenomenon exclusive to metazoans.<sup>62</sup> The ancestors of modern vascular plants diverged before the lineage leading to animals, fungi, and *Dictyostelids* did.<sup>29,30</sup> However, plants are known to carry out many of their signaling needs by means of serine/threonine receptor kinases, whereas receptor tyrosine kinase signaling appears to be more widespread in metazoans. It seems that plants and animals have solved their signaling needs according to similar yet distinct strategies. Still, the early appearance of putative phosphotyrosine signaling in evolution suggests that this might be an important, conserved, mechanism in plants and a number of plant-related species of ecological and economic importance such as red algae, chlorophyten green algae, and mosses.

### TKLs

The TKLs are present in only two of the fungal species: *C. neoformans* (2 TKLs) and *P. chrysosporium* (4 TKLs). The new observation made on the existence of putative TKLs in only two fungal species: *C. neoformans* (2 TKLs) and *P. chrysosporium* (4 TKLs); in the Apicomplexa *C. hominis* (4 TKLs), *P. yoelii* (4 TKLs), *P. falciparum* (6 TKLs); the amoeba *E. histolytica* (109); and in the algal species *C. merolae* (9 TKLs) and *T. pseudonana* (6 TKLs), suggests that TKLs might have been present at the very beginning of eukaryotic evolution. Although TKLs appear to be absent from most fungal species, it is reasonable to think that they were secondarily lost upon divergence and speciation. TKLs were found to be especially abundant in plants, with 776 TKLs in *A. thaliana* and around 1000 TKLs being encoded in the genomes of the two rice species, although its significance remains unknown. The unusually large number of TKLs found in the *E. histolytica* genome (109) relative to the number of genes encoded in its genome (9938), probably represents an

extreme metabolic adaptation of this single-celled organism to enhance its virulence in the environments it encounters.

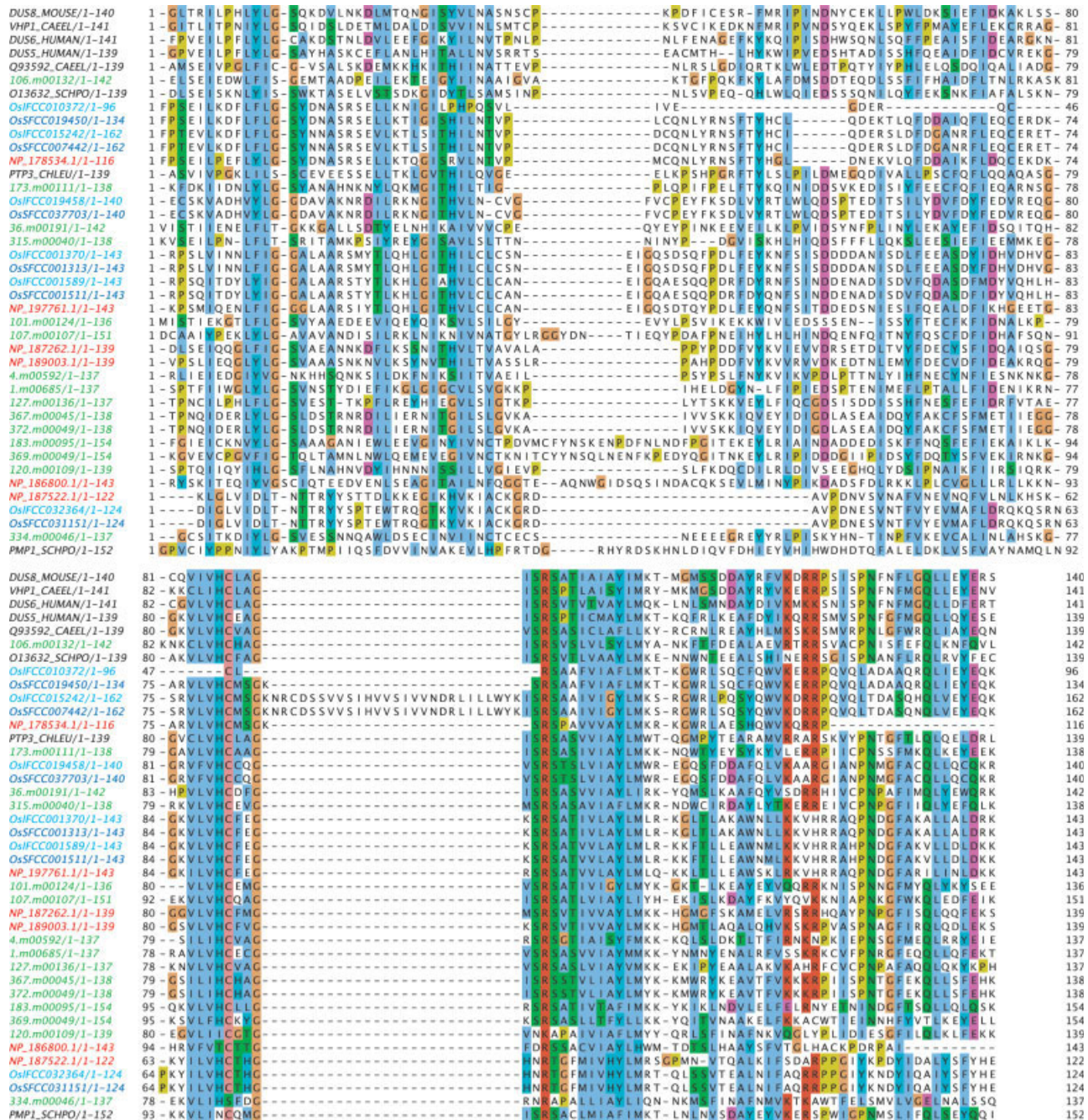
The understanding, both functional and evolutionary, of the primitive TKs and TKLs should shed light on their early functions and mechanisms of regulation, and will aid understanding of their subsequent diversification and incorporation into more sophisticated control circuits in animals.

### aPKs

Among the four families of aPKs, PIKKs appear in all the organisms examined here with no more than six PIKKs being found per genome. Their presence is indicative of their central importance in sensing DNA and RNA damage.<sup>66,67</sup> The PDHKs are present in all the organisms examined here with the exception of the Apicomplexan species, which possess divergent, tubular-shaped, mitochondria.<sup>68</sup> The PDHKs are known to regulate the pyruvate dehydrogenase complex (PDC), which catalyses the oxidative decarboxylation of pyruvate, linking the degradation of intracellular glycogen and extracellular glucose via glycolysis to the Krebs cycle.<sup>69</sup> Genes encoding enzymes of the TCA cycle have been predicted *in silico* for *P. falciparum*, but the parasite mitochondrion is likely to be divergent from mammalian mitochondria in other respects.<sup>68</sup> The RIO kinases, which seem to share the basic structural elements of the ePK fold<sup>8</sup> and whose function remains unknown, are present in all the eukaryotic lineages examined here, suggesting an essential function in eukaryotes. The copy number of RIO kinases per genome is never greater than three. The Alpha kinases, which also seem to share the ePK fold,<sup>70</sup> were absent from plants, one of the algal species (*C. merolae*), most fungi (except for two members in *N. crassa*) and from the Apicomplexa, but were present in *Entamoeba histolytica*.

