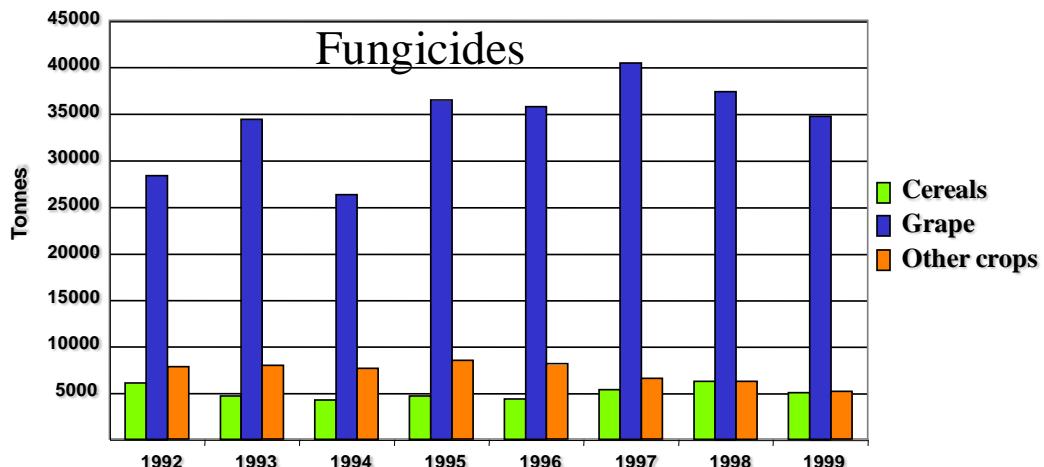
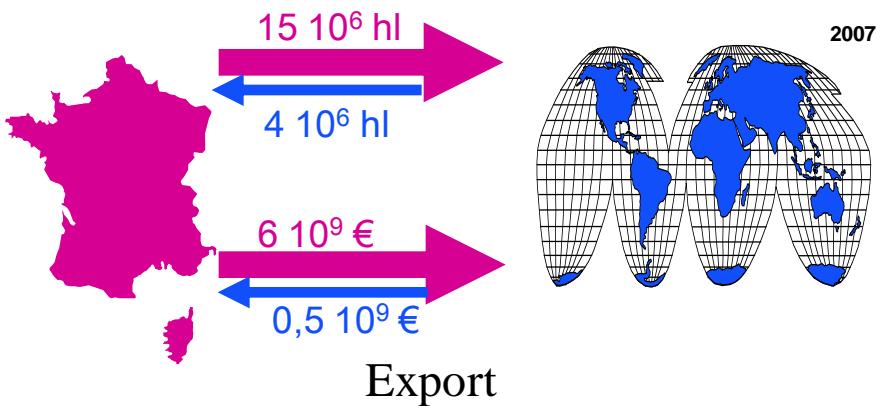
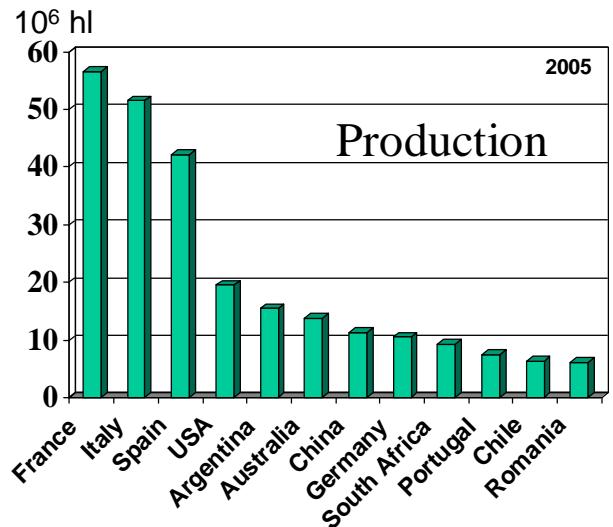


Structural and Functional Genomics in Grapevine through FLAGdb⁺⁺

Sandra Dèrozier, Cécile Guichard, Franck Samson, Jean-Philippe Tamby, Véronique Brunaud,
Vincent Thareau, Christophe Caron, Maria Tchoumakov, Roberto Bacilieri, Anne-Françoise Adam-Blondon
and Sébastien Aubourg

Why genomics on *Vitis vinifera* ?



$2n = 38$
 475 Mb

Genome sequencing



Franco-Italian consortium
Coord : A.-F. Adam-Blondon



France

Génoscope
INRA

7 X

Italy

Padova University 2.5 X
Udine University 2.5 X
Milano University

Sanger after WGS

8 X assembly (Jailon *et al.*, Nature 2007)
12 X assembly (2009)

Pinot Noir → PN40024 (93% homozygote)

Genome annotation strategy



ATAAACATGATCATGTTAACGTIGTA
ATCAAAACTGATTAAAAAGAACTAAT
GTCGT 'CTCA
GAAGT AGGG
TAAAC scaffolds GTGAA
AATGG TACGC
AGGATCCAAGCAAGCTTATTCAAGCTT
GACCTGATCACCCGTATTCATGAC
AACTGGAGGAAGCGGATACAGACGT

Masking of repeat sequences

Ab initio predictions
coding potential, signals

Sequence comparisons
similarities

Spliced alignments of
transcripts



GeneID
GlimmerHMM

EXOFISH

GeneWise
EST2Genome



SpliceMachine
EuGène IMM

BLASTX

GenomeThreader

training set

EMBL/GenBank/DDBJ

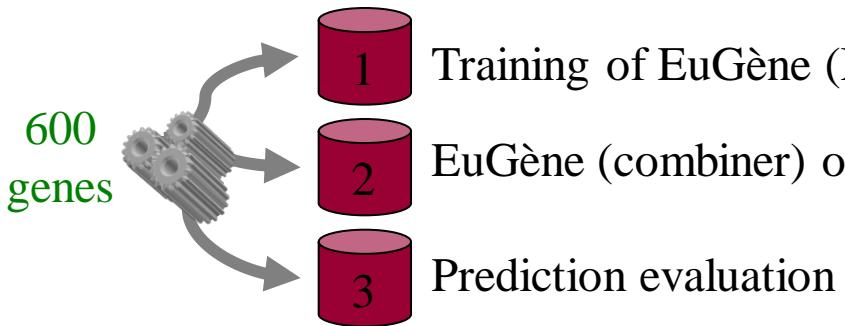
cDNA/EST resources

Integration and reconciliation

GAZE
EuGène

predicted gene models

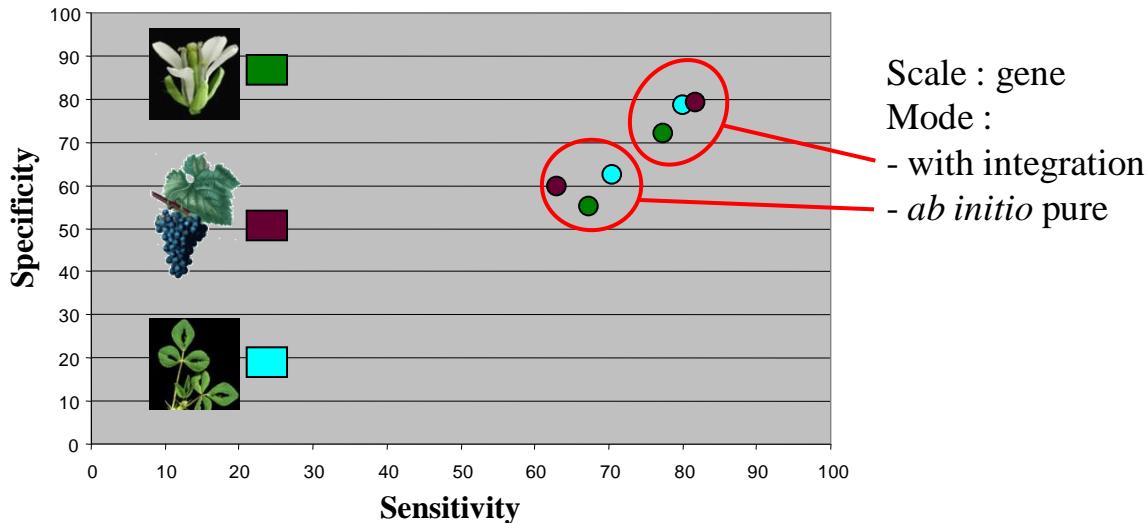
Training, evaluation and results



EuGène requires
high computation facilities:

migbole

(MIG)

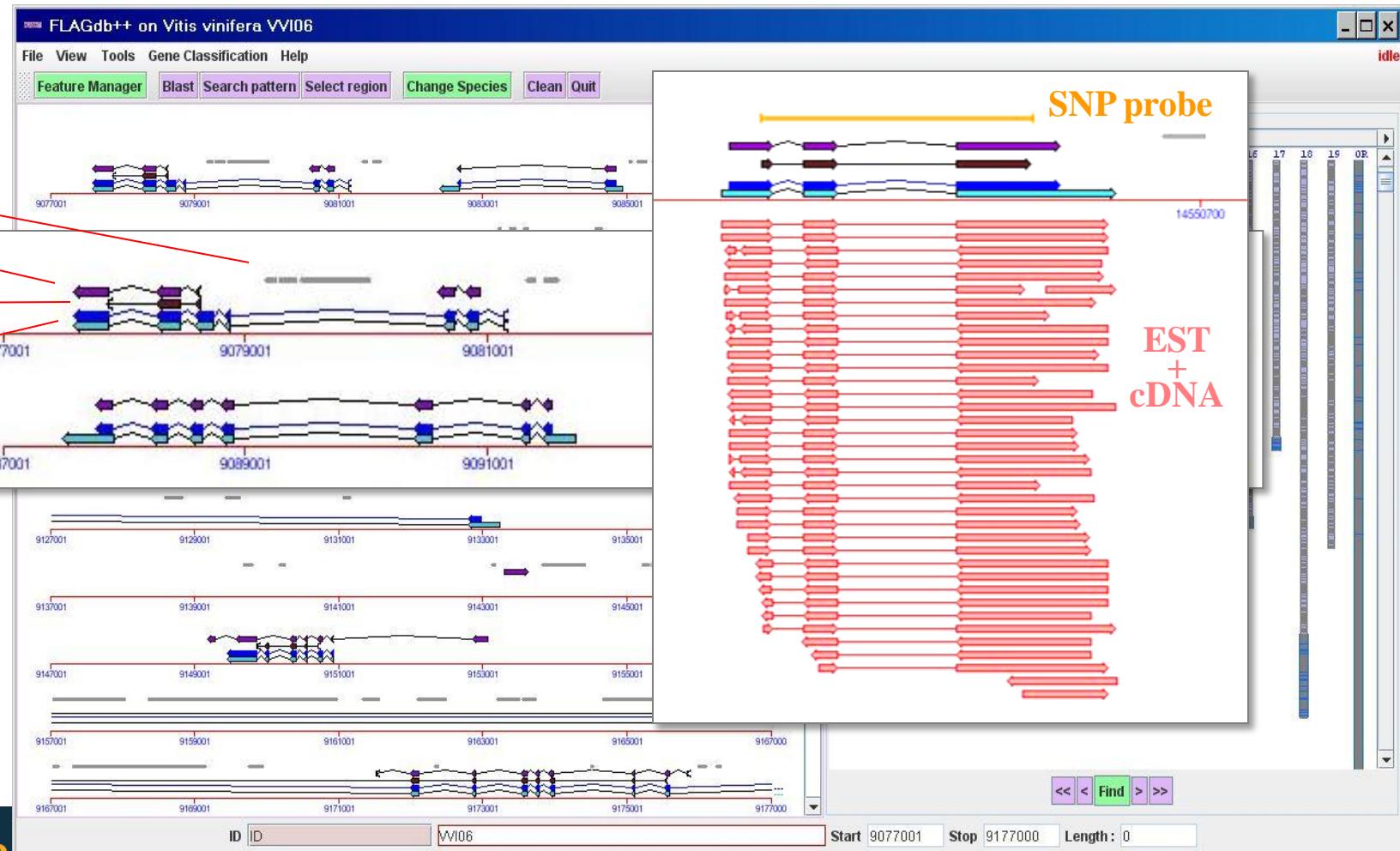


► EuGène predicts **44 414** coding genes in the 12x PN40024 genome
 (GAZE : 26 347 genes)

