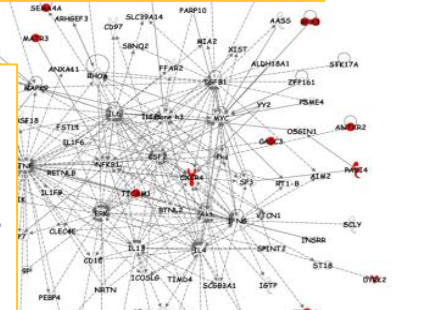
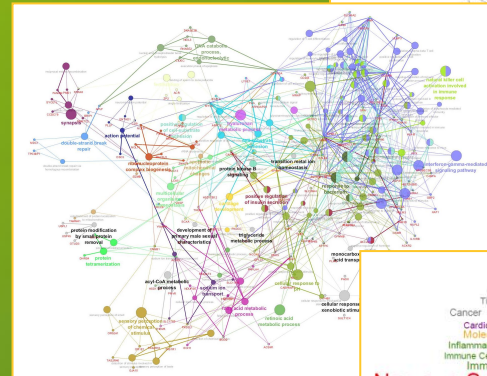
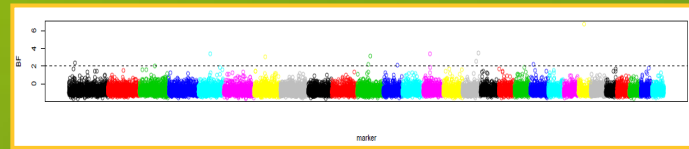


Biological interpretation of selection footprints



- Nervous System Development and Function**
- Cellular Compromise
 - Cellular Movement
 - Cell Death and Survival
 - Cellular Growth and Proliferation
 - Tissue Morphology
 - Cellular Morphology
 - Organ Morphology
 - Carbohydrate Metabolism
 - Inflammatory Disease
 - Tissue Development
 - Antimicrobial Response
 - Embryonic Development
 - Cancer
 - Organismal Functions
 - Call-To-Cell Signaling and Interaction
 - Cardiovascular Disease
 - DNA Replication, Recombination, and Repair
 - Molecular Transport
 - Reproductive System Development and Function
 - Inflammatory Response
 - Cellular Function and Maintenance
 - Cell Morphology
 - Immune Cell Trafficking
 - Cardiovascular System Development and Function
 - Immunological Disease
 - Endocrine System Development and Function
 - Cellular Assembly and Organization
 - Cellular System Development and Function
 - Cell Cycle
 - Small Molecule Biochemistry
 - Auditory and Vestibular System Development and Function
 - Metabolic Disease
 - Hematological System Development and Function
 - Dermatological Diseases and Conditions
 - Organismal Development
 - Hereditary Disorder
 - Organismal Injury and Abnormalities
 - Cellular Assembly and Organization
 - Endocrine System Disorders
 - Gastrointestinal Disease
 - Nucleic Acid Metabolism
 - Hematopoiesis
 - Reproductive System Disease
 - Lipid Metabolism
 - Drug Metabolism
 - Infectious Disease

FLORI Laurence
GABI, INRA, Jouy-en-josas
INTERTRYP, CIRAD, Montpellier



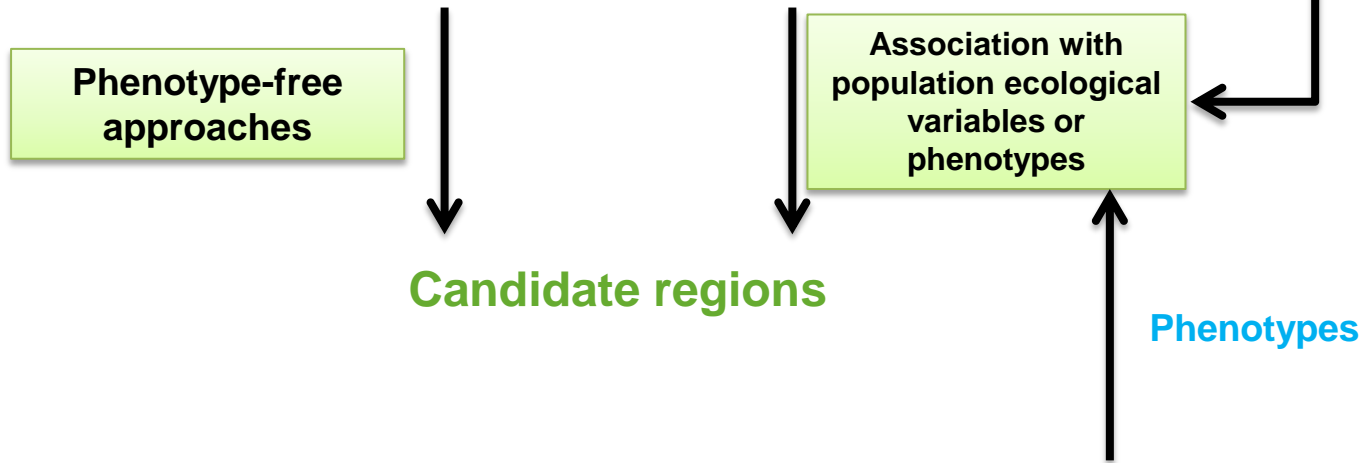
Overview



SNP chip
NGS (Pool and Ind-Seq)

Identification of footprints of
selection in the genome

Ecological variables
e.g. WorldClim DB



Overview

SNP chip
NGS (Pool and Ind-Seq)

Ecological variables
e.g. WorldClim DB

Identification of footprints of selection in the genome

Phenotype-free approaches

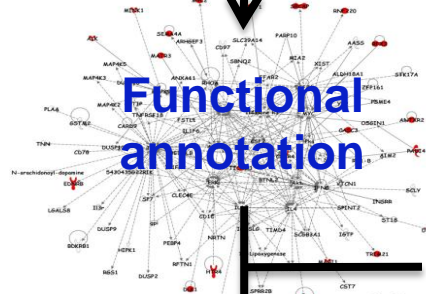
Association with population ecological variables or phenotypes

Candidate regions

Candidate genes

Systems biology Tools
(e.g. IPA, Cytoscape)

Principal selective pressures



Phenotyping RNA-Seq

Reverse ecology

Phenotypes

Whole Genome mutation screening (Seq)

Functional annotation of variants

Candidate causal mutations for adaptive traits



Functional annotation

- **Main goals**
 - Identify the main functions and biological pathways in which the genes are involved
 - Confirm the biological importance of the genes that have been identified as being positively selected
 - Strengthen the credibility of the positive selection model by the development of a sound scenario



Course outline

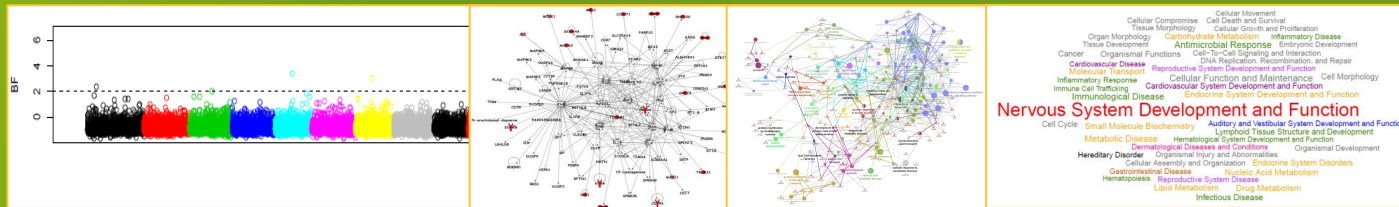
1. From regions under selection to biological interpretation

