

# LA REVOLUCIÓN DEL BIG DATA

¿Cómo podríamos usar estos datos para pronosticar y mejorar nuestras investigaciones?



UNIVERSIDAD  
DE ANTIOQUIA

Héctor Alejandro Rodríguez Cabal PhD

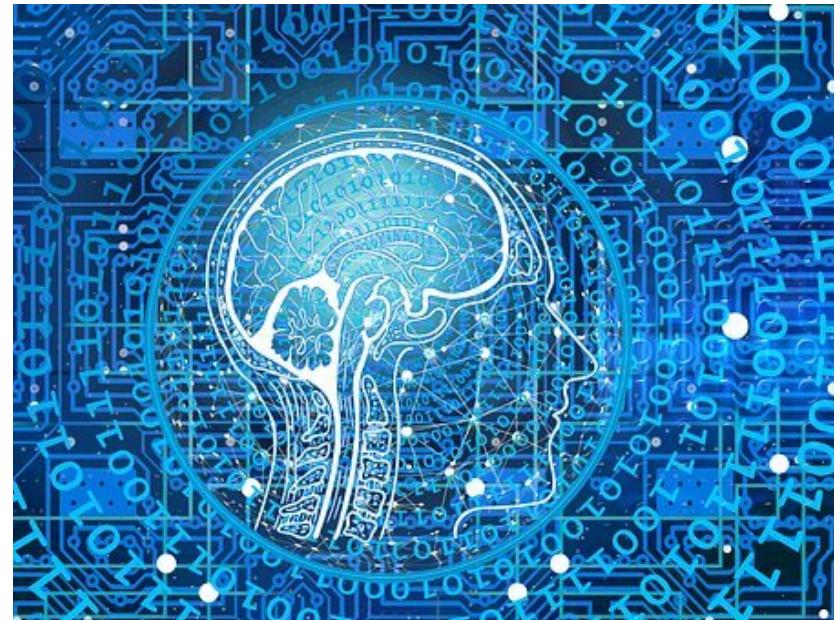
Docente FCEN

Coordinador del Grupo Agrobiotecnología



# Contenido

- Concepto de Big Data.
- Diferencias entre Big Data y Data Science.
- Big Data en el mundo.
- Big Data en Colombia.
- Big Data para la ciencia.
- Perspectivas.



# — ¿Qué es Big Data?



UNIVERSIDAD  
DE ANTIOQUIA



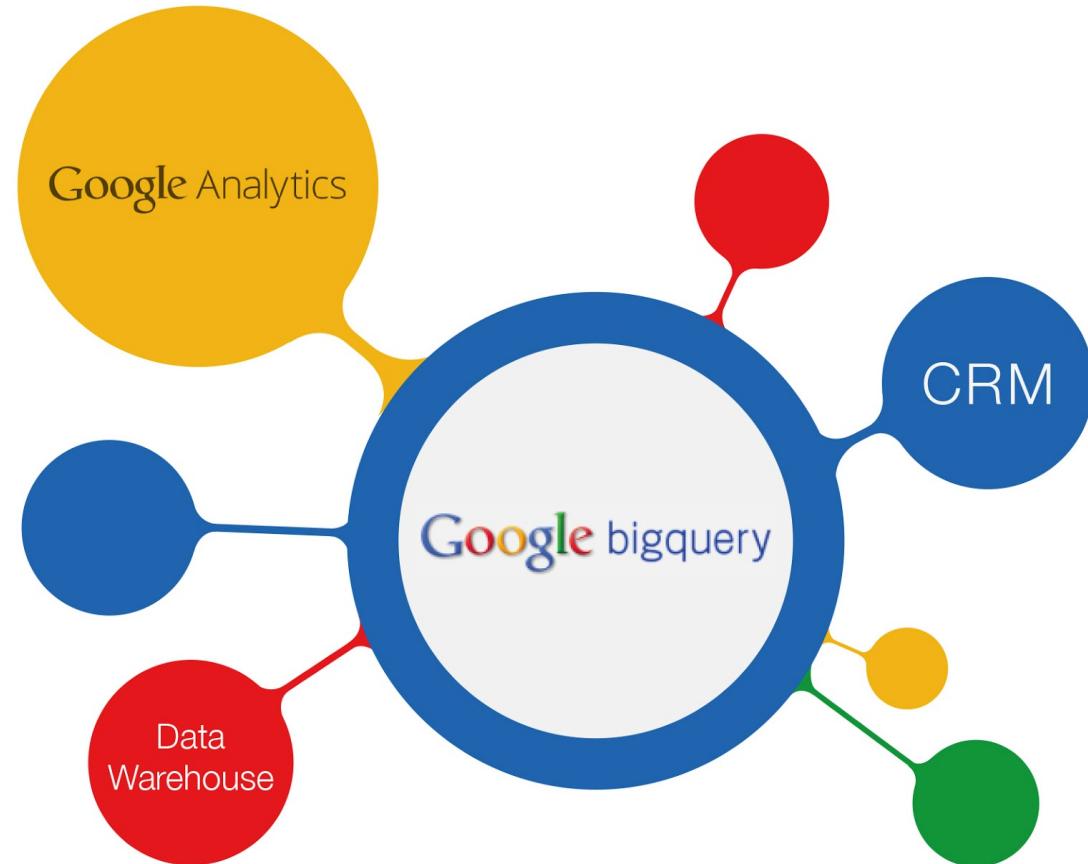
- Web logs.
- La identificación por radiofrecuencia (RFID) Los sensores incorporados en dispositivos.
- La maquinaria.
- Los vehículos.
- Las búsquedas en Internet.
- Las redes sociales como Facebook.
- Computadoras portátiles.
- Teléfonos inteligentes y otros teléfonos móviles.
- Dispositivos GPS.
- Registros de centros de llamadas.



# — ¿Por qué el Big Data es tan importante?

- El análisis de Big Data ayuda a las organizaciones a aprovechar sus datos y utilizarlos para identificar nuevas oportunidades.
  - **Reducción de coste.**
  - **Más rápido, mejor toma de decisiones.**
  - **Nuevos productos y servicios.**

# El as en la manga de big data Google

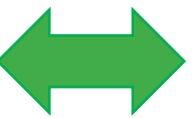


Si Google encuentra que de repente hay una cantidad creciente de usuarios que empieza a mirar determinado tipo de contenido concreto, es posible que esté detectando una nueva tendencia antes de que se convierta en popular.

"Con el big data somos capaces de decirle a un restaurante dónde es mejor abrir para tener éxito", Susana Voces de Deliveroo



# — Diferencias entre Big Data y Data Science.

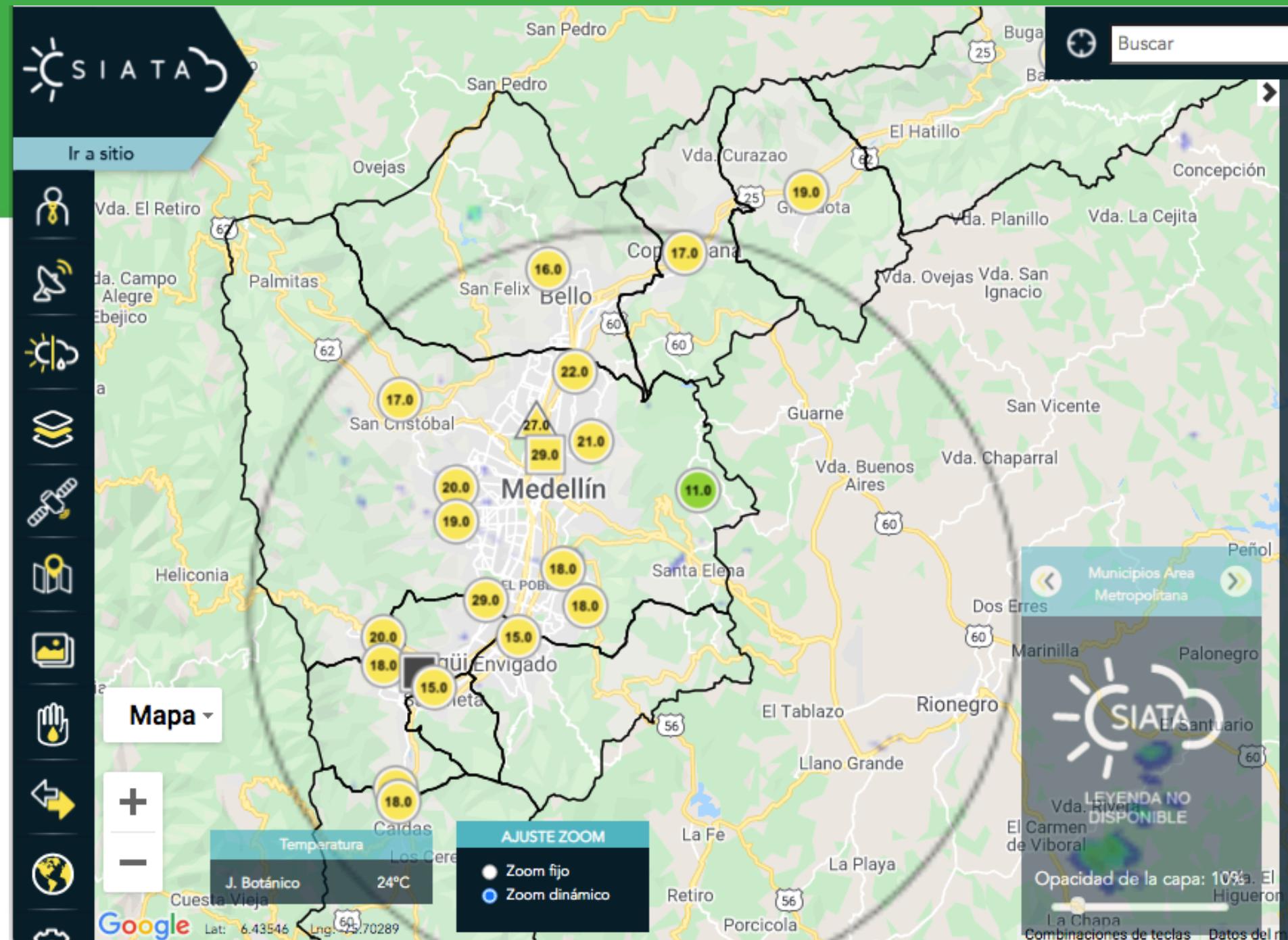


Big Data se distinguen por variedad,  
volumen y velocidad

El Data Science proporciona  
métodos o técnicas para analizarlos.

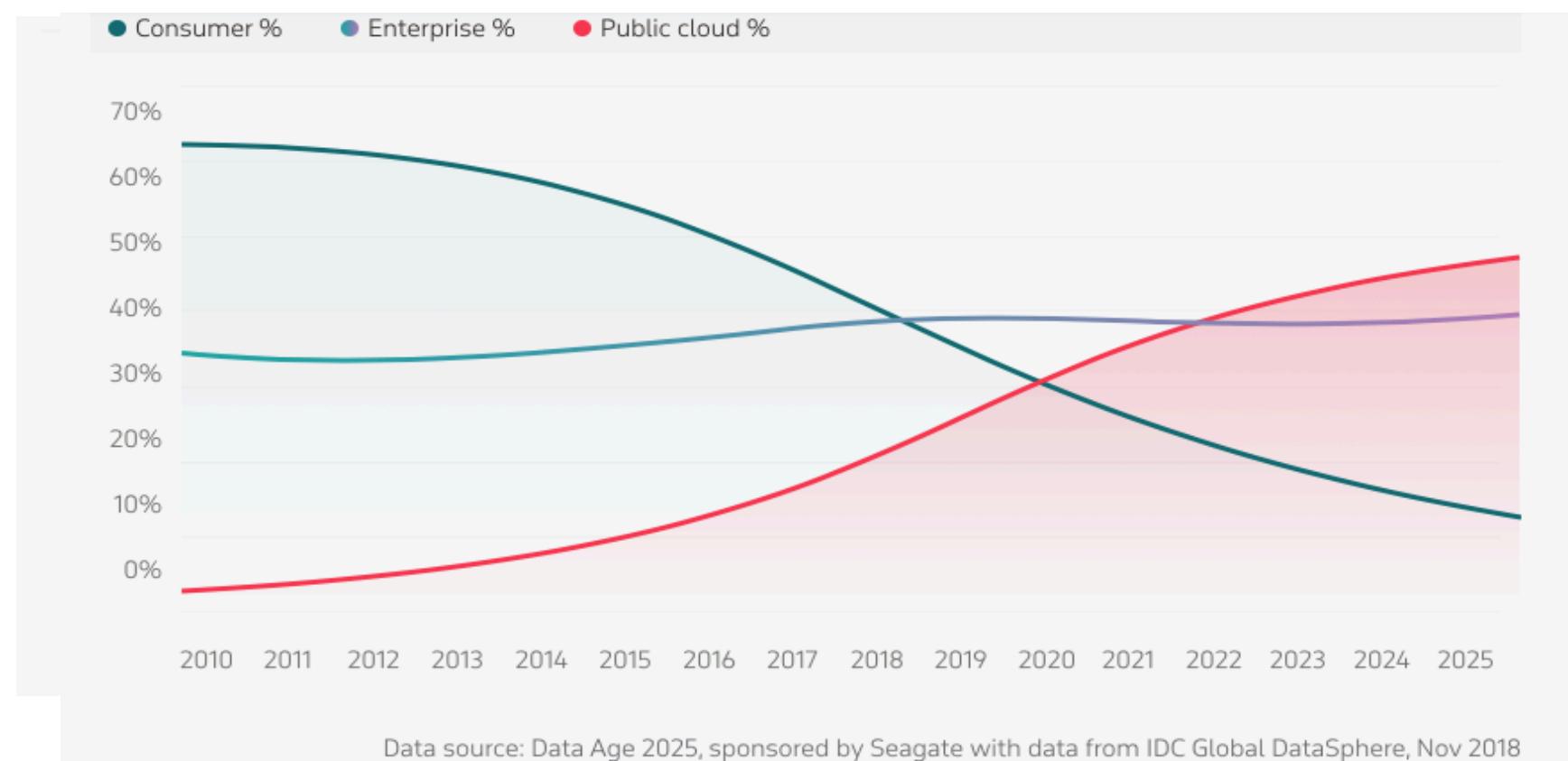
# Pirámide informacional





# Big data en el mundo

## 1. Los volúmenes de datos seguirán aumentando y migrando a la nube



# Big data en el mundo

## 2. El aprendizaje automático (machine learning) seguirá cambiando el panorama.

**“** Machine learning is becoming more sophisticated with every passing year. We are yet to see its full potential—beyond self-driving cars, fraud detection devices, or retail trends analyses.



**Wei Li**

Vice President and General Manager at Intel

**“** What fascinates me is combining big data with machine learning and especially natural language processing, where computers do the analysis by themselves to find things like new disease patterns.



