# Monopoly ANR project :

## Integrated gene network analyses between mono- and poly-ovulating species

S. Maman, G. Agoutin, J. Sarry, F. Plisson-Petit, S. Fabre, C. Genêt
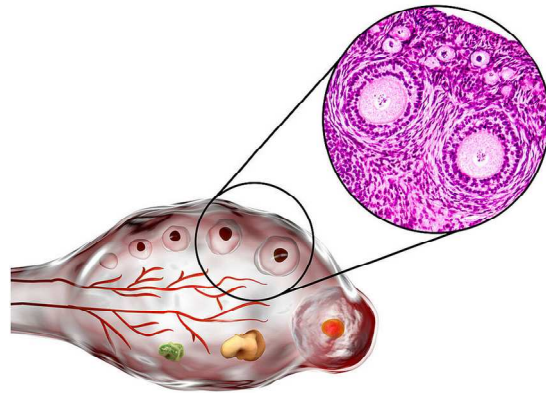
**GenROC team**

RÉPUBLIQUE FRANÇAISE
Liberté
Égalité
Fraternité

INRAE

GenPhySE
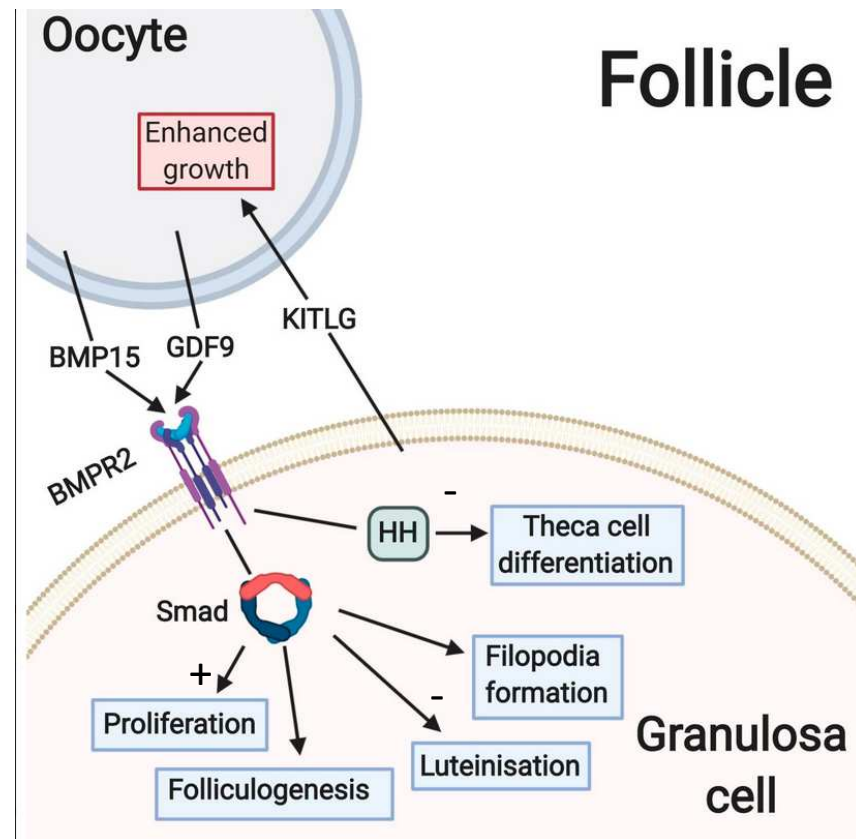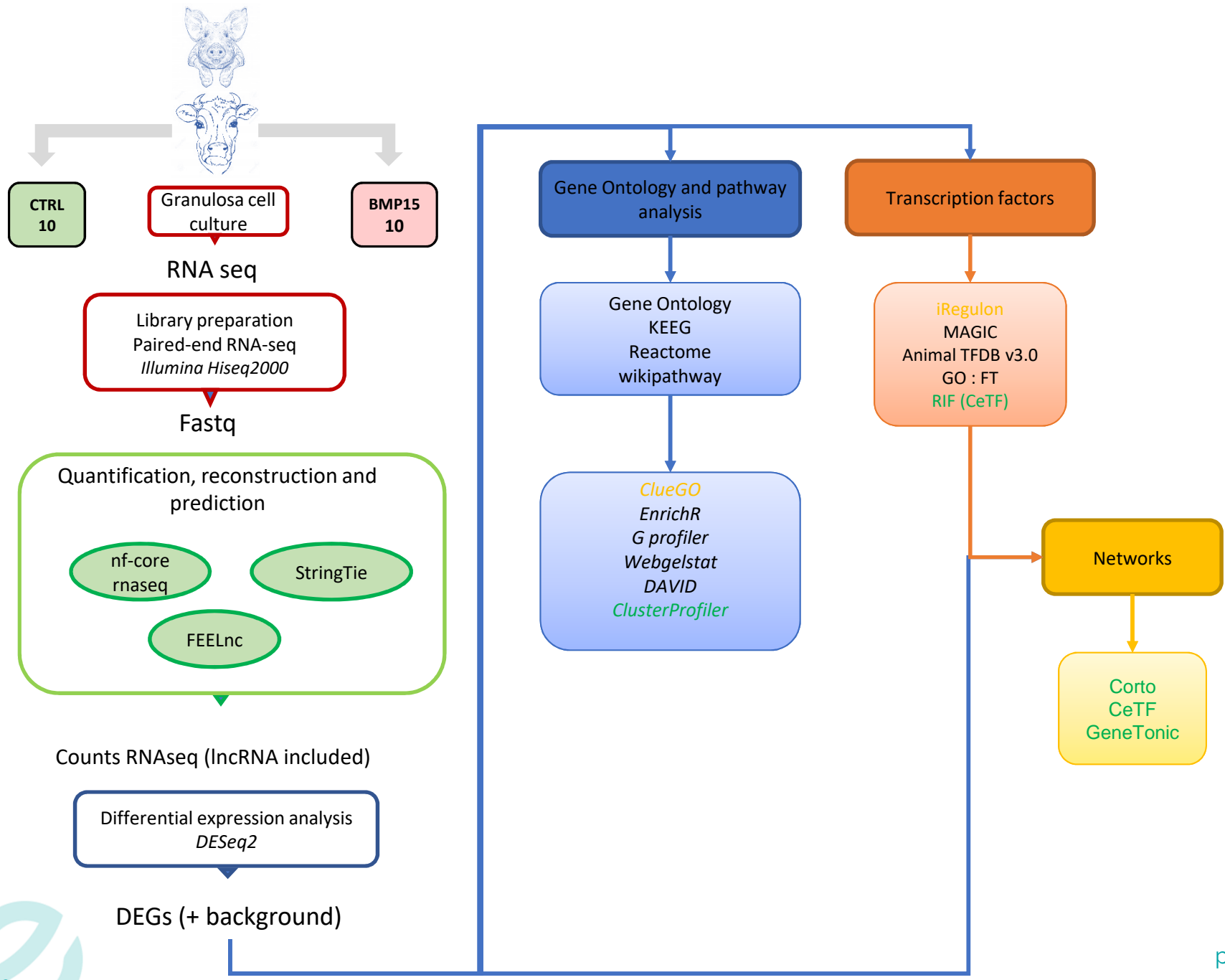Génétique Physiologie et Systèmes d'Elevage