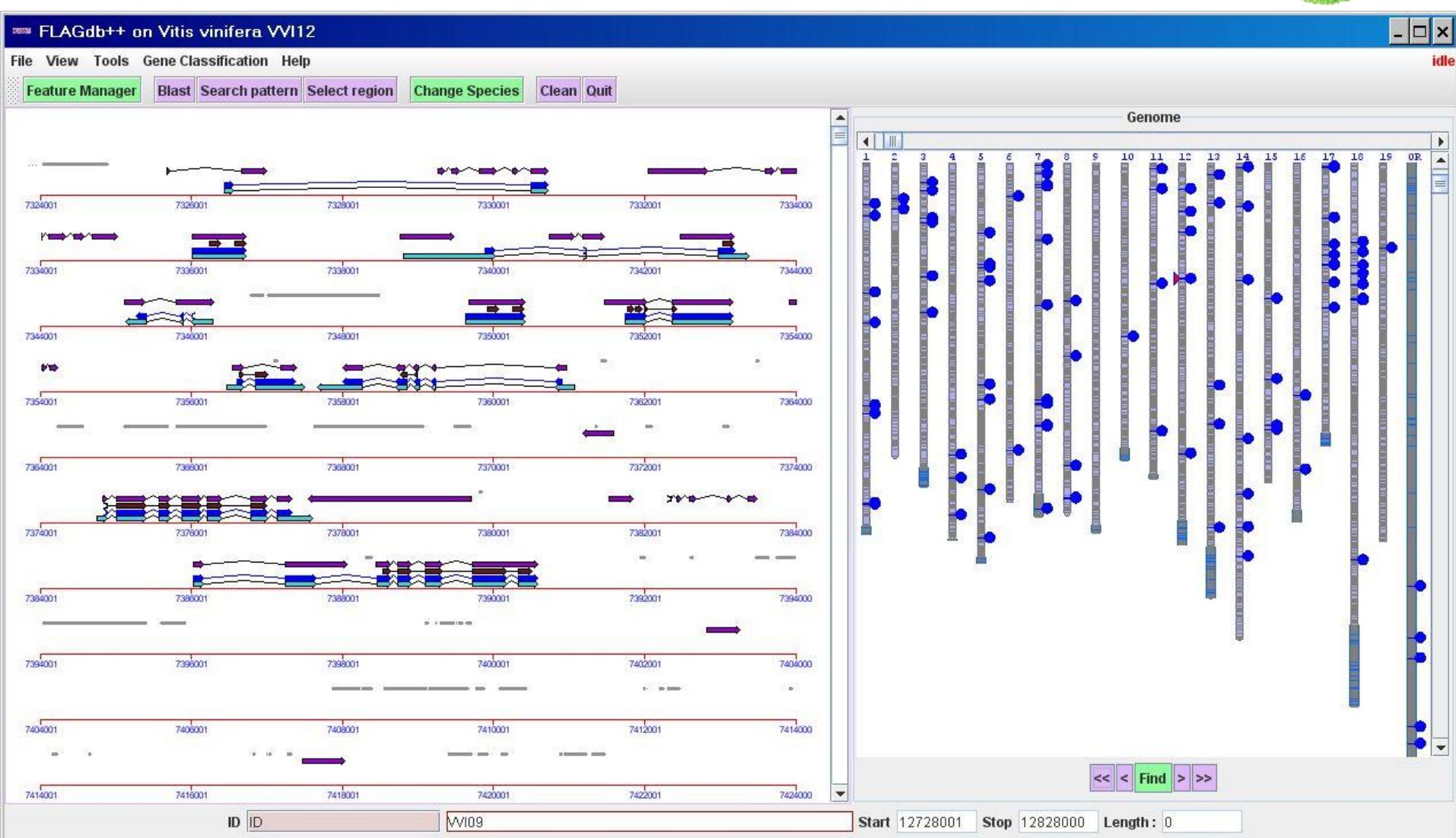
Vitis vinifera in FLAGdb⁺⁺



Relational database and JAVA application dedicated to integration and exploration of genomic data (experimental or predicted) in order to help the functional study of plant genes

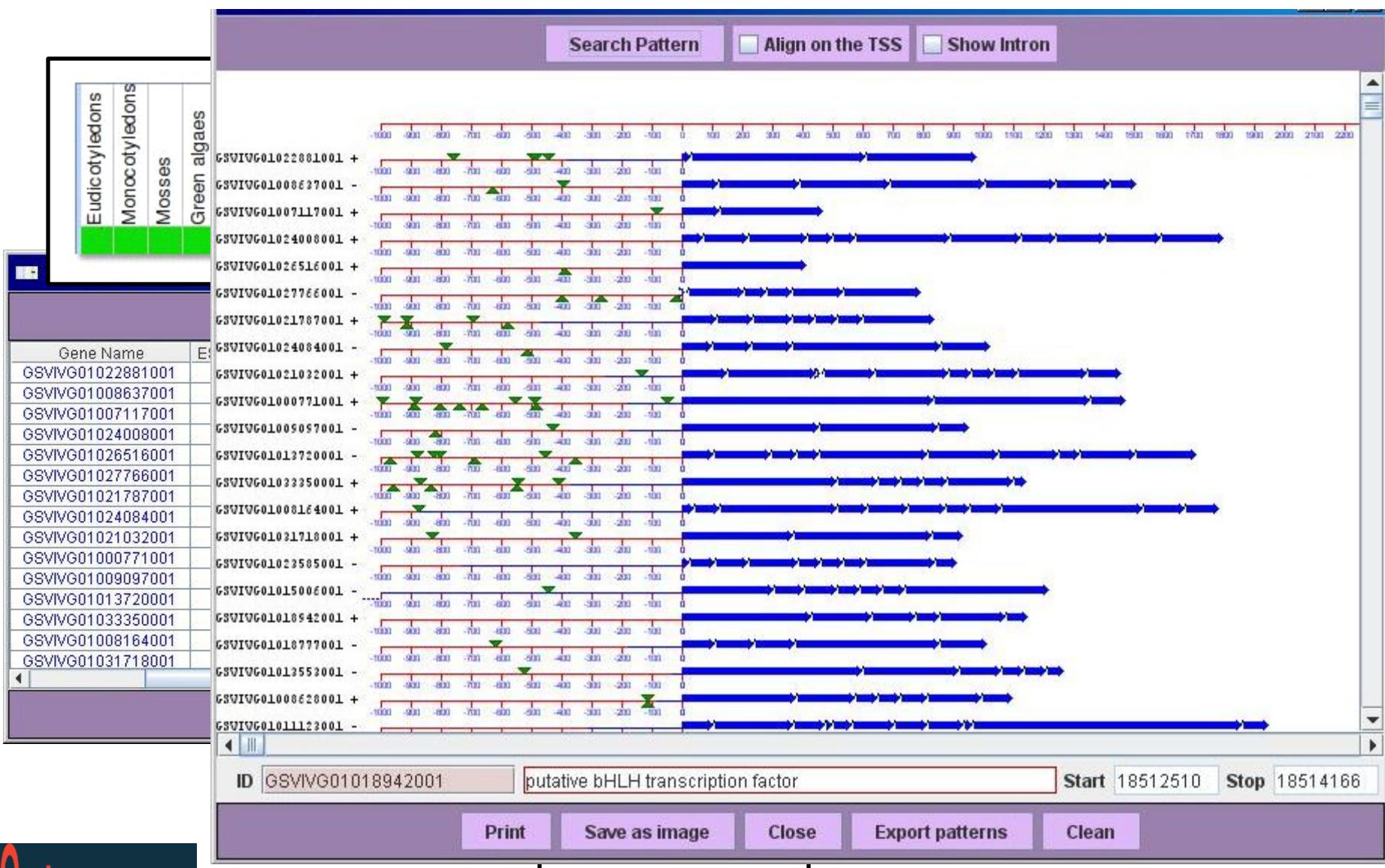


Vitis vinifera in FLAGdb⁺⁺



Easy access to gene families or large groups of genes

Vitis vinifera in FLAGdb⁺⁺



Vitis vinifera in FLAGdb⁺⁺



New tool dedicated to orthology relationships between the 4 plant genomes hosted in the database

Orthologues Map

Populus trichocarpa Oryza sativa Arabidopsis thaliana Vitis vinifera

Search Pattern Align on the TSS Show Intron

Clustal

Clustal it!

AT1G30670	AQSIIAARKRRRRI TEKTQELGKL IPGSQKHN-TAEMFNAAAKYVKFLQAQIEILQLKQT-
FGENESH4_PG....	TQSIAARERRRKITEKTRRELGFIPGGHKMN-TAEMFQAASKYVKFLQAQIGILELMGS-
GSVIVG010379...	AQSIIAARQRRRKITEKTQELGKLIPGGNKMN-TAEMFQAAFKYVKYLQAQVAILQLMGS-
OS12G31430	VQSIIARERRRRRISSKTAELSRLIPGAARMNSTAEMLQAAARHVRLLQAQVGMLALIHSS
AT2G34820	SQSIIAARGRRRRIA EKTHELGKLIPGGNKLN-TAEMFQAAAKYVKFLQSQVGILQLMQT-
ALIGN	*** *** * ***:***. **:***. ; * ***:*** :*: ; **: ***; : * * :

Feature information panel

ID AT2G34820 basic helix-loop-helix (bHLH) family protein Start 14696711 Stop 14697673 Length: 903

From grape genome to wine flavours...