- Annotation of candidate genes using systems biology tools
 - Functional enrichment analysis: Gene ontology, pathway analysis
 - Gene networks analysis
- Inferring the main selective pressures
- Interpretation and story-telling
- Including phenotypes and environmental covariates
- Identifying candidate mutations and prioritizing candidate genes

2. Interpretation of selection footprints: some examples

- Phenotype-free approaches
 - Manual functional annotation: Senepol cattle breed
 - Functional annotation using systems biology tools: french dairy cattle breeds, West-African cattle breeds, European bison/cattle
- Association with covariates
 - Phenotypes: dairy traits in French cattle breeds

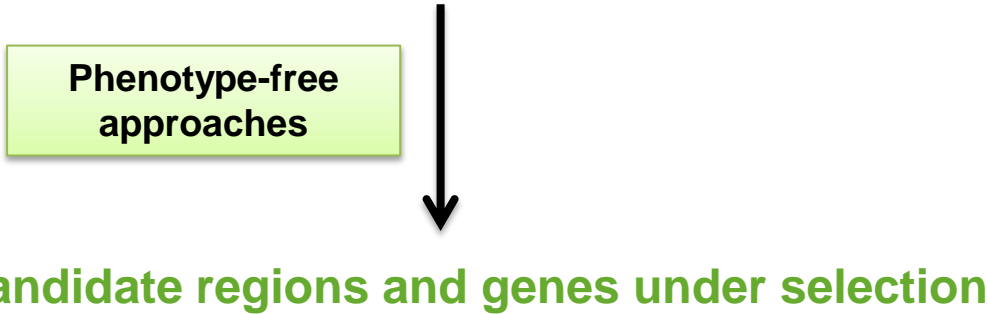
1. From regions under selection to biological interpretation



Interpretation of results obtained with phenotype-free approaches

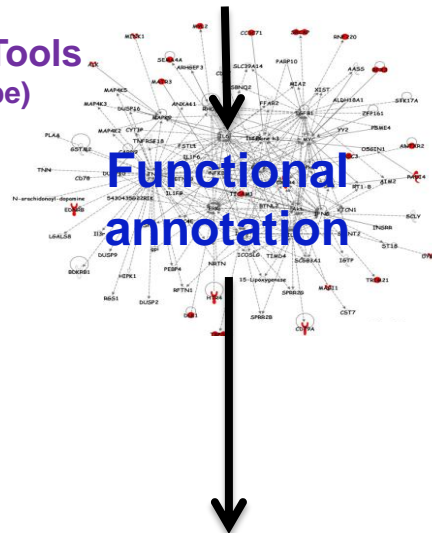
SNP chip
NGS (Pool and Ind-Seq)

Identification of footprints of selection in the genome



Systems biology Tools
(e.g. IPA, Cytoscape)

Principal selective pressures



Candidate causal mutations for adaptive traits

From regions to candidate genes

- **Identification of chromosomal regions under selection**
 - ⇒ **some arbitrary criteria**
 - Sliding windows (e.g. 1Mb, 500kb overlap)
 - Nb of SNPs with score > significant threshold
 - If several overlapping regions under selection: merging regions or choosing region with the highest peak and the highest proportion of significant SNPs
- **Mapping of regions on corresponding genome assembly on ucsc or ensembl (if available)**
- **List of all genes within regions (Refseq and Other species Refseq)**
- **Criteria to identify one or a few candidate genes per region**
 - e.g. proximity to the peak (e.g. 15-100kb from gene boundaries)

From candidate genes to biological interpretation



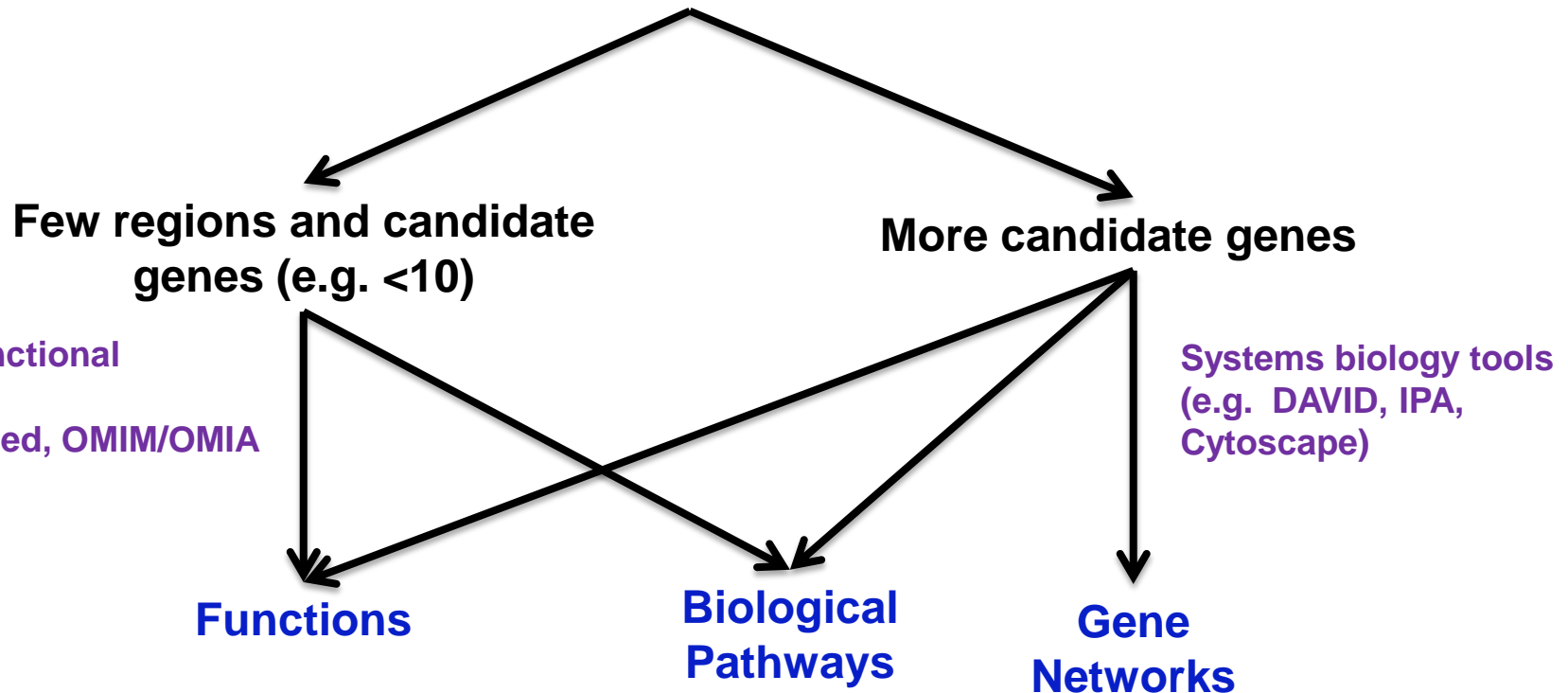
- **List of candidate genes with/without associated score**
 - Short list
 - Long list
- **Misleading gene names**
- **Possibly hundreds of papers describing gene functions**



How to find biological sense?

