

Ciclo Itinerante de Bioinformática y Exploración en Recursos Computacionales



TALLER HANDS-ON: EXPLORANDO HERRAMIENTAS DE ENSEMBL

Uso de herramientas de Ensembl y **Ensembl Genomes**

13:30 - 17:30 HS

- Facultad de Agronomía
- Salón de Cómputo (junto a Cantina)





UNIVERSIDAD

URUGUAY

Carla Filippi cfilippi@fagro.edu.uy

Instituto Bioinformático Europeo (EMBL-EBI, Hinxton UK)

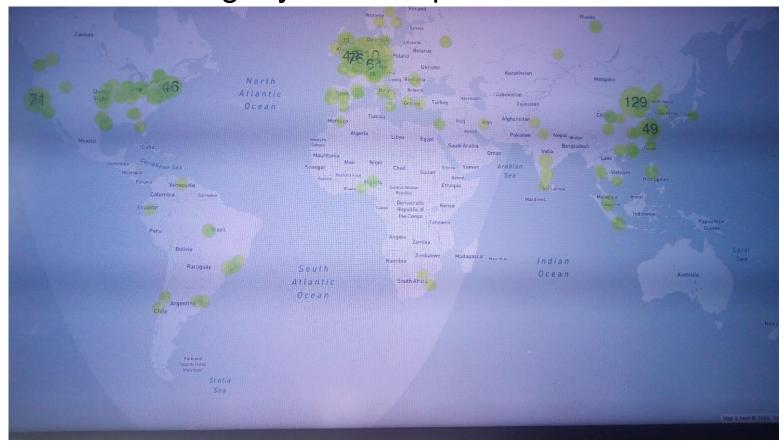




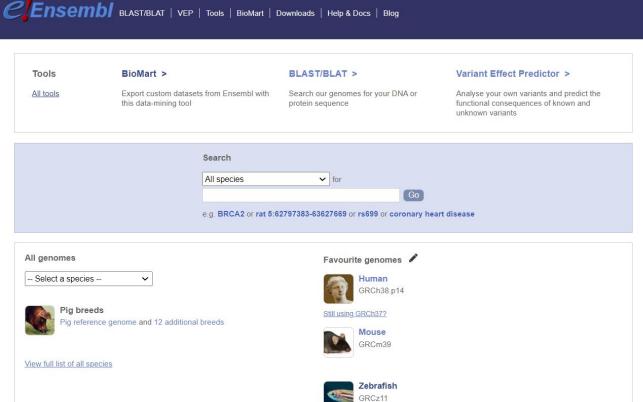




Poniendo a Uruguay en el mapa del EMBL-EBI









Ensembl is a genome browser for vertebrate genomes that supports research in comparative genomics, evolution, sequence variation and transcriptional regulation. Ensembl annotate genes, computes multiple alignments, predicts regulatory function and collects disease data. Ensembl tools include BLAST, BLAT, BioMart and the Variant Effect Predictor (VEP) for all supported species.

Ensembl Release 112 (May 2024)

- Many new fish genomes have been added to Ensembl
- Population frequency data is available for chicken, dog, goat and sheep through VEP
- Update to our current regulation annotation. The promoters now align with the 5' ends of known transcripts
- VEP will be updated to use the dbNSFP commercial data release

More release news

on our blog

Ensembl Rapid Release

New assemblies with gene and protein annotation every two weeks.

Note: species that already exist on this site will continue to be updated with the full range of annotations.



The Ensembl Rapid Release website provides annotation for recently produced, publicly available vertebrate and non-vertebrate genomes from biodiversity initiatives such as Darwin Tree of Life, the Vertebrate Genomes Project and the Earth BioGenome Project.