The results presented in Table XII also allow the comparison of protein kinase families in organisms that are related phylogenetically. Fungi are believed to have appeared approximately 1 billion years ago, and the divergence of the Basidiomycota and Ascomycota probably occurred about 968 million years ago.<sup>30</sup> The group of fungi examined here encompasses species from different lineages. Although the predicted gene number of the various fungal species ranges from 5283 for *C. glabrata* to 11,777 for *P. chrysosporium*, it is remarkable that the number of protein kinases in most families of ePKs is rather similar in most instances: around 20 for the AGCs, 2–5 for the CK1s, around 25 for the CMGCs, and around 14 for the STEs. The CAMKs of fungi present a wider range from 16 for *C. neoformans* to 44 for *S. cerevisiae*. Some fungal lineages are known to have undergone whole genome duplication events, followed by massive gene losses.<sup>32</sup> Therefore, direct comparisons between





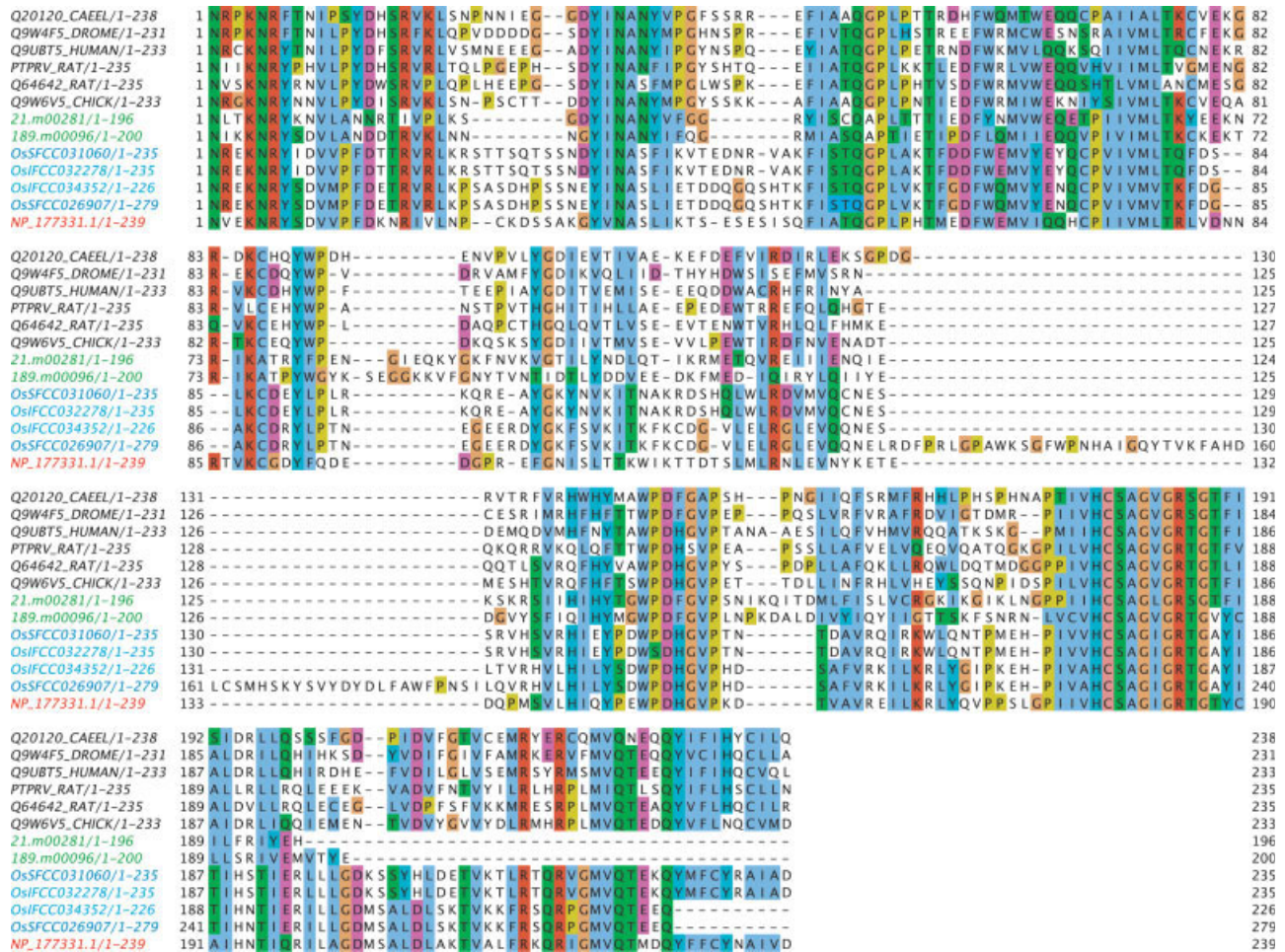
**Figure 4**

Multiple alignment of representative dual-specificity phosphatase catalytic domains with those identified for *E. histolytica* (15 members, sequence ids in green), *A. thaliana* (six members, sequence ids in red), *O. sativa* ssp. *Indica* (six members, sequence ids in light blue) and *O. sativa* ssp. *Japonica* (six members, sequence ids in dark blue). The sequence ids in black are representative examples of the Pfam signature PF00782. Dual-specificity phosphatases carry out the dephosphorylation of phosphoserine, phosphothreonine, and phosphotyrosine residues. The alignment was generated and displayed with Jalview.<sup>65</sup> [Color figure can be viewed in the online issue, which is available at [www.interscience.wiley.com](http://www.interscience.wiley.com).]

the protein kinase complements of the various fungal species is not strictly correct. However, the fact that some ePK families in fungi tend to be more compact than others in numbers might be an indication that in fungi some protein kinase families tolerate gene duplications

better than others, and that novel genes might be allowed to co-exist with the old copies, either contributing to the function of the former by being incorporated to existing signal transduction pathways without detrimental effects, or by evolving new functions (neo-functionalization). It





**Figure 5**

Alignment of representative tyrosine phosphatase sequences with those identified in *Entamoeba* (two members, sequence ids in green), *A. thaliana* (one member, sequence ids in red), and the two rice species (two members each, sequence ids in light blue for *O. sativa* ssp. *Indica*, and dark blue for *O. sativa* ssp. *japonica*). The sequence ids in black are representative examples of the Pfam signature PF00102. The presence of tyrosine phosphatases in the amoeba and plant species suggests that tyrosine phosphorylation signalling systems were present before the lineages leading to plants, animals, and amoebozoans split. The alignment was generated and displayed with Jalview.<sup>63</sup> [Color figure can be viewed in the online issue, which is available at [www.interscience.wiley.com](http://www.interscience.wiley.com).]

will be interesting to perform phylogenetic analysis of the fungal kinomes to determine what extent of each kinome is present in all fungi, which kinases are lineage-specific, and also whether the homologous kinase genes have different combinations of transcription factor binding sites that would promote differential transcription of homologous kinase genes, thereby accounting for the vast range of metabolic adaptations displayed by the fungal species considered here.