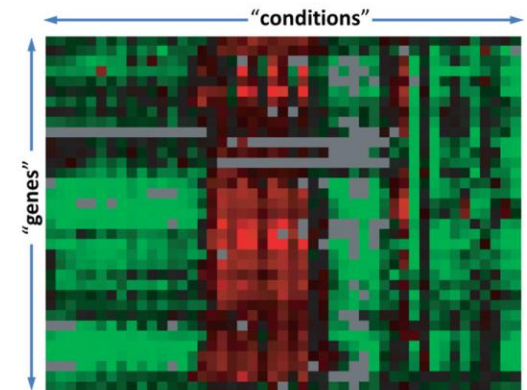
From candidate genes to biological interpretation

List of candidate genes with/without associated score
(one or few per region)



Annotation using systems biology tools

- **Systems biology**
- **Tools used to interpret results of transcriptomic and proteomic experiments**
- **e.g.: list of differentially expressed genes between different experimental conditions, different tissues**



Annotation using systems biology tools:

Functional enrichment analysis

- **Is the list of candidate genes statistically enriched for some functions or some biological pathways?**
- **Needs:**
 - Shared vocabulary: Gene Ontology, Biological pathways
 - Annotation of genes : association between terms and genes or gene products.
 - Statistical tests of enrichment

Annotation using systems biology tools

Functional enrichment analysis: Gene Ontology

- **Gene ontology (GO):** hierarchical vocabulary of terms describing genes and protein
- **Three GO represented by a root ontology term**
 - **cellular component** referring to the place in the cell (e.g. nucleus, ribosome)
 - **biological process** referring to a biological objective (e.g. signal transduction, alpha-glucoside transport)
 - **molecular function** describing activities (e.g. catalytic activity, Toll receptor binding)
- **Developed by a consortium** (geneontology.org)
=> collaborative effort
- **Widely used biological resource**
- **41 species**
- **Generic terms / Some species-specific terms**



Annotation using systems biology tools

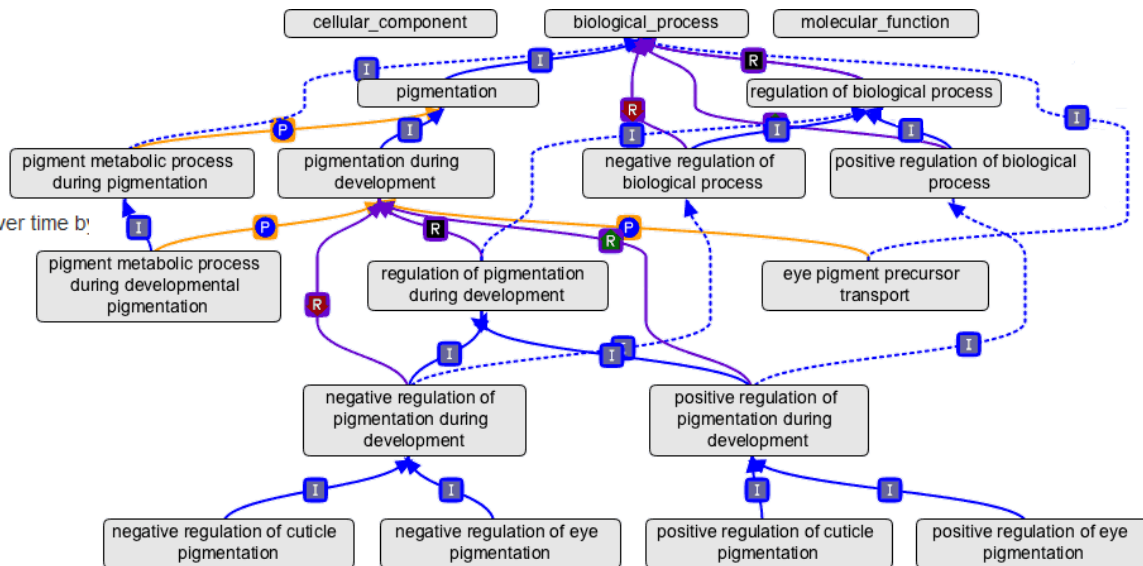
Functional enrichment analysis: Gene Ontology

- **GO structure:**
 - directed acyclic graph
 - types of relationships of child to parent: “is a” or “part of”
 - each term has defined relationships to one or more other terms.

Sample GO Term

The following is a GO term taken from the OBO format file.

- id: GO:0016049
- name: cell growth
- namespace: biological_process
- def: "The process in which a cell irreversibly increases in size over time b"
- subset: goslim_generic
- subset: goslim_plant
- subset: gosubset_prok
- synonym: "cell expansion" RELATED []
- synonym: "cellular growth" EXACT []
- synonym: "growth of cell" EXACT []
- is_a: GO:0009987 ! cellular process
- is_a: GO:0040007 ! growth
- relationship: part_of GO:0008361 ! regulation of cell size



geneontology.org

Annotation using systems biology tools

Functional enrichment analysis: Gene Ontology

- However some (commercial) tools developed their own ontology

e.g.: Ingenuity Pathway Analysis (IPA) => IPA ontology

three main categories of function

- Diseases and disorders
- Molecular and cellular functions
- Physiological system development and function

INGENUITY[®]
PATHWAY ANALYSIS

Annotation using systems biology tools

Functional enrichment analysis: Biological pathway

- **Biological pathway**
 - Biochemical engines responsible for the transduction of signals into output responses
 - A series of actions among molecules in a cell that leads to a certain product or a change in the cell.

- **Pathway building**
 - Process of identifying and integrating the molecules entities, interactions, and associated annotations
 - Contribute to the knowledgebase.
 - Can have either a data-driven or a knowledge-driven objective

Annotation using systems biology tools

Functional enrichment analysis: Biological pathway

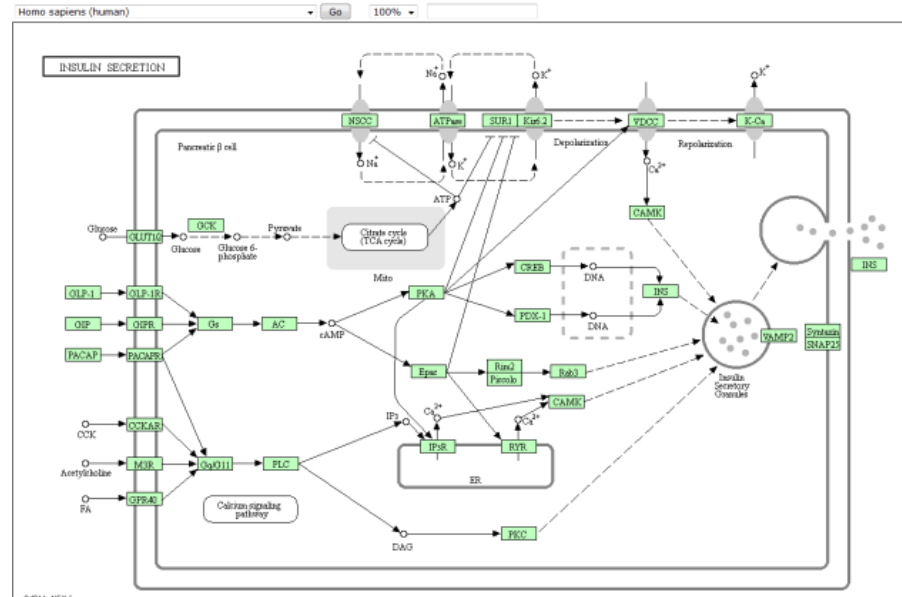
- **Pathway databases**
 - e.g. consortium effort



- Integration of several data sources
 - E.g.



