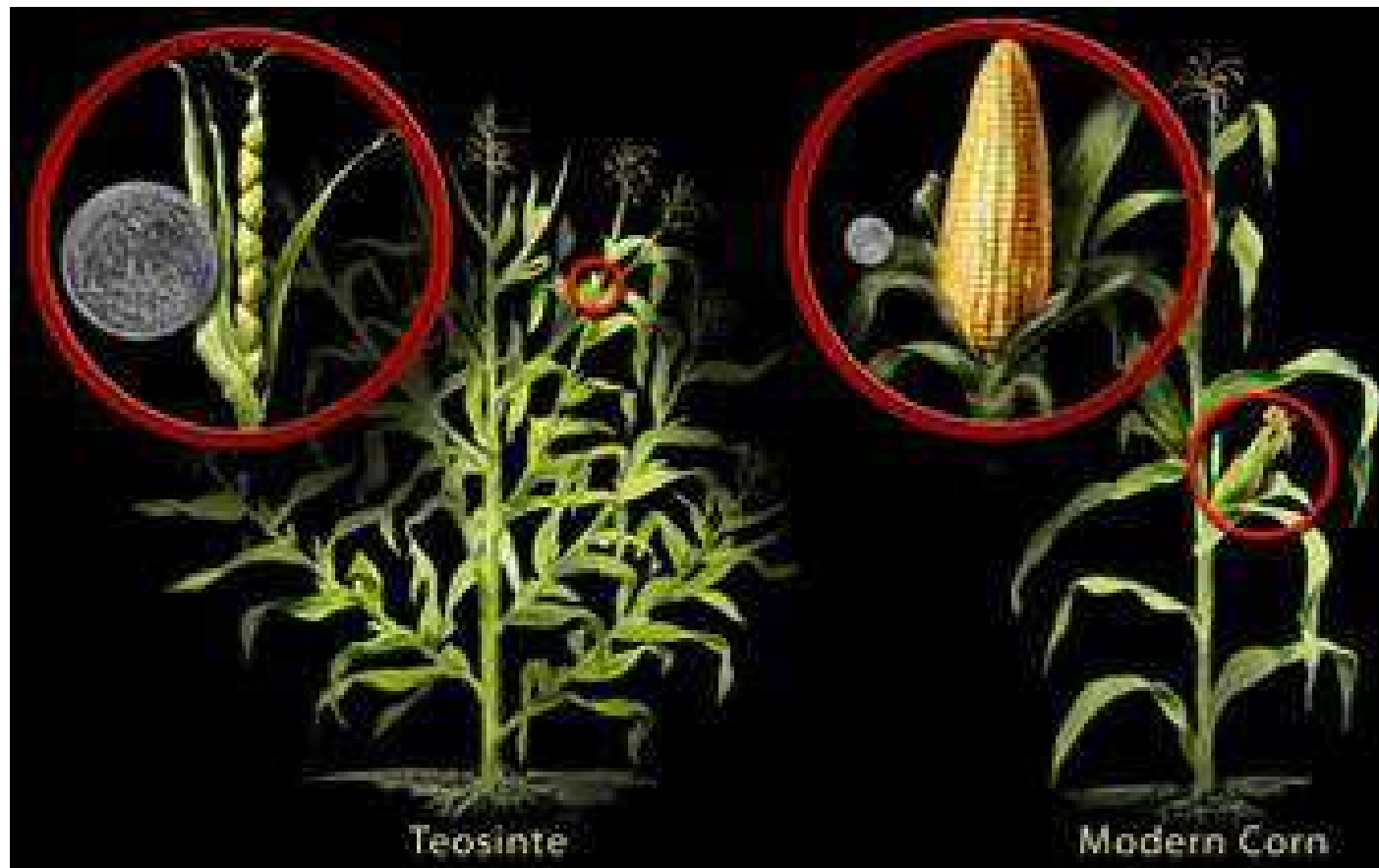
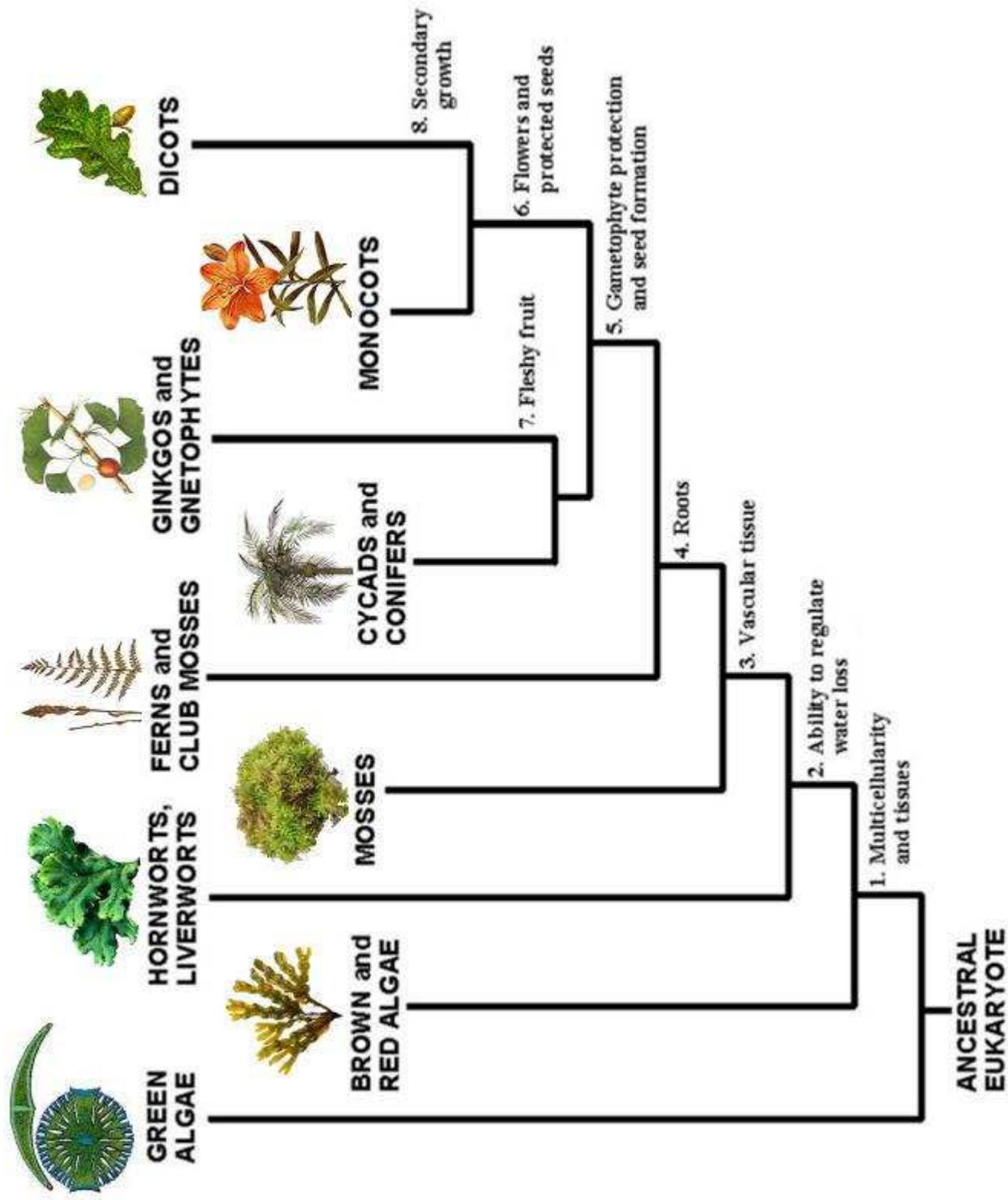


STUDIO DEI GENOMI VEGETALI

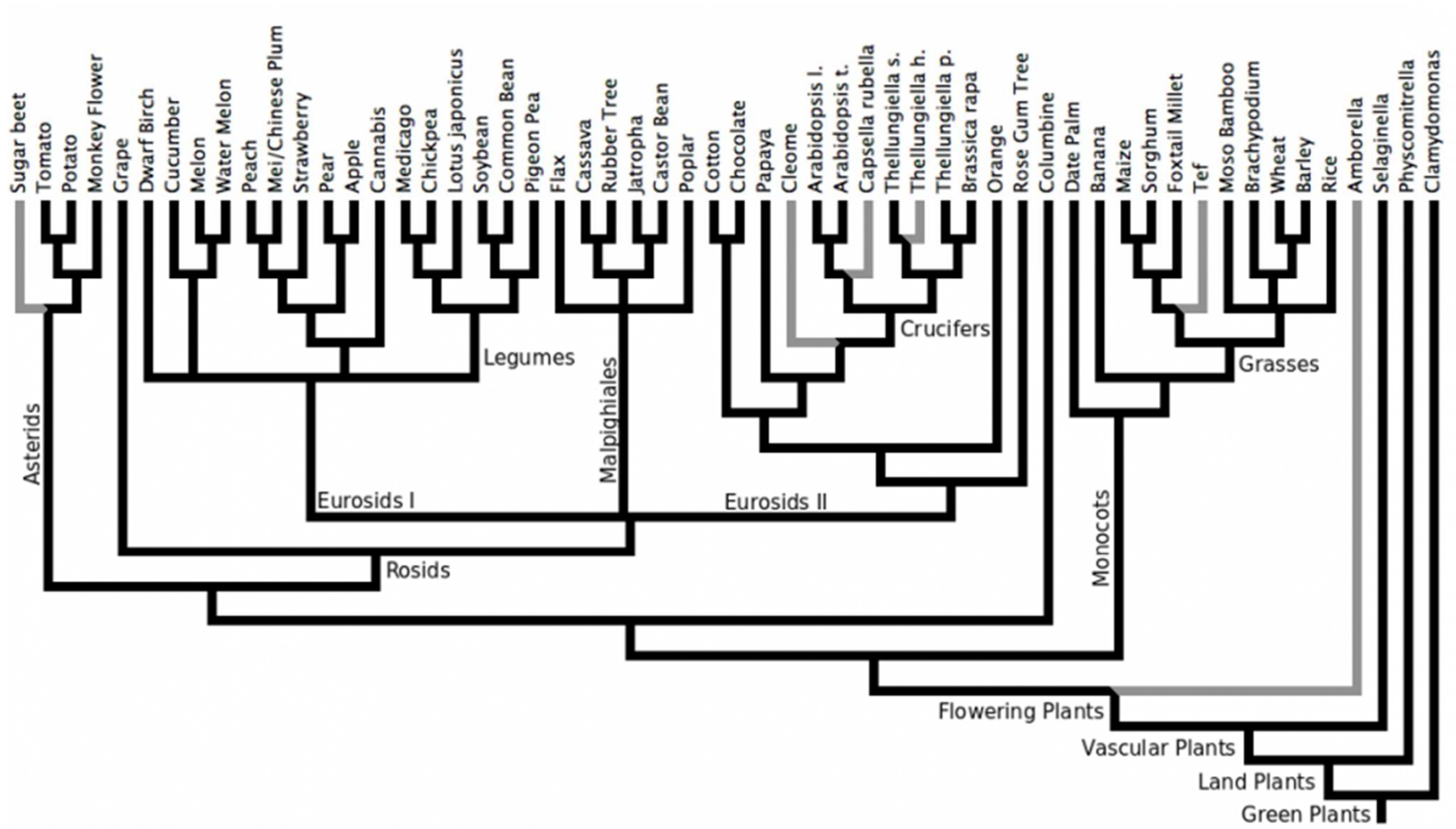
- **Identificazione di geni importanti per caratteri agronomici (produttività, resistenza a stress, proprietà nutrizionali)**
- **Comprensione dell'evoluzione delle piante**

Addomesticamento delle specie vegetali





http://genomeevolution.org/wiki/index.php/Sequenced_plant_genomes



Quali specie sequenziare?

- Impatto economico, sociale e scientifico
- Distanza filogenetica da altre specie sequenziate (-> nuove informazioni)
- Informazioni disponibili (mappe genetiche e fisiche)
- Capacità di persuasione dei ricercatori

Arabidopsis thaliana

Dicotiledone (*Brassicaceae*)

Piccolo genoma dipolide (C1 = 125 Mbp)

Trasformabile facilmente

5 cromosomi

Piccole dimensioni

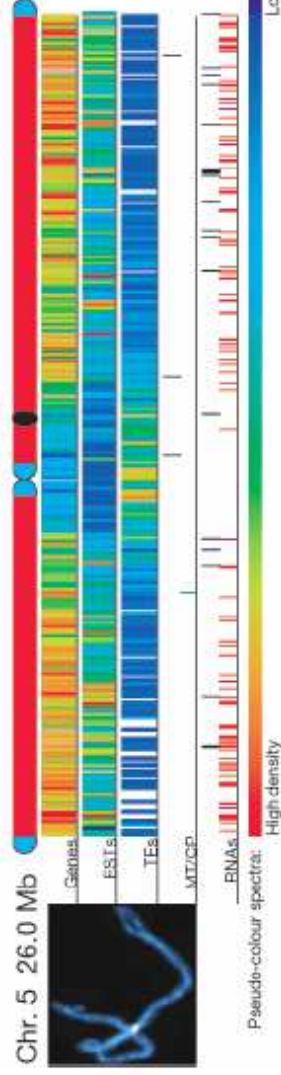
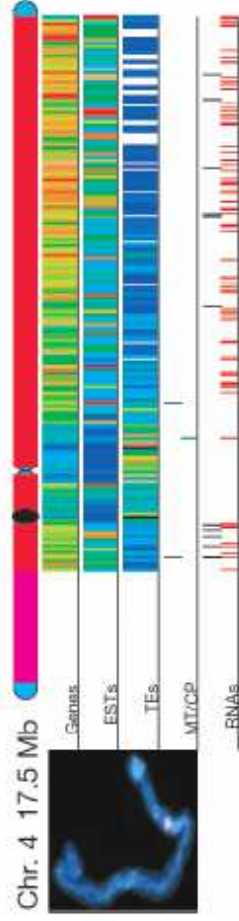
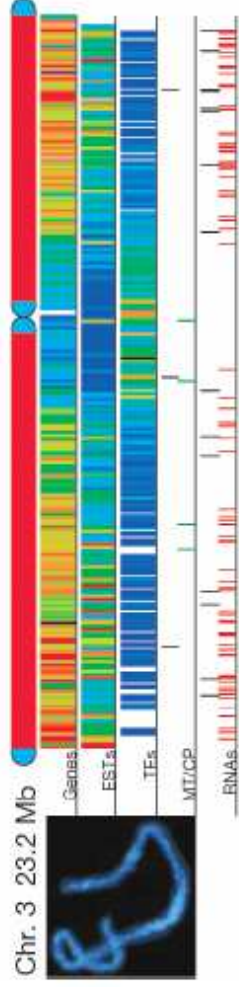
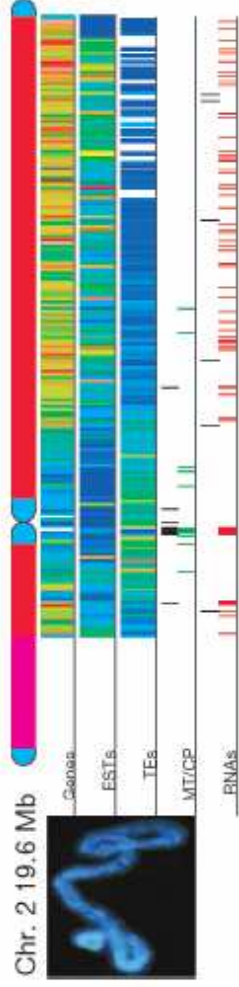
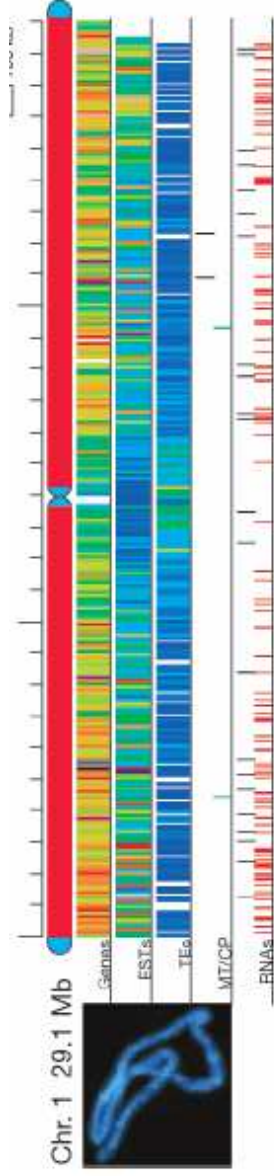
Ciclo vitale breve (2 mesi)

Primo genoma sequenziato

The Arabidopsis Genome Initiative (2000).

