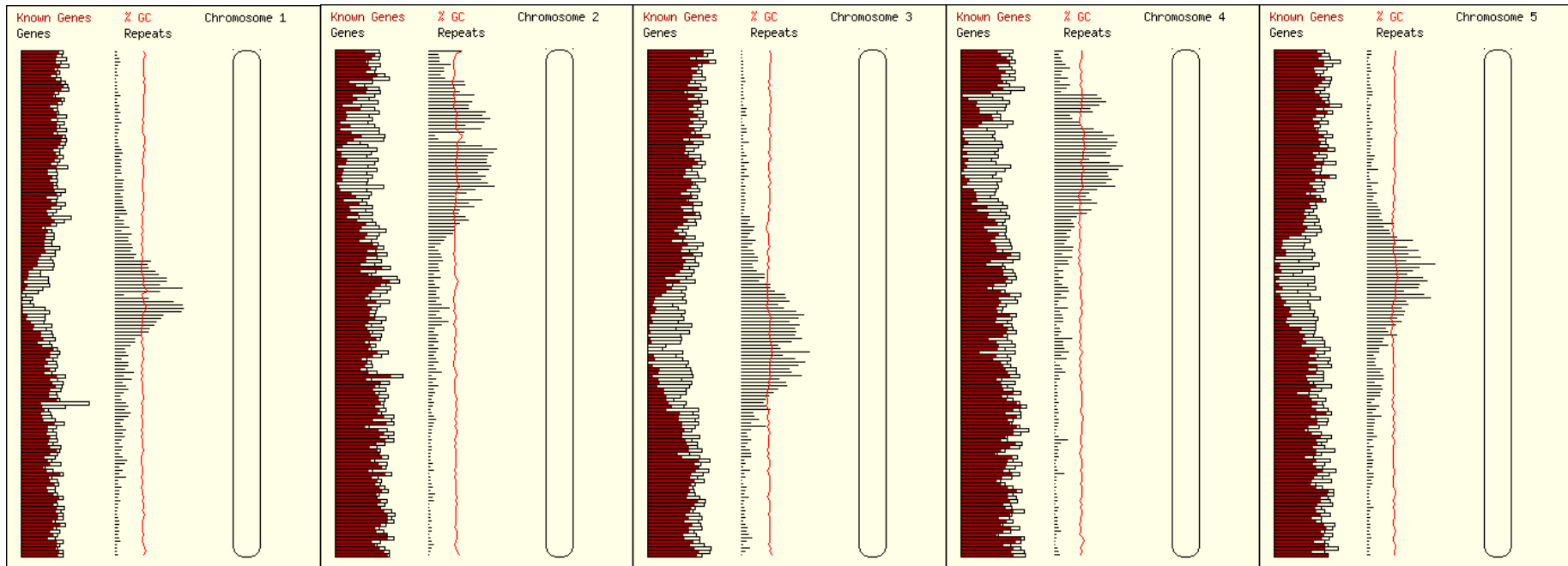


# Comparative Features of Multicellular Eukaryotic Genomes (2017)

(First three statistics from [www.ensembl.org](http://www.ensembl.org); other from original papers)

	<i>C. elegans</i>	<i>A. thaliana</i>	<i>D. melanogaster</i>	<i>M. musculus</i>	<i>H. sapiens</i>
Size (Mb)	103	136	143	3,482	3,555
# Protein-coding genes	20,362	27,655 (25,498 original est.)	13,918 (13,601 original est.)	22,598 (30,000 original est.)	20,338 (30,000 original est.)
Transcripts	58,941	55,157	34,749	131,195	200,310
Gene density (#/kb)	1/5	1/4.5	1/8.8	1/83	1/97
LINE/SINE (%)	0.4	0.5	0.7	27.4	33.6
LTR (%)	0.0	4.8	1.5	9.9	8.6
DNA Elements	5.3	5.1	0.7	0.9	3.1
Total repeats	6.5	10.5	3.1	38.6	46.4
Exons					
% genome size	27	28.8	24.0		
per gene	4.0	5.4	4.1	8.4	8.7
average size (bp)		250	506		
Introns					
% genome size		15.6			
average size (bp)		168			