**Bernard Marr**

Author, Big Data: Using SMART Big Data, Analytics and Metrics to Make Better Decisions and Improve Performance

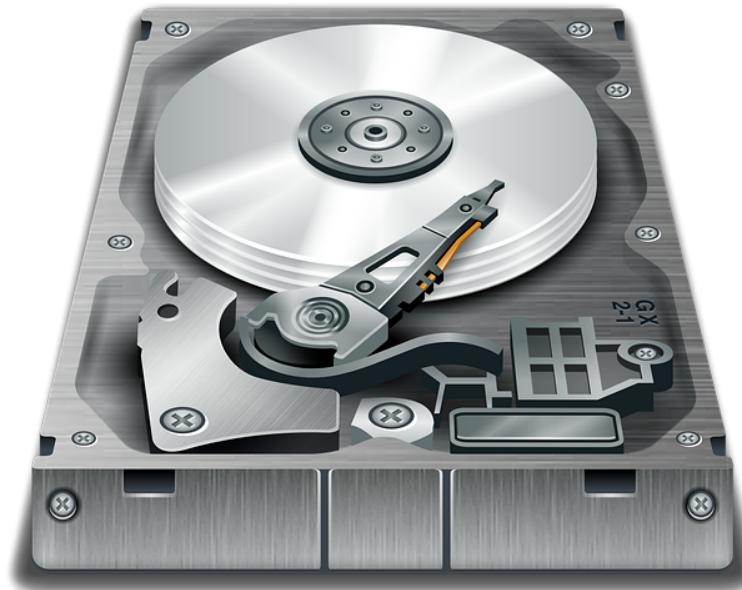
# Big data en el mundo

## 3. La privacidad seguirá siendo un tema candente

- **Brecha de habilidades de seguridad**, causada por la falta de oportunidades de educación y capacitación. Esta brecha crece constantemente y alcanzó los 3,5 millones de puestos de seguridad cibernética vacantes para 2021, según Cybercrime Magazine.
- **Evolución de los ciberataques**. Las amenazas utilizadas por los piratas informáticos están evolucionando y se vuelven más complejas cada día.
- **Adhesión irregular a las normas de seguridad**. Aunque los gobiernos están tomando medidas para estandarizar las regulaciones de protección de datos, siendo GDPR el ejemplo, la mayoría de las organizaciones aún ignoran los estándares de seguridad de datos.

# Big data en el mundo

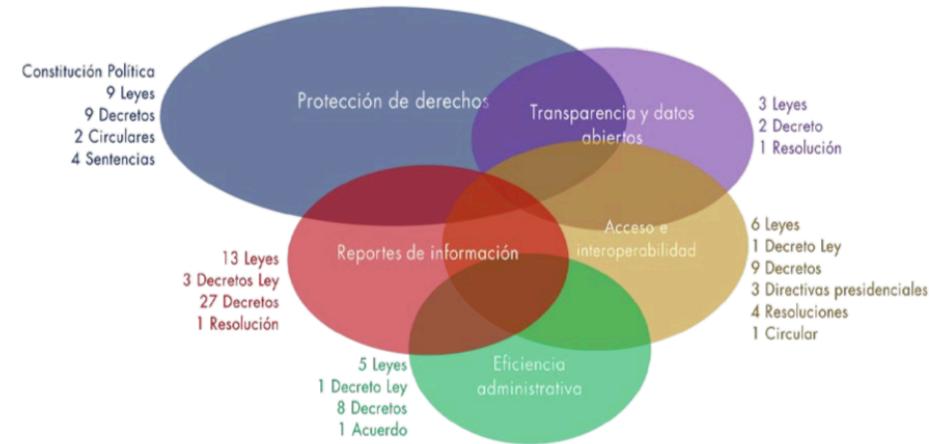
## 4. Los datos rápidos y procesables pasarán a primer plano



# Big Data en Colombia.

El Gobierno de Colombia ha realizado importantes esfuerzos para consolidar un marco de política que promueva el aprovechamiento de datos para la generación de valor social y económico en el país. La Política Nacional de Explotación de Datos y Big data - CONPES 3920 (2018) y la Política Nacional de Transformación Digital e Inteligencia Artificial CONPES 3975 (2019).

**Figura 1. Marco jurídico aplicable a los datos**



Fuente: DNP (2018).

# Big Data en Colombia.



**MinTIC expide el Plan Nacional de Infraestructura de Datos, que impulsará la transformación digital del Estado**

Nacional

## Plan Nacional de Infraestructura de Datos

Documento técnico y hoja de ruta

MinTIC, DNP, DAPRE  
Diciembre 2021

**PND**

# Enfoque de datos como infraestructura

- Permite abordar la infraestructura de datos, desde un enfoque de bien público, transversal a diversas actividades del Estado, necesario para crear otros bienes y servicios en la economía.



# Propósitos del Plan Nacional de Infraestructura de Datos

- Aumentar la reutilización de los datos que integran la infraestructura de datos
- Consolidar un sector público basado en datos
- Consolidar espacios de intercambio de datos que impulsen la innovación en el país
- Promover el desarrollo e integración de tecnologías emergentes
- Posicionar al modelo de gobernanza de los datos para la consolidación de una economía digital guiada por datos
- Posicionar al país como un referente en el uso de los datos para el desarrollo de la economía digital
- Construir un entorno de confianza pública para el aprovechamiento y protección de los datos

# Big Data en Colombia.

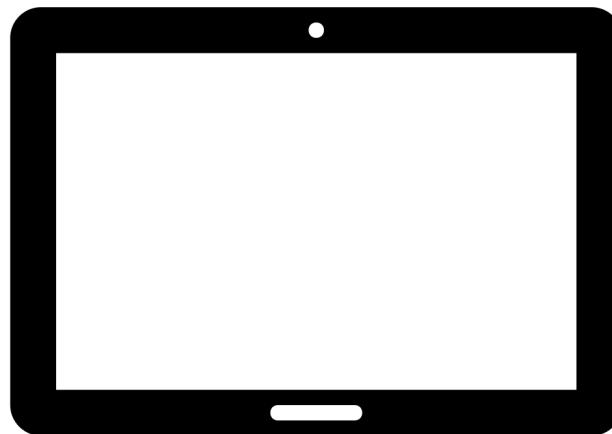


**Misión de Expertos de Inteligencia Artificial**

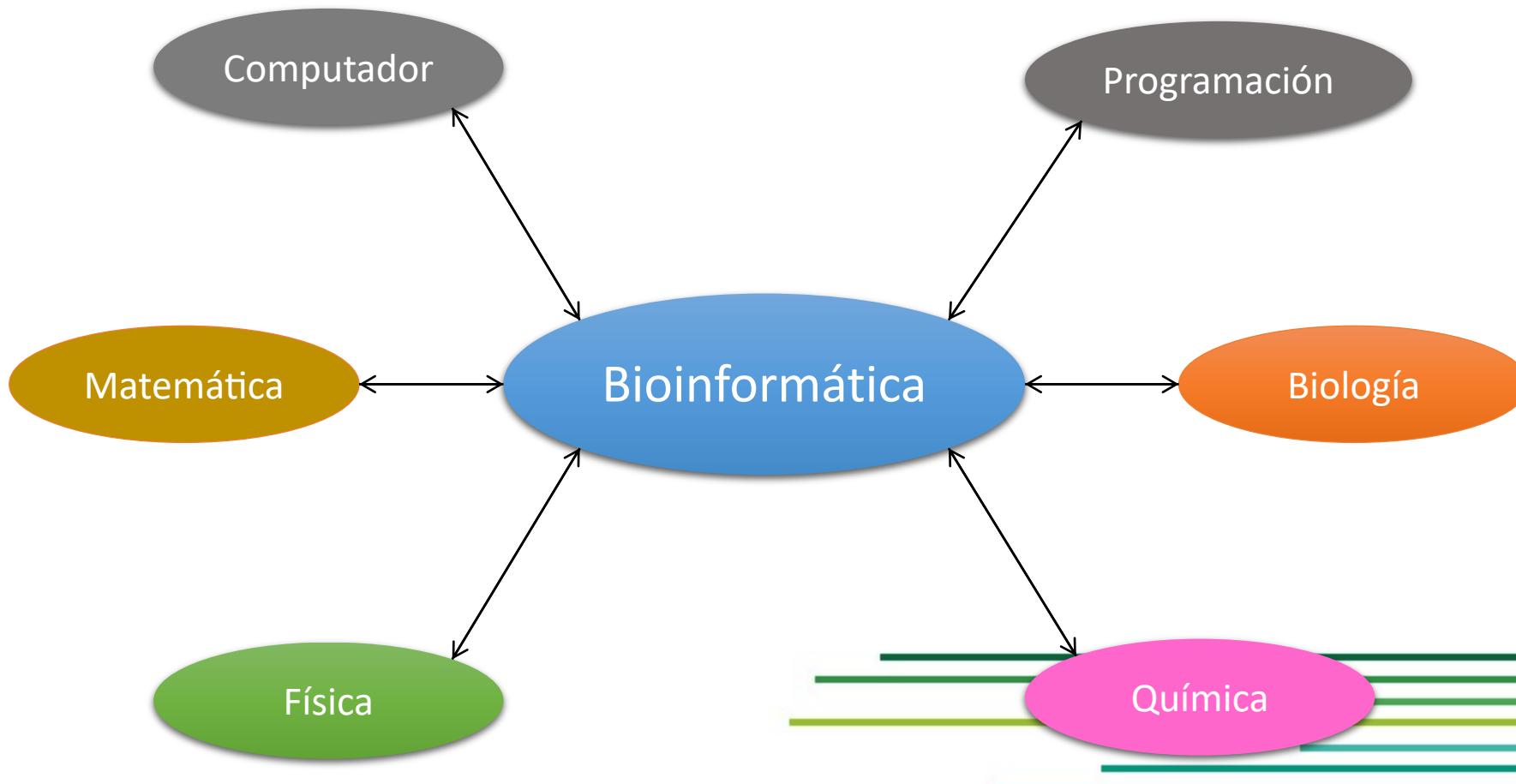


La Misión de Expertos en IA de Colombia es uno de los primeros grupos de expertos de IA en la región, y uno de los primeros en enfocarse en la generación de medidas para el desarrollo de políticas de educación y empleo para la cuarta revolución industrial.

# — ¿Podría Big Data impactar la educación?



# Big data para la Ciencia



# Big data para la Ciencia

## En proyectos a gran escala

- Se busca disminuir al máximo la intervención humana.
- Sistematización del proceso.
- Diseño de pipelines y nuevos scripts

# Big data para la Ciencia



National Library of Medicine

# THIRTY YEARS OF GROWTH: GenBank Sequences & NCBI Web Users

Sequences (Millions)

200

150

100

50

0

7.0

6.0

5.0

4.0

3.0

2.0

1.0

0.0

GenBank Sequences  
NCBI Web Users

1989 1992 1994 1996 1998 2000 2002 2004 2006 2008 2010 2012 2014 2016 2018 2019

PubMed Central

NIH Public Access

PubChem

Genome Reference  
Consortium

Genome-Wide  
Association Studies

Human Genome

PubMed

OMIM

Genomes

GenBank  
at NCBI

Entrez

BLAST

1000 Genomes

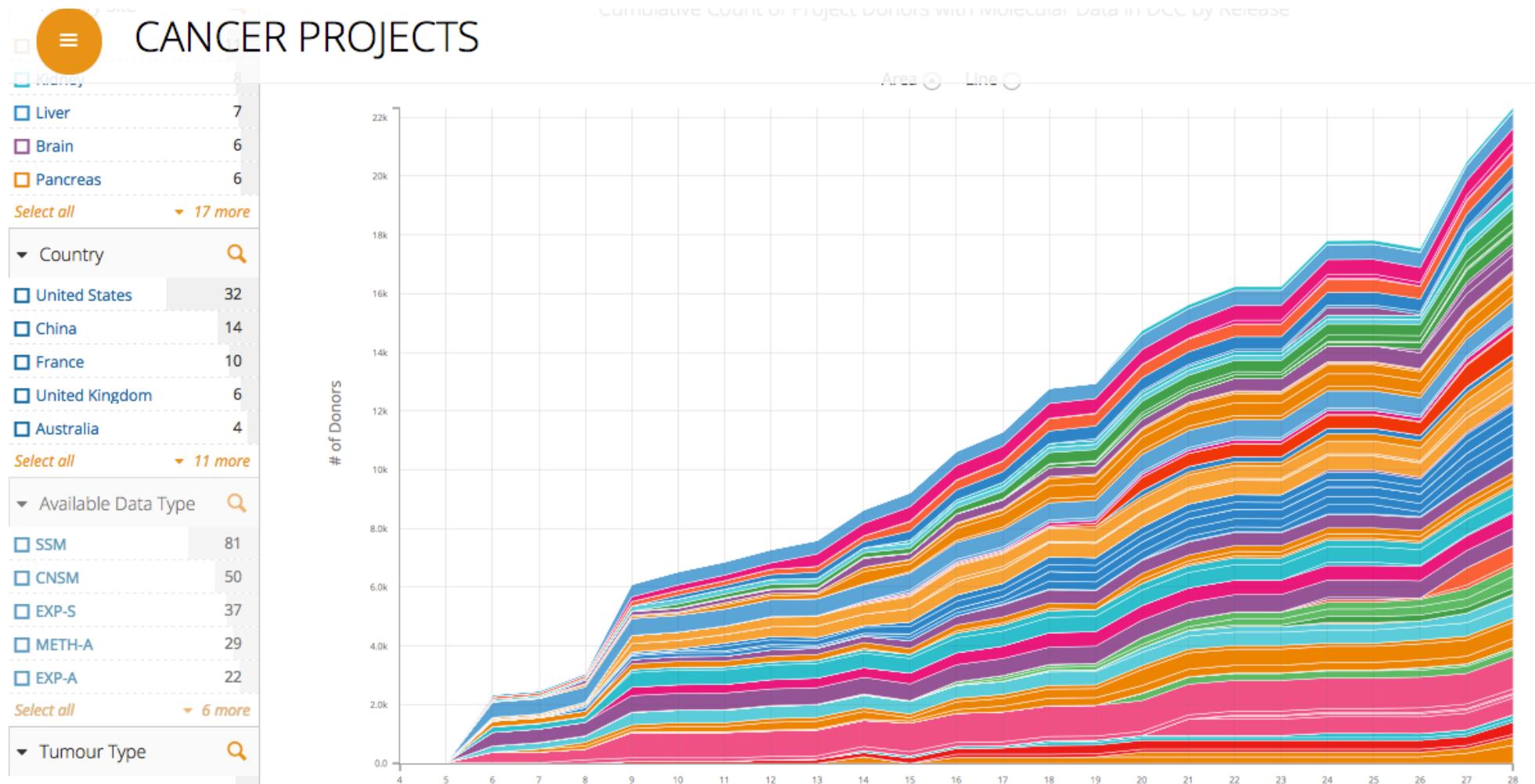
Genetic Testing  
Registry & ClinVar

MedGen  
PubReader

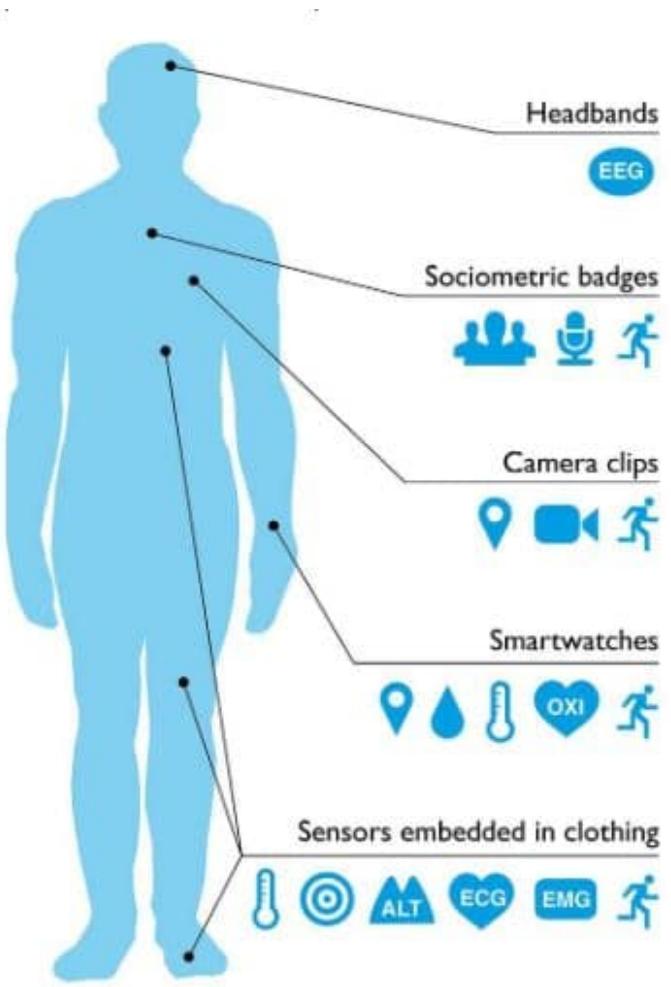
Food Pathogens  
Project

Antimicrobial  
Resistance

# Big data para la Ciencia



# Big data para la Ciencia



- Accelerometer:** Running figure icon
- Altimeter:** ALT icon
- Digital camera:** Video camera icon
- ECG:** Heart icon
- EMG:** Muscle icon
- EEG:** Brain icon
- Electrodermograph:** Water droplet icon
- Location GPS:** Location pin icon
- Microphone:** Microphone icon
- Oximeter:** Oximeter (OXI) icon
- Bluetooth proximity:** Two people icon
- Pressure:** Target icon
- Thermometer:** Thermometer icon

# Big data para la Ciencia

ADVERTISEMENT FEATURE Advertiser retains sole responsibility for the content of this article

## Startups and big data set to fuel genomic medicine in Korea

Genomic medicine is taking off in South Korea in a big way, thanks in large part to innovative startups and systems for collecting and analyzing big data.

