

# Analysis of the genome sequence of the flowering plant *Arabidopsis thaliana*

## The Arabidopsis Genome Initiative\*

\* Authorship of this paper should be cited as 'The Arabidopsis Genome Initiative'. A full list of contributors appears at the end of this paper

**The flowering plant *Arabidopsis thaliana* is an important model system for identifying genes and determining their functions. Here we report the analysis of the genomic sequence of *Arabidopsis*. The sequenced regions cover 115.4 megabases of the 125-megabase genome and extend into centromeric regions. The evolution of *Arabidopsis* involved a whole-genome duplication, followed by subsequent gene loss and extensive local gene duplications, giving rise to a dynamic genome enriched by lateral gene transfer from a cyanobacterial-like ancestor of the plastid. The genome contains 25,498 genes encoding proteins from 11,000 families, similar to the functional diversity of *Drosophila* and *Caenorhabditis elegans*— the other sequenced multicellular eukaryotes. *Arabidopsis* has many families of new proteins but also lacks several common protein families, indicating that the sets of common proteins have undergone differential expansion and contraction in the three multicellular eukaryotes. This is the first complete genome sequence of a plant and provides the foundations for more comprehensive comparison of conserved processes in all eukaryotes, identifying a wide range of plant-specific gene functions and establishing rapid systematic ways to identify genes for crop improvement.**

The plant and animal kingdoms evolved independently from unicellular eukaryotes and represent highly contrasting life forms. The genome sequences of *C. elegans*<sup>1</sup> and *Drosophila*<sup>2</sup> reveal that metazoans share a great deal of genetic information required for developmental and physiological processes, but these genome sequences represent a limited survey of multicellular organisms. Flowering plants have unique organizational and physiological properties in addition to ancestral features conserved between plants and animals. The genome sequence of a plant provides a means for understanding the genetic basis of differences between plants and other eukaryotes, and provides the foundation for detailed functional characterization of plant genes.

*Arabidopsis thaliana* has many advantages for genome analysis, including a short generation time, small size, large number of offspring, and a relatively small nuclear genome. These advantages promoted the growth of a scientific community that has investigated the biological processes of *Arabidopsis* and has characterized many genes<sup>3</sup>. To support these activities, an international collaboration (the Arabidopsis Genome Initiative, AGI) began sequencing the genome in 1996. The sequences of chromosomes 2 and 4 have been reported<sup>4,5</sup>, and the accompanying Letters describe the sequences of chromosomes 1 (ref. 6), 3 (ref. 7) and 5 (ref. 8).

Here we report analysis of the completed *Arabidopsis* genome sequence, including annotation of predicted genes and assignment of functional categories. We also describe chromosome dynamics and architecture, the distribution of transposable elements and other repeats, the extent of lateral gene transfer from organelles, and the comparison of the genome sequence and structure to that of other *Arabidopsis* accessions (distinctive lines maintained by single-seed descent) and plant species. This report is the summation of work by experts interested in many biological processes selected to illuminate plant-specific functions including defence, photomorphogenesis, gene regulation, development, metabolism, transport and DNA repair.

The identification of many new members of receptor families, cellular components for plant-specific functions, genes of bacterial origin whose functions are now integrated with typical eukaryotic components, independent evolution of several families of transcription factors, and suggestions of as yet uncharacterized metabolic pathways are a few more highlights of this work. The implications of these discoveries are not only relevant for plant

biologists, but will also affect agricultural science, evolutionary biology, bioinformatics, combinatorial chemistry, functional and comparative genomics, and molecular medicine.

## Overview of sequencing strategy

We used large-insert bacterial artificial chromosome (BAC), phage (P1) and transformation-competent artificial chromosome (TAC) libraries<sup>9–12</sup> as the primary substrates for sequencing. Early stages of genome sequencing used 79 cosmid clones. Physical maps of the genome of accession Columbia were assembled by restriction fragment 'fingerprint' analysis of BAC clones<sup>13</sup>, by hybridization<sup>14</sup> or polymerase chain reaction (PCR)<sup>15</sup> of sequence-tagged sites and by hybridization and Southern blotting<sup>16</sup>. The resulting maps were integrated (<http://nucleus/cshl.org/arabmaps/>) with the genetic map and provided a foundation for assembling sets of contigs into sequence-ready tiling paths. End sequence ([http://www.tigr.org/tdb/at/abe/bac\\_end\\_search.html](http://www.tigr.org/tdb/at/abe/bac_end_search.html)) of 47,788 BAC clones was used to extend contigs from BACS anchored by marker content and to integrate contigs.

Ten contigs representing the chromosome arms and centromeric heterochromatin were assembled from 1,569 BAC, TAC, cosmid and P1 clones (average insert size 100 kilobases (kb)). Twenty-two PCR products were amplified directly from genomic DNA and sequenced to link regions not covered by cloned DNA or to optimize the minimal tiling path. Telomere sequence was obtained from specific yeast artificial chromosome (YAC) and phage clones, and from inverse polymerase chain reaction (IPCR) products derived from genomic DNA. Clone fingerprints, together with BAC end sequences, were generally adequate for selection of clones for sequencing over most of the genome. In the centromeric regions, these physical mapping methods were supplemented with genetic mapping to identify contig positions and orientation<sup>17</sup>.

Selected clones were sequenced on both strands and assembled using standard techniques. Comparison of independently derived sequence of overlapping regions and independent reassembly sequenced clones revealed accuracy rates between 99.99 and 99.999%. Over half of the sequence differences were between genomic and BAC clone sequence. All available sequenced genetic markers were integrated into sequence assemblies to verify sequence contigs<sup>4–8</sup>. The total length of sequenced regions, which extend from either the telomeres or ribosomal DNA repeats to the 180-base-pair

(bp) centromeric repeats, is 115,409,949 bp (Table 1). Estimates of the unsequenced centromeric and rDNA repeat regions measure roughly 10 megabases (Mb), yielding a genome size of about 125 Mb, in the range of the 50–150 Mb haploid content estimated by different methods<sup>18</sup>. In general, features such as gene density, expression levels and repeat distribution are very consistent across the five chromosomes (Fig. 1), and these are described in detail in reports on individual chromosomes<sup>4–8</sup> and in the analysis of centromere, telomere and rDNA sequences.

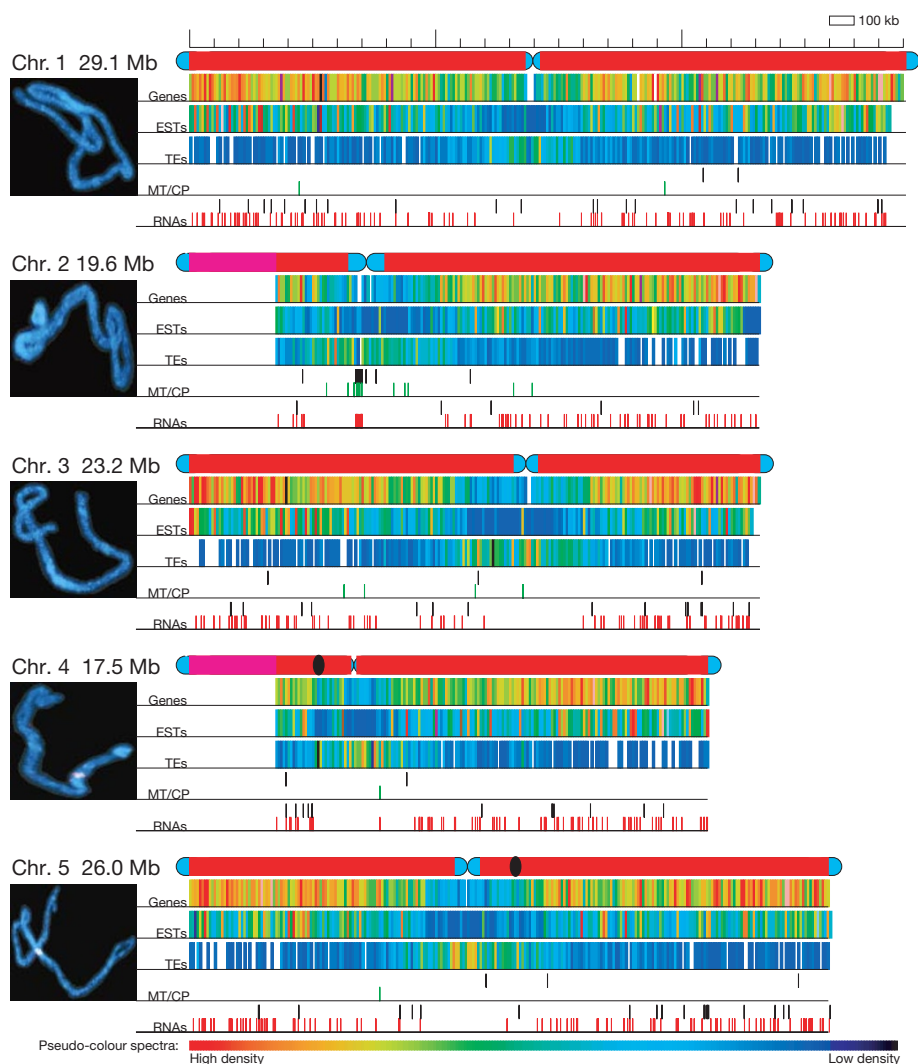
We used tRNAscan-SE 1.21 (ref. 19) and manual inspection to identify 589 cytoplasmic transfer RNAs, 27 organelle-derived tRNAs and 13 pseudogenes—more than in any other genome sequenced to date. All 46 tRNA families needed to decode all possible 61 codons were found, defining the completeness of the functional set. Several highly amplified families of tRNAs were found on the same strand<sup>6</sup>; excluding these, each amino acid is decoded by 10–41 tRNAs.

The spliceosomal RNAs (U1, U2, U4, U5, U6) have all been experimentally identified in *Arabidopsis*. The previously identified

sequences for all RNAs were found in the genome, except for U5 where the most similar counterpart was 92% identical. Between 10 and 16 copies of each small nuclear RNA (snRNA) were found across all chromosomes, dispersed as singletons or in small groups.

The small nucleolar RNAs (snoRNAs) consist of two subfamilies, the C/D box snoRNAs, which includes 36 *Arabidopsis* genes, and the H/ACA box snoRNAs, for which no members have been identified in *Arabidopsis*. U3 is the most numerous of the C/D box snoRNAs, with eight copies found in the genome. We identified forty-five additional C/D box snoRNAs using software (www.rna.wustl.edu/snoRNAdb/) that detects snoRNAs that guide ribose methylation of ribosomal RNA.

A combination of algorithms, all optimized with parameters based on known *Arabidopsis* gene structures, was used to define gene structure. We used similarities to known protein and expressed sequence tag (EST) sequence to refine gene models. Eighty per cent of the gene structures predicted by the three centres involved were completely consistent, 93% of ESTs matched gene models, and less than 1% of ESTs matched predicted non-coding regions, indicating



**Figure 1** Representation of the *Arabidopsis* chromosomes. Each chromosome is represented as a coloured bar. Sequenced portions are red, telomeric and centromeric regions are light blue, heterochromatic knobs are shown black and the rDNA repeat regions are magenta. The unsequenced telomeres 2N and 4N are depicted with dashed lines. Telomeres are not drawn to scale. Images of DAPI-stained chromosomes were kindly supplied by P. Fransz. The frequency of features was given pseudo-colour assignments, from red (high density) to deep blue (low density). Gene density ('Genes')

ranged from 38 per 100 kb to 1 gene per 100 kb; expressed sequence tag matches ('ESTs') ranged from more than 200 per 100 kb to 1 per 100 kb. Transposable element densities ('TEs') ranged from 33 per 100 kb to 1 per 100 kb. Mitochondrial and chloroplast insertions ('MT/CP') were assigned black and green tick marks, respectively. Transfer RNAs and small nucleolar RNAs ('RNAs') were assigned black and red ticks marks, respectively.

that most potential genes were identified. The sensitivity and selectivity of the gene prediction software used in this report has been comprehensively and independently assessed<sup>20</sup>.

The 25,498 genes predicted (Table 1) is the largest gene set published to date: *C. elegans*<sup>1</sup> has 19,099 genes and *Drosophila*<sup>2</sup> 13,601 genes. *Arabidopsis* and *C. elegans* have similar gene density, whereas *Drosophila* has a lower gene density; *Arabidopsis* also has a significantly greater extent of tandem gene duplications and segmental duplications, which may account for its larger gene set.

The rDNA repeat regions on chromosomes 2 and 4 were not sequenced because of their known repetitive structure and content. The centromeric regions are not completely sequenced owing to large blocks of monotonic repeats such as 5S rDNA and 180-bp repeats. The sequence continues to be extended further into centromeric and other regions of complex sequence.

### Characterization of the coding regions

To assess the similarities and differences of the *Arabidopsis* gene complement compared with other sequenced eukaryotic genomes, we assigned functional categories to the complete set of *Arabidopsis* genes. For chromosome 4 genes and the yeast genome, predicted functions were previously manually assigned<sup>5,21</sup>. All other predicted proteins were automatically assigned to these functional categories<sup>22</sup>, assuming that conserved sequences reflect common functional relationships.

The functions of 69% of the genes were classified according to sequence similarity to proteins of known function in all organisms; only 9% of the genes have been characterized experimentally (Fig. 2a). Generally similar proportions of gene products were predicted to be targeted to the secretory pathway and mitochondria in *Arabidopsis* and yeast, and up to 14% of the gene products are

**Table 1 Summary statistics of the *Arabidopsis* genome**

Feature	Value					
<b>(a) The DNA molecules</b>						
	Chr. 1	Chr. 2	Chr. 3	Chr. 4	Chr. 5	Σ
Length (bp)	29,105,111	19,646,945	23,172,617	17,549,867	25,953,409	115,409,949
Top arm (bp)	14,449,213	3,607,091	13,590,268	3,052,108	11,132,192	
Bottom arm (bp)	14,655,898	16,039,854	9,582,349	14,497,759	14,803,217	
Base composition (%GC)						
Overall	33.4	35.5	35.4	35.5	34.5	
Coding	44.0	44.0	44.3	44.1	44.1	
Non-coding	32.4	32.9	33.0	32.8	32.5	
Number of genes	6,543	4,036	5,220	3,825	5,874	25,498
Gene density (kb per gene)	4.0	4.9	4.5	4.6	4.4	
Average gene length (bp)	2,078	1,949	1,925	2,138	1,974	
Average peptide length (bp)	446	421	424	448	429	
Exons						
Number	35,482	19,631	26,570	20,073	31,226	13,2982
Total length (bp)	8,772,559	5,100,288	6,654,507	5,150,883	7,571,013	33,249,250
Average per gene	5.4	4.9	5.1	5.2	5.3	
Average size (bp)	247	259	250	256	242	
Introns						
Number	28,939	15,595	21,350	16,248	25,352	107,484
Total length (bp)	4,828,766	2,768,430	3,397,531	3,030,649	4,030,045	18,055,421
Average size (bp)	168	177	159	189	159	
Number of genes with ESTs (%)	60.8	56.9	59.8	61.4	61.4	
Number of ESTs	30,522	14,989	20,732	16,605	22,885	105,733
<b>(b) The proteome</b>						
Classification/function						
Total proteins	6,543	4,036	5,220	3,825	5,874	25,498
With INTERPRO domains	4,194	1,205	2,989	1,545	3,136	13,069
Genes containing at least one TM domain	64.1%	29.9%	57.8%	40.4%	53.4%	51.3%
Genes containing at least one SCOP domain	2,334	1,322	1,615	1,402	1,940	8,613
Genes containing at least one SCOP domain	35.7%	32.8%	30.9%	36.7%	33.0%	33.8%
Genes containing at least one SCOP domain	2,513	1,424	1,664	1,304	2,121	9,026
Genes containing at least one SCOP domain	38.4%	35.3%	31.9%	34.1%	36.1%	35.4%
With putative signal peptides						
Secretory pathway	1,242	675	877	659	1,014	4,467
>0.95 specificity	19.0%	16.7%	17.0%	17.2%	17.3%	17.6%
>0.95 specificity	1,146	632	813	632	964	4,167
>0.95 specificity	17.5%	15.7%	15.7%	16.5%	16.4%	16.4%
Chloroplast	866	535	754	532	887	3,574
>0.95 specificity	13.2%	13.2%	14.6%	13.9%	15.1%	14.0%
>0.95 specificity	602	290	420	298	475	2,085
>0.95 specificity	9.2%	7.2%	8.1%	7.8%	8.1%	8.2%
mitochondria	901	425	554	390	627	2,897
>0.95 specificity	13.8%	10.5%	10.7%	10.2%	10.7%	11.4%
>0.95 specificity	113	49	63	59	65	349
>0.95 specificity	1.7%	1.2%	1.2%	1.5%	1.1%	1.4%
Functional classification						
Cellular metabolism	1,188	620	745	588	868	4,009
Cellular metabolism	22.7%	23.3%	22.8%	22.9%	21.1%	22.5%
Transcription	880	474	566	335	763	3,018
Transcription	16.8%	17.8%	17.3%	13.1%	18.6%	16.9%
Plant defence	640	276	354	295	490	2,055
Plant defence	12.2%	10.4%	10.8%	11.5%	11.9%	11.5%
Signalling	573	296	356	210	420	1,855
Signalling	11.0%	11.1%	10.9%	8.2%	10.2%	10.4%
Growth	542	263	357	448	469	2,079
Growth	10.4%	9.9%	10.9%	17.5%	11.4%	11.7%
Protein fate	520	273	314	264	395	1,766
Protein fate	9.9%	10.2%	9.6%	10.3%	9.6%	9.9%
Intracellular transport	435	214	269	220	334	1,472
Intracellular transport	8.3%	8.9%	8.2%	8.6%	8.1%	8.3%
Transport	236	139	155	113	206	849
Transport	4.5%	5.2%	4.7%	4.4%	5.0%	4.8%
Protein synthesis	216	111	148	90	165	730
Protein synthesis	4.1%	4.2%	4.5%	3.5%	4.0%	4.1%
Total	5,230	2,666	3,264	2,563	4,110	17,833

The features of *Arabidopsis* chromosomes 1–5 and the complete nuclear genome are listed. Specialized searches used the following programs and databases: INTERPRO<sup>23</sup>; transmembrane (TM) domains by ALOM2 (unpublished); SCOP domain database<sup>21</sup>; functional classification by the PEDANT analysis system<sup>22</sup>. Signal peptide prediction (secretory pathway, targeted to chloroplast or mitochondria) was performed using TargetP<sup>22</sup> and <http://www.cbs.dtu.dk/services/TargetP/>.

