

A Gene Expression Map of the *Arabidopsis* Root

Kenneth Birnbaum,¹ Dennis E. Shasha,² Jean Y. Wang,³
Jee W. Jung,¹ Georgina M. Lambert,⁴ David W. Galbraith,⁴
Philip N. Benfey^{3*}

A global map of gene expression within an organ can identify genes with coordinated expression in localized domains, thereby relating gene activity to cell fate and tissue specialization. Here, we present localization of expression of more than 22,000 genes in the *Arabidopsis* root. Gene expression was mapped to 15 different zones of the root that correspond to cell types and tissues at progressive developmental stages. Patterns of gene expression traverse traditional anatomical boundaries and show cassettes of hormonal response. Chromosomal clustering defined some coregulated genes. This expression map correlates groups of genes to specific cell fates and should serve to guide reverse genetics.

In multicellular eukaryotes, development is a process regulated by differential gene expression whereby cells acquire specific fates. Thus, cell-type specific gene expression profiles will help reveal the determinants of cell fate. Expression analysis of single cell types alone from complex organs in both plants (1–3) and animals (4, 5) do not allow one to infer the relative expression levels among cell types. Gene expression patterns will be more informative when all or most cell types within an organ are analyzed.

We have developed a set of techniques to generate high-resolution spatial and temporal expression profiles throughout the *Arabidopsis* root. The method measures gene expression among cell types and tissues and along a developmental gradient, resulting in a digital readout with resolution approaching that of in situ hybridization (thus, we call the output a digital in situ). We chose the root to test this protocol because of its relatively simple radial organization and its mode of continuous development from a set of stem cells (6). An additional advantage of the root was the availability of well-characterized transgenic lines expressing green fluorescent protein (GFP) in specific cell populations.

For cell-type and tissue-specific expression, the roots of plants expressing GFP in specific cell types were dissociated into single cells by enzymatic digestion of their cell walls (protoplasting) [for example, (7)]. The GFP-expressing cells were isolated with the use of a fluorescence-activated cell sorter, and their mRNA was analyzed with the use of

microarrays. We used five separate GFP lines [expressing in stele, endodermis, endodermis plus cortex, epidermal atrichoblast cells, and lateral root cap (fig. S1)] that together captured all but two cell types produced from the primary meristem of the root (Fig. 1A). High-throughput techniques allowed the harvesting, protoplasting, and sorting of about 10 million cells in about 1.5 hours. Within this time period, specialized cells do not appear to undergo substantial changes in their transcriptional identity (8) (see below).

To obtain gene expression data along a temporal axis, we took advantage of the fact that the developmental stages of root cells are roughly correlated with distance from the apical meristem. Roots were dissected at stage-specific cellular landmarks along this longitudinal axis (Fig. 1A), and RNA from developmental stages was hybridized separately to microarrays. All samples were hybridized to the ATH1 GeneChip (Affymetrix, Santa Clara, CA), which contains probes for more than 22,000 *Arabidopsis* genes, covering about 90% of the genome. The three developmental stages analyzed corresponded to the same region of the root that was efficiently protoplasted. Thus, for each gene, expression signals from the three developmental stages (hereafter stages) were used to apportion expression of the cell type and tissue profiles (hereafter zones) forming 15 separate subzones (five cell types by three stages). For example, a gene with an expression signal of 200 in the stele that has 50% of its stage-specific expression in stage 1 (promeristem) would be given an expression signal of 100 for stele stage 1, and so on (Fig. 1A). The continuous growth of the root meristem allowed us to collect different cell types at essentially synchronized stages of development, making the root an ideal model for the digital reconstruction technique. Other organs can be profiled along a developmental

gradient with the use of stage-specific GFP marker lines or time series sampling for tissues with synchronized, determinate growth.

To determine the extent to which the protoplasting treatment altered gene expression, we compared expression in protoplasted roots with untreated roots. Roots were grown in the high-throughput conditions described above and split into two pools, one that was processed immediately for RNA extraction and a second that was treated with protoplasting enzymes and collected in increments over 1 hour to simulate the conditions of sorted cells. Labeled RNA from treated and untreated roots was then hybridized to microarrays, and expression for each gene was compared between the two treatments. Expression of individual transcripts was highly correlated in the two treatments in all three replicates ($r = 0.90$, $r = 0.92$, and $r = 0.92$), showing that the rapid protoplasting treatment did not dramatically change global gene expression profiles. We did, however, identify several hundred transcripts that were consistently induced by protoplasting treatment (table S1). These were removed from further analysis.