**Analysis of the genome sequence of the
flowering plant *Arabidopsis thaliana*. *Nature*,
408 (6814), 796-815 DOI: 10.1038/35048692**

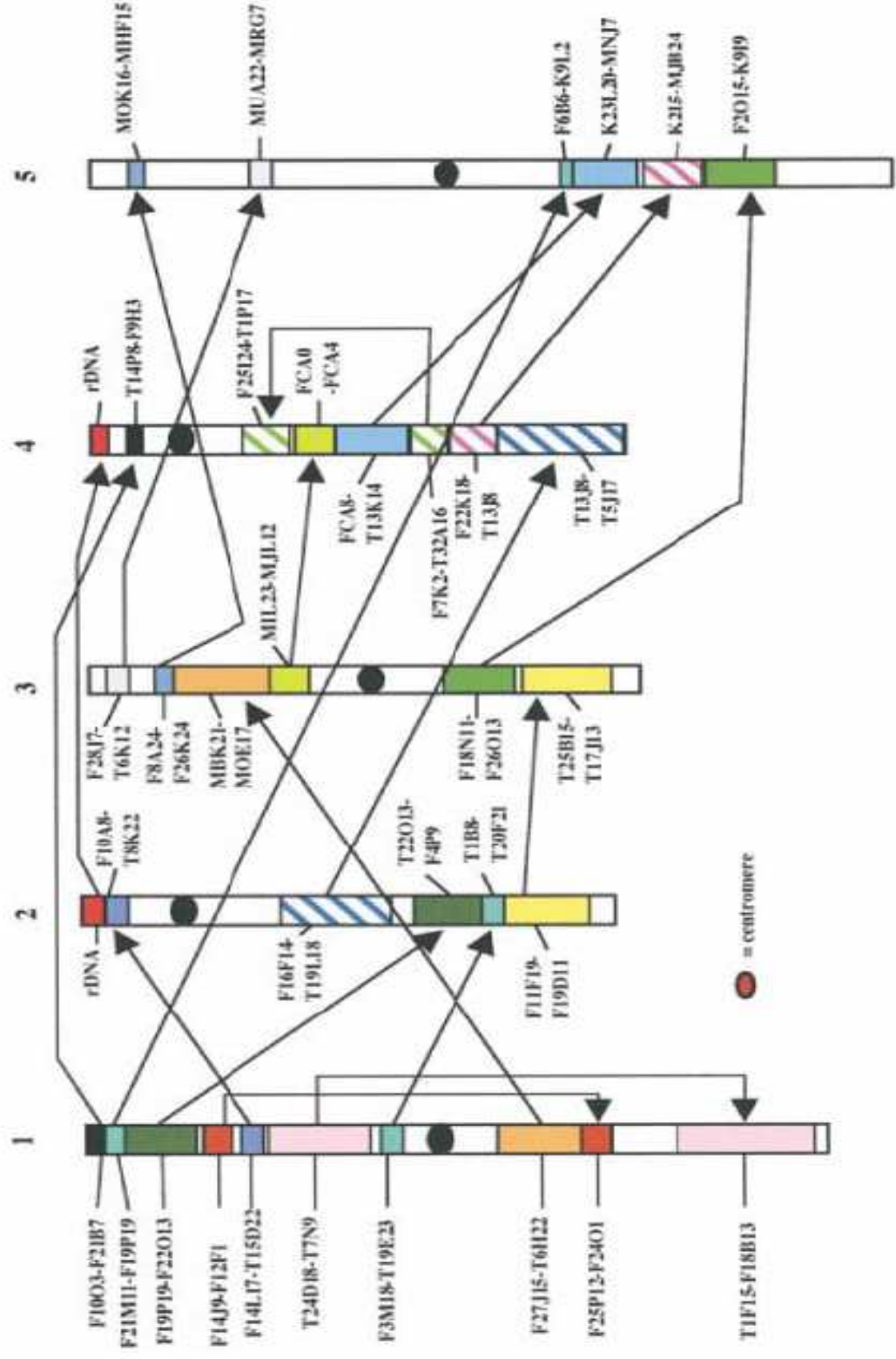




Arabidopsis genome sequence.
As published in 2000.
Nature 408, 796.

115 Mb of 125 Mb genome.
Gene annotation using
Expressed sequence tags (ESTs)
Homology with cloned plant genes
and genes of other organisms
Identified 25,500 genes.

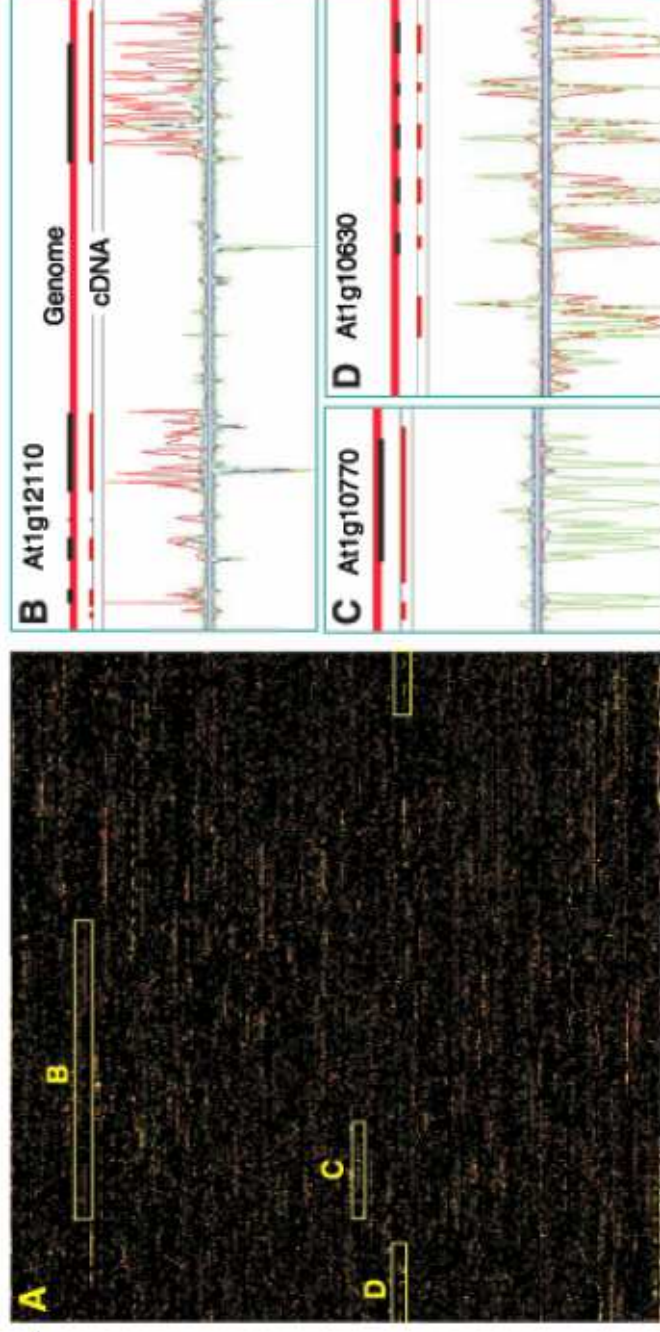
Large segments of the Arabidopsis genome are duplicated



Improved annotation of the *Arabidopsis* genome sequence

Initial sequence analysis relied heavily on expressed sequence tags (ESTs) and gene prediction programmes. Interpretation of genome sequence constantly improved.

Isolation of full-length cDNAs and use of whole genome oligonucleotide tiling arrays greatly improved the annotation. (32% of predicted gene models incorrect).



Arabidopsis genome represented on 12 oligonucleotide arrays. Each array 834,000 25-mer oligonucleotides. Is hybridized with cDNA made from 4 RNA populations; seedlings, roots, flowers, cultured cells. Sequences that hybridize are mapped directly onto genome sequence.

Enzymes involved in secondary metabolism

Arabidopsis genome contains many classes of enzymes involved in secondary metabolism that are required for the synthesis of specialized compounds.

An example, is the family of genes encoding the Cytochrome P450 monooxygenase enzymes.

Mammals, *C.elegans*, *Drosophila* – 80 – 105 genes.

Arabidopsis – 246 genes.

In plants these enzymes are required for the synthesis of compounds such as growth regulators (gibberellic acid, Brassinosteroid), carotenoids (protect cell from oxidative damage) and phenylpropanoids that are present in plant cell walls.

Transcription factors

Arabidopsis contains around 1500 genes encoding transcription factors (aprox. 5%)

Drosophila contains around 640 genes encoding transcription factors, around 4.5%.

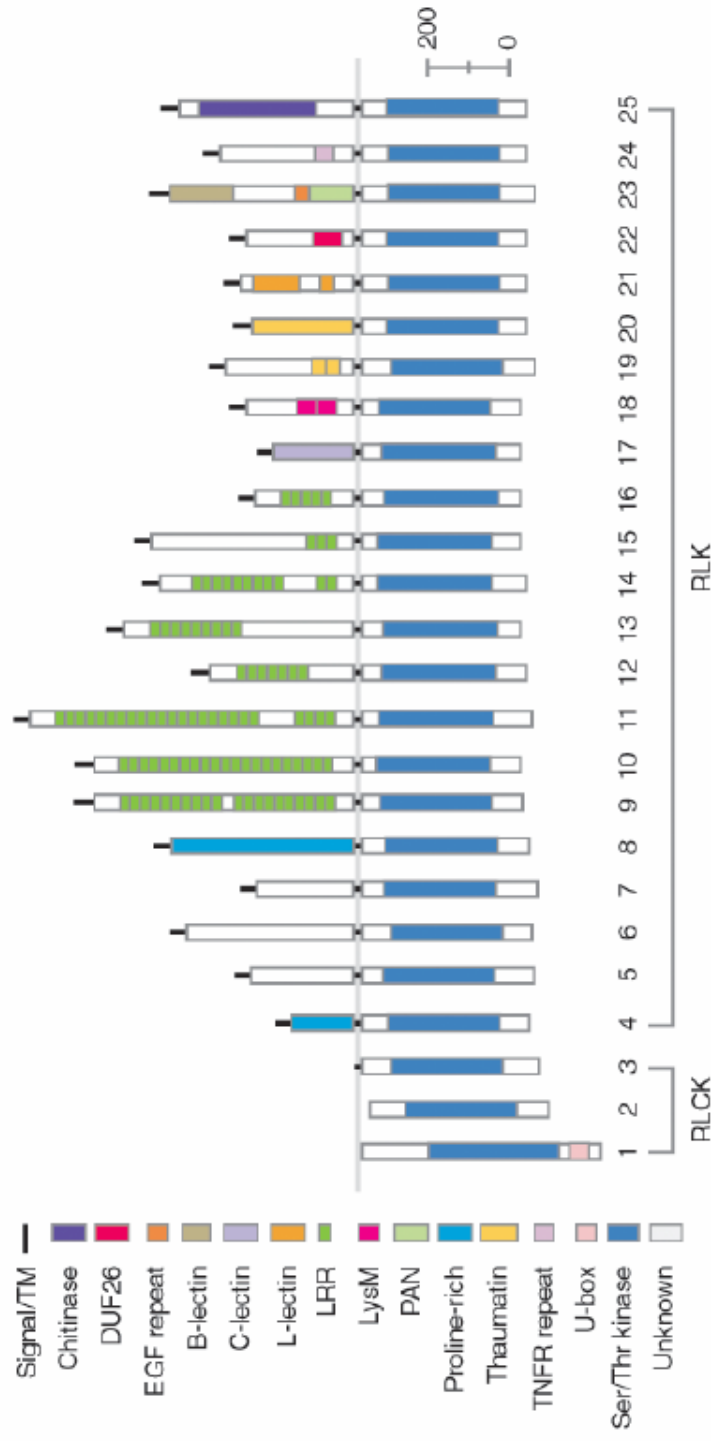
Many important animal transcription factor families are absent in plants, such as nuclear steroid receptors, NHR zinc finger proteins (252 in *C. Elegans*) and Fork head transcription factors (18 in *Drosophila*, 15 in *C.elegans*).

Each eukaryotic lineage has its own set of transcription factor families.