SIGENAE

# ❯ Context

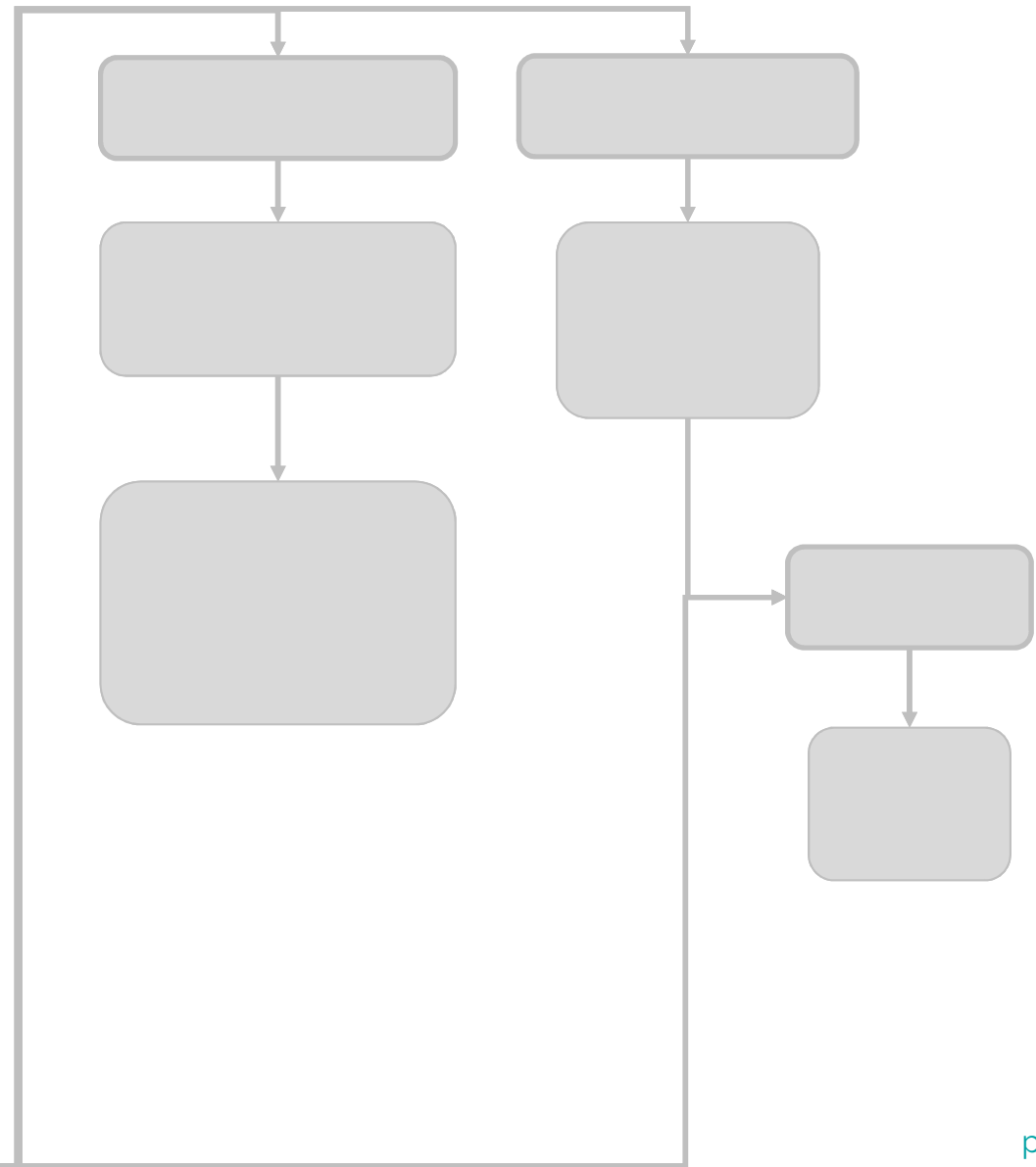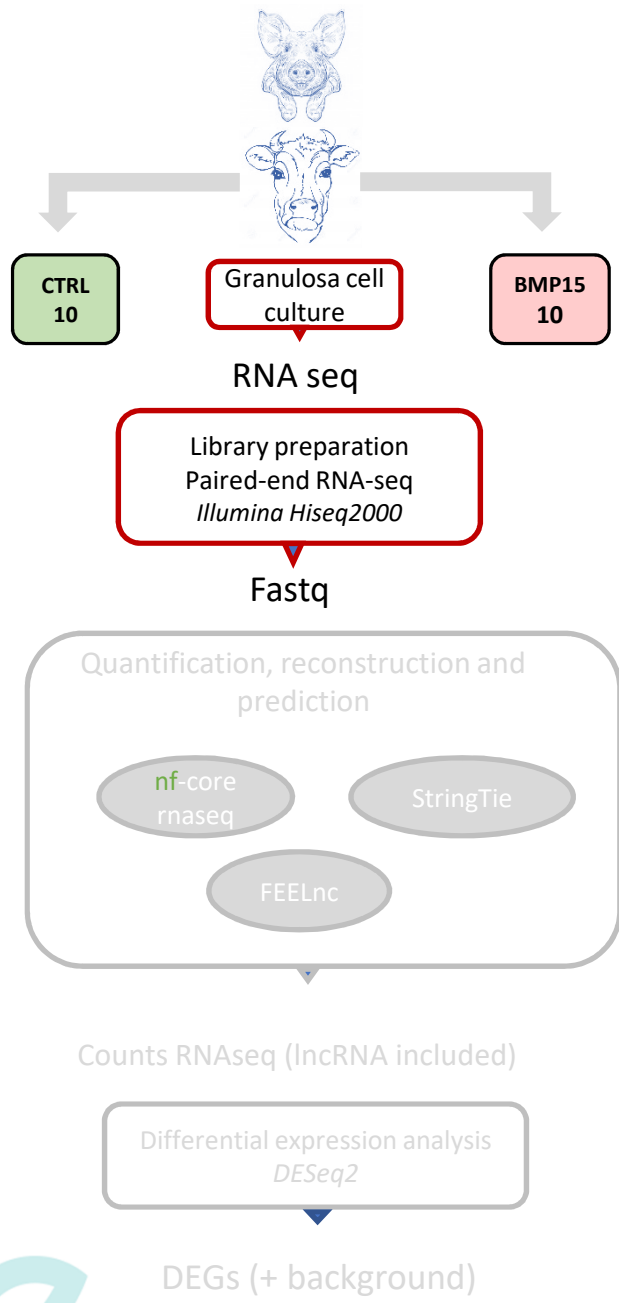Folliculogenesis is important for the development and maintenance of fertility

Ovarian microcosmos.



- BMP15 protein
- Secreted ligand of the TGF-β superfamily
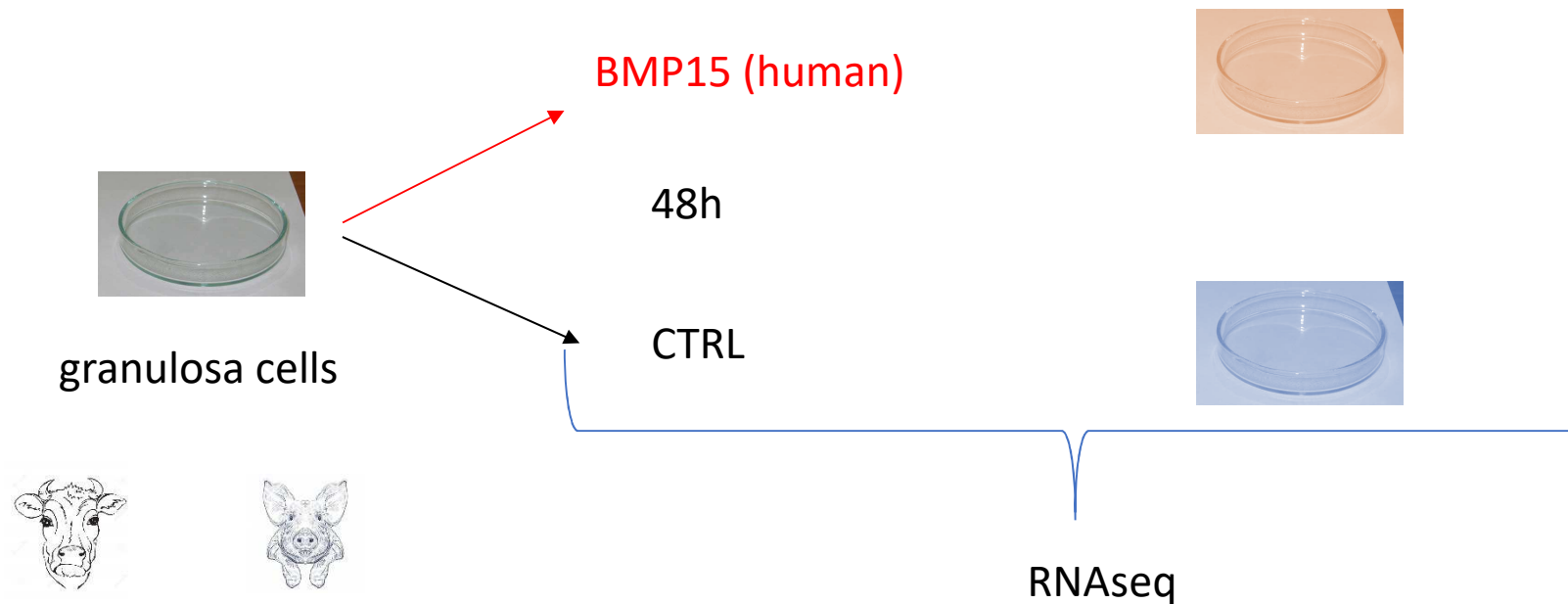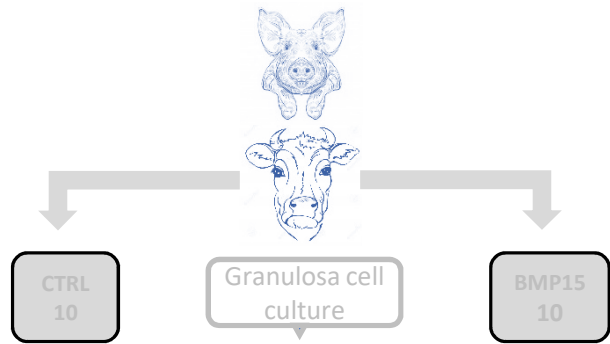- Involved in folliculogenesis
- Contributed to paracrine dialog

CTRL
10

Granulosa cell
culture

BMP15
10

RNA seq

Library preparation
Paired-end RNA-seq
*Illumina Hiseq2000*

Fastq

Quantification, reconstruction and
prediction

nf-core
rnaseq

StringTie

FEELnc

Counts RNAseq (lncRNA included)

Differential expression analysis
*DESeq2*

DEGs (+ background)

Gene Ontology and pathway
analysis

Gene Ontology
KEEG
Reactome
wikipathway

*ClueGO*
*EnrichR*
*G profiler*
*Webgelstat*
*DAVID*
*ClusterProfiler*

Transcription factors

*iRegulon*
MAGIC
Animal TFDB v3.0
GO : FT
RIF (CeTF)