To determine whether the relative abundance of individual RNA species as assayed by our microarray analysis mirrored the amounts found in the sorted cells, we used real-time reverse transcription polymerase chain reaction (RT-PCR) on RNA from the sorted cells on four transcripts with putative zone-specific expression. In each case, relative amounts of cell-type enrichment were closely replicated by the quantitative RT-PCR analysis (fig. S2).

We validated the digital in situ method by comparing its results for individual transcripts to previously documented expression patterns derived from promoter reporter constructs or in situ hybridization (e.g., Fig. 1, B to D). Root expression patterns generated by digital in situ matched documented expression patterns in 25 out of 26 cases (fig. S3). In an analysis of zone and stage data separately, only one transcript failed to match documented patterns. We then used the documented expression patterns as a training set to determine an expression level that represented a minimum transcript presence signal above noise (set at 75). At this threshold, the combined data had 19 false positives and 4 false negatives out of 390 measurements (15 subzones in 26 cases). To further test the method, we confirmed digital in situ patterns of two uncharacterized genes with in situ hybridization. Both matched the predicted digital patterns (e.g., Fig. 1D). Thus, the overall error rate was about 5.5%. Together the validation experiments showed that the digital in situ recreated root expression patterns at a high level of detail and accuracy (the full data set is available in table S2).

¹Department of Biology, ²Courant Institute of Mathematical Sciences, New York University, New York, NY 10003, USA. ³Department of Biology, Duke University, Box 91000, Durham, NC 27708, USA. ⁴Department of Plant Sciences, University of Arizona, Tucson, AZ 85721, USA.

*To whom correspondence should be addressed. E-mail: philip.benfey@duke.edu

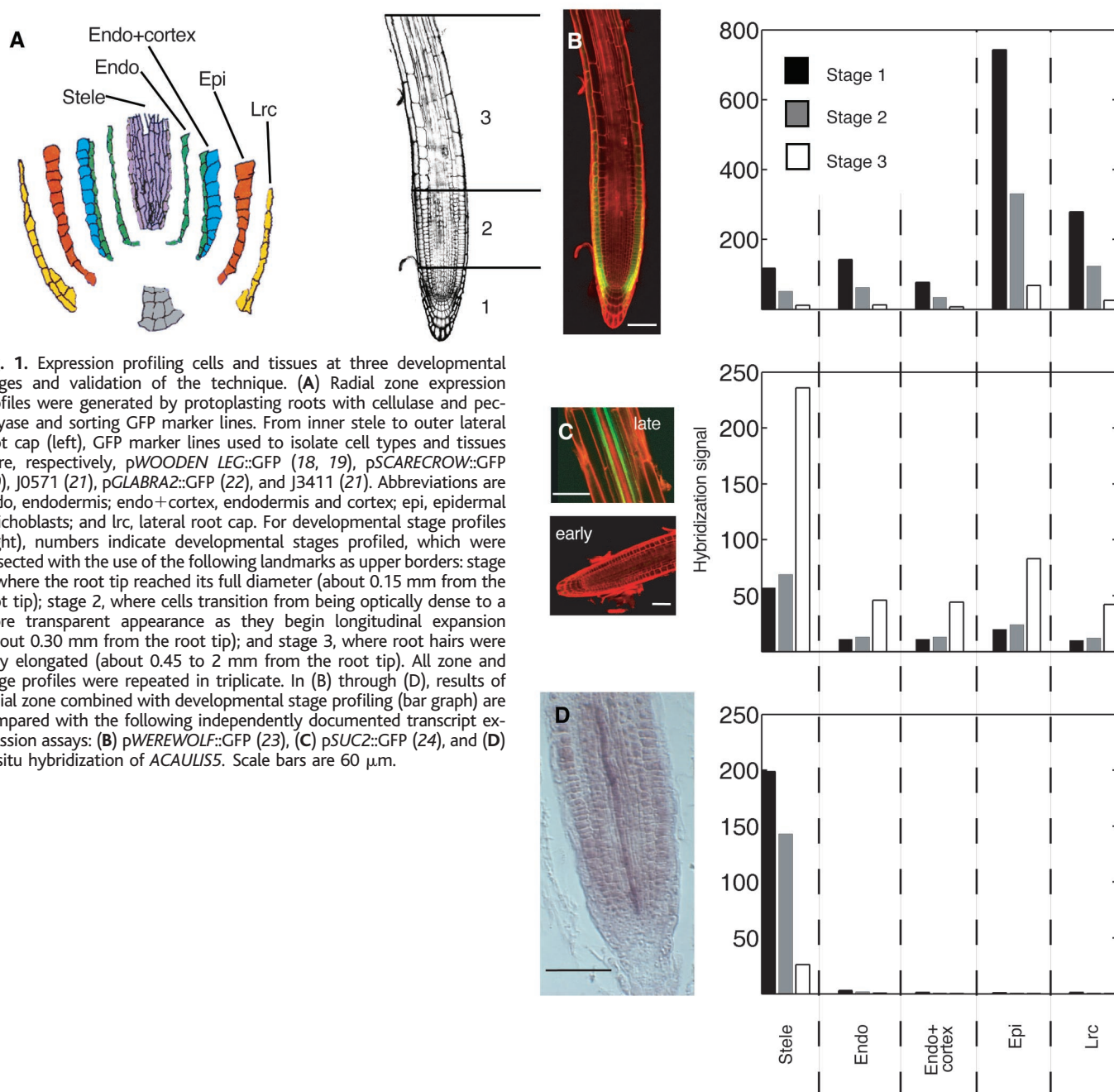
We next analyzed differential regulation and broad patterns of gene expression. We defined “differential regulation” as those genes with a maximum expression value four times higher than their minimal expression value (which we determined to be a highly conservative criteria for distinguishing true induction from expression fluctuations due to noise on the basis of genes with known root expression patterns). Out of 10,492 genes detected in the root, 5717 transcripts (~54%) were differentially regulated across subzones. Thus, the majority of detectable genes changed dramatically in expression across root subzones.

