# GENES EXPRESSION MEASUREMENT BY USING THE GENOME SEQUENCES OF *O. CUMANA* AND SUNFLOWER

**Stéphane Muños**
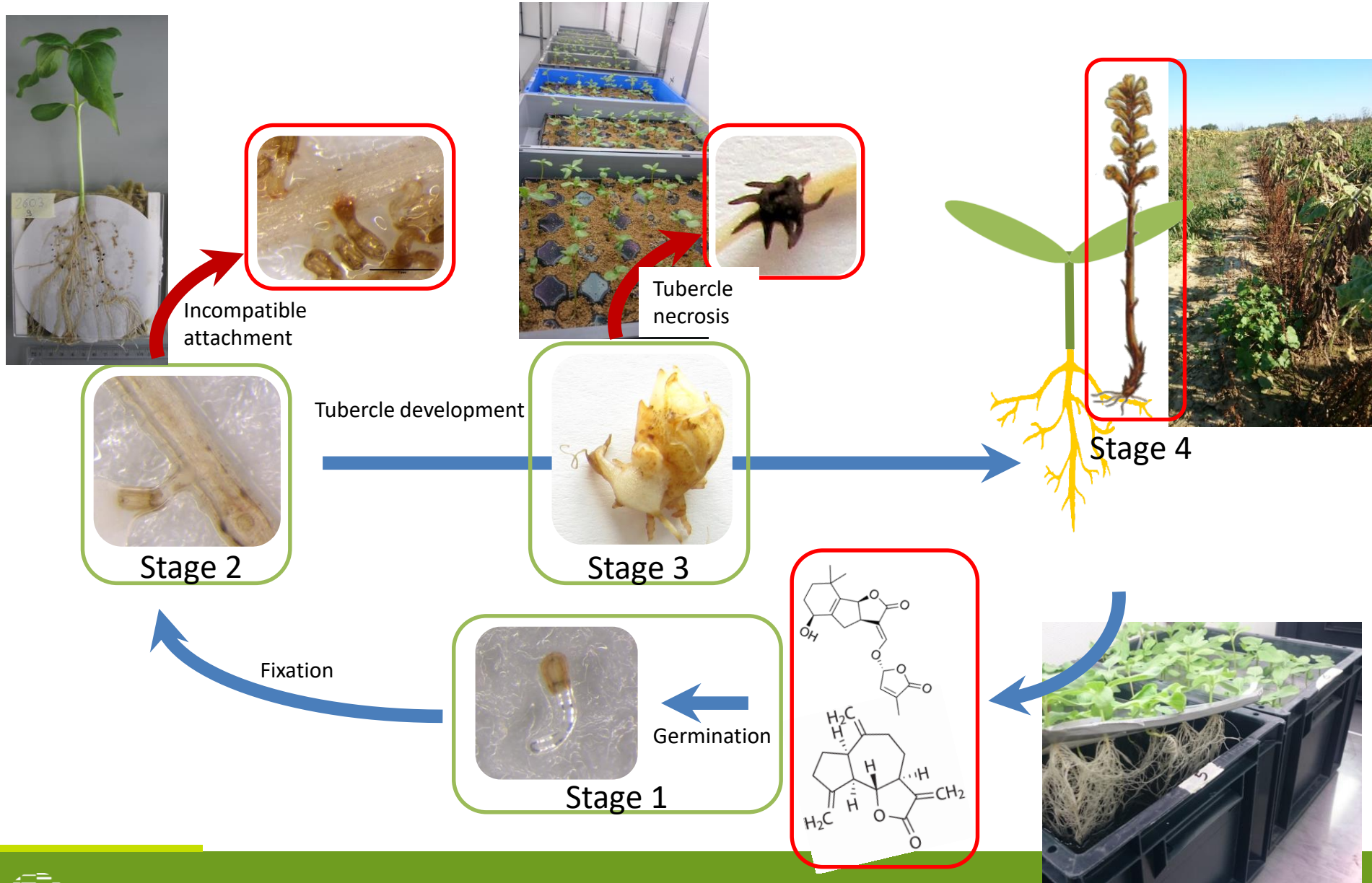**stephane.munos@inra.fr**
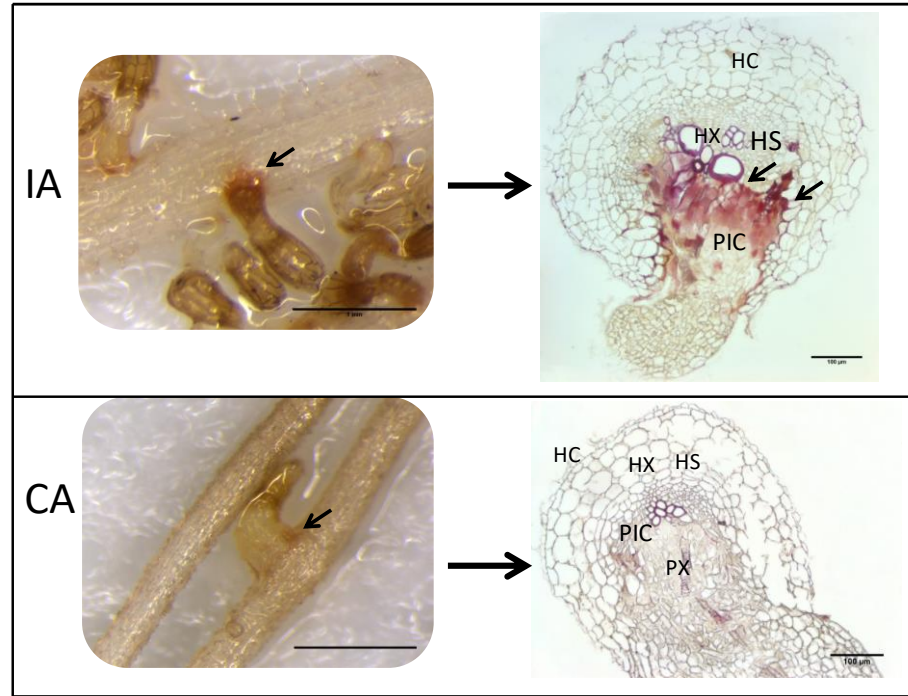**@stephane_munos**