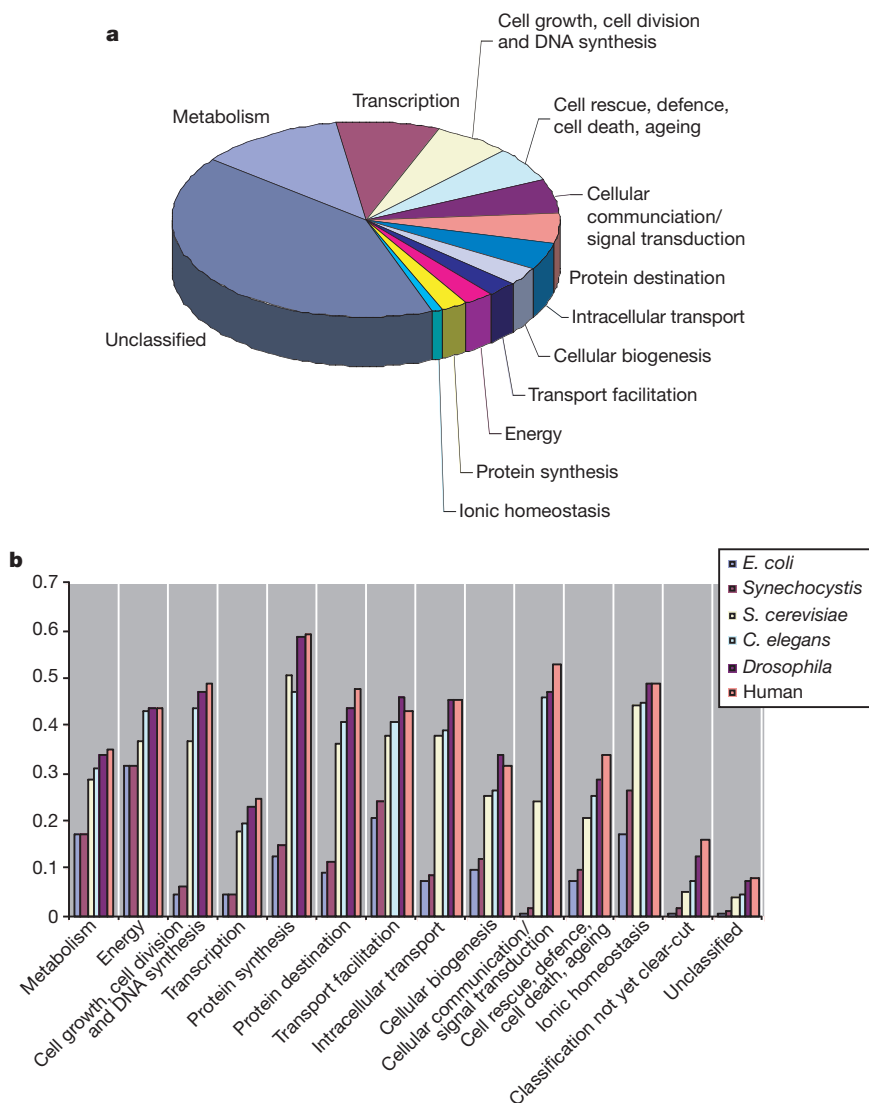
\* Default value.

likely to be targeted to the chloroplast (Table 1). The significant proportion of genes with predicted functions involved in metabolism, gene regulation and defence is consistent with previous analyses<sup>5</sup>. Roughly 30% of the 25,498 predicted gene products, (Fig. 2a), comprising both plant-specific proteins and proteins with similarity to genes of unknown function from other organisms, could not be assigned to functional categories.

To compare the functional categories in more detail, we compared data from the complete genomes of *Escherichia coli*<sup>23</sup>, *Synechocystis* sp.<sup>24</sup>, *Saccharomyces cerevisiae*<sup>21</sup>, *C. elegans*<sup>1</sup> and *Drosophila*<sup>2</sup>, and a non-redundant protein set of *Homo sapiens*, with the *Arabidopsis* genome data (Fig. 2b), using a stringent BLASTP threshold value of  $E < 10^{-30}$ . The proportion of *Arabidopsis* proteins having related counterparts in eukaryotic genomes varies by a factor of 2 to 3 depending on the functional category. Only 8–23% of *Arabidopsis* proteins involved in transcription have related genes in other eukaryotic genomes, reflecting the independent evolution of many plant transcription factors. In contrast, 48–60% of genes involved in protein synthesis have counterparts in the other eukaryotic genomes, reflecting highly

conserved gene functions. The relatively high proportion of matches between *Arabidopsis* and bacterial proteins in the categories ‘metabolism’ and ‘energy’ reflects both the acquisition of bacterial genes from the ancestor of the plastid and high conservation of sequences across all species. Finally, a comparison between unicellular and multicellular eukaryotes indicates that *Arabidopsis* genes involved in cellular communication and signal transduction have more counterparts in multicellular eukaryotes than in yeast, reflecting the need for sets of genes for communication in multicellular organisms.

Pronounced redundancy in the *Arabidopsis* genome is evident in segmental duplications and tandem arrays, and many other genes with high levels of sequence conservation are also scattered over the genome. Sequence similarity exceeding a BLASTP value  $E < 10^{-20}$  and extending over at least 80% of the protein length were used as parameters to identify protein families (Table 2). A total of 11,601 protein types were identified. Thirty-five per cent of the predicted proteins are unique in the genome, and the proportion of proteins belonging to families of more than five members is substantially higher in *Arabidopsis* (37.4%) than in *Drosophila* (12.1%) or



**Figure 2** Functional analysis of *Arabidopsis* genes. **a**, Proportion of predicted *Arabidopsis* genes in different functional categories. **b**, Comparison of functional categories between organisms. Subsets of the *Arabidopsis* proteome containing all proteins that fall into a common functional class were assembled. Each subset was searched against the complete set of translations from *Escherichia coli*, *Synechocystis* sp. PCC6803,

*Saccharomyces cerevisiae*, *Drosophila*, *C. elegans* and a *Homo sapiens* non-redundant protein database. The percentage of *Arabidopsis* proteins in a particular subset that had a BLASTP match with  $E \leq 10^{-30}$  to the respective reference genome is shown. This reflects the measure of sequence conservation of proteins within this particular functional category between *Arabidopsis* and the respective reference genome. y axis, 0.1 = 10%.

**Table 2 Proportion of genes in different organisms present as either singletons or in paralogous families**

	No of singletons and distinct gene families	Unique	Gene families containing				
			2 members	3 members	4 members	5 members	>5 members
<i>H. influenzae</i>	1,587	88.8%	6.8%	2.3%	0.7%	0.0%	1.4%
<i>S. cerevisiae</i>	5,105	71.4%	13.8%	3.5%	2.2%	0.7%	8.4%
<i>D. melanogaster</i>	10,736	72.5%	8.5%	3.4%	1.9%	1.6%	12.1%
<i>C. elegans</i>	14,177	55.2%	12.0%	4.5%	2.7%	1.6%	24.0%
<i>Arabidopsis</i>	11,601	35.0%	12.5%	7.0%	4.4%	3.6%	37.4%

The number of genes in the genomes of *Haemophilus influenzae*, *S. cerevisiae*, *Drosophila*, *C. elegans* and *Arabidopsis* that are present either as singletons or in gene families with two or more members are listed. To be grouped in a gene family, two genes had to show similarity exceeding a BLASTP value  $E < 10^{-20}$  and a FASTA alignment over at least 80% of the protein length. In column 1, the number of genes that are unique plus the number of gene families are listed. Columns 2 to 6 give the percentage of genes present as singletons or in gene families of  $n$  members.

*C. elegans* (24.0%). The absolute number of *Arabidopsis* gene families and singletons (types) is in the same range as the other multicellular eukaryotes, indicating that a proteome of 11,000–15,000 types is sufficient for a wide diversity of multicellular life. The proportion of gene families with more than two members is considerably more pronounced in *Arabidopsis* than in other eukaryotes (Fig. 3). As segmental duplication is responsible for 6,303 gene duplications (see below), the extent of tandem gene duplications accounts for a significant proportion of the increased family size. These features of the *Arabidopsis*, and presumably other plant genomes, may indicate more relaxed constraints on genome size in plants, or a more prominent role of unequal crossing over to generate new gene copies.

Conserved protein domains revealed more informative differences through INTERPRO<sup>25</sup> analysis of the predicted gene products from *Arabidopsis*, *S. cerevisiae*, *C. elegans* and *Drosophila*. Statistically over-represented domains, and those that are absent from the *Arabidopsis* genome, indicate domains that may have been gained or lost during the evolution of plants (Supplementary Information Table 1). Proteins containing the Pro-Pro-Arg repeat, which is involved in RNA stabilization and RNA processing, are over-represented as compared to yeast, fly and worm; 400 proteins containing this signature were detected in *Arabidopsis* compared with only 10 in total in yeast, *Drosophila* and *C. elegans*. Protein kinases and associated domains, 169 proteins containing a disease resistance protein signature, and the Toll/IL-1R (TIR) domain, a component of pathogen recognition molecules<sup>26</sup>, are also relatively abundant. This suggests that pathways transducing signals in response to pathogens and diverse environmental cues are more abundant in plants than in other organisms.

The RING zinc finger domain is relatively over-represented in *Arabidopsis* compared with yeast, *Drosophila* and *C. elegans*, whereas the F-box domain is over-represented as compared with yeast and *Drosophila* only. These domains are involved in targeting proteins to the proteasome<sup>27</sup> and ubiquitylation<sup>28</sup> pathways of protein degradation, respectively. In plants many processes such as hormone and defence responses, light signalling, and circadian rhythms and pattern formation use F-box function to direct negative regulators

to the ubiquitin degradation pathway. This mode of regulation appears to be more prevalent in plants and may account for a higher representation of the F box than in *Drosophila* and for the over-representation of the ubiquitin domain in the *Arabidopsis* genome. RING finger domain proteins in general have a role in ubiquitin protein ligases, indicating that proteasome-mediated degradation is a more widespread mode of regulation in plants than in other kingdoms.

Most functions identified by protein domains are conserved in similar proportions in the *Arabidopsis*, *S. cerevisiae*, *Drosophila* and *C. elegans* genomes, pointing to many ubiquitous eukaryotic pathways. These are illustrated by comparing the list of human disease genes<sup>29</sup> to the complete *Arabidopsis* gene set using BLASTP. Out of 289 human disease genes, 139 (48%) had hits in *Arabidopsis* using a BLASTP threshold  $E < 10^{-10}$ . Sixty-nine (24%) exceeded an  $E < 10^{-40}$  threshold, and 26 (9.3%) had scores better than  $E < 10^{-100}$  (Table 3). There are at least 17 human disease genes more similar to *Arabidopsis* genes than yeast, *Drosophila* or *C. elegans* genes (Table 3).

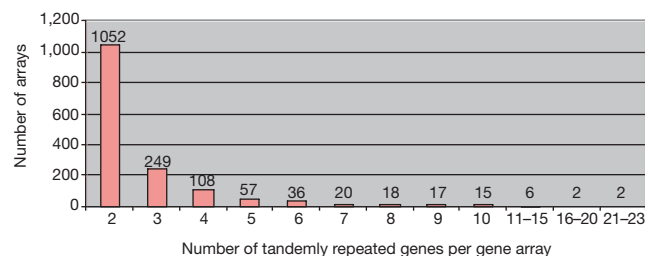
This analysis shows that, although numerous families of proteins are shared between all eukaryotes, plants contain roughly 150 unique protein families. These include transcription factors, structural proteins, enzymes and proteins of unknown function. Members of the families of genes common to all eukaryotes have undergone substantial increases or decreases in their size in *Arabidopsis*. Finally, the transfer of a relatively small number of cyanobacteria-related genes from a putative endosymbiotic ancestor of the plastid has added to the diversity of protein structures found in plants.

### Genome organization and duplication

The *Arabidopsis* genome sequence provides a complete view of chromosomal organization and clues to its evolutionary history. Gene families organized in tandem arrays of two or more units have been described in *C. elegans*<sup>1</sup> and *Drosophila*<sup>2</sup>. Analysis of the *Arabidopsis* genome revealed 1,528 tandem arrays containing 4,140 individual genes, with arrays ranging up to 23 adjacent members (Fig. 3). Thus 17% of all genes of *Arabidopsis* are arranged in tandem arrays.

Large segmental duplications were identified either by directly aligning chromosomal sequences or by aligning proteins and searching for tracts of conserved gene order. All five chromosomes were aligned to each other in both orientations using MUMmer<sup>30</sup>, and the results were filtered to identify all segments at least 1,000 bp in length with at least 50% identity (Supplementary Information Fig. 1). These revealed 24 large duplicated segments of 100 kb or larger, comprising 65.6 Mb or 58% of the genome. The only duplicated segment in the centromeric regions was a 375-kb segment on chromosome 4. Many duplications appear to have undergone further shuffling, such as local inversions after the duplication event.

We used TBLASTX<sup>5</sup> to identify collinear clusters of genes residing in large duplicated chromosomal segments. The duplicated regions encompass 67.9 Mb, 60% of the genome, slightly more than was



**Figure 3** Distribution of tandemly repeated gene arrays in the *Arabidopsis* genome. Tandemly repeated gene arrays were identified using the BLASTP program with a threshold of  $E < 10^{-20}$ . One unrelated gene among cluster members was tolerated. The histogram gives the number of clusters in the genome containing 2 to  $n$  similar gene units in tandem.

found in the DNA-based alignment (Fig. 4), and these data extend earlier findings<sup>4,5,31</sup>. The extent of sequence conservation of the duplicated genes varies greatly, with 6,303 (37%) of the 17,193 genes in the segments classified as highly conserved ( $E < 10^{-30}$ ) and a further 1,705 (10%) showing less significant similarity up to  $E < 10^{-5}$ . The proportion of homologous genes in each duplicated segment also varies widely, between 20% and 47% for the highly conserved class of genes. In many cases, the number of copies of a gene and its counterpart differ (for example, one copy on one chromosome and multiple copies on the other; see Supplementary Information Fig. 2); this could be due to either tandem duplication or gene loss after the segmental duplication.

What does the duplication in the *Arabidopsis* genome tell us about the ancestry of the species? Polyploidy occurs widely in plants and is proposed to be a key factor in plant evolution<sup>32</sup>. As the majority of the *Arabidopsis* genome is represented in duplicated (but not triplicated) segments, it appears most likely that *Arabidopsis*, like maize, had a tetraploid ancestor<sup>33</sup>. A comparative sequence analysis of *Arabidopsis* and tomato estimated that a duplication occurred ~112 Myr ago to form a tetraploid<sup>34</sup>. The degrees of conservation of the duplicated segments might be due to divergence from an ancestral autotetraploid form, or might reflect differences present in an allotetraploid ancestor. It is also possible, however, that several independent segmental duplication events took place instead of tetraploid formation and stabilization.

The diploid genetics of *Arabidopsis* and the extensive divergence of the duplicated segments have masked its evolutionary history. The determination of *Arabidopsis* gene functions must therefore be pursued with the potential for functional redundancy taken into account. The long period of time over which genome stabilization has occurred has, however, provided ample opportunity for the divergence of the functions of genes that arose from duplications.

### Comparative analysis of *Arabidopsis* accessions

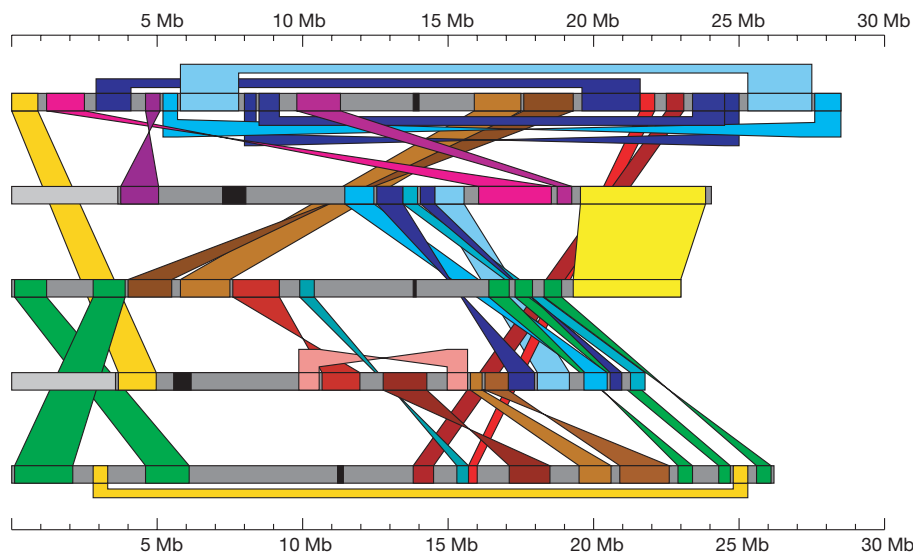
Comparing the multiple accessions of *Arabidopsis* allows us to identify commonly occurring changes in genome microstructure. It also enables the development of new molecular markers for genetic mapping. High rates of polymorphism between *Arabidopsis* accessions, including both DNA sequence and copy number of tandem arrays, are prevalent at loci involved in disease resistance<sup>35</sup>. This has been observed for other plant species, and such loci are thought to serve as templates for illegitimate recombination

to create new pathogen response specificities<sup>36</sup>. We carried out a comparative analysis between 82 Mb of the genome sequence of *Arabidopsis* accession Columbia (Col-0) and 92.1 Mb of non-redundant low-pass (twofold redundant) sequence data of the genomic DNA of accession Landsberg *erecta* (Ler). We identified two classes of differences between the sequences: single nucleotide polymorphisms (SNPs), and insertion-deletions (InDels). As we used high stringency criteria, our results represent a minimum estimate of numbers of polymorphisms between the two genomes.

In total, we detected 25,274 SNPs, representing an average density of 1 SNP per 3.3 kb. Transitions (A/T–G/C) represented 52.1% of the SNPs, and transversions accounted for the remainder: 17.3% for A/T–T/A, 22.7% for A/T–C/G and 7.9% for C/G–G/C. In total, we detected 14,570 InDels at an average spacing of 6.1 kb. They ranged from 2 bp to over 38 kilobase-pairs, although 95% were smaller than 50 bp. Only 10% of the InDels were co-located with simple sequence repeats identified with the program Sputnik. An analysis of 416 relative insertions greater than 250 bp in Col-0 showed that 30% matched transposon-related proteins, indicating that a substantial proportion of the large InDels are the result of transposon insertion or excision. Many InDels contained entire active genes not related to transposons. Half of such genes absent from corresponding positions in the Col-0 sequence were found elsewhere on the genome of *Ler*. This indicates that genes have been transferred to new genomic locations.

Gene structures are often affected by small InDels and SNPs. The positions of SNPs and InDels were mapped relative to 87,427 exons and 70,379 introns annotated in the Col-0 sequence. SNPs were found in exons, introns and intergenic regions at frequencies of 1 SNP per 3.1, 2.2 and 3.5 kb, respectively. The frequencies for InDels were 1 per 9.3, 3.1 and 4.3 kb, respectively. Polymorphisms were detected in 7% of exons, and alter the spliced sequences of 25% of the predicted genes. For InDels in exons, insertion lengths divisible by three are prevalent for small insertions (< 50 bp), indicating that many proteins can withstand small insertions or deletions of amino acids without loss of function.

Our analyses show that sequence polymorphisms between accessions of *Arabidopsis* are common, and that they occur in both coding and non-coding regions. We found evidence for the relocation of genes in the genome, and for changes in the complement of transposable elements. The data presented here are available at <http://www.arabidopsis.org/cereon/>.