The difficulty in producing a catalogue of the major families of eukaryotic protein kinases for an organism is highlighted by the lack of discussion of kinase-mediated signal transduction pathways in all but the genome papers of *E. histolytica*, *A. thaliana*, *N. crassa*, and *F. rubripes*. Even though the genome of *A. thaliana* has been carefully annotated, the knowledge website ([http://](http://www.arabidopsis.org)

[www.arabidopsis.org](http://www.arabidopsis.org)) features only 1031 open-reading frames annotated automatically as protein kinases. Application of Library 2 predicts that this plant contains 1301 putative protein kinases. The plant species examined here were found to have an unusually large number of kinases in comparison with, for example, the human genome. It must be taken into account that reported whole genome duplication events in plants<sup>45</sup> prevent the meaningful comparison between genome size, phenotypic complexity, and kinome size.

At the time of the publication of the *Arabidopsis* genome, the *ab initio* gene models were validated with known ESTs and protein sequences, the researchers finding that 93% of ESTs of *Arabidopsis* matched gene models, with less than 1% of ESTs matching noncoding regions.<sup>45</sup> With the rice species, only 61% of *ab initio*

gene models were found to have a high identity match with rice ESTs or full-length cDNAs. Whereas 71% of predicted rice proteins were found to have a homologue in the *Arabidopsis* proteome, 89.8% of the proteins from the *Arabidopsis* genome have a homologue in the rice proteome.<sup>71</sup> This suggests that caution must be exercised when considering the gene models of the rice species that do not have EST coverage or cannot be verified otherwise. Still, the predicted putative ePKs for the rice species (1276 and 1330, Table XII) agree with those predicted by the authors (between 1075 and 1425, depending on the InterPro signature used<sup>71</sup>).

The use of Library 2 provides a starting point for characterizing the kinomes of organisms that are representative of key points in the evolutionary history of eukaryotes. Detailed kinase family phylogenetic analyses will help identify orthologous subfamilies across the vast spectrum of eukaryotic life. This should illuminate the early stages of protein kinase evolution and help understanding of how gene loss, duplication, and innovation within each eukaryotic lineage correlate with specific functions. Protein kinase orthologue and paralogue identification is also an essential step in the process of drug target identification for establishing assay and animal models. Paralogues are known to introduce complications such as selectivity issues, pleiotropy, and functional redundancy of targets, all of which are critical to assessing the druggability of a particular protein kinase. This wealth of detailed kinomes will help extend Library 2 to classify automatically and reliably to the subfamily and sub-subfamily level. This will be particularly important for the larger families (e.g. CAMK), where it could help identify further subfamilies in their own right which might have been obscured hitherto by the limited number of kinomes available. Having enough information on what the key residues/regions of a particular subfamily are, coupled to structural information, will help highlighting those regions of the structure that are important for the subfamily, but which might not be captured appropriately by the HMM.

Library 2 has been used to assist in the analysis of the kinomes of the human parasitic protozoan *Trichomonas vaginalis*<sup>59</sup> and the human parasitic nematode *Brugia malayi* (manuscript submitted). A server is being developed that will integrate the ability to scan peptide sequences with Library 2, and display precalculated analyses of the kinomes of a number of species of biological, medical, and economic interest. This facility will be available under the URL: <http://www.compbio.dundee.ac.uk/kinomes/>. Library 2 is available from the authors upon request.

## CONCLUSIONS

In this article, a sensitive library of HMMs (Library 2) has been developed to identify and subclassify protein kinase catalytic domains into one of the 12 accepted families. The main conclusions of the work are:

1. Library 2 showed a protein kinase family misclassification rate of zero on the characterized kinomes of *H. sapiens*, *M. musculus*, *D. melanogaster*, *C. elegans*, *S. cerevisiae*, *P. falciparum*, and *D. discoideum*. Library 2 was able to retrieve all protein kinases from KinBase in a database search. In contrast, BLASTP could only retrieve on average 1/3 of protein kinases from the same database. The general kinase domain HMM PF00069 failed to identify between 2.61% and 10.09% of the ePKs of characterized kinomes. These properties make Library 2 a useful tool for database searching of kinases and their automatic classification into families.
2. 27/38 "Other" kinases of yeast, and their syntenic homologues in the fungus *A. gossypii*, were classified using Library 2 into the main ePK families; AGC,<sup>5</sup> CAMK,<sup>17</sup> CMGC,<sup>4</sup> and STE,<sup>1</sup> raising the family-level classification rate of the yeast ePKs from 66.96 to 90.43%.
3. The application of Library 2 to 21 eukaryotic genomes showed that five ePK (AGC, CAMK, CK1, CMGC, and STE) and two aPK (PIKK and RIO) protein kinase families were present in all eukaryotic genomes analyzed. These seven families are likely to be indispensable to all eukaryotic life forms. The families present in specific lines of descent (RGC, TK, TKL, Alpha, PDHK) are likely to be late innovations. Alternatively, their absence from a particular phylogenetic group might indicate that they have been lost secondarily.
4. No TKs were found in fungi. The plants *A. thaliana*, *O. sativa ssp. Indica*, and *O. sativa ssp. Japonica* were found to encode 2, 6, and 7 putative TKs, respectively. To our knowledge, TKs have never been predicted in plants before, although a tyrosine phosphorylation system is known to operate in *A. thaliana*.
5. The confirmation of the presence of seven putative TKs in *Entamoeba histolytica*, together with the putative TKs in plants, suggests that a primitive set of TKs might have predated the radiation of eukaryotic organisms.
6. The fungi *C. neoformans* and *P. chrysosporium* were found to encode 2 and 4 putative TKs, respectively. Putative TKs were also found in the Apicomplexa *C. hominis* (4), *P. yoelii* (4), *P. falciparum* (6), the amoeba *E. histolytica* (109), and in the algae *T. pseudonana* (6), and *C. merolae* (9). TKs were found to be especially abundant in plants: *A. thaliana* (776), *O. sativa ssp. Indica* (1010), *O. sativa ssp. Japonica* (969). These data suggest that (a) TKs predated the radiation of eukaryotic life forms; (b) TKs in plants probably carry out important and specific functions to plants; (c) TKs have been lost secondarily in some fungi.
7. The atypical protein kinase family PIKK is ubiquitous in all eukaryotic genomes examined here, with never more than six PIKKs per genome. Their universal presence reflects their importance in sensing DNA and RNA damage. RIO kinases were also found to be



ubiquitous, although their biological substrates remain unknown.