---

Produced by

**nature** research  
custom media



SEOUL  
NATIONAL  
UNIVERSITY





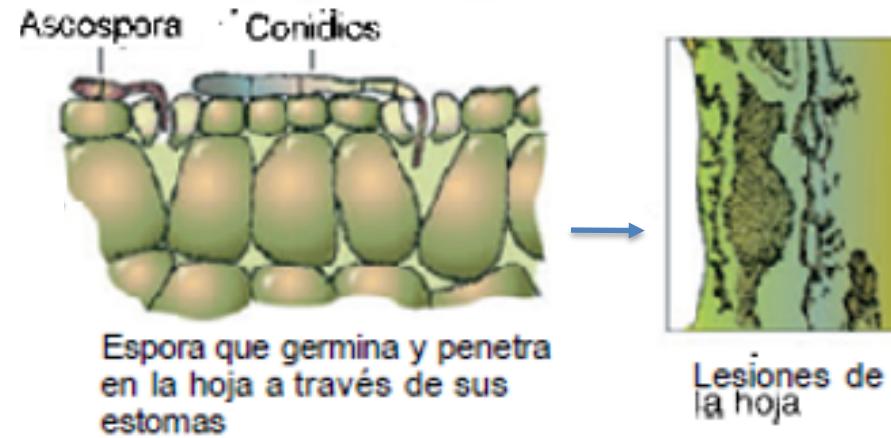
Variedad Susceptible  
Williams



Variedad resistente  
Calcutta 4



# Apoplastic fungus

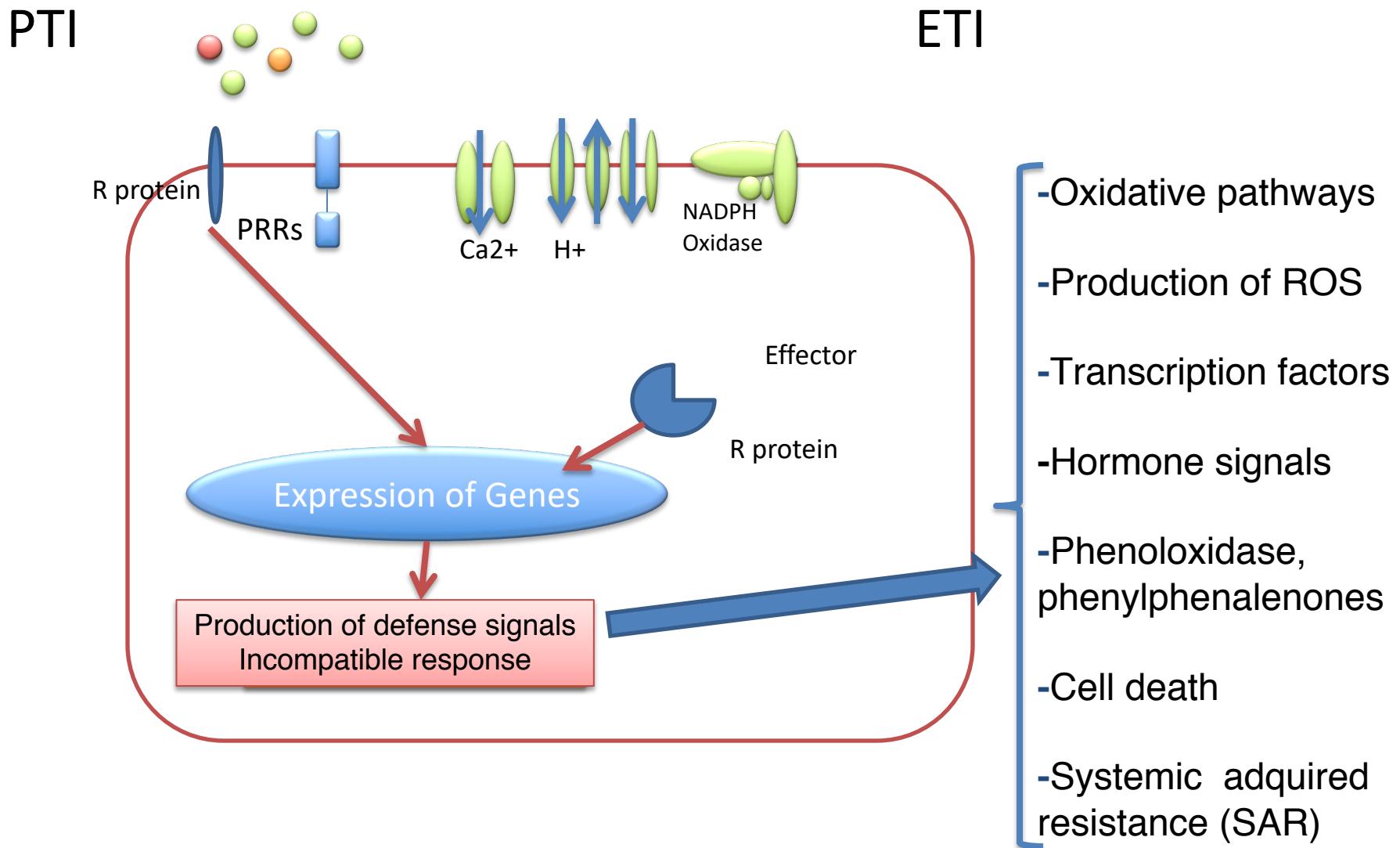


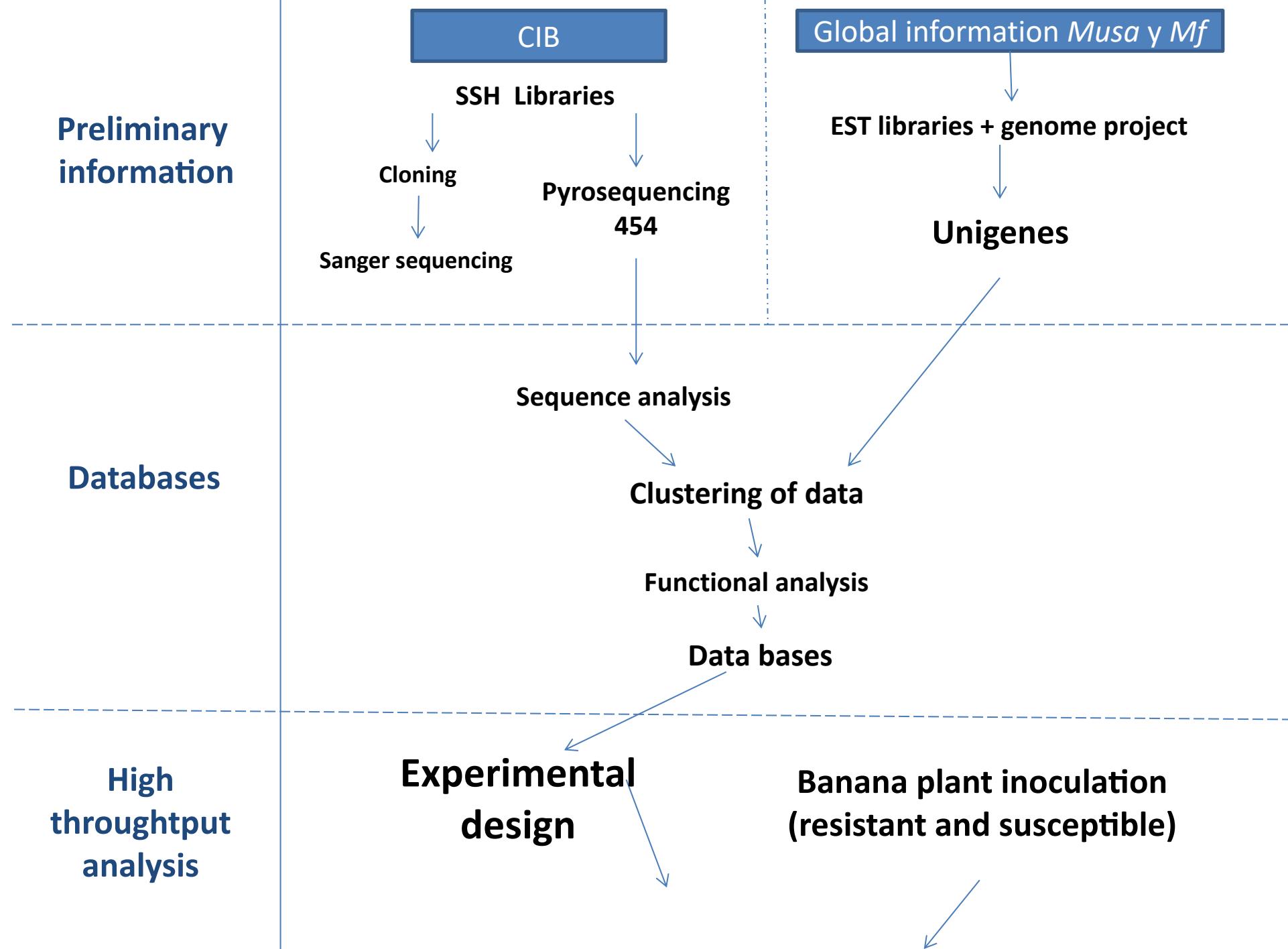
Biotrophic → Necrotrophic

10 to 15 days

Long biotrophic phase

# How Plants Defend Themselves?





High  
throughput  
analysis

Experimental  
design

Banana plant inoculation  
(resistant and susceptible)

Transcriptomic and  
Metabolomic

Data analysis

Analysis

Differential  
expression

12 hpi	18 hpi	24 hpi	48 hpi	72 hpi	144 hpi	9 dpi	15 dpi	30 dpi
control1	control1	control1	control1	control1	control1	control1	control1	control1
control2		control2		control2	control2		control2	
Control 3		Control 3		Control 3	Control 3		Control 3	
Replica 1	Replica 1	Replica 1	Replica 1	Replica 1	Replica 1	Replica 1	Replica 1	Replica 1
Replica 2		Replica 2		Replica 2	Replica 2		Replica 2	
Replica 3		Replica 3		Replica 3	Replica 3		Replica 3	
Contact	Contact	Penetration	Advance	Advance	HR (resistant)	Stage 0-1 susceptible	Stage 2-3 susceptible	Stage 4-5 susceptible

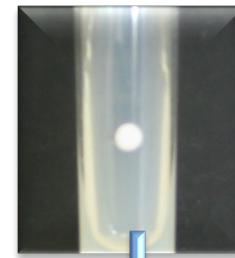
# Inoculation



Williams  
(Susceptible)



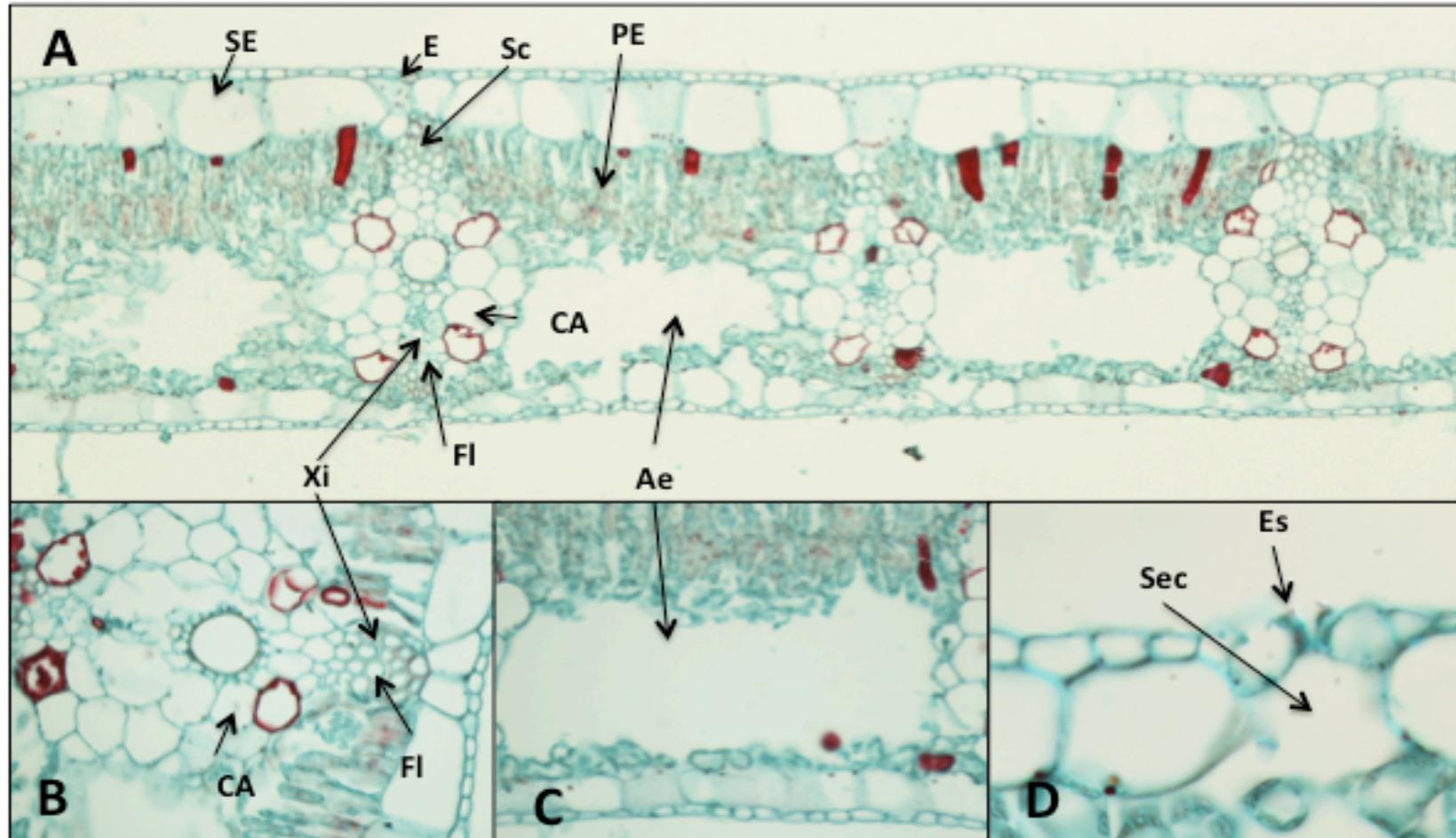
Calcutta  
(Resistente)