**Figure 4** Segmentally duplicated regions in the *Arabidopsis* genome. Individual chromosomes are depicted as horizontal grey bars (with chromosome 1 at the top), centromeres are marked black. Coloured bands connect corresponding duplicated

segments. Similarity between the rDNA repeats are excluded. Duplicated segments in reversed orientation are connected with twisted coloured bands. The scale is in megabases.

**Comparison of *Arabidopsis* and other plant genera**

Comparative genetic mapping can reveal extensive conservation of genome organization between closely related species<sup>37,38</sup>. The comparative analysis of plant genome microstructure reveals much about the evolution of plant genomes and provides unprecedented opportunities for crop improvement by establishing the detailed structures of, and relationships between, the genomes of crops and *Arabidopsis*.

The lineages leading to *Arabidopsis* and *Capsella rubella* (shepherd's purse) diverged between 6.2 and 9.8 Myr ago, and the gene content and genome organization of *C. rubella* is very similar to that of *Arabidopsis*<sup>39</sup>, including the large-scale duplications. Alignment of *Arabidopsis* complementary DNA and EST sequences with genomic DNA sequences of *Arabidopsis* and *C. rubella* showed conservation of exon length and intron positions. Coding sequences predicted from these alignments differed from the annotated *Arabidopsis* gene sequences in two out of five cases.

The ancestral lineages of *Arabidopsis* and the *Brassica* (cabbage and mustard) genera diverged 12.2–19.2 Myr ago<sup>40</sup>. *Brassica* genes show a high level of nucleotide conservation with their *Arabidopsis* orthologues, typically more than 85% in coding regions<sup>40</sup>. The structure of *Brassica* genomes resembles that of *Arabidopsis*, but with extensive triplication and rearrangement<sup>41</sup>, and extensive divergence of microstructure (Supplementary Information Fig. 3). The divergence between the genomes of *Arabidopsis* and *Brassica oleracea* is in striking contrast to that observed between *Arabidopsis* and *C. rubella*, although the time since divergence is only twofold greater. This accelerated rate of change in triplicated segments of the genome of *B. oleracea* indicates that polyploidy fosters rapid chromosomal evolution.

The *Arabidopsis* and tomato lineages diverged roughly 150 Myr ago, and comparative sequence analysis of segments of their genomes has revealed complex relationships<sup>34</sup>. Four regions of the *Arabidopsis* genome are related to each other and to one region in the tomato genome, suggesting that two rounds of duplication may

have occurred in the *Arabidopsis* lineage. The extensive duplication described here supports the proposal that the more recent of these duplications, estimated to have occurred ~112 Myr ago, was the result of a polyploidization event. The lineages of *Arabidopsis* and rice diverged ~200 Myr ago<sup>42</sup>. Three regions of the genome of *Arabidopsis* were related to each other and to one region in the rice genome, providing further evidence for multiple duplication events<sup>43,44</sup>.

The frequent occurrence of tandem gene duplications and the apparent deletion of single genes, or small groups of adjacent genes, from duplicated regions suggests that unequal crossing over may be a key mechanism affecting the evolution of plant genome microstructure. However, the segmental inversions and gene translocations in the genomes of both rice and *B. oleracea* that are not found in *Arabidopsis* indicate that additional mechanisms may be involved<sup>40</sup>.

**Integration of the three genomes in the plant cell**

The three genomes in the plant cell—those of the nucleus, the plastids (chloroplasts) and the mitochondria—differ markedly in gene number, organization and stability. Plastid genes are densely packed in an order highly conserved in all plants<sup>45</sup>, whereas mitochondrial genes<sup>46</sup> are widely dispersed and subjected to extensive recombination.

Organelar genomes are remnants of independent organisms—plastids are derived from the cyanobacterial lineage and mitochondria from the  $\alpha$ -Proteobacteria. The remaining genes in plastids include those that encode subunits of the photosystem and the electron transport chain, whereas the genes in mitochondria encode essential subunits of the respiratory chain. Both organelles contain sets of specific membrane proteins that, together with housekeeping proteins, account for 61% of the genes in the chloroplast and 88 % in the mitochondrion (Table 4). The balances are involved in transcription and translation.

The number of proteins encoded in the nucleus likely to be found

**Table 3 *Arabidopsis* genes with similarities to human disease genes**

Human disease gene	E value	Gene code	<i>Arabidopsis</i> hit
Darier–White, SERCA	$5.9 \times 10^{-272}$	T2711_16	Putative calcium ATPase
Xeroderma Pigmentosum, D-XPD	$7.2 \times 10^{-228}$	F15K9_19	Putative DNA repair protein
Xeroderma pigment, B-ERCC3	$9.6 \times 10^{-214}$	AT5g41360	DNA excision repair cross-complementing protein
Hyperinsulinism, ABCC8	$7.1 \times 10^{-188}$	F20D22_11	Multidrug resistance protein
Renal tubul. acidosis, ATP6B1	$1.0 \times 10^{-182}$	AT4g38510	Probable H <sup>+</sup> -transporting ATPase
HDL deficiency 1, ABCA1	$2.4 \times 10^{-181}$	At2g41700	Putative ABC transporter
Wilson, ATP7B	$7.6 \times 10^{-181}$	AT5g44790	ATP-dependent copper transporter
Immunodeficiency, DNA Ligase 1	$8.2 \times 10^{-172}$	T6D22_10	DNA ligase
Stargardt's, ABCA4	$2.8 \times 10^{-168}$	At2g41700	Putative ABC transporter
Ataxia telangiectasia, ATM	$3.1 \times 10^{-168}$	AT3g48190	Ataxia telangiectasia mutated protein AtATM
Niemann–Pick, NPC1	$1.2 \times 10^{-166}$	F7F22_1	Niemann–Pick C disease protein-like protein
Menkes, ATP7A	$1.1 \times 10^{-153}$	F2K11_17	ATP-dependent copper transporter, putative
HNPCC*, MLH1	$1.5 \times 10^{-150}$	AT4g09140	MLH1 protein
Deafness, hereditary, MYO15	$2.7 \times 10^{-150}$	At2g31900	Putative unconventional myosin
Fam, cardiac myopathy, MYH7	$6.5 \times 10^{-147}$	T1G11_14	Putative myosin heavy chain
Xeroderma Pigmentosum, F-XPF	$1.4 \times 10^{-146}$	AT5g41150	Repair endonuclease (gb)AAAF01274.1)
G6PD deficiency, G6PD	$7.6 \times 10^{-137}$	AT5g40760	Glucose-6-phosphate dehydrogenase
Cystic fibrosis, ABCC7	$2.3 \times 10^{-135}$	AT3g62700	ABC transporter-like protein
Glycerol kinase defic, GK	$7.9 \times 10^{-135}$	T21F11_21	Putative glycerol kinase
HNPCC, MSH3	$6.6 \times 10^{-134}$	AT4g25540	Putative DNA mismatch repair protein
HNPCC, PMS2	$5.1 \times 10^{-128}$	AT4g02460	No title
Zellweger, PEX1	$4.1 \times 10^{-125}$	AT5g08470	Putative protein
HNPCC, MSH6	$9.6 \times 10^{-122}$	AT4g02070	G/T DNA mismatch repair enzyme
Bloom, BLM	$4.4 \times 10^{-109}$	T19D16_15	DNA helicase isolog
Finnish amyloidosis, GSN	$2.2 \times 10^{-107}$	AT5g57320	Villin
Chediak–Higashi, CHS1	$5.8 \times 10^{-99}$	F10O3_11	Putative transport protein
Xeroderma Pigmentosum, G-XPG	$7.1 \times 10^{-89}$	AT3g28030	Hypothetical protein
Bare lymphocyte, ABCB3	$1.3 \times 10^{-84}$	AT5g39040	ABC transporter-like protein
Citrullinemia, type I, ASS	$3.2 \times 10^{-83}$	AT4g24830	Argininosuccinate synthase-like protein
Coffin–Lowry, RPS6KA3	$5.2 \times 10^{-81}$	AT3g08720	Putative ribosomal-protein S6 kinase (ATPK19)
Keratoderma, KRT9	$8.5 \times 10^{-81}$	AT3g17050	Unknown protein
Myotonic dystrophy, DM1	$1.4 \times 10^{-76}$	At2g20470	Putative protein kinase
Bartter's, SLC12A1	$1.6 \times 10^{-75}$	F26G16_9	Cation-chloride co-transporter, putative
Dents, CLCN5	$3.3 \times 10^{-74}$	AT5g26240	CLC-d chloride channel protein
Diaphanous 1, DAPH1	$1.9 \times 10^{-73}$	68069_m00158	Hypothetical protein
AKT2	$6.9 \times 10^{-72}$	AT3g08730	Putative ribosomal-protein S6 kinase (ATPK6)

in organelles was predicted using default settings on TargetP (Table 1). Many nuclear gene products that are targeted to either (or both) organelles were originally encoded in the organelle genomes and were transferred to the nuclear genome during evolutionary history. A large number also appear to be of eukaryotic origin, with functions such as protein import components, which were probably not required by the free-living ancestors of the endosymbionts.

To identify nuclear genes of possible organellar ancestry, we compared all predicted *Arabidopsis* proteins to all proteins from completed genomes including those from plastids and mitochondria (Supplementary Information Table 2). This search identified proteins encoded by the *Arabidopsis* nuclear genome that are most similar to proteins encoded by other species' organelle genomes (14 mitochondrial and 44 plastid). These represent organelle-to-nuclear gene transfers that have occurred sometime after the divergence of the organelle-containing lineages<sup>47</sup>. There is a great excess of nuclear encoded proteins most similar to proteins from the cyanobacteria *Synechocystis* (Supplementary Information Fig. 4; 806 *Arabidopsis* predicted proteins matching 404 different *Synechocystis* proteins, providing further evidence of a genome duplication). These 806 *Arabidopsis* predicted proteins, and many others of greatly diverse function, are possibly of plastid descent. Through searches against proteins from other cyanobacteria (with incompletely sequenced genomes), we identified 69 additional genes of possibly plastid descent. Only 25% of these putatively plastid-derived proteins displayed a target peptide predicted by TargetP, indicating potential cytoplasmic functions for most of these genes.

The difference between predicted plastid-targeted and predicted plastid-derived genes indicates that there is a probable overestimation by *ab initio* targeting prediction methods and a lack of resolution with respect to destination organelles, the possible extensive divergence of some endosymbiont-derived genes in the nuclear genome, the co-opting of nuclear genes for targeting to organelles, and cytoplasmic functions for cyanobacteria-derived proteins. Clearly more refined tools and extensive experimentation is required to catalogue plastid proteins.

The transfer of genes between genomes still continues (Supplementary Information Table 3). Plastid DNA insertions in the nucleus (17 insertions totalling 11 kb) contain full-length genes encoding proteins or tRNAs, fragments of genes and an intron as well as intergenic regions. Subsequent reshuffling in the nucleus is illustrated by the *atpH* gene, which was originally transferred completely, but is now in two pieces separated by 2 kb. The 13 small mitochondrial DNA insertions total 7 kb in addition to the large insertion close to the centromere of chromosome 2 (ref. 3). The high level of recombination in the mitochondrial genome may account for these events.

### Transposable elements

Transposons, which were originally identified in maize by Barbara McClintock, have been found in all eukaryotes and prokaryotes. A

subset of transposons replicate through an RNA intermediate (class I), whereas others move directly through a DNA form (class II). Transposons are further classified by similarity either between their mobility genes or between their terminal and/or internal motifs, as well as by the size and sequence of their target site. Internally deleted elements can often be mobilized in *trans* by fully functional elements.

Transposons in *Arabidopsis* account for at least 10% of the genome, or about one-fifth of the intergenic DNA. The *Arabidopsis* genome has a wealth of class I (2,109) and II (2,203) elements, including several new groups (1,209 elements; Supplementary Information Table 4). Mobile histories for many elements were obtained by identifying regions of the genome with significant similarity to 'empty' target sites (RESites) thus providing high-resolution information concerning the termini and target site duplications<sup>48,49</sup>. These regions were readily detected because of the propensity of transposons to integrate into repeats and because of duplications in the genome sequence. In several cases, genes appear to have been included as 'passengers' in transposable units<sup>48</sup>. In some cases, shared sequence similarity, coding capacity and RESites attest to recent activity of transposable elements in the *Arabidopsis* genome. Only about 4% of the complete elements identified correspond to an EST, however, suggesting that most are not transcribed.

Transposable elements found in many other plant genomes are well represented in *Arabidopsis*, including *copia*- and *gypsy*-like long terminal repeat (LTR) retrotransposons, long interspersal nuclear elements (LINEs); short interspersed nuclear elements (SINEs), *hobo/Activator/Tam3 (hAT)*-like elements, CACTA-like elements and miniature inverted-repeat transposable elements (MITES). Although usually small in size, some larger *Tourist*-like MITES contain open reading frames (ORFs) with similarity to the transposases of bacterial insertion sequences<sup>48</sup>. *Basho* and many *Mutator*-like elements (MULEs), first discovered in the *Arabidopsis* sequence, represent structurally unique transposons<sup>48-50</sup>. *Basho* elements have a target site preference for mononucleotide 'A' and wide distribution among plants<sup>48,51</sup>. MULEs exhibit a high level of sequence diversity and members of most groups lack long terminal inverted repeats (TIRs). Phylogenetic analysis of the *Arabidopsis* MURA-like transposases suggests that TIR-containing MULEs are more closely related to one another than to MULEs lacking TIRs<sup>49,52</sup>.

For many plants with large genomes, class I retrotransposons contribute most of the nucleotide content<sup>53</sup>. In the small *Arabidopsis* genome, class I elements are less abundant and primarily occupy the centromere. In contrast, *Basho* elements and class II transposons such as MITES and MULEs predominate on the periphery of pericentromeric domains (Fig. 5). In class II transposons, MULEs and CACTA elements are clustered near centromeres and heterochromatic knobs, whereas MITES and *hAT* elements have a less pronounced bias. The distribution pattern of transposable elements observed in *Arabidopsis* may reflect different types of pericentromeric heterochromatin regions and may be similar to those found in animals.

Numerous centromeric satellite repeats are located between each chromosome arm and have not yet been sequenced, but are represented in part by unanchored BAC contigs (R. Martienssen and M. Marra, unpublished data). End sequence suggests that these domains contain many more class I than class II elements, consistent with the distribution reported here (K. Lemcke and R. Martienssen, unpublished data). We do not know the significance of the apparent paucity of elements in telomeric regions and in the region flanking the rDNA repeats on chromosome 4 (but not on chromosome 2).

Overall, transposon-rich regions are relatively gene-poor and have lower rates of recombination and EST matches, indicating a correlation between low gene expression, high transposon density and low recombination<sup>51</sup>. The role of transposons in genome

**Table 4** General features of genes encoded by the three genomes in *Arabidopsis*

	Nucleus/cytoplasm	Plastid	Mitochondria
Genome size	125 Mb	154 kb	367 kb
Genome equivalent/cell	2	560	26
Duplication	60%	17%	10%
Number of protein genes	25,498	79	58
Gene order	Variable, but syntenic	Conserved	Variable
Density (kb per protein gene)	4.5	1.2	6.25
Average coding length	1,900 nt	900 nt	860 nt
Genes with introns	79%	18.4%	12%
Genes/pseudogenes	1/0.03	1/0	1/0.2-0.5
Transposons (% of total genome size)	14%	0%	4%



organization and chromosome structure can now be addressed in a model organism known to undergo DNA methylation and other forms of chromatin modification thought to regulate transposition<sup>52</sup>.

**rDNA, telomeres and centromeres**

Nucleolar organizers (NORs) contain arrays of unit repeats encoding the 18S, 5.8S and 25S ribosomal RNA genes and are transcribed by RNA polymerase I. Together with 5S RNA, which is transcribed by RNA polymerase III, these rRNAs form the structural and catalytic cores of cytoplasmic ribosomes. In *Arabidopsis*, the NORs juxtapose the telomeres of chromosomes 2 and 4, and comprise uninterrupted 18S, 5.8S and 25S units all orientated on the chromosomes in the same direction<sup>54</sup>. In contrast, the 5S rRNA genes are localized to heterogeneous arrays in the centromeric regions of chromosomes 3, 4 and 5 (ref. 55; and Fig. 6). Both NORs are roughly 3.5–4.0 megabase-pairs and comprise ~350–400 highly methylated rRNA gene units, each ~10 kb (ref. 54). The sequence between the euchromatic arms and NORs has been determined. Elsewhere in the genome, only one other 18S, 5.8S, 25S rRNA gene unit was identified in centromere 3. Although minor variations in sequence length and composition occur in the NOR repeats, these variants are highly clustered, supporting a model of sequence maintenance through concerted evolution<sup>55</sup>.

*Arabidopsis* telomeres are composed of CCCTAAA repeats and average ~2–3 kb (ref. 56). For TEL4N (telomere 4 North), consensus repeats are adjacent to the NOR; the remaining telomeres are typically separated from coding sequences by repetitive subtelomeric regions measuring less than 4 kb. Imperfect telomere-like arrays of up to 24 kb are found elsewhere in the genome, particularly

near centromeres. These arrays might affect the expression of nearby genes and may have resulted from ancient rearrangements, such as inversions of the chromosome arms.

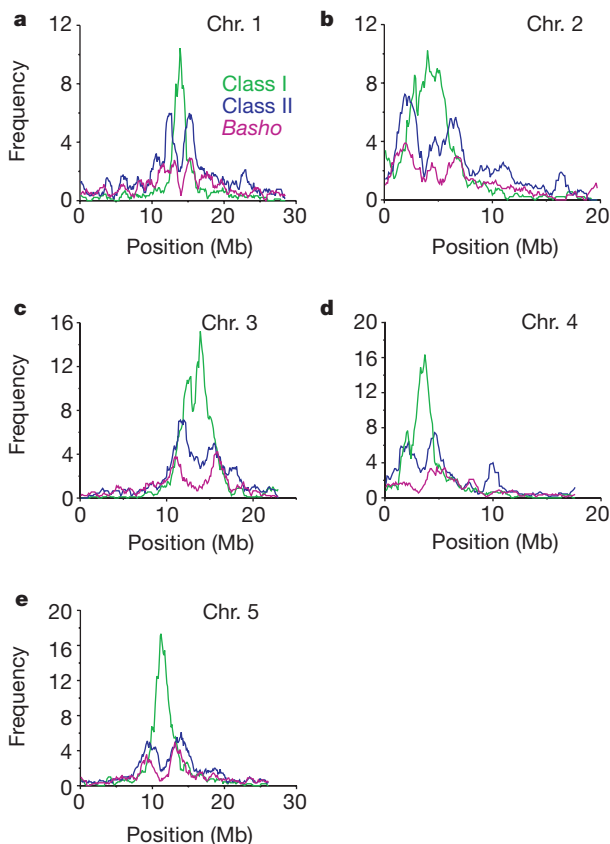
Centromere DNA mediates chromosome attachment to the meiotic and mitotic spindles and often forms dense heterochromatin. Genetic mapping of the regions that confer centromere function provided the markers necessary to precisely place BAC clones at individual centromeres<sup>17</sup>; 69 clones were targeted for sequencing, resulting in over 5 Mb of DNA sequence from the centromeric regions. The unsequenced regions of centromeres are composed primarily of long, homogeneous arrays that were characterized previously with physical<sup>57</sup> and genetic mapping<sup>17</sup> and contain over 3 Mb of repetitive arrays, including the 180-bp repeats and 5S rDNA<sup>51</sup> (Fig. 6).