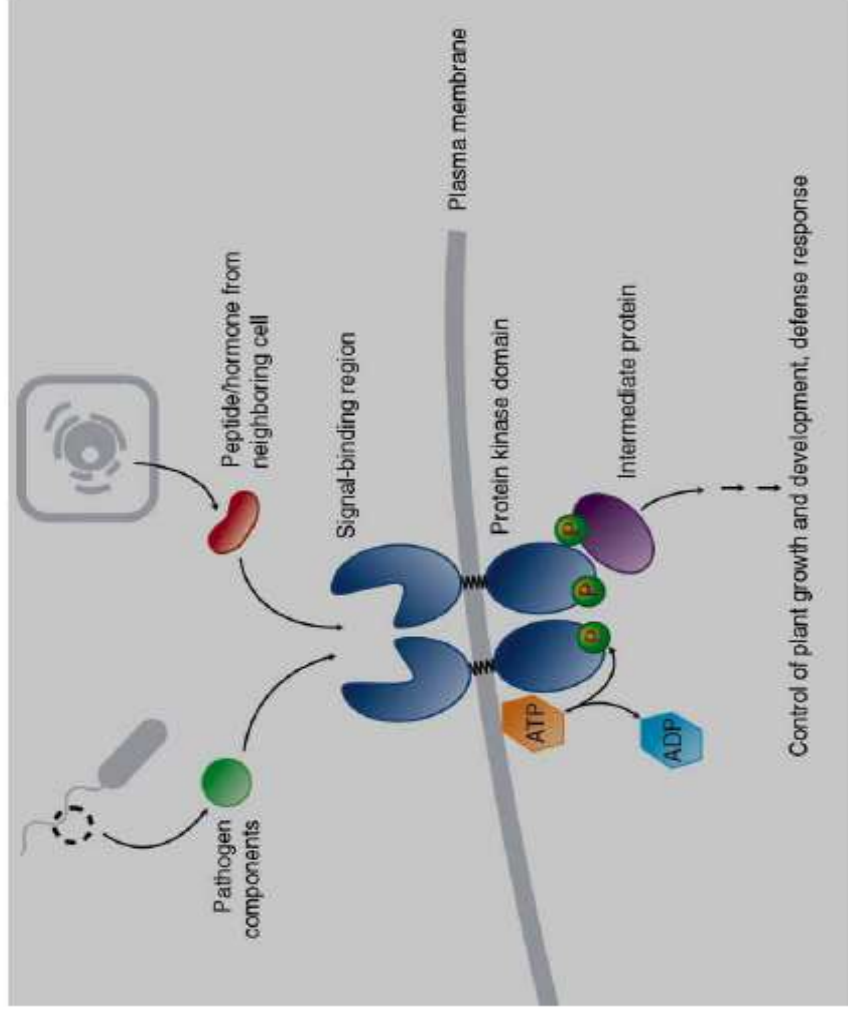
Receptor-like kinases

600 Arabidopsis genes encode receptor-like kinases predicted to be located in the membrane.

These are similar in domain organization to animal receptor tyrosine kinases, such as epidermal growth factor. However, are predicted to be serine/threonine kinases, and have divergent ligand binding domains.



Some Receptor Like Kinases have important functions, but for most their function is unknown



Many RLKs are of unknown function. However, some have defined roles:

**Brassinosteroid receptor
Clavata 1
Resistance to pathogens**

Other plant-specific processes

Hundreds of genes involved in photosynthesis

- light harvesting
- chlorophyll biosynthesis
- carbon dioxide fixation
- energy generating photosystems

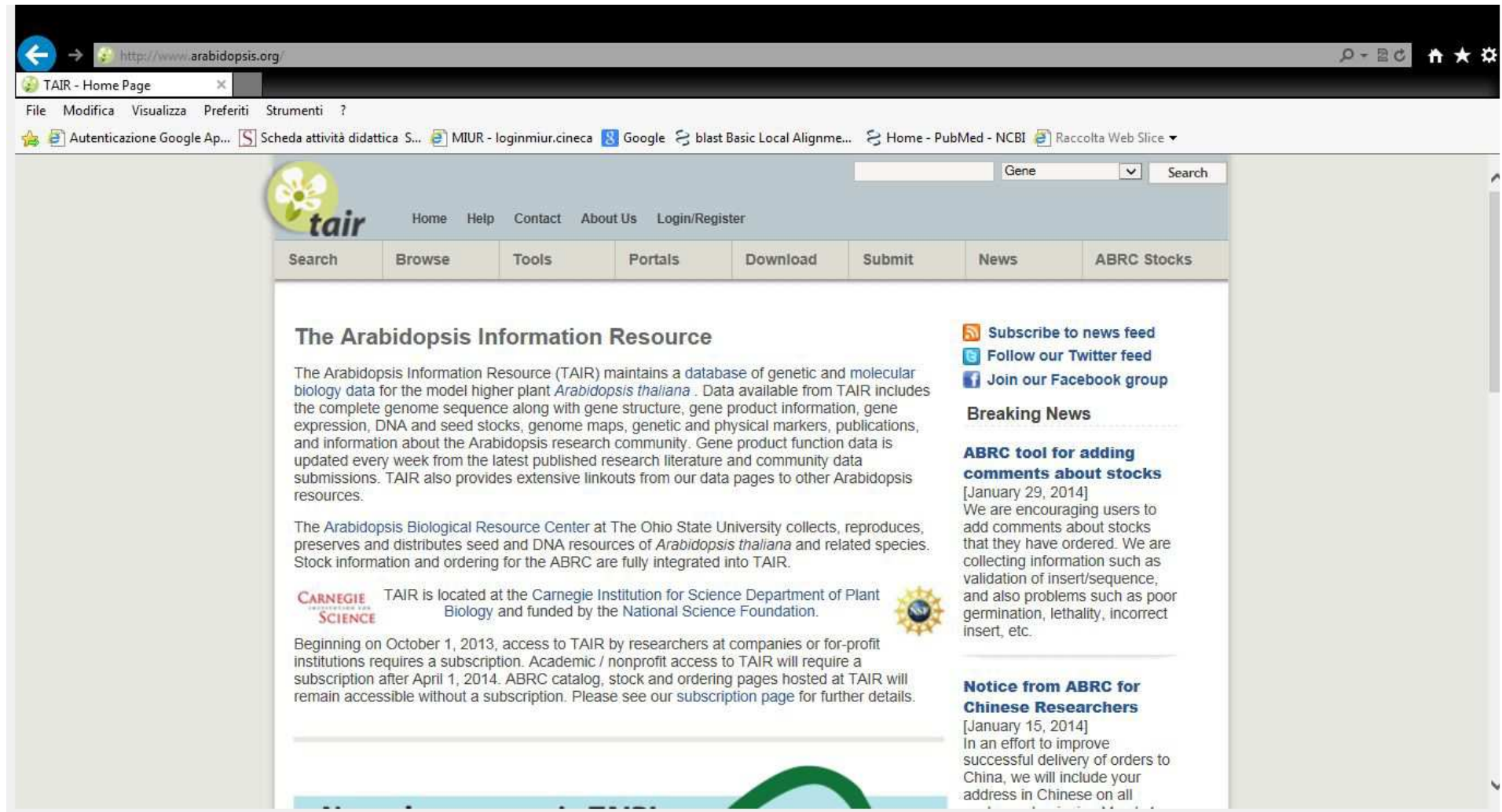
Transporters

- plants mainly use proton-type ATPases, whereas in animals transport is usely coupled with sodium ions via sodium-type ATPases.

**No major histocompatibility complex, however 100s of genes
Encoding nucleotide binding site leucine rich repeat proteins
Involved in pathogen resistance.**

The Arabidopsis Information Resource (TAIR)

<http://www.arabidopsis.org/>



The screenshot shows the TAIR website homepage in a browser window. The browser's address bar displays <http://www.arabidopsis.org/>. The page features a navigation menu with links for Home, Help, Contact, About Us, and Login/Register. Below this is a search bar with a dropdown menu labeled 'Gene' and a 'Search' button. A secondary navigation bar includes links for Search, Browse, Tools, Portals, Download, Submit, News, and ABRC Stocks. The main content area is titled 'The Arabidopsis Information Resource' and contains several paragraphs of text. On the right side, there are social media links for RSS, Twitter, and Facebook, followed by a 'Breaking News' section with two news items. The bottom of the page features a decorative green and blue graphic.

TAIR - Home Page

File Modifica Visualizza Preferiti Strumenti ?

Autenticazione Google Ap... | Scheda attività didattica S... | MIUR - loginmiur.cineca | Google | blast Basic Local Alignme... | Home - PubMed - NCBI | Raccolta Web Slice

Gene Search

Home Help Contact About Us Login/Register

Search Browse Tools Portals Download Submit News ABRC Stocks

The Arabidopsis Information Resource

The Arabidopsis Information Resource (TAIR) maintains a database of genetic and molecular biology data for the model higher plant *Arabidopsis thaliana*. Data available from TAIR includes the complete genome sequence along with gene structure, gene product information, gene expression, DNA and seed stocks, genome maps, genetic and physical markers, publications, and information about the Arabidopsis research community. Gene product function data is updated every week from the latest published research literature and community data submissions. TAIR also provides extensive linkouts from our data pages to other Arabidopsis resources.

The Arabidopsis Biological Resource Center at The Ohio State University collects, reproduces, preserves and distributes seed and DNA resources of *Arabidopsis thaliana* and related species. Stock information and ordering for the ABRC are fully integrated into TAIR.

CARNEGIE INSTITUTION FOR SCIENCE TAIR is located at the Carnegie Institution for Science Department of Plant Biology and funded by the National Science Foundation.

Beginning on October 1, 2013, access to TAIR by researchers at companies or for-profit institutions requires a subscription. Academic / nonprofit access to TAIR will require a subscription after April 1, 2014. ABRC catalog, stock and ordering pages hosted at TAIR will remain accessible without a subscription. Please see our subscription page for further details.

Subscribe to news feed

Follow our Twitter feed

Join our Facebook group

Breaking News

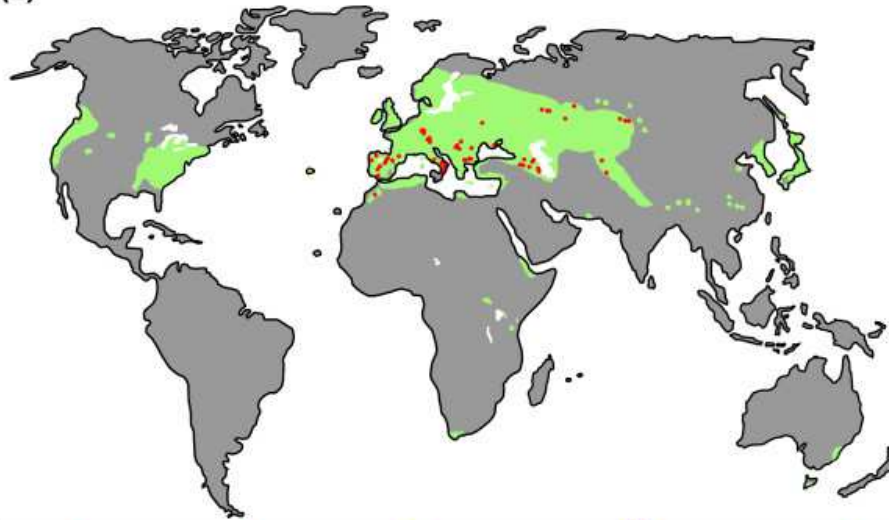
ABRC tool for adding comments about stocks
[January 29, 2014]
We are encouraging users to add comments about stocks that they have ordered. We are collecting information such as validation of insert/sequence, and also problems such as poor germination, lethality, incorrect insert, etc.

Notice from ABRC for Chinese Researchers
[January 15, 2014]
In an effort to improve successful delivery of orders to China, we will include your address in Chinese on all

1001 Genomes : A Catalog of *Arabidopsis thaliana* Genetic Variation

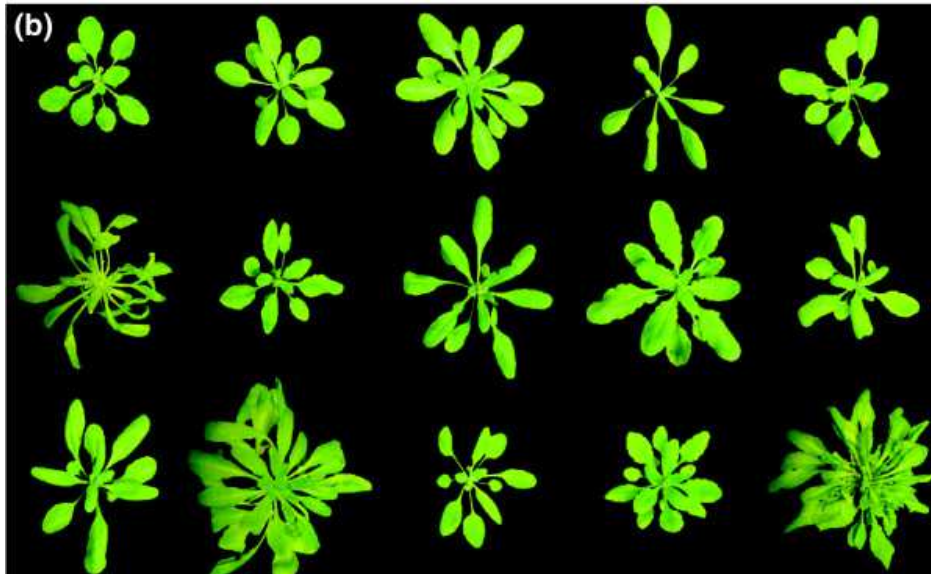
<http://1001genomes.org/accessions.html>

(a)



818 accessioni sequenziate e
rilasciate al 13-3-2014

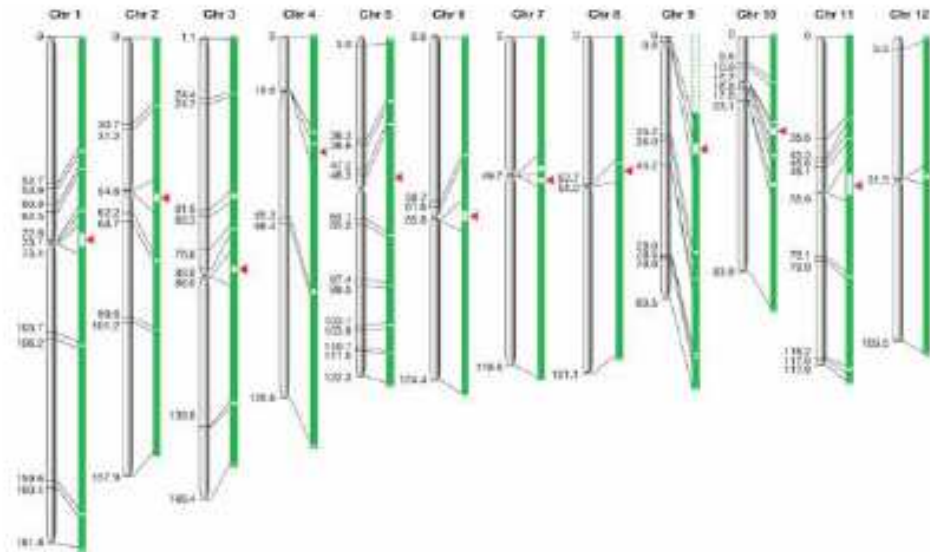
(b)



Weigel and Mott *Genome Biology*
2009, 10:107
doi:10.1186/gb-2009-10-5-107

RISO (*Oryza sativa*)

The rice genome



389 Mbp, 12 Chromosomes
Finished sequence (95% coverage)

37,544 / 22,840 [61%] genes (FGENESH models / cDNA-supported)

71% / 89% have homologs in Arabidopsis

90% of Arabidopsis genes have homologs in rice

2,859 genes not present in Arabidopsis

80,127 polymorphic loci *japonica* / *indica*

"With a large number of proteins of unknown function, the most interesting differences between the genome content of [Arabidopsis and rice] remain to be discovered."