## ACKNOWLEDGMENTS

The authors thank Professor Grahame Hardie and Dr. David Martin for many insightful discussions, and Dr. Jonathan Monk for computer support. The authors also wish to thank the reviewers for their insightful comments.

## REFERENCES

- Cohen P. The regulation of protein function by multisite phosphorylation—a 25 year update. *Trends Biochem Sci* 2000;25:596–601.
- Cohen P. The role of protein phosphorylation in human health and disease. The Sir Hans Krebs medal lecture. *Eur J Biochem* 2001;268:5001–5010.
- Cohen P. Protein kinases—the major drug targets of the twenty-first century? *Nat Rev Drug Discov* 2002;1:309–315.
- Manning G, Plowman GD, Hunter T, Sudarsanam S. Evolution of protein kinase signaling from yeast to man. *Trends Biochem Sci* 2002;27:514–520.
- Hanks SK, Hunter T. Protein kinases 6. The eukaryotic protein kinase superfamily: kinase (catalytic) domain structure and classification. *FASEB J* 1995;9:576–596.
- Manning G, Whyte DB, Martinez R, Hunter T, Sudarsanam S. The protein kinase complement of the human genome. *Science* 2002;298:1912–1934.
- Angermayr M, Bandlow W. The general regulatory factor Reb1p controls basal, but not Gal4p-mediated, transcription of the GCY1 gene in yeast. *Mol Gen Genet* 1997;256:682–689.
- LaRonde-LeBlanc N, Wlodawer A. Crystal structure of *A. fulgidus* Rio2 defines a new family of serine protein kinases. *Structure* 2004;12:1585–1594.
- Yamaguchi H, Matsushita M, Nairn AC, Kuriyan J. Crystal structure of the atypical protein kinase domain of a TRP channel with phosphotransferase activity. *Mol Cell* 2001;7:1047–1057.
- Wu CH, Apweiler R, Bairoch A, Natale DA, Barker WC, Boeckmann B, Ferro S, Gasteiger E, Huang H, Lopez R, Magrane M, Martin MJ, Mazumder R, O'Donovan C, Redaschi N, Suzek B. The universal protein resource (UniProt): an expanding universe of protein information. *Nucleic Acids Res* 2006;34 (database issue):D187–D191.
- Plowman GD, Sudarsanam S, Bingham J, Whyte D, Hunter T. The protein kinases of *Caenorhabditis elegans*: a model for signal transduction in multicellular organisms. *Proc Natl Acad Sci USA* 1999;96:13603–13610.
- Hunter T, Plowman GD. The protein kinases of budding yeast: six score and more. *Trends Biochem Sci* 1997;22:18–22.
- Morrison DK, Murakami MS, Cleghon V. Protein kinases and phosphatases in the *Drosophila* genome. *J Cell Biol* 2000;150:F57–F62.
- Caenepeel S, Charyczak G, Sudarsanam S, Hunter T, Manning G. The mouse kinome: discovery and comparative genomics of all mouse protein kinases. *Proc Natl Acad Sci USA* 2004;101:11707–11712.
- Ward P, Equinet L, Packer J, Doerig C. Protein kinases of the human malaria parasite *Plasmodium falciparum*: the kinome of a divergent eukaryote. *BMC Genomics* 2004;5:79.
- Goldberg JM, Manning G, Liu A, Fey P, Pilcher KE, Xu Y, Smith JL. The dictyostelium kinome—analysis of the protein kinases from a simple model organism. *PLoS Genet* 2006;2:e38.
- Altschul SF, Madden TL, Schaffer AA, Zhang J, Zhang Z, Miller W, Lipman DJ. Gapped BLAST and PSI-BLAST: a new generation of protein database search programs. *Nucleic Acids Res* 1997;25:3389–3402.
- Bateman A, Coin L, Durbin R, Finn RD, Hollich V, Griffiths-Jones S, Khanna A, Marshall M, Moxon S, Sonnhammer EL, Studholme DJ, Yeats C, Eddy SR. The Pfam protein families database. *Nucleic Acids Res* 2004;32 (database issue):D138–D141.
- Hanks SK. Genomic analysis of the eukaryotic protein kinase superfamily: a perspective. *Genome Biol* 2003;4:111.
- Krogh A, Brown M, Mian IS, Sjolander K, Haussler D. Hidden Markov models in computational biology. Applications to protein modeling. *J Mol Biol* 1994;235:1501–1531.
- Park J, Karplus K, Barrett C, Hughey R, Haussler D, Hubbard T, Chothia C. Sequence comparisons using multiple sequences detect three times as many remote homologues as pairwise methods. *J Mol Biol* 1998;284:1201–1210.
- Madera M, Vogel C, Kummerfeld SK, Chothia C, Gough J. The SUPERFAMILY database in 2004: additions and improvements. *Nucleic Acids Res* 2004;32 (database issue):D235–D239.
- Gough J, Karplus K, Hughey R, Chothia C. Assignment of homology to genome sequences using a library of hidden Markov models that represent all proteins of known structure. *J Mol Biol* 2001;313:903–919.
- Brown D, Krishnamurthy N, Dale JM, Christopher W, Sjolander K. Subfamily hmms in functional genomics. *Pac Symp Biocomput* 2005:322–333.
- Gudi R, Bowker-Kinley MM, Kedishvili NY, Zhao Y, Popov KM. Diversity of the pyruvate dehydrogenase kinase gene family in humans. *J Biol Chem* 1995;270:28989–28994.
- Barton GJ, Sternberg MJ. Evaluation and improvements in the automatic alignment of protein sequences. *Protein Eng* 1987;1:89–94.
- Eddy SR. Profile hidden Markov models. *Bioinformatics* 1998;14:755–763.
- Baldauf SL. The deep roots of eukaryotes. *Science* 2003;300:1703–1706.
- Baldauf SL, Doolittle WF. Origin and evolution of the slime molds (Mycetozoa). *Proc Natl Acad Sci USA* 1997;94:12007–12012.
- Hedges SB, Blair JE, Venturi ML, Shreeve JL. A molecular timescale of eukaryote evolution and the rise of complex multicellular life. *BMC Evol Biol* 2004;4:2.
- Dietrich FS, Voegeli S, Brachat S, Lerch A, Gates K, Steiner S, Mohr C, Pohlmann R, Luedi P, Choi S, Wing RA, Flavier A, Gaffney TD, Philippsen P. The *Ashbya gossypii* genome as a tool for mapping the ancient *Saccharomyces cerevisiae* genome. *Science* 2004;304:304–307.
- Kellis M, Birren BW, Lander ES. Proof and evolutionary analysis of ancient genome duplication in the yeast *Saccharomyces cerevisiae*. *Nature* 2004;428:617–624.
- Armbrust EV, Berges JA, Bowler C, Green BR, Martinez D, Putnam NH, Zhou S, Allen AE, Apt KE, Bechner M, Brzezinski MA, Chaal BK, Chiovitti A, Davis AK, Demarest MS, Detter JC, Glavina T, Goodstein D, Hadi MZ, Hellsten U, Hildebrand M, Jenkins BD, Jurka J, Kapitonov VV, Kroger N, Lau WW, Lane TW, Larimer FW, Lippmeier JC, Lucas S, Medina M, Montsant A, Obornik M, Parker MS, Palenik B, Pazour GJ, Richardson PM, Rynearson TA, Saito MA, Schwartz DC, Thamtrakoln K, Valentin K, Vardi A, Wilkerson FP, Rokhsar DS. The genome of the diatom *Thalassiosira pseudonana*: ecology, evolution, and metabolism. *Science* 2004;306:79–86.
- Matsuzaki M, Misumi O, Shin IT, Maruyama S, Takahara M, Miyagishima SY, Mori T, Nishida K, Yagisawa F, Nishida K, Yoshida Y, Nishimura Y, Nakao S, Kobayashi T, Momoyama Y, Higashiyama T, Minoda A, Sano M, Nomoto H, Oishi K, Hayashi H, Ohta F, Nishizaka S, Haga S, Miura S, Morishita T, Kabeya Y, Terasawa K, Suzuki Y, Ishii Y, Asakawa S, Takano H, Ohta N, Kuroiwa H, Tanaka K, Shimizu N, Sugano S, Sato N, Nozaki H, Ogasawara N, Kohara Y, Kuroiwa T. Genome sequence of the ultrasmall unicellular red alga *Cyanidioschyzon merolae* 10 D. *Nature* 2004;428:653–657.
- Holt RA, Subramanian GM, Halpern A, Sutton GG, Charlab R, Nusskern DR, Wincker P, Clark AG, Ribeiro JM, Wides R, Salzberg SL,