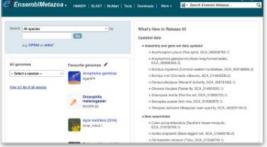
Rapid Release news P on our blog

Ensembl Genomes (https://ensemblgenomes.org/)

EnsemblPlants



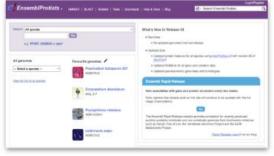
EnsemblMetazoa



EnsemblBacteria



EnsemblProtists

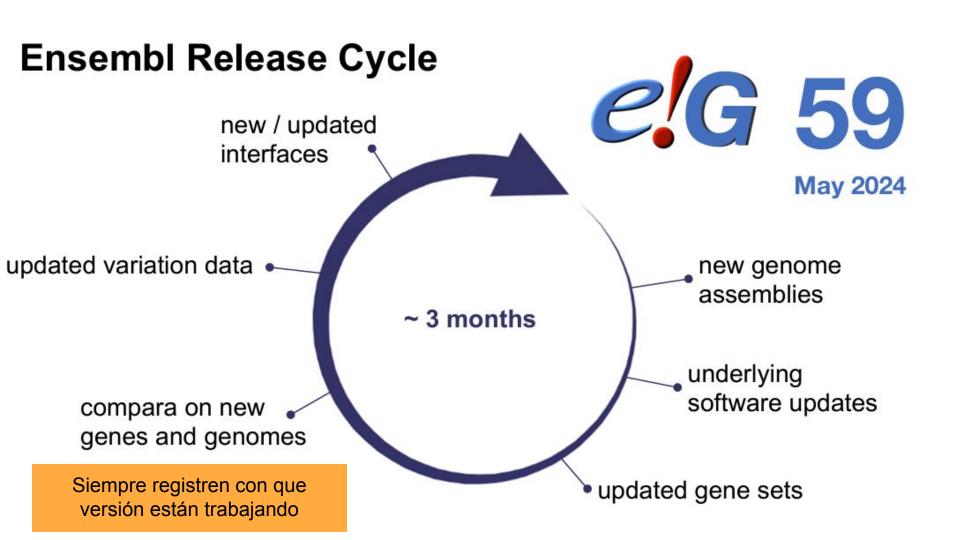


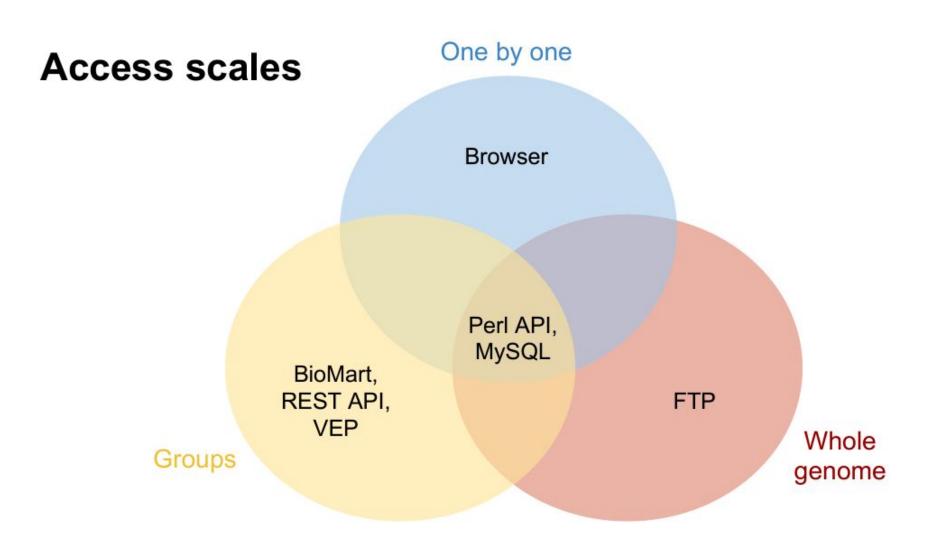
EnsemblFungi



EnsemblCOVID-19







Herramientas de Ensembl (y Ensembl genomes)



Tools

BioMart >

All tools

Export custom datasets from Ensembl with this data-mining tool BLAST/BLAT >

Search our genomes for your DNA or protein sequence

Variant Effect Predictor >

Analyse your own variants and predict the functional consequences of known and unknown variants

Hilo conductor del taller**



**No hay vino, solo café

Situación 1

RNAseq de Vitis vinifera (vid)

Caracter: respuesta a infección con hongo (grey mold, causado por *Botrytis cinerea*)

Tratamientos: control (0h), 72h y 120h post infección.