INRA Occitanie-Toulouse

**ISBS, July 3rd, 2018, Bucharest**

# Biology cycle of *Orobanche Cumana and resistance mechanisms*



Incompatible attachment

Tubercle necrosis

Tubercle development

Stage 2

Stage 3

Stage 4

Fixation

Germination

Stage 1

# Prevent connection to the vascular system in sunflower root :
## *a key resistance mechanism*



Incompatible attachment

Stage 2

IA

CA

Barrier formation: **no connection** to the vascular system of the host

**Connection** to the vascular system of the host

**Cell wall modifications seem to be involved both compatible and incompatible attachments.**
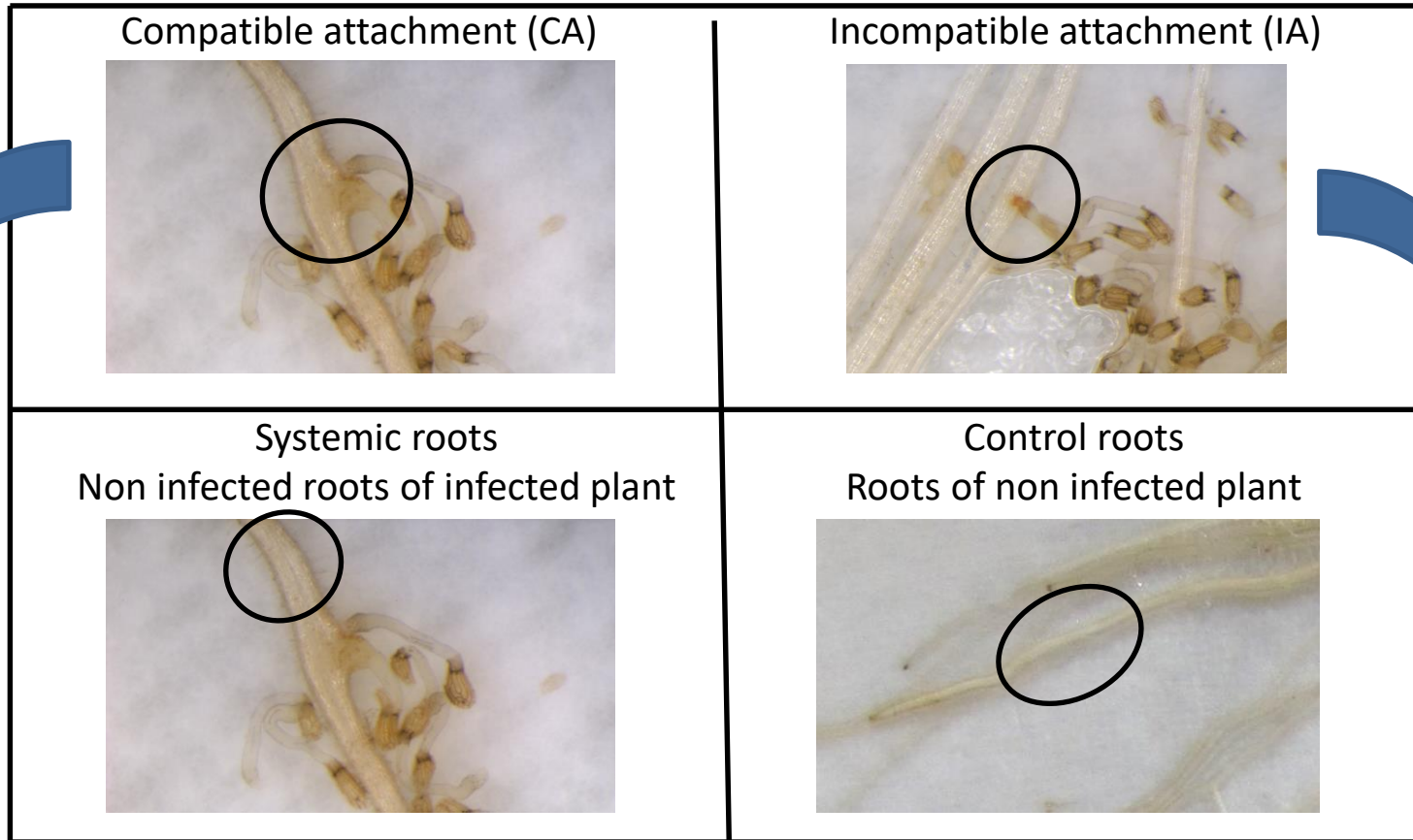**Which genes are involved?**

# Understanding the molecular and cellular mechanisms involved in the incompatible attachments could help in identifying new candidate genes for the resistance

| Genotype | IA rate | # of Healthy tubercles | # of necrotic tubercles | Total number of tubercles |
|---|---|---|---|---|
| LR1 | **66.5** | 2.4 | 0 | 2.4 |
| HA89 | **7.1** | 3.4 | 0 | 3.4 |
| RIL334 | **73.4** | 0.2 | 0.2 | 0.4 |
| 2603 | **0** | 13 | 0 | 13 |