Networks

Corto
CeTF
GeneTonic

12/05/2022

Web ; Cytoscape ; R package

CTRL
10

Granulosa cell
culture

BMP15
10

RNA seq

Library preparation
Paired-end RNA-seq
*Illumina Hiseq2000*

Fastq

Quantification, reconstruction and
prediction

nf-core
rnaseq

StringTie

FEELnc

Counts RNAseq (lncRNA included)

Differential expression analysis
*DESeq2*

DEGs (+ background)

12/05/2022

Web ; Cytoscape ; R package

# Monopoly ANR project (2010 S. Fabre)

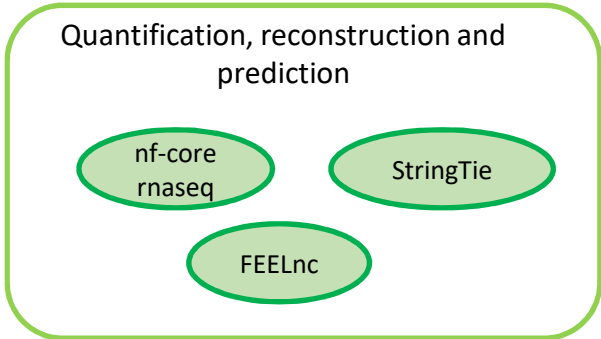Decipher BMP15 response in granulosa cells of cow and sow (mono vs poly-ovulating species)



BMP15 (human)

48h

CTRL

granulosa cells

RNAseq

20 samples by species (10 CTRL vs 10 BMP15)

40 sequencing lanes (0,94<Pearson Correlation<0,98)

INRAe

CTRL
10

Granulosa cell culture

BMP15
10

RNA seq

Library preparation
Paired-end RNA-seq
*Illumina Hiseq2000*

Fastq

Quantification, reconstruction and prediction

nf-core rnaseq

StringTie

FEELnc

Counts RNAseq (lncRNA included)

Differential expression analysis
*DESeq2*

DEGs (+ background)

12/05/2022

Web ; Cytoscape ; R package

Join nf-core

# nf-core 🍎

A community effort to collect a curated set of analysis pipelines built using Nextflow.

**VIEW PIPELINES**

| Search | Search |

### For facilities

Highly optimised pipelines with excellent reporting. Validated releases ensure reproducibility.

### For users

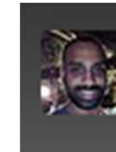Portable, documented and easy to use workflows. Pipelines that you can trust.

### For developers

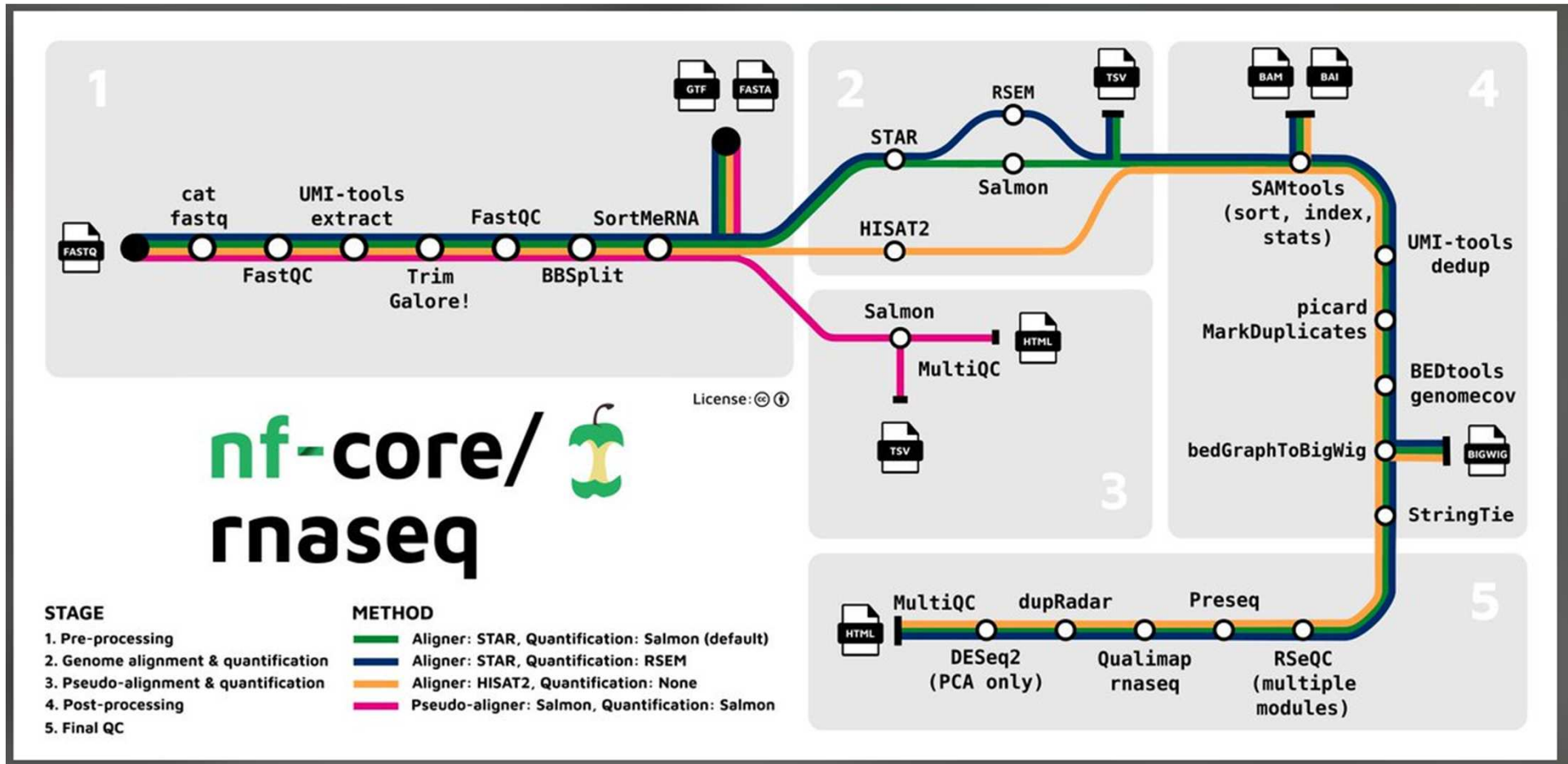Companion templates and tools help to validate your code and simplify common tasks.

nf-core is published in Nature Biotechnology!
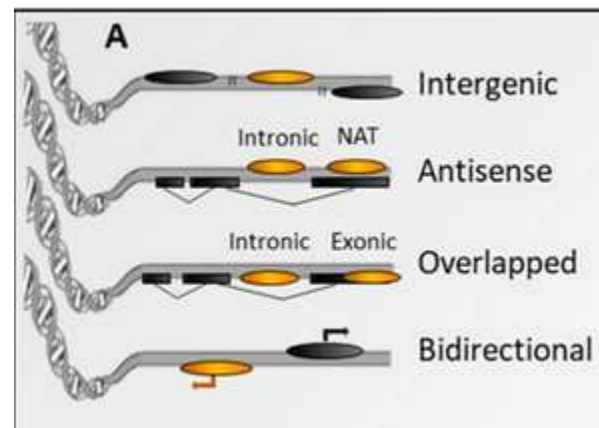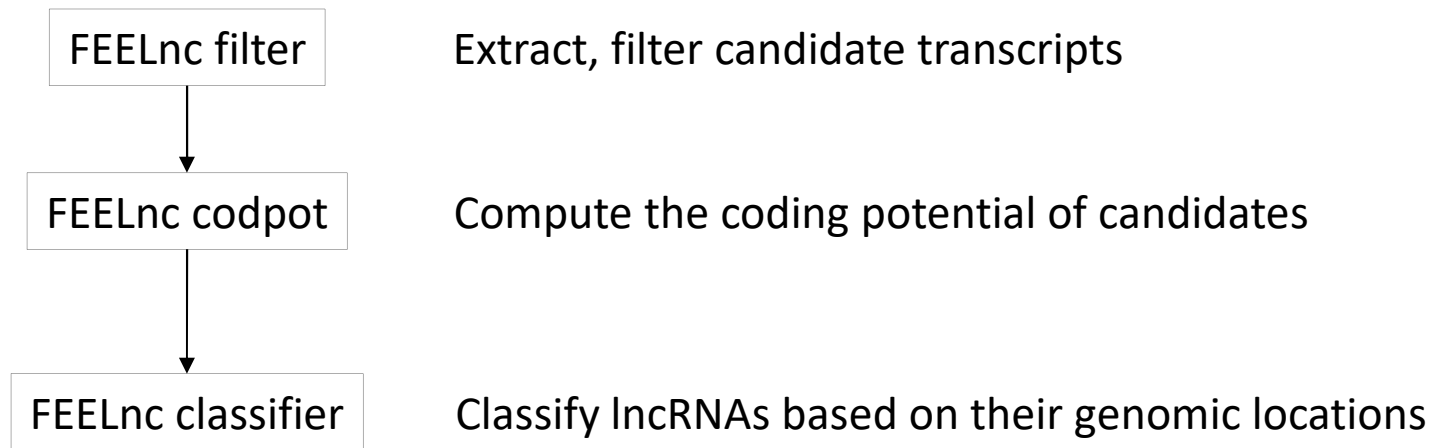*Nat Biotechnol* **38**, 276–278 (2020).