To identify dominant expression patterns, we first coded the expression of all genes into binary values, scoring expression signals ei-

ther above (1) or below (0) the previously defined detection cutoff of 75. This discretization of the data focuses on expression fluctuations at the level of detection. Although over 1500 different patterns were observed, 10 patterns accounted for 20% of all genes that were not either ubiquitous or entirely below detection levels ($n = 7480$ binary regulated genes). Many other patterns deviated by only one or two bits (out of 15) from the 10 major patterns, indicating that many genes fell into a few patterns. To corroborate the number of major patterns in the data, we applied principal component analysis (PCA) to the 5717 differentially expressed genes (9). The analysis showed that eight eigenvalues accounted for 85% of the variation in the

data, and each of the eight explained at least 4% of the variation (10). Thus, PCA, which uses the full range of expression values, indicated that about eight dominant patterns summarized much of the way genes varied among the 15 root zones. Moreover, many of the trends identified by the eigenvectors in PCA were similar to cell-type enrichment patterns revealed by the binary analysis.

Because the expression profiles fit into 8 to 10 major patterns and we sought to summarize expression patterns broadly, we used K-means clustering to group all 5717 differentially expressed genes into eight groups (11). Many of the same specific patterns detected in PCA and binary analysis emerged in the K-means clustering, which partitions the



REPORTS

Fig. 2. Global expression map depicting major patterns of gene activity in the *Arabidopsis* root. Each of the 15 (5 × 3) subgrids contains all 5712 genes that varied by at least a factor of four between any two average signals in the expression profiles. The 15 subgrids are superimposed on their position on one-half of the bilaterally symmetrical root, with radial zones representing cell and tissue types and stages representing stages of cell development. Transcripts are placed in order of the LED in which they fell and then ordered within LEDs by peak expression value. Genes are placed in the same coordinates in each subgrid so that numbers on the vertical axis mark the position of the same cluster at the three developmental stages. Colors indicate the magnitude of expression signals on a log scale. These values are based on microarray hybridization signals, which have no units. White boxes indicate the subzone(s) where genes in each of the eight LEDs (designated by the number) reach their peak expression. Numbers at right identify the position of each LED in each stage for stage comparisons.