EST	E0	E1/2	E2	E3	E3	E3	E3/4
T	0 - 72 hpi	144 hpi (6 dpi)	216 hpi (9 dpi)	360 hpi (15 dpi)	552 hpi (23 dpi)	720 hpi (30 dpi)	960 hpi (40 dpi)
CAL							
EST	E0	E1	E2	E3	E4	E5	E6
T	0 - 216 hpi	240 hpi (10 dpi)	312 hpi (13 dpi)	360 hpi (15 dpi)	552 hpi (23 dpi)	720 hpi (30 dpi)	960 hpi (40 dpi)
Will							

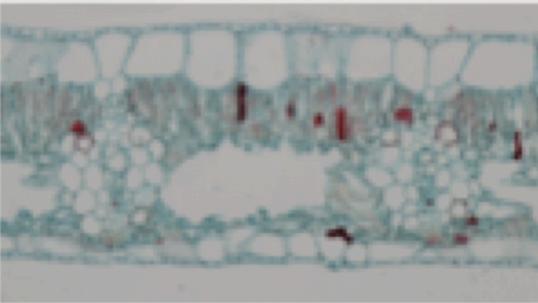
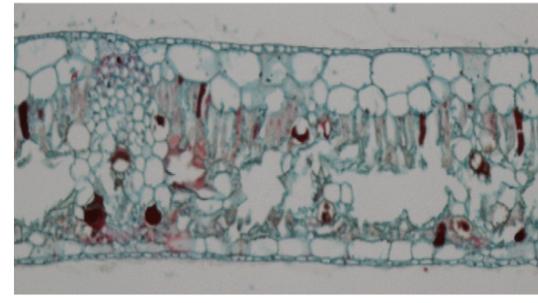
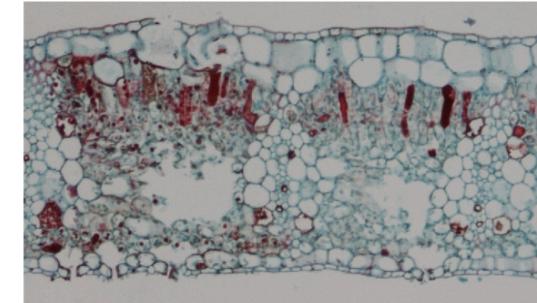
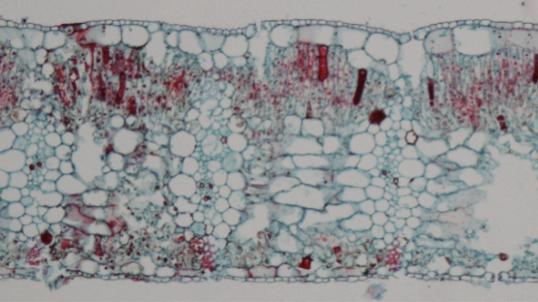
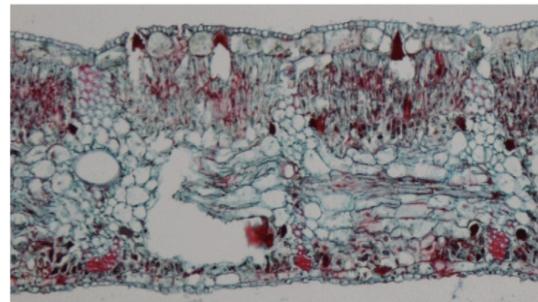
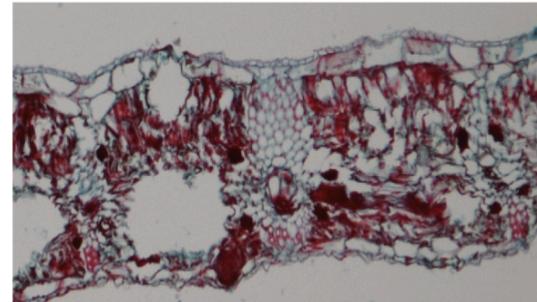
Est= Estadio; T= tiempo

# Caracterización morfoanatómica de una hoja sana de Williams



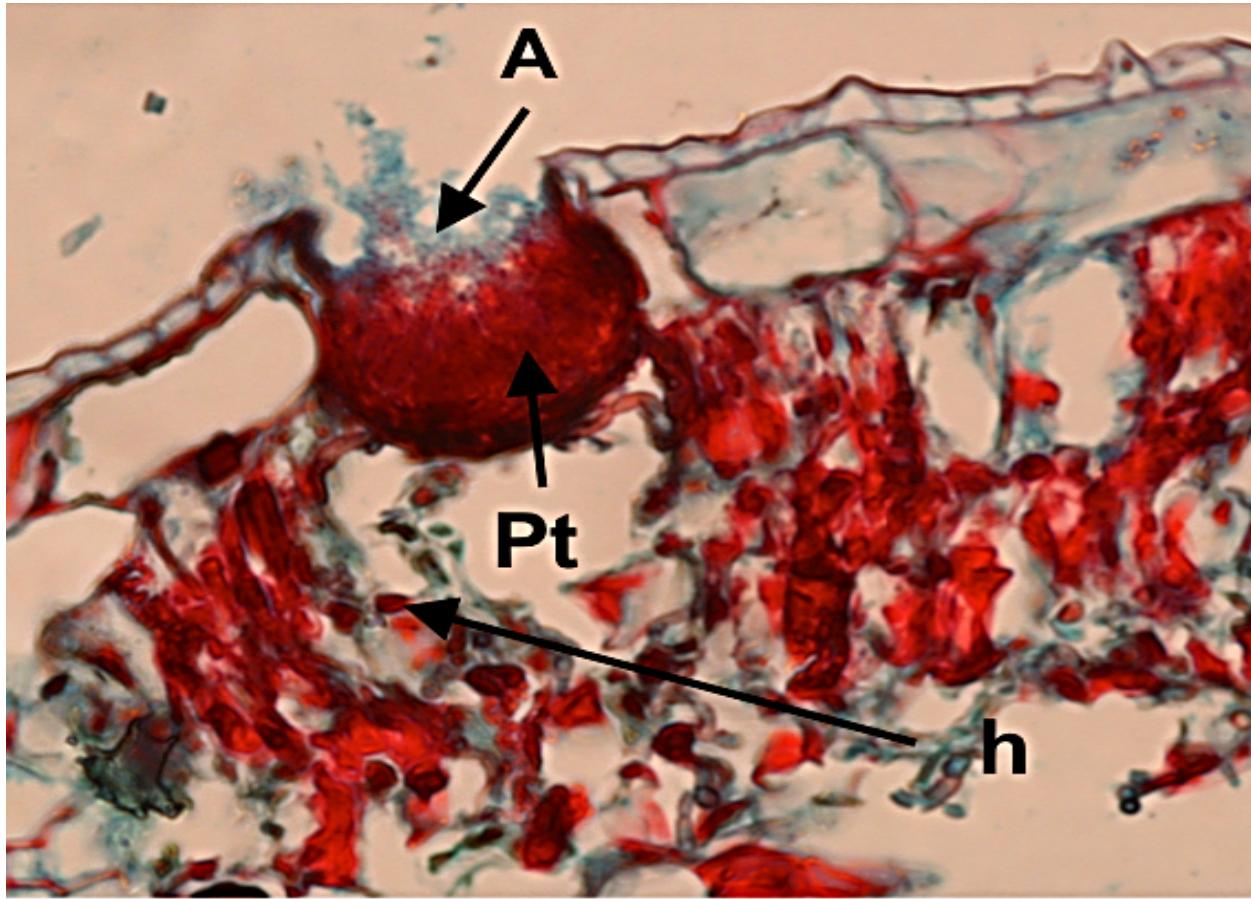
SE= capa Sub-Epidermal; PE= células del Parénquima de Empalizada; E= Epidermis; Sc= células del Esclerénquima; CA=Células Acompañantes; Sec= cámara Sub-estomatal; Ae= Espacios aeríferos; Fl= Floema; Es= Estoma; Xi= Vasos del Xilema.

# Caracterización morfoanatómica de hojas de la variedad Williams después de la inoculación con el hongo.

EST	E1	E2	E3
T	240 hpi (10 dpi)	312 hpi (13 dpi)	360 hpi (15 dpi)
Will			
EST	E4	E5	E6
T	552 hpi (23 dpi)	720 hpi (30 dpi)	960 hpi (40 dpi)
Will			

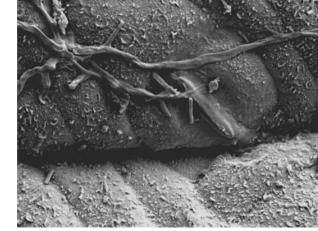
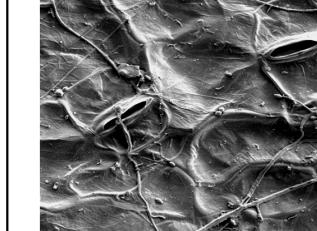
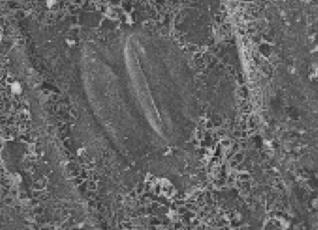
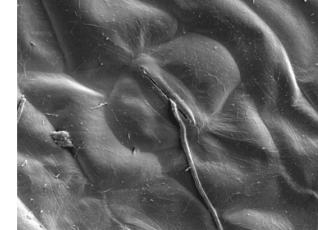
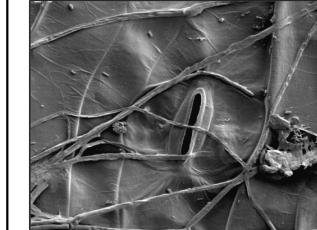
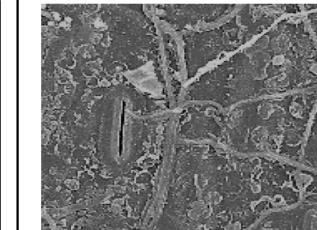
EST= estadío; T= Tiempo

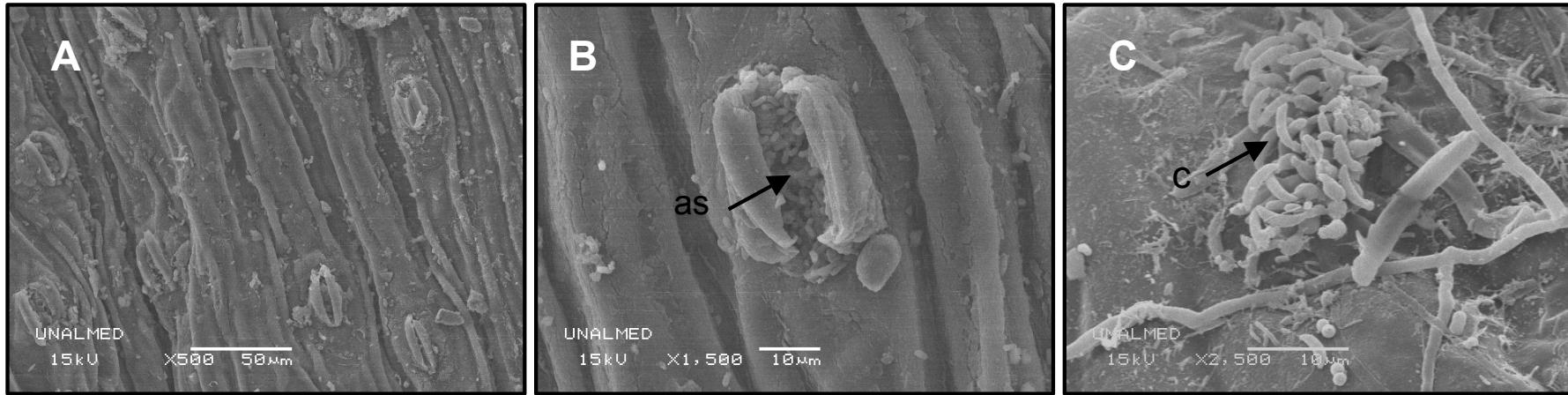
## Pseudotecio liberando ascosporas de *M. fijiensis*



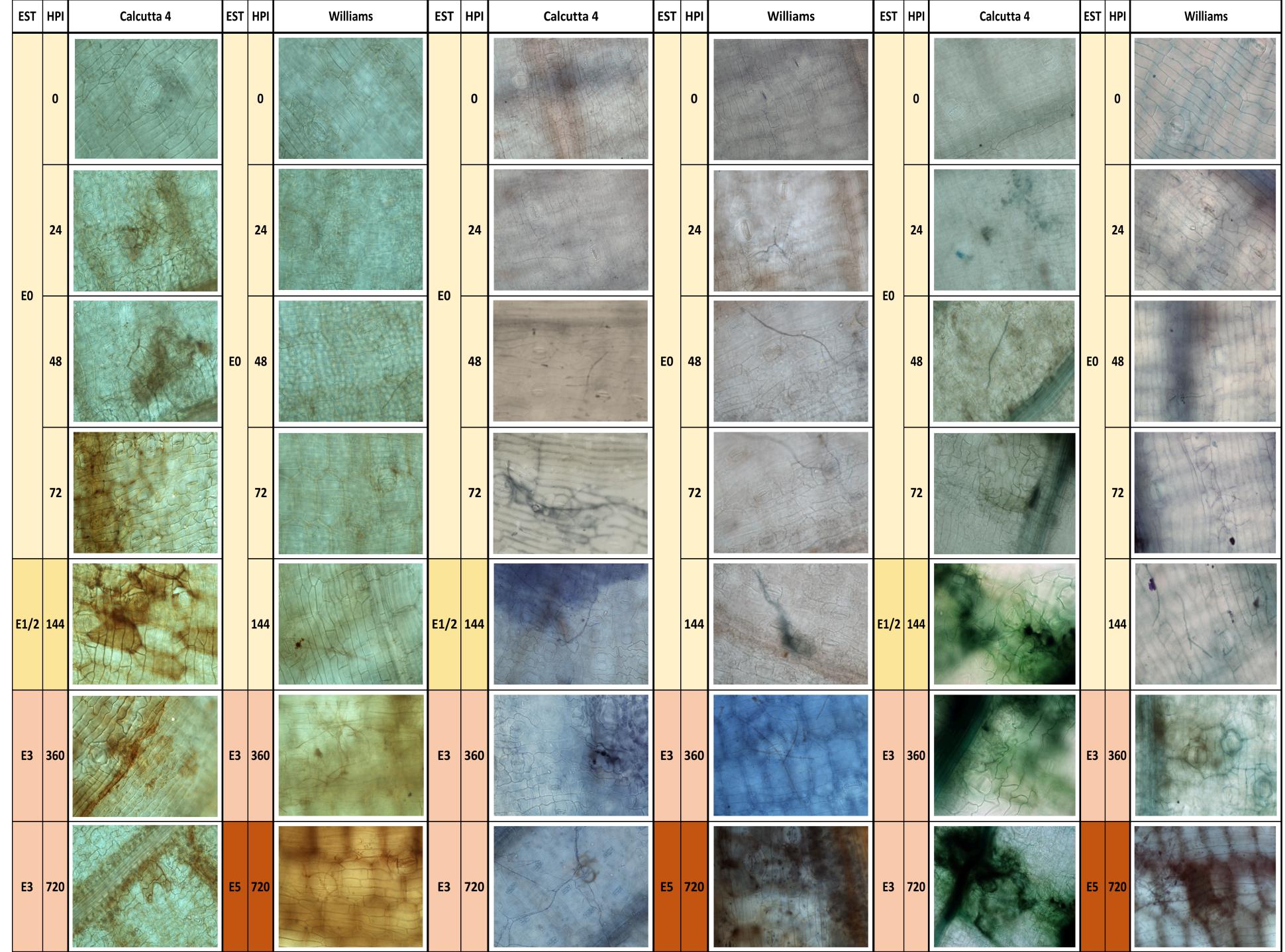
A= ascoporas; Pt= Pseudotecio; h= hifas dentro del tejido