INRAⓔ

12/05/2022

# Pipeline nf-core rnaseq

nf-core 🍎

hpatel



MONOPOLY : BLUE method

INRAe

# Feelnc (prediction only !!)

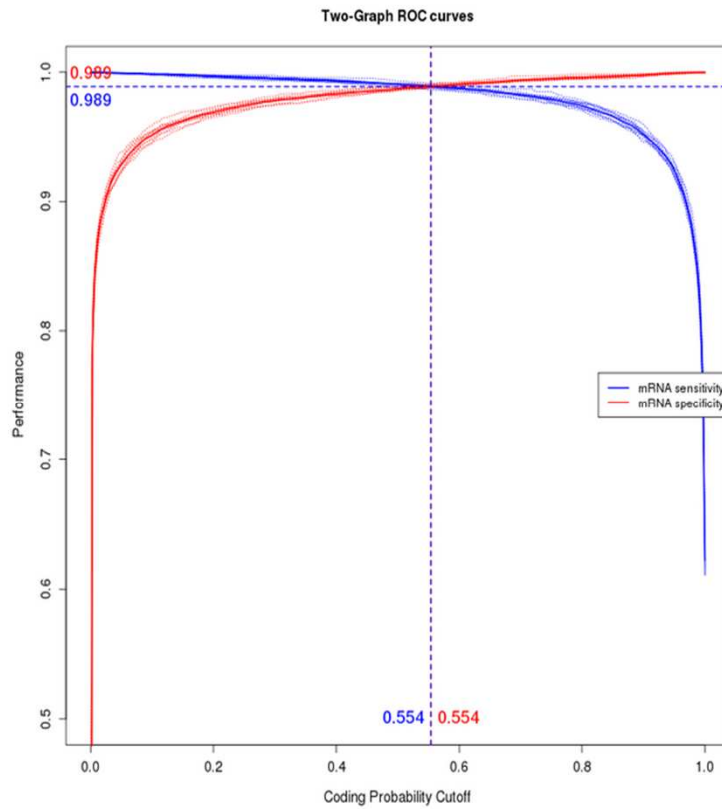| | |
|---|---|
| FEELnc filter | Extract, filter candidate transcripts |
| ↓ | |
| FEELnc codpot | Compute the coding potential of candidates |
| ↓ | |
| FEELnc classifier | Classify lncRNAs based on their genomic locations |



Salviano-Silva et al., 2018, *Non-coding RNA*

INRAe

FEELnc Wucher et al., 2017, *NAR*   p. 9

# Results FEELnc



intergenic

shuffle

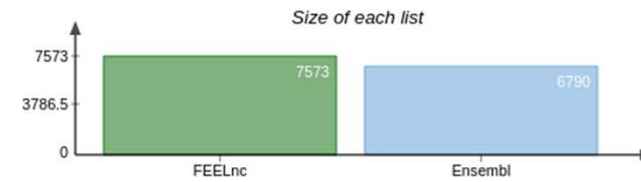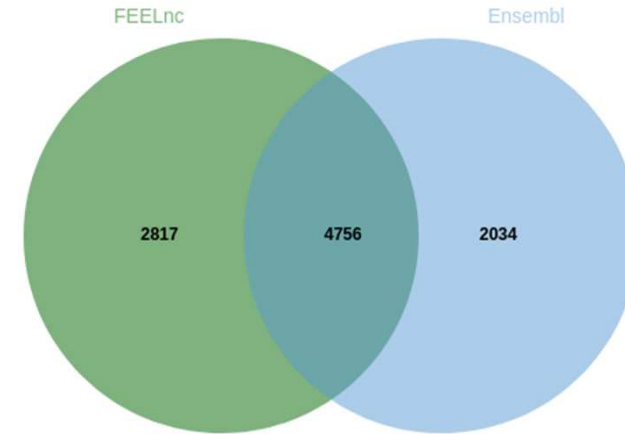lncRNAs      mRNAs

Venn diagram comparing the lncRNAs already known by ENSEMBL and those predicted by FEELnc

12/05/2022

INRAe

CTRL
10

Granulosa cell
culture

BMP15
10

RNA seq

Library preparation
Paired-end RNA-seq
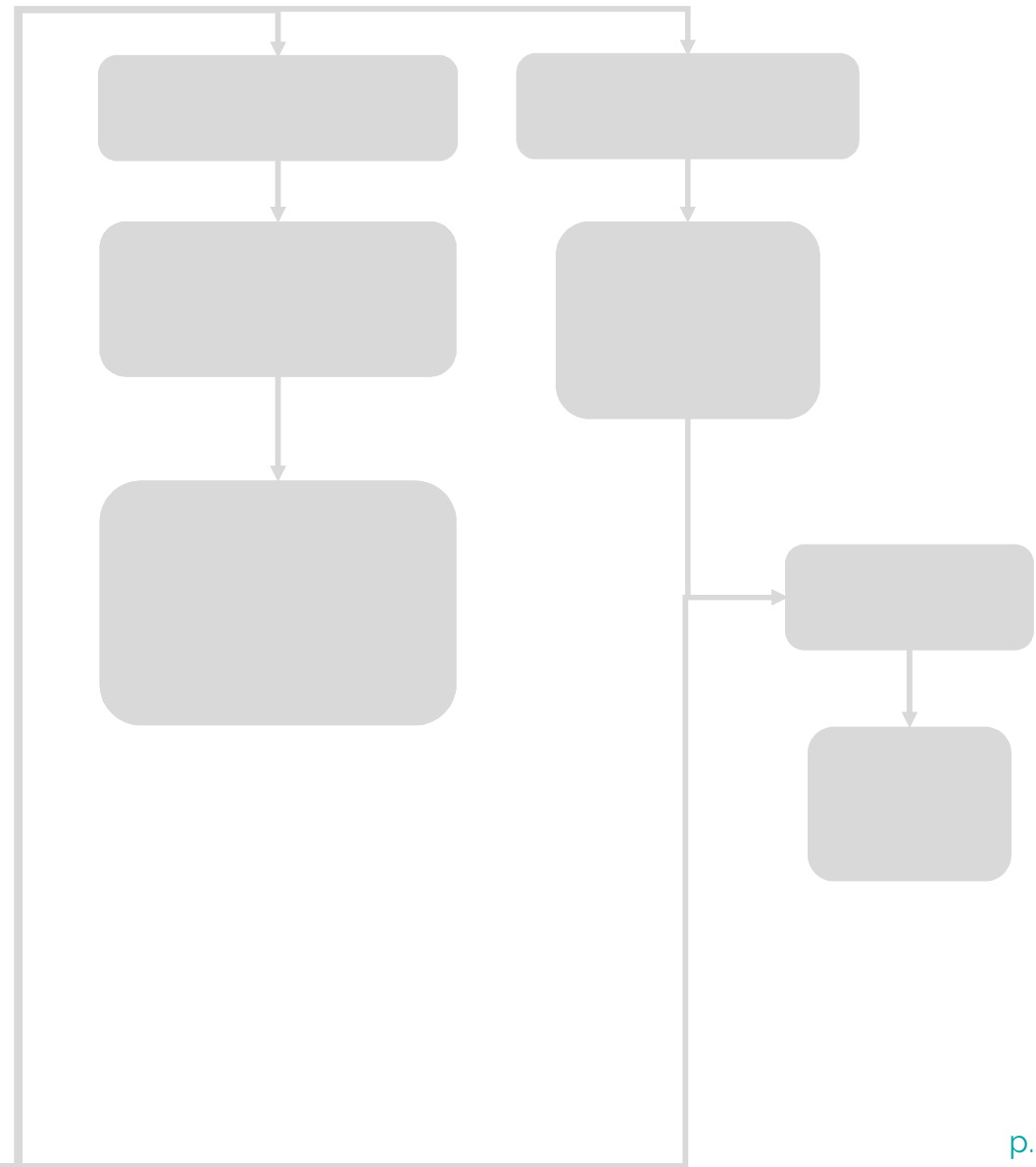*Illumina Hiseq2000*

Fastq

Quantification, reconstruction and
prediction

nf-core
rnaseq

StringTie

FEELnc

Counts RNAseq (lncRNA included)

Differential expression analysis
DESeq2

DEGs (+ background)