CCAGTGTGCAGCTTGCAAGCGCCAAACGTGATGCCCTGCTTCAGTTCAAAGACGCTAACG
TGCTGTGGTGGAGTATGACCAGGCACTCATGACCTGAAGACCCCTGCTCTACACTACTTIG
AGGAGGCTGAATGGGATGGGTTGCAAGCTGACATACGCCCTCGTGTGCGGGTGGACC
CTGACGGGCCCTGCAAACTGCCTATGGCACAGGGCTGGTAGTTCTGCTTCCCAC
GAGAGAGCTGGCTGAGGAACATGAGGGCTCATGGCGAAGGGCAGAGGTCTAGCTTCTGC
CCAGCTACATCATGGTACGGGCTCTGGATGAGAACTACTAACATCATGGACTGCAGT
TCCCTGCTGGCTACTATGAGCCCCACCTGCTTATCTGTTGAGCCCAACAGACTTGGCCAG
GGCGGTGGCTGTGAGGCAGGACACGTGCTCATTGTCGCTGAACATCACACAGA
AAGTCCATCCAGTCATCTGGTCCCACAGCTGGCTTGAACCTGGCCCTGGCTG
TGCCCAAGGCGATAGGTGGGGTAGTGTCTTCGCTCAACTCACTGGTAGCTGAACCCAGA
GTGTTCCCCCATAGGTGGCTCTCAATAGTCTCACACCAGGCACACCCTGGCTTCCATTGC
GTACCCAAGAGGGTGTACCGGATCACCTAGACTGGCCACAGGGCCCTCATCTTATGACA
AGATGGTCATCTCCCTCAAGGGTGGTGAGATTATGTGTCACCTGACGGCATGC
GAAGTGTCCGAGCGTTCACTTGACAAGGCAGCTGCTAGTGTCTTACACCCAGTGGTCA
CAATGGAGCCGGATACTGTCTCTAGGCTCTGCCCTGGCAATTCCCTCTCCCTCAAGTACA
CGGAGAAGCTGCAGGAGCCCCCAGCCAGCTCCGTCGGAGGCTGCTGACAAGGAAGACCTGC

a gene family story



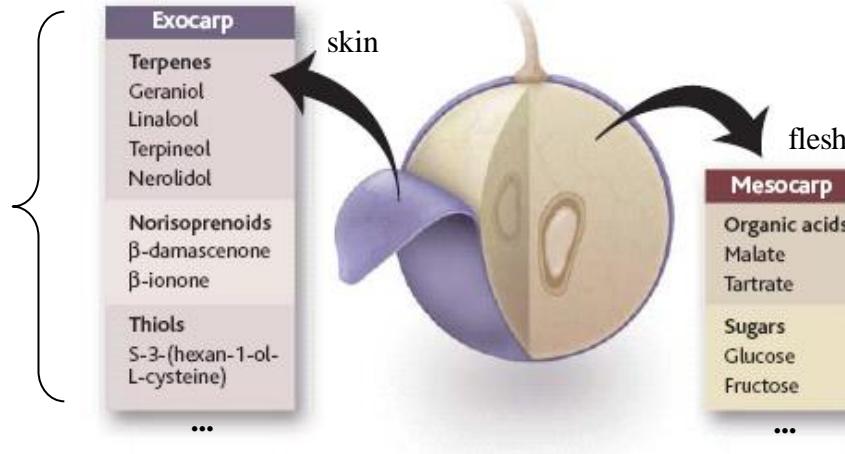
Many volatile organic compounds ► wine flavour (taste and aroma)

Stored as sugar or amino acid conjugates in vacuoles

glycosidases & peptidases

from grape, yeast and/or industrial enzymes

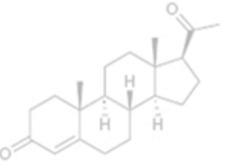
Volatilization of compounds essential for the flavour perception



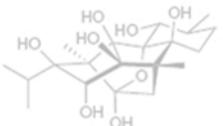
Lund and Bohlmann (2006) Science

The terpenoid world

(22 000 different molecules !)

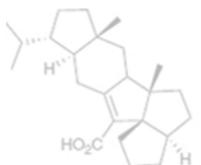


Gewürztraminer

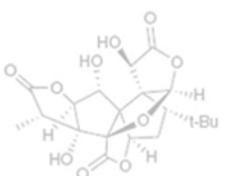


linalool
geraniol
nerol
citronellol
 α -terpineol
...

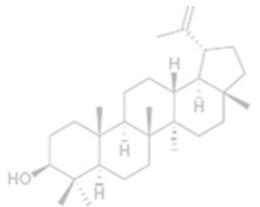
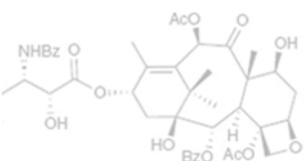
Distinctive floral character



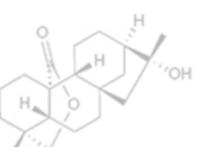
Girard et al. (2002) Am. J. Enol. Vitic.



Cabernet Sauvignon
anthers and pollen grains

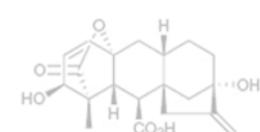
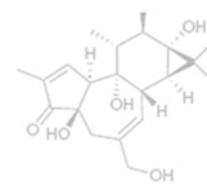
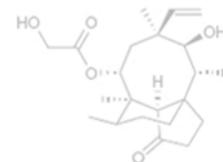


(+)-valencene
(-)-7-epi- α -selinene
...

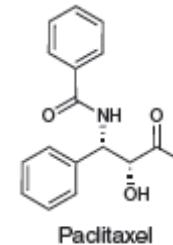
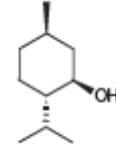
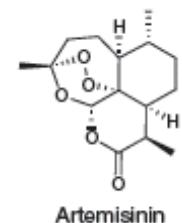
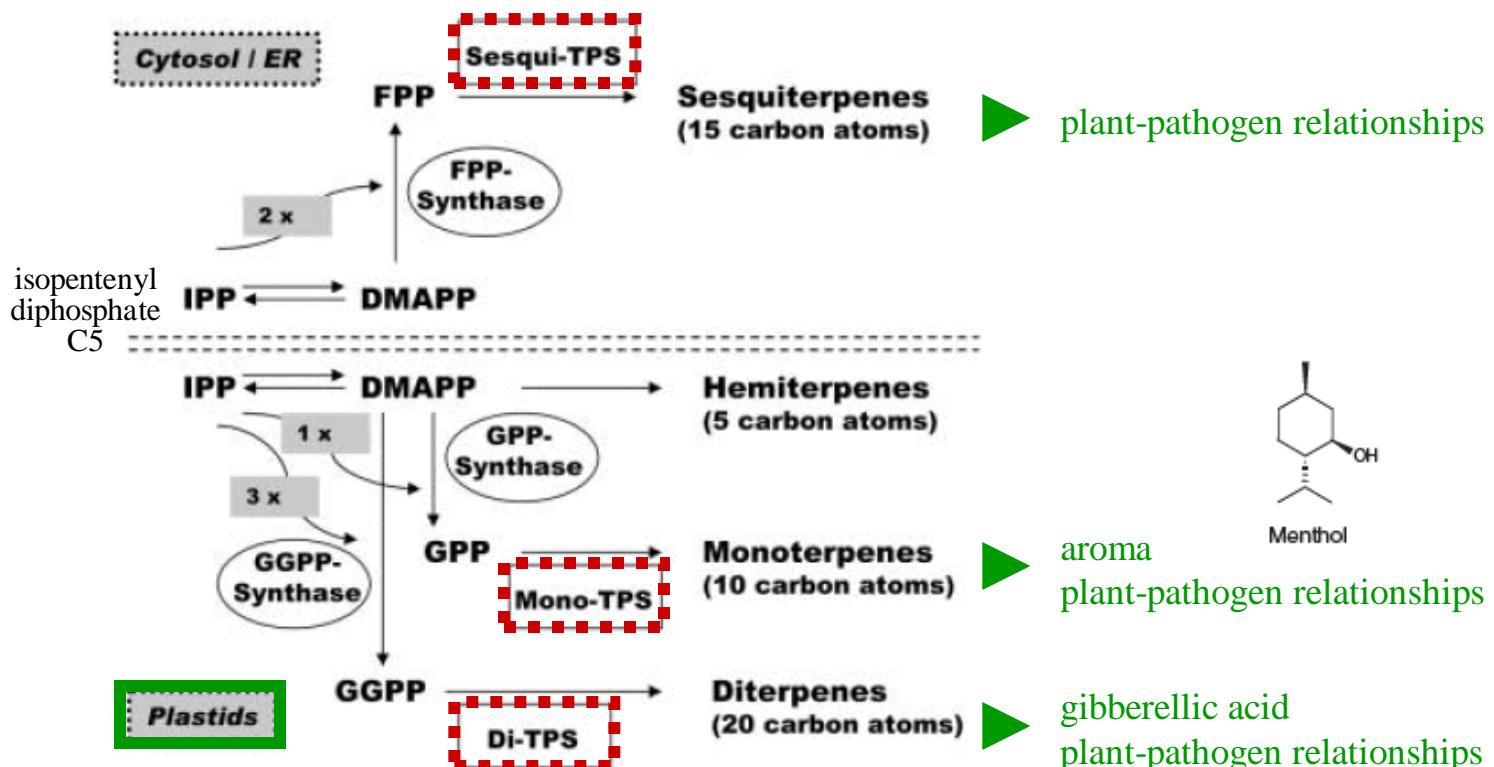


Protect reproductive tissues against pathogens
Allow pollinators to locate flowers

Martin et al. (2009) PNAS



The terpenoid biosynthesis : TPS family



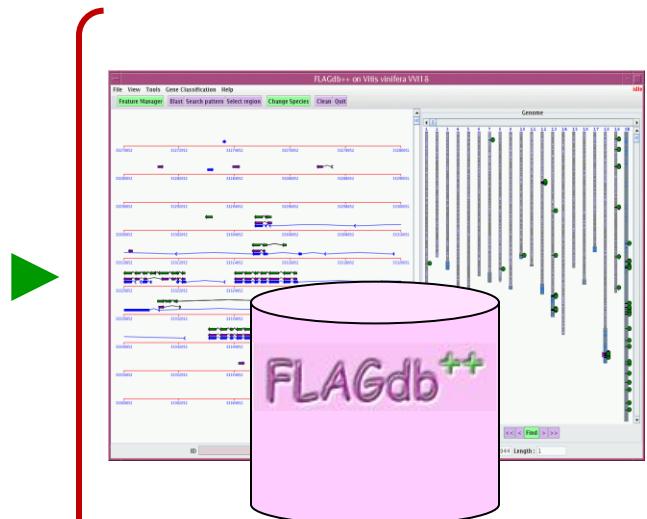
Expert annotation process



Known TPS proteins (Arabidopsis, Pinus, Citrus...)

TBLASTN
Vitis 12x genome
PN40024

152 loci TPS-like

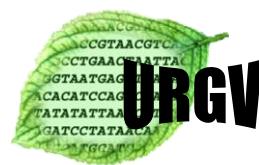


- ✓ GAZE predictions
 - ✓ EUGENE predictions
 - ✓ PFAM motifs
 - ✓ EST
 - ✓ cDNA
 - ✓ BLASTX results
 - ✓ TPS knowledge
- From
GB/EMBL
&
URGV

Clean and exhaustive structural annotation of the Vitis TPS gene family

- Detect and correct erroneous annotations and missing genes
- Discriminate between pseudogenes, partial and complete genes

The Vitis TPS family

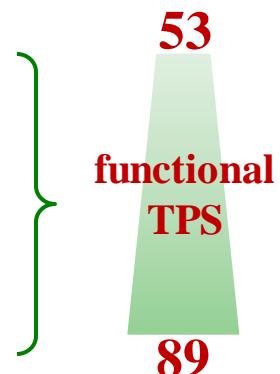


- GAZE fails to detect 54 of them
- EUGENE fails to detect 12 of them

- GAZE predicts perfectly 12 TPS (23%)
- EUGENE predicts perfectly 43 TPS (81%)

152 loci homologous to known TPS :

- **53 full** and perfect TPS genes
- **16 complete** TPS genes with only one punctual problem of sequence (sequencing error ?)
- **20 partial** TPS genes disrupted by unsequenced gap in the 12x assembly
- **63 pseudogenes** with numerous frameshifts, stop codons and/or large deletion(s)