*O. cumana* race F used for infection

# Sampling for transcriptomic analysis



Compatible attachment (CA)

Incompatible attachment (IA)

Systemic roots
Non infected roots of infected plant

Control roots
Roots of non infected plant

**Genes from both broomrape and sunflower are expressed in the attachments samples**

**RNASeq experiment is the main method to measure the expression of thouthands of genes simultaneously.**

**But reference sequences of all genes are needed**

# Genomics of parasitic plants and of their hosts

| Feature | *Mimulus guttatus* | *Triphysaria versicolor* | *Striga hermonthica* | *Orobanche aegyptiaca* |
|---|---|---|---|---|
| Nutrition | Autotrophic | Hemiparasite | Hemiparasite | Holoparasite |
| Dependence on host | Free living | Facultative | Obligate | Obligate |
| Genome size (Mb/1C) | 430 | 1975 | 1672 | 3900 |
| Chromosome number (2N) | 28 | 22 | 38 | 24 |
| Hosts with abundant sequence information (model hosts) | N/A | *Arabidopsis*, *Medicago*, tomato | Maize, rice, sorghum | *Arabidopsis*, tobacco, tomato |

Loss of photosynthesis

Origin of terminal haustorium (obligate parasitism)

Origin of lateral haustorium (parasitism)

TRENDS in Plant Science

**Illustration from Westwood *et al.*, 2010**

**Parasitic Plant Genome Project**
**http://ppgp.huck.psu.edu**

Transcriptome libraries using 454 sequencing (Roche)

Fully sequenced (Hellsten *et al.*, 2013)
https://phytozome.jgi.doe.gov

Sanger shotgun sequencing

# scaffolds: 1507
Sequence size total: 312.7 Mb
(72.7% of the genome, 7.3% N)
N50: 21.2Mb (7 scaffolds)

33573 protein-coding transcripts

# Genomics of parasitic plants and their hosts

| Feature | *Mimulus guttatus* | *Triphysaria versicolor* | *Striga hermonthica* | *Orobanche aegyptiaca* | *Orobanche cumana* |
|---|---|---|---|---|---|
| Nutrition | Autotrophic | Hemiparasite | Hemiparasite | Holoparasite | Holoparasite |
| Dependence on host | Free living | Facultative | Obligate | Obligate | Obligate |
| Genome size (Mb/1C) | 430 | 1975 | 1672 | 3900 | 1420 |
| Chromosome number (2N) | 28 | 22 | 38 | 24 | 38 |
| Hosts with abundant sequence information (model hosts) | N/A | *Arabidopsis, Medicago,* tomato | Maize, rice, sorghum | *Arabidopsis,* tobacco, tomato | sunflower |

→Fully sequenced

Fully sequenced

**Illustration from Westwood *et al.*, 2010**

Loss of photosynthesis

Origin of terminal haustorium (obligate parasitism)

Origin of lateral haustorium (parasitism)

*TRENDS in Plant Science*

**Genomic resources available make *O. cumana*–sunflower a unique pathosystem**

# The HeliOr project

A public-private joint French-Spanish consortium

*Sunflower team:*
**Sunflower genetics and genomics**
*O. cumana* **diversity**
*Bioinformatics team:*
**Genome assembly and annotation**
**SNP calling**
**Informatic tools for biologists**

*O. cumana* **biology**

**Sequencing platform**

**Sunflower genetics**
*O. cumana* **diversity**

**BAC libraries**
**Genomic tools (optical map...)**

**Advices to farmers**
**List of infected fields in France**

**Sunflower genetics**
*O. cumana* **genetics and diversity**

Fundings:

# Sunflower and *Orobanche cumana* in Spain

- **Race F has been the most virulent race until recently in both main areas of occurrence of sunflower broomrape: Cuenca (CU) and the Guadalquivir Valley: two distinct genetic pools**

Castilla y León 261,983 has

Castilla La Mancha 198,252 has

Andalucía 354,767 has

VA
CU
BA
CR
AB
CO
HU
SE
CA

INSTITUTO DE AGRICULTURA SOSTENIBLE **IAS**

**IN23**
Used for genome sequencing (selfed three times to increase homozygosity)

INRA SCIENCE & IMPACT

# *O.cumana*: a nightmare not only for sunflower.
## For bioinformaticians too!!!

J. Gouzy



**A lot of repeats (like sunflower, i.e. 33% of the genome) but longer than in sunflower**

# Long read sequencing using PacBio RSII (Pacific Biosciences)

## P6-C4: Read Length Performance



Half of data in reads: > 14 kb
Top 5% of reads: > 24 kb
Maximum read length: > 40 kb
Data per SMRT® Cell: 500 Mb – 1 Gb (in 4 hours)

**PacBio produces sequences longer than the known repeats**

P6-C4, 4-hr movie, 20-kb BluePippin™ size-selected *E. coli* library (1 SMRT Cell)

# PacBio data

N. Pouilly

- All data produced at GeT-Plage (INRA)

PacBio *RS* II

- Produced from October 2015 to February 2016

  - 100X depth expected
  - 126 SMRT Cells (mean: 1.19 Gb/SMRT Cell)