*Arabidopsis* centromeres, like those of many higher eukaryotes, contain numerous repetitive elements including retroelements, transposons, microsatellites and middle repetitive DNA<sup>17</sup>. These repeats are rare in the euchromatic arms and often most abundant in pericentromeric DNA. The repeats, affinity for DNA-binding dyes, dense methylation patterns and inhibition of homologous recombination indicate that the centromeric regions are highly heterochromatic, and such regions are generally viewed as very poor environments for gene expression. Unexpectedly, we found at least 47 expressed genes encoded in the genetically defined centromeres of *Arabidopsis* (<http://preuss.bsd.uchicago.edu/arabidopsis.genome.html>). In several cases, these genes reside on islands of unique sequence flanked by repetitive arrays, such as 180-bp or 5S rDNA repeats. Among the genes encoded in the centromeres are members of 11 of the 16 functional categories that comprise the proteome. The centromeres are not subject to recombination; consequently, genes residing in these regions probably exhibit unique patterns of molecular evolution.

The function of higher eukaryotic centromeres may be specified by proteins that bind to centromere DNA, by epigenetic modifications, or by secondary or higher order structures. A pairwise comparison of the non-repetitive portions of all five centromeres showed they share limited (1–7%) sequence similarity. Forty-one families of small, conserved centromere sequences (AtCCS, see <http://preuss.bsd.uchicago.edu/arabidopsis.genome.html>) are enriched in the centromeric and pericentromeric regions and differ from sequences found in the centromeres of other eukaryotes. Molecular and genetic assays will be required to determine whether these conserved motifs nucleate *Arabidopsis* centromere activity. Apart from the AtCCS sequences, most centromere DNA is not shared between chromosomes, complicating efforts to derive clear evolutionary relationships. In contrast, genetic and cytological assays indicate that homologous centromeres are highly conserved among *Arabidopsis* accessions, albeit subject to rearrangements such as inversions to form knobs<sup>5,58,59</sup> and insertions<sup>4</sup>. Further investigation of centromere DNA promises to yield information on the evolutionary forces that act in regions of limited recombination, as well as an improved understanding of the role of DNA sequence patterns in chromosome segregation.

**Membrane transport**

Transporters in the plasma and intracellular membranes of *Arabidopsis* are responsible for the acquisition, redistribution and compartmentalization of organic nutrients and inorganic ions, as well as for the efflux of toxic compounds and metabolic end products, energy and signal transduction, and turgor generation. Previous genomic analyses of membrane transport systems in *S. cerevisiae* and *C. elegans* led to the identification of over 100 distinct families of membrane transporters<sup>60,61</sup>. We compared membrane transport processes between *Arabidopsis*, animals, fungi and prokaryotes, and identified over 600 predicted membrane transport systems in *Arabidopsis* (<http://www-biology.ucsd.edu/~ipaulsen/transport/>), a similar number to that of *C. elegans*



**Figure 5** Distribution of class I, II and *Bashed* transposons in *Arabidopsis* chromosomes. The frequency of class I retroelements (green), class II DNA transposons (blue) and *Bashed* elements (purple) are shown at 100-kb intervals along the five chromosomes (a–e) of *Arabidopsis*.

(~700 transporters) and over twofold greater than either *S. cerevisiae* or *E. coli* (~300 transporters).

We compared the transporter complement of *Arabidopsis*, *C. elegans* and *S. cerevisiae* in terms of energy coupling mechanisms (Fig. 7a). Unlike animals, which use a sodium ion P-type ATPase pump to generate an electrochemical gradient across the plasma membrane, plants and fungi use a proton P-type ATPase pump to form a large membrane potential (-250 mV)<sup>62</sup>. Consequently, plant secondary transporters are typically coupled to protons rather than to sodium<sup>63</sup>. Compared with *C. elegans*, *Arabidopsis* has a surprisingly high percentage of primary ATP-dependent transporters (12% and 21% of transporters, respectively), reflecting increased numbers of P-type ATPases involved in metal ion transport and ABC ATPases proposed to be involved in sequestering unusual metabolites and drugs in the vacuole or in other intracellular compartments. These processes may be necessary for pathogen defence and nutrient storage.

About 15% of the transporters in *Arabidopsis* are channel proteins, five times more than in any single-celled organism but half the number in *C. elegans* (Fig. 7b). Almost half of the *Arabidopsis* channel proteins are aquaporins, and *Arabidopsis* has 10-fold more Mfamily major intrinsic protein (MIP) family water channels than any other sequenced organism. This abundance emphasizes the importance of hydraulics in a wide range of plant processes, including sugar and nutrient transport into and out of the vasculature, opening of stomatal apertures, cell elongation and epinastic movements of leaves and stems. Although *Arabidopsis* has a diverse range of metal cation transporters, *C. elegans* has more, many of which function in cell-cell signalling and nerve signal transduction. *Arabidopsis* also possesses transporters for inorganic anions such as phosphate, sulphate, nitrate and chloride, as well as for metal cation channels that serve in signal transduction or cell homeostasis. Compared with other sequenced organisms, *Arabidopsis* has 10-fold more predicted peptide transporters, primarily of the proton-dependent oligopeptide transport (POT) family, emphasizing the importance of peptide transport or indicating that there is broader substrate specificity than previously realized. There are nearly 1,000 *Arabidopsis* genes encoding Ser/Thr protein kinases, suggesting that peptides may have an important role in plant signalling<sup>64</sup>.

Virtually no transporters for carboxylates, such as lactate and pyruvate, were identified in the *Arabidopsis* genome. About 12% of the transporters were predicted to be sugar transporters, mostly consisting of paralogues of the MFS family of hexose transporters. Notably, *S. cerevisiae*, *C. elegans* and most prokaryotes use APC family transporters as their principle means of amino-acid

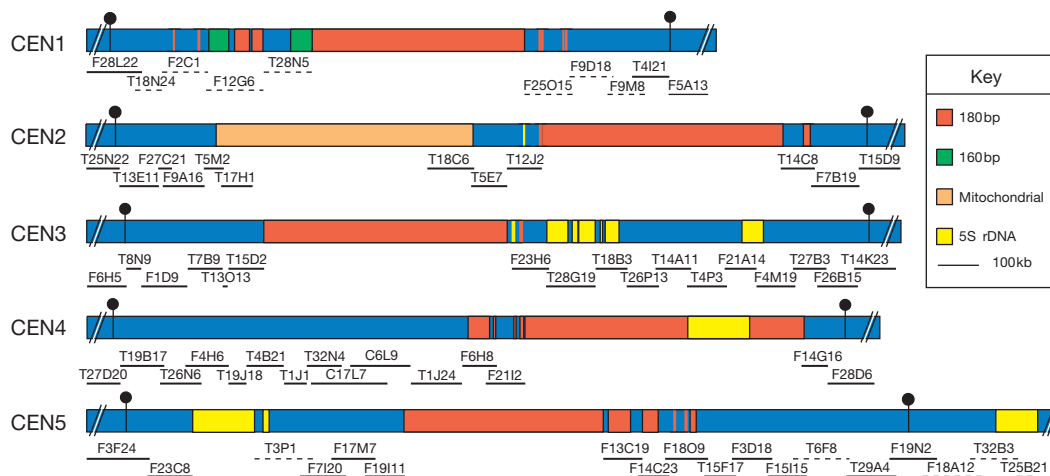
transport, but *Arabidopsis* appears to rely primarily on the AAAP family of amino-acid and auxin transporters. More than 10% of the transporters in *Arabidopsis* are homologous to drug efflux pumps; these probably represent transporters involved in the sequestration into vacuoles of xenobiotics, secondary metabolites, and breakdown products of chlorophyll.

Surprisingly, *Arabidopsis* has close homologues of the human ABC TAP transporters of antigenic peptides for presentation to the major histocompatibility complex (MHC). In *Arabidopsis*, these transporters may be involved in peptide efflux, or more speculatively, in some form of cell-recognition response. *Arabidopsis* also has 10-fold more members of the multi-drug and toxin extrusion (MATE) family than any other sequenced organism; in bacteria, these transporters function as drug efflux pumps. Curiously, *Arabidopsis* has several homologues of the *Drosophila* RND transporter family Patched protein, which functions in segment polarity, and more than ten homologues of the *Drosophila* ABC family eye pigment transporters. In plants, these are presumably involved in intracellular sequestration of secondary metabolites.

### DNA repair and recombination

DNA repair and recombination pathways have many functions in different species such as maintaining genomic integrity, regulating mutation rates, chromosome segregation and recombination, genetic exchange within and between populations, and immune system development. Comparing the *Arabidopsis* genome with other species<sup>65</sup> indicates that *Arabidopsis* has a similar set of DNA repair and recombination (RAR) genes to most other eukaryotes. The pathways represented include photoreactivation, DNA ligation, non-homologous end joining, base excision repair, mismatch excision repair, nucleotide excision repair and many aspects of DNA recombination (Supplementary Information Table 5). The *Arabidopsis* RAR genes include homologues of many DNA repair genes that are defective in different human diseases (for example, hereditary breast cancer and non-polyposis colon cancer, xeroderma pigmentosum and Cockayne's syndrome).

One feature that sets *Arabidopsis* apart from other eukaryotes is the presence of additional homologues of many RAR genes. This is seen for almost every major class of DNA repair, including recombination (four RecA), DNA ligation (four DNA ligase I), photoreactivation (one class II photolyase and five class I photolyase homologues) and nucleotide excision repair (six RPA1, two RPA2, two Rad25, three TFB1 and four Rad23). This is most striking for genes with probable roles in base excision repair. *Arabidopsis* encodes 16 homologues of DNA base glycosylases (enzymes that



**Figure 6** Predicted centromere composition. Genetically defined centromere boundaries are indicated by filled circles; fully and partially assembled BAC sequences are represented by solid and dashed black lines, respectively. Estimates of repeat sizes within

the centromeres were derived from consideration of repeat copy number, physical mapping and cytogenetic assays.

recognize abnormal DNA bases and cleave them from the sugar-phosphate backbone)—more than any other species known. This includes several homologues of each of three families of alkylation damage base glycosylases: two of the *S. cerevisiae* MPG; six of the *E. coli* TagI; and two of the *E. coli* AlkA. *Arabidopsis* also encodes three homologues of the apurinic-apyrimidinic (AP) endonuclease Xth. AP endonucleases continue the base excision repair started by glycosylases by cleaving the DNA backbone at abasic sites.

Evolutionary analysis indicates that some of the extra copies of RAR genes in *Arabidopsis* originated through relatively recent gene duplications—because many of the sets of genes are more closely related to each other than to their homologues in any other species. As duplication is frequently accompanied by functional divergence, the duplicate (paralogous) genes may have different repair specificities or may have evolved functions that are outside RAR functions (as is the case for two of the five class I photolyase homologues, which function as blue-light receptors). In most cases, it is not known whether the paralogous gene copies have different functions. The presence of multiple paralogues might also allow functional redundancy or a greater repair or recombination capacity.

The multiplicity of RAR genes in *Arabidopsis* is also partly due to the transfer of genes from the organellar genomes to the nucleus. Repair gene homologues that appear to be of chloroplast origin (Supplementary Information Tables 2 and 5) include the recombination proteins RecA, RecG and SMS, two class I photolyase homologues, Fpg, two MutS2 proteins, and the transcription-repair coupling factor Mfd. Two of these (RecA and Fpg) are involved in RAR functions in the plastid, suggesting that the others may be as well. The finding of an Mfd orthologue of cyanobacterial descent is surprising. In *E. coli*, Mfd couples nucleotide excision repair carried out by UvrABC to transcription, leading to the rapid repair of DNA damage on the transcribed strand of transcribed genes<sup>66</sup>. The absence of orthologues of UvrABC in *Arabidopsis* renders the function of Mfd difficult to predict. The presence of Mfd but not UvrABC has been reported for only one other species, a bacterial endosymbiont of the pea aphid.

Other nuclear-encoded *Arabidopsis* DNA repair gene homologues are evolutionarily related to genes from  $\alpha$ -Proteobacteria, and thus may be of mitochondrial descent. In particular, the six homologues of the alkyl-base glycosylase TagI appear to be the result of a large expansion in plants after transfer from the mitochondrial genome. Whether any of these TagI homologues function in the repair and maintenance of mitochondrial DNA has not been determined. More detailed phylogenetic analysis may reveal additional *Arabidopsis* RAR genes to be of organellar ancestry.

There are some notable absences of proteins important for RAR in other species, including alkyltransferases, MSH4, RPA3 and many components of TFIIH (TFB2, TFB3, TFB4, CCL1, Kin28). Nevertheless, *Arabidopsis* shows many similarities to the set of DNA repair genes found in other eukaryotes, and therefore offers an experimental system for determining the functions of many of these proteins, in part through characterization of mutants defective in DNA repair<sup>67</sup>.

### Gene regulation

Eukaryotic gene expression involves many nuclear proteins that modulate chromatin structure, contribute to the basal transcription machinery, or mediate gene regulation in response to developmental, environmental or metabolic cues. As predicted by sequence similarity, more than 3,000 such proteins may be encoded by the *Arabidopsis* genome, suggesting that it has a comparable complexity of gene regulation to other eukaryotes. *Arabidopsis* has an additional level of gene regulation, however, with DNA methylation potentially mediating gene silencing and parental imprinting.

Plants have evolved several variations on chromatin remodelling proteins, such as the family of HD2 histone deacetylases<sup>68</sup>. Although *Arabidopsis* possesses the usual number of SNF2-type chromatin

remodelling ATPases, which regulate the expression of nearly all genes, there are significant structural differences between yeast and metazoan SNF2-type genes and their orthologues in *Arabidopsis*. DDM1, a member of the SNF2 superfamily, and MOM1, a gene with similarity to the SNF2 family, are involved in transcriptional gene silencing in *Arabidopsis*. MOM1 has no clear orthologue in fungal or metazoan genomes.

Consistent with its methylated DNA, *Arabidopsis* possesses eight DNA methyltransferases (DMTs). Two of the three types are orthologous to mammalian DMT<sup>69</sup> whereas one, chromomethyltransferase<sup>70</sup>, is unique to plants. No DMTs are found in yeast or *C. elegans*, although two DMT-like genes are found in *Drosophila*<sup>71</sup>. *Arabidopsis* also encodes eight proteins with methyl-DNA-binding domains (MBDs). Despite lacking methylated DNA, *Drosophila* encodes four MBD proteins and *C. elegans* has two. These differences in chromatin components are likely to reflect important differences in chromatin-based regulatory control of gene expression in eukaryotes (Supplementary Information Table 6; <http://Ag.Arizona.Edu/chromatin/chromatin.html>).

The *Arabidopsis* genome encodes transcription machinery for the three nuclear DNA-dependent RNA polymerase systems typical of eukaryotes (Supplementary Information Table 6). Transcription by RNA polymerases II and III appears to involve the same machinery as is used in other eukaryotes; however, most transcription factors for RNA polymerase I are not readily identified. Only two polymerase I regulators (other than polymerase subunits and TATA-binding protein) are apparent in *Arabidopsis*, namely homologues of yeast RRN3 and mouse TTF-1. All eukaryotes examined to date have distinct genes for the largest and second largest subunits of polymerase I, II and III. Unexpectedly, *Arabidopsis* has two genes encoding a fourth class of largest subunit and second-largest subunit (Supplementary Information Fig. 5). It will be interesting to determine whether the atypical subunits comprise a polymerase that has a plant-specific function. Four genes encoding single-subunit plastid or mitochondrial RNA polymerases have been identified in *Arabidopsis* (Supplementary Information Table 6). Genes for the bacterial  $\beta$ -,  $\beta'$ - and  $\alpha$ -subunits of RNA polymerase are also present, as are homologues of various  $\sigma$ -factors, and these proteins may regulate chloroplast gene expression. Mutations in the *Sde-1* gene, encoding RNA-dependent RNA polymerase (RdRp), lead to defective post-transcriptional gene silencing<sup>72</sup>. We also identified five more closely related RdRp genes.

Our analysis, using both similarity searches and domain matches, has identified 1,709 proteins with significant similarity to known classes of plant transcription factors classified by conserved DNA-binding domains. This analysis used a consistent conservative threshold that probably underestimates the size of families of diverse sequence. This class of protein is the least conserved among all classes of known proteins, showing only 8–23% similarity to transcription factors in other eukaryotes (Fig. 2b). This reduced similarity is due to the absence of certain classes of transcription factors in *Arabidopsis* and large numbers of plant-specific transcription factors. We did not detect any members of several widespread families of transcription factors, such as the REL (Rel-like DNA-binding domain) homology region proteins, nuclear steroid receptors and forkhead-winged helix and POU (Pit-1, Oct- and Unc-8b) domain families of developmental regulators. Conversely, of 29 classes of *Arabidopsis* transcription factors, 16 appear to be unique to plants (Supplementary Information Table 6). Several of these, such as the AP2/EREBP-RAV, NAC and ARF-AUX/IAA families, contain unique DNA-binding domains, whereas others contain plant-specific variants of more widespread domains, such as the DOF and WRKY zinc-finger families and the two-repeat MYB family.

Functional redundancy among members of large families of closely related transcription factors in *Arabidopsis* is a significant potential barrier to their characterization<sup>73</sup>. For example, in the

SHATTERPROOF and SEPALLATA families of MADS box transcription factors, all genes must be defective to produce visible mutant phenotypes<sup>74,75</sup>. These functionally redundant genes are found on the segmental duplications described above. Our analyses, together with the significant sequence similarity found in large families of transcription factors such as the R2R3-repeat MYB and WRKY families, suggest that strategies involving overexpression will be important in determining the functions of members of transcription factor families.

*Arabidopsis* has two or over three times more transcription factors than identified in *Drosophila*<sup>29</sup> or *C.elegans*<sup>1</sup>, respectively. The significantly greater extent of segmental chromosomal and local tandem duplications in the *Arabidopsis* genome generates larger gene families, including transcription factors. The partly overlapping functions defined for a few transcription factors are also likely to be much more widespread, implicating many sequence-related transcription factors in the same cellular processes. Finally, the expanded number of genes involved in metabolism, defence and environmental interaction in *Arabidopsis* (Fig. 2a), which have few counterparts in *Drosophila* and *C. elegans*, all require additional numbers and classes of transcription factors to integrate gene function in response to a vast range of developmental and environmental cues.