Datos (listas de genes diferencialmente expresados, DEG) disponibles en GitHub RSG



https://doi.org/10.3389/fpls.2023.1127206

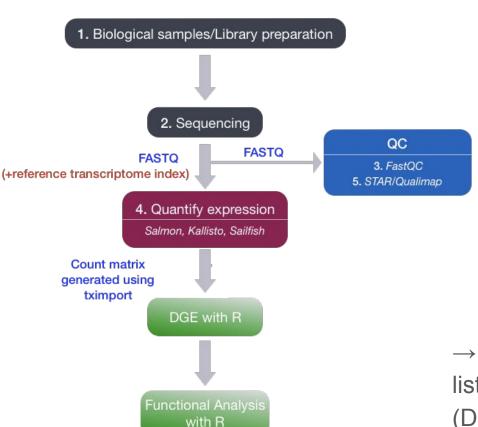
SRA ID: PRJNA788159

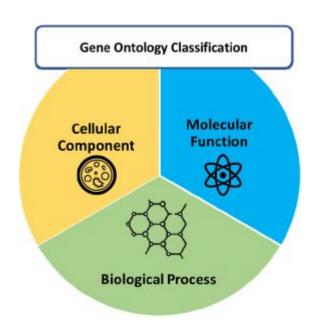
Johann Georg Sturm, 1796

Vamos a estar trabajando con datos de Vitis vinifera

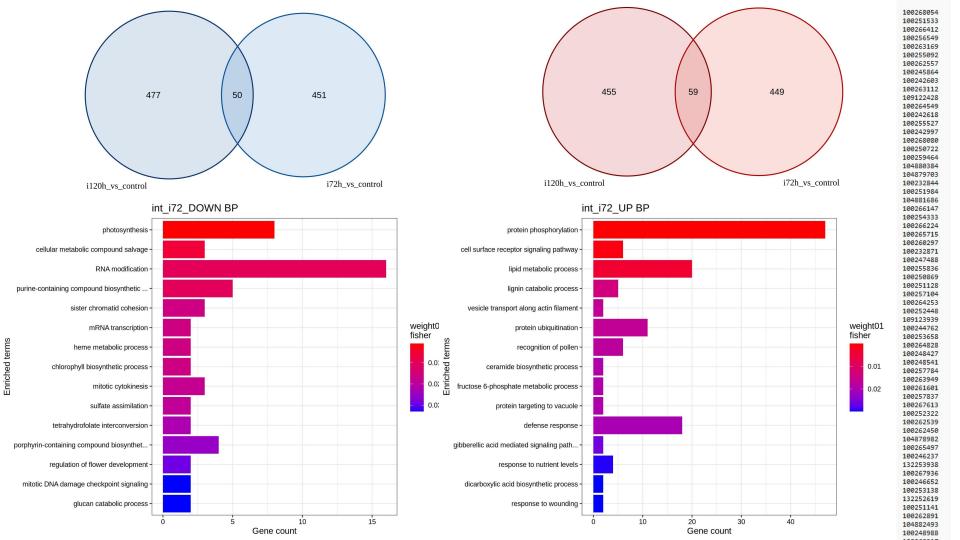
- Vayan al portal de Ensembl Genomes que corresponda
- Registren la versión actual de Ensembl Genomes
- Encuentren el sitio para V. vinifera
- Registren el nombre del assembly (genoma de referencia), su tamaño en pb, número de cromosomas y cantidad de genes anotados (pista: More information and statistics)

Breves lineamientos RNAseq





→ Su bioinformático amigo les envía las listas de genes diferencialmente expresados (DEGs) y GOs/KEGG





Herramienta 1: BioMart





Use this data-mining tool to export custom datasets from Ensembl.

https://plants.ensembl.org/biomart/martview/



Step 1: Dataset

Choose the database and species.

Step 2: Filters

Narrow down the dataset.

Specify your output and what to print on your table.

Step 3: Attributes

Choose the format of your results and export.