Web ; Cytoscape ; R package

# PCA : paired samples

# Differential analysis

DESeq2 Shrunken log2FC, paired samples

Apeglm package

# Identification of DEGs

15 164 expressed genes    (rowsum(counts(dds)>1)>10    16 135 expressed genes

4 521 DEGs    padj<0,05    7 909 DEGs

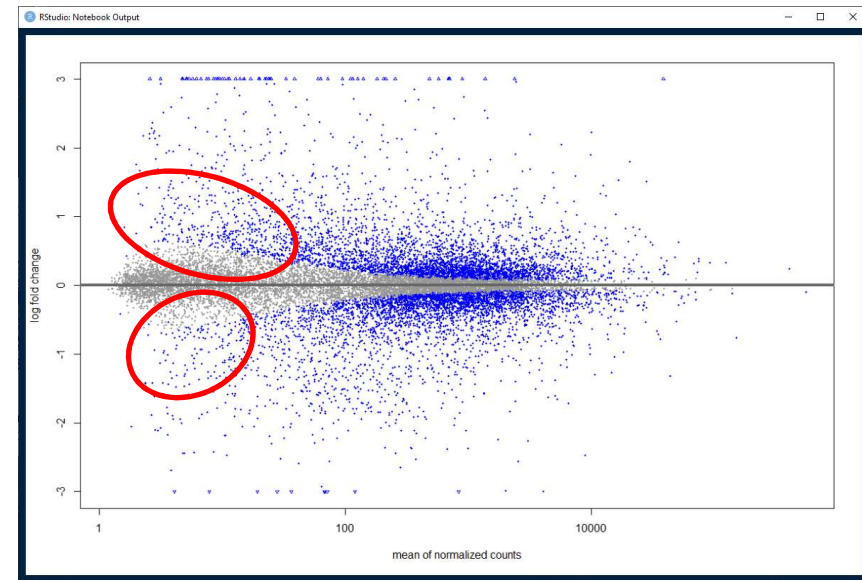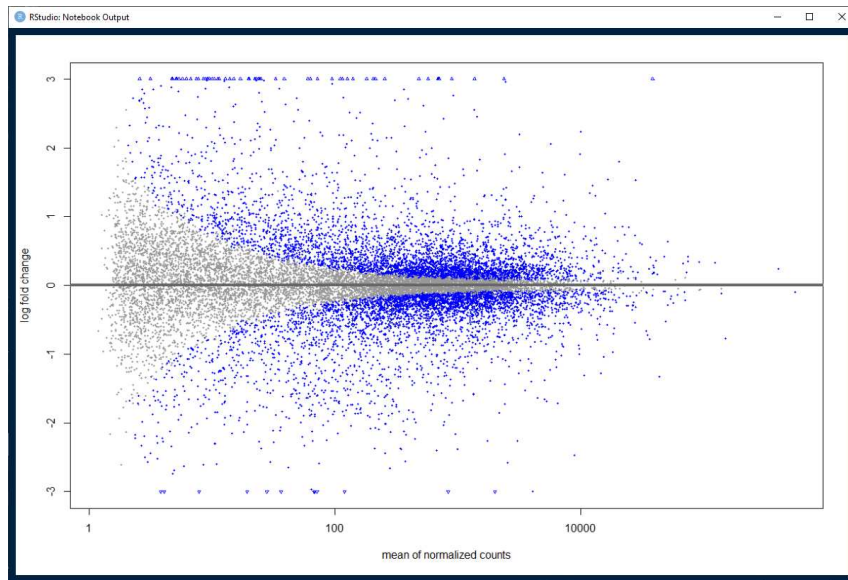240 up ; 203 down    padj<0,05 et |Log2FC|>1    487 up ; 521 down

443 DEGs    1 008 DEGs

INRAe

Granulosa cell culture

RNA seq

Library preparation
Paired-end RNA-seq
*Illumina Hiseq2000*

Fastq

Quantification, reconstruction and prediction

nf-core rnaseq

StringTie

FEELnc

Counts RNAseq (lncRNA included)

Differential expression analysis
*DESeq2*

DEGs (+ background)

Gene Ontology and pathway analysis

Gene Ontology
KEGG
Reactome
wikipathway

*EnrichR*
*g:profiler*
*Webgelstat*
*DAVID*
*ClueGO*
*ClusterProfiler*
*GeneTonic*
*topGO*
*GOstats*

12/05/2022

Web ; Cytoscape ; R package

# GeneOntology – functional enrichment

ORA : Over representation analysis (input = DEGs), is based on a hypergeometric test (Fisher's exact test) (Boyle et al., 2004).

GSEA : Gene Set Enrichment Analysis (All expressed genes, ranked), is based on an enrichment score (Subramanian et al., 2005)

*DAVID*
*EnrichR*
*g:profiler*
*Webgelstat*
*ClueGO*
*GeneTonic*
*topGO*
*Gostats*
*ClusterProfiler*

DAVID : https://david.ncifcrf.gov/

Enrichr: https://amp.pharm.mssm.edu/Enrichr/

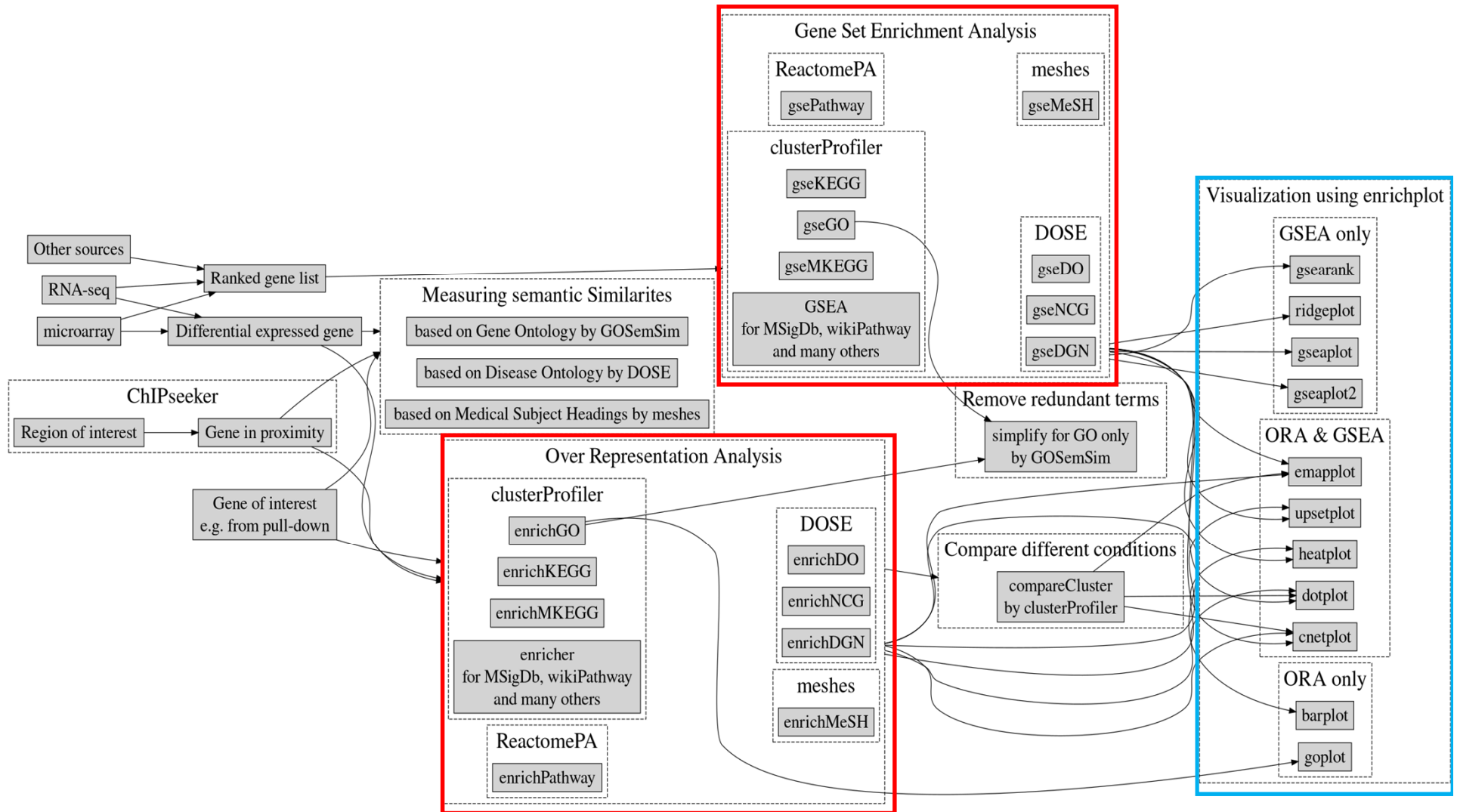g:profiler: https://biit.cs.ut.ee/gprofiler/

Webgelstat : http://www.webgestalt.org/

TopGO and fgsea R packages

Package *GeneTonic* (Marini et al., 2021)

Package *ClusterProfiler* (Wu et al, 2021)