International Rice Genome Sequencing Project, Nature (2005)

GENOMICA COMPARATIVA

- Analisi e confronto di genomi di specie diverse
- Fornisce informazioni sull'evoluzione delle specie e sulla funzione di geni e sequenze non codificanti
- Es.: funzione di un gene dedotta dallo studio di geni ortologhi in specie modello

GENOMICA COMPARATIVA

Cosa si analizza?

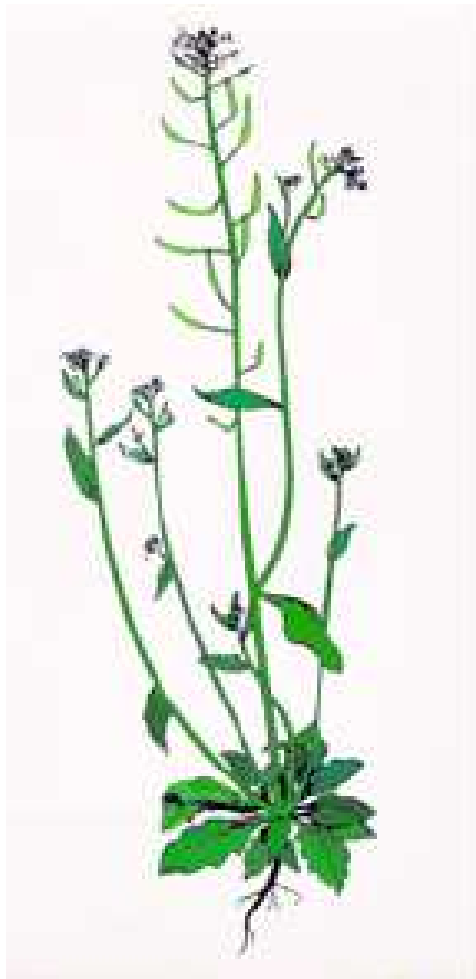
- Similarità di sequenza
- Localizzazione cromosomica dei geni
- Lunghezza e numero esoni
- Quantità di DNA non codificante
- Conservazione di regioni cromosomiche

Ostacoli al sequenziamento di specie coltivate

- **Dimensioni**
- **DNA ripetitivo**
- **Poliploidia**

Dimensioni del genoma

Arabidopsis: 125 Mb



***Fritillaria assyriaca*: 125 Gb!**



ILLUMINA GENOME ANALYZER

Permette il sequenziamento in parallelo di un numero massiccio di frammenti genomici

-> 1 milione di basi sequenziate per volta!



DNA ripetitivo

**Responsabile per gran parte della
variabilità nelle dimensioni del
genoma vegetale**

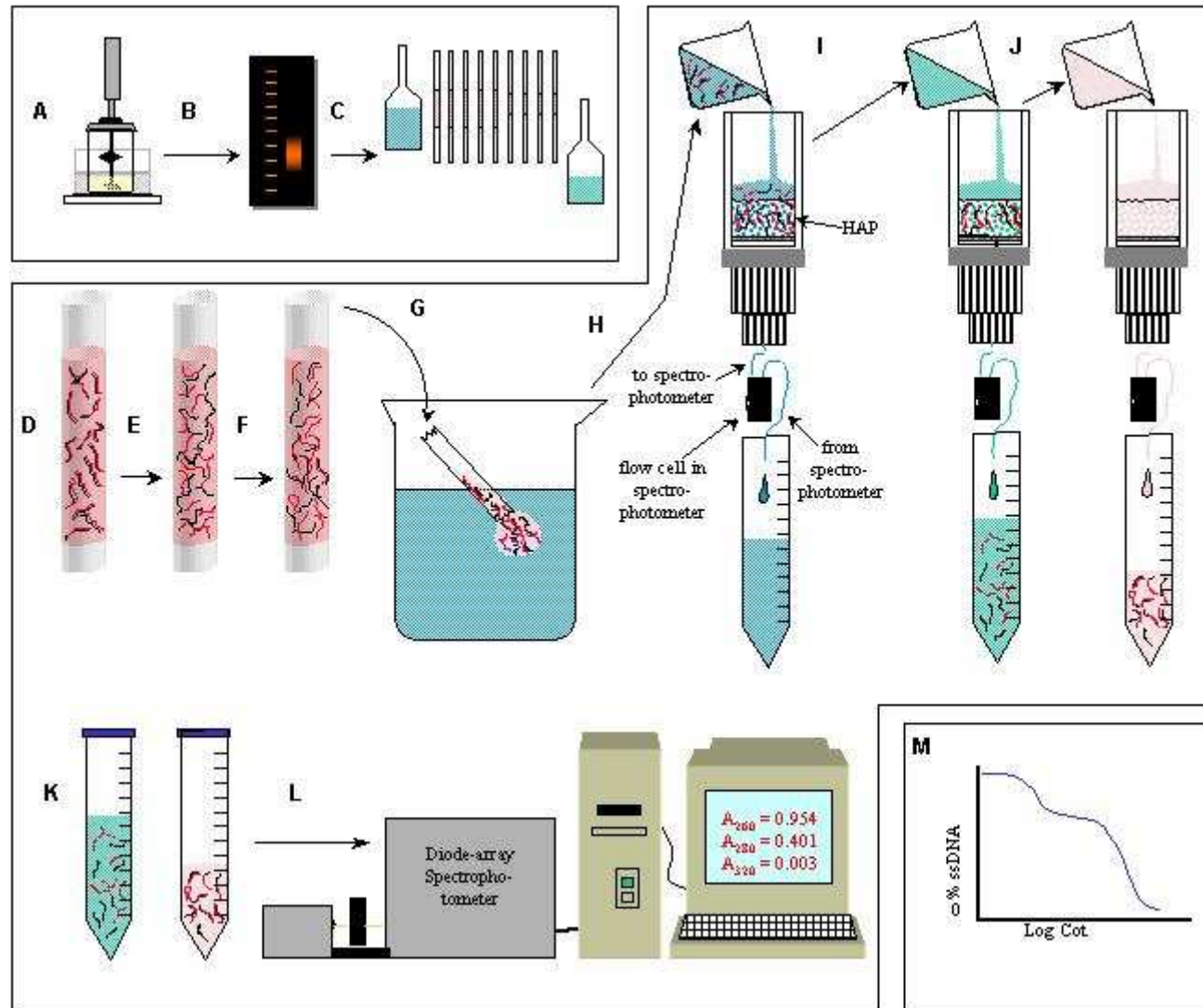
**Complica l'assemblamento delle
sequenze**

**Sequenze non-ridondanti nel genoma: da
13% (cipolla) a 77% (pomodoro)**

DNA ripetitivo

N.B.: le piante hanno più DNA ripetitivo degli animali, e copie individuali possono avere meno mutazioni per distinguerle, perchè più recenti

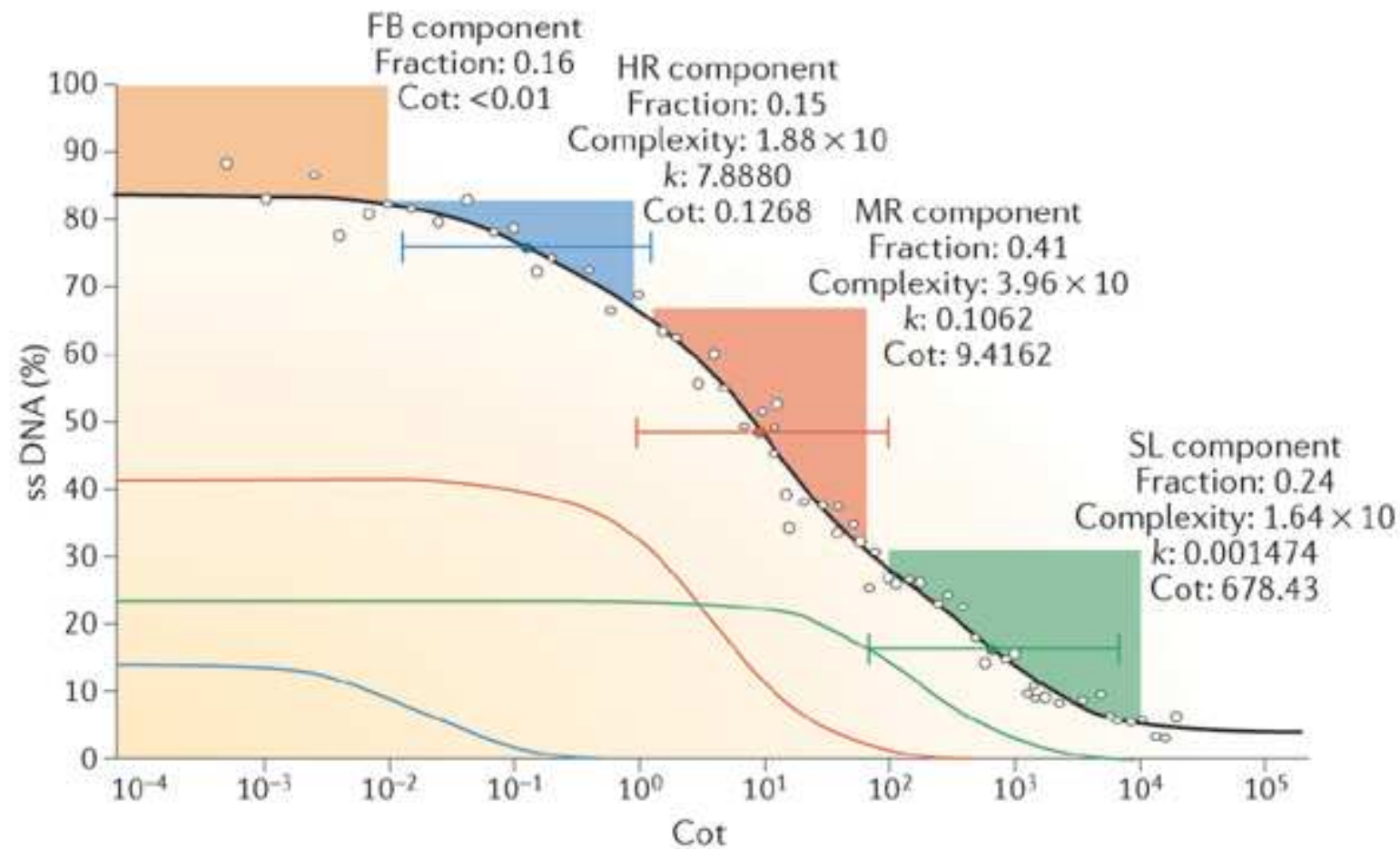
Cinetica di riassociazione



Cinetica di riassociazione

- Fornisce il valore Cot , cioè il prodotto fra la concentrazione dei nucleotidi (C_0) ed il tempo di riassociazione (normalizzato per la conc. di cationi nel tampone)
- La cromatografia su colonna di idrossiapatite (che lega il dsDNA) permette di isolare la frazione di DNA che si riassocia ad un particolare valore di Cot .
- Più il DNA è ripetitivo, più basso sarà il suo valore Cot

Cinetica di riassociazione



Clonaggio basato sul valore Cot (CBCS)

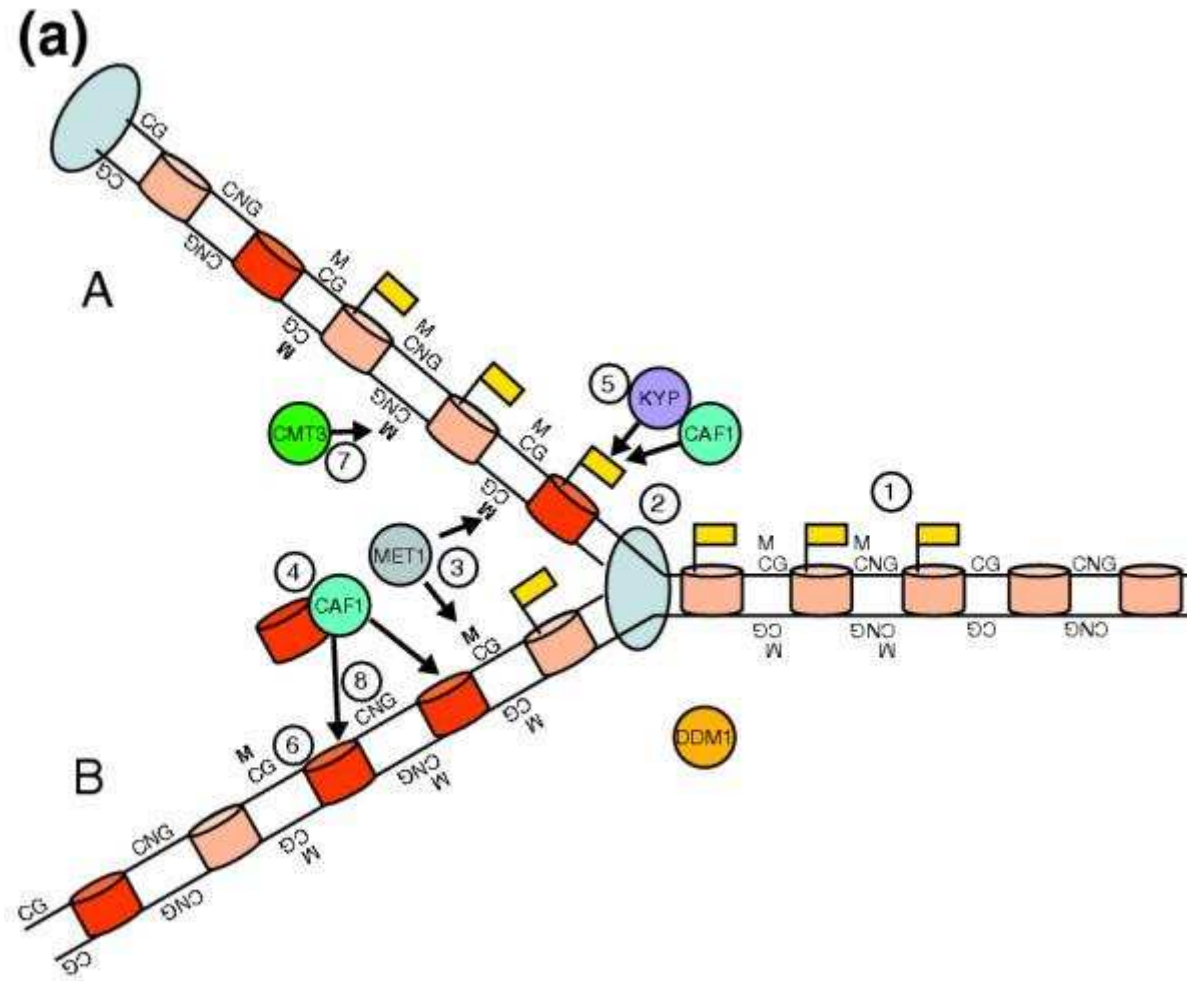
L'analisi Cot permette di isolare specificamente frazioni più o meno ripetitive

DNA meno ripetitivo viene sequenziato

-> più facile da assemblare in contigs

-> maggiore percentuale di geni

il DNA più ricco in geni è ipometilato rispetto a quello non codificante (inclusa una parte di DNA ripetitivo)



Methylation filtration (MF)