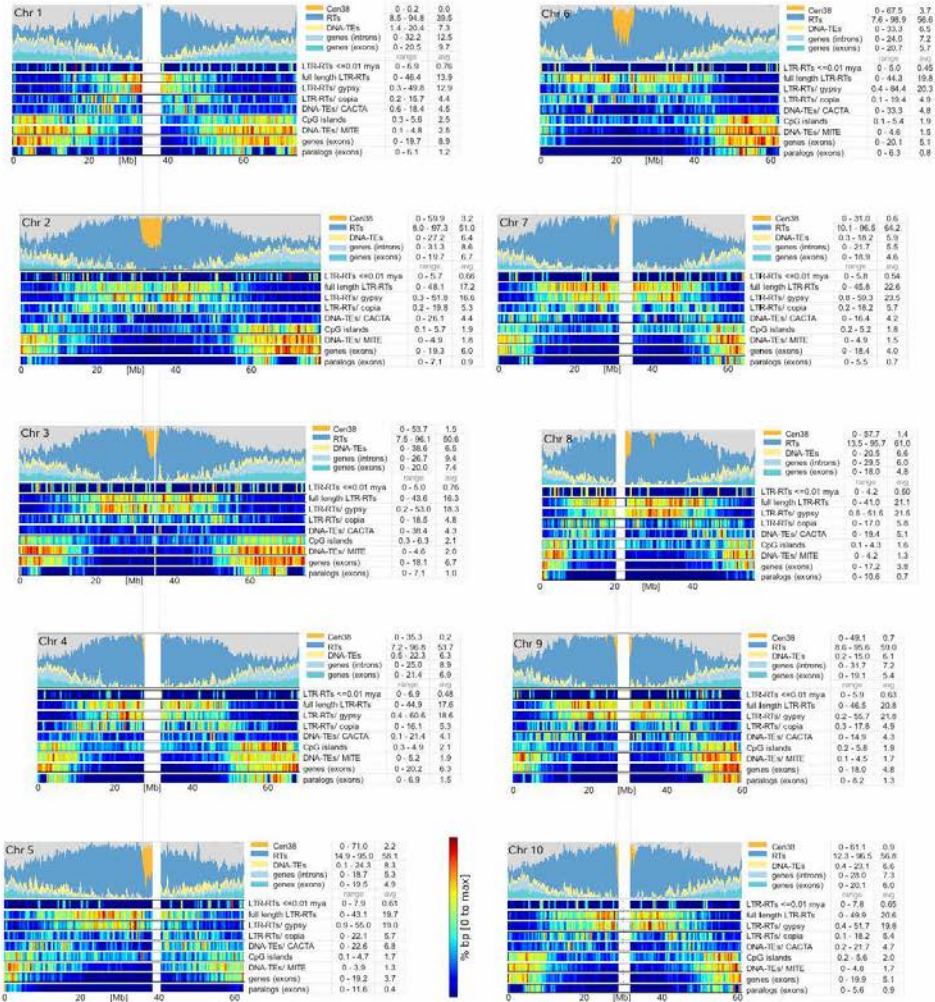
# Arabidopsis Chromosome Structures



# Sorghum Whole Genome Details

doi: 10.1038/nature07723

SUPPLEMENTARY INFORMATION



# Characterizing the Proteome

## The Protein World

- Sequencing has defined
  - Many, many proteins
- How can we use this data to:
  - Define genes in new genomes
  - Look for evolutionarily related genes
  - Follow evolution of genes
    - Mixing of domains to create new proteins
  - Uncover important subsets of genes that
    - That deep phylogenies
      - Plants vs. animals
      - Placental vs. non-placental animals
      - Monocots vs. dicots plants
- Common nomenclature needed
  - Ensure consistency of interpretations

## InterPro (<http://www.ebi.ac.uk/interpro/>)

- Integrated documentation resource for protein families, domains and functional sites
- Nucleic Acids Research (2017) 45:D190-D199 (database issue)
- A database of protein families, domains, and functional sites in which identifiable features found in known proteins can be applied to unknown protein sequences
- A database that collapses information from 11 other databases

- Current Release: 28 September 2017
  - 33,283 total entries
    - 14,223 InterPro entries
      - Match to 33,283 GO functions
  - Based on 90,606,305 UniProtKD proteins
  - Updated six times since November 2016

## InterPro Entries Defined as:

### Family (20,229 on 9/17)

- Evolutionarily related proteins
- Share architecture base on:
  - Domain signature
  - Repeat signature
- Multiple signatures can define a family
- 7/2003 release: 6.416 families
- Examples
  - IPR000003 Retinoid X Receptor
  - IPR000006 Vertebrate metallothionein
  - IPR000007 Tubby
  - IPR000009 Protein phosphate 2A regulatory subunit PR55
  - IPR000010 Cysteine protease inhibitor

## Example Families

### Cytochrome P450

The cytochrome P450 enzymes constitute a superfamily of haem-thiolate proteins. P450 enzymes usually act as terminal oxidases in multicomponent electron transfer chains, called P450-containing monooxygenase systems and are involved in metabolism of a plethora of both exogenous and endogenous compounds. P450-containing monooxygenase systems primarily fall into two major classes: bacterial/mitochondrial (type I), and microsomal (type II). All P450 enzymes can be categorised into two main groups, the so-called B- and E-classes: P450 proteins of prokaryotic 3-component systems and fungal P450nor (CYP55) belong to the B-class; all other known P450 proteins from distinct systems are of the E-class [ [1](#) ].

### Disease resistance protein

Plants are attacked by a range of phytopathogenic organisms, including viruses, mycoplasma, bacteria, fungi, nematodes, protozoa and parasites. Resistance to a pathogen is manifested in several ways and is often correlated with a hypersensitive response (HR), localised induced cell death in the host plant at the site of infection [ [1](#) ]. The induction of the plant defense response that leads to HR is initiated by the plants recognition of specific signal molecules (elicitors) produced by the pathogen; R genes are thought to encode receptors for these elicitors. RPS2, N and L6 genes confer resistance to bacterial, viral and fungal pathogens.

## Domain (8,735 on 9/17)

- Independent structural unit
- Can be associated with other
  - Domains
  - Repeat
- Proteins with the domain are evolutionarily related
- 7/2003 release: 1,902 domains
- Examples
  - IPR000001 Kringle
  - IPR000002 Cdc20/Fizzy
  - IPR000005 Helix-turn-helix, AraCtype
  - IPR000008 C2 domain
  - IPR000014 PAS domain