# Statistics of contig sequences after assembly of the raw data

| Steps | NUM | MAX (Mp) | N50 (Mp) | NUM >=N50 | MEAN (bp) | MEDIAN | Total (Gb) |
|---|---|---|---|---|---|---|---|
| Raw data (subreads) 126 SMRT Cells | 13.2M | 85.05 | | | | | **149.9** |
| Corrected reads (CANU) | 7.04M | 55.53 | 13.98 | 2.01M | 10651 | 9777 | **75.06** |
| | NUM | MAX (Mb) | N50 (Mb) | NUM >=N50 | MEAN (Mb) | MEDIAN (Mb) | Total (Gb) |
| Genome assembly (CANU) | 905 | 16.88 | 3.57 | 107 | 1.53 | 6.49 | 1.388 |
| Remove spurious + Sequence based Scaffolding + polishing (QUIVER) | **793** | 16.98 | **4.21** | 96 | 1.74 | 7.43 | **1.380** |

/2

INRA
SCIENCE & IMPACT

# From contigs to chromosomes sequences!

## Using genetic map and optical map!

# Diversity analysis in *O. cumana*



M. Coque

| | Total |
|---|---|
| FRANCE | 3 |
| HUNGARY | 1 |
| ROMANIA | 1 |
| SPAIN | 5 |
| UKRAINE | 2 |
| **Total** | **12** |

**Exome capture from the 12 populations : 362285 SNPs**

**1536 SNPs selected to maximize the diversity of the whole set**

# A segregating population for the first genetic of *O. cumana*



B. Pérez-Vich   L. Velasco

IN12

INA

IN12          F1          INA

IN23 used for genome sequencing

# The first genetic map of
# *O. cumana*

X. Grand

1536 SNPs + 168 SSR x 91 F2 and parental lines have been genotyped

509 SNPs + 18 SSR were polymorphic and did not show any distortion of their segregation in the full population

Genetic map built using CarthaGène software (INRA) with a high stringency

**1479cM**
**28 linkage groups for the 19 chromosomes**

# The first genetic map of
# *O. cumana*

# Colinearity between the genetic map and the genome

Almost perfect colinearity between physical map and genetic map.

But strong variations between physical distances and genetic distances according to the regions.

**This genetic map enabling the anchoring of 95 contigs (593Mb) from the 256 contigs representing 90% of the genome assembly.**

**161 contigs remain unmapped to anchor 90% of the genome!**

INRA
SCIENCE & IMPACT

# ReSequencing the 2 parental lines of the segregating population

A. Calderon

The genome of the 2 parental lines used to produce the F2 segregating populations have been resequenced (HiSeq, Illumina)



40 912 accurate polymorphic SNPs between the two parental lines

# Improving the genetic map

L. Hu

Not enough contigs anchored to the genetic map (95 contigs)

40 912 polymorphic SNPs between the two parental lines

13511 SNPs located on 145 contigs from the 161 contigs that remain unmapped

278 SNPs will be genotyped in the full segregating population to anchor the 145 more contigs

SNP1                                    SNP2

contig

**All these data should anchor 90% of the genome assembly**

# Transcriptomic data to annotate the *genome*

P. Delavault

**60 RNASeq libraries sequenced** (data obtained on 26 April 2016)

Corresponding to **20 broomrape development stages** from seeds to flowering (3 replicates/stage)

# Use of an automatic annotation pipeline

E. Sallet

- EuGene Plant pipeline (egn-ep)

Total number of genes: **55726**

**Number of protein coding genes: 46447**
Mean gene length (bp): 3568.75
Per cent genes with introns : 63  **genes with 5' UTR: 82%    genes with 3' UTR: 83%**
- Exons      Mean number per gene: 3.63    Mean length (bp): 472.95    GC %:  40.94
- Introns   Mean number per gene: 2.63    Mean length (bp): 706.35    GC %    35.40
-CDS        Mean length (bp): 732.32           GC %  45.25

**Experimental expression evidence added!**

**Number of non protein coding genes      9279**

\* Tephra – repeat annotation - https://github.com/sestaton/tephra

        11 456  genes (8 256 protein coding + 3 200 ncRNA) are fully included in repeat_region annotated by tephra

# Assessment of the annotation quality using « BUSCO » (a set of conserved genes in plants)

| | Complete and unique | Complete but Duplicated | % complete | partial | Missing |
|---|---|---|---|---|---|
| Sunflower | 1153 | 196 | 93.7% | 21 | 70 |
| *O. cumana* | 999 | 47 | **72.7%** | 72 | 322 |

A lot of missing genes!
Due to the biology of *O. cumana*?

# All genome information are available in a Genome Browser

https://www.heliagene.org

S. Carrère

# All genome information are available in a Genome Browser



MTA for the access

# Biochimical informations from the genome sequence

L. Cottret

**Automatic process to build the metabolic database**



Schlapfer *et al.,* Plant Physiology. 2017.