Clonaggio del DNA genomico totale in ceppi di *E. coli* che degradano il DNA metilato -> sequenziamento dei cloni e assemblaggio in contigs

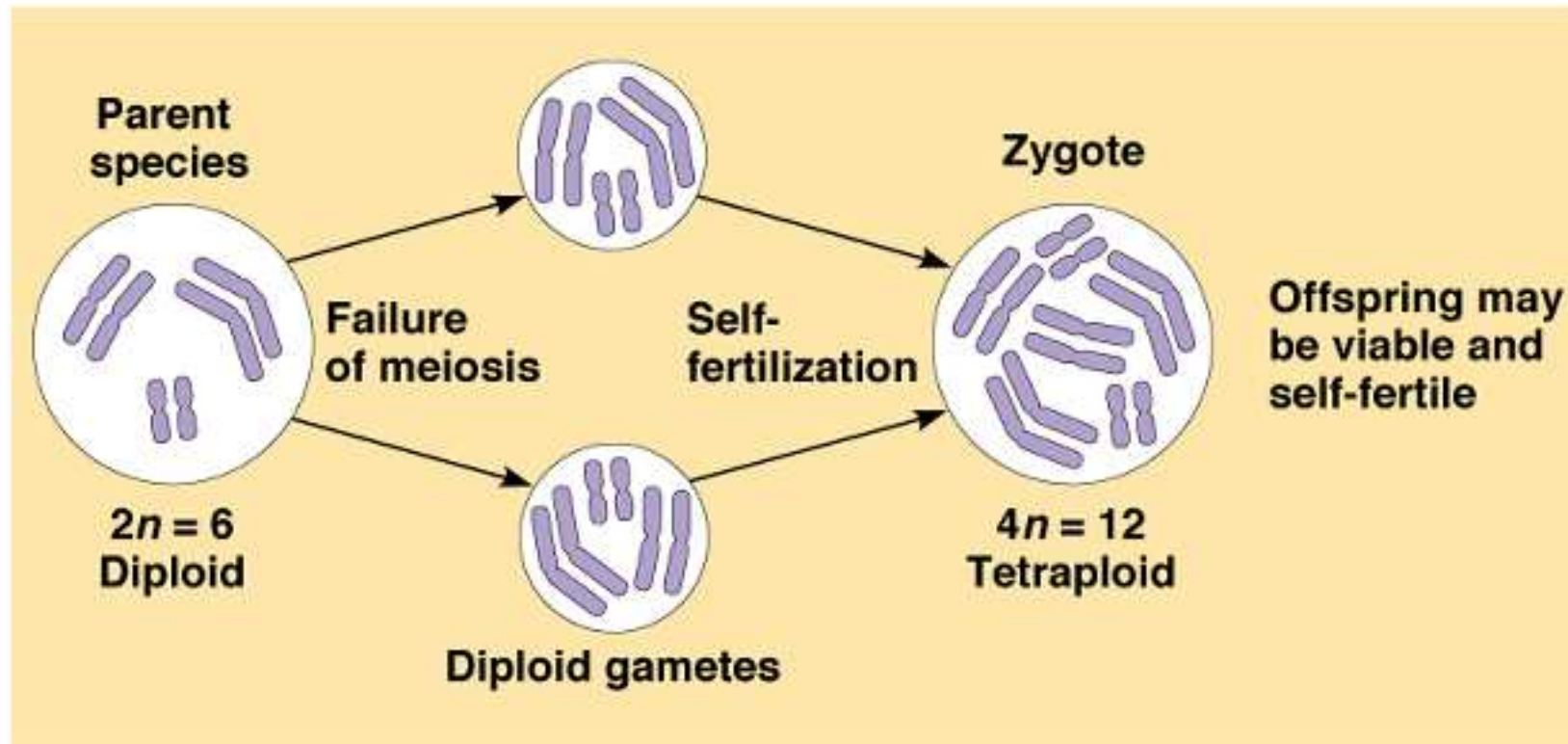
Svantaggio: non sempre il DNA codificante è ipometilato (es. metilazione indotta da stress, o in colture cellulari)

POLIPLOIDIA

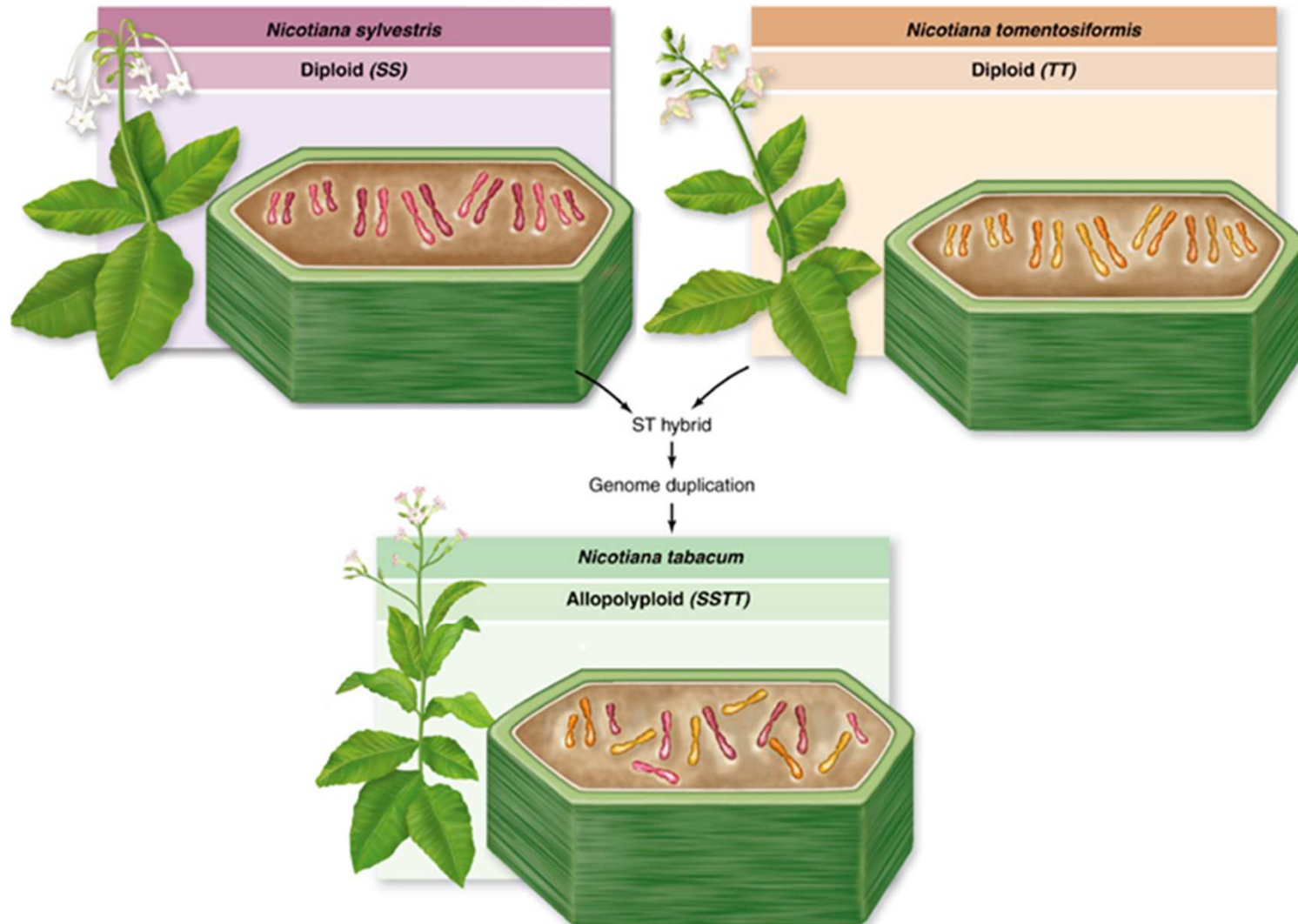
- Duplicazione del genoma in una specie (**autopoliploidia**), attraverso errore meiotico (4 copie di ogni cromosoma)
- Ibridazione di due specie diverse (**allopoliploidia**)

AUTOPOLIPLOIDIA

(es. Canna da zucchero, patata, erba medica, caffè)



ALLOPOLIPLIIDIA: TABACCO



ALLOPOLIPLIIDIA: FRUMENTO

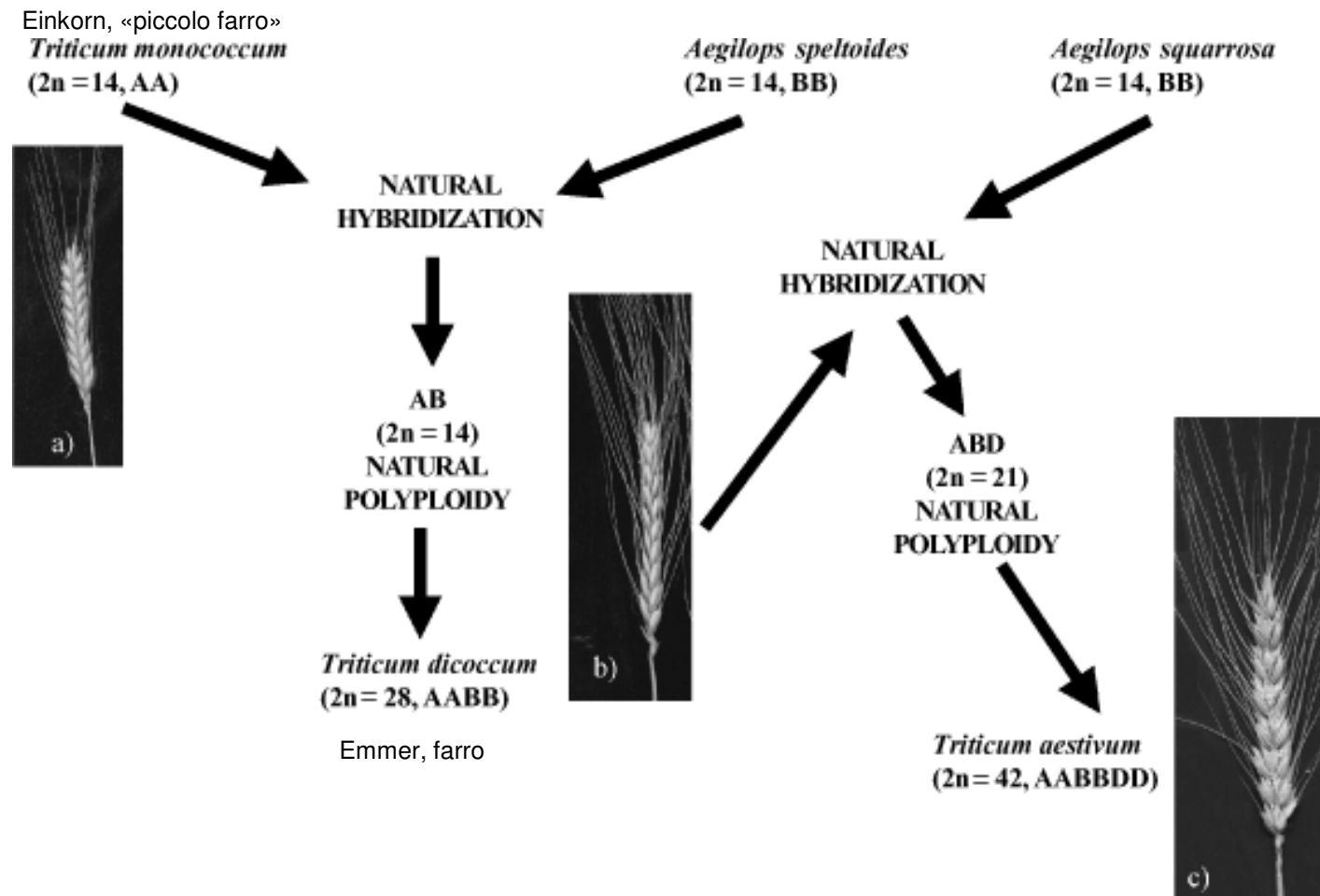
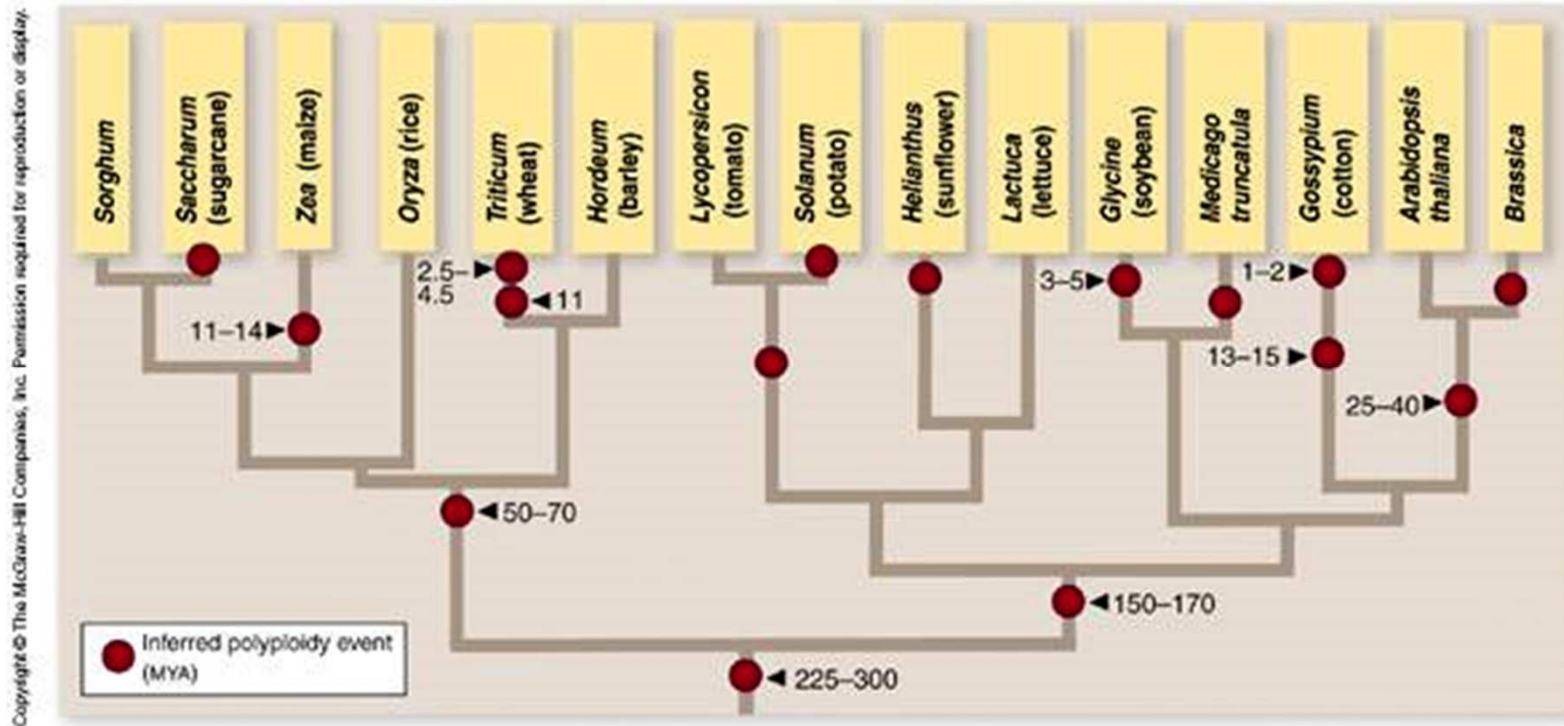


Figure 1 - Synoptic chart of cultivated wheats evolution: the diploid (2n = 14, AA) forms of *Triticum monococcum* (a) were naturally pollinated by weed species, possible *Aegilops speltoides* (2n = 14, BB?), in about 10,000 B.C. primitive farms. The subsequent genome duplication of hybrids by natural polyploidy gave rise to several wild and cultivated tetraploid species (2n = 28, AABB) like *Triticum dicoccum* (b) and *Triticum durum* (Figure 2a); again, the natural pollination of the tetraploid *T. dicoccum* (b) by another weed species, *Aegilops squarrosa* (2n = 14, DD) gave rise to the hexaploid (2n = 42, AABBDD) species (c).