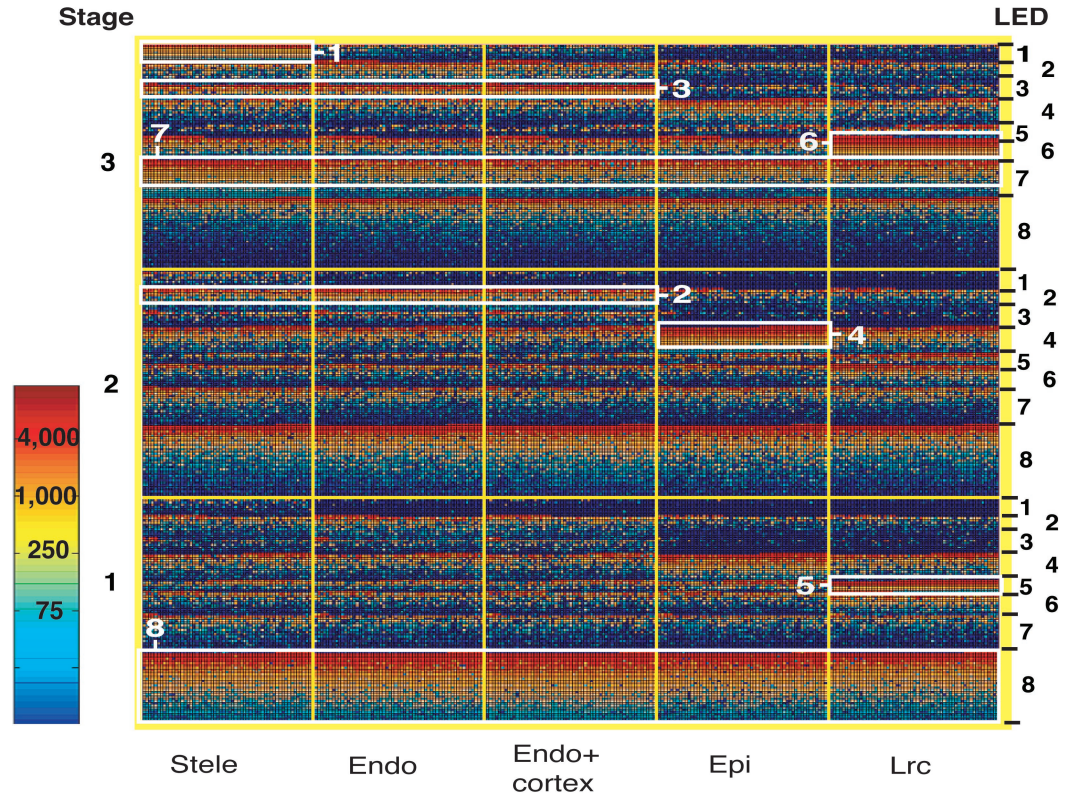
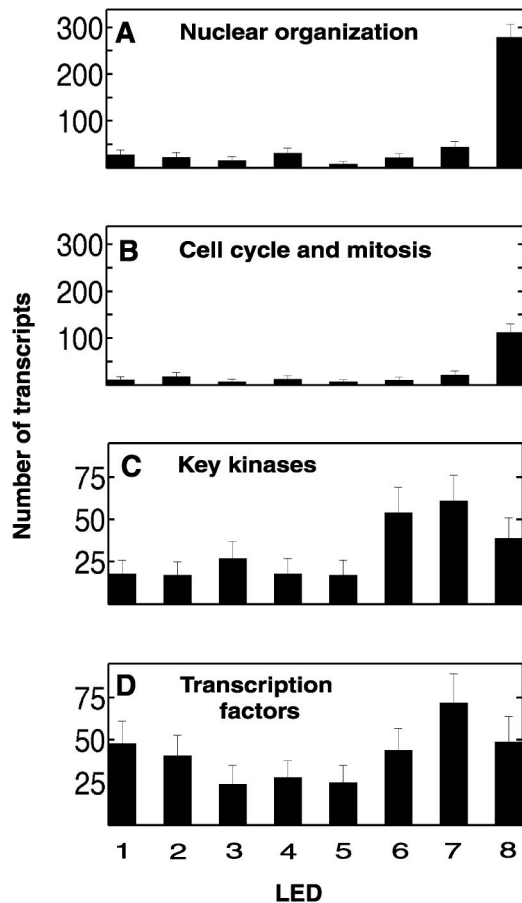