- commercial database





Functional enrichment analysis

Statistical tests

- **Which functional category or biological pathway is more prevalent in the gene list than expected by chance?**
- **Tests**
 - Fisher's exact test (one-tailed, right)
 - Gene Set Enrichment Analysis
 - Correction for multiple testing

Functional enrichment analysis

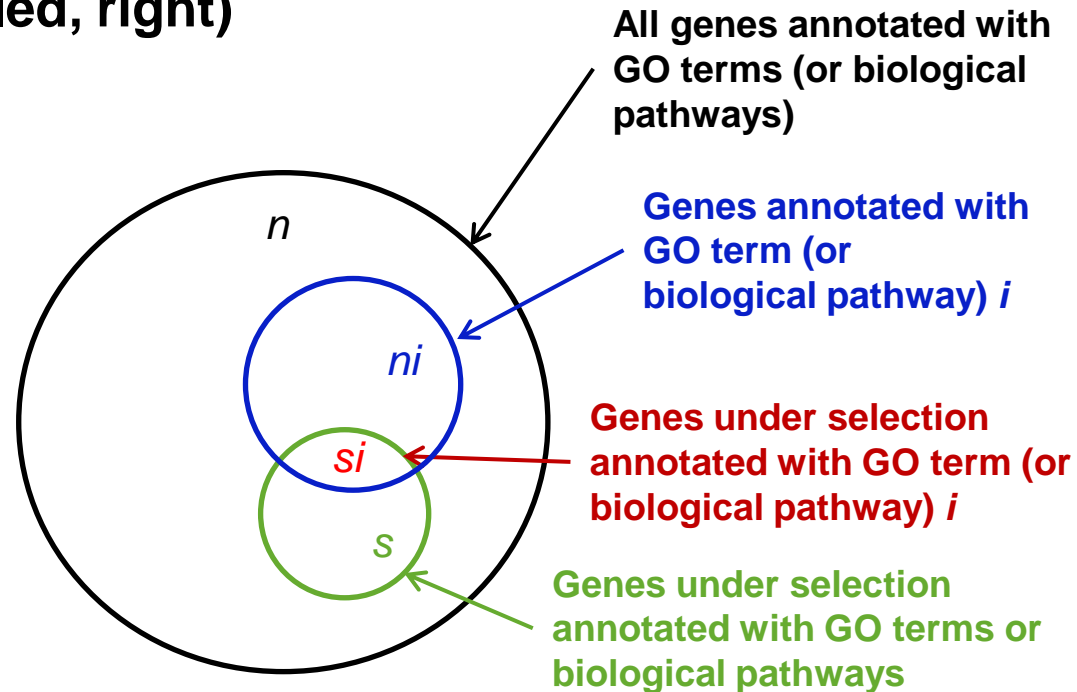
Statistical tests

- Fisher's exact test (one-tailed, right)

2x2 contingency table

	User genes	All genes tested or genome
In pathway	s_i	ni
Not in pathway	$s-s_i$	$n-ni$

Is s_i/s more than random chance comparing to the background of ni/n ?



- Modified Fisher's exact test (one-tailed, right): e.g. EASE score

More conservative: uses s_i-1 instead of s_i

Functional enrichment analysis

Statistical tests

- **Gene Set Enrichment analysis (GSEA)** [Subramanian et al, PNAS, 2005](#)
 - Determines if an *a priori* defined set of genes are statistically significant (presumably concordantly different) between two biological states.
 - Sets of genes can be those within a pathway, biological process, etc.
 - Statistical significance determined by permutation (shuffling of the data)
 - Strategy
 - The genes are ordered on the basis of the parameter from the statistical test
 - For each gene set compute an enrichment score (ES) is computed (i.e a measure of how relevant or associated a biological process is for discerning the difference between the two biological states)
 - Create a running sum of a normalized Kolmogorov-Smirnov (non-parametric test) statistic.
 - Permute the class labels a large nb of times, each time recording the maximum ES over all gene sets.
 - Compare the observed ES score to the distribution of the ES scores from the permuted data.
 - Test the hypothesis that no gene set is associated with the class distinction



Functional enrichment analysis

Statistical tests

- **Gene Set Enrichment analysis (GSEA)**
 - A user-supplied ranked list of genes could also be used.
 - It determines whether *a priori* defined sets of genes show statistically significant enrichment at either end of the ranking.
 - A statistically significant enrichment indicates that the biological activity (e.g., biomolecular pathway) characterized by the gene set is correlated with the user-supplied ranking



Functional enrichment analysis

Statistical tests

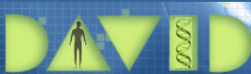
- Which functional category or biological pathway is more prevalent in the gene list than expected by chance?
- **Fisher's exact test** (parametric)
- **Gene Set Enrichment Analysis** (non-parametric)
- **Multiple testing correction:**
 - Bonferroni correction or false discovery rate (FDR)

Functional enrichment analysis

Tools

- DAVID
- EASE
- AmiGO
- GeneGo MetaCore
- GOMiner
- BiNGO & ClueGO integrated with Cytoscape
- sigPathway & GOSTat (R/Bioconductor based)
- FuncAssociate
- FatiGO
- GOEAST
- TopGO
- Gene Set Analysis (GSA)
- GSEAPreranked
- KEGG
- Wikipathway
- Reactome
- Genemania in Cytoscape
- PANTHER
- InnateDB
- STRING
- Ingenuity Pathway Analysis

Functional enrichment analysis Tools


DAVID Bioinformatics Resources 6.7
National Institute of Allergy and Infectious Diseases (NIAID), NIH

[Home](#) | [Start Analysis](#) | [Shortcut to DAVID Tools](#) | [Technical Center](#) | [Downloads & APIs](#) | [Term of Service](#) | [Why DAVID?](#) | [About Us](#)

Shortcut to DAVID Tools

▶ **Functional Annotation**

Gene-annotation enrichment analysis, functional annotation clustering, BioCarta & KEGG pathway mapping, gene-disease association, homologue match, ID translation, literature match and [more](#)

▶ **Gene Functional Classification**

Provide a rapid means to reduce large lists of genes into functionally related groups of genes to help unravel the biological content captured by high throughput technologies. [More](#)

▶ **Gene ID Conversion**

Convert list of gene ID/accessions to others of your choice with the most comprehensive gene ID mapping repository. The ambiguous accessions in the list can also be determined semi-automatically. [More](#)

▶ **Gene Name Batch Viewer**

Display gene names for a given gene list; Search functionally related genes within your list or not in your list; Deep links to enriched detailed information. [More](#)


Recommending: A [paper](#) published in *Nature Protocols* describes step-by-step procedure to use DAVID!

Welcome to DAVID 6.7


2003 - 2016

The Database for Annotation, Visualization and Integrated Discovery (DAVID) v6.7 is an update to the sixth version of our original web-accessible programs. DAVID now provides a comprehensive set of functional annotation tools for investigators to understand biological meaning behind large list of genes. For any given gene list, DAVID tools are able to:


- Identify enriched biological themes, particularly GO terms
- Discover enriched functional-related gene groups
- Cluster redundant annotation terms
- Visualize genes on BioCarta & KEGG pathway maps
- Display related many-genes-to-many-terms on 2-D view.
- Search for other functionally related genes not in the list
- List interacting proteins
- Explore gene names in batch
- Link gene-disease associations
- Highlight protein functional domains and motifs
- Redirect to related literatures
- Convert gene identifiers from one type to another.
- And more



Screen Shot 1



Screen Shot 2



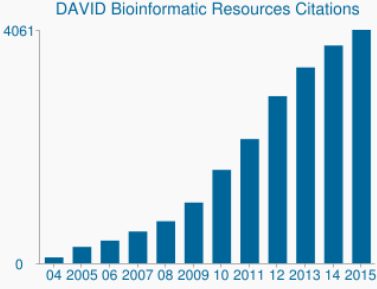
Screen Shot 3

What's Important in DAVID?