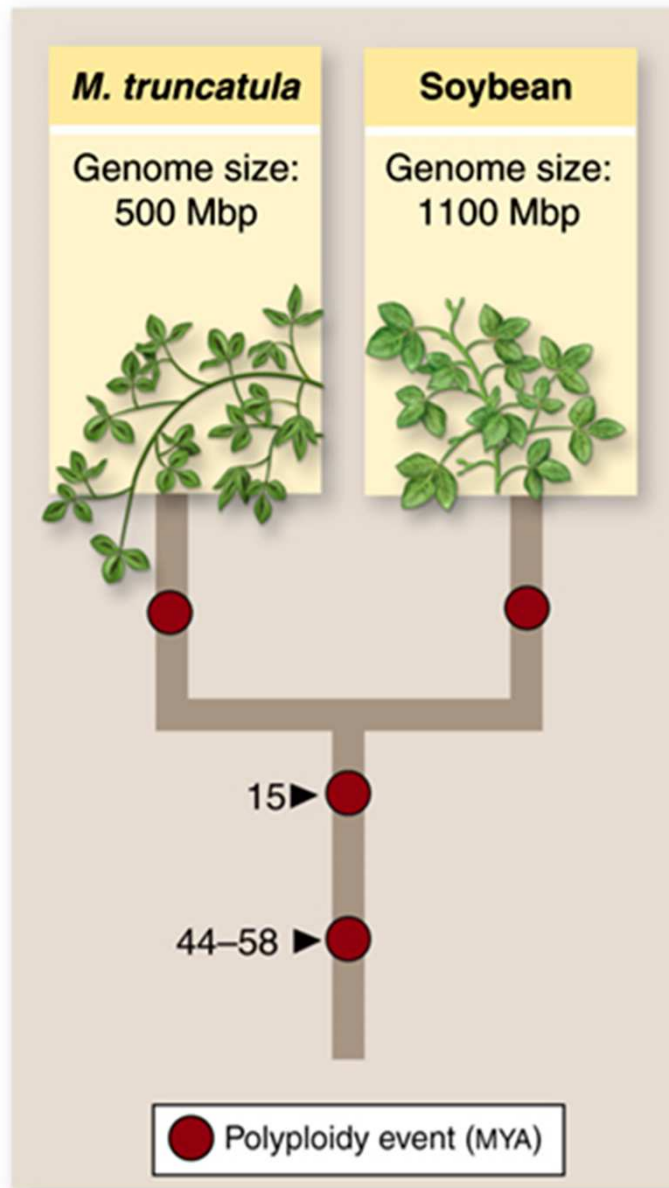
La poliploidia guida lo studio dell'evoluzione dei genomi

- **Paleopoliploidia:** confronto degli eventi di ploidizzazione
 - Divergenza di sequenze duplicate
 - Presenza o assenza di coppie di geni duplicati in seguito a ibridazione

Evoluzione dei genomi

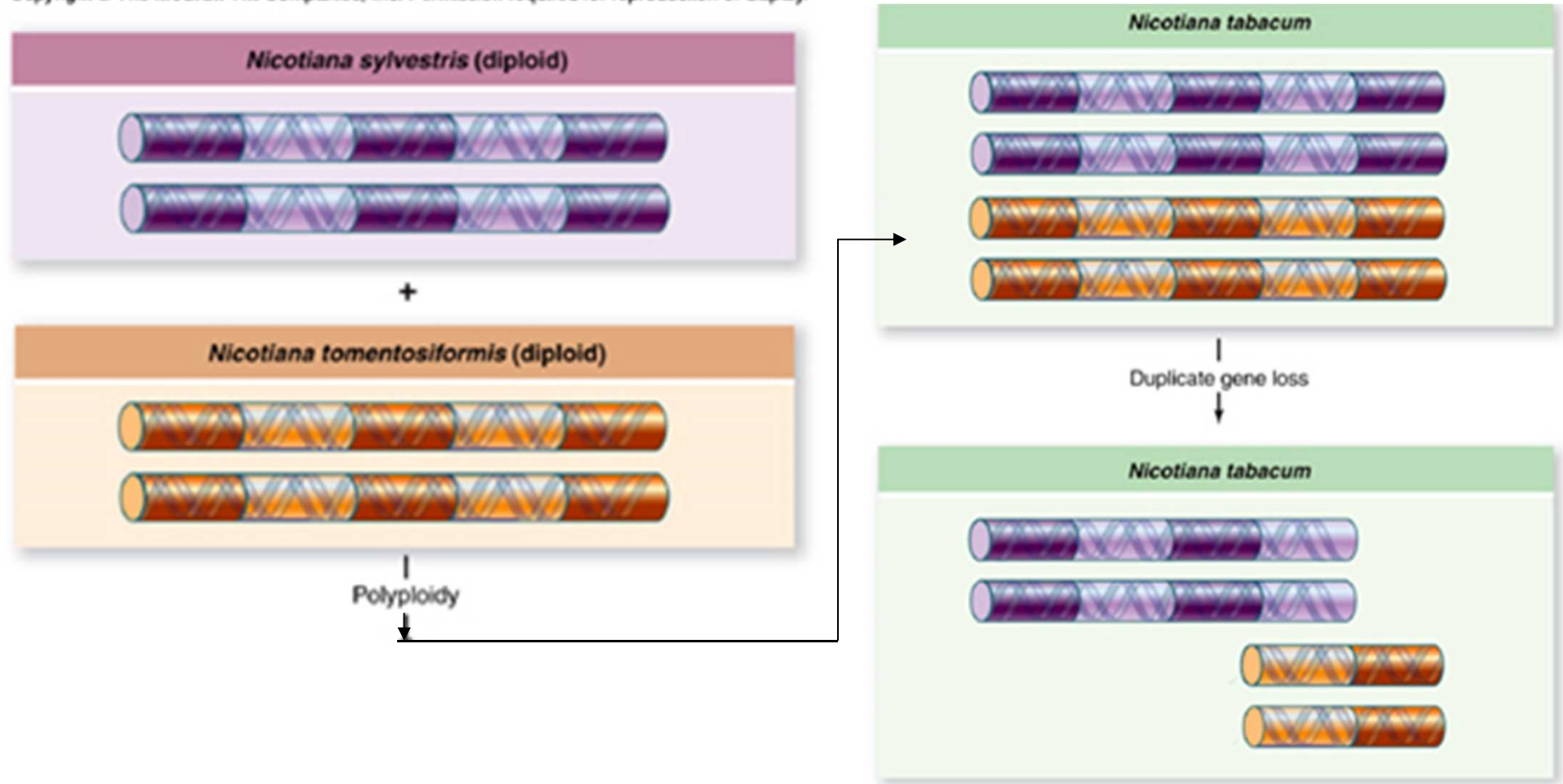


La poliploidia è diffusa nelle piante e ha avuto origini multiple durante l'evoluzione



Riduzione delle dimensioni del genoma

- Destino dei geni duplicati
 - Perdita di funzione per mutazione
 - Nuove funzioni
 - Suddivisione delle funzioni tra le due copie



Perdita di geni duplicati-> problema anche per identificare geni ortologhi in specie diverse

POLIPLOIDIA E SEQUENZIAMENTO DEI GENOMI

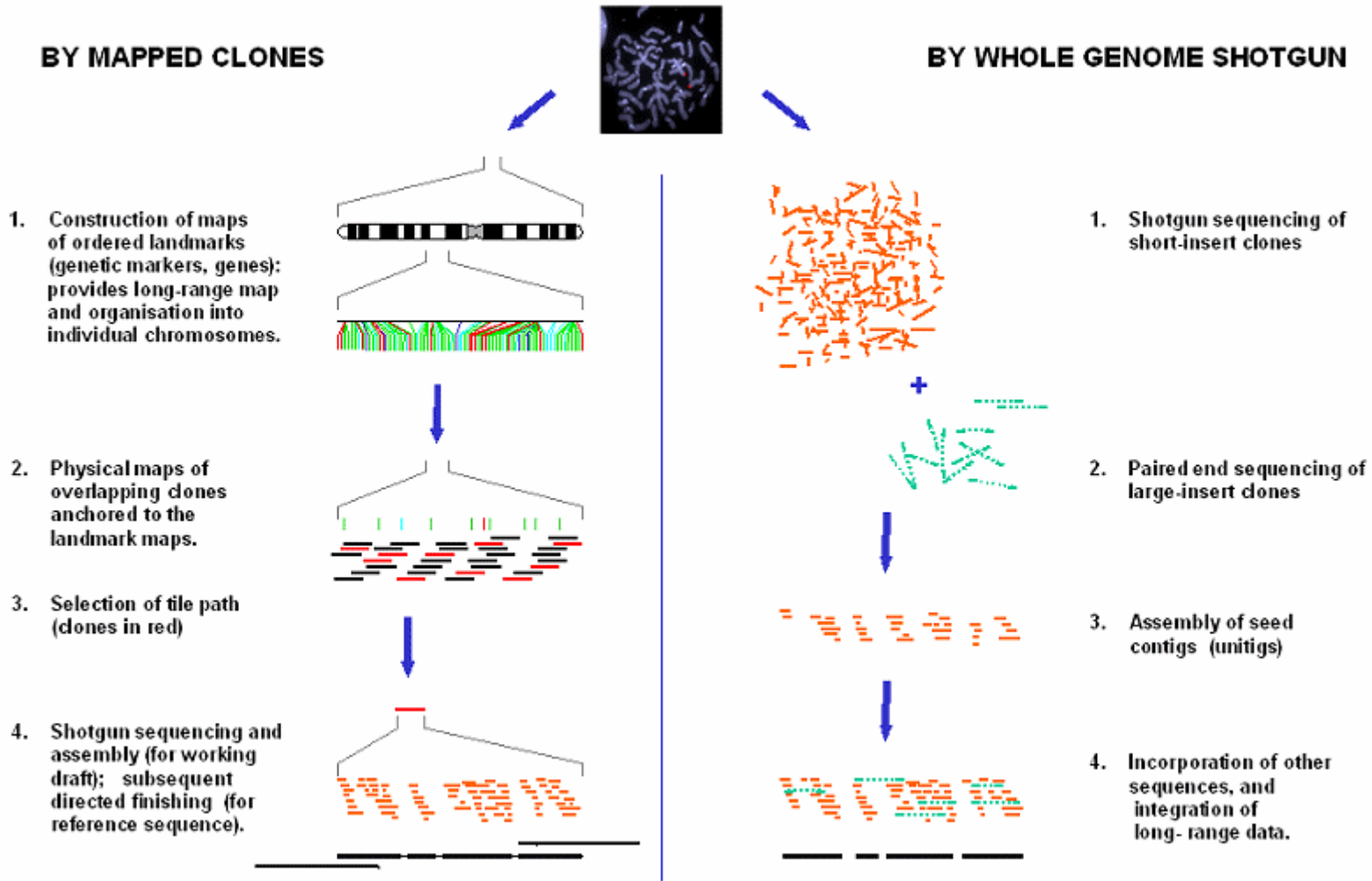
Molte specie autopoliploidi sono intolleranti all'INBREEDING, e hanno alti livelli di eterozigosità, importanti per la produttività

-> problema nell'assemblaggio dei contigs (più alleli diversi per ogni gene)

Negli alloploiploidi i cromosomi duplicati hanno subito sufficiente divergenza per non appaiarsi tra loro -> le sequenze delle coppie geniche sono distinguibili

N.B.: tutte le angiosperme sono PALEOPOLIPLOIDI, ma i geni "paleologhi" sono normalmente ben differenziati

Sequenziamento “whole-genome shotgun” o “clone-by-clone”?



“whole-genome shotgun”

Vantaggi

Rapido
Meno costoso
Utile per sequenziare regioni
refrattarie alla mappatura
fisica (es. regioni ripetitive)

Svantaggi

Assemblaggio complicato se ci
sono molte regioni ripetitive
In autoploididi, non distingue
aplotipi diversi di geni identici

“clone-by-clone”

Vantaggi

Delimita l'incertezza a intervalli
piccoli (100Kb)

Un allele alla volta -> no problema
di eterozigosità

Svantaggi

Costo dell'assemblaggio della
library e dell'ordinamento dei
contigs

EST = Expressed Sequence Tags

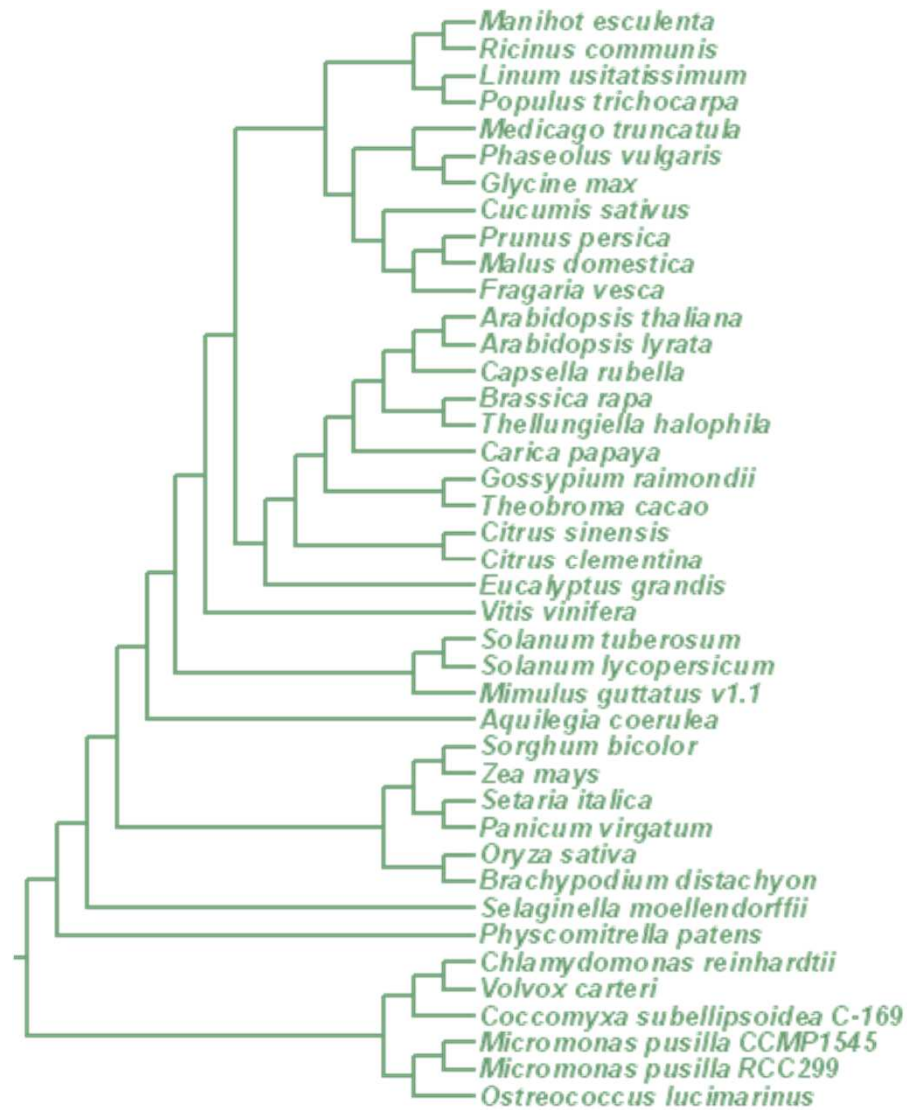
Creata sequenziando l'estremità 5' e/o 3' di mRNA isolati a caso e convertiti in cDNA (di solito 200–900 nt)