# Metabolic database



Summary of *Orobanche cernua*, Subspecies cumana, version 1.0
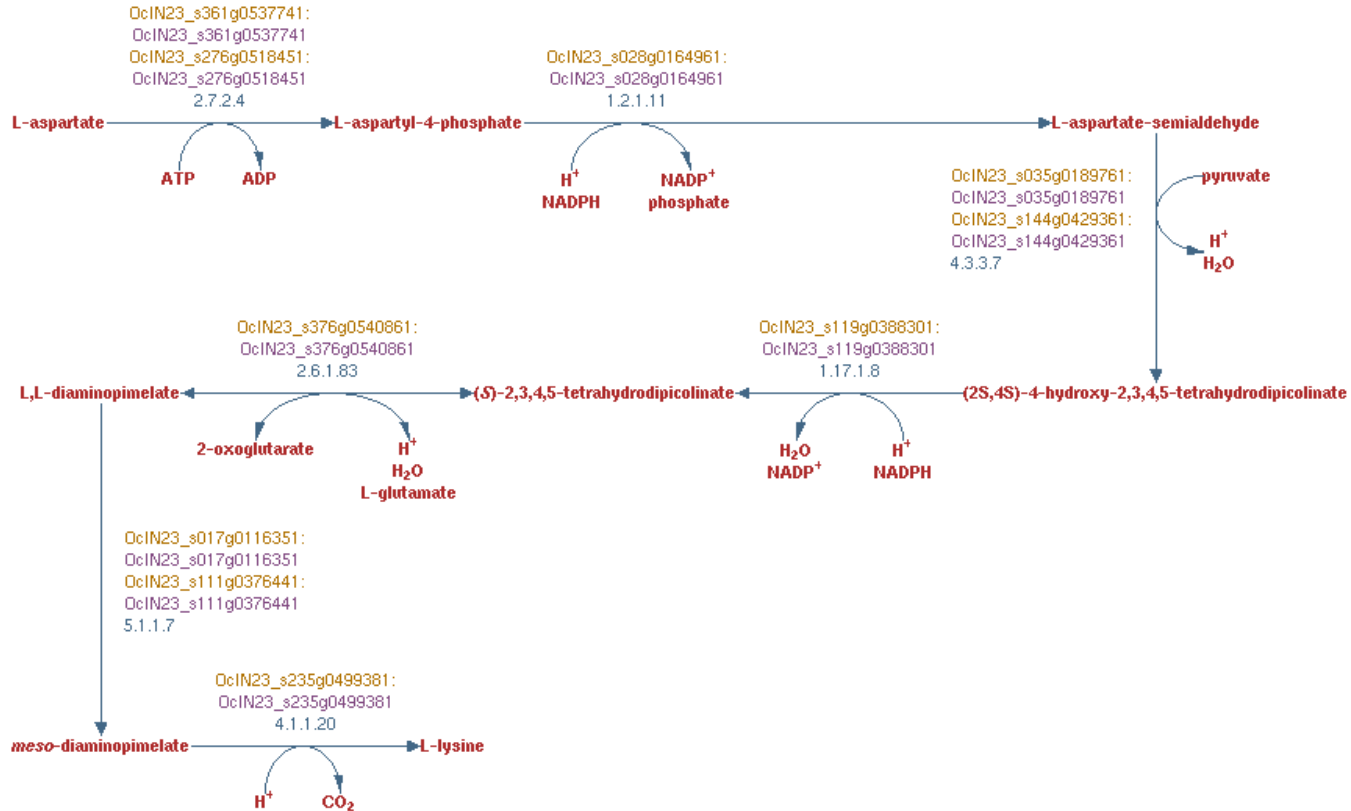
447 pathways

3458 reactions

1517 classified into pathways

Restricted access but collaborations are welcome

# Link between pathways, reactions, enzymes and genes



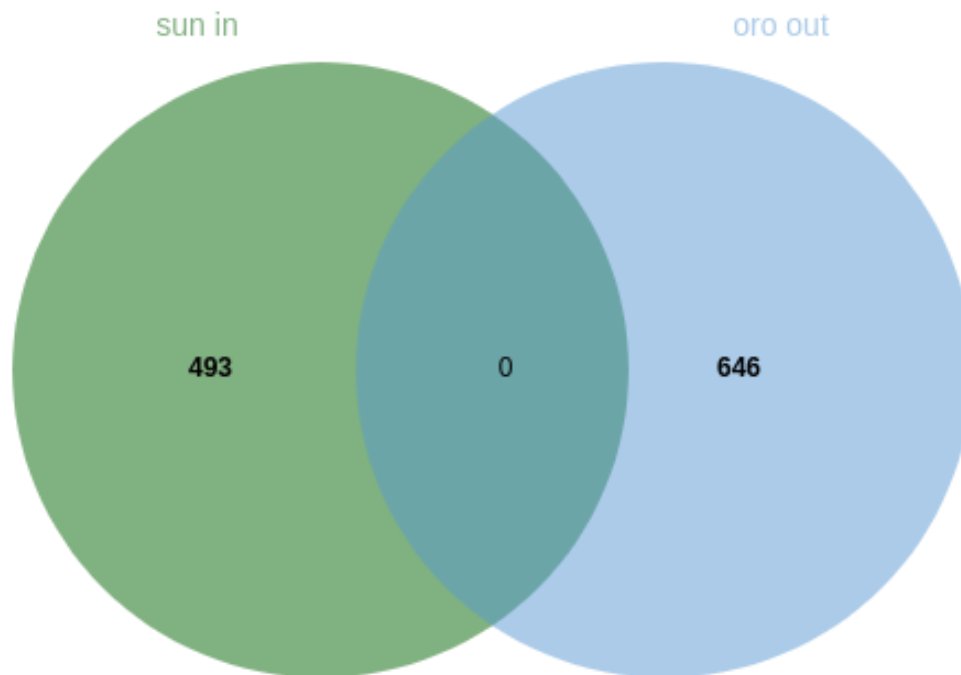Orobanche cernua cumana Pathway: L-lysine biosynthesis VI

# Comparison of sunflower and *O. cumana* pathways

# Intersection between Orobanche input metabolites and sunflower output metabolites

# Intersection between broomrape output metabolites and sunflower input metabolites

# Use of the annotated genome for RNASeq analysis

We used the annotated mRNA sequences from XRQ (sunflower) and IN23 (*O. cumana*) to map the RNASeq data