- [Current \(v 6.7\) release note](#)
- [New requirement to cite DAVID](#)
- [IDs of Affy Exon and Gene arrays supported](#)
- [Novel Classification Algorithms](#)
- [Pre-built Affymetrix and Illumina backgrounds](#)
- [User's customized gene background](#)
- [Enhanced calculating speed](#)

Statistics of DAVID

DAVID Bioinformatic Resources Citations



Year	Citations
04	~100
05	~200
06	~300
07	~400
08	~500
09	~600
10	~800
11	~1000
12	~1300
13	~1700
14	~2200
15	~2800

- [> 21,000 Citations](#)
- Average Daily Usage: ~2,600 gene lists/sublists from ~800 unique researchers.
- Average Annual Usage: ~1,000,000 gene lists/sublists from >5,000 research institutes world-wide

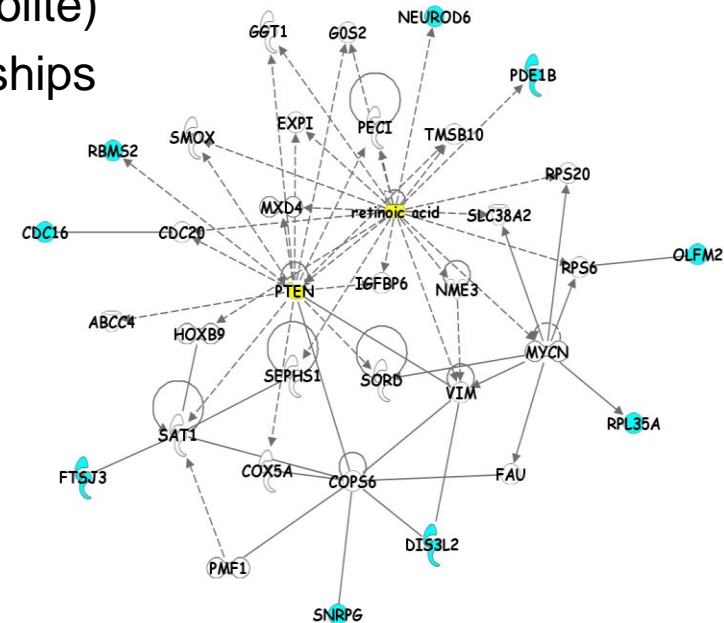
Functional enrichment analysis

Limits

- **Function enrichment analysis (GO)**
 - Structure of GO: difficult to determine which level of hierarchy is most responsible for statistical enrichment
 - The most enriched terms are often broad functional categories, not very informative
- **Pathway analysis**
 - The majority of genes have not been assigned to a pathway
 - Bias towards well-studied signalling pathways
 - Informative about information already known
- **Statistical tests**
 - Choice of the reference gene set
- **Species of interest:** some tools are species specific

Gene network analysis

- Biological processes are mainly controlled by complex networks of molecular interactions
- Network
 - graph in which an entity (molecule or metabolite) is represented by a **node** and entities' relationships by **edges** between nodes
 - Not restricted to one type of nodes or edges
 - Based on publicly available data such as experimentally-supported interactions (e.g. protein-protein interactions: PPI)



Gene network analysis

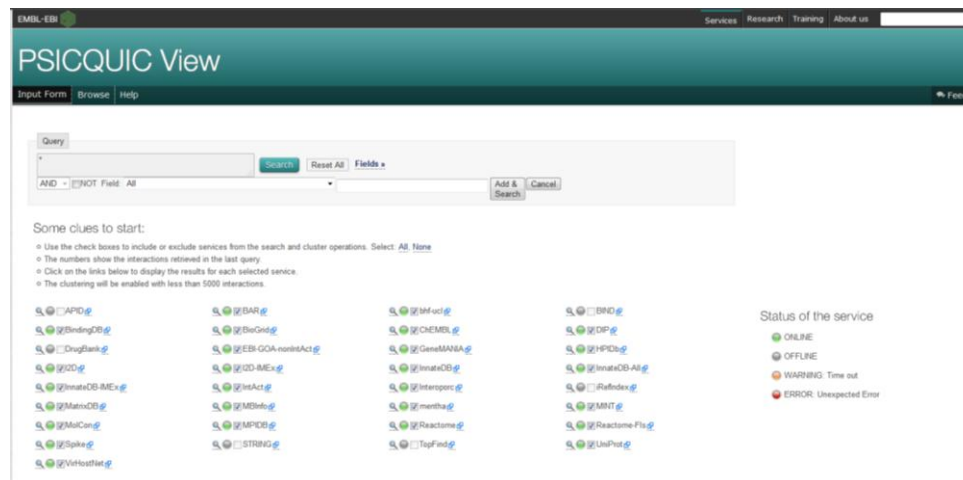
Molecular interactions databases

- **Different type of interactions:**
 - physical (e.g. PPI, protein-DNA)
 - regulatory (e.g. miRNA-mRNA)
 - biochemicals interactions (e.g. phosphorylation)
- **Experimentally-validated interaction data obtained by:**
 - Primary database: Text-mining from peer-reviewed literature
 - Manual curation from peer-reviewed literature
 - Meta-databases :integration of different sources
- **Some databases included *in silico*-predicted interactions**
- **Limited overlap between interaction information of primary databases**
- **Need integration of information provided by several primary databases**

Gene network analysis

Molecular interactions databases

- **Need integration of information provided by several primary databases:** e.g.
 - Web service PSICQUICK (integrated in Cytoscape and in Bioconductor/R)



- Commercial: IPA, manual curation by experts
- **The meaning of an edge could vary as it integrates several different types of interaction**

Gene network analysis

A complementary approach to enrichment analysis methods

- **More data-driven**
- **Interactome** (molecular interactions within a biological system) **available for several species**
- **Less biased towards well-studied pathways**
- **Additional possible integration of information associated with nodes or edges**
- **Software development for network visualization**



Gene network analysis

Limits

- **Integration of different type of interactions**
 - The meaning of an edge could vary.
- **Level of confidence associated with interactions**
 - Some techniques used to predict PPI in large-scale studies are associated with a high false positive rate (e.g. Yeast-2 Hybrid)
 - For more focused studies, biases toward well-studied biological processes
- **Interactions are context-specific**
- **Species**

Gene network analysis

Tools

- **Cytoscape (open source):**

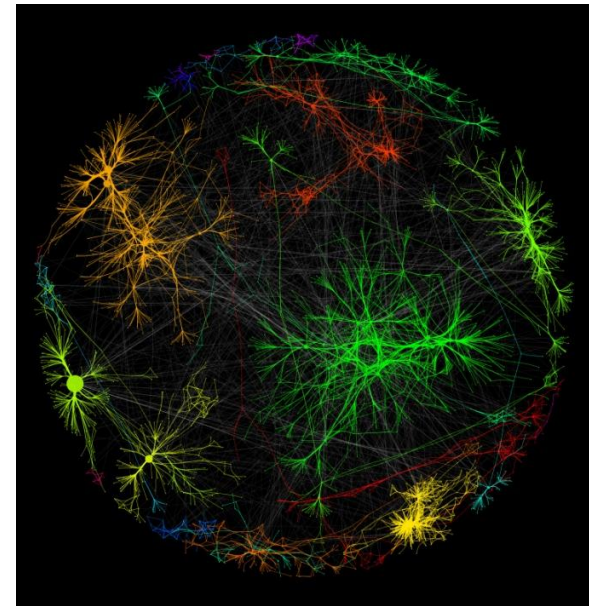


- open source software platform for visualizing complex networks and integrating these with any type of attribute data.