### Cellular organization

Plant cells differ from animal cells in many features such as plastids, vacuoles, Golgi organization, cytoskeletal arrays, plasmodesmata linking cytoplasm of neighbouring cells, and a rigid polysaccharide-rich extracellular matrix—the cell wall. Because the cell wall maintains the position of a cell relative to its neighbours, both changes in cell shape and organized cell divisions, involving cytoskeleton reorganization and membrane vesicle targeting, have major roles in plant development. Plant cytokinesis is also unique in that the partitioning membrane is formed *de novo* by vesicle fusion. We compared the *Arabidopsis* genome with those of *C. elegans*,

*Drosophila* and yeast to glimpse the genetic basis of plant-cell-specific features.

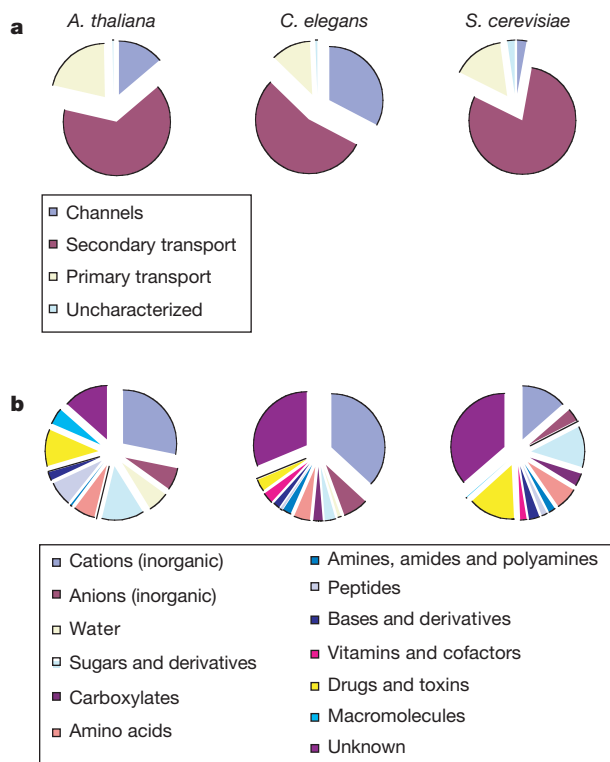
The principal components of the plant cytoskeleton are microtubules (MTs) and actin filaments (AFs); intermediate filaments (IFs) have not been described in plants. *Arabidopsis* appears to lack genes for cyokeratin or vimentin, the main components of animal IFs, but has several variants of actin,  $\alpha$ - and  $\beta$ -tubulin. The *Arabidopsis* genome also encodes homologues of chaperones that mediate the folding of tubulin and actin polypeptides in yeast and animal cells, such as the prefoldin and cytosolic chaperonin complexes and tubulin-folding cofactors. The dynamic stability of MTs and AFs is influenced by MT-associated proteins and actin-binding proteins, respectively, several of which are encoded by *Arabidopsis* genes. These include the MT-severing ATPase katanin, AF-cross-linking/bundling proteins, such as fimbrins and villins, and AF-disassembling proteins, such as profilin and actin-depolymerizing factor/cofilin. The *Arabidopsis* proteome appears to lack homologues of proteins that, in animal cells, link the actin cytoskeleton across the plasma membrane to the extracellular matrix, such as integrin, talin, spectrin,  $\alpha$ -actinin, vitronectin or vinculin. This apparent lack of ‘anchorage’ proteins is consistent with the different composition of the cell wall and with a prominence of cortical MTs at the expense of cortical AFs in plant cells.

Plant-specific cytoskeletal arrays include interphase cortical MTs mediating cell shape, the preprophase band marking the cortical site of cell division, and the phragmoplast assisting in cytokinesis<sup>76</sup>. Although plant cells lack structural counterparts of the yeast spindle pole body and the animal centrosome, *Arabidopsis* has homologues of core components of the MT-nucleating  $\gamma$ -tubulin ring complex, such as  $\gamma$ -tubulin, Spc97/hGCP2 and Spc98/hGCP3. *Arabidopsis* has numerous motor molecules, both kinesins and dyneins with associated dynactin complex proteins, which are presumably involved in the dynamic organization of MTs and in transporting cargo along MT tracks. There are also myosin motors that may be involved in AF-supported organelle trafficking. Essential features of the eukaryotic cytoskeleton appear to be conserved in *Arabidopsis*.

The *Arabidopsis* genome encodes homologues of proteins involved in vesicle budding, including several ARFs and ARF-related small G-proteins, large but not small ARF GEFs (adenosine ribosylation factor on guanine nucleotide exchange factor), adapter proteins, and coat proteins of the COP and non-COP types. *Arabidopsis* also has homologues of proteins involved in vesicle docking and fusion, including SNAP receptors (SNAREs), N-ethylmaleimide-sensitive factor (NSF) and Cdc48-related ATPases, accessory proteins such as Sec1 and soluble NSF attachment protein (SNAP), and Rab-type GTPases. The large number of *Arabidopsis* SNAREs can be grouped by sequence similarity to yeast and animal counterparts involved in specific trafficking pathways, and some have been localized to the trans-Golgi and the pre-vacuolar pathway<sup>77</sup>. *Arabidopsis* also has a receptor for retention of proteins in the endoplasmic reticulum, a cargo receptor for transport to the vacuole and several phragmoplastins related to animal dynamin GTPases. Thus, plant cells appear to use the same basic machinery for vesicle trafficking as yeast and animal cells.

Animal cells possess many functionally diverse small G-proteins of the Ras superfamily involved in signal transduction, AF reorganization, vesicle fusion and other processes. Surprisingly, *Arabidopsis* appears to lack genes for G-proteins of the Ras, Rho, Rac and Cdc42 subfamilies but has many Rab-type G-proteins involved in vesicle fusion and several Rop-type G-proteins, one of which has a role in actin organization of the tip-growing pollen tube<sup>78</sup>. The significance of this divergent amplification of different subfamilies of small G-proteins in plants and animals remains to be determined.

*Arabidopsis* possesses cyclin-dependent kinases (CDKs), including a plant-specific Cdc2b kinase expressed in a cell-cycle-dependent manner, several cyclin subtypes, including a D-type cyclin that



**Figure 7** Comparison of the transport capabilities of *Arabidopsis*, *C. elegans* and *S. cerevisiae*. Pie charts show the percentage of transporters in each organism according to bioenergetics (a) and substrate specificity (b).

mediates cytokinin-stimulated cell-cycle progression<sup>79</sup>, a retinoblastoma-related protein and components of the ubiquitin-dependent proteolytic pathway of cyclin degradation. In yeast and animal cells, chromosome condensation is mediated by condensins, sister chromatids are held together by cohesins such as Scc1, and metaphase–anaphase transition is triggered by separin/Esp1 endopeptidase proteolysis of Scc1 on APC-mediated degradation of its inhibitor, securin/Psd1. Related proteins are encoded by the *Arabidopsis* genome. Thus, the basic machinery of cell-cycle progression, genome duplication and segregation appears to be conserved in plants. By contrast, entry into M phase, M-phase progression and cytokinesis seem to be modified in plant cells. *Arabidopsis* does not appear to have homologues of Cdc25 phosphatase, which activates Cdc2 kinase at the onset of mitosis, or of polo kinase, which regulates M-phase progression in yeast and animals. Conversely, plant-specific mitogen-activated protein (MAP) kinases appear to be involved in cytokinesis.

Cytokinesis partitions the cytoplasm of the dividing cell. Yeast and animal cells expand the membrane from the surface towards the centre in a cleavage process supported by septins and a contractile ring of actin and type II myosin. By contrast, plant cytokinesis starts in the centre of the division plane and progresses laterally. A transient membrane compartment, the cell plate, is formed *de novo* by fusion of Golgi-derived vesicles trafficking along the phragmoplast MTs<sup>80</sup>. Consistent with the unique mode of plant cytokinesis, *Arabidopsis* appears to lack genes for septins and type II myosin. Conversely, cell-plate formation requires a cytokinesis-specific syntaxin that has no close homologue in yeast and animals. Although syntaxin-mediated membrane fusion occurs in animal cytokinesis and cellularization, the vesicles are delivered to the base of the cleavage furrow. Thus, the plant-specific mechanism of cell division is linked to conserved eukaryotic cell-cycle machinery.

Two main conclusions are suggested by this comparative analysis. First, *Arabidopsis* and eukaryotic cells have common features related to intracellular activities, such as vesicle trafficking, cytoskeleton and cell cycle. Second, evolutionarily divergent features, such as organization of the cytoskeleton and cytokinesis, appear to relate to the plant cell wall.

## Development

The regulation of development in *Arabidopsis*, as in animals, involves cell–cell communication, hierarchies of transcription factors, and the regulation of chromatin state; however, there is no reason to suppose that the complex multicellular states of plant and animal development have evolved by elaborating the same general processes during the 1.6 billion years since the last common unicellular ancestor of plants and animals<sup>81,82</sup>. Our genome analyses reflect the long, independent evolution of many processes contributing to development in the two kingdoms.

Plants and animals have converged on similar processes of pattern formation, but have used and expanded different transcription factor families as key causal regulators. For example, segmentation in insects and differentiation along the anterior–posterior and limb axes in mammals both involve the spatially specific activation of a series of homeobox gene family members. The pattern of activation is causal in the later differentiation of body and limb axis regions. In plants the pattern of floral whorls (sepals, petals, stamens, carpels) is also established by the spatially specific activation of members of a family of transcription factors, but in this instance the family is the MADS box family. Plants also have homeobox genes and animals have MADS box genes, implying that each lineage invented separately its mechanism of spatial pattern formation, while converging on actions and interactions of transcription factors as the mechanism. Other examples show even greater divergence of plant and animal developmental control. Examples are the AP2/EREBP and NAC families of transcription factors, which have important roles in flower and meristem development; both families are so far found

only in plants (Supplementary Information Table 6).

A similar story can be told for cell–cell communication. Plants do not seem to have receptor tyrosine kinases, but the *Arabidopsis* genome has at least 340 genes for receptor Ser/Thr kinases, belonging to many different families, defined by their putative extracellular domains (Supplementary Information Table 7). Several families have members with known functions in cell–cell communication, such as the CLV1 receptor involved in meristem cell signalling, the S-glycoprotein homologues involved in signalling from pollen to stigma in self-incompatible *Brassica* species, and the BRI1 receptor necessary for brassinosteroid signalling<sup>83</sup>. Animals also have receptor Ser/Thr kinases, such as the transforming growth factor- $\beta$  (TGF- $\beta$ ) receptors, but these act through SMAD proteins that are absent from *Arabidopsis*. The leucine-rich repeat (LRR) family of *Arabidopsis* receptor kinases shares its extracellular domain with many animal and fungal proteins that do not have associated kinase domains, and there are at least 122 *Arabidopsis* genes that code for LRR proteins without a kinase domain. Other *Arabidopsis* receptor kinase families have extracellular domains that are unfamiliar in animals. Thus, evolution is modular, and the plant and animal lineages have expanded different families of receptor kinases for a similar set of developmental processes.

Several *Arabidopsis* genes of developmental importance appear to be derived from a cyanobacteria-like genome (Supplementary Information Table 2), with no close relationship to any animal or fungal protein. One salient example is the family of ethylene receptors; another gene family of apparent chloroplast origin is the phytochromes—light receptors involved in many developmental decisions (see below). Whereas the land plant phytochromes show clear homology to the cyanobacterial light receptors, which are typical prokaryotic histidine kinases, the plant phytochromes are histidine kinase paralogues with Ser/Thr specificity<sup>84</sup>. Similarly to the ethylene receptors, the proteins that act downstream of plant phytochrome signalling are not found in cyanobacteria, and thus it appears that a bacterial light receptor entered the plant genome through horizontal transfer, altered its enzymatic activity, and became linked to a eukaryotic signal transduction pathway. This infusion of genes from a cyanobacterial endosymbiont shows that plants have a richer heritage of ancestral genes than animals, and unique developmental processes that derive from horizontal gene transfer.

## Signal transduction

Being generally sessile organisms, plants have to respond to local environmental conditions by changing their physiology or redirecting their growth. Signals from the environment include light and pathogen attack, temperature, water, nutrients, touch and gravity. In addition to local cellular responses, some stimuli are communicated across the plant body, with plant hormones and peptides acting as secondary messengers. Some hormones, such as auxin, are taken up into the cell, whereas others, such as ethylene and brassinosteroids, and the peptide CLV3, act as ligands for receptor kinases on the plasma membrane. No matter where the signal is perceived by the cell, it is transduced to the nucleus, resulting in altered patterns of gene expression.

Comparative genome analysis between *Arabidopsis*, *C. elegans* and *Drosophila* supports the idea that plants have evolved their own pathways of signal transduction<sup>85</sup>. None of the components of the widely adopted signalling pathways found in vertebrates, flies or worms, such as Wingless/Wnt, Hedgehog, Notch/lin12, JAK/STAT, TGF- $\beta$ /SMADs, receptor tyrosine kinase/Ras or the nuclear steroid hormone receptors, is found in *Arabidopsis*. By contrast, brassinosteroids are ligands of the BRI1 Ser/Thr kinase, a member of the largest recognizable class of transmembrane sensors encoded by 340 receptor-like kinase (RLK) genes in the *Arabidopsis* genome (Supplementary Information Table 7). With a few notable exceptions, such as CLV1, the types of ligands sensed by RLKs are

completely unknown, providing an enormous future challenge for plant biologists. G-protein-coupled receptors (GPCRs)/ seven-transmembrane proteins are an abundant class of proteins in mammalian genomes, instrumental in signal transduction. INTERPRO detected 27 GPCR-related domains in *Arabidopsis* (Supplementary Information Table 1), although there is no direct experimental evidence for these. *Arabidopsis* contains a family of 18 seven-transmembrane proteins of the mildew resistance (MIO) class, several of which are involved in defence responses. Notably, only single G $\alpha$  (GPA1) and G $\beta$  (AGB1) subunits are found in *Arabidopsis*, both previously known<sup>86</sup>.

Although cyclic GMP has been proposed to be involved in signal transduction in *Arabidopsis*<sup>87</sup>, a protein containing a guanylate cyclase domain was not identified in our analyses. Nevertheless, cyclic nucleotide-binding domains were detected in various proteins, indicating that cNMPs may have a role in plant signal transduction. Thus, although cNMP-binding domains appear to have been conserved during evolution, cNMP synthesis in *Arabidopsis* may have evolved independently.

We were unable to identify a protein with significant similarity to known G $\gamma$  subunits, but recent biochemical studies suggest that a protein with this functional capacity is likely to be present in plant cells (H. Ma, personal communication). Therefore, there is potential for the formation of only a single heterotrimeric G-protein complex; however, its functional interaction with any of the potential GPCR-related proteins remains to be determined.

Modules of cellular signal pathways from bacteria and animals have been combined and new cascades have been innovated in plants. A pertinent example is the response to the gaseous plant hormone ethylene<sup>88</sup>. Ethylene is perceived and its signal transmitted by a family of receptors related to bacterial-type two-component histidine kinases (HKs). In bacteria, yeast and plants, these proteins sense many extracellular signals and function in a His-to-Asp phosphorelay network<sup>89</sup>. In turn, these proteins physically interact with the genetically downstream protein CTR1, a Raf/MAPKKK-related kinase, revealing the juxtaposition of bacterial-type two-component receptors and animal-type MAP kinase cascades. Unlike animals, however, *Arabidopsis* does not seem to have a Ras protein to activate the MAP kinase cascade. MAP kinases are found in abundance in *Arabidopsis*: we identified ~20, a higher number than in any other eukaryote. As potentially counteracting components, we found ~70 putative PP2C protein phosphatases. Although this group is largely uncharacterized functionally, several members are related to ABI1/ABI2, key negative regulators in the signalling pathway for the plant hormone abscisic acid. Additional components of the His-to-Asp phosphorelay system were also found in *Arabidopsis*, including authentic response regulators (ARRs), pseudoresponse regulators (PRRs) and phosphotransfer intermediate protein (HPT)<sup>90</sup>. We found 11 HKs in the proteome (3 new), 16 RRs (2 new) and 8 PRRs (2 new). The biological roles of most ARR, PRR and HPTs are largely unknown, but several have been found to have diverse functions in plants, including transcriptional activation in response to the plant hormone cytokinin<sup>91</sup>, and as components of the circadian clock<sup>92</sup>.

Plants seem to have evolved unique signalling pathways by combining a conserved MAP kinase cascade module with new receptor types. In many cases, however, the ligands are unknown. Conversely, some known signalling molecules, such as auxin, are still in search of a receptor. Auxin signalling may represent yet another plant-specific mode of signalling, with protein degradation through the ubiquitin-proteasome pathway preceding altered gene expression. With many *Arabidopsis* genes encoding components of the ubiquitin-proteasome pathway, elimination of negative regulators may be a more widespread phenomenon in plant signalling.

### Recognizing and responding to pathogens

Plants are constantly exposed to pests, parasites and pathogens and

have evolved many defences. In mammals, polymorphism for parasite recognition encoded in the MHC genes contributes to resistance. In plants, disease resistance (*R*) genes that confer parasite recognition are also extremely polymorphic. This polymorphism has been proposed to restrict parasites, and its absence may explain the breakdown of resistance in crop monocultures<sup>93</sup>. In contrast to MHC genes, plant resistance genes are found at several loci, and the complete genome sequence enables analysis of their complement and structure. Parasite recognition by resistance genes triggers defence mechanisms through various signalling molecules, such as protein kinases and adapter proteins, ion fluxes, reactive oxygen intermediates and nitric oxide. These halt pathogen colonization through transcriptional activation of defence genes and a form of programmed cell death called the hypersensitive response<sup>94</sup>. The *Arabidopsis* genome contains diverse resistance genes distributed at many loci, along with components of signalling pathways, and many other genes whose role in disease resistance has been inferred from mutant phenotypes.

Most resistance genes encode intracellular proteins with a nucleotide-binding (NB) site typical of small G proteins, and carboxy-terminal LRRs<sup>95</sup>. Their amino termini either carry a TIR domain, or a putative coiled coil (CC). There are 85 TIR–NB–LRR resistance genes at 64 loci, and 36 CC–NB–LRR resistance genes at 30 loci. Some NB–LRR resistance genes express neither obvious TIR nor CC domains at their N termini. This potential class is present seven times, at six loci. There are 15 truncated TIR–NB genes that lack an LRR at 10 loci, often adjacent to full TIR–NB–LRR genes. There are also six CC–NB genes, at five loci. These truncated products may function in resistance. Intriguingly, two TIR–NB–LRR genes carry a WRKY domain, found in transcription factors that are implicated in plant defence, and one of these also encodes a protein kinase domain.

Resistance gene evolution may involve duplication and divergence of linked gene families<sup>36</sup>; however, most (46) resistance genes are singletons; 50 are in pairs, 21 are in 7 clusters of 3 family members, with single clusters of 4, 5, 7, 8 and 9 members, respectively. Of the non-singletons, ~60% of pairs are in direct repeats, and ~40% are in inverted repeats. Resistance genes are unevenly distributed between chromosomes, with 49 on chromosome 1; 2 on chromosome 2; 16 on chromosome 3; 28 on chromosome 4; and 55 on chromosome 5.