## Example domains

### Protein kinase-like; Protein kinase

Protein kinases ( [IPR000719](#) ) catalyze the phosphotransfer reaction fundamental to most signalling and regulatory processes in the eukaryotic cell [ [1](#) ]. The catalytic subunit contains a core that is common to both serine/threonine and tyrosine protein kinases. The catalytic domain contains the nucleotide-binding site and the catalytic apparatus in an inter-lobe cleft. Structurally it shares functional and structural similarities with the ATP-grasp fold, which is found in enzymes that catalyse the formation of an amide bond, and with PIPK (phosphoinositol phosphate kinase). The three-dimensional fold of the protein kinase catalytic domain is similar to domains found in several other proteins. These include the catalytic domain of actin-fragmin kinase, an atypical protein kinase that regulates the F-actin capping activity in plasmodia [ [2](#) ]; the catalytic domain of phosphoinositide-3-kinase (PI3K), which phosphorylates phosphoinositides and as such is involved in a number of fundamental cellular processes such as apoptosis, proliferation, motility and adhesion [ [3](#) ]; the catalytic domain of the MHCK/EF2 kinase, an atypical protein kinase that includes the TRP (transient channel potential) calcium-channel kinase involved in the modulation of calcium channels in eukaryotic cells in response to external signals [ [4](#) ]; choline kinase, which catalyses the ATP-dependent phosphorylation of choline during the biosynthesis of phosphatidylcholine [ [5](#) ]; and 3',5'-aminoglycoside phosphotransferase type IIIa, a bacterial enzyme that confers resistance to a range of aminoglycoside antibiotics [ [6](#) ].

## Repeat (280 on 9/17)

- Region not expected to form a globular structure on their own
  - WD40
    - Six to eight copies needed for globular domain to form
- 7/2003 release: 163 repeats
- Examples
  - IPR000033 Low-density lipoprotein receptor, YWTD repeat
  - IPR000102 Neuraxin/MAP1B repeat
  - IPR000127 Ubiquitin-activating enzyme repeat
  - IPR000225 Armadillo repeat
  - IPR000258 Bacterial ice-nucleation proteins octamer repeat

## Example Repeats

### Leucine-rich repeat

Leucine-rich repeats (LRR) consist of 2-45 motifs of 20-30 amino acids in length that generally folds into an arc or horseshoe shape [ [1](#) ]. LRRs occur in proteins ranging from viruses to eukaryotes, and appear to provide a structural framework for the formation of protein-protein interactions [ [2](#) ]. Proteins containing LRRs include tyrosine kinase receptors, cell-adhesion molecules, virulence factors, and extracellular matrix-binding glycoproteins, and are involved in a variety of biological processes, including signal transduction, cell adhesion, DNA repair, recombination, transcription, RNA processing, disease resistance, apoptosis, and the immune response.

Sequence analyses of LRR proteins suggested the existence of several different subfamilies of LRRs. The significance of this classification is that repeats from different subfamilies never occur simultaneously and have most probably evolved independently. It is, however, now clear that all major classes of LRR have curved horseshoe structures with a parallel beta sheet on the concave side and mostly helical elements on the convex side. At least six families of LRR proteins, characterized by different lengths and consensus sequences of the repeats, have been identified. Eleven-residue segments of the LRRs (LxxLxLxxN/CxL), corresponding to the  $\beta$ -strand and adjacent loop regions, are conserved in LRR proteins, whereas the remaining parts of the repeats (herein termed variable) may be very different. Despite the differences, each of the variable parts contains two half-turns at both ends and a "linear" segment (as the chain follows a linear path overall), usually formed by a helix, in the middle. The concave face and the adjacent loops are the most common protein interaction surfaces on LRR proteins. 3D structure of some LRR proteins-ligand complexes show that the concave surface of LRR domain is ideal for interaction with alpha-helix, thus supporting earlier conclusions that the elongated and curved LRR structure provides an outstanding framework for achieving diverse protein-protein interactions [ [2](#) ]. Molecular modeling suggests that the conserved pattern LxxLxL, which is shorter than the previously proposed LxxLxLxxN/CxL is sufficient to impart the characteristic horseshoe curvature to proteins with 20- to 30-residue repeats [ [3](#) ].



## PPR repeat

Pentatricopeptide repeat proteins are characterised by the presence of a tandem array of repeats, where the number of PPR motifs controls the affinity and specificity of the PPR protein for RNA. These proteins occur predominantly in plants, where they appear to play essential roles in RNA/DNA metabolism in mitochondria and chloroplasts [ [1](#) ]. It has been suggested that each of the highly variable PPR proteins is a gene-specific regulator of plant organellar RNA metabolism. PPR proteins may also play a role in organelle biogenesis, probably via binding to organellar transcripts [ [2](#) ]. Examples of PPR repeat-containing proteins include PET309 [P32522](#) , which may be involved in RNA stabilisation [ [3](#) ], and crp1, which is involved in RNA processing [ [4](#) ]. The repeat is associated with a predicted plant protein [O49549](#) that has a domain organization similar to the human BRCA1 protein.

## Site (812 on 9/17)

- **Post-translational modification**
  - Modifies primary protein structure
  - Glycosylation
  - Splicing
  - Phosphorylation
- 7/2003 release: 20 post-translational modifications
- Examples
  - IPR000042 N-glycosylation site
  - IPR000134 Amidation site
  - IPR000152 Aspartic acid and asparagin hydroxylation site
  - IPR000220 Tyrosine kinase phosphorylation site
  - IPR000338 N-myristoylation site

## Binding site (76 on 9/17)

- Binds chemical compounds
  - Compounds not substrates
  - Cofactors required for protein function
- 7/2003 release: 20 binding sites
- Examples
  - IPR000205 NAD-binding site
  - IPR000214 Formamidopyrimidine-DNA glycolase, Zn-binding site
  - IPR000345 Cytochrome c heme-binding site
  - IPR000634 Serine/threonine dhydrate, pyridoxal-phosphate-binding site
  - IPR001216 Cysteine synthase/cystathionine beta-synthase P-phosphate binding site