- a lot of Apps (plugins) available (most of them freely)

Cytoscape Apps store:

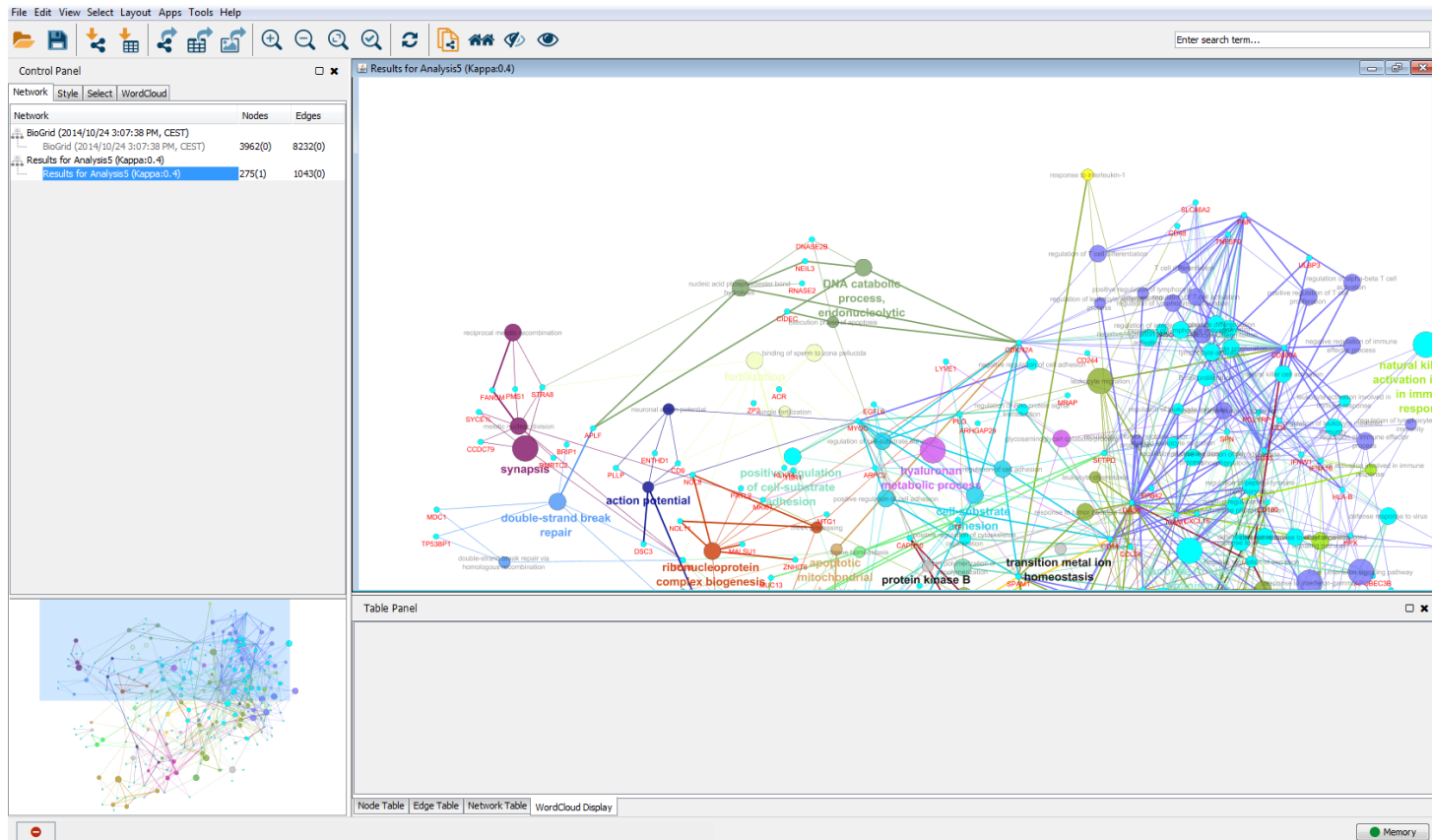
<http://apps.cytoscape.org/apps/all>



Gene network analysis

Tools

- Cytoscape interface

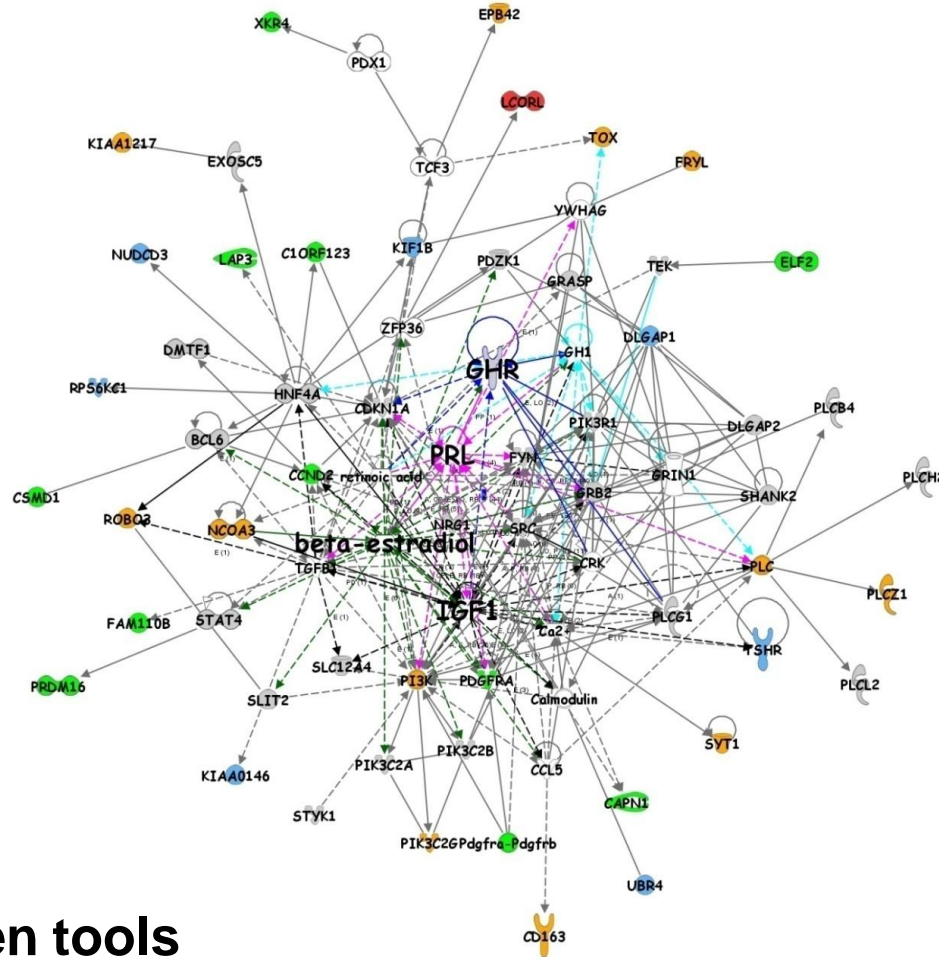


Gene network analysis

Tools

- Ingenuity Pathway Analysis (commercial)