Fig. 3. Overrepresentation of functional classes in expression domains. Numbered LEDs on the x axis are defined by their subzones of peak expression: (1) stele stage 3, (2) stele and ground tissue stage 2, (3) stele and ground tissue stage 3, (4) epidermis stage 2, (5) lateral root cap stage 1, (6) all zones peaking in lateral root cap stage 3, (7) all zones stage 3, (8) all zones stage 1. Genes annotated with the designations of nuclear organization (A) and cell cycle and mitosis (B) are overrepresented in LED 8. Genes annotated as key kinases (C) are overrepresented in LEDs 6 and 7, and genes annotated as specific transcription factors (D) are overrepresented in LED 7. Attached bars are 95% confidence levels generated by bootstrapping the sample 1000 times with the observed frequency of positive scores. Annotations are from the Munich Information Center for Protein Sequences functional categories database, except for transcription factors, which was a list of 1600 transcription factors compiled from several sources.



data into discrete expression groups and also uses the full range of expression values. The agreement of these different analyses supports the robustness of the major expression patterns found (table S3).

To visualize the patterns, we mapped the expression of all the differentially regulated genes onto a grid representing the physical root (Fig. 2). The grid was organized into 15 subgrids corresponding to the subzones of the root, and each subgrid contained the expression signal of all 5717 genes, with their expression intensity depicted on a color scale. Genes were laid down in the subgrids by clusters with the same gene at the same coordinates in each subgrid. Thus, groups of genes with coordinated induction and repression are apparent as broad color bands that represent spatial and temporal domains. These clusters are not the only expression patterns we observed among root-expressed genes but rather represent large-scale trends in expression. We call these clusters, which are composed of sets of genes whose expression is enriched in one or more subzones, localized expression domains (LEDs) (table S2).

We expected LEDs with the most dramatic biases in gene functional categories to correlate strictly with cell maturity. Thus, stage-specific patterns offered an opportunity to test whether LEDs organized genes into functionally coherent groups. For example,

Fig. 4. Putative hormone activity centers. Ovals show the regions encompassed by LEDs with an overrepresentation of hormone-related genes: (A) auxin, (B) JA, and (C) GA. Colors correspond to LED 1 (red, stele stage 3), LED 2 (green; stele, endodermis, and cortex stage 2), LED 4 (yellow, epidermis stage 2), and LED 5 (blue, lateral root cap stage 1).

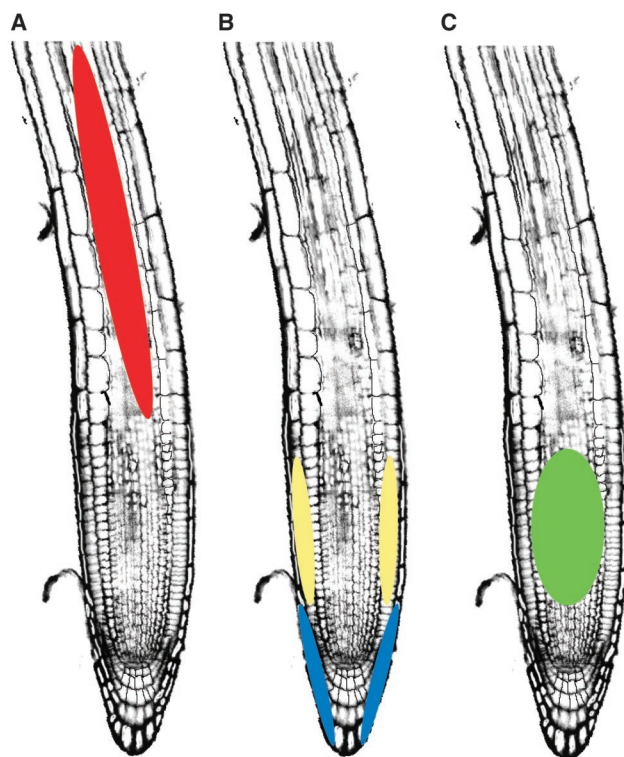


Table 1. Root expression patterns among members of five transcription factor families. LEDs are the same as those defined in Fig. 2. The number of members of each family in each LED is represented. The category "Other" is composed of all 4775 genes that showed relatively constant root expression, changing in expression by less than a factor of four between any 2 of the 15 subzones.