## • Active site (132 on 9/17)

- Catalytic sites of protein
- Site where substrate binds
- 7/2003 release: 26 active sites
- Examples
  - IPR000126 Serine protease V8, active site
  - IPR000138 Hydroxymethylglutaryl-coenzyme A lyase, active site
  - IPR000180 Peptidase M19, renal dipeptidase, active site
  - IPR000189 Prokaryotic transglycosylase, active site
  - IPR000590 Hydroxymethylglutaryl-coenzyme A synthase, active site

## Conserved Site (687 on 9/17)

- A site motif specific defined by PROSITE that characterizes a protein or protein family
- Not an PTM, binding site, or active site
- 9/2/08 release: 166 conserved sites
- Examples
  - IPR000035 Alkylbase DNA glycosidase, conserved site
  - IPR000059 NUDIX hydrolase, conserved site
  - IPR000132 Nitrilase/cyanide hydratase, conserved site
  - IPR000283 NADH dehydrogenase 75 kDa subunit, conserved sit
  - IPR000291 D-alanine--D-alanine ligase/VANA/B/C, conserved site

## InterPro Strategy









([https://www.ebi.ac.uk/interpro/about.html#about\\_08](https://www.ebi.ac.uk/interpro/about.html#about_08))







- Uses nomenclature developed by other databases
- Collapses those nomenclatures to define each InterPro entry

## The InterPro Consortium

(<https://www.ebi.ac.uk/interpro/about.html>)

The following databases make up the InterPro Consortium:

	<p><u>CATH-Gene3D</u> database describes protein families and domain architectures in complete genomes. Protein families are formed using a Markov clustering algorithm, followed by multi-linkage clustering according to sequence identity. Mapping of predicted structure and sequence domains is undertaken using hidden Markov models libraries representing CATH and Pfam domains. CATH-Gene3D is based at University College, London, UK.</p>
	<p><u>CDD</u> is a protein annotation resource that consists of a collection of annotated multiple sequence alignment models for ancient domains and full-length proteins. These are available as position-specific score matrices (PSSMs) for fast identification of conserved domains in protein sequences via RPS-BLAST. CDD content includes NCBI-curated domain models, which use 3D-structure information to explicitly define domain boundaries and provide insights into sequence/structure/function relationships, as well as domain models imported from a number of external source databases.</p>
	<p><u>MobiDB</u> offers a centralized resource for annotations of intrinsic protein disorder. The database features three levels of annotation: manually curated, indirect and predicted. The different sources present a clear tradeoff between quality and coverage. By combining them all into a consensus annotation, MobiDB aims at giving the best possible picture of the "disorder landscape" of a given protein of interest.</p>
	<p><u>HAMAP</u> stands for High-quality Automated and Manual Annotation of Proteins. HAMAP profiles are manually created by expert curators. They identify proteins that are part of well-conserved proteins families or subfamilies. HAMAP is based at the SIB Swiss Institute of Bioinformatics, Geneva, Switzerland.</p>
	<p><u>PANTHER</u> is a large collection of protein families that have been subdivided into functionally related subfamilies, using human expertise. These subfamilies model the divergence of specific functions within protein families, allowing more accurate association with function, as well as inference of amino acids important for functional specificity. Hidden Markov models (HMMs) are built for each family and subfamily for classifying additional protein sequences. PANTHER is based at at University of Southern California, CA, US.</p>
	<p><u>Pfam</u> is a large collection of multiple sequence alignments and hidden Markov models covering many common protein domains. Pfam is based at EMBL-EBI, Hinxton, UK.</p>
	<p><u>PIRSF</u> protein classification system is a network with multiple levels of sequence diversity from superfamilies to subfamilies that reflects the evolutionary relationship of full-length proteins and domains. PIRSF is based at the Protein Information Resource, Georgetown University Medical Centre, Washington DC, US.</p>
	<p><u>PRINTS</u> is a compendium of protein fingerprints. A fingerprint is a group of conserved motifs used to characterise a protein family or domain. PRINTS is</p>

	based at the University of Manchester, UK.
	<u>ProDom</u> protein domain database consists of an automatic compilation of homologous domains. Current versions of ProDom are built using a novel procedure based on recursive PSI-BLAST searches. ProDom is based at PRABI Villeurbanne, France.
	<u>PROSITE</u> is a database of protein families and domains. It consists of biologically significant sites, patterns and profiles that help to reliably identify to which known protein family a new sequence belongs. PROSITE is base at the Swiss Institute of Bioinformatics (SIB), Geneva, Switzerland.
	<u>SFLD</u> (Structure-Function Linkage Database) is a hierarchical classification of enzymes that relates specific sequence-structure features to specific chemical capabilities.
	<u>SMART</u> (a Simple Modular Architecture Research Tool) allows the identification and annotation of genetically mobile domains and the analysis of domain architectures. SMART is based at at EMBL, Heidelberg, Germany.
<i>Superfamily</i> 	<u>SUPERFAMILY</u> is a library of profile hidden Markov models that represent all proteins of known structure. The library is based on the SCOP classification of proteins: each model corresponds to a SCOP domain and aims to represent the entire SCOP superfamily that the domain belongs to. SUPERFAMILY is based at the University of Bristol, UK.
	<u>TIGRFAMs</u> is a collection of protein families, featuring curated multiple sequence alignments, hidden Markov models (HMMs) and annotation, which provides a tool for identifying functionally related proteins based on sequence homology. TIGRFAMs is based at the J. Craig Venter Institute, Rockville, MD, US.