- Lofthus B, Yandell M, Majoros WH, Rusch DB, Lai Z, Kraft CL, Abril JF, Anthouard V, Arensburger P, Atkinson PW, Baden H, de Berardinis V, Baldwin D, Benes V, Biedler J, Blass C, Bolanos R, Boscos D, Barnstead M, Cai S, Center A, Chaturverdi K, Christophides GK, Chrystal MA, Clamp M, Cravchik A, Curwen V, Dana A, Delcher A, Dew I, Evans CA, Flanigan M, Grundschober-Freimoser A, Friedli L, Gu Z, Guan P, Guigo R, Hillenmeyer ME, Hladun SL, Hogan JR, Hong YW, Hoover J, Jaillon O, Ke Z, Kodira C, Kokoza E, Koutsos A, Letunic I, Levitsky A, Liang Y, Lin JJ, Lobo NF, Lopez JR, Malek JA, McIntosh TC, Meister S, Miller J, Mobarry C, Mongin E, Murphy SD, O'Brochta DA, Pfannkoch C, Qi R, Regier MA, Remington K, Shao H, Sharakhova MV, Sitter CD, Shetty J, Smith TJ, Strong R, Sun J, Thomasova D, Ton LQ, Topalis P, Tu Z, Unger ME, Walenz B, Wang A, Wang J, Wang M, Wang X, Woodford KJ, Wortman JR, Wu M, Yao A, Zdobnov EM, Zhang H, Zhao Q, Zhao S, Zhu SC, Zhimulev I, Coluzzi M, della Torre A, Roth CW, Louis C, Kalush F, Mural RJ, Myers EW, Adams MD, Smith HO, Broder S, Gardner MJ, Fraser CM, Birney E, Bork P, Brey PT, Venter JC, Weissenbach J, Kafatos FC, Collins FH, Hoffman SL. The genome sequence of the malaria mosquito *Anopheles gambiae*. Science 2002;298:129–149.
36. Dehal P, Satou Y, Campbell RK, Chapman J, Degnan B, De Tomaso A, Davidson B, Di Gregorio A, Gelpke M, Goodstein DM, Harafuji N, Hastings KE, Ho I, Hotta K, Huang W, Kawashima T, Lemaire P, Martinez D, Meinertzhagen IA, Necula S, Nonaka M, Putnam N, Rash S, Saiga H, Satake M, Terry A, Yamada L, Wang HG, Awazu S, Azumi K, Boore J, Branno M, Chin-Bow S, DeSantis R, Doyle S, Francino P, Keys DN, Haga S, Hayashi H, Hino k, Imai KS, Inaba K, Kano S, Kobayashi K, Kobayashi M, Lee BI, Makabe KW, Manohar C, Matassi G, Medina M, Mochizuki Y, Mount S, Morishita T, Miura S, Nakayama A, Nishizaka S, Nomoto H, Ohta F, Oishi K, Rigoutsos I, Sano M, Sasaki A, Sasakura Y, Shoguchi E, Shin-i T, Spagnuolo A, Stainier D, Suzuki MM, Tassy O, Takatori N, Tokuoka M, Yagi K, Yoshizaki F, Wada S, Zhang C, Hyatt PD, Larimer F, Detter C, Doggett N, Glavina T, Hawkins T, Richardson P, Lucas S, Kohara Y, Levine M, Satoh N, Rokhsar DS. The draft genome of *Ciona intestinalis*: insights into chordate and vertebrate origins. Science 2002;298:2157–2167.
37. Jaillon O, Aury JM, Brunet F, Petit JL, Stange-Thomann N, Mauceli E, Bouneau L, Fischer C, Ozouf-Costaz C, Bernot A, Nicaud S, Jaffe D, Fisher S, Lutfalla G, Dossat C, Segurens B, Dasilva C, Salanoubat M, Levy M, Boudet N, Castellano S, Anthouard V, Jubin C, Castelli V, Katinka M, Vacherie B, Biemont C, Skalli Z, Cattolico L, Poulain J, De Berardinis V, Cruaud C, Duprat S, Brottier P, Coutanceau JP, Gouzy J, Parra G, Lardier G, Chapple C, McKernan KJ, McEwan P, Bosak S, Kellis M, Volff JN, Guigo R, Zody MC, Mesirov J, Lindblad-Toh K, Birren B, Nusbaum C, Kahn D, Robinson-Rechavi M, Laudet V, Schachter V, Quetier F, Saurin W, Scarpelli C, Wincker P, Lander ES, Weissenbach J. Genome duplication in the teleost fish *Tetraodon nigroviridis* reveals the early vertebrate proto-karyotype. Nature 2004;431:946–957.
38. Aparicio S, Chapman J, Stupka E, Putnam N, Chia JM, Dehal P, Christoffels A, Rash S, Hoon S, Smit A, Gelpke MD, Roach J, Oh T, Ho IY, Wong M, Detter C, Verhoef F, Predki P, Tay A, Lucas S, Richardson P, Smith SF, Clark MS, Edwards YJ, Doggett N, Zharkikh A, Tavtigian SV, Pruss D, Barnstead M, Evans C, Baden H, Powell J, Glusman G, Rowen L, Hood L, Tan YH, Elgar G, Hawkins T, Venkatesh B, Rokhsar D, Brenner S. Whole-genome shotgun assembly and analysis of the genome of *Fugu rubripes*. Science 2002;297:1301–1310.
39. Dujon B, Sherman D, Fischer G, Durrens P, Casaregola S, Lafontaine I, De Montigny J, Marck C, Neuveglise C, Talla E, Goffard N, Frangeul L, Aigle M, Anthouard V, Babour A, Barbe V, Barnay S, Blanchin S, Beckerich JM, Beyne E, Bleykrasten C, Boisrame A, Boyer J, Cattolico L, Confanioli F, De Daruvar A, Despons L, Fabre E, Fairhead C, Ferry-Dumazet H, Groppi A, Hantraye F, Hennequin C, Jauniaux N, Joyet P, Kachouri R, Kerrest A, Koszul R, Lemaire M, Lesur I, Ma L, Muller H, Nicaud JM, Nikolski M, Oztas S, Ozier-Kalogeropoulos O, Pellenz S, Potier S, Richard GF, Straub ML, Suleau A, Swennen D, Tekcia F, Wesolowski-Louvel M, Westhof E, Wirth B, Zeniou-Meyer M, Zivanovic I, Bolotin-Fukuhara M, Thierry A, Bouchier C, Caudron B, Scarpelli C, Gaillardin C, Weissenbach J, Wincker P, Souciet JL. Genome evolution in yeasts. Nature 2004;430:35–44.