Transcription factor family	LED								Other
	1	2	3	4	5	6	7	8	
MYB	4	6	8	2	0	2	4	2	7
AP2-like	2	4	1	2	3	1	5	5	16
WRKY	1	3	0	3	1	12	2	1	14
HD-ZIP	5	1	0	1	0	1	2	0	0
bHLH	4	3	0	1	1	1	1	0	7

LED 8, which spanned all radial zones but peaked in stages 1 and 2, was significantly overrepresented with genes involved in nuclear organization, cell cycle, and mitosis (Fig. 3, A and B), as expected of rapidly dividing cells in this broad meristematic zone. Other subcategories that showed the same pattern were organization of cytoplasm, mitochondrial organization, and ribosomal proteins (12). LED 7, which spanned all radial zones but peaked in stage 3, contained significantly more kinases and transcription factors than most other LEDs (Fig. 3, C and D), delineating an apparent extensive signaling network associated with cell maturation across cell types.

We were particularly interested in LEDs that spanned only a subset of radial zones because of their potential to define cell and tissue specificity. Taking as an example genes that regulate or respond to plant hormones, we found these genes in many LEDs.

However, three LEDs showed an aggregation of genes with known or putative roles in auxin, gibberellic acid (GA), or jasmonic acid (JA) pathways, suggesting the possibility of localized signaling centers.

LED 1 (451 genes), which contains the 4.3% of root-expressed genes enriched in the stele in stage 3 (Fig. 4A), contained 7 of the 49 auxin-related genes expressed in the root (14%) (table S4). This was the only significant grouping of auxin-related genes ($P = 0.015$), because most such genes were contained within the nondifferentially regulated group (4775 genes) or in the two large stage-specific LEDs (7 and 8). Thus, genes related to auxin function show expression in the root center, where auxin is known to affect vascular development (13).

Gene expression related to JA function characterized the outer layers of the root. Of 47 JA-related genes expressed in the root, 9 were expressed in LED 4 (672

genes), which peaked in the epidermis in stage 2 (Fig. 4B, yellow) and 10 were expressed in LED 5 (299 genes), which peaked in the lateral root cap in stage 1 (Fig. 4B, blue). The two LEDs comprised only 6.4% and 2.8% of root-expressed genes but contained 19% and 21% of the root-expressed JA-related genes, respectively (table S4). They were the only two significant aggregations of JA-related genes ($P < 0.01$ and $P < 0.0002$). These LEDs included many genes implicated in pathogen response, such as myrosinase-binding proteins, lectins, and CYP79B2 (14, 15). These LEDs may reflect constitutive JA-induced defense responses in the outer layers of the root (roots in our experiments were grown under sterile conditions) or may signify a function for JA in coordinating root development.

Another aggregation of hormone-related genes was found in LED 2 (587 genes; Fig. 4C), which comprised only 5.6% of the root-expressed genes but contained 9 of the 29 GA-related genes expressed in the root (31%) (table S4). Similarly, this represented the only significant aggregation of GA-related genes ($P < 0.0002$). LEDs 2 and 3 do not conform to classically defined tissue boundaries, and they may identify new developmental domains not apparent by morphological classification. Together these data lead us to hypothesize the existence of JA and GA hormone signaling centers with the potential to serve as local organizing centers in root development.

We also examined the expression of transcription factors (TFs) to explore the complexity of regulatory mechanisms in the root. From a list of 1411 transcription factors in *Arabidopsis* on the ATH1 microarray, we detected 577 in the root, of which 331 were differentially regulated. Most of the TF families examined had members in multiple LEDs, indicating a divergence in their developmental roles. Most gene families examined also had multiple paralogs in the same LED (Table 1), identifying small groups of potentially redundant TFs that could be targeted for multiple mutant analyses. One extreme case was the WRKY TF family; 12 of the 37 members found in the root were contained in LED 6, indicating a possible specialized role for these TFs in cell maturation. Most of the differentially regulated TFs found are as-yet uncharacterized.