# **Pfam Database**

**<http://pfam.xfam.org/>**

## **The protein world is vast**

- GenBank RefSeq
- Proteins 45,166,402 proteins

## **Pfam – the Protein Family Database**

The new standard for defining what families a protein belongs to  
Helpful to determine the potential function of a protein

## **Pfam Families**

- Entry types
  - Family
    - A group of proteins that are related by structure
  - Domain
    - A short amino acid sequence
  - Repeat
    - A short amino acid sequence that functions when a group of the sequences are clustered
  - Motif
    - A short amino acid that defines a specific feature

## **Two families**

- Pfam-A
  - High reliable family domain structure
  - Repeated over many proteins and considered to be “real” and associated with function
  - Annotated
- Pfam-B (discontinued in 2016 release)
  - Family domains of lower quality
  - Fewer proteins contain the family domain
  - Not annotated

## Common bean FLS2 ortholog Pfam structure

Sequence search results

Show the detailed description of this results page.

We found **3** Pfam-A matches to your search sequence (**all** significant). You did not choose to search for Pfam-B matches.



Show the search options and sequence that you submitted.

[Return](#) to the search form to look for Pfam domains on a new sequence.

### Output from HMMR3

(Bold lines are those outputted from the Pfam site. Pfam just returns “high” quality hits.)

seq_id	alignment_start	alignment_end	envelope_start	envelope_end	hmm_acc	hmm_name	type	hmm_start	hmm_end	hmm_length	bit_score	E-value
<b>Phvul.002G196200</b>	<b>31</b>	<b>68</b>	<b>28</b>	<b>68</b>	<b>PF08263</b>	<b>LRRNT_2</b>	<b>Family</b>	<b>4</b>	<b>43</b>	<b>43</b>	<b>44.8</b>	<b>7.60E-12</b>
Phvul.002G196200	120	156	120	164	PF12799	LRR_4	Family	1	36	44	25.6	5.90E-06
Phvul.002G196200	218	253	217	259	PF12799	LRR_4	Family	3	37	44	23.3	3.10E-05
Phvul.002G196200	264	324	264	324	PF13855	LRR_8	Repeat	1	61	61	23.9	2.30E-05
Phvul.002G196200	337	373	336	381	PF12799	LRR_4	Family	2	37	44	18.6	0.00096
Phvul.002G196200	385	421	384	427	PF12799	LRR_4	Family	2	37	44	21	0.00016
Phvul.002G196200	432	471	432	481	PF12799	LRR_4	Family	1	39	44	22.6	5.00E-05
Phvul.002G196200	481	518	480	526	PF12799	LRR_4	Family	2	38	44	21.7	0.0001
Phvul.002G196200	578	638	578	638	PF13855	LRR_8	Repeat	3	61	61	29.8	3.20E-07
<b>Phvul.002G196200</b>	<b>654</b>	<b>687</b>	<b>650</b>	<b>695</b>	<b>PF12799</b>	<b>LRR_4</b>	<b>Family</b>	<b>5</b>	<b>37</b>	<b>44</b>	<b>21.8</b>	<b>9.30E-05</b>
Phvul.002G196200	701	759	699	759	PF13855	LRR_8	Repeat	3	61	61	35.5	5.50E-09
<b>Phvul.002G196200</b>	<b>870</b>	<b>1153</b>	<b>867</b>	<b>1156</b>	<b>PF00069</b>	<b>Pkinase</b>	<b>Domain</b>	<b>4</b>	<b>256</b>	<b>260</b>	<b>145.2</b>	<b>1.70E-42</b>

[The Pfam protein families database](#): R.D. Finn, A. Bateman, J. Clements, P. Coghill, R.Y. Eberhardt, S.R. Eddy, A. Heger, K. Hetherington, L. Holm, J. Mistry, E.L.L. Sonnhammer, J. Tate, M. Punta Nucleic Acids Research (2014) Database Issue 42:D222-D230



# Phytozome: Integrates Multiple Data Into Display

## *Phaseolus vulgaris* FLS2 gene output

Phytozome v12.1: Gene

Secure | <https://phytozome.jgi.doe.gov/pz/portal.html#?gene?method=0&crow=1&search=1&detail=1&searchText=transcriptid:37176694>

JGI Phytozome 12 THE PLANT GENOMICS RESOURCE

JGI HOME LOG IN

Species Tools Info Download Help Cart Subscribe

Previous view Help with this page

Actions

- Revise query
- Launch Jalview
- Find related ...
- Add to cart
- Composite family

My Data (0)

- View cart
- Add to cart
- Upload user data
- Send to BioMart
- Send to PhytoMine
- Get from PhytoMine
- Quick download
- Delete data

Settings

- Species display
- Family filter
- Homolog filter

### Gene Phvul.002G196200

**Gene Info**

**Organism** Phaseolus vulgaris

**Transcript Name** Phvul.002G196200.1 (primary)

**Location** Chr02:36081370..36086171 forward

**Alias** Phvul.002G196200.v1.0 Phvul.002G196200.1.v1.0

**Description** (1 of 1) K13420 - LRR receptor-like serine/threonine-protein kinase FLS2 (FLS2)

**Links** [B](#) [M](#)

Functional Annotation Genomic Sequences Protein Homologs Gene Ancestry Expression