- > veloce e poco costoso
- > scoperta geni nuovi
- > marcatori per mappatura
- > base per futuri progetti di sequenziamento genomico
- > parziale copertura della porzione codificante del genoma

GENOMICA COMPARATIVA

- Analisi e confronto di genomi di specie diverse
- Fornisce informazioni sull'evoluzione delle specie e sulla funzione di geni e sequenze non codificanti
- Es.: funzione di un gene dedotta dallo studio di geni ortologhi in specie modello

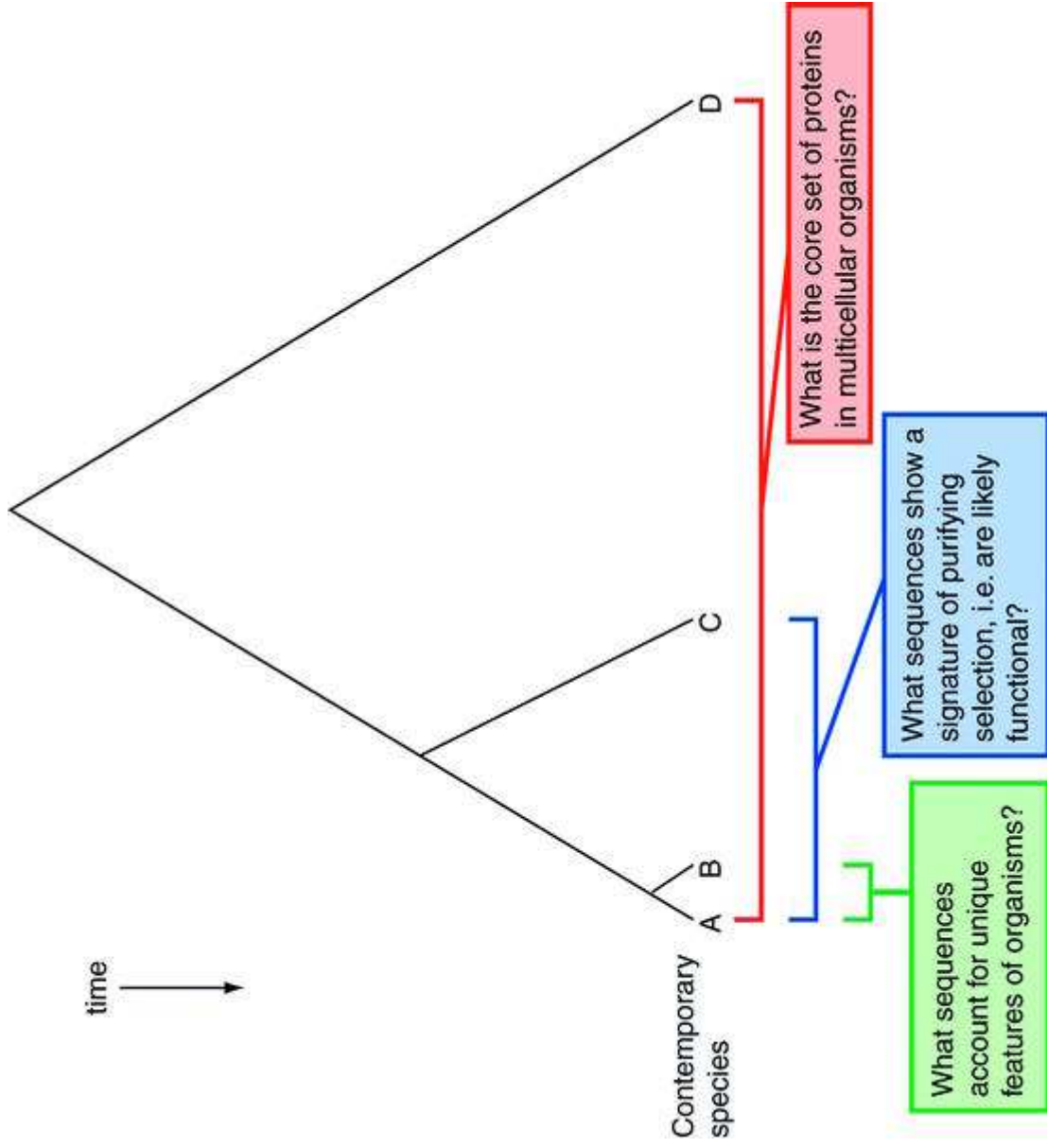
SEQUENCED AND ANNOTATED GREEN PLANT GENOMES



GENOMICA COMPARATIVA

Cosa si analizza?

- Similarità di sequenza
- Localizzazione cromosomica dei geni
- Lunghezza e numero esoni
- Quantità di DNA non codificante
- Conservazione di regioni cromosomiche



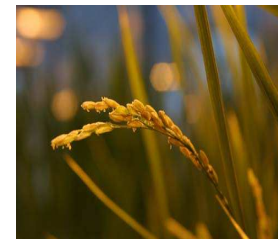
Predizione della funzione di un gene a partire dalla sequenza di geni in altre specie

Gene con funzione ignota



Trasferimento di annotazione

Specie modello

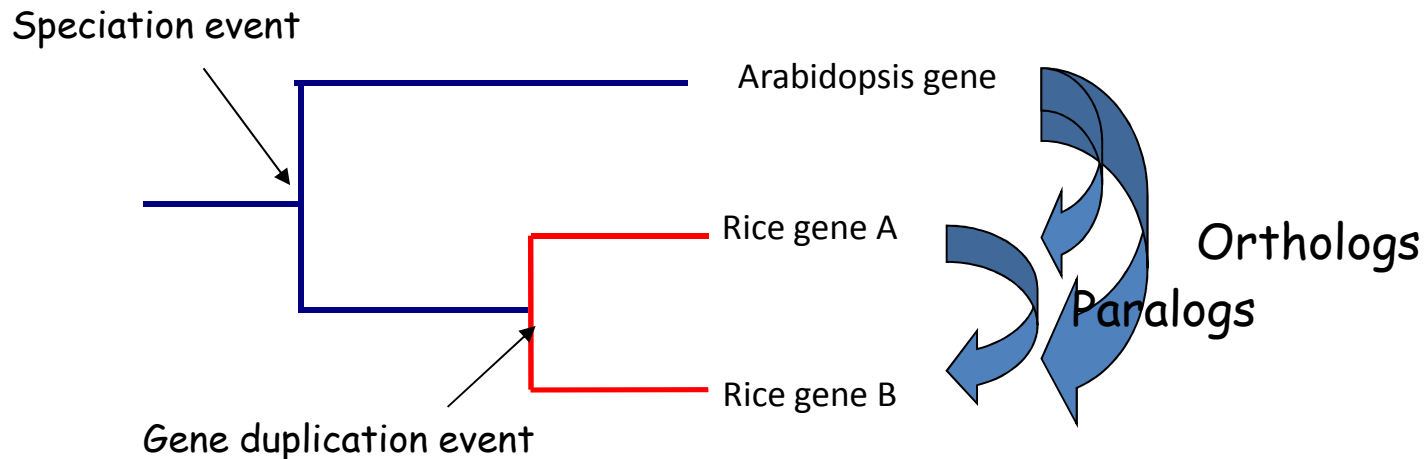


Geni omologhi

Gene con funzione X

Geni omologhi

- **Geni ortologhi** sono geni omologhi che discendono dall'ultimo ancestore comune attraverso speciazione
- Molto probabilmente codificano per proteine con funzione simile



- **Geni paraloghi** sono geni omologhi che si sono evoluti per duplicazione e possono codificare proteine con funzioni più divergenti
- **Geni inparaloghi**: geni ortologhi che hanno subito duplicazione

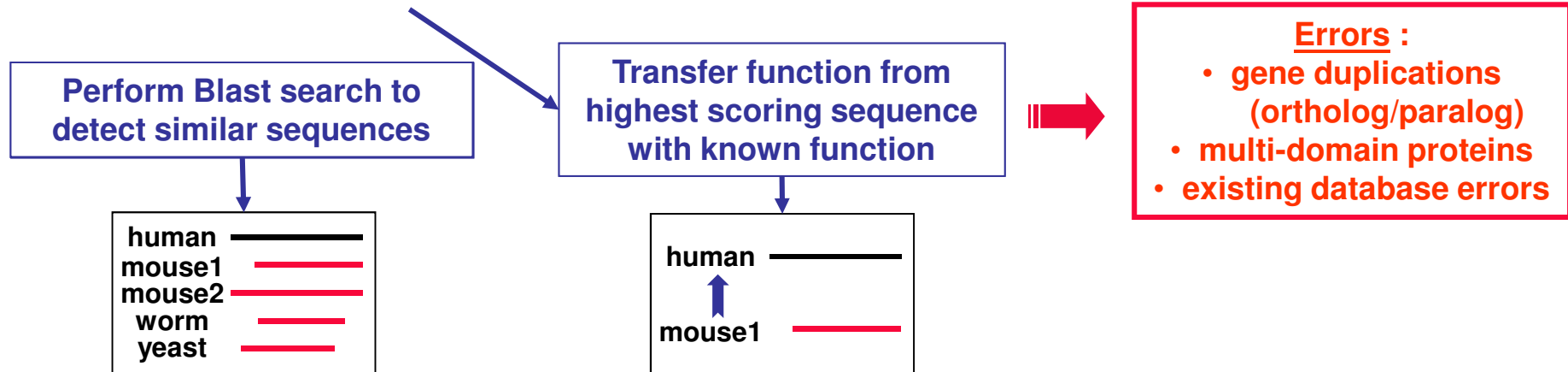
Come predire l'omologia?

Similarità e omologia non sono la stessa cosa!

Geni simili si assomigliano sulla base di un'osservazione empirica

Geni omologhi sono geneticamente correlati (fatto storico: hanno antenato comune)

Metodo classico : annotazione funzionale basata sulla somiglianza (Blast)



Predizione dell'omologia sulla base della similarità

Es. BLAST

Vantaggi:

- Facile
- Veloce
- Direttamente sul genoma completo

Svantaggi:

- Come stabilire la soglia di E-value per trasferire l'annotazione del gene da una specie all'altra?

Due sequenze possono presentare similarità senza essere evolutivamente correlate!

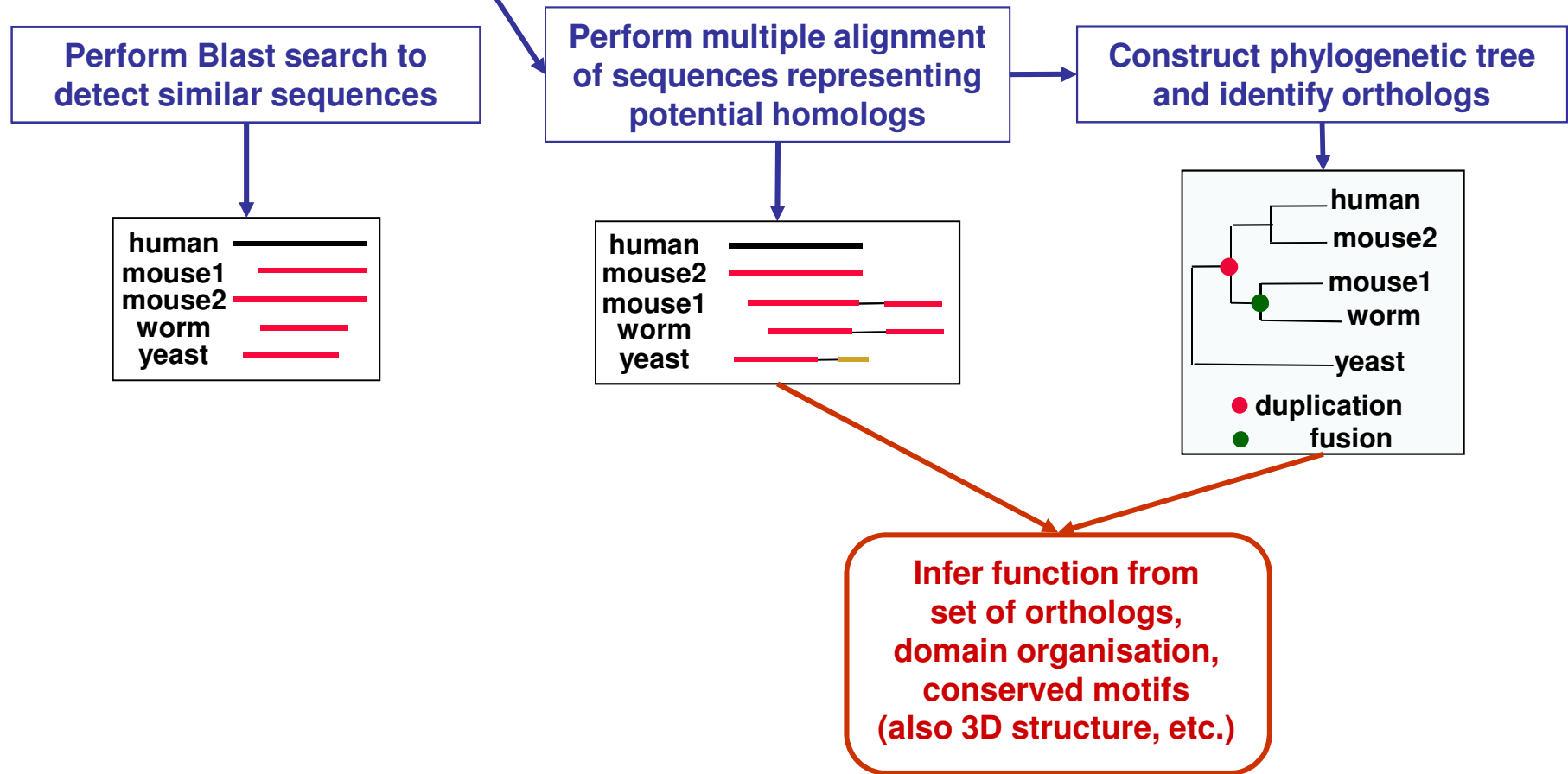
- Non identifica eventi di duplicazione genica