In other plant species, resistance genes encode both transmembrane receptors for secreted pathogen products and protein kinases, and some other classes are also found. The *Cf* genes in tomato encode extracellular LRRs with a transmembrane domain and short cytoplasmic domain. Mutation in an *Arabidopsis* homologue, *CLAVATA2*, results in enlarged meristems, but to date no resistance function has been assigned to the 30 *Arabidopsis* *CLV2* homologues. *CLAVATA1*, a transmembrane LRR kinase, is also required for meristem function. *Xa21*, a rice LRR-kinase, confers *Xanthomonas* resistance, and the *Arabidopsis* *FLS2* LRR kinase confers recognition of flagellin. It has been proposed that *CLV1* and *CLV2* function as a heterodimer; perhaps this is also true for *Xa21*, *FLS2* and *Cf* proteins. There are 174 LRR transmembrane kinases in *Arabidopsis*, with only *FLS2* assigned a role in resistance. A unique resistance gene, beet *Hs1pro-1*, which confers nematode resistance, has two *Arabidopsis* homologues.

The tomato Pto Ser/Thr kinase acts as a resistance protein in conjunction with an NB–LRR protein, so similar kinases might do the same for *Arabidopsis* NB–LRR proteins. There are 860 Ser/Thr kinases in the *Arabidopsis* sequence. Fifteen of these share 50% identity over the Pto-aligned region. The Toll pathway in *Drosophila* and mammals regulates innate immune responses through LRR/TIR domain receptors that recognize bacterial lipopolysaccharides<sup>96</sup>. Pto is highly homologous to *Drosophila* PELLE and mammalian IRAK protein kinases that mediate the TIR pathway.

Additional genes have been defined that are required for resistance by our analysis of the genome sequence. The *ndr1* mutation defines a gene required by the CC–NB–LRR gene *RPS2* and *RPM1*. *NDR1* is 1 of 28 *Arabidopsis* genes that are similar both to each other and to the tobacco *HIN1* gene that is transcriptionally induced early during the hypersensitive response. *EDS1* is a gene required for TIR–NB–LRR function, and like *PAD4*, encodes a protein with a putative lipase motif. *EDS1*, *PAD4* and a third gene comprise the *EDS1/PAD4* family. The *NPR1/NIM1/SAI1* gene is required for systemic acquired resistance, and we found five additional *NPR1* homologues. Recessive mutations at both the barley Mlo and *Arabidopsis* *LSD1* loci confer broad-spectrum resistance and derepress a cell-death program. There are at least 18 Mlo family members that resemble heterotrimeric GPCRs in *Arabidopsis*, and only two *LSD1* homologues.

One of the earliest responses to pathogen recognition is the production of reactive oxygen intermediates. This involves a specialized respiratory burst oxidase protein that transfers an electron across the plasma membrane to make superoxide. *Arabidopsis* encodes eight apparently functional *gp91* homologues, called *Atrboh* genes. Unlike *gp91*, they all carry an ~300 amino-acid N-terminal extension carrying an EF-hand  $\text{Ca}^{2+}$ -binding domain. In mammals, activation of the respiratory oxidative burst complex in the neutrophil, which includes *gp91*, requires the action of Rac proteins. As no Rac or Ras proteins are found in *Arabidopsis*, members of the large *rop* family of G proteins may carry this out. Similarly, we did not detect any *Arabidopsis* homologues of other mammalian respiratory burst oxidase components (p22, p47, p67, p40).

There are no clear homologues of many mammalian defence and cell-death control genes. Although nitric oxide production is involved in plant defence, there is no obvious homologue of nitric oxide synthase. Also absent are apparent homologues of the REL domain transcription factors involved in innate immunity in both *Drosophila* and mammals. We found no similarity to proteins involved in regulating apoptosis in animal cells, such as classical caspases, *bcl2/ced9* and baculovirus p35. There are, however, 36 cysteine proteases. There are also eight homologues of a newly defined metacaspase family<sup>97</sup>, two of which, along with *LSD1*, have a clear GATA-type zinc-finger.

### Photomorphogenesis and photosynthesis

Because nearly all plants are sessile and most depend on photosynthesis, they have evolved unique ways of responding to light. Light serves as an energy source, as well as a trigger and modulator of complex developmental pathways, including those regulated by the circadian clock. Light is especially important during seedling emergence, where it stimulates chlorophyll production, leaf development, cotyledon expansion, chloroplast biogenesis and the coordinated induction of many nuclear- and chloroplast-encoded genes, while at the same time inhibiting stem growth. The goal of this process, called photomorphogenesis, is the establishment of a body plan that allows the plant to be an efficient photosynthetic machine under varying light conditions<sup>98</sup>. The signal transduction cascade leading to light-induced responses begins with the activation of photoreceptors. Next, the light signal is transduced via positively and negatively acting nuclear and cytoplasmic proteins, causing activation or derepression of nuclear and chloroplast-encoded photosynthetic genes and enabling the plant to establish optimal photoautotrophic growth. Although genetic and biochemical studies have defined many of the components in this process, the genome sequence provides an opportunity to identify comprehensively *Arabidopsis* genes involved in photomorphogenesis and the establishment of photoautotrophic growth. We identified at least 100 candidate genes involved in light perception and signalling, and 139 nuclear-encoded genes that potentially function in photosynthesis.

The roles have been described of only 35 of the 100 candidate photomorphogenic genes (Supplementary Information Table 8). All of the light photoreceptors had been discovered previously, including five red/far-red absorbing phytochromes (PHYA–E), two blue/ultraviolet-A absorbing cryptochromes (CRY1 and CRY2), one blue-absorbing phototropin (NPH1) and one NPH1-like (or NPL1). In contrast, we uncovered many new proteins similar to the photomorphogenesis regulators COP/DET/FUS, PKS1, PIF3, NDPK2, SPA1, FAR1, GIGANTEA, FIN219, HY5, CCA1, ATHB-2, ZEITLUPE, FKF1, LKP1, NPH3 and RPT2.

Both the phytochromes and NPH1 contain chromophores for light sensing coupled to kinase domains for signal transmission. Phytochromes have an N-terminal chromophore-binding domain, two PAS domains, and a C-terminal Ser/Thr kinase domain<sup>99</sup>, whereas NPH1 has two LOV domains (members of the PAS domain superfamily) for flavin mononucleotide binding and a C-terminal Ser/Thr kinase domain<sup>100</sup>. PAS domains potentially sense changes in light, redox potential and oxygen energy levels, as well as mediating protein–protein interactions<sup>99,100</sup>. We searched for uncharacterized proteins with the combination of a kinase domain and either a phytochrome chromophore-binding site or PAS domains. Although we found no new phytochrome-like genes, we did identify four predicted proteins that contain PAS and kinase domains (Supplementary Information Fig. 6). These proteins share 80% amino-acid identity, but, unlike NPH1 and NPL1, have only one PAS domain. The combination of potential signal sensing and transmitting domains makes it tempting to speculate that these proteins may be receptors for light or other signals.

Our screen included searches for components of photosynthetic reaction centres and light-harvesting complexes, enzymes involved in  $\text{CO}_2$  fixation and enzymes in pigment biosynthesis. We identified 11 core proteins of photosystem I, including the eukaryotic-specific components PsaG and PsaH<sup>101</sup>, and 8 photosystem II proteins, including a single member (*psbW*) of the photosystem II core. We also found 26 proteins similar to the Chlorophyll-a/b binding proteins (8 Lhca and 18 Lhcb). Of the seven subunits of the cytochrome *b<sub>6</sub>f* complex (PetA–D, PetG, PetL, PetM), only one (*PetC*) was found in the nuclear genome, whereas the remainder are probably encoded in the chloroplast. Similarly, of the nine subunits of the chloroplast ATP synthase complex, three are encoded in the nucleus, including the II-,  $\gamma$ - and  $\delta$ -subunits; the remaining subunits (I, III, IV,  $\alpha$ ,  $\beta$ ,  $\epsilon$ ) are encoded in the chloroplast<sup>102</sup>. Ten genes were related to the soluble components of the electron transfer chain, including two plastocyanins, five ferredoxins and three ferredoxin/NADP oxidoreductases. Forty genes are predicted to have a role in  $\text{CO}_2$  fixation, including all of the enzymes in the Calvin–Benson cycle. For pigment biosynthesis, 16 genes in chlorophyll biosynthesis and 31 genes in carotenoid biosynthesis were found (Supplementary Information Table 8). Our analyses have identified several potential components of the light perception pathway, and have revealed the complex distribution of components of the photosynthetic apparatus between nuclear and plastid genomes.

### Metabolism

*Arabidopsis* is an autotrophic organism that needs only minerals, light, water and air to grow. Consequently, a large proportion of the genome encodes enzymes that support metabolic processes, such as photosynthesis, respiration, intermediary metabolism, mineral acquisition, and the synthesis of lipids, fatty acids, amino acids, nucleotides and cofactors<sup>103</sup>. With respect to these processes, *Arabidopsis* appears to contain a complement of genes similar to those in the photoautotrophic cyanobacterium *Synechocystis*<sup>45</sup>, but, whereas *Synechocystis* generally has a single gene encoding an enzyme, *Arabidopsis* frequently has many. For example, *Arabidopsis* has at least seven genes for the glycolytic enzyme pyruvate kinase, with an

additional five for pyruvate kinase-like proteins. Whatever the reason for this high level of redundancy, it varies from gene to gene in the same pathway; the 11 enzymes of glycolysis are encoded by up to 51 genes that are present in as few as one or as many as eight copies. Similarly, of the 59 genes encoding proteins involved in glycerolipid metabolism, 39 are represented by more than one gene<sup>104</sup>. Genome duplication and expansion of gene families by tandem duplication have contributed to this diversity.

This high degree of apparent structural redundancy does not necessarily imply functional redundancy. For instance, although there are seven genes for serine hydroxymethyltransferase, a mutation in the gene for the mitochondrial form completely blocks the photorespiratory pathway<sup>105</sup>. Although there are 12 genes for cellulose synthase, mutations in at least 2 of the 12 confer distinct phenotypes because of tissue-specific gene expression<sup>106</sup>.

The metabolome of *Arabidopsis* differs from that of cyanobacteria, or of any other organism sequenced to date, by the presence of many genes encoding enzymes for pathways that are unique to vascular plants. In particular, although relatively little is known about the enzymology of cell-wall metabolism, more than 420 genes could be assigned probable roles in pathways responsible for the synthesis and modification of cell-wall polymers. Twelve genes encode cellulose synthase, and 29 other genes encode 6 families of structurally related enzymes thought to synthesize other major polysaccharides<sup>106</sup>. Roughly 52 genes encode polygalacturonases, 20 encode pectate lyases and 79 encode pectin esterases, indicating a massive investment in modifying pectin. Similarly, the presence of 39  $\beta$ -1,3-glucanases, 20 endoxyloglucan transglycosylases, 50 cellulases and other hydrolases, and 23 expansins reflects the importance of wall remodelling during growth of plant cells. Excluding ascorbate and glutathione peroxidases, there are 69 genes with significant similarity to known peroxidases and 15 laccases (diphenol oxidases). Their presence in such abundance indicates the importance of oxidative processes in the synthesis of lignin, suberin and other cell-wall polymers. The high degree of apparent redundancy in the genes for cell-wall metabolism might reflect differences in substrate specificity by some of the enzymes.

The high degree of apparent redundancy in the genes for cell wall metabolism might reflect differences in substrate specificity by some of the enzymes. It is already known that cell types have different wall compositions, which may require that the relevant enzymes be subject to cell-type-specific transcriptional regulation. Of the 40 or so cell types that plants make, almost all can be identified by unique features of their cell wall<sup>107</sup>. A large number of genes involved in wall metabolism have yet to be defined. Although more than 60 genes for glycosyltransferases can be found in the genome sequence, most of these are probably involved in protein glycosylation or metabolite catabolism and do not seem to be adequate to account for the polysaccharide complexity of the wall. For instance, at least 21 enzymes are required just to produce the linkages of the pectic polysaccharide RGII, and none of these enzymes has been identified at present. Thus, if these and related enzymes involved in the synthesis of other cell-wall polymers are also represented by multiple genes, a substantial number of the genes of currently unknown function may be involved in cell-wall metabolism.

Higher plants collectively synthesize more than 100,000 secondary metabolites. Because flowering plants are thought to have similar numbers of genes, it is apparent that a great deal of enzyme creation took place during the evolution of higher plants. An important factor in the rapid evolution of metabolic complexity is the large family of cytochrome P450s that are evident in *Arabidopsis* (Supplementary Information Table 1). These enzymes represent a superfamily of haem-containing proteins, most of which catalyse NADPH- and O<sub>2</sub>-dependent hydroxylation reactions. Plant P450s participate in myriad biochemical pathways including those devoted to the synthesis of plant products, such as phenylpropanoids, alkaloids, terpenoids, lipids, cyanogenic glycosides and

glucosinolates, and plant growth regulators, such as gibberellins, jasmonic acid and brassinosteroids. Whereas *Arabidopsis* has ~286 P450 genes, *Drosophila* has 94, *C. elegans* has 73 and yeast has only 3. This low number in yeast indicates that there are few reactions of basic metabolism that are catalysed by P450s. It seems likely that many animal P450s are involved in detoxification of compounds from food plant sources. The role of endogenous enzymes is poorly understood; only a few dozen P450 enzymes from plants have been characterized to any extent. The discrepancy between the number of known P450-catalysed reactions and the number of genes suggests that *Arabidopsis* produces a relatively large number of metabolites that have yet to be identified.

In addition to the large number of cytochrome P450s, *Arabidopsis* has many other genes that suggest the existence of pathways or processes that are not currently known. For instance, the presence of 19 genes with similarity to anthranilate *N*-hydroxycinnamoyl/benzoyl transferase is currently inexplicable. This enzyme is involved in the synthesis of dianthramide phytoalexins in Caryophyllaceae and Gramineae. No phytoalexins of this class have been described in *Arabidopsis* as yet. Similarly, the presence of 12 genes with sequence similarity to the berberine bridge enzyme, ((S)-reticuline:oxygen oxidoreductase (methylene-bridge-forming); EC 1.5.3.9), and 13 genes with similarity to tropinone reductase, suggests that *Arabidopsis* may have the ability to produce alkaloids. In other plants, the berberine bridge enzyme transforms reticuline into scoulerine, a biosynthetic precursor to a multitude of species-specific protopine, protoberberine and benzophenanthridine alkaloids. The discovery of these and many other intriguing genes in the *Arabidopsis* genome has created a wealth of new opportunities to understand the metabolic and structural diversity of higher plants.

### Concluding remarks

The twentieth century began with the rediscovery of Mendel's rules of inheritance in pea<sup>108</sup>, and it ends with the elucidation of the complete genetic complement of a model plant, *Arabidopsis*. The analysis of the completed sequence of a flowering plant reported here provides insights into the genetic basis of the similarities and differences of diverse multicellular organisms. It also creates the potential for direct and efficient access to a much deeper understanding of plant development and environmental responses, and permits the structure and dynamics of plant genomes to be assessed and understood.

*Arabidopsis*, *C. elegans* and *Drosophila* have a similar range of 11,000–15,000 different types of proteins, suggesting this is the minimal complexity required by extremely diverse multicellular eukaryotes to execute development and respond to their environment. We account for the larger number of gene copies in *Arabidopsis* compared with these other sequenced eukaryotes with two possible explanations. First, independent amplification of individual genes has generated tandem and dispersed gene families to a greater extent in *Arabidopsis*, and unequal crossing over may be the predominant mechanism involved. Second, ancestral duplication of the entire genome and subsequent rearrangements have resulted in segmental duplications. The pattern of these duplications suggests an ancient polyploidy event, and mutant analysis indicates that at least some of the many duplicate genes are functionally redundant. Their occurrence in a functionally diploid genetic model came as a surprise, and is reminiscent of the situation in maize, an ancient segmental allotetraploid. The remarkable degree of genome plasticity revealed in the large-scale duplications may be needed to provide new functions, as alternative promoters and alternative splicing appear to be less widely used in plants than they are in animals. Apart from duplicated segments, the overall chromosome structure of *Arabidopsis* closely resembles that of *Drosophila*; transposons and other repetitive sequences are concentrated in the heterochromatic regions surrounding the centromere,



whereas the euchromatic arms are largely devoid of repetitive sequences. Conversely, most protein-coding genes reside in the euchromatin, although a number of expressed genes have been identified in centromeric regions. Finally, *Arabidopsis* is the first methylated eukaryotic genome to be sequenced, and will be invaluable in the study of epigenetic inheritance and gene regulation.

Unlike most animals, plants generally do not move, they can perpetuate indefinitely, they reproduce through an extended haploid phase, and they synthesize all their metabolites. Our comparison of *Arabidopsis*, bacterial, fungal and animal genomes starts to define the genetic basis for these differences between plants and other life forms. Basic intracellular processes, such as translation or vesicle trafficking, appear to be conserved across kingdoms, reflecting a common eukaryotic heritage. More elaborate intercellular processes, including physiology and development, use different sets of components. For example, membrane channels, transporters and signalling components are very different in plants and animals, and the large number of transcription factors unique to plants contrasts with the conservation of many chromatin proteins across the three eukaryotic kingdoms. Unexpected differences between seemingly similar processes include the absence of intracellular regulators of cell division (Cdc25) and apoptosis (Bcl-2). On the other hand, DNA repair appears more highly conserved between plants and mammals than within the animal kingdom, perhaps reflecting common factors such as DNA methylation. Our analysis also shows that many genes of the endosymbiotic ancestor of the plastid have been transferred to the nucleus, and the products of this rich prokaryotic heritage contribute to diverse functions such as photoautotrophic growth and signalling.

The sequence reported here changes the fundamental nature of plant genetic analysis. Forward genetics is greatly simplified as mutations are more conveniently isolated molecularly, but at the same time extensive gene duplications mean that functional redundancy must be taken into account. At a biochemical level, the specificity conferred by nucleotide sequence, and the completeness of the survey allow complex mixtures of RNA and protein to be resolved into their individual components using micro-arrays and mass spectrometry. This specificity can also be used in the parallel analysis of genome-wide polymorphisms and quantitative traits in natural populations<sup>109</sup>. Looking ahead, the challenge of determining the function of the large set of predicted genes, many of which are plant-specific, is now a clear priority, and multinational programs have been initiated to accomplish this goal using site-selected mutagenesis among the the necessary tools<sup>110</sup>. Finally, productive paths of crop improvement, based on enhanced knowledge of *Arabidopsis* gene function, will help meet the challenge of sustaining our food supply in the coming years.