## Preliminary results, more analysis need to be performed to mak then more accurate!

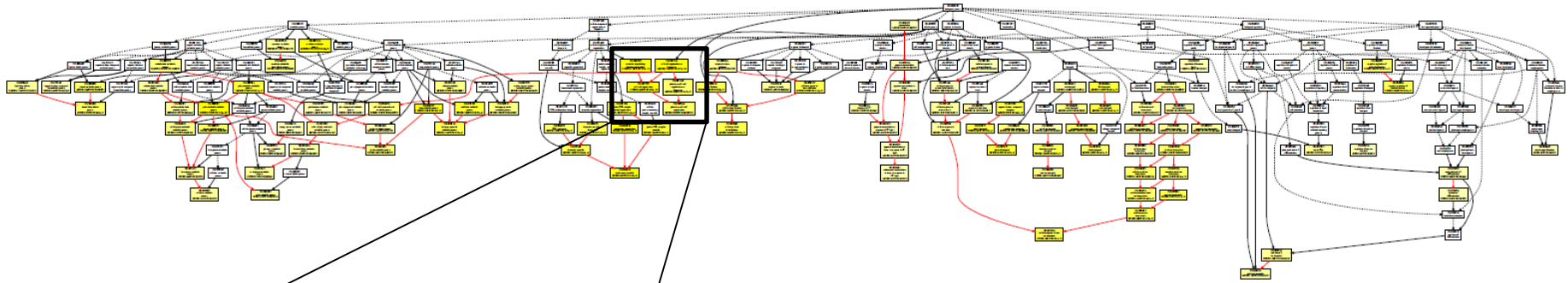# Sunflower DEG during incompatible attachment

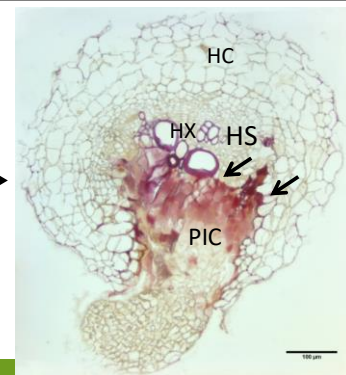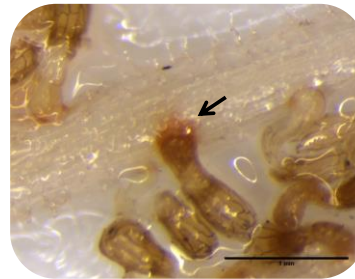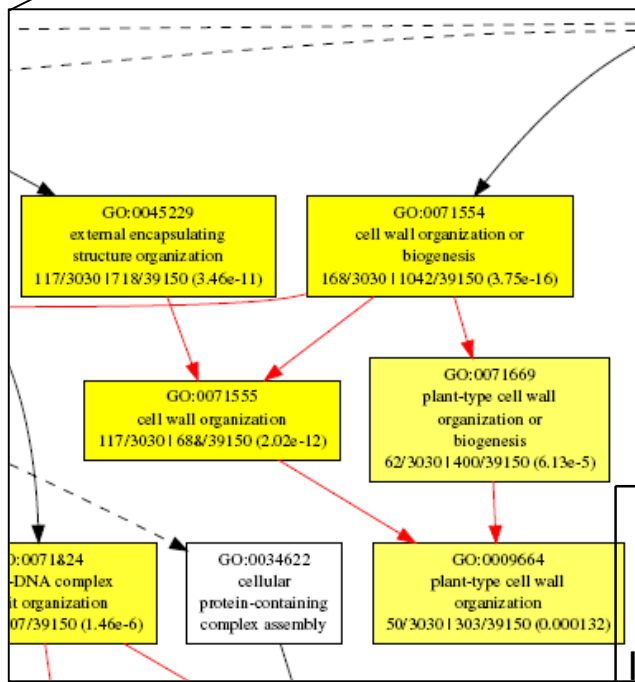Incompatible Attachments

LR1: 3444 DEG

Control

Gene Ontology terms enrichment analysis using GOEAST

A more complex response during incompatible attachment

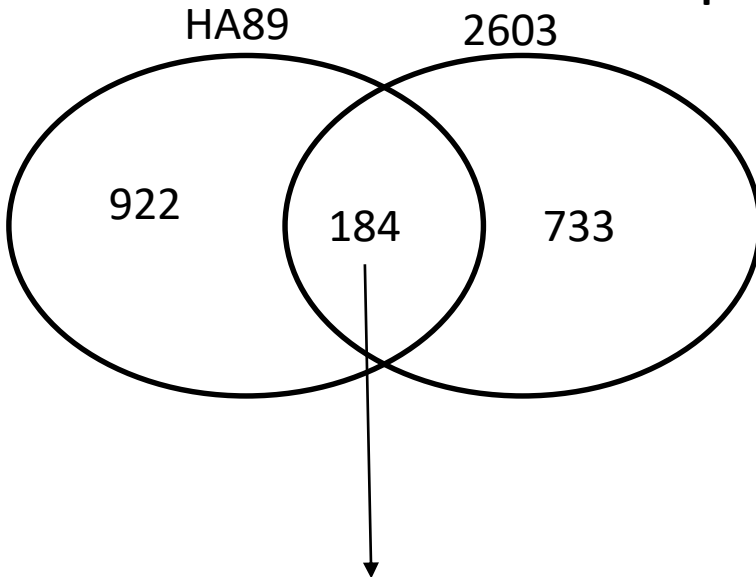DEG involved in cell wall biogenesis are enriched to prevent connection of the vascular system of the host