# Seguimiento del hongo *M. fijiensis* en las variedades de Calcutta 4 y Williams mediante SEM

EST	E0			E1/2	E3
HPI	0	12	24	144	360
Cal					
EST	E0			E3	
HPI	0	12	24	144	360
Will					



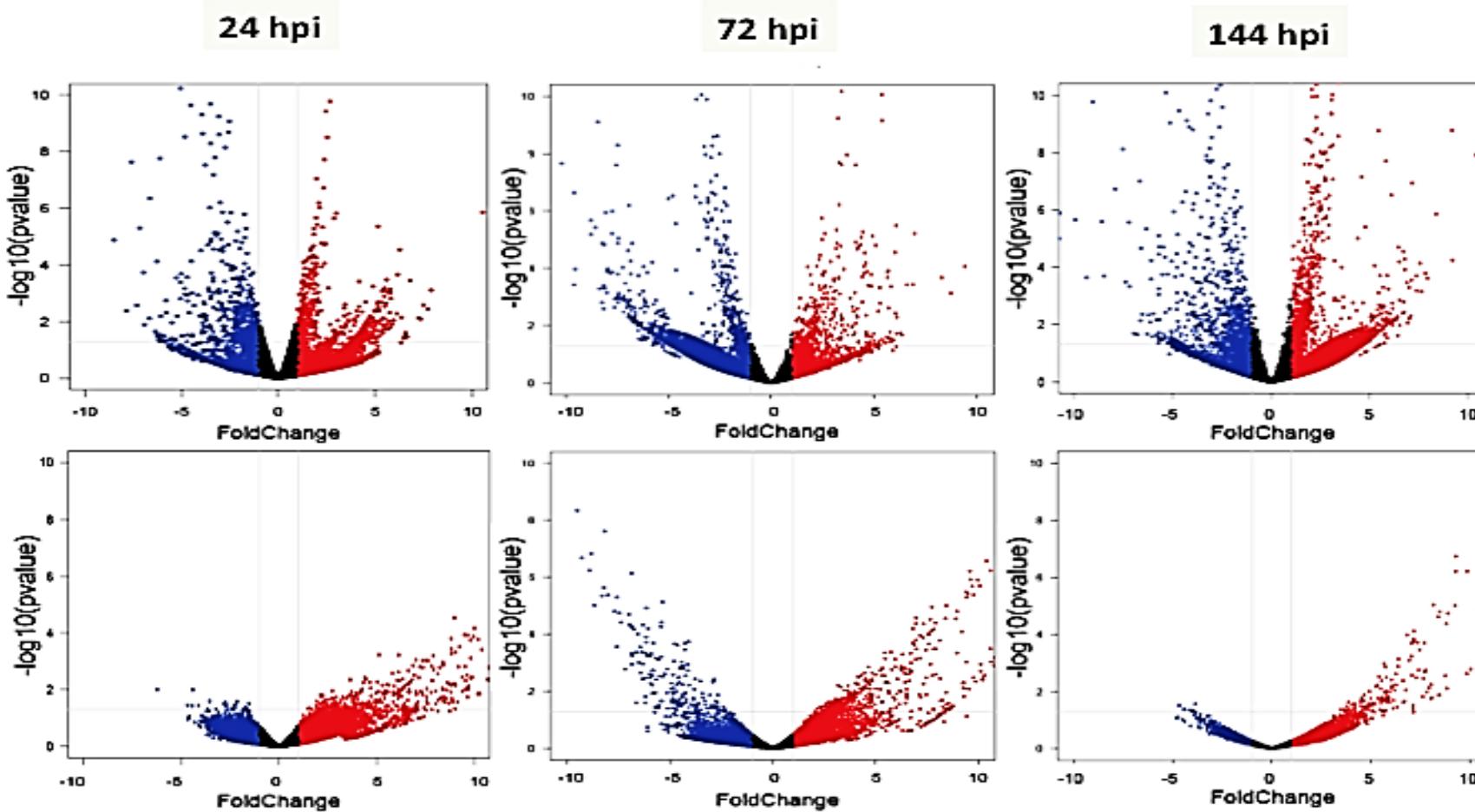
**Figura 2.6.** Fotografías de hojas de Williams durante estadíos avanzados de la Sigatoka Negra. A) Mancha madura durante estadio 5 de la enfermedad; B) Pseudotecio con ascosporas al interior de un estoma; C) Esporodonquios con conidias en mancha joven de la enfermedad. As= ascosporas; c= conidias.



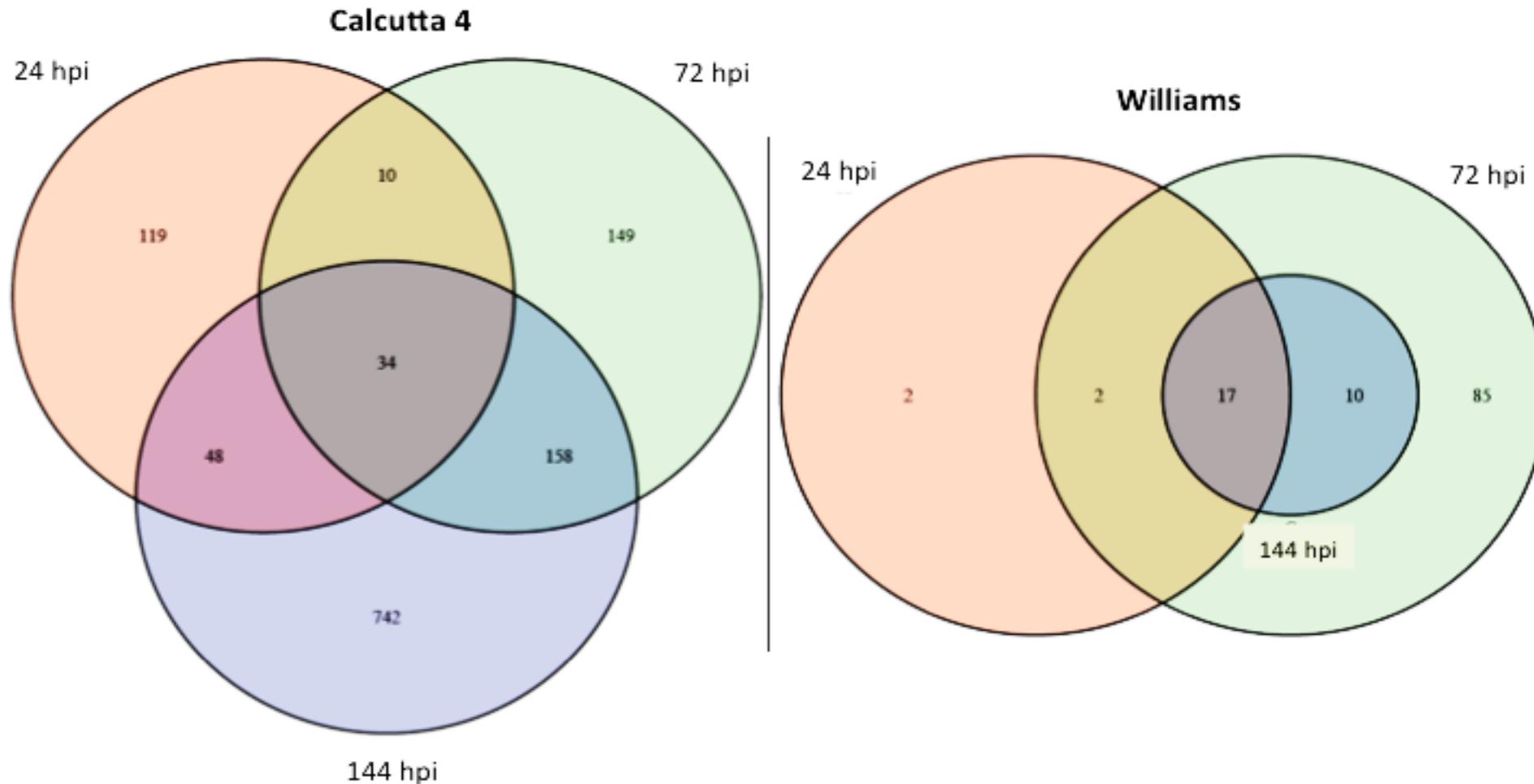
**Perfiles transcriptómicos de genotipos de *Musa acuminata* revelan genes relacionados con la defensa en la interacción con *Mycosphaerella fijiensis***

# Volcano plot

Calcutta 4



## Diagrama de Venn de los genes que se comparten entre los tiempos por variedad

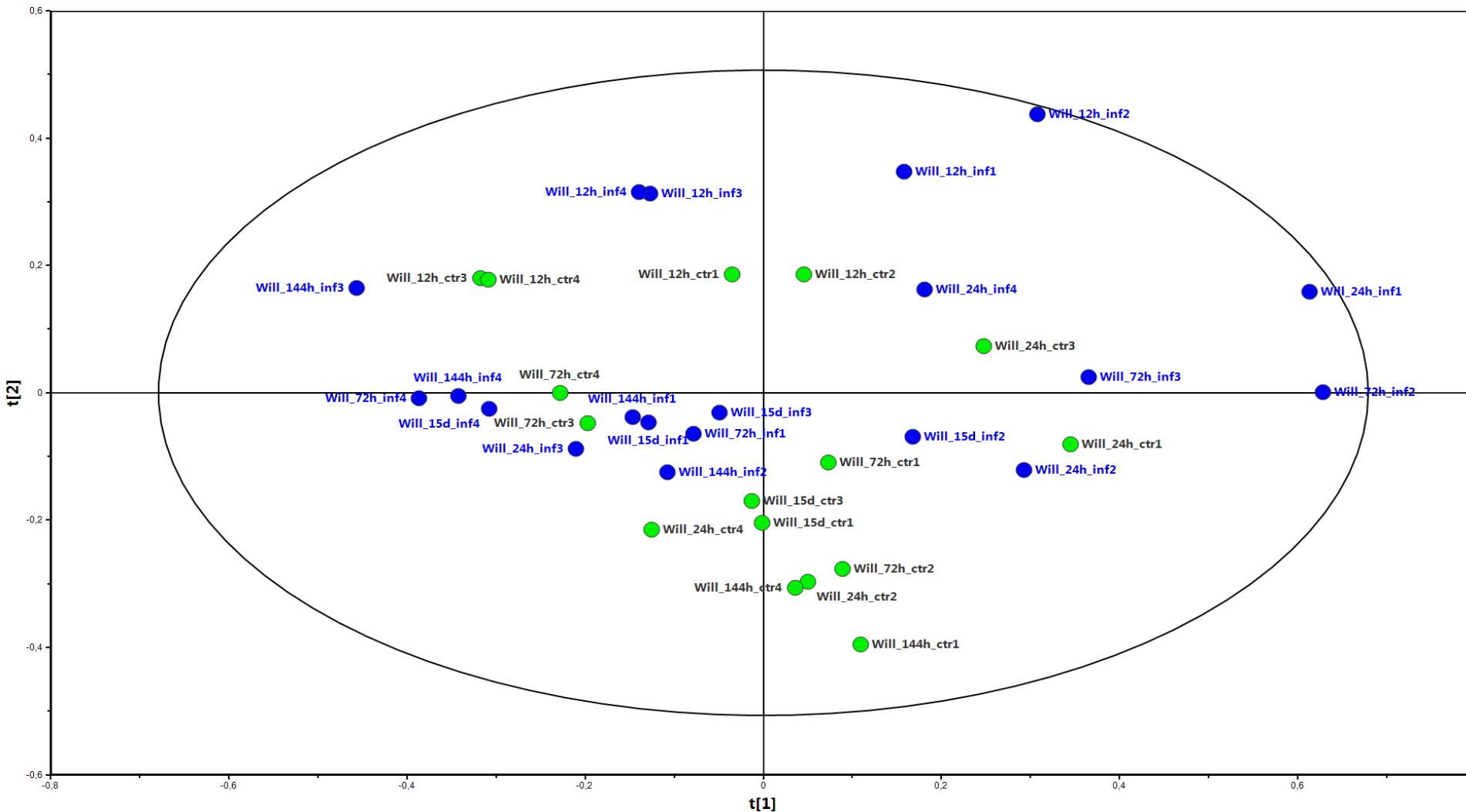


# Transducción de señales hormonales.



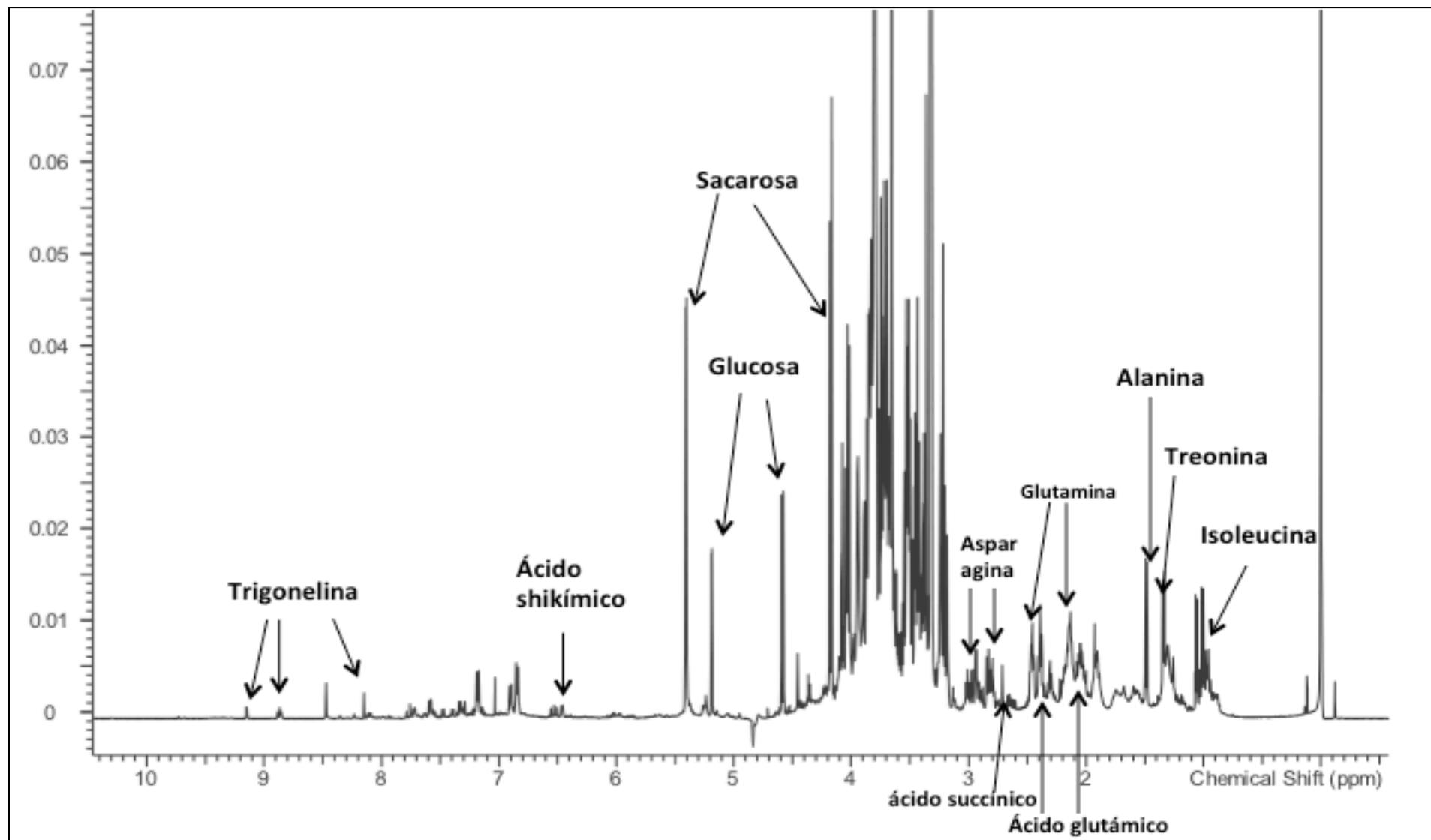
**Perfil metabolico de hojas de banano  
durante la interacción con el hongo  
*Mycosphaerella fijiensis***