We also found evidence for chromosomal clustering (4, 16, 17) of the coregulated sets of genes within LEDs. Genes in four of the eight LEDs were significantly clustered on the chromosome as compared to random gene sets [e.g., (7)]—LED 1 (47 observed, 31 expected, $P < 0.035$), LED 3 (44 observed, 23 expected, $P < 0.0014$), LED 4 (86 observed, 63 expected, $P < 0.0278$),

and LED 8 (541 observed, 456 expected, $P < 0.0006$)—providing evidence in plants for a link between genome organization and gene regulation.

Together these data provide an organ expression map, revealing putative localized hormone-response domains and a complex pattern of regulatory genes that could mediate primary developmental cues. These data should help identify candidate genes involved in pattern formation and cell specificity in the root, which is a model for organogenesis. The expression map will also facilitate both computational and experimental methods aimed at decoding regulatory mechanisms in the root. Thus, these results can now be used to explore how the hundreds of different expression patterns they reveal are established and interpreted at the cellular level to generate a complex organ.

References and Notes

1. N. M. Kerk, T. Ceserani, S. L. Tausta, I. M. Sussex, T. M. Nelson, *Plant Physiol.* **132**, 27 (2003).
 2. T. Asano *et al.*, *Plant J.* **32**, 401 (2002).

3. D. Milioni, P. E. Sado, N. J. Stacey, K. Roberts, M. C. McCann, *Plant Cell* **14**, 2813 (2002).
 4. P. J. Roy, J. M. Stuart, J. Lund, S. K. Kim, *Nature* **418**, 975 (2002).
 5. H. Jasper *et al.*, *Dev. Cell* **3**, 511 (2002).
 6. P. N. Benfey, J. W. Schiefelbein, *Trends Genet.* **10**, 84 (1994).
 7. Materials and methods are available as supporting material on Science Online.
 8. J. Sheen, *Plant Physiol.* **127**, 1466 (2001).
 9. J. Quackenbush, *Nature Rev. Genet.* **2**, 418 (2001).
 10. The program Cluster was used in the analysis and downloaded from <http://rana.lbl.gov/EisenSoftware.htm>.
 11. M. B. Eisen, P. T. Spellman, P. O. Brown, D. Botstein, *Proc. Natl. Acad. Sci. U.S.A.* **95**, 14863 (1998).
 12. K. Birnbaum *et al.*, unpublished data.
 13. T. Berleth, J. Mattsson, *Curr. Opin. Plant Biol.* **3**, 406 (2000).
 14. U. Wittstock, B. A. Halkier, *Trends Plant Sci.* **7**, 263 (2002).
 15. L. L. Murdock, R. E. Shade, *J. Agric. Food Chem.* **50**, 6605 (2002).
 16. B. A. Cohen, R. D. Mitra, J. D. Hughes, G. M. Church, *Nature Genet.* **26**, 183 (2000).
 17. H. Caron *et al.*, *Science* **291**, 1289 (2001).
 18. A. P. Mahonen *et al.*, *Genes Dev.* **14**, 2938 (2000).
 19. M. Bonke, S. Thitamadee, A. P. Mahonen, M. T. Hauser, Y. Helariutta, *Nature*, in press.
 20. J. W. Wysocka-Diller, Y. Helariutta, H. Fukaki, J. E. Malamy, P. N. Benfey, *Development* **127**, 595 (2000).
 21. The plant line was generated by the Haseloff labora-

tory (www.plantsci.cam.ac.uk/Haseloff/Home.html). The lines were obtained through the Arabidopsis Information Resource (www.arabidopsis.org/).

22. Y. Lin, J. Schiefelbein, *Development* **128**, 3697 (2001).
 23. M. M. Lee, J. Schiefelbein, *Cell* **99**, 473 (1999).
 24. E. Truernit, N. Sauer, *Planta* **196**, 564 (1995).
 25. We thank J. Malamy for valuable ideas on the protoplasting technique; H. Petri, K. Gordon, and J. Hirst for assistance in cell sorting; H. Dressman and the Duke Microarray Core Facility for assistance with microarrays; A. Pekka Mähönen and Y. Helariutta for use of the pWOL::GFP line and M. Cilia and D. Jackson for the pSUC2::GFP line, both before publication; M. Levesque for valuable discussions; and G. Sena and T. Navy for photos. This work was supported by NSF grants MCB-020975 (P.N.B. and D.E.S.), DBI-9813360 (D.W.G.), DBI-0211857 (D.W.G.), and a Small Grant for Exploratory Research (P.N.B. and D.E.S.). The NIH supported K.B. with a postdoctoral fellowship grant (5 F32 GM20716-03).