**Come trovare in una specie un gene ortologo ad
un gene noto in un'altra specie?**



FILOGENOMICA

Phylogeny-based inference



Predizione dell'omologia sulla base della filogenesi

Vantaggi:

- Efficiente per identificare duplicazioni (paraloghi e ortologhi)

Svantaggi:

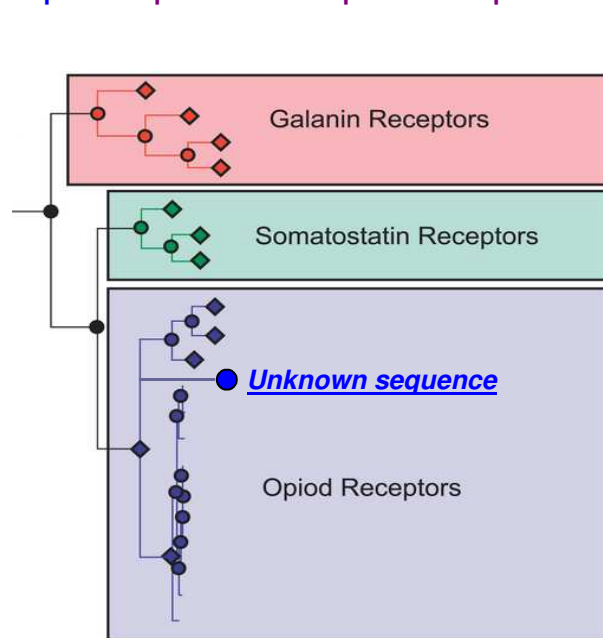
- Lento
- Richiede raggruppamento dei geni in famiglie

Metodi correnti

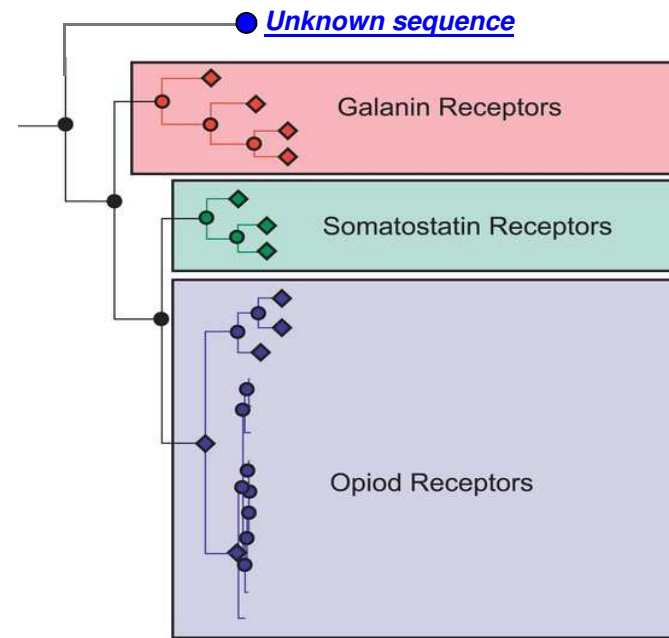
- RIO e Orthostrapper : solo per 1900 famiglie di geni vegetali (Pfam)
- GOST (usa GreenPhylDB family : 6420 famiglie geniche vegetali)

- Tree-based orthology: build a phylogenetic tree of a group of genes and compare gene tree to species tree to define speciation, duplication events
 - Resampled Inference of Orthologs (RIO) (Zmasek and Eddy, 2002)
 - Orthostrapper (Storm and Sonnhammer, 2002)
 - Levels Of Orthology From Trees (LOFT) (Van de Heijden et al, 2007)

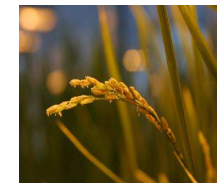
■ Example: G protein-coupled receptors



Prediction: Opioid receptor



More general prediction: GPCR of unknown specificity



Due specie modello

Inizialmente *Oryza sativa* e *Arabidopsis thaliana*:

- Genoma completo
- Alta qualità dell'annotazione (TAIR release 7, TIGR release 5)
- Evidenze funzionali disponibili

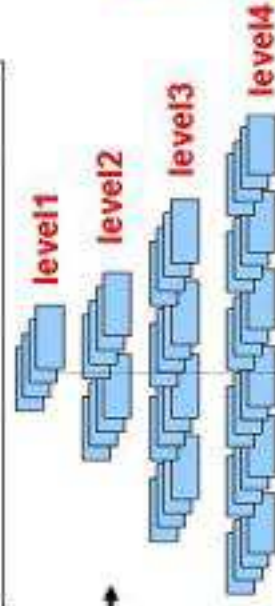
Nel tempo integra altre specie

<http://www.greenphyl.org/cgi-bin/index.cgi>

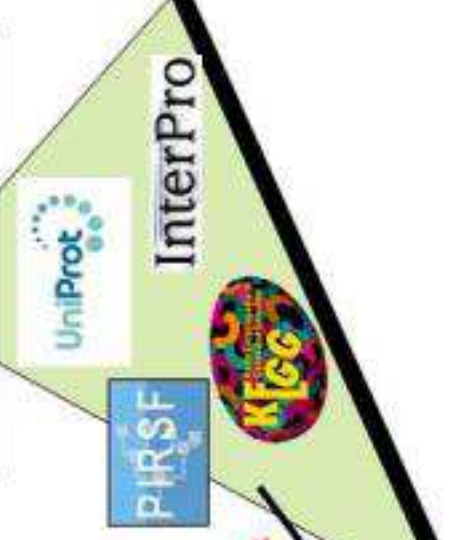
Full genome

```
>Seq1  
MESNYSGVVMQYDVSEFLPTSIPDLOFGVPSSDF  
>Seq2  
DURMDQYYRQPSIMVFDGDHHSPPA  
>Seq3  
DEIDSEMITLLKYYVNGLLMEESELAEKGGIFYDSLALR  
>Seq4  
QTEMLDQVYSDSQTQSSIPNNSI  
>Seq5  
...
```

Gene family clustering

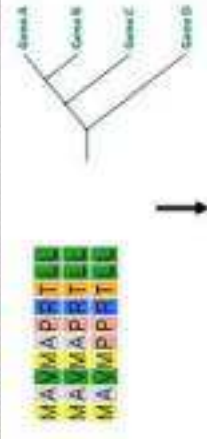


Cross references



Cluster annotation

Phylogenetic analyses of validated families



Phylogenomic inference: orthologs and paralogs identification

<u><i>Amborella trichopoda</i></u>	<u><i>Gossypium raimondii</i></u>	<u><i>Picea abies</i></u>
<u><i>Arabidopsis thaliana</i></u>	<u><i>Hordeum vulgare</i></u>	<u><i>Populus trichocarpa</i></u>
<u><i>Brachypodium distachyon</i></u>	<u><i>Lotus japonicus</i></u>	<u><i>Ricinus communis</i></u>
<u><i>Cajanus cajan</i></u>	<u><i>Malus domestica</i></u>	<u><i>Selaginella moellendorffii</i></u>
<u><i>Carica papaya</i></u>	<u><i>Manihot esculenta</i></u>	<u><i>Setaria italica</i></u>
<u><i>Chlamydomonas reinhardtii</i></u>	<u><i>Medicago truncatula</i></u>	<u><i>Solanum lycopersicum</i></u>
<u><i>Cicer arietinum</i></u>	<u><i>Musa acuminata</i></u>	<u><i>Solanum tuberosum</i></u>
<u><i>Citrus sinensis</i></u>	<u><i>Musa balbisiana</i></u>	<u><i>Sorghum bicolor</i></u>
<u><i>Coffea canephora</i></u>	<u><i>Oryza sativa</i></u>	<u><i>Theobroma cacao</i></u>
<u><i>Cucumis sativus</i></u>	<u><i>Ostreococcus tauri</i></u>	<u><i>Vitis vinifera</i></u>
<u><i>Cyanidioschyzon merolae</i></u>	<u><i>Phaseolus vulgaris</i></u>	<u><i>Zea mays</i></u>
<u><i>Elaeis guineensis</i></u>	<u><i>Phoenix dactylifera</i></u>	
<u><i>Glycine max</i></u>	<u><i>Physcomitrella patens</i></u>	



Family_name (e.g. GRAS)

Family ID: 20923
 Family name: Pollen Allergen/Expansin Superfamily

Synonym(s)

Cross-reference(s)

Curation status



Not available

Phylogenetic analyses

(a)

Family structure Family composition Protein domains Protein list Phylogenomic analysis Chromosome position

Clustering level Family id (number of sequences)

1	20923 (812)
2	24999 (689) 25676 (26) 26339 (66)
3	30583 (515) 31138 (117) 31147 (26) 31717 (57) 31804 (66)
4	36255 (410) 36521 (194) 36760 (28) 37177 (65) 37759 (61)

(b)

Family structure Family composition Protein domains Protein list Phylogenomic analysis Chromosome position

InterPro family

IPR	Annotation	Type	%	Occurrence	Specificity
IPR007118	Expansin/Loi pl	Family	80	(653)	Y

Other InterPro signatures

IPR	Annotation	Type	%	Occurrence	Specificity
IPR014734	Pollen allergen, N-terminal	Domain	57	(466)	N
IPR009009	Barwin-related endoglucanase	Domain	58	(472)	N
IPR005132	Rare lipoprotein A	Domain	89	(722)	N
IPR007117	Pollen allergen/expansin, C-terminal	Domain	82	(667)	N
IPR007112	Expansin 45, endoglucanase-like	Domain	88	(714)	N

Domain architecture

Identified using MEME suite

[see sequence motifs](#)

Representative InterPro domains - Consensus schema (alpha version)

Show phylogenetic Tree

Tools View as Text Font Size Options Type Help

- Phylocon
- Dyna Hide
- Rollover
- Show Internal Data
- Taxonomic Colorize
- Annotation Colorize
- Collapse Branches
- Use Branch Width

Display Data:

- Node Names
- Taxonomic Code
- Taxonomic Matrix
- Prot/Genes Symbols
- Prot/Genes Name
- Prot/Genes Acc.
- Annotation
- Binary Characters
- Binary Char Counts
- Domains
- Confidence Value
- Event

Sequence relations to display
(type) [orthology]

- Relation confidence filter

Click on Node to:

[collapse] [expand]

Zoom:

Y+ [Y-]

X+ [X-]

Y+ [Y-]

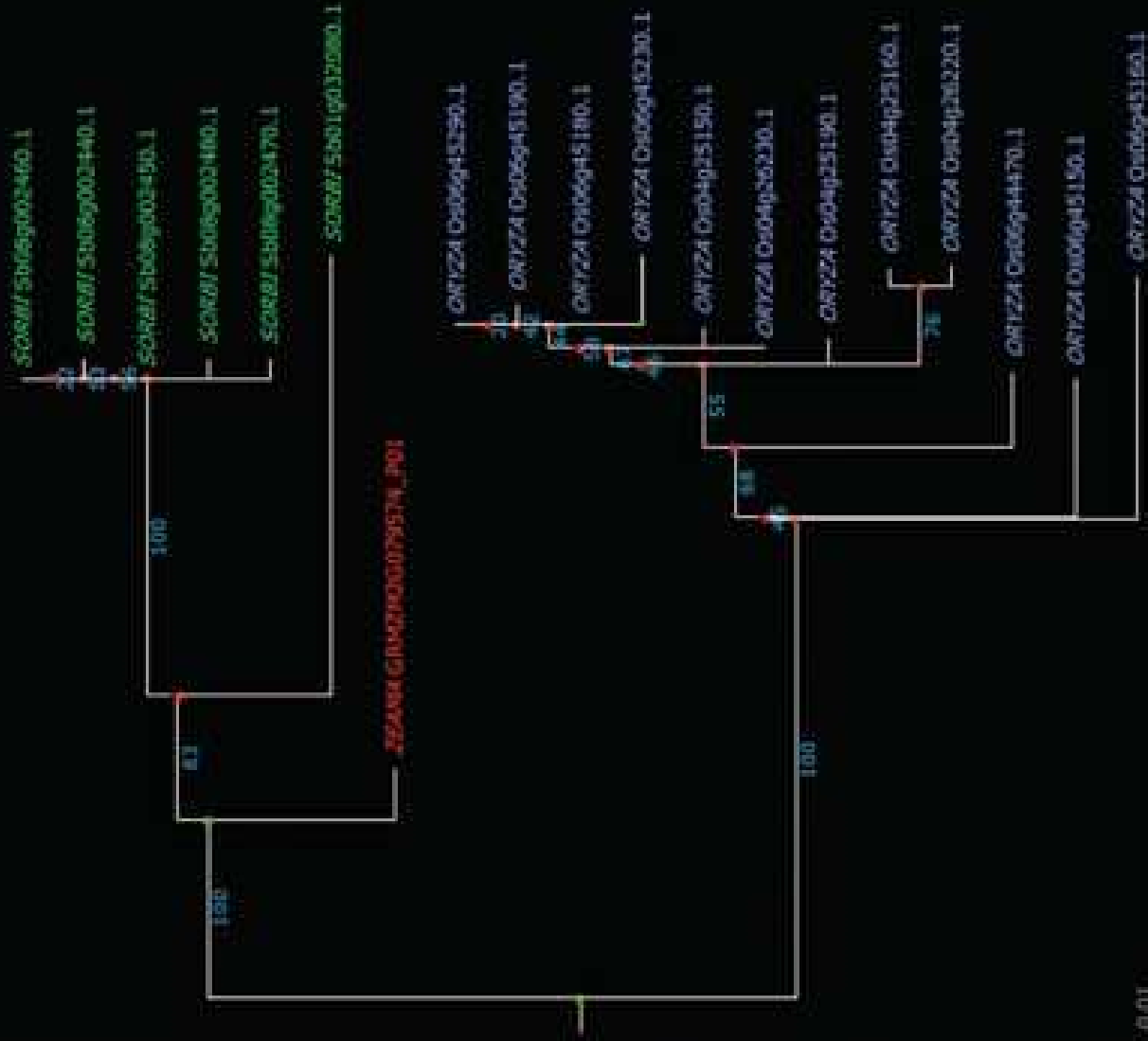
Back to Rooted Tree

Order Subtrees

Uncollapse All

Search:

[Search Input Field]



(e)

InParanoid

- <http://inparanoid.sbc.su.se>
- Database per identificare geni ortologhi e inparaloghi tra specie diverse di eucarioti (animali, piante, funghi, protisti)