## Análisis de componentes principales de la región 0.5-9.0 ppm del espectro de $^1\text{H}$ NMR (500 MHz)

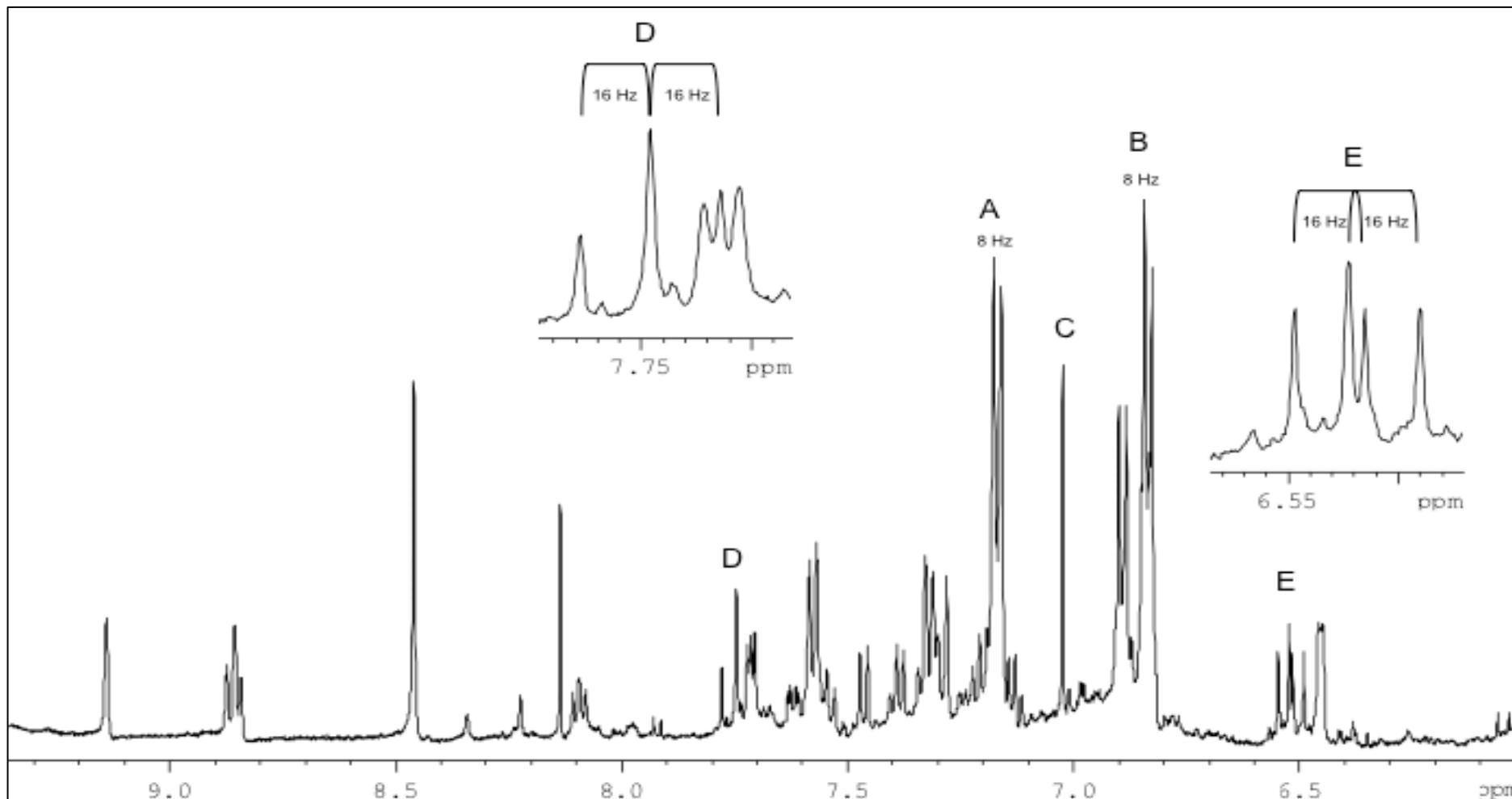


## Variedad susceptible Williams

# Metabolitos identificados en los espectros de $^1\text{H}$ NMR

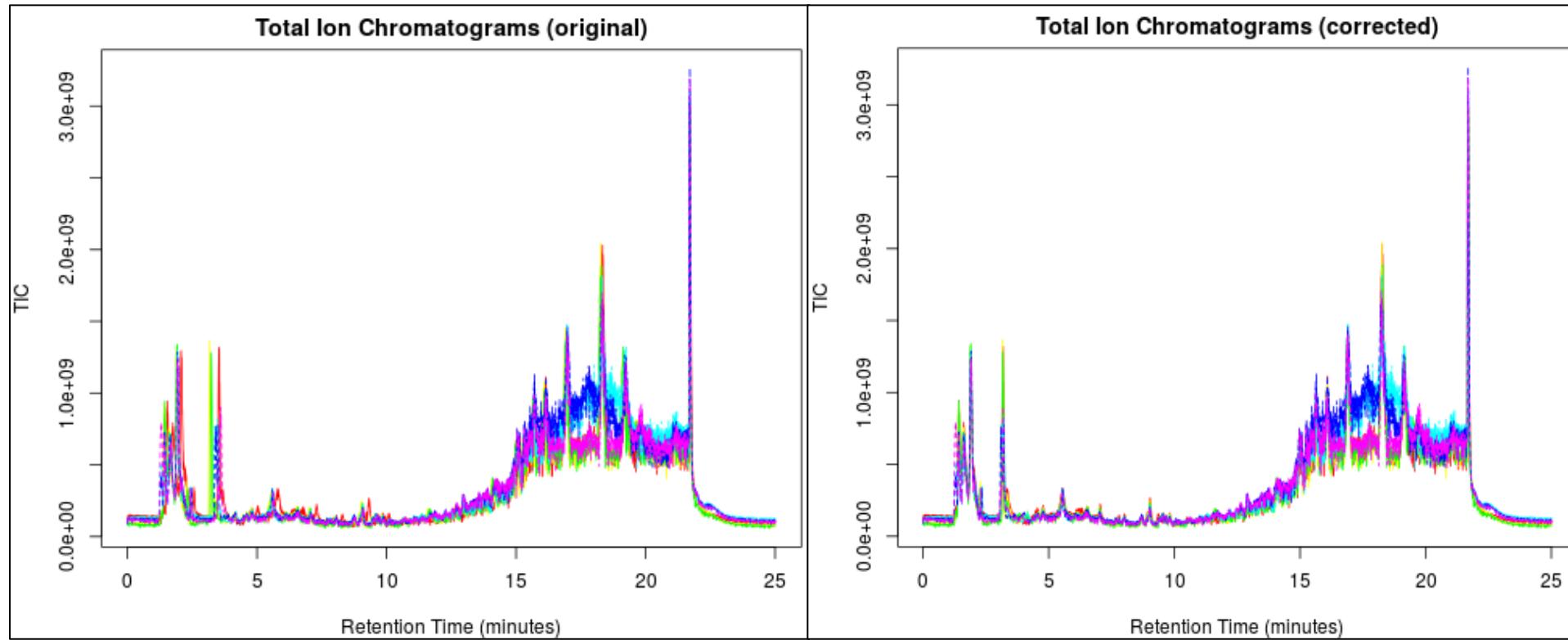


**Señales relacionadas con fenilfenalenonas (A-E),  
evidenciando la presencia de dichos compuestos en las  
muestras**

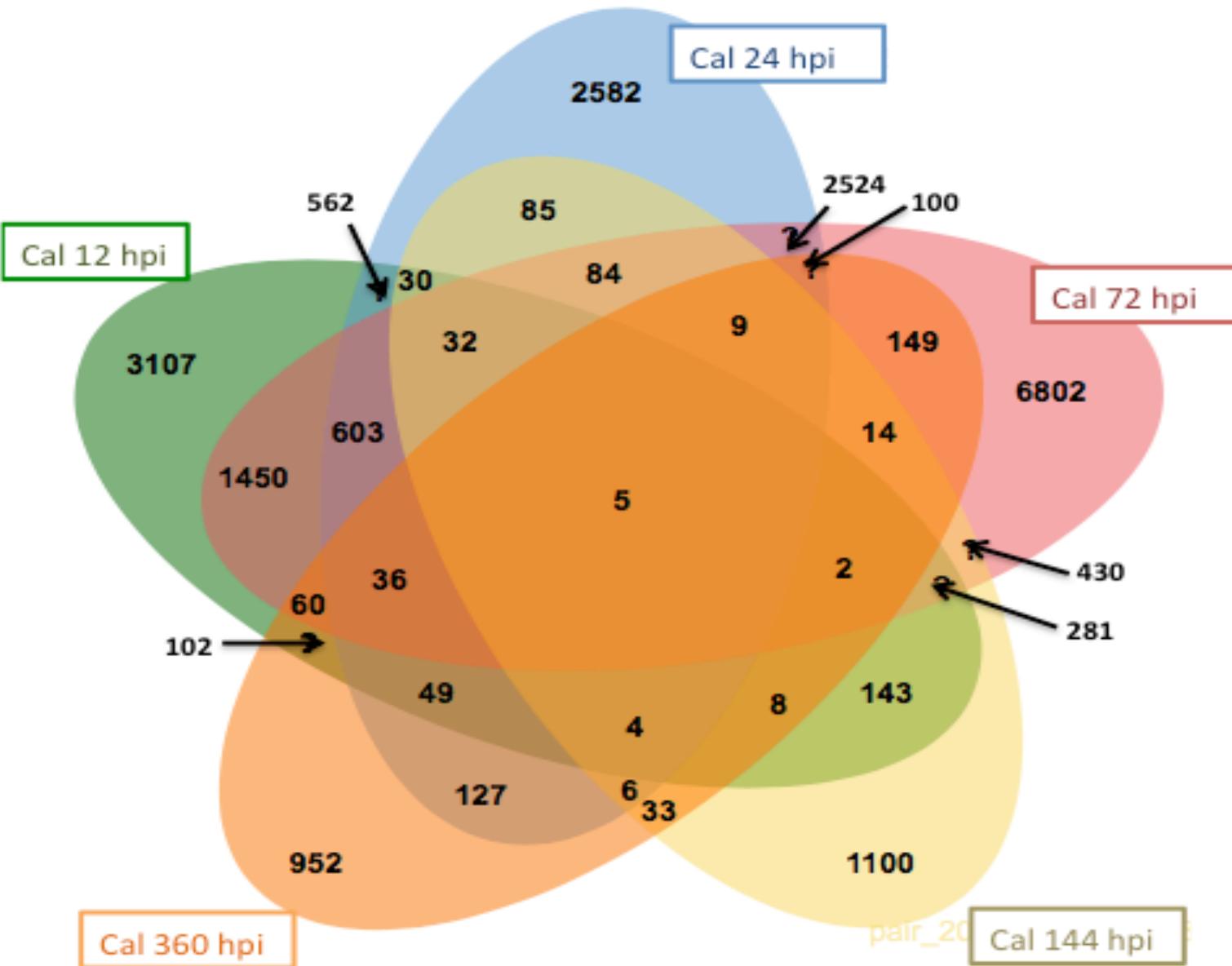


**Análisis no supervisado de metabolitos diferencialmente expresados en plantas de banano durante la inoculación con *M. fijiensis* por UPLC-ESI-MS**