Supporting Online Material

www.sciencemag.org/cgi/content/full/302/5652/1956/DC1

Materials and Methods

Figs. S1 to S3

Tables S1 to S4

4 August 2003; accepted 15 October 2003

Inferring Nonneutral Evolution from Human-Chimp-Mouse Orthologous Gene Trios

Andrew G. Clark,¹ Stephen Glanowski,³ Rasmus Nielsen,² Paul D. Thomas,⁴ Anish Kejariwal,⁴ Melissa A. Todd,² David M. Tanenbaum,⁵ Daniel Civello,⁶ Fu Lu,⁵ Brian Murphy,³ Steve Ferriera,³ Gary Wang,³ Xianqun Zheng,⁵ Thomas J. White,⁶ John J. Sninsky,⁶ Mark D. Adams,^{5*} Michele Cargill^{6†}

Even though human and chimpanzee gene sequences are nearly 99% identical, sequence comparisons can nevertheless be highly informative in identifying biologically important changes that have occurred since our ancestral lineages diverged. We analyzed alignments of 7645 chimpanzee gene sequences to their human and mouse orthologs. These three-species sequence alignments allowed us to identify genes undergoing natural selection along the human and chimp lineage by fitting models that include parameters specifying rates of synonymous and nonsynonymous nucleotide substitution. This evolutionary approach revealed an informative set of genes with significantly different patterns of substitution on the human lineage compared with the chimpanzee and mouse lineages. Partitions of genes into inferred biological classes identified accelerated evolution in several functional classes, including olfaction and nuclear transport. In addition to suggesting adaptive physiological differences between chimps and humans, human-accelerated genes are significantly more likely to underlie major known Mendelian disorders.

Although the human genome project will allow us to compare our genome to that of other primates and discover features that are uniquely human, there is no guarantee that such features are responsible for any of our unique biological attributes. To identify genes and biological processes that have been most altered by our recent evolutionary divergence from other primates, we need to fit the data to models of sequence divergence that allow us to distinguish between diver-

gence caused by random drift and divergence driven by natural selection. Early observations of unexpectedly low levels of protein divergence between humans and chimpanzees led to the hypothesis that most of the evolutionary changes must have occurred at the level of gene regulation (*1*). Recently, much more extensive efforts at DNA sequencing in nonhuman primates has confirmed the very close evolutionary relationship between humans and chimps (*2*), with an

average nucleotide divergence of just 1.2% (*3–5*). The role of protein divergence in causing morphological, physiological, and behavioral differences between these two species, however, remains unknown.

Here we apply evolutionary tests to identify genes and pathways from a new collection of more than 200,000 chimpanzee exonic sequences that show patterns of divergence consistent with natural selection along the human and chimpanzee lineages.

To construct the human-chimp-mouse alignments, we sequenced PCR amplifications using primers designed to essentially all human exons from one male chimpanzee, resulting in more than 20,000 human-chimp gene alignments spanning 18.5 Mb (*6–8*). To identify changes that are specific to the divergence in the human lineage, we compared the human-chimp aligned genes to their mouse ortholog. Inference of orthology involved a combination of reciprocal best matches and syntenic evidence between human and mouse gene annotations (*9, 10*). This genome-wide set of orthologs underwent a series of filtering steps to remove ambiguities, orthologs with little sequence data, and genes with suspect annotation (*6*). The filtered ortholog set was compared to

¹Molecular Biology and Genetics, ²Biological Statistics and Computational Biology, Cornell University, Ithaca, NY 14853, USA. ³Applied Biosystems, 45 West Gude Drive, Rockville, MD 20850, USA. ⁴Protein Informatics, Celera Genomics, 850 Lincoln Centre Drive, Foster City, CA 94404, USA. ⁵Celera Genomics, 45 West Gude Drive, Rockville, MD 20850, USA. ⁶Celera Diagnostics, 1401 Harbor Bay Parkway, Alameda, CA 94502, USA.

*Present address: Department of Genetics, Case Western Reserve University, 10900 Euclid Avenue, Cleveland, OH 44106, USA.

†To whom correspondence should be addressed. E-mail: michele_cargill@celeradiagnostics.com