Barrier formation: **no connection** to the vascular system of the host

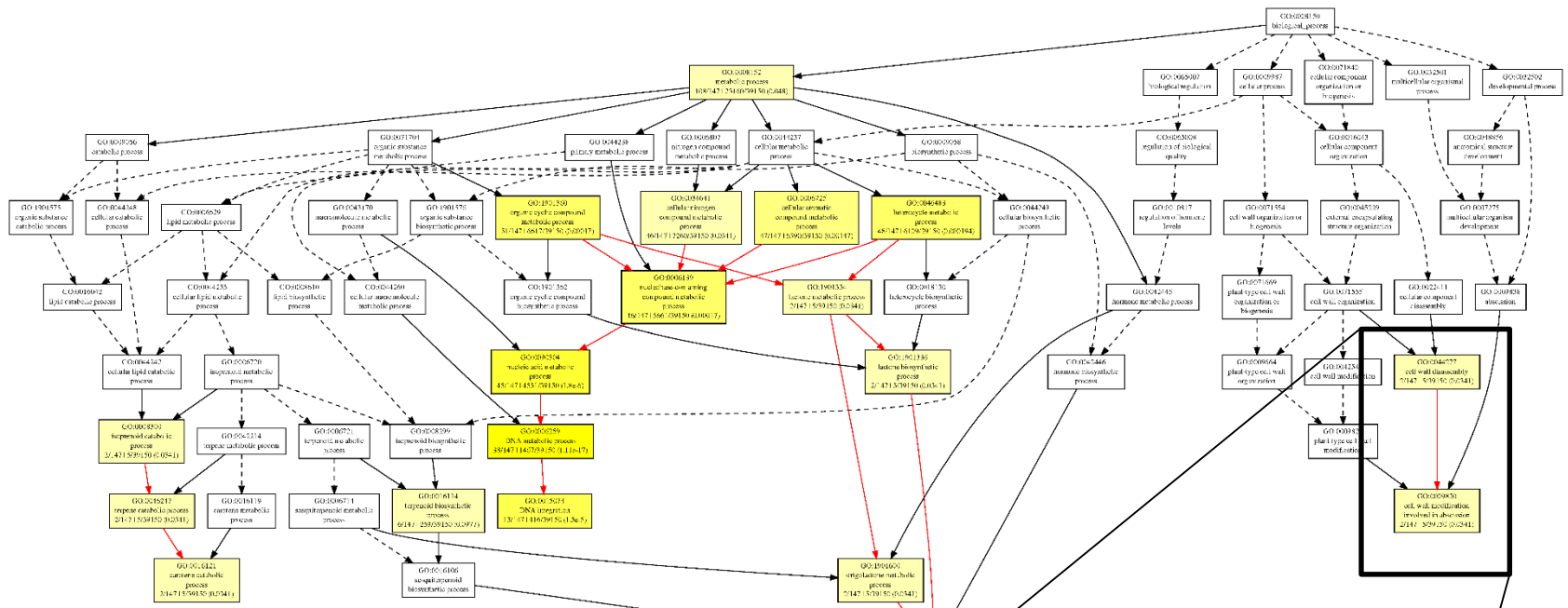# Sunflower DEG during compatible attachment

HA89: 1106 DEG
CA 2603: 917 DEG

HA89          2603

922      184      733

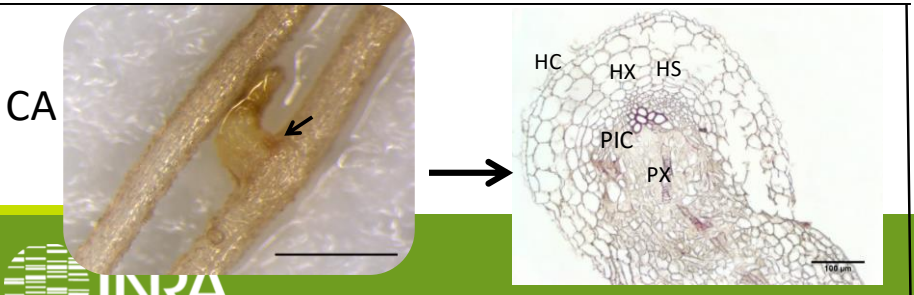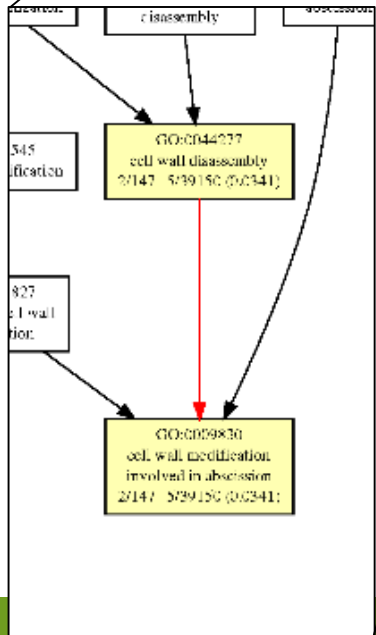**Compatible Attachments**



**Control**

**Gene Ontology terms enrichment analysis using GOEAST**

DEG involved in Cell wall disassembly are enriched.