#### Protein domain view

Functional annotations for this locus

ID	Type	Description
<input type="checkbox"/> PTHR27000	PANTHER	FAMILY NOT NAMED
<input type="checkbox"/> PTHR27000:SF16	PANTHER	LRR RECEPTOR-LIKE SERINE/THREONINE-PROTEIN KINASE FLS2
<input type="checkbox"/> PF00069	PFAM	Protein kinase domain
<input type="checkbox"/> PF00560	PFAM	Leucine Rich Repeat
<input type="checkbox"/> PF08263	PFAM	Leucine rich repeat N-terminal domain
<input type="checkbox"/> PF13855	PFAM	Leucine rich repeat
<input type="checkbox"/> 2.7.11.1	EC	Non-specific serine/threonine protein kinase
<input type="checkbox"/> KOG1187	KOG	Serine/threonine protein kinase
<input type="checkbox"/> GO:0004672	GO	Catalysis of the phosphorylation of an amino acid residue in a protein, usually according to the reaction: a protein + ATP = a phosph...
<input type="checkbox"/> GO:0005515	GO	Interacting selectively and non-covalently with any protein or protein complex (a complex of two or more proteins that may include ...
<input type="checkbox"/> GO:0005524	GO	Interacting selectively and non-covalently with ATP, adenosine 5'-triphosphate, a universally important coenzyme and enzyme regul...

Contact Disclaimer  
Accessibility / Section 508 Statement  
©1997-2017 The Regents of the University of California

U.S. DEPARTMENT OF ENERGY Office of Science

Start Desktop 1:50 PM 10/4/2017

**InterPro statistics for sequenced eukaryote genomes. (10/12/03 data from InterPro; no longer available from InterPro)**

	# Proteins in proteome	Proteins with InterPro matches (%)	Number of signatures	Number of InterPro entries (%)
<i>A. thaliana</i>	26,070	20,685 (70.3)	5629	3119 (36.5)
<i>C. elegans</i>	21,128	15,090 (71.4)	5402	2956 (34.5)
<i>D. melanogaster</i>	15,455	11,267 (72.9)	5536	3035 (35.5)
<i>E. cuniculi</i> <sup>1</sup>	1,908	1,258 (65.9)	1765	927 (10.9)
<i>G. theta</i> <sup>2</sup>	451	279 (61.9)	596	295 (5.3)
<i>H. sapiens</i>	29,638	22,036 (74.4)	7296	4163 (48.7)
<i>M. musculus</i>	21,041	16,652 (79.1)	7123	4053 (47.4)
<i>S. cerevisiae</i>	6,202	4,379 (70.6)	4436	2347 (27.5)
<i>S. pombe</i>	4,988	3,889 (78.0)	4839	2335 (27.3)

<sup>1</sup>*Encephalitozoon cuniculi*: microsporidian; primitive eukaryote; lack mitochondria

<sup>2</sup>*Guillarida theta* enslaved nucleus: nucleus containing chloroplasts are found in some unicellular algae.

# InterPro Data

## Arabidopsis Top 15 Entries (10/2003)

<b>Rank</b>	<b>No. proteins (% proteome)</b>	<b>InterPro Entry (type)</b>
1	1050 (4.0)	Protein kinase (family)
2	796 (3.1)	Serine/threonine protein kinase active site (active site)
3	648 (2.5)	Cyclin-like F-box (domain)
4	518 (2.0)	Leucine-rich repeat (repeat)
5	491 (1.9)	Zn-finger, RING (domain)
6	454 (1.7)	PPR repeat (Repeat)
7	345 (1.3)	Myb DNA-binding protein (domain)
8	342 (1.3)	Tyrosine protein kinase (family)
9	332 (1.3)	Leucine-rich repeat, plant specific (repeat)
10	314 (1.2)	AAA ATPase (domain)
11	300 (1.2)	RNA-binding region RNP-1 (RNA recognition motif) (domain)
12	267 (1.0)	G-protein beta WD40 repeat (repeat)
13	252 (1.0)	Cytochrome P450 (family)
14	233 (0.9)	Zn-finger, CCHC type (domain)
15	227 (0.9)	E-class P450, group I (family)

## Arabidopsis Top 15 Families (10/2003)

<b>Rank</b>	<b>No. proteins</b>	<b>InterPro Entry</b>
1	1050	Protein kinase
2	252	Cytochrome P450
3	157	Disease resistance protein
4	114	2OG-Fe(II) oxygenase superfamily
5	110	Retrotransposon gag protein
6	110	Lipolytic enzyme, G-DE-S-L family
7	109	UDG-glucuronosyl/UDP-glucosyl transferase
8	108	No apical meristem (NAM) protein
9	106	Ras GTPase superfamily
10	96	Haem peroxidase
11	91	Short-chain dehydrogenase/reductase SDR
12	91	General substrate transporter
13	85	Mitochondrial substrate carrier
14	84	Major facilitator superfamily
15	81	Transcriptional factor B3

## Arabidopsis Top 15 Domains (10/2003)

<b>Rank</b>	<b>Proteins matched (Matches/genome)</b>	<b>InterPro Entry (type)</b>
1	648 (1487)	Cyclin-like F-box
2	491 (1287)	Zn-finger, RING
3	345 (1535)	Myb DNA-binding domain
4	314 (390)	AAA ATPase
5	300 (1302)	RNA-binding region RNP-1 (RNA recognition motif
6	233 (923)	Zn-finger, CCHC type
7	224 (647)	Basic helix-loop-helix dimerization domain bHLH
8	222 (1546)	Calcium-binding EF-hand
9	213 (223)	F-box protein interaction domain
10	210 (213)	Esterase/lipase/thioesterase
11	196 (852)	Zn-finger,C2H2 type
12	164 (472)	Zn-finger-like, PHD finger
13	159 (192)	FBD
14	154 (157)	NB-ARC domain
15	153 (158)	Integrase, catalytic domain