INGENUITY
PATHWAY ANALYSIS



- Links between tools

Gene network analysis

Properties and features of gene networks

- **Hubs**

- **i.e. High degree nodes**

- Node degree: number of interactions/edges that a node has

- Important for the network structure

- Central

- **Less targeted by selection**

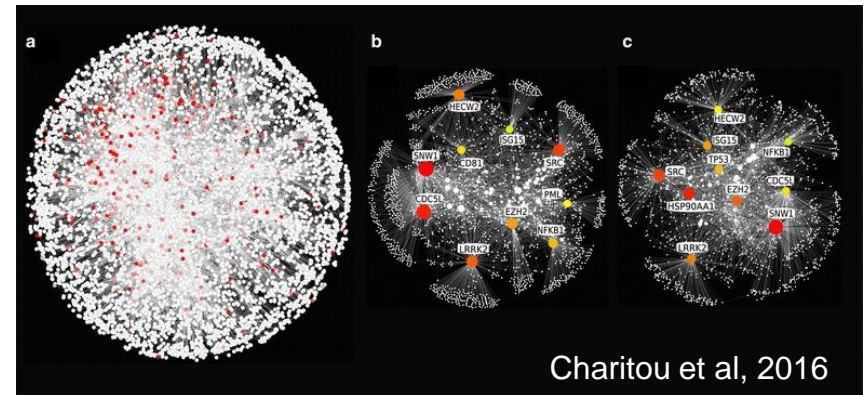
- **Bottlenecks**

- **i.e. Nodes with a high betweenness centrality**

- nodes which are the crossroads of many shortest paths

- Distance between nodes: minimum number of steps between them

- **i.e. Major bridges**

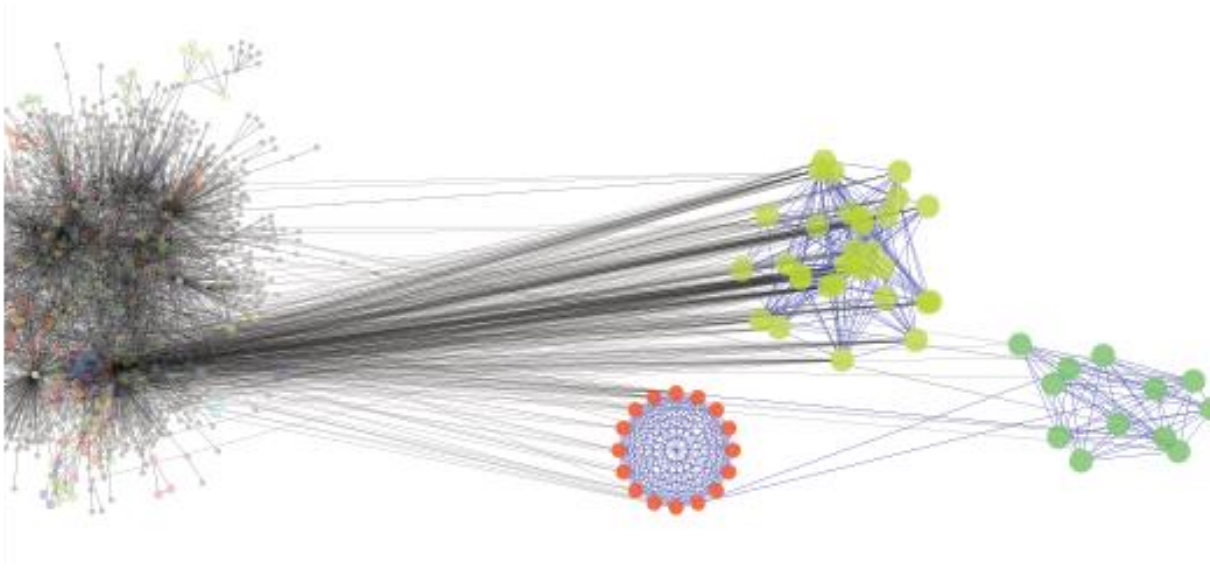


Gene network analysis

Properties and features of gene networks

- **Modules**

- i.e. group of molecules that preferentially interact with each other (sub-networks)
- Tend to be enriched for common biological functions or diseases





Gene network analysis

Properties and features of gene networks

- **Tools to identify gene networks properties**
 - **Network Analyst**
 - **Cytoscape Apps:**
 - CytoHubba
 - jActiveModules

- **Cannot be done with some tools (e.g. IPA)**
 - **Export network in Cytoscape format (need sometimes a specific licence)**

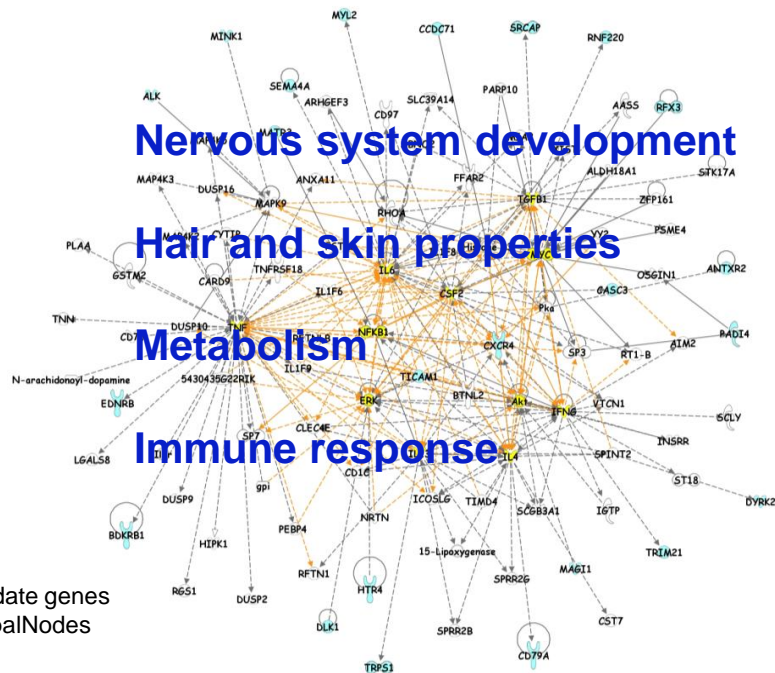
Gene network and selection

- **Position of selected genes selected within PPI network**
 - Kim et al, 2007
 - Hubs tend to be constrained and under negative selection
 - Observed positive selection at the network periphery
 - Qian et al, 2014
 - Signal of selection tend to be underrepresented for genes with fewer neighbors
 - Tend to enrich at subcentral position of the PPI network
 - Molecules of high centrality could be under strong evolutionary constraints
 - Molecules at periphery may not contribute enough to phenotypic effect
- **Interactions between recently selected genes** (Qian et al, 2014)
 - Closer interaction among genes under selection consistent with the effect of coselection

Inferring the main selective pressure

- Gene network reconstruction and identification of the main functions targeted by selection (IPA, Cytoscape)
- Exemple in West-African cattle breeds

Gautier *et al*, 2009, BMC Genomics
Flori *et al*, 2009, PLoS One
Flori *et al*, 2014, Mol. Ecol



Selective pressures

Climatic conditions

Drought and food shortage

Breeders' choices

Pathogens

Interpretation and over-interpretation

- **Goals of functional annotation of selection footprints**
 - Confirm the biological importance of the genes that have been identified as being positively selected .
 - Strengthen the credibility of the positive selection model by a development of a sound scenario
- **Sometimes a misleading approach prone to storytelling**
 - The majority of genes in a genome have important biological functions
 - easy to find connections between key genes using GO annotation and literature mining tools
 - Ex: Pavlidis et al, MBE, 2012

Integrating additional information: QTL location

- QTL databases

Release 29
(Apr 29, 2016)

Animal QTLdb

The Animal Quantitative Trait Loci (QTL) Database (Animal QTLdb) strives to collect all publicly available trait mapping data, i.e. QTL (phenotype/expression, eQTL), candidate gene and association data (GWAS), and copy number variations (CNV) mapped to livestock animal genomes, in order to facilitate locating and comparing discoveries within and between species. New data and database tools are continually developed to align various trait mapping data to map-based genome features such as annotated genes.

Many scientific journals require or recommend that any original QTL/association data be deposited into a public database before a paper may be accepted for publication. We provide user/curator accounts for direct data submission, and supply users with a data summary link to facilitate the manuscript review process.

- Cattle QTL**
There are **71,448** QTLs from **672** publications curated into the database. Those QTLs represent **505** different traits (see [data summary](#) for details).
- Chicken QTL**
There are **5,462** QTLs from **245** publications curated into the database. Those QTLs represent **336** different traits (see [data summary](#) for details).
- Horse QTL**
There are **1,139** QTLs from **61** publications curated into the database. Those QTLs represent **34** different traits (see [data summary](#) for details).
- Pig QTL**
There are **16,108** QTLs from **537** publications curated into the database. Those QTLs represent **599** different traits (see [data summary](#) for details).
- Rainbow Trout QTL**
There are **127** QTLs from **10** publications curated into the database. Those QTLs represent **14** different traits (see [data summary](#) for most recent updates).
- Sheep QTL**
There are **1,173** QTLs from **114** publications curated into the database. Those QTLs represent **204** different traits (see [data summary](#) for most recent updates).