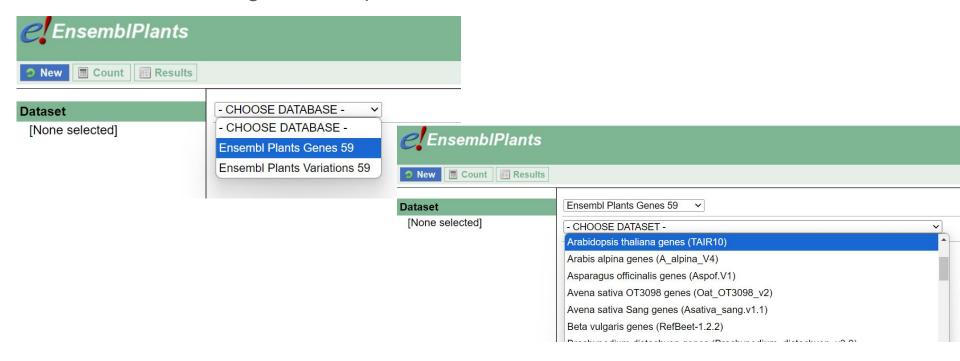
Step 4: Results





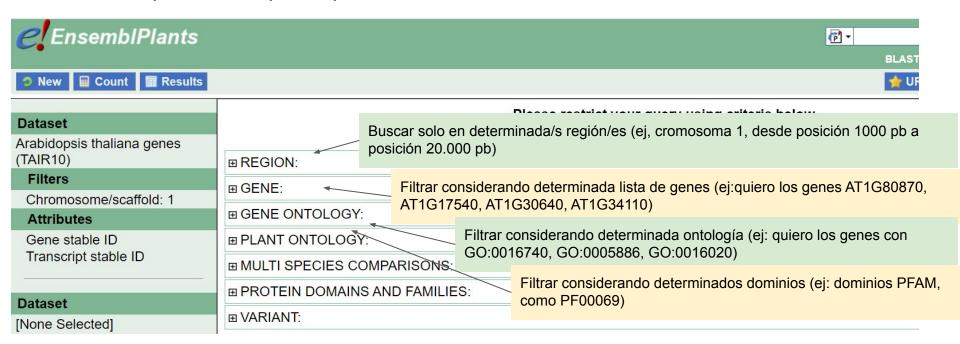


Use this data-mining tool to export custom datasets from Ensembl.



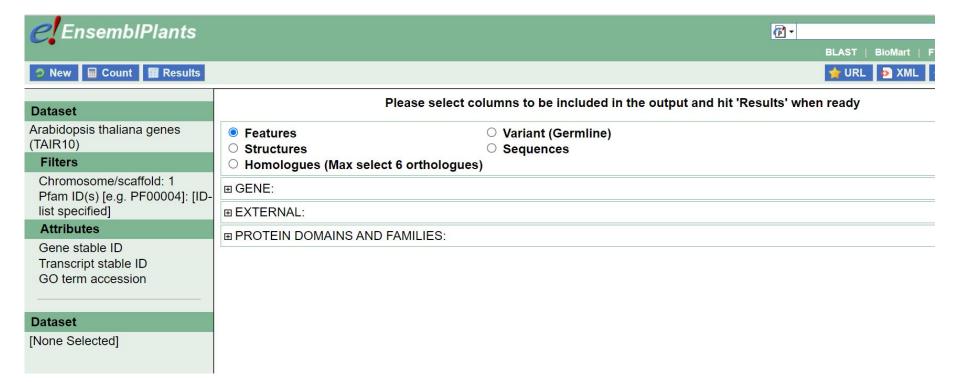
Herramienta 1: BioMart

Filters → que filtros aplicar para extraer datos de interés?



Herramienta 1: BioMart

Attributes → que quiero recuperar? (GO? secuencias? dominios? IDs?)



coffee



Para esta actividad tendremos listas de genes diferencialmente expresados (UP: *upregulados* y DOWN, *downregulados*) → descargar de Github

Datos útiles:

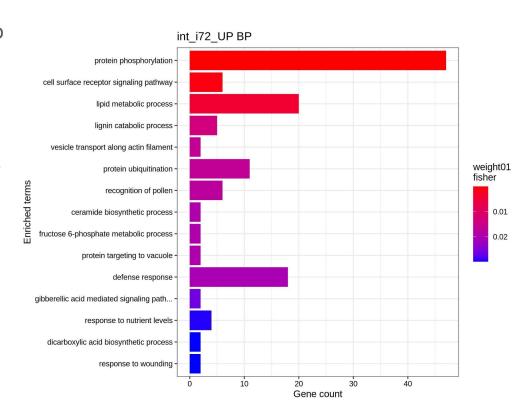
- El análisis se hizo con referencia de NCBI, por lo que los IDs de las listas están en formato NCBI (anteriormente Entrez gene ID)
- Vamos a usar Biomart para responder algunas preguntas (pueden inventarse otras)