## Arabidopsis Top 15 Repeats (10/2003)

Rank	No. proteins	InterPro Entry
1	518	Leucine-rich repeat
2	454	PPR repeat
3	267	G-protein beta WD-40 repeat
4	136	TPR repeat
5	117	Ankyrin
6	115	Kelch repeat
7	70	Armadillo repeat
8	62	Parallel beta-helix repeat
9	28	Paired amphipathic helix
10	26	Bacterail transferase hexapeptide repeat
11	26	Leucine-rich repeat, cysteine-containing
12	24	Protein of unknown function DUF321
13	20	Prenyltransferase/squalene oxidase
14	18	Late embryogenesis abundant protein
15	17	HEAT repeat

## Are their plant specific families?

Rank	InterPro entry	<i>A. thaliana</i> # proteins	<i>C. elegans</i> # proteins	<i>D. melanogaster</i> # proteins	<i>H. sapiens</i> # proteins	<i>S. cerevisiae</i> # proteins
1	Protein kinase	1050	528	303	694	118
2	Cytochrome P450	252	83	88	90	3
3	<b><i>Disease resistance protein</i></b>	<b>157</b>	<b>1</b>	<b>0</b>	<b>4</b>	<b>2</b>
4	2OG-Fe(II) oxygenase superfamily	114	11	26	19	1
5	<b><i>Retrotransposon gag protein</i></b>	<b>110</b>	<b>5</b>	<b>3</b>	<b>4</b>	<b>3</b>
6	Lipolytic enzyme, G-DE-S-L family	110	11	2	5	1
7	UDG-glucuronosyl/UDP-glucosyl transferase	109	76	35	29	0
8	<b><i>No apical meristem (NAM) protein</i></b>	<b>108</b>	<b>0</b>	<b>0</b>	<b>0</b>	<b>0</b>
9	Ras GTPase superfamily	106	82	86	118	34
10	Haem peroxidase	96	26	16	32	4
11	Short-chain dehydrogenase/reductase SDR	91	92	60	80	13
12	General substrate transporter	91	96	118	77	54
13	Mitochondrial substrate carrier	85	65	74	93	38
14	Major facilitator superfamily	84	128	123	106	670
15	<b><i>Transcriptional factor B3</i></b>	<b>81</b>	<b>0</b>	<b>0</b>	<b>0</b>	<b>0</b>

## Are their plant specific domains? [No. proteins (Matches/proteome)]

Rank	InterPro entry	<i>A. thaliana</i> # proteins	<i>C. elegans</i> # proteins	<i>D. melanogaster</i> # proteins	<i>H. sapiens</i> # proteins	<i>S. cerevisiae</i> # proteins
1	Cyclin-like F-box	648 (1487)	441 (957)	37 (78)	86 (204)	14 (32)
2	Zn-finger, RING	491 (1287)	175 (504)	150 (376)	387 (1194)	40 (108)
3	Myb DNA-binding domain	345 (1535)	39 (112)	50 (103)	92 (276)	21 (91)
4	AAA ATPase	314 (390)	124 (173)	131 (193)	176 (251)	89 (125)
5	RNA-binding region RNP-1 (RNA recognition motif)	300 (1302)	131 (193)	176 (251)	346 (1809)	58 (307)
6	Zn-finger, CCHC type	233 (923)	47 (210)	27 (111)	56 (308)	14 (89)
7	Basic helix-loop-helix dimerization domain bHLH	224 (647)	101 (271)	93 (307)	168 (635)	9 (36)
8	Calcium-binding EF-hand	222 (1546)	133 (813)	134 (873)	316 (1945)	17 (102)
<b>9</b>	<b><i>F-box protein interaction domain</i></b>	<b>213 (223)</b>	<b>0 (0)</b>	<b>0 (0)</b>	<b>0 (0)</b>	<b>0 (0)</b>
10	Esterase/lipase/thioesterase	210 (213)	134 (142)	148 (149)	121 (121)	37 (37)
11	Zn-finger,C2H2 type	196 (852)	251 (2042)	397 (7249)	962 (29,334)	54 (381)
12	Zn-finger-like, PHD finger	164 (472)	46 (146)	51 (203)	126 (529)	16 (57)
<b>13</b>	<b><i>FBD</i></b>	<b>159 (192)</b>	<b>1 (1)</b>	<b>2 (2)</b>	<b>1 (1)</b>	<b>0 (0)</b>
<b>14</b>	<b><i>NB-ARC domain</i></b>	<b>154 (157)</b>	<b>0 (0)</b>	<b>0 (0)</b>	<b>0 (0)</b>	<b>0 (0)</b>
15	Integrase, catalytic domain	153 (158)	30 (31)	6 (6)	15 (15)	40 (40)



## Are their plant specific repeats?

Rank	InterPro entry	<i>A. thaliana</i> # proteins	<i>C. elegans</i> # proteins	<i>D. melanogaster</i> # proteins	<i>H. sapiens</i> # proteins	<i>S. cerevisiae</i> # proteins
1	Leucine-rich repeat	518	85	127	269	11
<b>2</b>	<b><i>PPR repeat</i></b>	<b>454</b>	<b>1</b>	<b>4</b>	<b>5</b>	<b>2</b>
3	G-protein beta WD-40 repeat	267	170	195	405	100
4	TPR repeat	136	64	92	189	33
5	Ankyrin	117	116	104	306	21
6	Kelch repeat	115	17	22	97	6
7	Armadillo repeat	70	7	11	47	2
<b>8</b>	<b><i>Parallel beta-helix repeat</i></b>	<b>62</b>	<b>3</b>	<b>2</b>	<b>9</b>	<b>1</b>
9	Paired amphipathic helix	28	1	3	3	1
10	Bacterail transferase hexapeptide repeat	26	7	5	8	4
11	Leucine-rich repeat, cysteine-containing	26	7	10	21	3
<b>12</b>	<b><i>Protein of unknown function DUF321</i></b>	<b>24</b>	<b>0</b>	<b>0</b>	<b>0</b>	<b>0</b>
13	Prenyltransferase/squalene oxidase	20	3	3	4	4
<b>14</b>	<b><i>Late embryogenesis abundant protein</i></b>	<b>18</b>	<b>4</b>	<b>1</b>	<b>2</b>	<b>0</b>
15	HEAT repeat	17	14	6	30	10