40. Lofthus BJ, Fung E, Roncaglia P, Rowley D, Amedeo P, Bruno D, Vamathevan J, Miranda M, Anderson IJ, Fraser JA, Allen JE, Bosdet IE, Brent MR, Chiu R, Doering TL, Donlin MJ, D'Souza CA, Fox DS, Grinberg V, Fu J, Fukushima M, Haas BJ, Huang JC, Janbon G, Jones SJ, Koo HL, Krzywinski MI, Kwon-Chung, Lengeler KB, Maiti R, Marra MA, Marra RE, Mathewson CA, Mitchell TG, Perteau M, Riggs FR, Salzberg SL, Schein JE, Shvartsbeyn A, Shin H, Shumway M, Specht CA, Suh BB, Tenney A, Utterback TR, Wickes BL, Wortman JR, Wye NH, Kronstad JW, Lodge JK, Heitman J, Davis RW, Fraser CM, Hyman RW. The genome of the basidiomycetous yeast and human pathogen *Cryptococcus neoformans*. Science 2005;307:1321–1324.
41. Galagan JE, Calvo SE, Borkovich KA, Selker EU, Read ND, Jaffe D, FitzHugh W, Ma LJ, Smirnov S, Purcell S, Rehman B, Elkins T, Engels R, Wang S, Nielsen CB, Butler J, Endrizzi M, Qui D, Ianakiev P, Bell-Pedersen D, Nelson MA, Werner-Washburne M, Selitrennikoff CP, Kinsey JA, Braun EL, Zelter A, Schulte U, Kothe GO, Jedd G, Mewes W, Staben C, Marcotte E, Greenberg D, Roy A, Foley K, Naylor J, Strange-Thomann N, Barrett R, Gnerre S, Kamal M, Kamvyselis M, Mauceli E, Bielke C, Rudd S, Frishman D, Krysstofova S, Rasmussen C, Metznerberg RL, Perkins DD, Kroken S, Cogoni C, Macino G, Catcheside D, Li W, Pratt RJ, Osmani SA, DeSouza CP, Glass L, Orbach MJ, Berglund JA, Voelker R, Yarden O, Plamann M, Seiler S, Dunlap J, Radford A, Aramayo R, Natvig DO, Alex LA, Mannhaupt G, Ebbold DJ, Freitag M, Paulsen I, Sachs MS, Lander ES, Nusbaum C, Birren B. The genome sequence of the filamentous fungus *Neurospora crassa*. Nature 2003;422:859–868.
42. Martinez D, Larrondo LF, Putnam N, Gelpke MD, Huang K, Chapman J, Helfenbein KG, Ramaiya P, Detter JC, Larimer F, Coutinho PM, Henrissat B, Berka R, Cullen D, Rokhsar D. Genome sequence of the lignocellulose degrading fungus *Phanerochaete chrysosporium* strain RP78. Nat Biotechnol 2004;22:695–700.
43. Wood V, Gwilliam R, Rajandream MA, Lyne M, Lyne R, Stewart A, Sgouros J, Peat N, Hayles J, Baker S, Basham D, Bowman S, Brooks K, Brown D, Brown S, Chillingworth T, Churcher C, Collins M, Connor R, Cronin A, Davis P, Feltwell T, Fraser A, Gentles S, Goble A, Hamlin N, Harris D, Hidalgo J, Hodgson G, Holroyd S, Hornsby T, Howarth S, Huckle EJ, Hunt S, Jagels K, James K, Jones L, Jones M, Leather S, McDonald S, McLean J, Mooney P, Moule S, Mungall K, Murphy I, Niblett D, Odell C, Oliver K, O'Neil S, Pearson D, Quail MA, Rabinowitsch E, Rutherford K, Rutter S, Saunders D, Seeger K, Sharp S, Skelton J, Simmonds M, Squares R, Squares S, Stevens K, Taylor K, Taylor RG, Tivey A, Walsh S, Warren T, Whitehead S, Woodward J, Volckaert G, Aert R, Robben J, Grymonprez B, Weltjens I, Vanstreels E, Rieger M, Schafer M, Muller-Auer S, Gabel C, Fuchs M, Dusterhoft A, Fritzc C, Holzer E, Moestl D, Hilbert H, Borzym K, Langer I, Beck A, Lehrach H, Reinhardt R, Pohl TM, Eger P, Zimmermann W, Wedler H, Wambutt R, Purnelle B, Goffeau A, Cadieu E, Dreano S, Gloux S, Lelaure V, Mottier S, Galibert F, Aves SJ, Ziang Z, Hunt C, Moore K, Hurst SM, Lucas M, Rochet M, Gaillardin C, Tallada VA, Garzon A, Thode G, Daga RR, Cruzado L, Jimenez J, Sanchez M, del Rey F, Benito J, Dominguez A, Revuelta JL, Moreno S, Armstrong J, Forsburg SL, Cerutti L, Lowe T, McCombie WR, Paulsen I, Potashkin J, Shpakovski GV, Ussery D, Barrell BG, Nurse P. The genome sequence of *Schizosaccharomyces pombe*. Nature 2002;415:871–880.
44. Gibbs RA, Weinstock GM, Metzker ML, Muzny DM, Sodergren EJ, Scherer S, Scott G, Steffen D, Worley KC, Burch PE, Okwuonu G, Hines S, Lewis L, DeRamo C, Delgado O, Dugan-Rocha S, Miner G, Morgan M, Hawes A, Gill R, Celera, Holt RA, Adams MD,