**INRA℮**
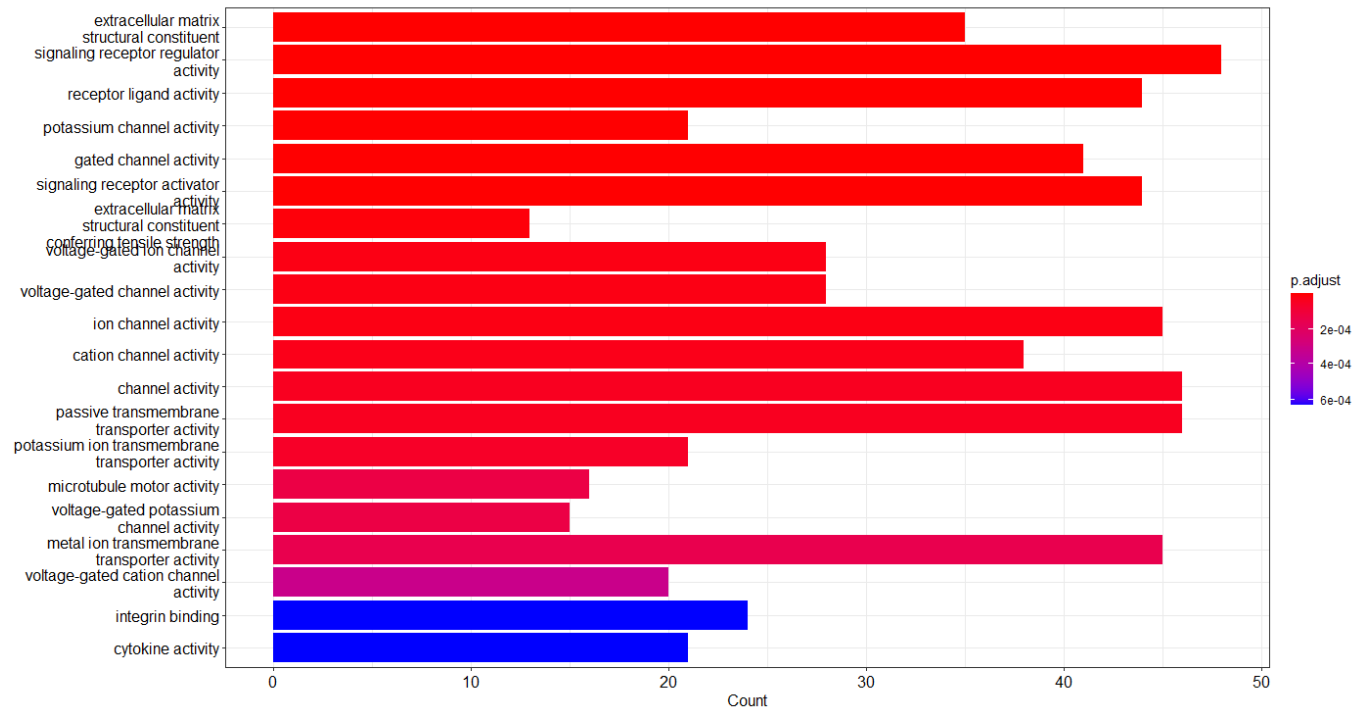
# GO plot
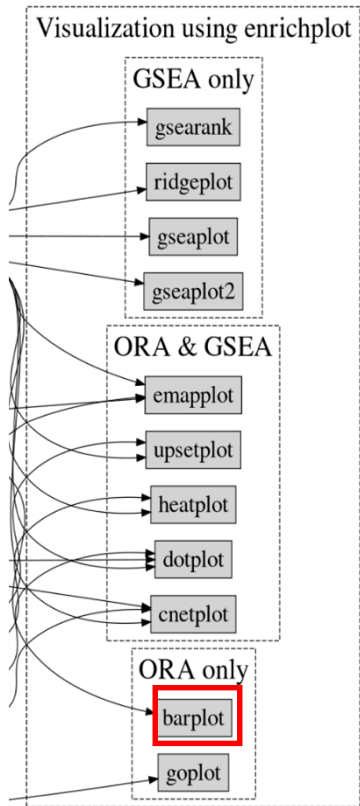
# Bar plot

# Dot plot

# Cnet plot (circular)

# Cnet plot

12/05/2022

# Compare Cluster



KEGG pathway

INRAе

# Compare cluster



BP GSEA

INRAe

Inhibition cell cycle and tissue development, inhibition cell survival; increased apoptosis

Granulosa cell culture

RNA seq

Library preparation
Paired-end RNA-seq
*Illumina Hiseq2000*

Fastq

Quantification, reconstruction and prediction

nf-core rnaseq

StringTie

FEELnc

Counts RNAseq (lncRNA included)

Differential expression analysis
*DESeq2*

DEGs (+ background)

Transcription factors

iRegulon
MAGIC
Animal TFDB v3.0
GO : FT
RIF (CeTF)

Web ; Cytoscape ; R package

**❯ Transcriptions factors :**

Raw view :

Iregulon : Database de 10 000 FTs + 1 120 CHiP-seq datasets (Janky et al., 2014)

iRegulon
MAGIC
Animal TFDB v3.0
GO : FT
RIF (CeTF)

MAGIC  Mining Algorithm for GenetIc Controllers: 2312 CHiP-seq tracks 684 FT in 588 cell lines (ENCODE data) (Roopra 2020 *Plos Comput Biol) (cluster)*

**Gene Regulatory Networks (GRN) →** Detailed view

*Bos taurus: 1396 TFs and 935 TF Cofactors*

*Sus scrofa: 1490 TFs and 937 TF Cofactors*

Gene Ontology : DNA-binding transcription factor activity

Granulosa cell
culture

RNA seq

Library preparation
Paired-end RNA-seq
*Illumina Hiseq2000*

Fastq

Quantification, reconstruction and
prediction

Counts RNAseq (lncRNA included)

Differential expression analysis
*DESeq2*

DEGs (+ background)

Transcription factors

Networks

Corto
CeTF
GeneTonic

Web ; Cytoscape ; R package

## Gene Regulatory Networks reconstruction

Several inference algorithms

- Information theory (co-expression)

- Boolean networks

- Differential equations models

- Bayesian models

- Neural models

## PCIT approach : Partial Correlation with Information Theory

For every trio of genes in x, y and z

1) The three first-order partial correlation coefficients are computed -> the strength of the linear relationship between x and y that is independent of (uncorrelated with) z.

   Obtention d'un seuil « local » pour capturer les associations significatives

2) Data Processing Inequality (DPI) or theorem of Information. Theory which states that 'no clever manipulation of the data can improve the inference that can be made from the data' (Cover and Thomas 2012)

Packages R CeTF, Corto



Gene co-expression → Module formation

Expression / Samples

Reverter et Chan, 2008 *Bioinformatics*

INRAⒺ

# Package CeTF : identification of crucial FTs

RIF Regulatory Information Factors algorithm (Reverter et al., 2010 *Bioinformatics*) → MSTN

Counts normalisés de DEGs
**Cibles**

Liste de FTs
**Régulateurs**

iRegulon
MAGIC
Animal TFDB v3.0
GO : FT
**RIF (CeTF)**

Calcul de la corrélation de la co-expression de chaque TF
avec les DEGs pour les 2 conditions

RIF1: high score to TFs that are most differentially co-expressed, highly abundant.
High RIF1 = FTs « constants »

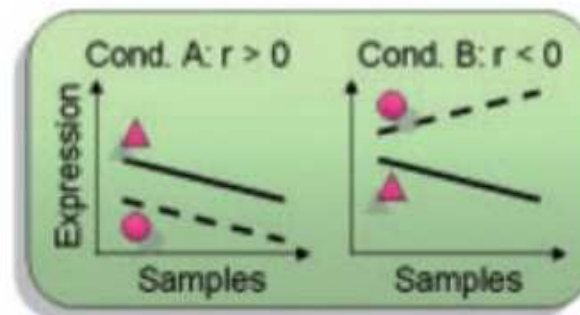RIF2: high score to TFs whose expression can predict better the abundance of DEGs



**INRAe**

Reverter et al., 2010, *bioinformatics*

# RIF Results

| TF | avgexpr | RIF1 | RIF2 |
|---|---|---|---|
| RARG | 4.4308315 | -2.362989 | -0.4269886 |
| STAT2 | 5.5116914 | 2.347202 | 0.3046848 |
| CIB1 | 6.4784370 | -2.079084 | 0.4312936 |
| RORA | 2.8853325 | -1.983128 | -0.2532182 |
| MALT1 | 5.3576710 | 3.138360 | -0.5968538 |
| CREB3 | 5.9433637 | -2.617405 | 0.3039166 |
| HEY1 | 0.6291504 | 2.012688 | 0.4760702 |
| COPS5 | 5.9009044 | 2.583354 | -1.1730226 |
| ATF6 | 8.2596171 | 2.067810 | -0.6565662 |
| GATA5 | 0.5430692 | 2.918902 | 0.4725577 |

26 TFs RIF1

| TF | avgexpr | RIF1 | RIF2 |
|---|---|---|---|
| ELK4 | 5.9506900 | 5.0683310 | -3.121338 |
| DLX2 | 1.4871928 | 0.8961663 | -2.845057 |
| ZNF35 | 3.4982397 | 1.6921604 | -2.709549 |
| TICAM1 | 4.4586161 | 0.8491830 | -2.689220 |
| FEZF2 | 0.2689711 | -0.5681378 | -2.648282 |
| UFL1 | 6.2945395 | 1.5768689 | -2.400373 |
| SP7 | 2.4673947 | 2.3311699 | -2.390159 |
| KLF13 | 7.9588770 | -0.3132815 | -2.286667 |
| ZHX3 | 5.3036326 | -0.5644726 | -2.175302 |
| RARA | 6.1011870 | 2.5207260 | -2.148701 |