## Examples of Arabidopsis Genes With Known Phenotypic Functions (Same structure; many functions)

<i>Family</i>	<i>Gene</i>	<i>Phenotype</i>
Protein kinase	<i>YDA</i> (At1g63700)	Defective embryo and seedling
Cytochrome P450	<i>REF8</i> (At2g40890)	Reduced phenylpropanoid metabolites and lignin
Disease resistance (LRR)	<i>RPM1</i> (At3g07040)	Resistance to bacterial pathogens
Glycosyl Transferase	<i>QUA1</i> (At3g25410)	Dwarf plant reduced pectin
<b><i>Domain</i></b>		
Myb-transcription factor	<i>AS1</i> (At2g37630)	Altered leaf morphology
	<i>CPC</i> (At2g46410)	Reduced root hairs
	<i>LAF1</i> (At4g25560)	Elongated hypocotyls in far-red light
Basic helix-loop-helix	<i>HFR1</i> (At1g02340)	Elongated hypocotyls in far-red light
	<i>SPT</i> (At4g36930)	Altered carpel development
<b><i>Repeats</i></b>		
Leucine-rich repeat	<i>BRL1</i> (At1g55610)	Altered vascular cell differentiation
	<i>RPP1</i> (At3g44480)	Fungal disease resistance
Ankyrin	<i>EMB506</i> (At5g40160)	Embryo defective
	<i>CAO</i> (At2g47450)	Pale, chlorotic plants
TPR	<i>HBT</i> (At2g20000)	Altered root meristem development
	<i>SPY</i> (At3g11540)	Spindly plants
WD40	<i>TTG1</i> (At5g24520)	Yellow seed coat

from: Meinke, D.W., L.K. Meinke, T.C. Showalter, A.M. Schissel, L.A. Mueller, and I. Tzafrir. 2003. A sequence-based map of Arabidopsis genes with mutant phenotypes. *Pl. Physiol.* 131:409-418.

## **Gene Ontology (GO) (<http://www.geneontology.org/>)**

“The goal of the Gene Ontology Consortium is to produce a controlled vocabulary that be applied to all organisms even as knowledge of gene and protein roles in cells is accumulating and changing.”

### **Major Categories**

#### ***Molecular Function Ontology***

“The action characteristic of a gene product.”

#### ***Biological Process Ontology***

A phenomenon marked by changes that lead to a particular result, mediated by one or more gene products.

#### ***Cellular Component Ontology***

“The part of a cell of which a gene product is a component; for purpose of GO includes the extracellular environment of cells; a gene product may be a component of one or more parts of a cell; this term includes gene products that are parts of macromolecular complexes, by the definition that all members of a complex normally copurify under all except extreme conditions.”

## GO Annotations for Eukaryotic Species

<b>Species</b>	<b>Biological process</b>	<b>Molecular function</b>	<b>Cellular component</b>	<b>Total genes associated</b>
<i>A. thaliana</i>	6,712	7,967	13,532	18,572
<i>C. elegans</i>	5,115	5,755	3,057	6,925
<i>D. melanogaster</i>	4,439	6,795	3,942	7,938
<i>H. sapiens</i>	16,643	18,772	13,880	20,687
<i>S. cerevisiae</i>	6,646	6,434	6,435	6,448

Note: Each of these annotations was performed by a different group.

## Lineage Specific Expansion

(<http://www.ncbi.nlm.nih.gov/COG/new/sholse.cgi>)

Uses the “Clusters of Orthologous Genes”

Determines how many genes are specific to only one species

<b>Species</b>	<b>LSEs</b>
<i>A. thaliana</i>	1,458
<i>C. elegans</i>	845
<i>D. melanogaster</i>	303
<i>H. sapiens</i>	1,373
<i>S. cerevasiae</i>	132

### *Arabidopsis* LSE

<b>Category</b>	<b>Unique</b>	<b>Total</b>
Signal transduction mechanisms (T)	34	1080
Defense mechanisms (V)	15	337

## ***Arabidopsis* Example**

- Disease resistance proteins
  - Shared structure among each protein
    - Domains
      - Coiled-coil receptor site
      - TIR:Toll-interleukin receptor site
      - NBS: nucleotide binding site
      - LRR: leucine rich repeat
    - Structure shared among many plant resistance genes
  - Lineage specific for *Arabidopsis* relative to other sequenced genome
    - But not lineage specific for plants in general
  - Number
    - 105

## **Other *Arabidopsis* Examples**

- Receptor protein kinase containing LRR repeats
  - Gene expression regulator
    - Involved in phosphorylating other proteins
  - Number
    - 227
- bHL
  - **b**asic **H**elix **L**oop **H**elix
    - Transcriptional regulator
  - Number
  - 143

## Human Example

- Olfactory receptor
  - Smell
  - Number
    - 570
  
- Immunoglobulin heavy chain like
  - Immune system protein
    - Defense response
  - Number
    - 156