Arabidopsis



32 full TPS
8 pseudogenes

Rice



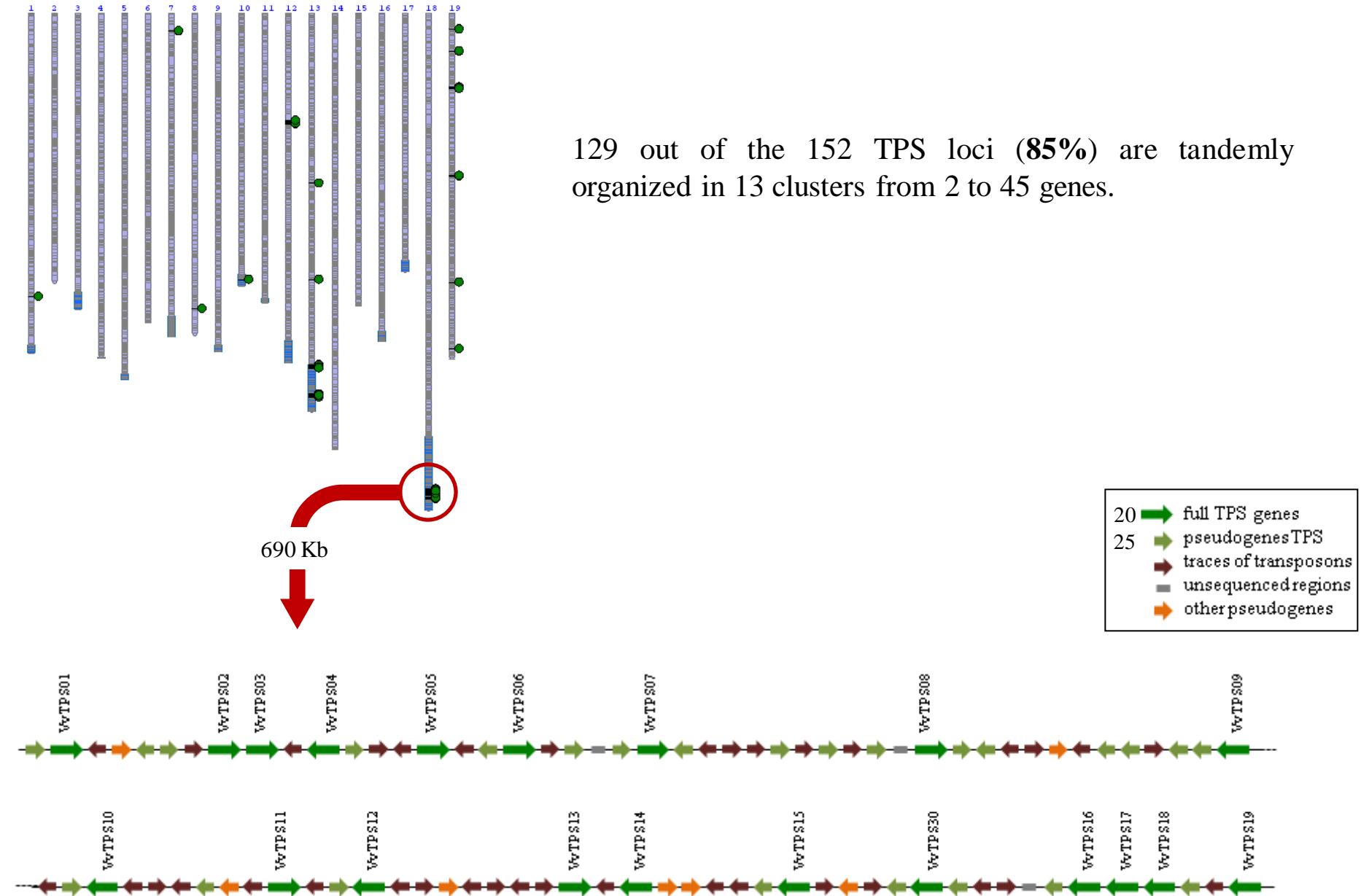
48 TPS loci
including pseudogenes

Poplar



45 TPS loci
including pseudogenes
and partial genes

Topological organization of the family

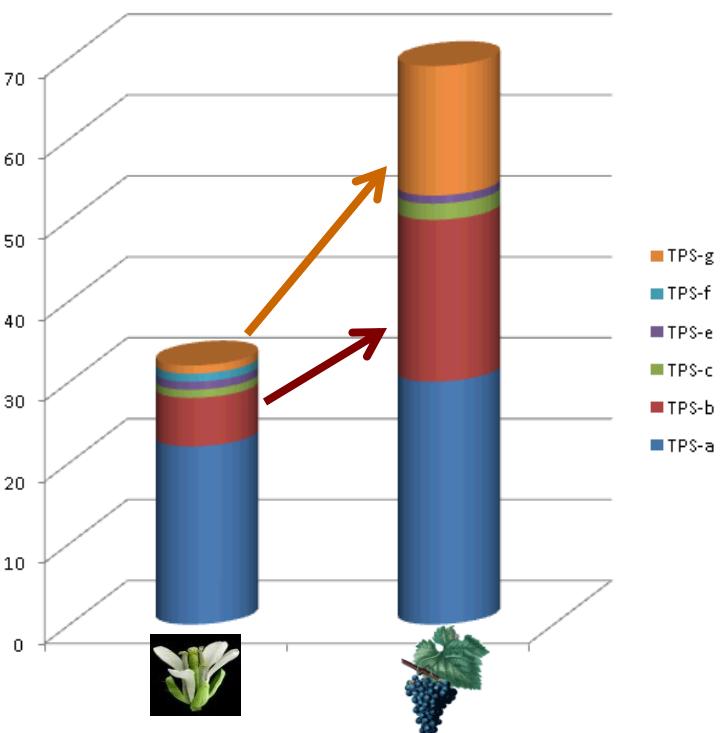
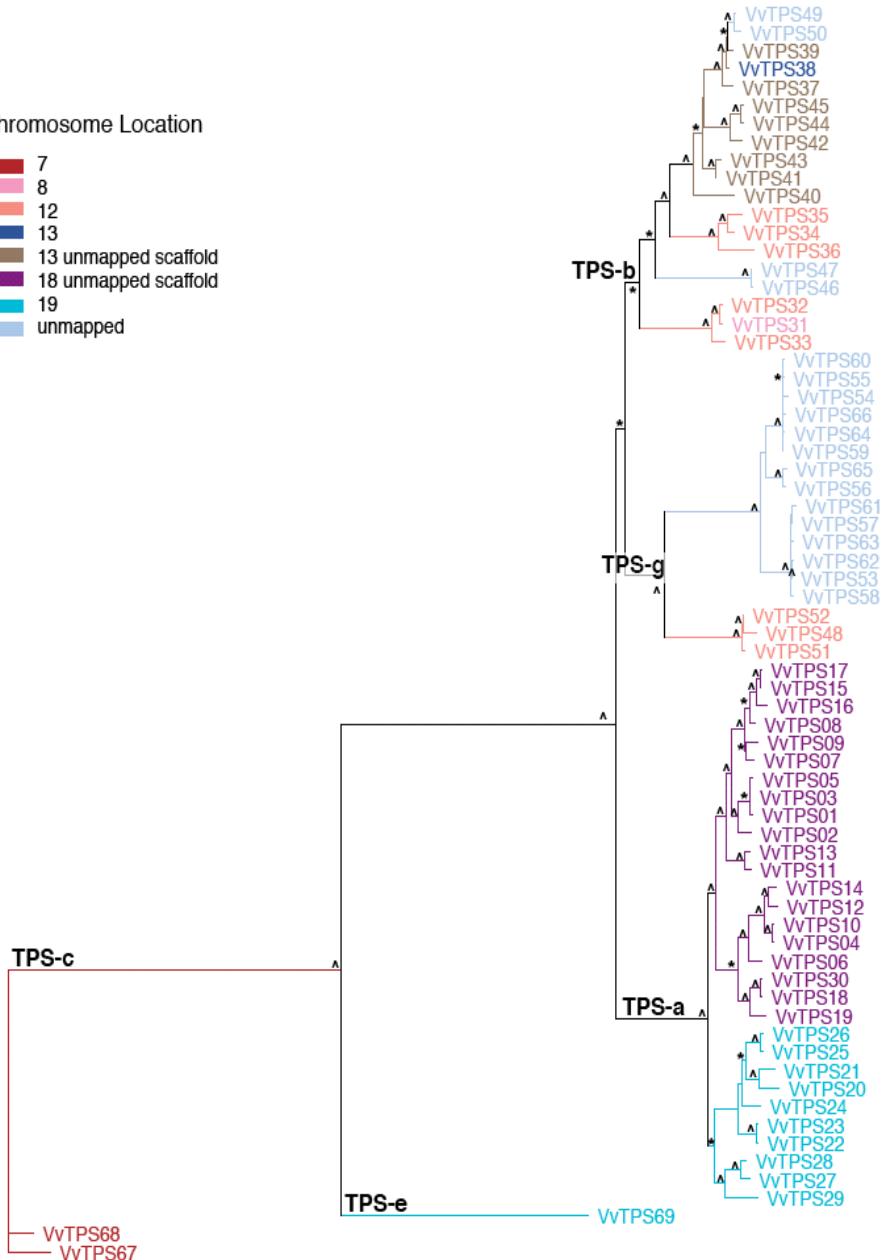


Vitis TPS classification



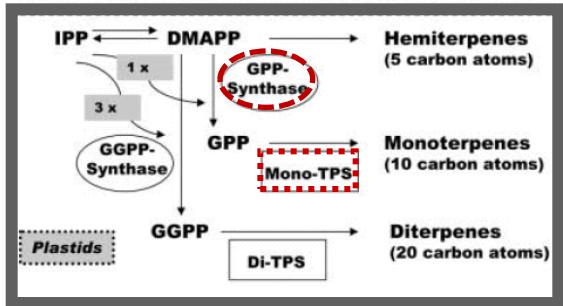
Chromosome Location

- 7
- 8
- 12
- 13
- 13 unmapped scaffold
- 18 unmapped scaffold
- 19
- unmapped



► high diversity
of monoterpenes
in PN40024 ?

Geranyl diphosphate synthases (GPPS) family



In *Arabidopsis thaliana* :
homodimeric GPPS

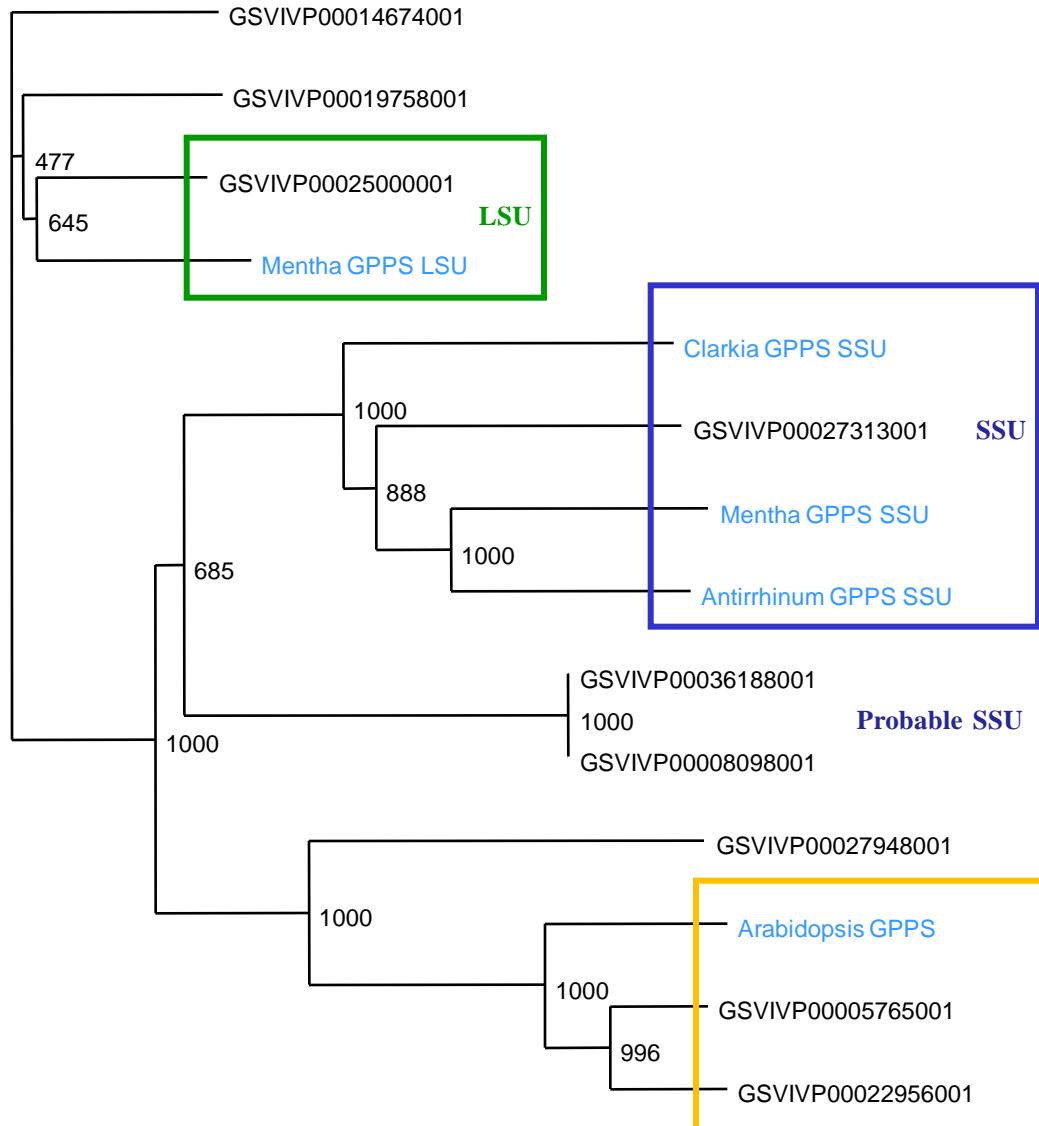
In *Mentha piperita* and *Clarkia breweri* :
heterodimeric GPPS