- 455 59 449 172h_vs_control 172h_vs_control
- 1. En Biomart, seleccionar el dataset correspondiente (i.e. *Vitis*)
- Ingresar como filtro nuestra lista de genes upregulados a 72h post infección (cuyos IDs están codificados en formato NCBI) unique_72h_UP.txt
- 3. Obtener para cada gen de la lista el ID usado por la comunidad científica para esa especie (puntualmente para Vitis, se llaman "Vitvi_XXX")
- 4. Obtener las secuencias de las proteínas codificadas por estos genes (formato FASTA, aminoácidos).

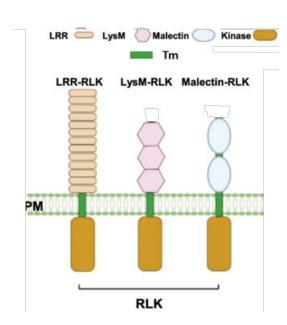
(5. opcional) Pueden mejorar los encabezados de esos FASTA haciéndolos más informativos. Ej, agregar cromosoma, start y end, etc)

6. Tenemos este gráfico de enriquecimiento GO, donde se ven los procesos biológicos enriquecidos en los genes de nuestra lista de DEGs upregulados a las 72h post infección con el hongo. Vemos que "protein fosforilation" es el término más representado.

Averiguar qué genes de mi lista tienen ese GO BP (NOTA: el código es GO:0006468. Para buscar términos GO, ir a https://www.ebi.ac.uk/QuickGO/).



- 7. Averiguar cuántos genes de mi lista codifican para quinasas. (pista: ver dominios, ejemplo para la base de datos PFAM, el dominio quinasa es PF00069)
- 8. Identificar cuáles de estas quinasas DEG tienen dominios transmembrana, como posibles RLK (receptor like kinase, proteínas que participan en procesos de defensa a patógenos en plantas).
- 9. Recuperar en formato FASTA (nucleótidos) la región río arriba de cada una de estas quinasas (1000 pb rio arriba)





Situación 3

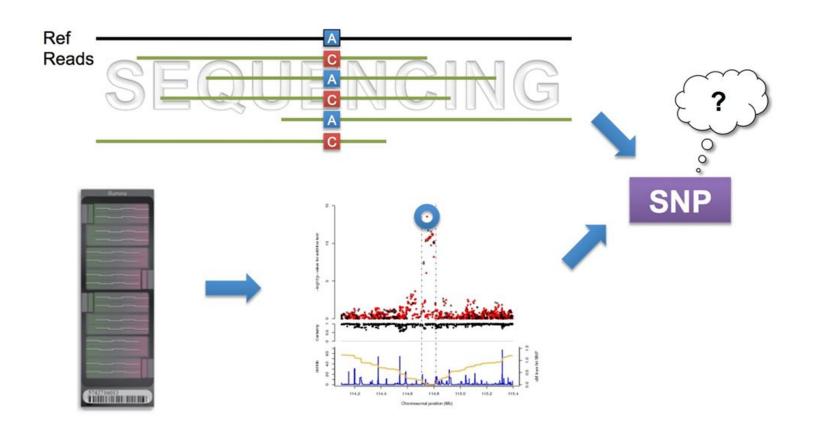
Búsqueda de polimorfismos en Saccharomyces cerevisiae

Tenemos datos de polimorfismos de tipo SNP (polimorfismo de nucleótido simple)

Queremos saber si afectan genes y sus consecuencias



Gráficamente...



The Variant Effect Predictor

(Similar a SNPeff, 10.4161/fly.19695)

Data input:

Variant coordinates

VCF

1 30609607 . A C . . . 3 6658706 . A T . . . 2 35790350 . GC G . . .