*Note added in proof:* at the time of publication 17 centromeric BACs and 5 sequence gaps in chromosome arms are being sequenced. □

## Methods

The three centres used similar annotation approaches involving in silico gene-finding methods, comparison to EST and protein databases, and manual reconciliation of that data. Gene finding involved three steps: (1) analysis of BAC sequences using a computational gene finder; (2) alignment of the sequence to the protein and EST databases; (3) assignment of functions to each of the genes. Genscan<sup>111</sup>, GeneMark.HMM<sup>112</sup>, Xgrail<sup>113</sup> Genefinder (P. Green, unpublished software) and GlimmerA<sup>114</sup> were used to analyse BAC sequences. All of these systems were specially trained for *Arabidopsis* genes. Splice sites were predicted using NetGene2<sup>115</sup>, Splice Predictor<sup>116</sup> and GeneSplicer (M. Perlea and S. Salzberg, unpublished software). For the second step, BACs were aligned to ESTs and to the *Arabidopsis* gene index<sup>117</sup> using programs such as DDS/GAP2<sup>118</sup> or BLASTN<sup>119</sup>. Segmental duplications were analysed and displayed using a modified version of DIALIGN2 (ref. 120).

Received 20 October; accepted 15 November 2000.

1. The *C. elegans* Sequencing Consortium. Sequence and analysis of the genome of *C. elegans*. *Science* **282**, 2012–2018 (1998).
2. Adams, M. D. The genome sequence of *Drosophila melanogaster*. *Science* **287**, 2185–2195 (2000).
3. Meinke, D. W., Cherry, J. M., Dean, C., Rounsley, S. D. & Koornneef, M. *Arabidopsis thaliana*: a model plant for genome analysis. *Science* **282**, 662–665 (1998).

4. Lin, X. *et al.* Sequence and analysis of chromosome 2 of the plant *Arabidopsis thaliana*. *Nature* **402**, 761–768 (1999).
5. Mayer, K. *et al.* Sequence and analysis of chromosome 4 of the plant *Arabidopsis thaliana*. *Nature* **402**, 769–777 (1999).
6. Theologis, A. *et al.* Sequence and analysis of chromosome 1 of the plant *Arabidopsis thaliana*. *Nature* **408**, 816–820 (2000).
7. Salanoubat, M. *et al.* Sequence and analysis of chromosome 3 of the plant *Arabidopsis thaliana*. *Nature* **408**, 820–822 (2000).
8. Tabata, S. *et al.* Sequence and analysis of chromosome 5 of the plant *Arabidopsis thaliana*. *Nature* **408**, 820–822 (2000).
9. Choi, S. D., Creelman, R., Mullet, J. & Wing, R. A. Construction and characterisation of a bacterial artificial chromosome library from *Arabidopsis thaliana*. *Weeds World* **2**, 17–20 (1995).
10. Mozo, T., Fischer, S., Shizuya, H. & Altmann, T. Construction and characterization of the IGF *Arabidopsis* BAC library. *Mol. Gen. Genet.* **258**, 562–570 (1998).
11. Lui, Y. -G., Mitsukawa, N., Vazquez-Tello, A. & Whittier, R. F. Generation of a high-quality P1 library of *Arabidopsis* suitable for chromosome walking. *Plant J.* **7**, 351–358 (1995).
12. Lui, Y. -G. *et al.* Complementation of plant mutants with large genomic DNA fragments by a transformation-competent artificial chromosome vector accelerates positional cloning. *Proc. Natl Acad. Sci. USA* **96**, 6535–6540 (1999).
13. Marra, M. *et al.* A map or sequence analysis of the *Arabidopsis thaliana* genome. *Nature Genet.* **22**, 265–270 (1999).
14. Mozo, T. *et al.* A complete BAC-based physical map of the *Arabidopsis thaliana* genome. *Nature Genet.* **22**, 271–275 (1999).
15. Sato, S. *et al.* Structural analysis of *Arabidopsis thaliana* chromosome 5. I. Sequence features of the 1.6 Mb regions covered by twenty physically assigned P1 clones. *DNA Res.* **4**, 215–230 (1997).
16. Bent, E., Johnson, S. & Bancroft, I. BAC representation of two low-copy regions of the genome of *Arabidopsis thaliana*. *Plant J.* **13**, 849–855 (1998).
17. Copenhaver, G. P. *et al.* Genetic definition and sequence analysis of *Arabidopsis* centromeres. *Science* **286**, 2468–2474 (1999).
18. Meyerowitz, E. M. & Somerville, C. R. *Arabidopsis* (Cold Spring Harbor Laboratory Press, Cold Spring Harbor, New York, 1994).
19. Lowe, T. M. & Eddy, S. R. tRNAscan-SE: A program for improved detection of transfer RNA genes in genomic sequence. *Nucleic Acids Res.* **25**, 955–964 (1997).
20. Pavy, N. *et al.* Evaluation of gene prediction software using a genomic data set: application to *Arabidopsis thaliana* sequences. *Bioinformatics* **15**, 887–900 (1999).
21. Mewes, H. W. *et al.* Overview of the yeast genome. *Nature* **387** (Suppl.) 7–65 (1997).
22. Frishman, D. *et al.* Functional and structural genomics using PEDANT. *Bioinformatics* (in the press).
23. Blattner, F. R. *et al.* The complete genome sequence of *Escherichia coli* K-12. *Science* **277**, 1453–1462 (1997).
24. Kotani, H. & Tabata, S. Lessons from the sequencing of the genome of a unicellular cyanobacterium, *Synechocystis* SP. PCC6803. *Annu. Rev. Plant Physiol. Plant Mol. Biol.* **49**, 151–171 (1998).
25. Apweiler, R. *et al.* INTERPRO (<http://www.ebi.ac.uk/interpro/>). Collaborative Computer Project 11 Newsletter no. 10 (Cambridge, 2000).
26. Bent, A. F. *et al.* RPS2 of *Arabidopsis thaliana* a leucine-rich repeat class of plant disease resistance genes. *Science* **265**, 1856–1860 (1994).
27. Skowryra, D. *et al.* F box proteins are receptors that recruit phosphorylated substrates to the SCF ubiquitin-ligase complex. *Cell* **91**, 209–219 (1997).
28. Joazeiro, C. A. P. & Weissman, A. M. RING finger proteins: mediators of ubiquitin ligase activity. *Cell* **102**, 549–552 (2000).
29. Rubin, G. M. *et al.* Comparative genomics of the eukaryotes. *Science* **287**, 2204–2215 (2000).
30. Delcher, A. L. *et al.* Alignment of whole genomes. *Nucleic Acids Res.* **27**, 2369–2376 (1999).
31. Blanc, G. *et al.* Extensive duplication and reshuffling in the *Arabidopsis* genome. *Plant Cell* **12**, 1093–1102 (2000).
32. Wendel, J. F. Genome evolution in polyploids. *Plant Mol. Biol.* **42**, 225–249 (2000).
33. Gaut, B. S. & Doebley, J. F. DNA sequence evidence for the segmental allotetraploid origin of maize. *Proc. Natl Acad. Sci. USA* **94**, 6809–6814 (1997).
34. Ku, H. -M., Vision, T., Liu, J. & Tanksley, S. D. Comparing sequenced segments of the tomato and *Arabidopsis* genomes: Large-scale duplication followed by selective gene loss creates a network of synteny. *Proc. Natl Acad. Sci. USA* **97**, 9121–9126 (2000).
35. Noel, L. *et al.* Pronounced intraspecific haplotype divergence at the RPP5 complex disease resistance locus of *Arabidopsis*. *Plant Cell* **11**, 2099–2111 (1999).
36. Ellis, J., Dodds, P. & Pryor, T. Structure, function, and evolution of plant disease resistance genes. *Trends Plant Sci.* **3**, 278–284 (2000).
37. Tanksley, S. D. *et al.* High density molecular linkage maps of the tomato and potato genomes. *Genetics* **132**, 1141–1160 (1992).
38. Moore, G., Devos, K. M., Wang, Z. & Gale, M. D. Grasses, line up and form a circle. *Curr. Biol.* **5**, 737–739 (1995).
39. Acarkan, A., Rossberg, M., Koch, M. & Schmidt, R. Comparative genome analysis reveals extensive conservation of genome organisation for *Arabidopsis thaliana* and *Capsella rubella*. *Plant J.* **23**, 55–62 (2000).
40. Cavell, A., Lydiate, D., Parkin, I., Dean, C. & Trick, M. A 30 centimorgan segment of *Arabidopsis thaliana* chromosome 4 has six collinear homologues within the *Brassica napus* genome. *Genome* **41**, 62–69 (1998).
41. O'Neill, C. & Bancroft, I. Comparative physical mapping of segments of the genome of *Brassica oleracea* var *alboglabra* that are homologous to sequenced regions of the chromosomes 4 and 5 of *Arabidopsis thaliana*. *Plant J.* **23**, 233–243 (2000).
42. Wolfe, K. H., Gouy, M., Yang, Y. -W., Sharp, P. M. & Li, W. -H. Date of the monocot-dicot divergence estimated from the chloroplast DNA sequence data. *Proc. Natl Acad. Sci. USA* **86**, 6201–6205 (1989).
43. van Dodeweerd, A. -M. *et al.* Identification and analysis of homologous segments of the genomes of rice and *Arabidopsis thaliana*. *Genome* **42**, 887–892 (1999).
44. Mayer, K. Sequence level analysis of homologous segments of the genomes of rice and *Arabidopsis thaliana*. *Genome Res.* (submitted).
45. Sato, S. Complete structure of the chloroplast genome of *Arabidopsis thaliana*. *DNA Research* **6**, 283–290 (1999).
46. Unsel, M., Marienfeld, J., Brandt, P. & Brennicke, A. The mitochondrial genome in *Arabidopsis*

- thaliana* contains 57 genes in 366,924 nucleotides. *Nature Genet.* **15**, 57–61 (1997).
47. Palmer, J. D. et al. Dynamic evolution of plant mitochondrial genomes: mobile genes and introns and highly variable mutation rates. *Proc. Natl Acad. Sci. USA* **97**, 6960–6966 (2000).
  48. Le, Q. -H. et al. Transposon diversity in *Arabidopsis thaliana*. *Proc. Natl Acad. Sci. USA* **97**, 7376–7381 (2000).
  49. Yu, Z., Wright, S. & Bureau, T. *Mutator*-like elements (MULEs) in *Arabidopsis thaliana*: Structure, diversity and evolution. *Genetics* (in the press).
  50. Feschotte, C. & Mouches, C. Evidence that a family of miniature inverted-repeat transposable elements (MITEs) from the *Arabidopsis thaliana* genome has arisen from a pogo-like DNA transposon. *Mol. Biol. Evol.* **17**, 730–737 (2000).
  51. Martienssen, R. Transposons, DNA methylation and gene control. *Trends Genet.* **14**, 263–264 (1998).
  52. Singer, T., Yordan, C. & Martienssen, R. Robertson's *Mutator* transposons in *Arabidopsis* are regulated by the chromatin-remodeling gene *Decrease in DNA Methylation (DDM1)*. *Genes Dev.* (in the press).
  53. SanMiguel, P. et al. Nested retrotransposons in the intergenic regions of the maize genome. *Science* **274**, 765–768 (1996).
  54. Copenhaver, G. P. & Pikaard, C. S. Two-dimensional RFLP analyses reveal megabase-sized clusters of rRNA gene variants in *Arabidopsis thaliana*, suggesting local spreading of variants as the mode for gene homogenization during concerted evolution. *Plant J.* **9**, 273–282 (1996).
  55. Franz, P. et al. Cytogenetics for the model system *Arabidopsis thaliana*. *Plant J.* **13**, 867–876 (1998).
  56. Richards, E. J. & Ausubel, F. M. Isolation of a higher eukarotic telomere from *Arabidopsis thaliana*. *Cell* **53**, 127–136 (1988).
  57. Round, E. K., Flowers, S. K. & Richards, E. J. *Arabidopsis thaliana* centromere regions: genetic map positions and repetitive DNA structure. *Genome Res.* **7**, 1045–1053 (1997).
  58. The CSHL/WUGSC/PEB *Arabidopsis* Sequencing Consortium. The complete sequence of a heterochromatic island from a higher eukaryote. *Cell* **100**, 377–386 (2000).
  59. Franz, P. F. et al. Integrated cytogenetic map of chromosome arm 4S of *A. thaliana*: Structural organization of heterochromatic knob and centromere region. *Cell* **100**, 367–376 (2000).
  60. Paulsen, I. T., Nguyen, L., Sliwinski, M. K., Rabus, R. & Saier, M. H. Jr Microbial genome analyses: comparative transport capabilities in eighteen prokaryotes. *J. Mol. Biol.* **301**, 75–101 (2000).
  61. Paulsen, I. T., Sliwinski, M. K., Nelissen, B., Goffeau, A. & Saier, M. H. Jr Unified inventory of established and putative transporters encoded within the complete genome of *Saccharomyces cerevisiae*. *FEBS Lett.* **430**, 116–125 (1998).
  62. Hirsch, R. E., Lewis, B. D., Spalding, E. P. & Sussman, M. R. A role for the AKT1 potassium channel in plant nutrition. *Science* **280**, 918–921 (1998).
  63. Slayman, C. L. & Slayman, C. W. Depolarization of the plasma membrane of *Neurospora* during active transport of glucose: evidence for a proton-dependent cotransport system. *Proc. Natl Acad. Sci. USA* **71**, 1035–1039 (1974).
  64. Ryan, C. A. & Pearce, G. Systemin: a polypeptide signal for plant defensive genes. *Annu. Rev. Cell. Dev. Biol.* **14**, 1–17 (1998).
  65. Eisen, J. A. & Hanawalt, P. C. A phylogenomic study of DNA repair genes, proteins, and processes. *Mutat. Res.* **435**, 171–213 (1999).
  66. Selby, C. P. & Sancar, A. Structure and function of transcription-repair coupling factor. Structural domains and binding properties. *J. Biol. Chem.* **270**, 4882–4889 (1995).
  67. Britt, A. B. Molecular genetics of DNA repair in higher plants. *Trends Plant Sci.* **4**, 20–25 (1999).
  68. Dangel, M. Response to Aravind, L. & Koonin, E. V. Second Family of Histone Deacetylases. *Science* **280**, 1167 (1998).
  69. Cao, X. et al. Conserved plant genes with similarity to mammalian de novo DNA methyltransferases. *Proc. Natl Acad. Sci. USA* **97**, 4979–4984 (2000).
  70. Henikoff, S. & Comai, L. A DNA methyltransferase homologue with a chromodomain exists in multiple polymorphic forms in *Arabidopsis*. *Genetics* **149**, 307–318 (1998).
  71. Hung, M. -S. et al. *Drosophila* proteins related to vertebrate DNA (5-cytosine) methyltransferases. *Proc. Natl Acad. Sci. USA* **96**, 11940–11945 (1999).
  72. Dalmay, T., Hamilton, A. J., Rudd, S., Angell, S. & Baulcombe, D. C. An RNA-dependent-RNA polymerase in *Arabidopsis* is required for post transcriptional gene silencing mediated by a transgene but not by a virus—the truth. *Cell* **101**, 543–553 (2000).
  73. Riechmann, J. L. & Ratcliffe, O. J. A genomic perspective on plant transcription factors. *Curr. Opin. Plant Biol.* **3**, 423–434 (2000).
  74. Liljegren, S. J. et al. SHATTERPROOF MADS-box genes control seed dispersal in *Arabidopsis*. *Nature* **404**, 766–770 (2000).
  75. Pelaz, S. et al. B and C floral organ identity functions require SEPALLATA MADS-box genes. *Nature* **405**, 200–203 (2000).
  76. Canaday, J., Stoppin-Mellet, V., Mutterer, J., Lambert, A. M. & Schmit, A. C. Higher plant cells: gamma-tubulin and microtubule knob in the absence of centrosomes. *Gamma. Res. Technol.* **49**, 487–495 (2000).
  77. Bassham, D. C. & Raikhel, N. V. Unique features of the plant vacuolar sorting machinery. *Curr. Opin. Cell Biol.* **12**, 491–495 (2000).
  78. Zheng, Z. L. & Yang, Z. The Rop GTPase switch turns on polar growth in pollen. *Trends Plant Sci.* **5**, 298–303 (2000).
  79. den Boer, B. G. & Murray, J. A. Triggering the cell cycle in plants. *Trends Cell Biol.* **10**, 245–250 (2000).
  80. Heese, M., Mayer, U. & Jurgens, G. Cytokinesis in flowering plants: cellular process and developmental integration. *Curr. Opin. Plant Biol.* **1**, 486–491 (1998).
  81. Meyerowitz, E. M. Plants, animals, and the logic of development. *Trends Genet.* **15**, M65–M68 (1999).
  82. Wang, D. Y. C. et al. Divergence time estimates for the early history of animal phyla and the origin of plants, animals and fungi. *Proc. R. Soc. Lond. B Biol.* **266**, 63–171 (1999).
  83. Torii, K. Receptor kinase activation and signal transduction in plants: an emerging picture. *Curr. Opin. Plant Biol.* **3**, 362–367 (2000).
  84. Yeh, K. C. & Lagarias, J. C. Eukaryotic phytochromes: Light-regulated serine/threonine protein kinases with histidine kinase ancestry. *Proc. Natl Acad. Sci. USA* **95**, 13976–13981 (1998).
  85. McCarty, D. R. & Chory, J. Conservation and innovation in plant signaling pathways. *Cell* **103**, 201–211 (2000).
  86. Weiss, C. A., Garnaat, C., Mukai, K., Hu, Y. & Ma, H. Molecular cloning of cDNAs from maize and *Arabidopsis* encoding a G protein beta subunit. *Proc. Natl Acad. Sci. USA* **91**, 9554–9558 (1994).
  87. Bowler, C. et al. Cyclic GMP and calcium mediate phytochrome phototransduction. *Cell* **77**, 73–81 (1994).
  88. Stepanova, A. & Ecker, J. R. Ethylene signaling: from mutants to molecules. *Curr. Opin. Plant Biol.* **3**, 353–360 (2000).
  89. Urao, T., Yamaguchi-Shinozaki, K. & Shinozaki, K. Two-component systems in plant signal transduction. *Trends Plant Sci.* **5**, 67–74 (2000).
  90. Makino, S. et al. Genes encoding pseudo-response regulators: Insight into His-to-Asp phosphorylation and circadian rhythm in *Arabidopsis thaliana*. *Plant Cell Physiol.* **41**, 791–803 (2000).
  91. D'Agostino, I. B. & Kieber, J. J. Phosphorelay signal transduction: the emerging family of plant response regulators. *Trends Biol. Sci.* **24**, 452–456 (1999).
  92. Strayer, C. et al. Cloning of the *Arabidopsis* clock gene TOC1, an autoregulatory response regulator homologue. *Science* **289**, 768–771 (2000).
  93. Stahl, E. A. & Bishop, J. G. Plant-Pathogen arms races at the molecular level. *Curr. Opin. Plant Biol.* **3**, 299–304 (2000).
  94. McDowell, J. M. & Dangl, J. L. Signal transduction in the plant innate immune response. *Trends Biochem. Sci.* **25**, 79–82 (2000).
  95. Van der Biezen, E. A. & Jones, J. D. Plant disease-resistance proteins and the gene-for-gene concept. *Trends Biochem. Sci.* **23**, 454–456 (1998).
  96. Belvin, M. P. & Anderson, K. V. A conserved signaling pathway: the *Drosophila* toll-dorsal pathway. *Annu. Rev. Cell. Dev. Biol.* **12**, 393–416 (1996).
  97. Uren, A. G. et al. Identification of paracaspases and metacaspases: Two ancient families of caspase-like proteins, one of which plays a key role in MALT lymphoma. *Mol. Cell* **6**, 961–967 (2000).
  98. Fankhauser, C. & Chory, J. Light control of plant development. *Annu. Rev. Cell. Dev. Biol.* **13**, 203–229 (1997).
  99. Briggs, W. R. & Hua, E. Blue-light photoreceptors in higher plants. *Annu. Rev. Cell. Dev. Biol.* **15**, 33–62 (1999).
  100. Christie, J. M., Salomon, M., Nozue, K., Wada, M. & Briggs, W. R. LOV (light, oxygen, or voltage) domains of the blue-light photoreceptor phototropin (nph1): binding sites for the chromophore flavin mononucleotide. *Proc. Natl Acad. Sci. USA* **96**, 8779–8783 (1999).
  101. Golbeck, J. H. Structure and function of photosystem I. *Annu. Rev. Plant Physiol. Plant Mol. Biol.* **43**, 293–324 (1992).
  102. Maier, R. M., Neckermann, K., Igloi, G. L. & Kossel, H. Complete sequence of the maize chloroplast genome: gene content, hotspots of divergence and fine tuning of genetic information by transcript editing. *J. Mol. Biol.* **251**, 614–28 (1995).
  103. Buchanan, B. B., Grusissem, W. & Jones, R. L. in *Biochemistry and Molecular Biology of Plants* 1367 (Am. Soc. Plant Physiol., Rockville, Maryland, 2000).
  104. Mekhedov, S., Martinez de Ilarduya, O. & Ohlrogge, J. Toward a functional catalog of the plant genome. A survey of genes for lipid biosynthesis. *Plant Physiol.* **122**, 389–401 (2000).
  105. Somerville, C. R. & Ogren, W. L. Photorespiration deficient mutants of *Arabidopsis thaliana* lacking mitochondrial serine transhydroxymethylase activity. *Plant Physiol.* **67**, 666–671 (1981).
  106. Richmond, T. & Somerville, C. R. The cellulose synthase superfamily. *Plant Physiol* **124**, 495–499 (1999).
  107. Carpita, N. Vergara C. A recipe for cellulose. *Science* **279**, 672–673 (1998).
  108. De Vries, H. Sur la loi de disjonction des hybrides. *C. R. Acad. Sci. Paris* **130**, 845–847 (1900).
  109. Alonso-Blanco, C. & Koornneef, M. Naturally occurring variation in *Arabidopsis*: an underexploited resource for plant genetics. *Trends Plant Sci.* **5**, 1360–1385 (1999).
  110. Chory, J. Functional genomics and the virtual plant. A blueprint for understanding how plants are built and how to improve them. *Plant Physiology* **123**, 423–425 (2000).
  111. Burge, C. & Karlin, S. Prediction of complete gene structures in human genomic DNA. *J. Mol. Biol.* **268**, 78–94 (1997).
  112. Lukashin, A. V. & Borodovsky, M. GeneMark.hmm: new solutions for gene finding. *Nucleic Acids Res.* **26**, 1107–1115 (1998).
  113. Uberbacher, E. C. & Mural, R. J. Locating protein-coding regions in human DNA sequences by a multiple sensor-neural network approach. *Proc. Natl Acad. Sci. USA* **88**, 11261–11265 (1991).
  114. Salzberg, S. L., Pertea, M., Delcher, A. L., Gardiner, M. J. & Tettelin, H. Interpolated Markov models for eukaryotic gene finding. *Genomics* **59**, 24–31 (1999).
  115. Hebsgaard, S. M. et al. Splice site prediction in *Arabidopsis thaliana* DNA by combining local and global sequence information. *Nucleic Acids Res.* **24**, 3439–3452 (1996).
  116. Brendel, V. & Kleffe, J. Prediction of locally optimal splice sites in plant pre-mRNA with applications to gene identification in *Arabidopsis thaliana* genomic DNA. *Nucleic Acids Res.* **26**, 4748–4757 (1998).
  117. Quackenbush, J., Liang, F., Holt, I., Pertea, G. & Upton, J. The TIGR gene indices: reconstruction and representation of expressed gene sequences. *Nucleic Acids Res.* **28**, 141–145 (2000).
  118. Huang, X., Adams, M. D., Zhou, H. & Kerlavage, A. R. A tool for analyzing and annotating genomic sequences. *Genomics* **46**, 37–45 (1997).
  119. Altschul, S. F. et al. Basic local alignment search tool. *J. Mol. Biol.* **215**, 403–410 (1990).
  120. Morgenstern, B. DIALIGN2: improvement of the segment-to-segment approach to multiple sequence alignment. *Bioinformatics* **15**, 211–218 (1999).
  121. Murzin, A. G., Brenner, S. E., Hubbard, T. & Chothia, C. SCOP: a structural classification of proteins database for the investigation of sequences and structures. *J. Mol. Biol.* **247**, 536–540 (1995).
  122. Emanuelsson, O., Nielsen, H., Brunak, S. & von Heijne, G. Predicting subcellular localization of proteins based on their N-terminal amino acid sequence. *J. Mol. Biol.* **300**, 1005–1016 (2000).