LSU + SSU

In PN40024, there are orthologous
genes for the 2 forms of GPPS



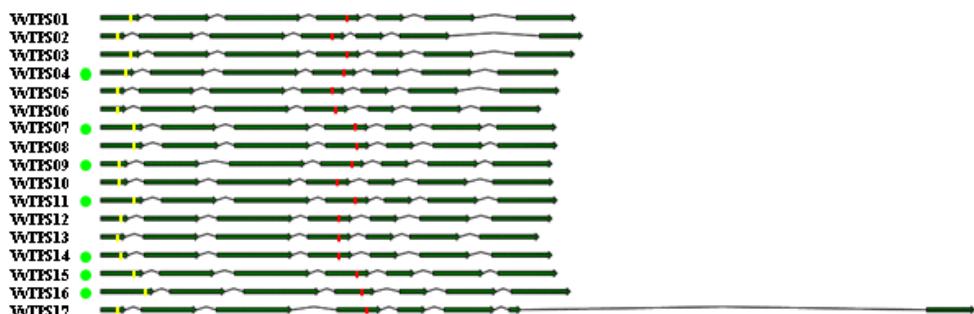
high diversity
of monoterpenes



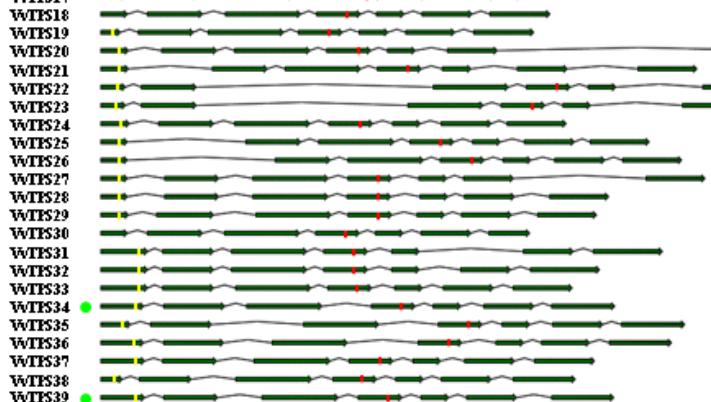
TPS gene structures



TPS-a



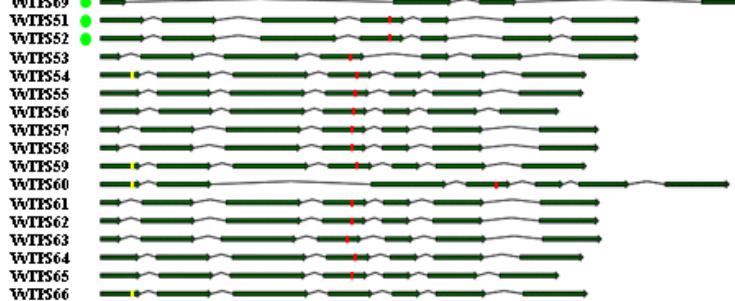
TPS-b



TPS-c
TPS-e



TPS-g

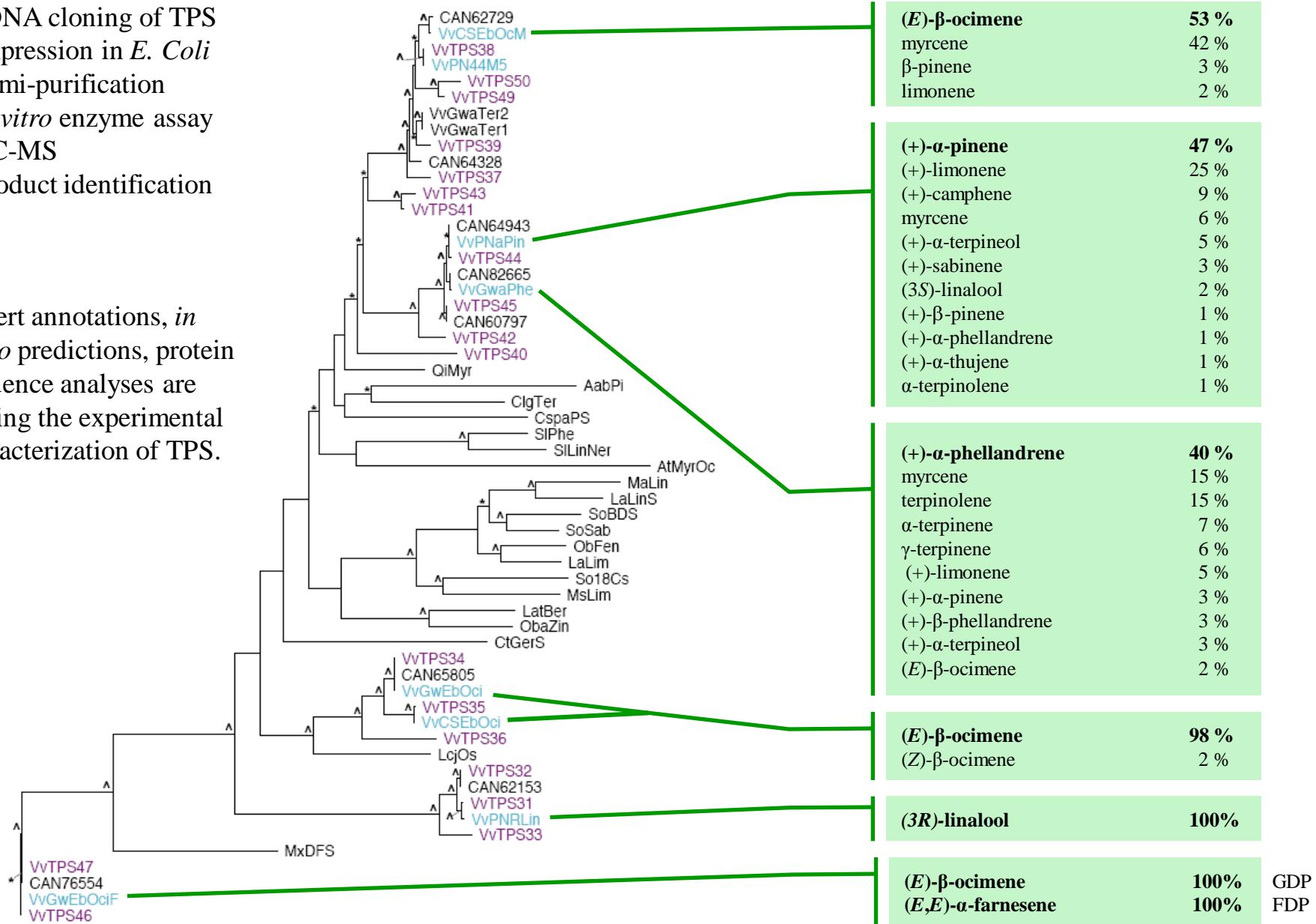


- predicted plastidial targeting peptide
mono- and diterpene synthases
- RR_x(8)W motif
isomerization-cyclization reaction
- DDxxD motif
metal ion binding for cleavage
of prenyl diphosphate substrates

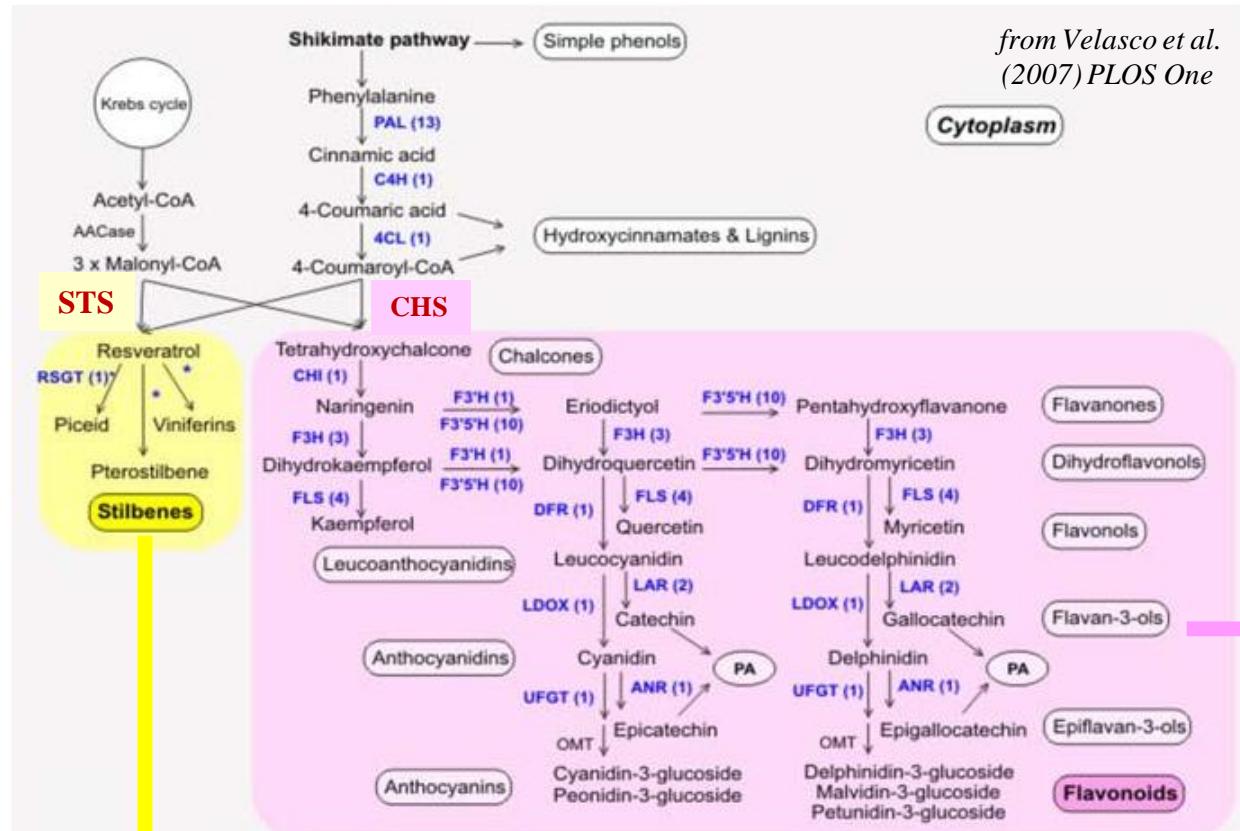
Functional characterization of VvTPS genes

1. cDNA cloning of TPS
2. Expression in *E. Coli*
3. Semi-purification
4. *In vitro* enzyme assay
5. GC-MS
6. Product identification

Expert annotations, *in silico* predictions, protein sequence analyses are driving the experimental characterization of TPS.