29 TFs RIF2

Zscore 1,96 pval<0,05

| | TF | avgexpr | RIF1 | RIF2 |
|---|---|---|---|---|
| 1 | ZNF423 | 5.261707 | 3.43867 | -1.28137 |

| | TF | avgexpr | RIF1 | RIF2 |
|---|---|---|---|---|
| 1 | NTRK1 | 3.732381 | 1.2677245 | 1.807717 |
| 2 | LRRC7 | 2.767650 | 0.0279897 | -1.881063 |
| 3 | ARNTL2 | 6.167500 | -0.4328875 | 1.866842 |

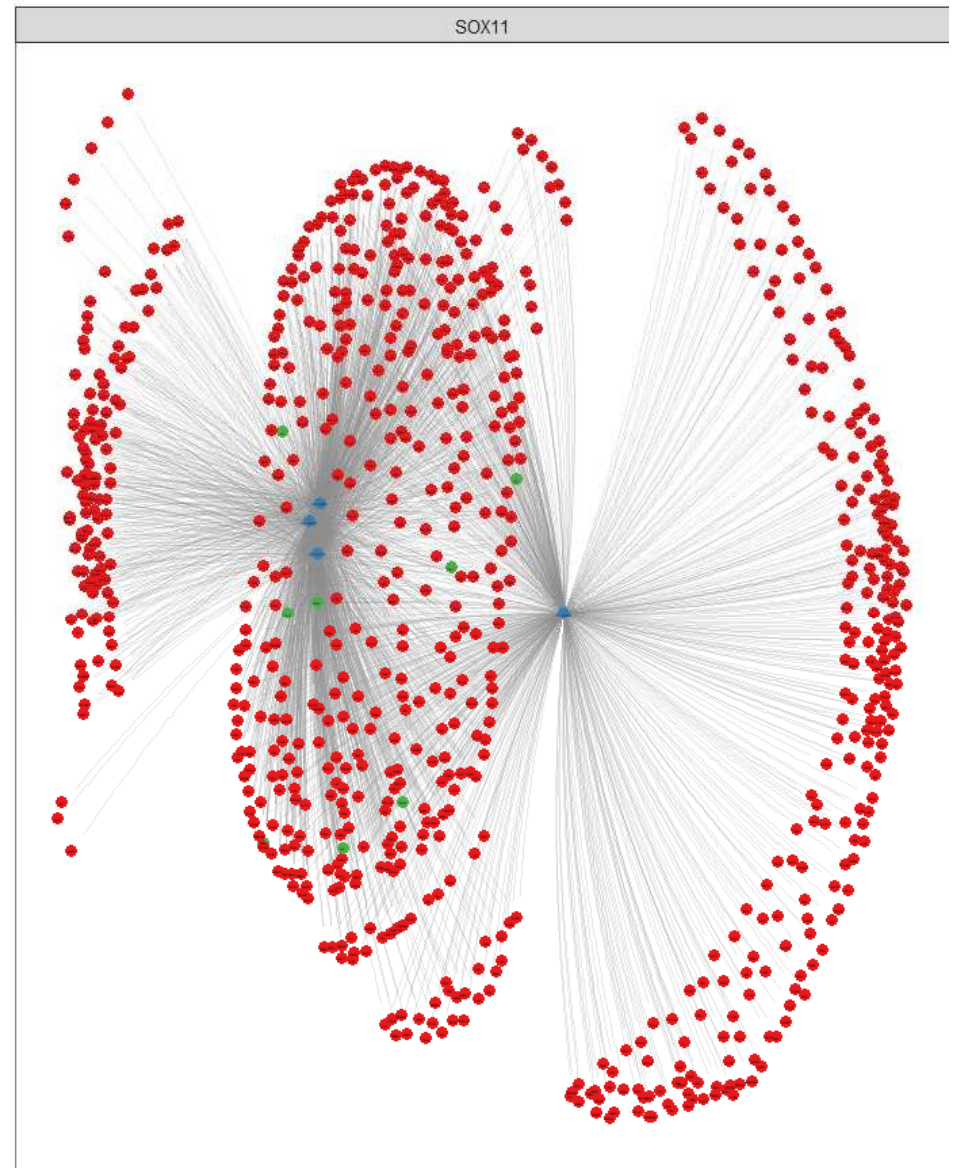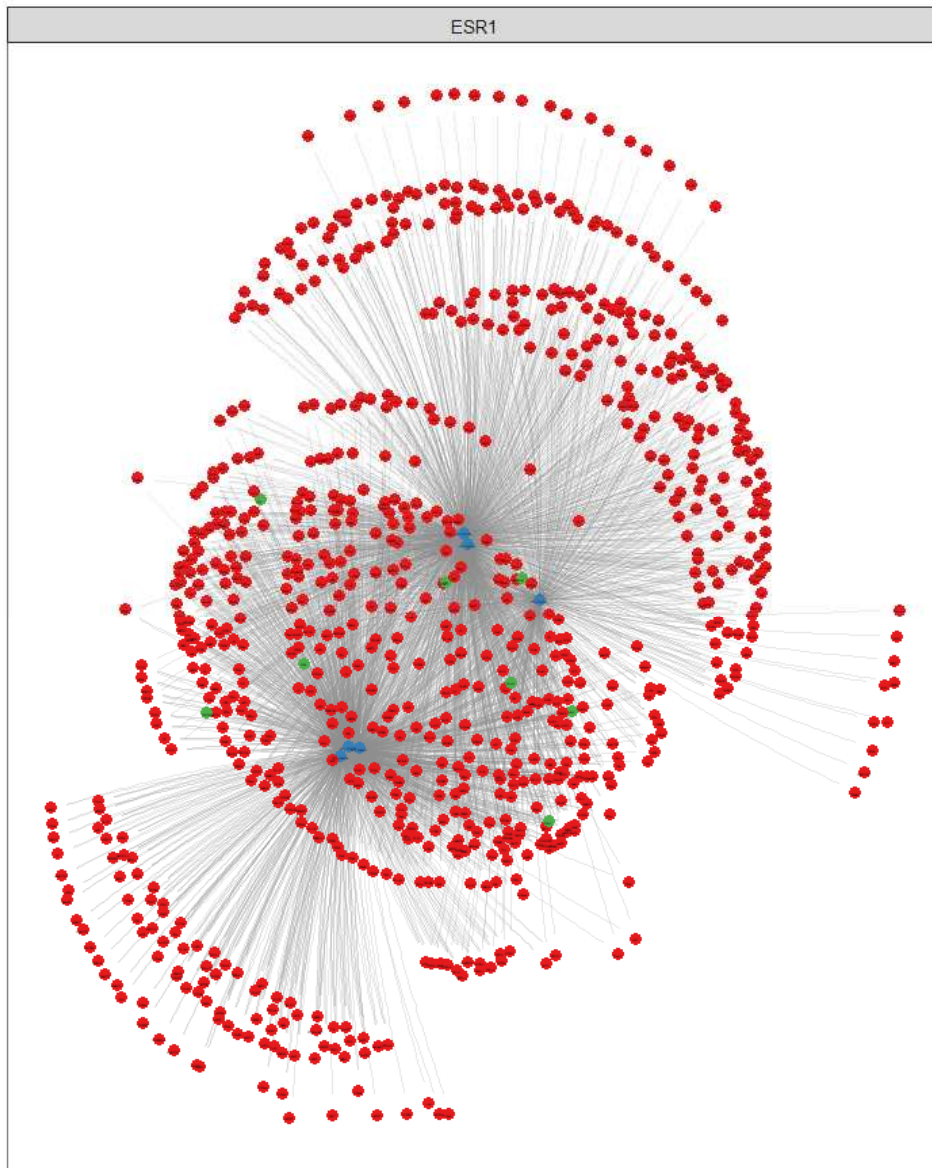Zscore 1,65 pval<0,10

Packages R CeTF

INRAe

# FT ESR1 and its targets



Smear Plot for ESR1 and its targets

ESR1     SOX11

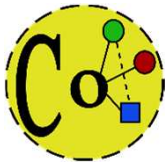Type    ● gene    ● pathway    ● TF

## Package CeTF : Pros and cons

🙂 RIF algorithm in R

🙁 2 conditions, little used, maintainer not responding, graphic interface.