- Amanatides PG, Baden-Tillson H, Barnstead M, Chin S, Evans CA, Ferriera S, Fosler C, Glodek A, Gu Z, Jennings D, Kraft CL, Nguyen T, Pfannkoch CM, Sitter C, Sutton GG, Venter JC, Woodage T, Smith D, Lee HM, Gustafson E, Cahill P, Kana A, Doucette-Stamm L, Weinstock K, Fachtel K, Weiss RB, Dunn DM, Green ED, Blakesley RW, Bouffard GG, De Jong PJ, Osoegawa K, Zhu B, Marra M, Schein J, Bosdet I, Fjell C, Jones S, Krzywinski M, Mathewson C, Siddiqui A, Wye N, McPherson J, Zhao S, Fraser CM, Shetty J, Shatsman S, Geer K, Chen Y, Abramson S, Nierman WC, Haylak PH, Chen R, Durbin KJ, Egan A, Ren Y, Song XZ, Li B, Liu Y, Qin X, Cawley S, Worley KC, Cooney AJ, D'Souza LM, Martin K, Wu JQ, Gonzalez-Garay ML, Jackson AR, Kalafus KJ, McLeod MP, Milosavljevic A, Virk D, Volkov A, Wheeler DA, Zhang Z, Bailey JA, Eichler EE, Tuzun E, Birney E, Mongin E, Ureta-Vidal A, Woodward C, Zdobnov E, Bork P, Suyama M, Torrents D, Alexandersson M, Trask BJ, Young JM, Huang H, Wang H, Xing H, Daniels S, Gietzen D, Schmidt J, Stevens K, Vitt U, Wingrove J, Camara F, Mar Alba M, Abril JF, Guigo R, Smit A, Dubchak I, Rubin EM, Couronne O, Poliakov A, Hubner N, Ganten D, Goesele C, Hummel O, Kreidler T, Lee YA, Monti J, Schulz H, Zimdahl H, Himmelbauer H, Lehrach H, Jacob HJ, Bromberg S, Gullings-Handley J, Jensen-Seaman MP, Kwitek AE, Lazar J, Pasko D, Tonelato PJ, Twigger S, Ponting CP, Duarte JM, Rice S, Goodstadt L, Beatson SA, Emes RD, Winter EE, Webber C, Brandt P, Nyakatura G, Adetobi M, Chiaromonte F, Elnitski L, Eswara P, Hardison RC, Hou M, Kolbe D, Makova K, Miller W, Nekrutenko A, Riemer C, Schwartz S, Taylor J, Yang S, Zhang Y, Lindpaintner K, Andrews TD, Caccamo M, Clamp M, Clarke L, Curwen V, Durbin R, Eyras E, Searle SM, Cooper GM, Batzoglu S, Brudno M, Sidow A, Stone EA, Venter JC, Payseur BA, Bourque G, Lopez-Otin C, Puente XS, Chakrabarti K, Chatterji S, Dewey C, Pachter L, Bray N, Yap VB, Caspi A, Tesler G, Pevzner PA, Haussler D, Roskin KM, Baertsch R, Clawson H, Furey TS, Hinrichs AS, Karolchik D, Kent WJ, Rosenbloom KR, Trumbower H, Weirauch M, Cooper DN, Stenson PD, Ma B, Brent M, Arumugam M, Shteynberg D, Copley RR, Taylor MS, Riethman H, Mudunuri U, Peterson J, Guyer M, Felsenfeld A, Old S, Mockrin S, Collins F. Genome sequence of the Brown Norway rat yields insights into mammalian evolution. *Nature* 2004; 428:493–521.
45. Initiative AG. Analysis of the genome sequence of the flowering plant *Arabidopsis thaliana*. *Nature* 2000;408:796–815.
  46. Yu J, Hu S, Wang J, Wong GK, Li S, Liu B, Deng Y, Dai L, Zhou Y, Zhang X, Cao M, Liu J, Sun J, Tang J, Chen Y, Huang X, Lin W, Ye C, Tong W, Cong L, Geng J, Han Y, Li L, Li W, Hu G, Huang X, Li W, Li J, Liu Z, Li L, Liu J, Qi Q, Liu J, Li L, Li T, Wang X, Lu H, Wu T, Zhu M, Ni P, Han H, Dong W, Ren X, Feng X, Cui P, Li X, Wang H, Xu X, Zhai W, Xu Z, Zhang J, He S, Zhang J, Xu J, Zhang K, Zheng X, Dong J, Zeng W, Tao L, Ye J, Tan J, Ren X, Chen X, He J, Liu D, Tian W, Tian C, Xia H, Bao Q, Li G, Gao H, Cao T, Wang J, Zhao W, Li P, Chen W, Wang X, Zhang Y, Hu J, Wang J, Liu S, Yang J, Zhang G, Xiong Y, Li Z, Mao L, Zhou C, Zhu Z, Chen R, Hao B, Zheng W, Chen S, Guo W, Li G, Liu S, Tao M, Wang J, Zhu L, Yuan L, Yang H. A draft sequence of the rice genome (*Oryza sativa* L. ssp. *indica*). *Science* 2002;296: 79–92.
  47. Goff SA, Ricke D, Lan TH, Presting G, Wang R, Dunn M, Glazebrook J, Sessions A, Oeller P, Varma H, Hadley D, Hutchison D, Martin C, Katagiri F, Lange BM, Moughamer T, Xia Y, Budworth P, Zhong J, Miguel T, Paszkowski U, Zhang S, Colbert M, Sun WL, Chen L, Cooper B, Park S, Wood TC, Mao L, Quail P, Wing R, Dean R, Yu Y, Zharkikh A, Shen R, Sahasrabudhe S, Thomas A, Cannings R, Gutin A, Pruss D, Reid J, Tavtigian S, Mitchell J, Eldredge G, Scholl T, Miller RM, Bhatnagar S, Adey N, Rubano T, Tusneem N, Robinson R, Feldhaus J, Macalma T, Oliphant A, Briggs S. A draft sequence of the rice genome (*Oryza sativa* L. ssp. *japonica*). *Science* 2002;296:92–100.
  48. Xu P, Widmer G, Wang Y, Ozaki LS, Alves JM, Serrano MG, Puiu D, Manque P, Akiyoshi D, Mackey AJ, Pearson WR, Dear PH, Bankier AT, Peterson DL, Abrahamsen MS, Kapur V, Tzipori S, Buck GA. The genome of *Cryptosporidium hominis*. *Nature* 2004;431:1107–1112.
  49. Loftus B, Anderson I, Davies R, Alsmark UC, Samuelson J, Amedeo P, Roncaglia P, Berriman M, Hirt RP, Mann BJ, Nozaki T, Suh B, Pop M, Duchene M, Ackers J, Tannich E, Leippe M, Hofer M, Bruchhaus I, Willhoeft U, Bhattacharya A, Chillingworth T, Churcher C, Hance Z, Harris B, Harris D, Jagels K, Moule S, Mungall K, Ormond D, Squares R, Whitehead S, Quail MA, Rabinowitsch E, Norbertczak H, Price C, Wang Z, Guillen N, Gilchrist C, Stroup SE, Bhattacharya S, Lohia A, Foster PG, Sicheritz-Ponten T, Weber C, Singh U, Mukherjee C, El-Sayed NM, Petri WA Jr, Clark CG, Embley TM, Barrell B, Fraser CM, Hall N. The genome of the protist parasite *Entamoeba histolytica*. *Nature* 2005;433: 865–868.
  50. Carlton JM, Angiuoli SV, Suh BB, Kooij TW, Perteau M, Silva JC, Ermolaeva MD, Allen JE, Selengut JD, Koo HL, Peterson JD, Pop M, Kosack DS, Shumway MF, Bidwell SL, Shallom SJ, van Aken SE, Riedmuller SB, Feldblyum TV, Cho JK, Quackenbush J, Sedegah M, Shoaibi A, Cummings LM, Florens L, Yates JR, Raine JD, Sinden RE, Harris MA, Cunningham DA, Preiser PR, Bergman LW, Vaidya AB, van Lin LH, Janse CJ, Waters AP, Smith HO, White OR, Salzberg SL, Venter JC, Fraser CM, Hoffman SL, Gardner MJ, Carucci DJ. Genome sequence and comparative analysis of the model rodent malaria parasite *Plasmodium yoelii yoelii*. *Nature* 2002;419:512–519.
  51. Alonso A, Sasín J, Bottini N, Friedberg I, Friedberg I, Osterman A, Godzik A, Hunter T, Dixon J, Mustelin T. Protein tyrosine phosphatases in the human genome. *Cell* 2004;117:699–711.
  52. Hunter T. A thousand and one protein kinases. *Cell* 1987;50:823–829.
  53. Mustelin T, Feng GS, Bottini N, Alonso A, Kholod N, Birle D, Merlo J, Huynh H. Protein tyrosine phosphatases. *Front Biosci* 7: d85–142, 2002.
  54. Mustelin T, Abraham RT, Rudd CE, Alonso A, Merlo JJ. Protein tyrosine phosphorylation in T cell signaling. *Front Biosci* 2002;7: d918–d969.
  55. Rudrabhatla P, Reddy MM, Rajasekharan R. Genome-wide analysis and experimentation of plant serine/threonine/tyrosine-specific protein kinases. *Plant Mol Biol* 2006;60:293–319.
  56. Barizza E, Lo Schiavo F, Terzi M, Filippini F. Evidence suggesting protein tyrosine phosphorylation in plants depends on the developmental conditions. *FEBS Lett* 1999;447:191–194.
  57. Kameyama K, Kishi Y, Yoshimura M, Kanzawa N, Sameshima M, Tsuchiya T. Tyrosine phosphorylation in plant bending. *Nature* 2000;407:37.
  58. Shiu SH, Li WH. Origins, lineage-specific expansions, and multiple losses of tyrosine kinases in eukaryotes. *Mol Biol Evol* 2004;21:828–840.
  59. Carlton JM, Hirt RP, Silva JC, Delcher AL, Schatz M, Zhao Q, Wortman JR, Bidwell SL, Alsmark UC, Besteiro S, Sicheritz-Ponten T, Noel CJ, Dacks JB, Foster PG, Simillion C, Van de Peer Y, Miranda-Saavedra D, Barton GJ, Westrop GD, Muller S, Dessi D, Fiori PL, Ren Q, Paulsen I, Zhang H, Bastida-Corcuera FD, Simoes-Barbosa A, Brown MT, Hayes RD, Mukherjee M, Okumura CY, Schneider R, Smith AJ, Vanacova S, Villalvazo M, Haas BJ, Perteau M, Feldblyum TV, Utterback TR, Shu CL, Osoegawa K, de Jong PJ, Hrđy J, Horvathova L, Zubacova Z, Dolezal P, Malik SB, Logsdon JM Jr, Henze K, Gupta A, Wang CC, Dunne RL, Upcroft JA, Upcroft P, White O, Salzberg SL, Tang P, Chiu CH, Lee YS, Embley TM, Coombs GH, Mottram JC, Tachezy J, Fraser-Liggett CM, Johnson PJ. Draft genome sequence of the sexually transmitted pathogen *Trichomonas vaginalis*. *Science* 2007;315:207–212.
  60. King N, Carroll SB. A receptor tyrosine kinase from choanoflagellates: molecular insights into early animal evolution. *Proc Natl Acad Sci USA* 2001;98:15032–15037.
  61. King N, Hittinger CT, Carroll SB. Evolution of key cell signaling and adhesion protein families predates animal origins. *Science* 2003;301:361–363.