Other wine constituents...



Colour and astringent
features of wines

Resveratrol is involved in health benefits : 'French paradox'
Prevents cancer and cardiovascular disease.
Stilbenes are also known as fungicides.

► 48 Stilbene Synthase genes !



Acknowledgements



French-Italian Public Consortium for Vitis
Anne-Françoise Adam-Blondon
Olivier Jaillon, Patrick Wincker *et al.*



Franck Samson
Cécile Guichard
Sandra Derozier
Jean-Philippe Tamby
Véronique Brunaud



Jérôme Gouzy
Thomas Schiex



Joerg Bohlmann
Diane Martin



Christophe Caron



Philippe Hugueney



flagdb@evry.inra.fr

aubourg@evry.inra.fr

<http://urgv.evry.inra.fr/FLAGdb>

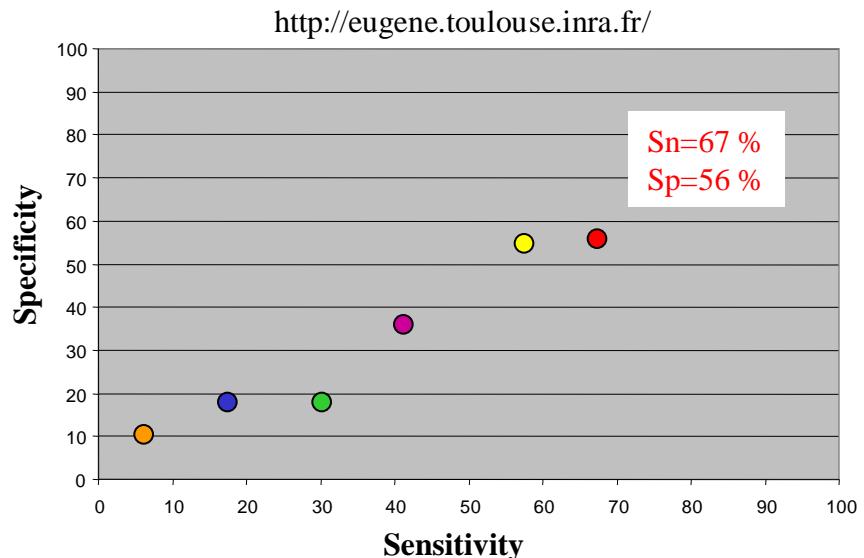


Why EuGène ?

An eukaryotic gene finder that combines multiple sources of evidence.
 T. Schiex *et al.* (BIA, INRA Toulouse)



- FgenesP
- GenScan
- GlimmerA
- GeneMark.hmm
- FgeneSH
- EuGène



S. Aubourg *et al.* (URGV, Evry)



S. Rombauts *et al.* (PSB, Gent)



J. Gouzy *et al.* (LIPM, Toulouse)



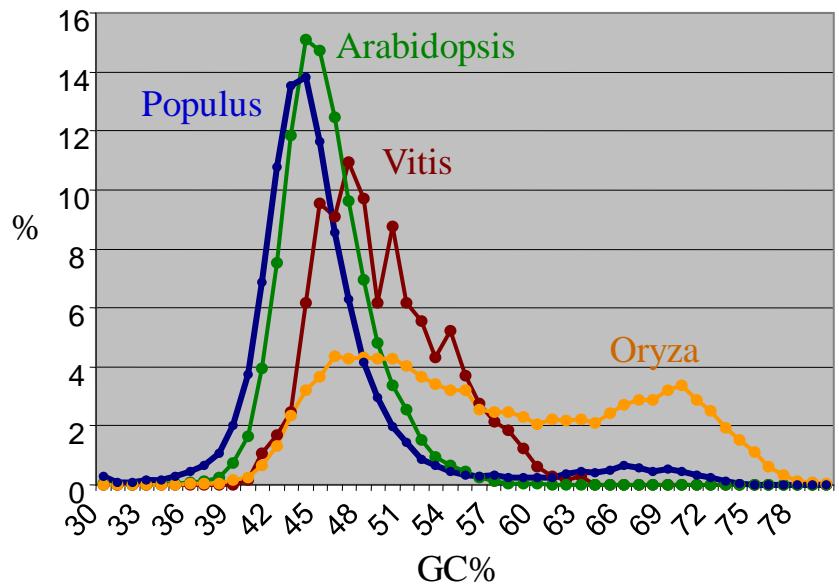
J. Amselem *et al.* (URGI, Versailles)

...

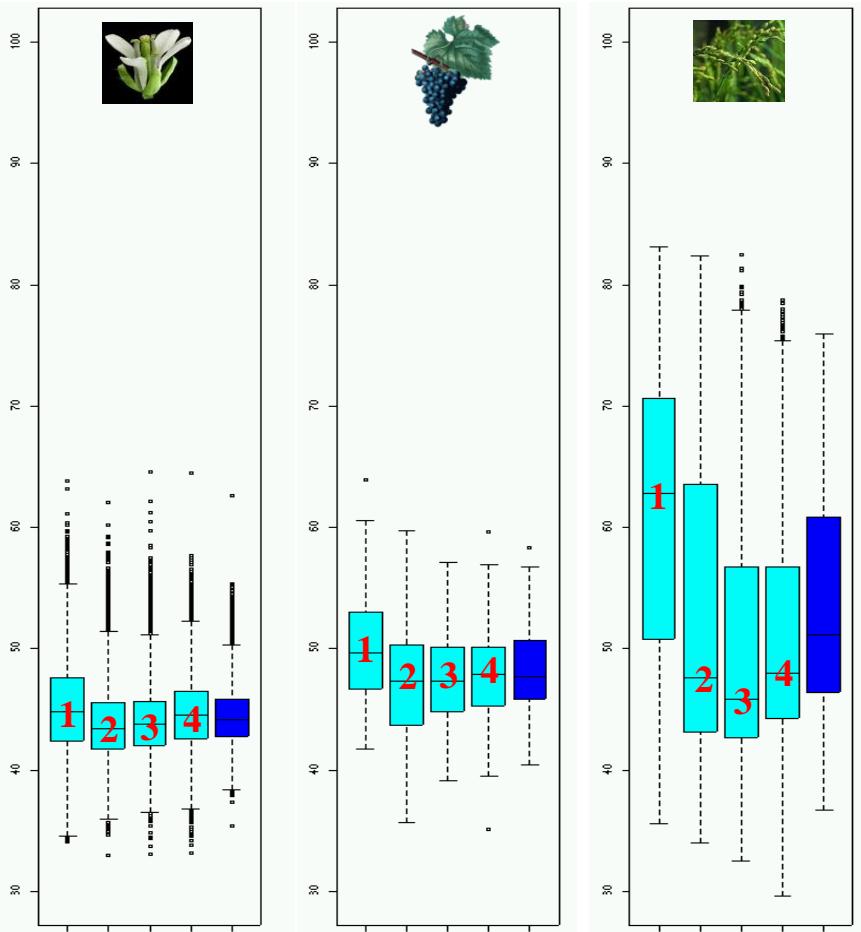
Gene training set



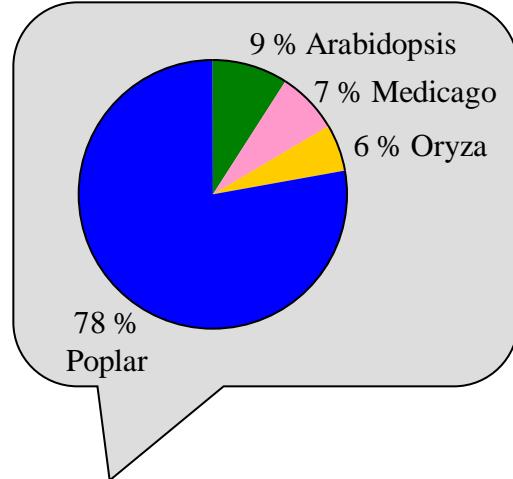
Set of 600 experimentally characterized complete genes (full-length mRNA cognate to genomic DNA) ‘representative’ of the whole genome + 4500 splicing sites



GC content



Integration and optimization steps



Tested formulas (plug-ins):

- 1- IMM + SpliceMachine
- 2- IMM + SpliceMachine + BLASTX
- 3- IMM + SpliceMachine + BLASTX + TBLASTX
- 4- IMM + SpliceMachine + BLASTX + TBLASTX + GenomeThreader
- 5- IMM + SpliceMachine + BLASTX + GenomeThreader