The interaction (*O. cumana?*) seems to target cell wall degradation gene to enable *O. cumana* to connect to the vascular system of the host

CA

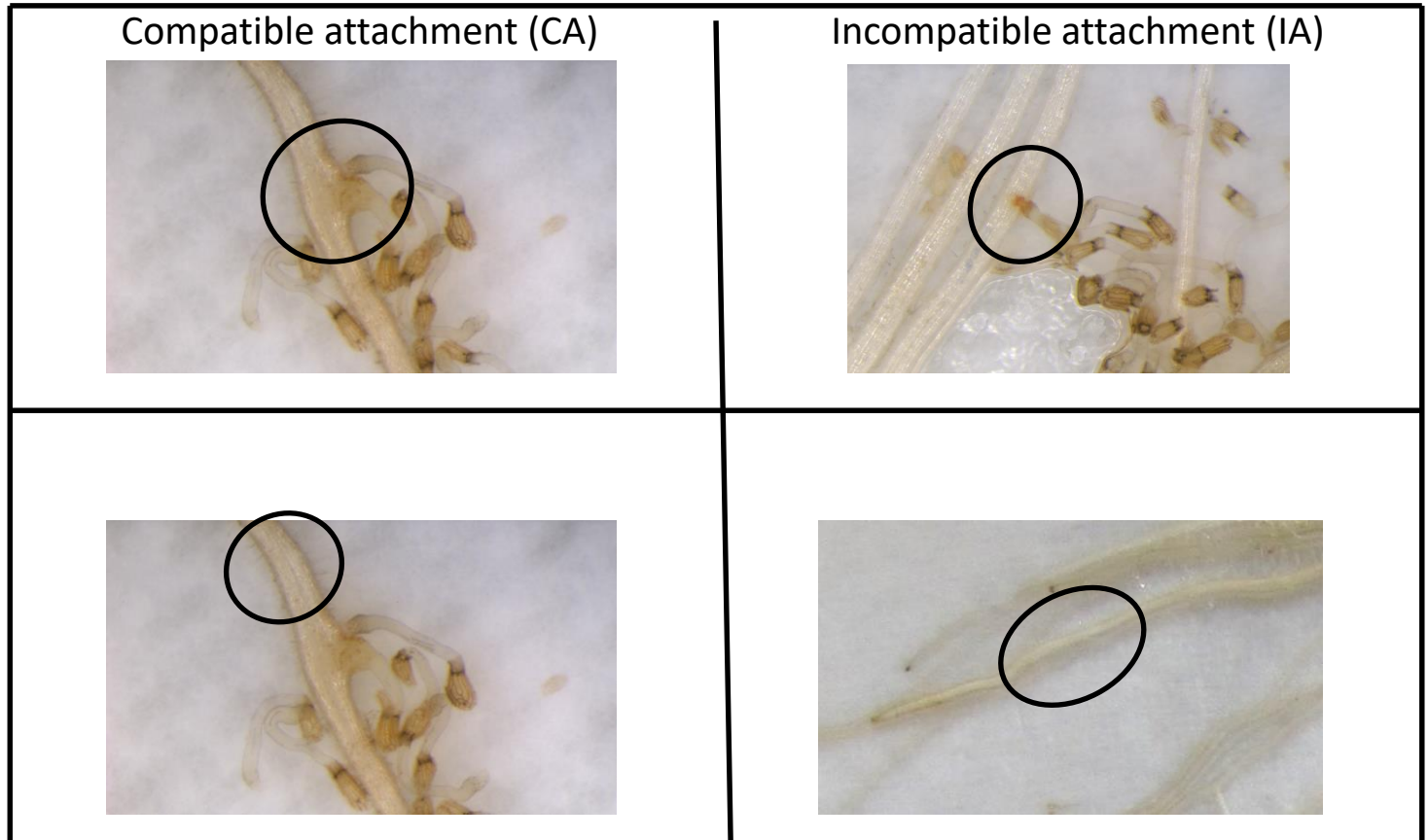**Connection** to the vascular system of the host

# *O. cumana* DEG

IA vs CA LR1: 15874

IA vs CA Other Resistant line : 742



| Compatible attachment (CA) | Incompatible attachment (IA) |

# *O. cumana* DEG

# Summary

A high quality genome sequence of *O. cumana* produced

[www.heliagene.org](www.heliagene.org)

Usefull for functional and genetic analysis

# Many Thanks to Collaborators

**IAS-CSIC:**
Álvaro Calderón González
Begoña Pérez-Vich
Leonardo Velasco

**LBPV, Nantes University :**
Philippe Delavault
Marc-Marie Lechat
Philippe Simier

**Biogemma:**
Clotilde Claudel
Marie Coque
Sébastien Faure
Xavier Grand
Nicolas Ribière

**CNRGV**:
Hélène Bergès
Stéphane Cauet
Céline Jéziorski
William Marande

**Get-Plage:**
Cécile Donadieu
Olivier Bouchez
Maarten Pirson

**Terre Inovia:**
Christophe Jestin

**LIPM:**

Julia Bazerque
Nicolas Blanchet
Marie-Claude Boniface
Fanny Bonnafous
Sébastien Carrère
Olivier Catrice
Mireille Chabaud
Ludovic Cottret
Alexandra Dühnen
Pauline Duriez
Louise Gody
Florie Gosseau
Jérôme Gouzy
Luyang Hu

Nicolas Langlade
Marion Larroque
Ludovic Legrand
Johann Louarn
Anne-Sophie Lubrano
Brigitte Mangin
Gwenola Marage
Prune Pegot-Espagnet
Nicolas Pouilly
David Rengel
Erika Sallet
Camille Tapy