Supplementary information is available on Nature's World-Wide Web site (<http://www.nature.com>) or as paper copy from the London editorial office of Nature.

## Acknowledgements

This work was supported by the National Science Foundation (NSF) Cooperative Agreements (funded by the NSF, the US Department of Agriculture (USDA) and the US Department of Energy (DOE)), the Kazusa DNA Research Institute Foundation, and by the European Commission. Additional support from the USDA, Ministère de la Recherche, GSF-Forschungszentrum f. Umwelt u. Gesundheit, BMBF (Bundesministerium f. Bildung, Forschung und Technologie), the BBSRC (Biotechnology and Biological

Sciences Research Council) and the Plant Research International, Wageningen, is also gratefully acknowledged. The authors wish to thank E. Magnien, D. Nasser and J. D. Watson for their continual support and encouragement.

Correspondence and requests for materials should be addressed to The Arabidopsis Genome Initiative (e-mail: genomeanalysis@tigr.org or genomeanalysis@gsf.de).

### Genome Sequencing Groups

**Samir Kaul, Hean L. Koo, Jennifer Jenkins, Michael Rizzo, Timothy Rooney, Luke J. Tallon, Tamara Feldblyum, William Nierman, Maria-Ines Benito, Xiaoying Lin, Christopher D. Town, J. Craig Venter & Claire M. Fraser**

*The Institute for Genomic Research, 9712 Medical Centre Drive, Rockville, Maryland 20850, USA*

**Satoshi Tabata, Yasukazu Nakamura, Takakazu Kaneko, Shusei Sato, Erika Asamizu, Tomohiko Kato, Hirokazu Kotani & Shigemi Sasamoto**

*Kazusa DNA Research Institute, 1532-3 Yana, Kisarazu, Chiba 292, Japan*

**Joseph R. Ecker<sup>1\*†</sup>, Athanasios Theologis<sup>2\*</sup>, Nancy A. Federspiel<sup>3\*†</sup>, Curtis J. Palm<sup>3</sup>, Brian I. Osborne<sup>2</sup>, Paul Shinn<sup>1</sup>, Aaron B. Conway<sup>3</sup>, Valentina S. Vysotskaia<sup>2</sup>, Ken Dewar<sup>1</sup>, Lane Conn<sup>3</sup>, Catherine A. Lenz<sup>2</sup>, Christopher J. Kim<sup>1</sup>, Nancy F. Hansen<sup>3</sup>, Shirley X. Liu<sup>2</sup>, Eugen Buehler<sup>1</sup>, Hootan Altafi<sup>3</sup>, Hitomi Sakano<sup>2</sup>, Patrick Dunn<sup>1</sup>, Bao Lam<sup>3</sup>, Paul K. Pham<sup>2</sup>, Qimin Chao<sup>1</sup>, Michelle Nguyen<sup>3</sup>, Guixia Yu<sup>2</sup>, Huaming Chen<sup>1</sup>, Audrey Southwick<sup>3</sup>, Jeong Mi Lee<sup>2</sup>, Molly Miranda<sup>3</sup>, Mitsue J. Toriumi<sup>2</sup> & Ronald W. Davis<sup>3</sup>**

*1, Plant Science Institute, Department of Biology, University of Pennsylvania, Philadelphia, Pennsylvania 19104 USA; 2, Plant Gene Expression Center/USDA-U.C.Berkeley, 800 Buchanan Street, Albany, California 94710, USA; 3, Stanford Genome Technology Center, 855 California Avenue, Palo Alto, California 94304, USA. \* These authors contributed equally to this work. † Present addresses: The Salk Institute for Biological Studies, 10010 North Torrey Pines Road, La Jolla, California 92037, USA (J.R.E.); Exelixis, Inc., 170 Harborway, P.O. Box 511, South San Francisco, California 94083-0511, USA (N.A.F)*

**European Union Chromosome 4 and 5 Sequencing Consortium: R. Wambutt<sup>1</sup>, G. Murphy<sup>2</sup>, A. Düsterhöft<sup>3</sup>, W. Stiekema<sup>4</sup>, T. Pohl<sup>5</sup>, K.-D. Entian<sup>6</sup>, N. Terryn<sup>7</sup> & G. Volckaert<sup>8</sup>**

*1, AGOWA GmbH, Glienicke Weg 185, D-12489 Berlin, Germany; 2, John Innes Centre, Colney Lane, Norwich NR4 7UH, UK; 3, QIAGEN GmbH, Max-Volmer-Str. 4, D-40724 Hilden, Germany; 4, Greenomics, Plant Research International, Droevendaalseleeg 1, NL 6700, AA Wageningen, The Netherlands; 5, GATC GmbH, Fritz-Arnold-Strasse 23, D-78467 Konstanz, Germany; 6, SRD GmbH, Oberurseler Str. 43, Oberursel 61440, Germany; 7, Department for Plant Genetics, (VIB), University of Gent, K.L. Ledeganckstraat 35, B-9000 Gent, Belgium; 8, Katholieke Universiteit Leuven, Laboratory of Gene Technology, Kardinaal Mercierlaan 92, B-3001 Leuven, Belgium*

**European Union Chromosome 3 Sequencing Consortium: M. Salanoubat<sup>1</sup>, N. Choise<sup>1</sup>, M. Rieger<sup>2</sup>, W. Ansoerge<sup>3</sup>, M. Unseld<sup>4</sup>, B. Fartmann<sup>5</sup>, G. Valle<sup>6</sup>, F. Artiguenave<sup>1</sup>, J. Weissenbach<sup>1</sup> & F. Quetier<sup>1</sup>**

*1, Genoscope and CNRS FRE2231, 2 rue G. Crémieux, 91057 Evry Cedex, France; 2, Genotype GmbH Angelhofweg 39, D-69259 Wilhelmsfeld, Germany; 3, European Molecular Biology Laboratory, Biochemical Instrumentation Program, Meyerhofstr. 1, D-69117 Heidelberg, Germany; 4, LION Bioscience AG, Im Neuenheimer Feld 515-517, 69120 Heidelberg, Germany; 5, MWG-Biotech AG, Anzinger Strasse 7a, 85560 Ebersberg, Germany; 6, CRIBI, Università di Padova, via G. Colombo 3, Padova 35131, Italy*

**The Cold Spring Harbor and Washington University Genome Sequencing Center Consortium: Richard K. Wilson<sup>1</sup>, Melissa de la Bastide<sup>2</sup>, M. Sekhon<sup>1</sup>, Emily Huang<sup>2</sup>, Lori Spiegel<sup>2</sup>, Lidia Gnoj<sup>2</sup>, K. Pepin<sup>1</sup>, J. Murray<sup>1</sup>, D. Johnson<sup>1</sup>, Kristina Habermann<sup>2</sup>, Neilay Dedhia<sup>2</sup>, Larry Parnell<sup>2</sup>, Raymond Preston<sup>2</sup>, L. Hillier<sup>1</sup>, Ellison Chen<sup>3</sup>, M. Marra<sup>2</sup>, Robert Martienssen<sup>4</sup> & W. Richard McCombie<sup>2</sup>**

*1, Washington University Genome Sequencing Center, Washington University in St Louis School of Medicine, 4444 Forest Park Blvd., St. Louis, Missouri 63108 USA; 2, Lita Annenberg Hazen Genome Center, Cold Spring Harbor Laboratory, Cold Spring Harbor, New York 11724, USA; 3, Celera Genomics, 850 Lincoln Center Drive, Foster City, California 94494, USA; 4, Plant Biology Group, Cold Spring Harbor Laboratory, Cold Spring Harbor, New York 11724, USA*

### Genome Analysis Group

**Klaus Mayer<sup>1\*</sup>, Owen White<sup>2\*</sup>, Michael Bevan<sup>3</sup>, Kai Lemcke<sup>1</sup>, Todd H. Creasy<sup>2</sup>, Cord Bielke<sup>2</sup>, Brian Haas<sup>1</sup>, Dirk Haase<sup>1</sup>, Rama Maiti<sup>2</sup>, Stephen Rudd<sup>1</sup>, Jeremy Peterson<sup>2</sup>, Heiko Schoof<sup>1</sup>, Dimitrij Frishman<sup>1</sup>, Burkhard Morgenstern<sup>1</sup>, Paulo Zaccaria<sup>1</sup>, Maria Ermolaeva<sup>2</sup>, Mihaela Perteau<sup>2</sup>, John Quackenbush<sup>2</sup>, Natalia Volfovskiy<sup>2</sup>, Dongying Wu<sup>2</sup>, Todd M. Lowe<sup>4</sup>, Steven L. Salzberg<sup>2</sup> & Hans-Werner Mewes<sup>1</sup>**

*1, GSF-Forschungszentrum f. Umwelt u. Gesundheit, Munich Information Center for Protein Sequences, am Max-Planck-Institut f. Biochemie, Am Klopferspitz 18a, D-82152, Germany; 2, The Institute for Genomic Research, 9712 Medical Center Drive, Rockville, Maryland 20850, USA; 3, Molecular Genetics Department, John Innes Centre, Colney Lane, Norwich NR4 7UH, UK; 4, Dept Genetics, Stanford University Medical School, Stanford, California 94305-5120, USA. \* These authors contributed equally to this work*

### Contributing Authors

**Comparative analysis of the genomes of *A. thaliana* accessions. S. Rounsley, D. Bush, S. Subramaniam, I. Levin & S. Norris**

*Cereon Genomics LLC, 45 Sidney St, Cambridge, Massachusetts 02139, USA*

**Comparative analysis of the genomes of *A. thaliana* and other genera. R. Schmidt<sup>1</sup>, A. Acarkan<sup>1</sup> & I. Bancroft<sup>2</sup>**

*1, Max-Delbrück-Laboratorium in der Max-Planck-Gesellschaft, Carl-von-Linné-Weg 10, 50829 Cologne, Germany; 2, Brassicas and Oilseeds Research Department, John Innes Centre, Norwich NR4 7UJ, UK*

**Integration of the three genomes in the plant cell: the extent of protein and nucleic acid traffic between nucleus, plastids and mitochondria. F. Quetier<sup>1</sup>, A. Brennicke<sup>2</sup> & J. A. Eisen<sup>3</sup>.**

1, Genoscope, Centre Nationale de Sequencage, 2 rue Gaston Cremieux, CP 5706, 91057 Evry Cedex, France; 2, Molekulare Botanik, Universität Ulm, 89069 Ulm, Germany; 3, The Institute for Genomic Research, 9712 Medical Centre Drive, Rockville, Maryland 20850, USA

**Transposable elements. T. Bureau<sup>1</sup>, B.-A. Legault<sup>1</sup>, Q.-H. Le<sup>1</sup>, N. Agrawal<sup>1</sup>, Z. Yu<sup>1</sup> & R. Martienssen<sup>2</sup>**

1, McGill University, Dept of Biology, 1205 rue Dr Penfield, Montreal, Quebec, H3A 1B1, Canada; 2, Plant Biology Group, Cold Spring Harbor Laboratory, Cold Spring Harbor, New York 11724, USA

**rDNA, telomeres and centromeres. G. P. Copenhaver<sup>1</sup>, S. Luo<sup>1</sup>, C. S. Pikaard<sup>2</sup> & D. Preuss<sup>1</sup>**

1, Howard Hughes Medical Institute, The University of Chicago, 1103 East 57th Street, Chicago, Illinois, USA; 2, Biology Department, Washington University in St Louis, St Louis, Missouri 63130, USA

**Membrane transport. I. T. Paulsen<sup>1</sup> & M. Sussman<sup>2</sup>**

1, The Institute for Genomic Research, 9712 Medical Centre Drive, Rockville, Maryland 20850, USA; 2, University of Wisconsin Biotechnology Center, 425 Henry Mall, Madison, Wisconsin 53706, USA

**DNA repair and recombination. A. B. Britt<sup>1</sup> & J. A. Eisen<sup>2</sup>**

1, Section of Plant Biology, University of California, Davis, California 95616, USA; 2, The Institute for Genomic Research, 9712 Medical Centre Drive, Rockville, Maryland 20850, USA

**Gene regulation. D. A. Selinger<sup>1</sup>, R. Pandey<sup>1</sup>, D. W. Mount<sup>2</sup>, V. L. Chandler<sup>1</sup>, R. A. Jorgensen<sup>1</sup> & C. Pikaard<sup>3</sup>**

1, Department of Plant Sciences, University of Arizona, 303 Forbes Hall; and 2, Department of Molecular and Cellular Biology, University of Arizona, Tucson, Arizona 85721, USA; 3, Biology Department, Washington University in St Louis, St Louis, Missouri 63130, USA

**Cellular organization. G. Juergens**

Entwicklungsgenetik, ZMBP-Centre for Plant Molecular Biology, auf der Morgenstelle 1, Tuebingen D-72076, Germany

**Development. E. M. Meyerowitz.**

Division of Biology, California Institute of Biology, Pasadena, California 91125, USA

**Signal transduction. J. R. Ecker<sup>1</sup> & A. Theologis<sup>2</sup>.**

1, The Salk Institute for Biological Studies, 10010 North Torrey Pines Road, La Jolla, California 92037, USA; 2, Plant Gene Expression Center/USDA-UC Berkeley, 800 Buchanan Street, Albany, California 94710, USA

**Recognition of and response to pathogens. J. Dangi<sup>1</sup> & J. D. G. Jones<sup>2</sup>**

1, Biology Department, Coker Hall, University of North Carolina, Chapel Hill, North Carolina 27599, USA; 2, Sainsbury Laboratory, John Innes Centre, Colney Lane, Norwich NR4 7UJ, UK

**Photomorphogenesis and photosynthesis. M. Chen & J. Chory**

Howard Hughes Medical Institute and Plant Biology Laboratory, The Salk Institute, 10010 North Torrey Pines Road, La Jolla, California 92037, USA

**Metabolism. C. Somerville**

Carnegie Institution, 260 Panama Street, Stanford, California 94305, USA