# Cromatograma de iones totales antes (A) y después (B) de la corrección del tiempo de retención



# Diagrama de Venn para comparar las características identificadas por cada tiempo en la variedad Calcutta 4



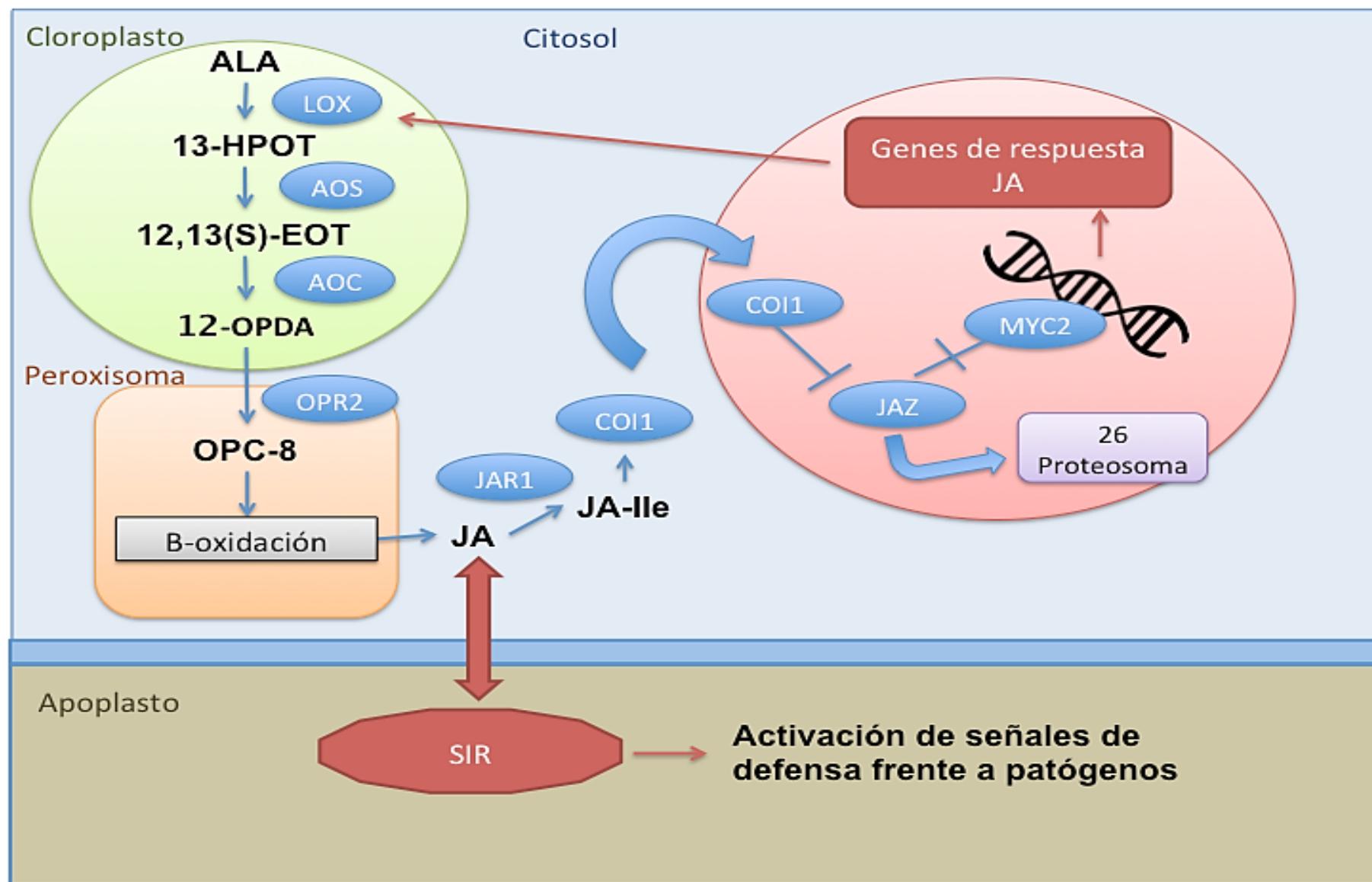
# Principales 30 rutas identificadas para Calcutta 4 a 72 hpi con el hongo *M. fijiensis*

sucrose biosynthesis II	"SOBRE"	primario	5	10	9.7e-3
sucrose degradation II (sucrose synthase)	"SOBRE"	primario	4	8	4.0e-3
Calvin-Benson-Bassham cycle	"SOBRE"	primario	4	12	2.1e-2
glycolysis I (from glucose 6-phosphate)	"SOBRE"	primario	4	11	1.5e-1
gluconeogenesis I	"SOBRE"	primario	4	11	1.5e-1
trehalose degradation II (trehalase)	"SOBRE"	primario	3	4	5.9e-4
guanosine nucleotides degradation I	"SOBRE"	primario	3	5	9.4e-4
pentose phosphate pathway (non-oxidative branch)	"SOBRE"	primario	3	6	1.5e-3
mannitol biosynthesis	"SOBRE"	primario	3	4	9.6e-3
Rubisco shunt	"SOBRE"	primario	3	11	1.5e-2
chlorophyllide abiosynthesis I (aerobic, light-dependent)	"SUB"	primario	3	13	2.2e-1

# Principales 30 rutas identificadas para Calcutta 4 a 72 hpi con el hongo *M. fijiensis*

Ruta	Regulación	tipo de metabolismo	Metabolitos Putativos	Metabolitos Ruta	Valor p
flavonoid biosynthesis	"SOBRE"	secundario	8	13	3.3e-4
kaempferol glycoside biosynthesis	"SOBRE"	secundario	5	14	9.5e-4
tetrahydrofolate biosynthesis II	"SOBRE"	secundario	4	18	1.8e-2
luteolin biosynthesis	"SOBRE"	secundario	3	6	1.5e-3
rosmarinic acid biosynthesis II	"SOBRE"	secundario	3	8	4.0e-3
coumarin biosynthesis (via 2-coumarate)	"SOBRE"	secundario	3	8	4.0e-3
flavonol biosynthesis	"SOBRE"	secundario	3	8	4.0e-3
monolignol glucosides biosynthesis	"SOBRE"	secundario	3	10	9.7e-3
traumatin and (Z)-3-hexen-1-yl acetate biosynthesis	"SOBRE"	secundario	3	10	9.7e-3
carotenoid cleavage	"SOBRE"	secundario	3	21	1.9e-1
phenylpropanoid biosynthesis	"SUB"	secundario	3	17	3.8e-1

# Reconstrucción de la ruta putativa de biosíntesis y transducción de señales de JA en la variedad Calcutta 4 durante la interacción con *Mycosphaerella fijiensis*



# Big data para la Ciencia



**Banana Genome Hub**

Home    Genomes ▾    JBrowse ▾    Search ▾    Download    About ▾    Partners ▾

## Contributions

## About the Banana Genome Hub

The **Banana Genome Hub (BGH)** is a community website that federates genomic data on banana which developed by [Cirad](#) and [the Alliance of Bioversity International and CIAT](#) supported by the [South Green Bioinformatics platform](#). The BGH is based on Drupal and [Tripal](#) and related modules to facilitate the integration between various systems (Jbrowse, [Galaxy](#), [Gigwa](#)) for plant genome analysis.

The Banana Genome Hub (BGH) belongs to the South Green genomes hubs that are part of the [Elixir France Service Delivery plan](#). The BGH has been registered in [bio.tools](#) and [FAIRsharing](#) to facilitate interoperability.



## Data contributions

The BGH intends to centralise data produced by the community on banana genomics. This information can be retrieved after publications by the team or by request from partners.

We aknowledge contributions from



Ellos encuentran mas de 1800 especies de microbios incluyendo 150 especies nuevas de bacterias y mas de 1.2 millones de nuevos genes.



Collecting sea water samples from the Sargasso Sea for the whole genome shotgun sequencing of microbial populations. (Photo supplied by The Center for the Advancement of Genomics.)

# Molecular Ecology

(is everything everywhere ? Does diversity matter for ecosystem function ? )



## Biogeochemistry

(who is where and what are they doing ?)

## Molecular Evolution

(Distribution/Diversity of ecologically meaningful genes/alleles ?)

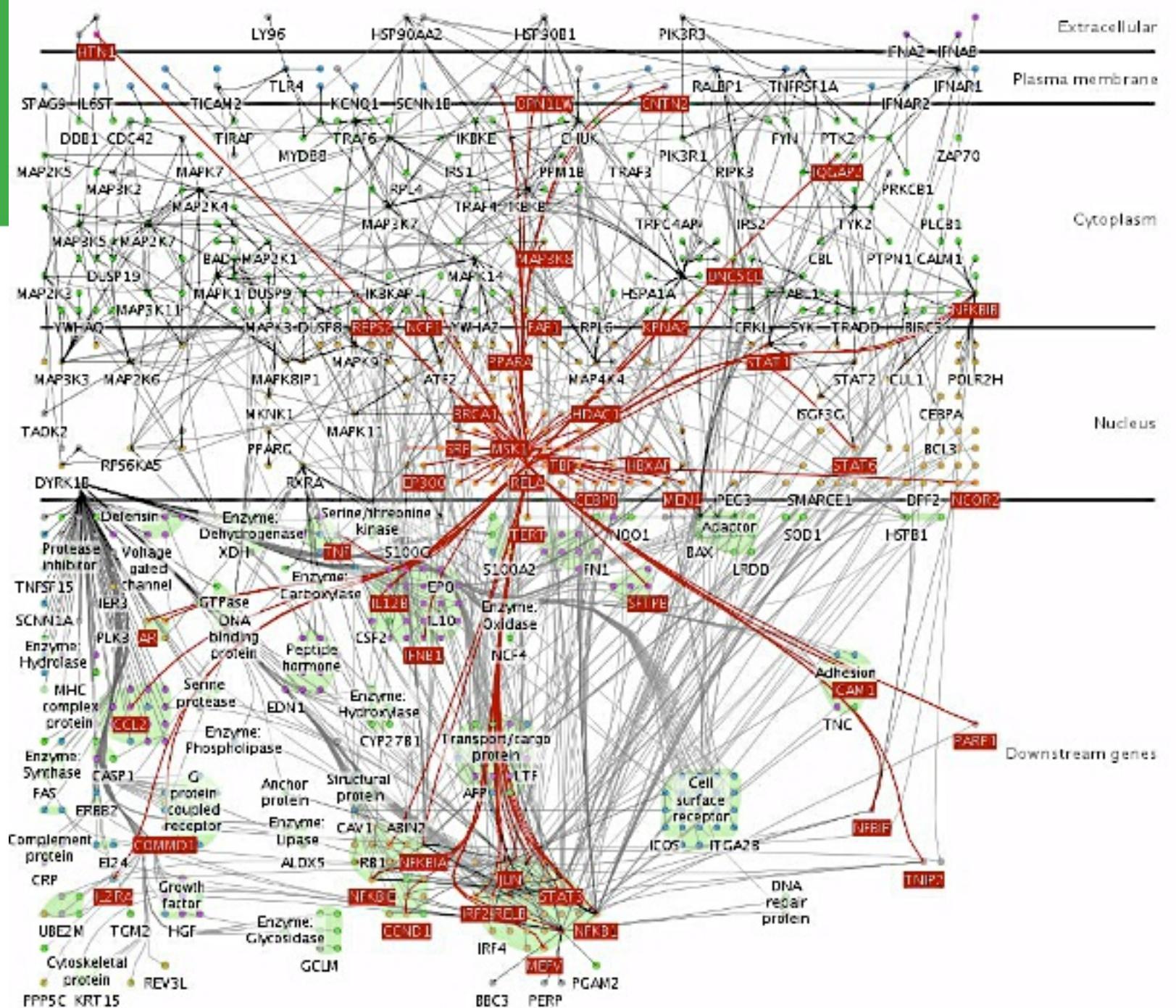
## Environmental Systems Biology

## Genome Biology

(Evolution/origin of gene families, major signaling/regulatory/metabolic mechanisms, duplications/deletions, constraints, organization, structure, size, etc.)

J. Craig Venter Comparative/Functional Genomics

(Lineage/taxa specific features, transcriptomics, protein function)

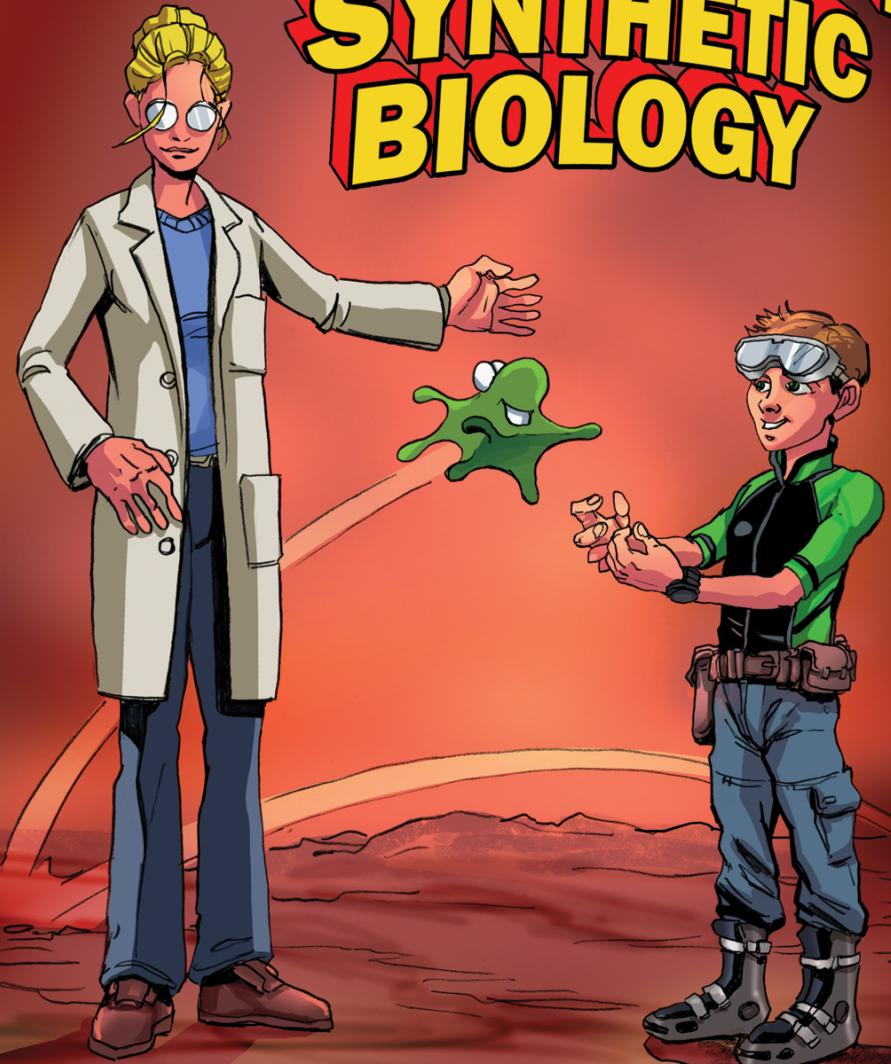


# ¿Cual es el siguiente paso?

¿Como puedo utilizar toda esta información?

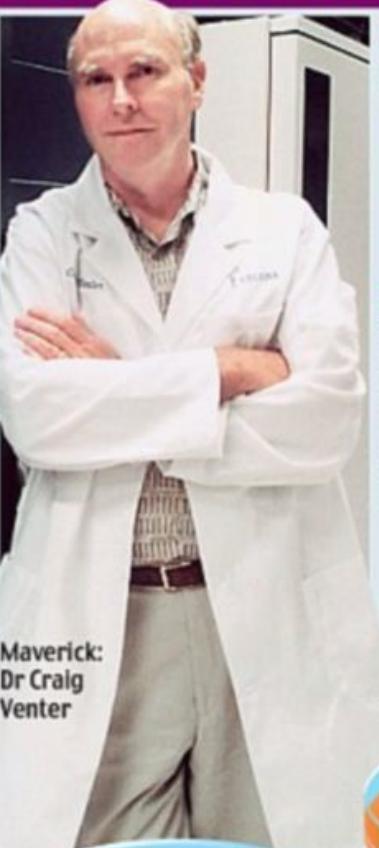
NO. 1  
NOV 2005

# ADVENTURES IN SYNTHETIC BIOLOGY



STORY: DREW ENDY ISADORA DEESE  
THE MIT SYNTHETIC BIOLOGY WORKING GROUP  
ART: CHUCK WADEY [WWW.CHUCKWADEY.COM](http://WWW.CHUCKWADEY.COM)

# HOW TO MAKE ARTIFICIAL LIFE



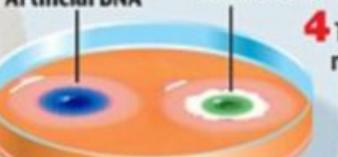
1 Entire DNA of Mycoplasma mycoides, a bug that usually infects goats, is decoded.



2 Researchers buy fragments of DNA from a mail order catalogue. Each of the four bottles of chemicals contains a section of the code.

3 The fragments are put into yeast, which 'stitches' them together, gradually building a synthetic copy of the original DNA.

Artificial DNA      Natural DNA



4 The artificial DNA is put into a recipient bacterium, which then grows and divides, creating two daughter cells, one with the artificial DNA and one with the natural DNA.