Corto : (Correlation Tool): an R package to generate correlation-based DPI networks. Mercatelli et al., 2020 Bioinformatics
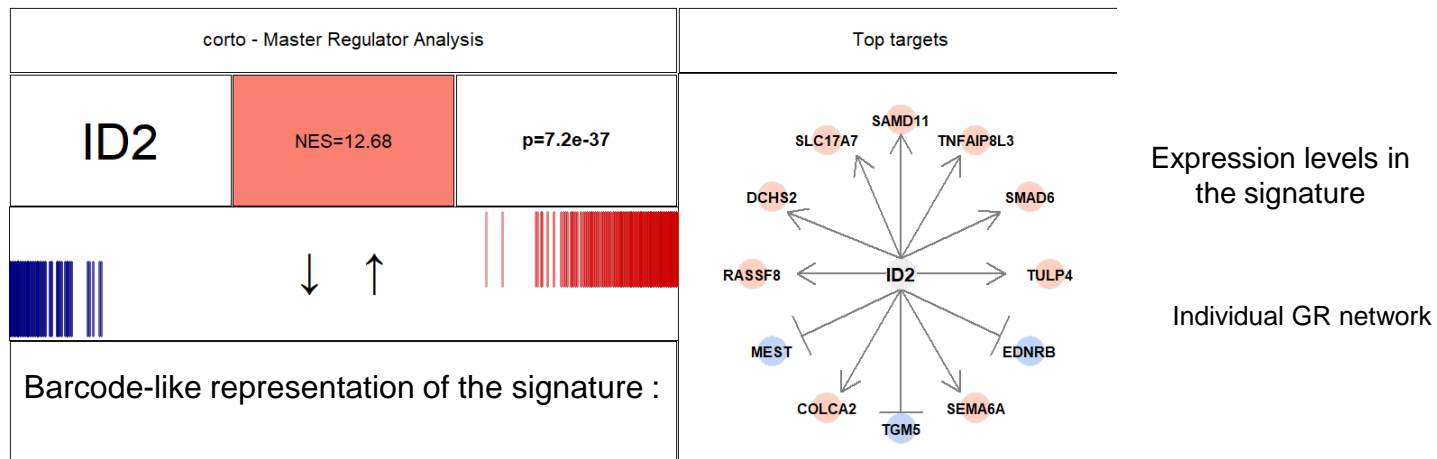
😃 Maintainer ++

INRAぇ

# Corto

1) Identification of TFs expressed in the dataset

```
 [1] "FOXM1"   "E2F7"    "TCF19"   "E2F1"    "MYBL2"   "E2F8"    "ZNF367"
 [8] "ZNF704"  "SMAD7"   "PRRX2"   "ATOH8"   "DLX2"    "TBX2"    "SP7"
[15] "DRGX"    "DLX5"    "ID3"     "NR5A2"   "HIVEP1"  "NFIB"    "ZNF182"
[22] "TEAD4"   "RTKN2"   "SFRP5"   "GATA5"   "ID2"     "SOX11"   "PBX4"
[29] "GATA4"   "ID1"     "E2F2"    "MYOCD"   "RGCC"    "ZFHX4"   "DLX3"
[36] "VDR"     "ZBTB16"  "DLX6"    "NR4A1"   "ADCY8"   "HAND2"   "SNAI2"
[43] "SHOX2"   "ATF3"    "TRAIP"   "ZNF423"  "WNT5A"   "TCF7"    "ESR2"
[50] "EN1"     "SOX9"    "MKX"     "THRB"    "MYCN"    "MSX2"    "NKX2-8"
>
```

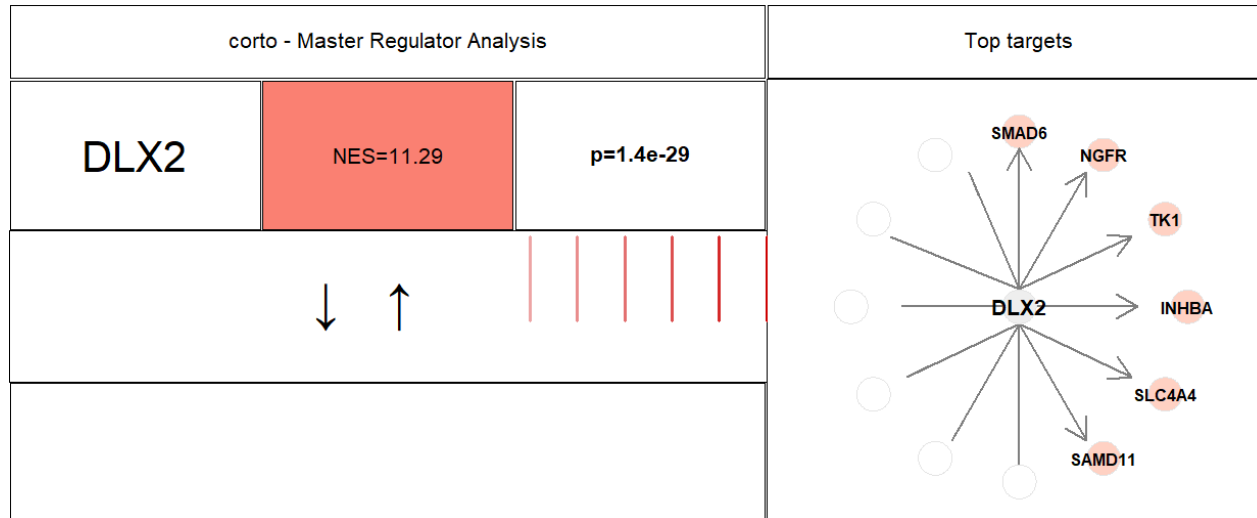2) Gene network inference (optimized pairwise correlation, DPI and bootstrapping)

Master Regulator Analysis : TF networks vs a signature (= 2 gene expression matrices)



Expression levels in the signature

Individual GR network

Barcode-like representation of the signature :

corresponding P value (based on 10 000 permutation tests and signature sample shuffling)

| corto - Master Regulator Analysis | | | Top targets |
|---|---|---|---|
| DLX2 | NES=11.29 | p=1.4e-29 | |

The transparency of each bar is associated to the value of the target in the signature.

INRAe

## ❯ What else ?

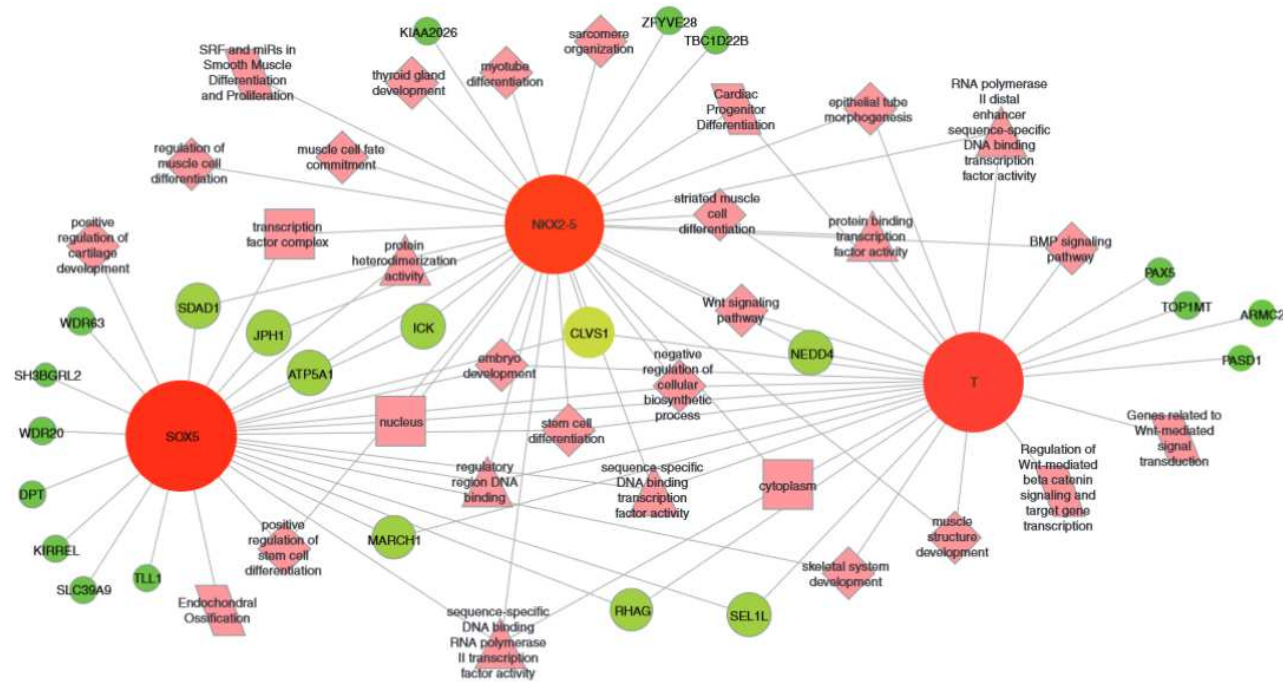- From R to Cytoscape or Gephi :



**Fig. 3.** Transcription factors network. Transcription factors (in red), with related genes (green), pathways (parallelogram) and gene ontology: molecular functions (triangle), biological process (diamond) and cellular component (rectangle).

- Add lncRNAs : Test TAGADA pipeline

## Perspectives : TAGADA pipeline (GeneSwitch project)

Another pipeline using Nextflow for count matrix generation and lncRNA prediction

|  | Nextflow nf-core rnaseq + FEELnc + StringTie | Nextflow TAGADA |
|---|---|---|
| Quantification | RSEM | Stringtie |
| Merge of GTFs by sample | Outside, StringTie Merge | Inside,T-merge |
| transcript quantification | Outside, Stringtie | Inside, Stringtie |
| LncRNA prediction | Outside, FEELnc Choice of the method intra WF | Inside, FEELnc Choice of the methode upstream WF |

+ Filters on transcripts « de novo » are also different.

INRAe

## ❯ What else ?

- GRN with lncRNAs

- Use co-expression networks to identify potential trans-targeted genes.

- Validate the trans target with lncTar (Li et al., 2014, *brief. Bioinform) :*

  Predict lncRNA-RNA interactions based on secondary structure

**INRAe**

## Thanks

TAGADA pipeline

Sarah Djebali

Cervin Guyomar

Cyril Kurylo

Sylvain Foissac