Variant ID



VEP output:

Genes, transcripts affected

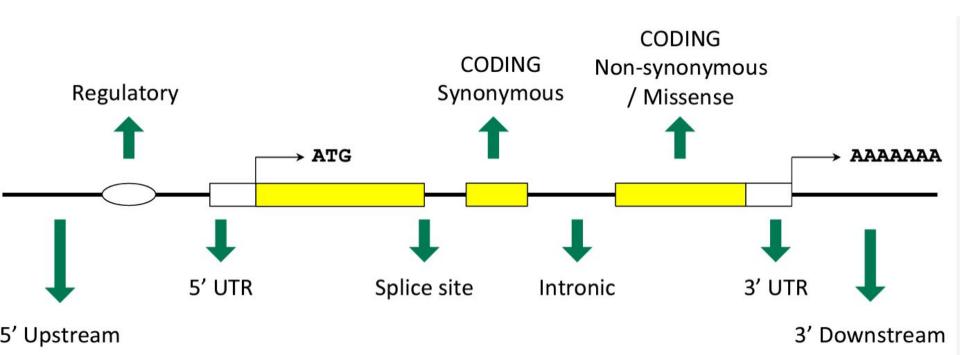
Pathogenicity

4 **

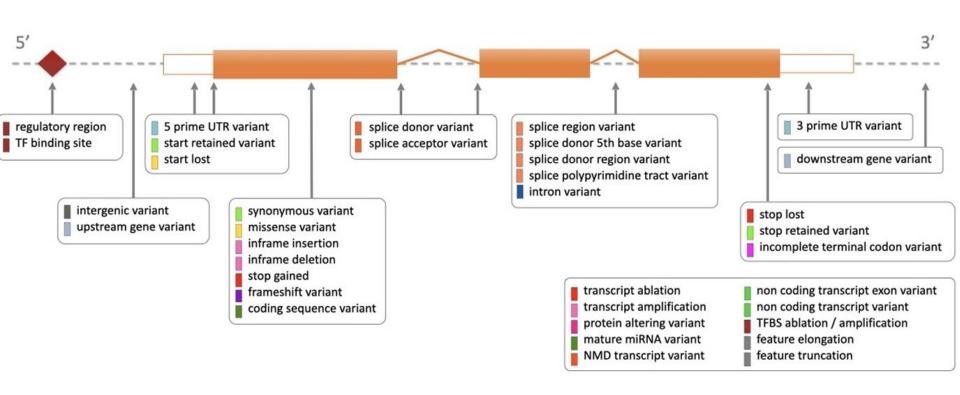
Frequency data Regulatory and splicing consequences

PubMed citations

Consecuencias de las variantes



Sequence Ontology (SO)



Vamos a trabajar con datos de *S. cerevisiae* → vayan al portal del Ensembl Genomes correspondiente.

- Encuentren el sitio para S. cerevisiae
- Registren el nombre del assembly (genoma de referencia), su tamaño en pb, número de cromosomas y cantidad de genes anotados (pista: More information and statistics)

- Usando VEP, predecir el efecto de las variantes disponibles en el VCF (descargar de GitHub)
- 2. Cuantas variantes tiene nuestro archivo?
- 3. Cuales son las consecuencias más frecuentes?
- 4. Cuántas variantes generan un "missense_variant"? (pista: usar filtros, consequences)
- 5. Cuantas variantes son de impacto alto (HIGH)? Cuáles son sus consecuencias? (Ojo que los filtros son sensibles a mayúsculas/ minúsculas)
- 6. Hay variantes que afectan el gen de la PDC1 (Pyruvate Decarboxylase 1)? (pista: su geneID/Symbol es PDC1, y está en cromosoma XII: posiciones 232390-234081)



Actividad 5 (integrando todo lo visto)

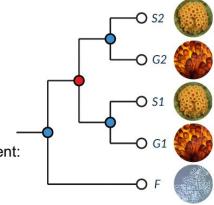
No encontramos variantes que afectan el gen de la PDC1... pero quizá a alguno de sus parálogos? (Si es que tiene)

Orthologues:

- Genes emerged through a speciation event:
 - e.g. G1 and S1; G2 and F; S2 and F
- 1-to-1: G1 and S1
- 1-to-many: F and S1, S2, G1, G2

Paralogues:

- Genes emerged through a duplication event:
 - e.g. G1 and G2, S1 and S2
- Within species: G1 and G2
- Between species: G1 and S2



Buscar los parálogos de PDC1. Si hay varios, quedarse con el de mayor % de similitud (pista: BioMart) → luego VEP p/ variantes, y sus consecuencias)

Exploración libre de Ensembl Browser (ejemplo, para el parálogo que identificamos recién)