- SwissProt
proteomes from :
- Arabidopsis
- Oryza
- Poplar
- Medicago

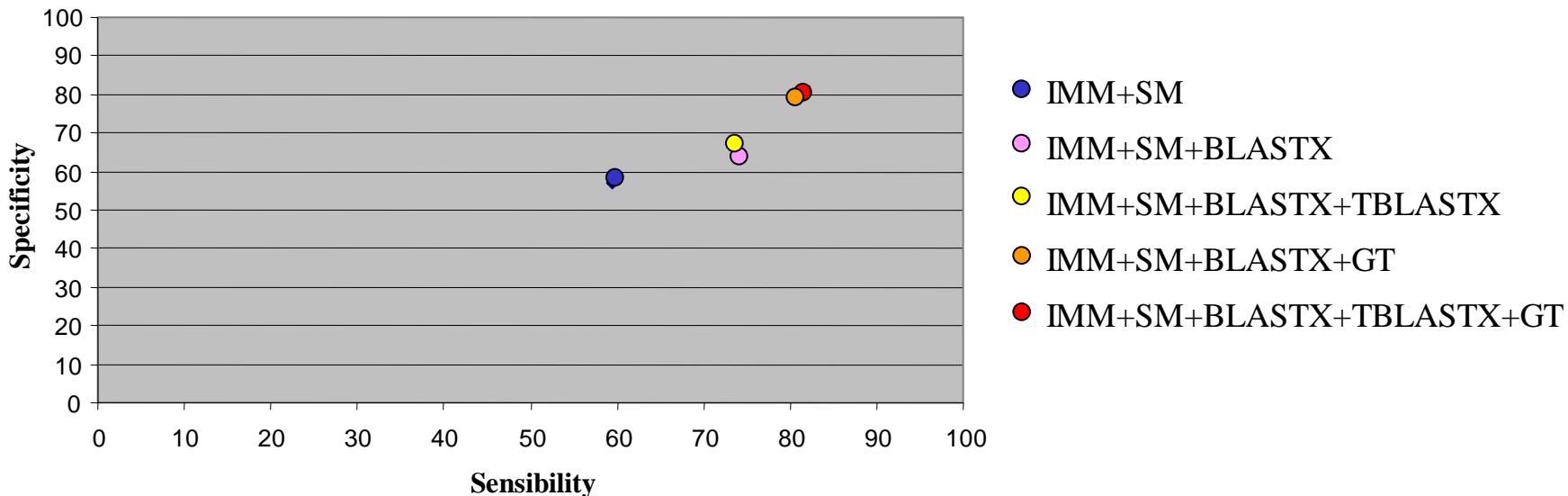


Poplar
genome

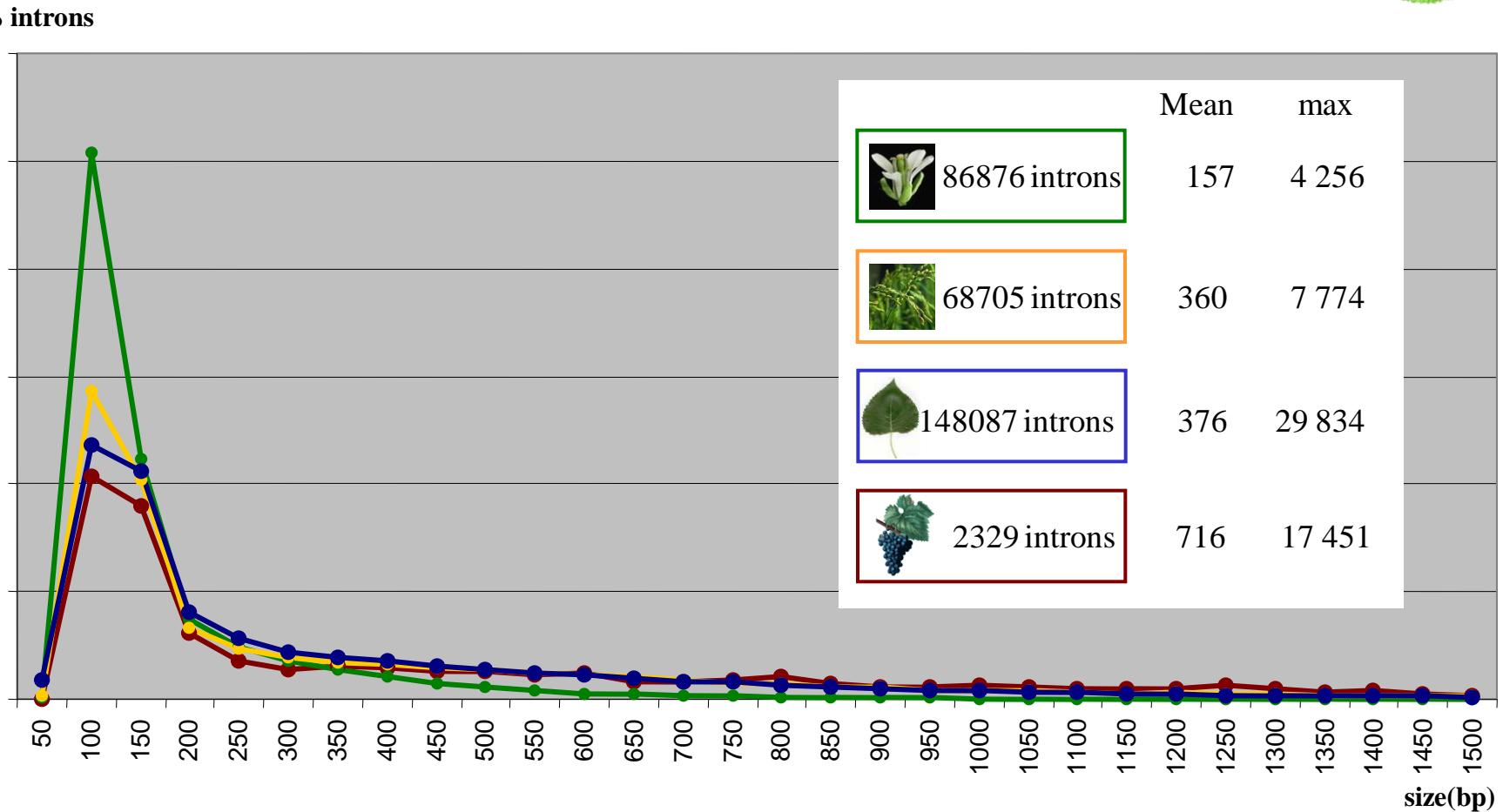


EST
cDNA

EUGENE training results

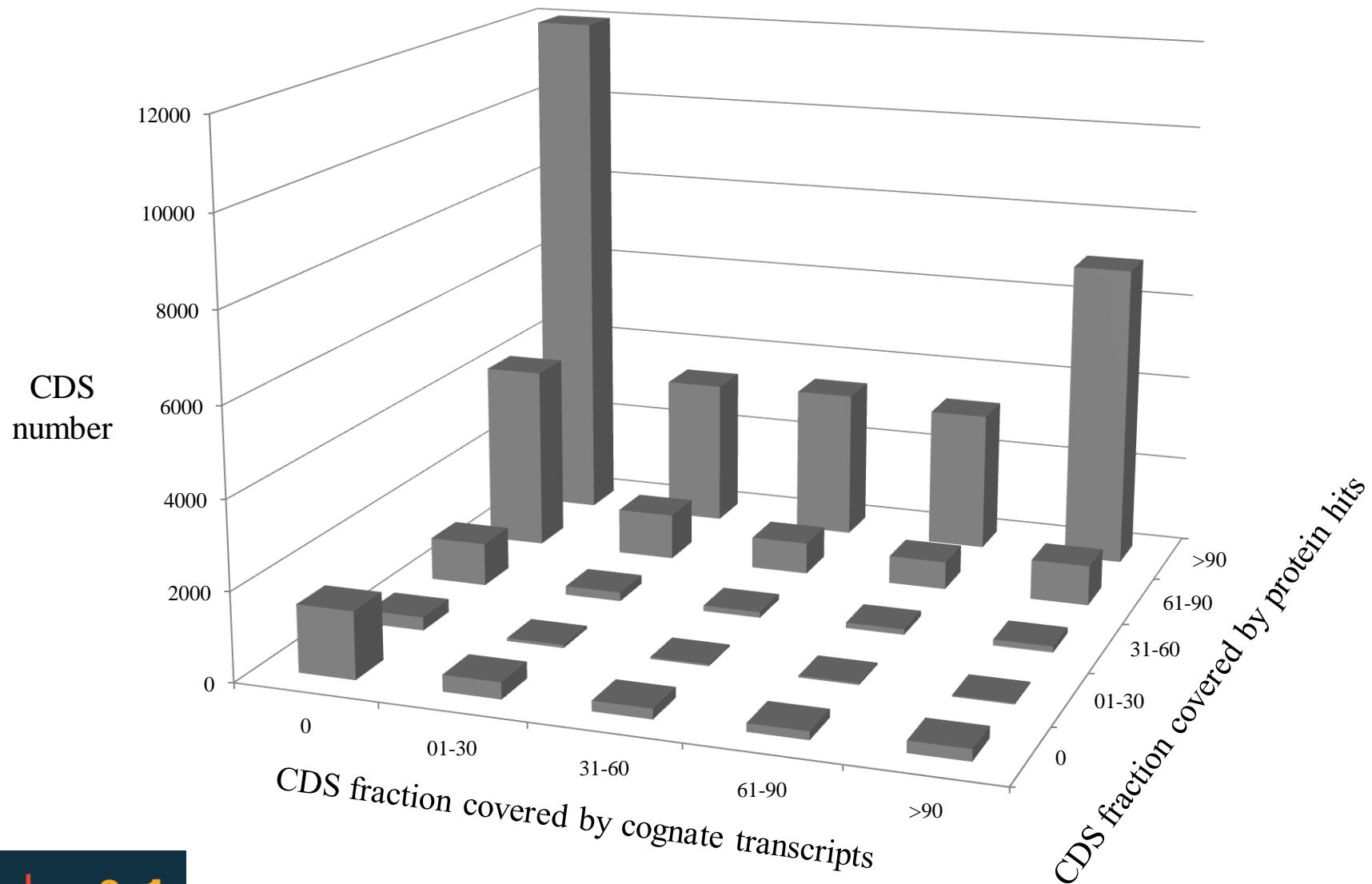


Intron size (training set)

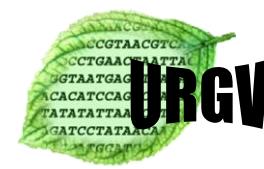


Long introns (12% include LINEs elements)

EuGène results

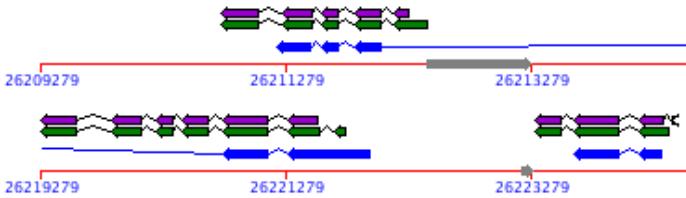


EuGène vs GAZE (TPS family)



Out of 152 TPS loci :

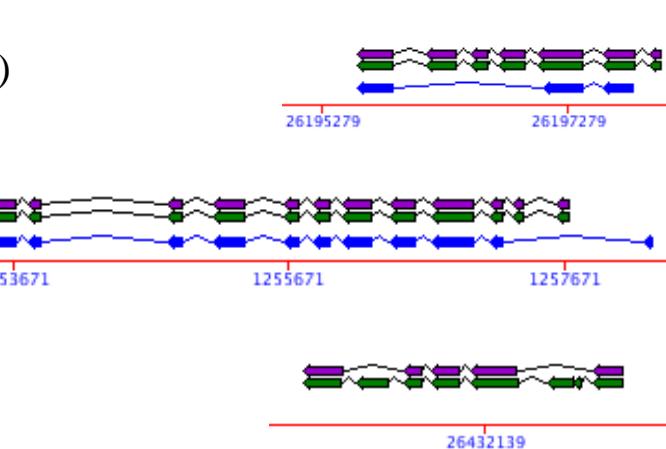
- GAZE fails to detect 54 of them (including 6 full TPS genes)
- EUGENE fails to detect 12 of them



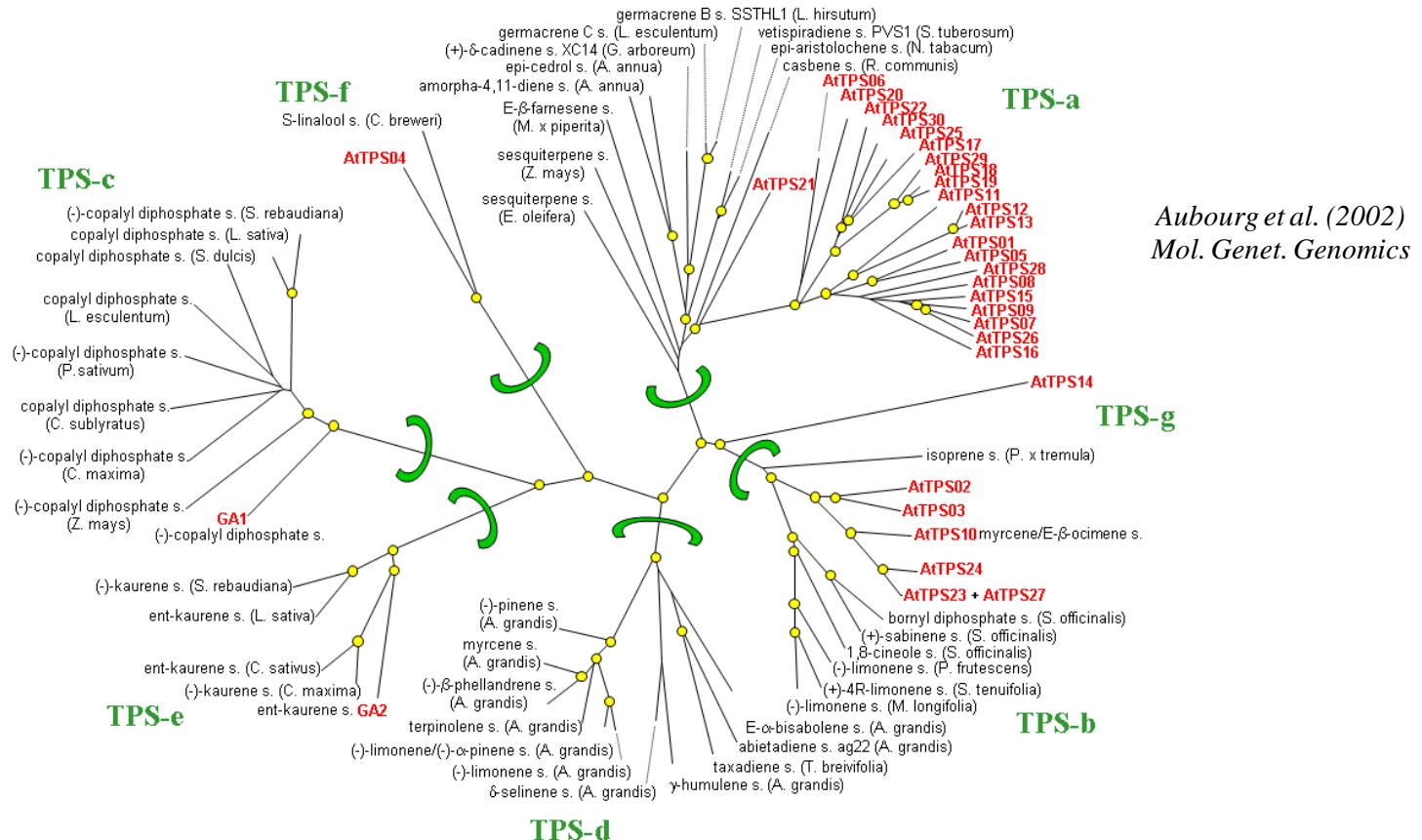
Out of 53 complete and perfect TPS genes:

(partial and disturbed genes can not be well predicted)

- GAZE predicts perfectly 12 TPS 23%
- EUGENE predicts perfectly 43 TPS 81%



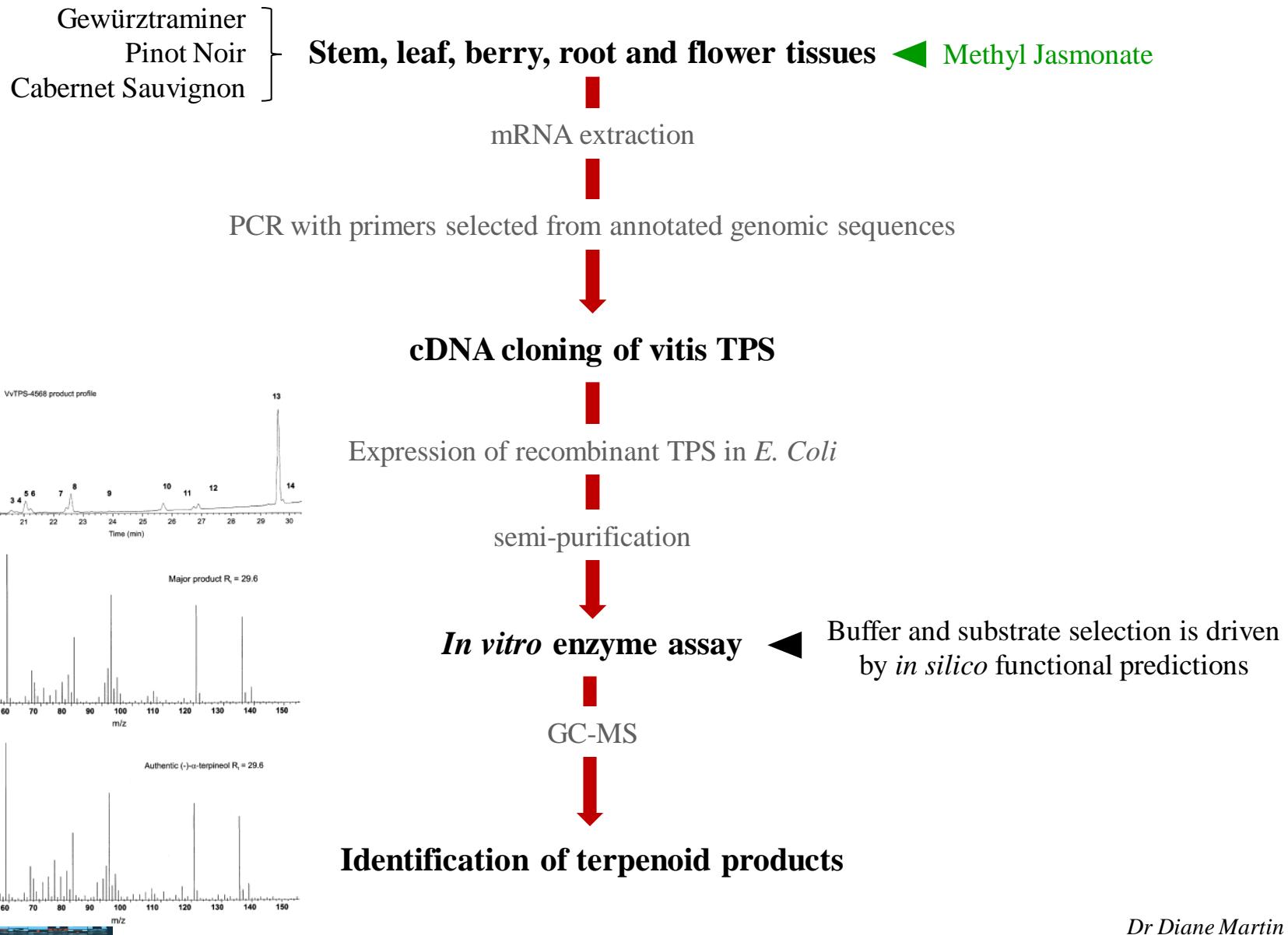
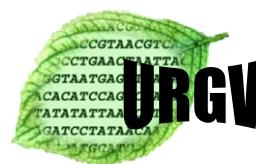
TPS classification

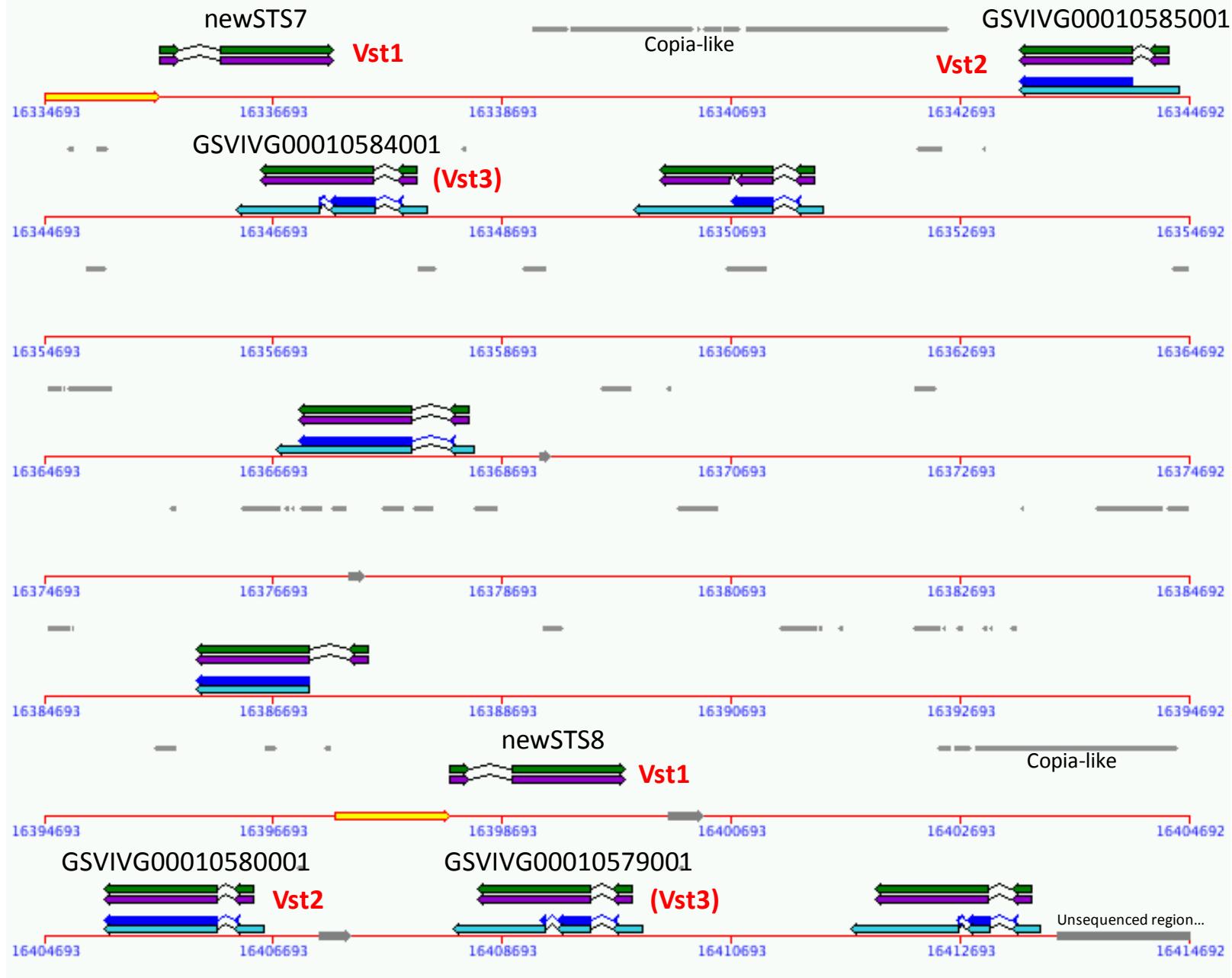


Aubourg et al. (2002)
Mol. Genet. Genomics

-
- TPS-a** Sesquiterpene and Diterpene Synthases
 - TPS-b** Monoterpene Synthases
 - TPS-c** Diterpene Synthases (copalyl diphosphate synthases)
 - TPS-d** TPS from conifers
 - TPS-e** Diterpene Synthases (ent-kaurene synthases)
 - TPS-f** Monoterpene Synthases (linalool synthases)
 - TPS-g** unknown TPS

Functional characterization of TPS genes





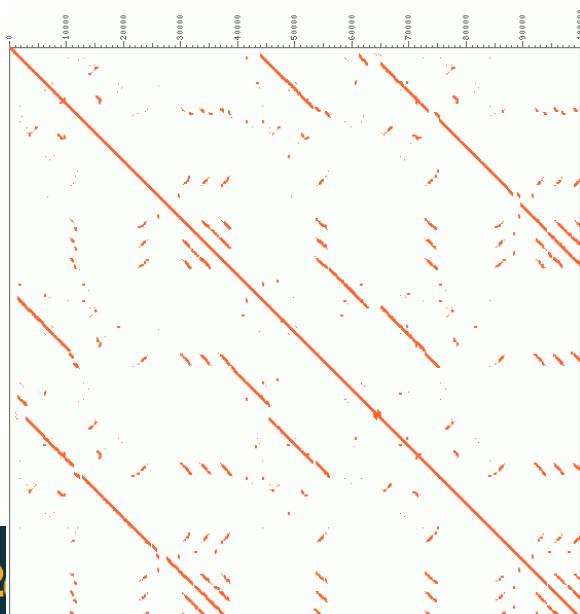
Stilbene and Chalcone Synthase families



Proteome and translated genome have been screened with sequences of 3 previously characterized Stilbene Synthases :

► 62 loci identified and expertised :

- 3 CHS genes highly expressed and 11 CHS-like
- 48 STS genes (discriminated with specific residues) distributing on 2 clusters (chromosomes 10 and 16). At least 73% of them are expressed. Only 2 grape STS have been experimentally characterized. All the STS share between 91% and 99.7% of identity.



Out of 62 STS/CHS genes (2 exons):

- GAZE fails to detect 36 of them and only 2 are well predicted !
- EUGENE fails to detect only one and 40 are well predicted.