62. Williams JG, Zvelebil M. SH2 domains in plants imply new signalling scenarios. *Trends Plant Sci* 2004;9:161–163.
63. Schlessinger J, Lemmon MA. SH2 and PTB domains in tyrosine kinase signaling. *Sci STKE* 2003;2003:RE12.
64. Benes CH, Wu N, Elia AE, Dharia T, Cantley LC, Soltoff SP. The C2 domain of PKC $\delta$  is a phosphotyrosine binding domain. *Cell* 2005;121:271–280.
65. Clamp M, Cuff J, Searle SM, Barton GJ. The Jalview Java alignment editor. *Bioinformatics* 2004;20:426–427.
66. Abraham RT. The ATM-related kinase, hSMG-1, bridges genome and RNA surveillance pathways. *DNA Repair (Amst)* 2004;3:919–925.
67. Abraham RT, Tibbetts RS. Cell biology. Guiding ATM to broken DNA. *Science* 2005;308:510–511.
68. Krungkrai J. The multiple roles of the mitochondrion of the malarial parasite. *Parasitology* 2004;129(Part 5):511–524.
69. Holness MJ, Sugden MC. Regulation of pyruvate dehydrogenase complex activity by reversible phosphorylation. *Biochem Soc Trans* 2003;31(Part 6):1143–1151.
70. Drennan D, Ryazanov AG.  $\alpha$ -Kinases: analysis of the family and comparison with conventional protein kinases. *Prog Biophys Mol Biol* 2004;85:1–32.
71. Project IRGS. The map-based sequence of the rice genome. *Nature* 2005;436:793–800.