5 Antibiotics in the petri dish kill the bacterium with the natural DNA, leaving the one with the synthetic DNA to multiply.



6 Within just a few hours, all traces of the recipient bug are wiped out and bugs with artificial DNA thrive. New life has been created.

7 Possible uses are bugs capable of producing clean fuels and sucking carbon dioxide out of the atmosphere. Also microbes capable of mopping up oil slicks (above) or generating drugs, including the flu vaccine

# Big data para la Ciencia



American College of  
Neuropsychopharmacology

[www.nature.com/npp](http://www.nature.com/npp)

ARTICLE

OPEN

Check for updates

## Genome-wide meta-analysis of alcohol use disorder in East Asians

Hang Zhou , Rasmon Kalayasiri<sup>3,4,5</sup>, Yan Sun<sup>6</sup>, Yaira Z. Nuñez<sup>1,2</sup>, Hong-Wen Deng , Xiang-Ding Chen<sup>8</sup>, Amy C. Justice<sup>9,10</sup>, Henry R. Kranzler , Suhua Chang , Lin Lu , Jie Shi , Kittipong Sanichwankul<sup>14</sup>, Apiwat Mutirangura<sup>5</sup>, Robert T. Malison<sup>1,16</sup> and Joel Gelernter

© The Author(s) 2022

Alcohol use disorder (AUD) is a leading cause of death and disability worldwide. Genome-wide association studies (GWAS) have identified ~30 AUD risk genes in European populations, but many fewer in East Asians. We conducted GWAS and genome-wide meta-analysis of AUD in 13,551 subjects with East Asian ancestry, using published summary data and newly genotyped data from five cohorts: (1) electronic health record (EHR)-diagnosed AUD in the Million Veteran Program (MVP) sample; (2) DSM-IV diagnosed alcohol dependence (AD) in a Han Chinese–GSA (array) cohort; (3) AD in a Han Chinese–Cyto (array) cohort; and (4) two AD Thai cohorts. The MVP and Thai samples included newly genotyped subjects from ongoing recruitment. In total, 2254 cases and 11,297 controls were analyzed. An AUD polygenic risk score was analyzed in an independent sample with 4464 East Asians (Genetic Epidemiology Research in Adult Health and Aging (GERA)). Phenotypes from survey data and ICD-9-CM diagnoses were tested for association with the AUD PRS. Two risk loci were detected: the well-known functional variant rs1229984 in *ADH1B* and rs3782886 in *BRAP* (near the *ALDH2* gene locus) are the lead variants. AUD PRS was significantly associated with days per week of alcohol consumption ( $\beta = 0.43$ ,  $SE = 0.067$ ,  $p = 2.47 \times 10^{-10}$ ) and nominally associated with pack years of smoking ( $\beta = 0.09$ ,  $SE = 0.05$ ,  $p = 4.52 \times 10^{-2}$ ) and ever vs. never smoking ( $\beta = 0.06$ ,  $SE = 0.02$ ,  $p = 1.14 \times 10^{-2}$ ). This is the largest GWAS of AUD in East Asians to date. Building on previous findings, we were able to analyze pleiotropy, but did not identify any new risk regions, underscoring the importance of recruiting additional East Asian subjects for alcohol GWAS.

# Big data para la Ciencia

## ARTICLE

<https://doi.org/10.1038/s41586-018-0043-0>

# Renewing Felsenstein's phylogenetic bootstrap in the era of big data

F. Lemoine<sup>1,2</sup>, J.-B. Domelevo Entfellner<sup>3,4</sup>, E. Wilkinson<sup>5</sup>, D. Correia<sup>1</sup>, M. Dávila Felipe<sup>1</sup>, T. De Oliveira<sup>5,6</sup> & O. Gascuel<sup>1,7\*</sup>

Felsenstein's application of the bootstrap method to evolutionary tree inference has been a cornerstone of phylogenetic inferences. However, increasing numbers of sequences and phylogenies based on hundreds or thousands of taxa are becoming available, and the bootstrap tends to yield very low supports, especially on deep branches. We propose a bootstrap in which the presence of inferred branches in replications is compared to the binary presence or absence index used in Felsenstein's original approach. Our method does not induce falsely supported branches. The application of our method reveals their phylogenetic signals, whereas Felsenstein's bootstrap fails to do so.

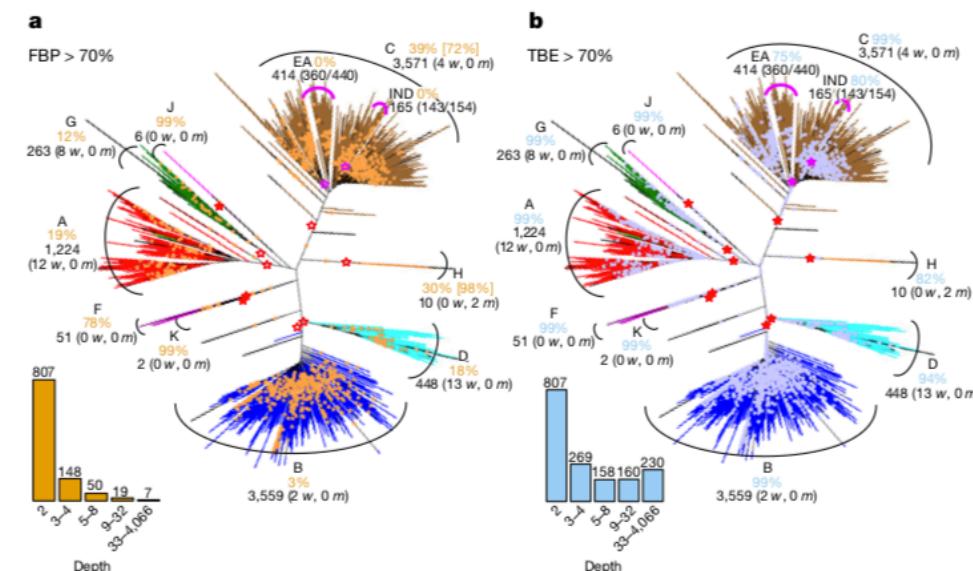
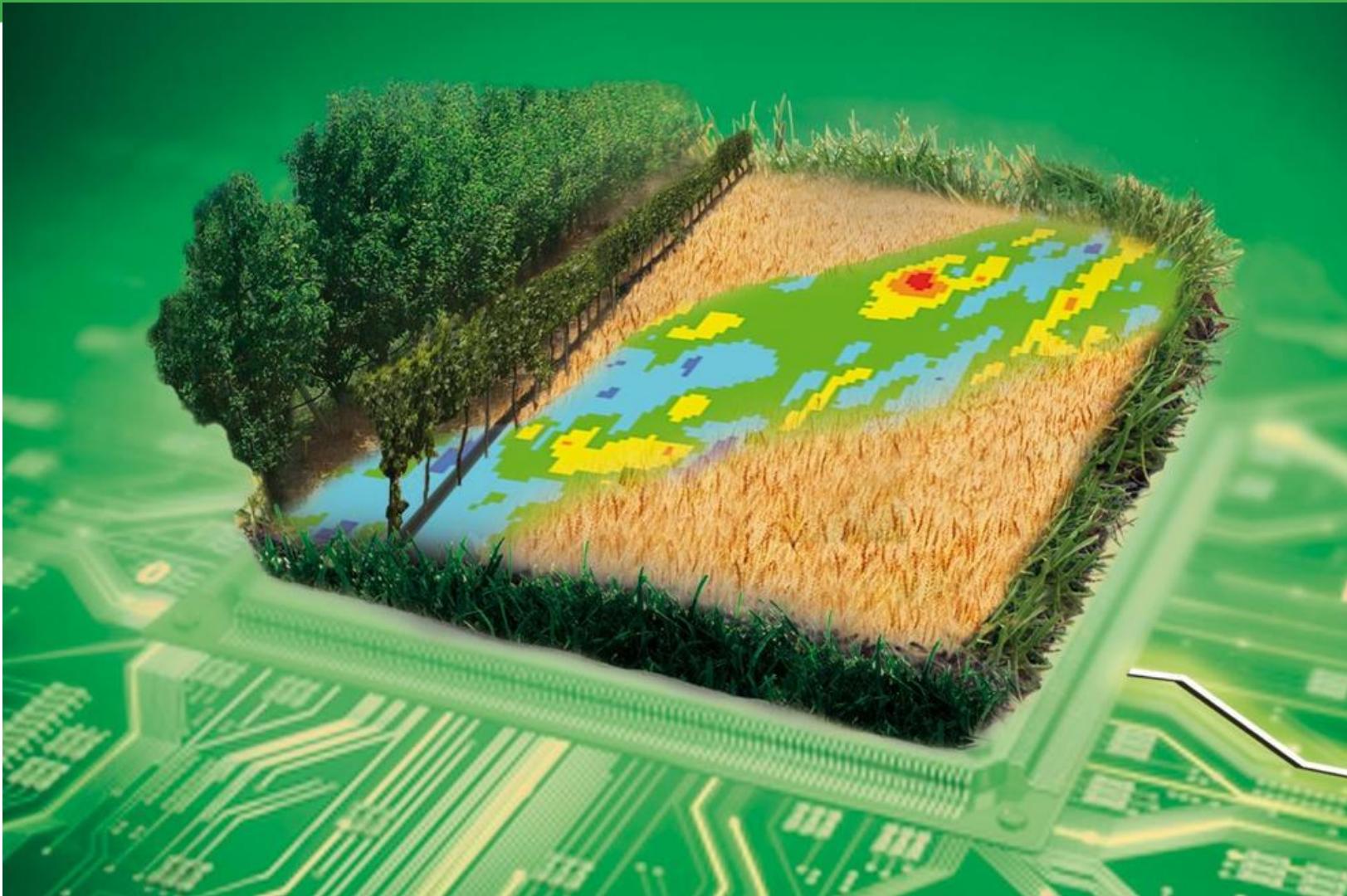


Fig. 1 | FBP and TBE bootstrap supports on the same phylogeny with 9,147 HIV-1M *pol* sequences. a, FBP. b, TBE. Subtypes are colourized;

supports larger than 70%, which are shown in square brackets. The same approach is applied to the C sub-epidemics in India (IND) and East Africa

# Big data para la Ciencia



# Big data para la Ciencia

Article | Open Access | Published: 13 October 2021

## From calibration to parameter learning: Harnessing the scaling effects of big data in geoscientific modeling

Wen-Ping Tsai, Dapeng Feng, Ming Pan, Hylke Beck, Kathryn Lawson, Yuan Yang, Jiangtao Liu & Chaopeng Shen 

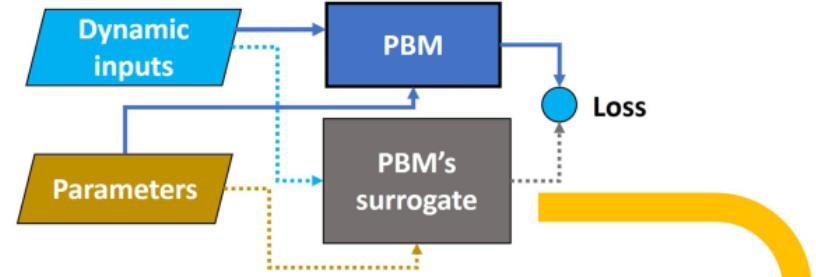
Nature Communications 12, Article number: 5988 (2021) | [Cite this article](#)

8313 Accesses | 3 Citations | 64 Altmetric | [Metrics](#)

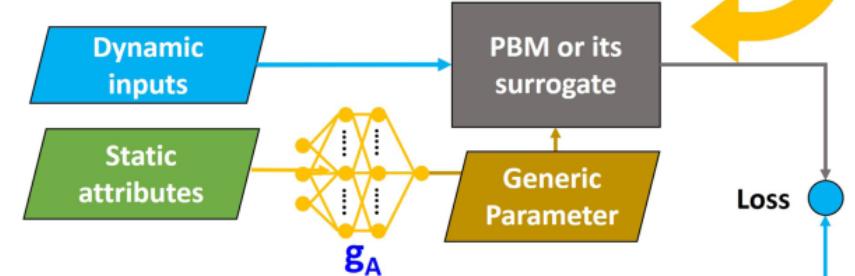
### Abstract

The behaviors and skills of models in many geosciences (e.g., hydrology and ecosystem sciences) strongly depend on spatially-varying parameters that need calibration. A well-calibrated model can reasonably propagate information from observations to unobserved variables via model physics, but traditional calibration is highly inefficient and results in non-unique solutions. Here we propose a novel differentiable parameter learning (dPL) framework that efficiently learns a global mapping between inputs (and optionally responses) and parameters. Crucially, dPL exhibits beneficial scaling curves not previously demonstrated to geoscientists: as training data increases, dPL achieves better performance, more physical coherence, and better generalizability (across space and uncalibrated

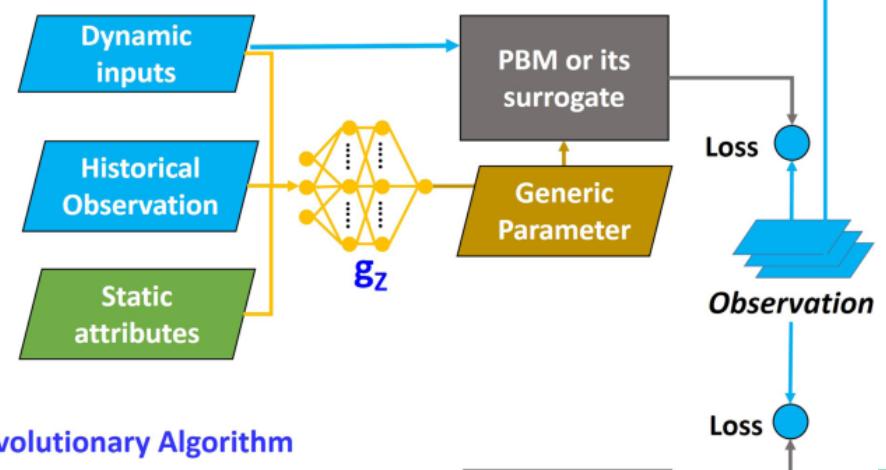
(a) PBM or PBM's surrogate (optional)



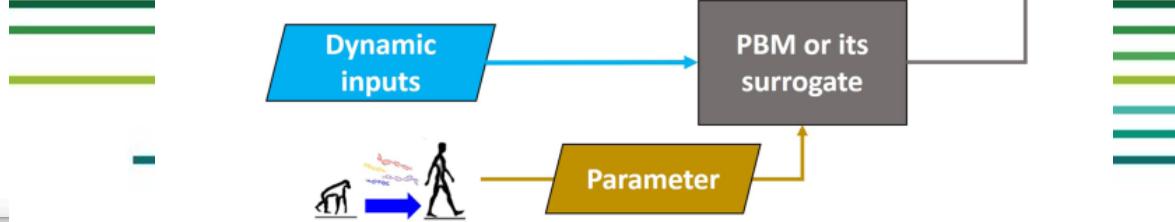
(b) dPL  $g_A$  framework (if historical observations are unavailable)



(c) dPL  $g_z$  framework (if historical observations are available)



(d) Evolutionary Algorithm



# Perspectivas

Muchas Gracias por su atención ....



1 8 0 3

# UNIVERSIDAD DE ANTIOQUIA



@UdeA



@UdeA



@universidaddeantioquia