Database summary:
Publications: 1,639
Species: 6
Traits: 1,692
QTL: 94,457

Data Alliances: ENCB, THOMSON REUTERS, UCSC, Ensembl

Features: CNV, QTLs, GWAS, Maps, Genes

This graph is partly adopted from the RGD with kind permission.

CattleQTLdb
Browse
Search
View Maps

Browse the Cattle QTLdb

Option 1 : by Chromosomes

- Chromosome 1
- Chromosome 11
- Chromosome 21
- Chromosome 2
- Chromosome 12
- Chromosome 22
- Chromosome 3
- Chromosome 13
- Chromosome 23
- Chromosome 4
- Chromosome 14
- Chromosome 24
- Chromosome 5
- Chromosome 15
- Chromosome 25
- Chromosome 6
- Chromosome 16
- Chromosome 26
- Chromosome 7
- Chromosome 17
- Chromosome 27
- Chromosome 8
- Chromosome 18
- Chromosome 28
- Chromosome 9
- Chromosome 19
- Chromosome 29
- Chromosome 10
- Chromosome 20
- Chromosome X

Option 2 : by Trait Classes

1. Health Traits
2. Meat and Carcass Traits
3. Milk Traits
4. Production Traits
5. Reproduction Traits
6. Exterior Traits

Cattle QTLdb at a Glance
— Counts by Chromosomes

Chromosome	QTLs Found
1	1867
2	1865
3	1485
4	5264
5	2489
6	3937
7	1758
8	917
9	1103
10	1390
11	2212
12	887
13	1540
14	3790
15	914
16	1077
17	1367
18	1030
19	1351
20	2149
21	939
22	739
23	823
24	579
25	544
26	3163
27	712
28	556
29	684
X	25560

Association with phenotypes and environmental covariates

SNP chip
NGS (Pool and Ind-Seq)

Identification of footprints of selection in the genome

Ecological variables
e.g. WorldClim DB

Phenotype-free approaches

Association with population ecological variables or phenotypes

Candidate regions and genes under selection

Systems biology Tools
(e.g. IPA, Cytoscape)

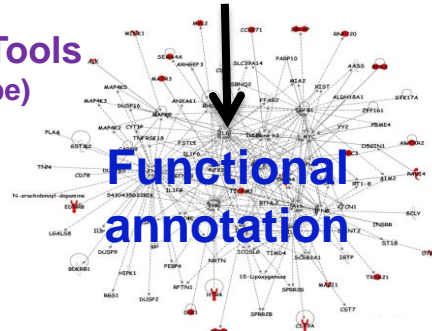
Phenotypes

Reverse ecology

Phenotyping
RNA-Seq

Principal selective pressures

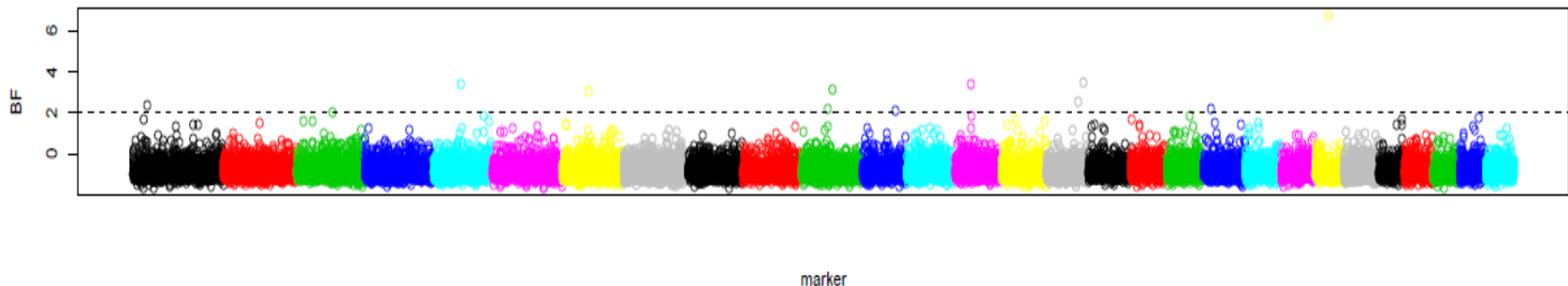
Functional annotation



Association with phenotypes and environmental covariates

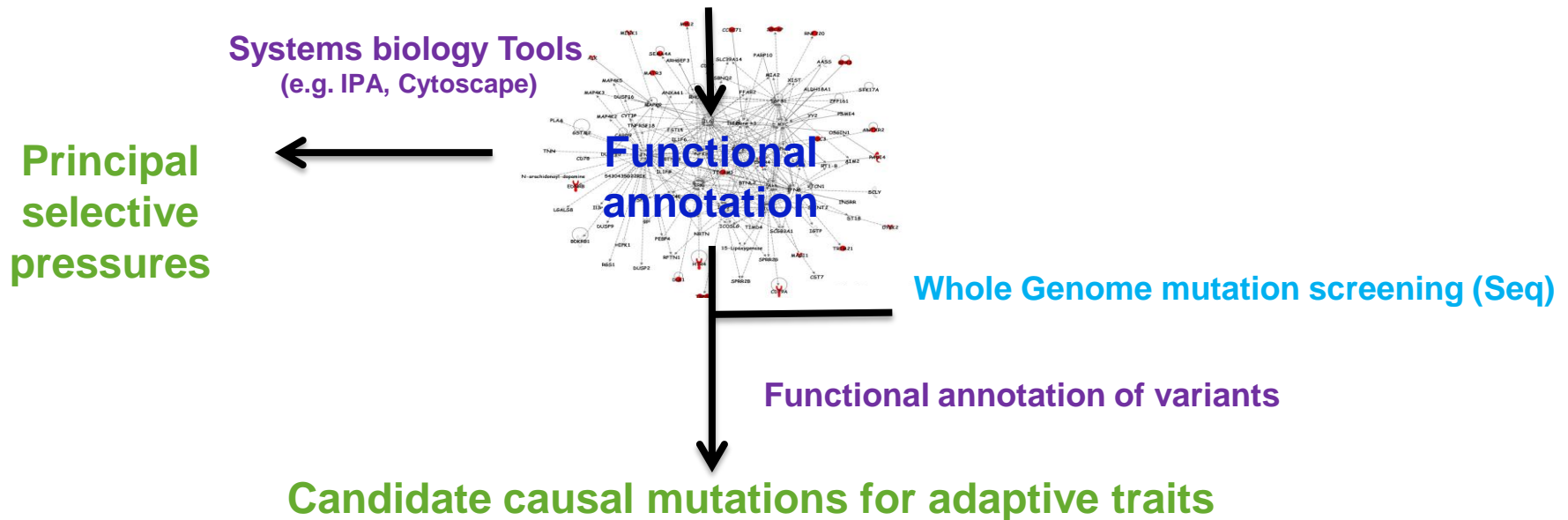
- Identifying loci underlying local adaptation using correlations between allele frequencies and ecological population variables or phenotypes
 - Multi-dimensional methods (PCA, Laloë et al, in prep)
 - Bayesian model-based approaches (e.g. Baypass; Gautier, 2015)
 - LFMM (Frichot et al, MBE, 2013)

e.g Mediterranean cattle breeds (GALIMED project): SNPs contribution to the genetic variability according to mean temperature



Identifying and annotating candidate mutations and prioritizing candidate genes

Candidate regions and genes under selection



Identifying candidate mutations

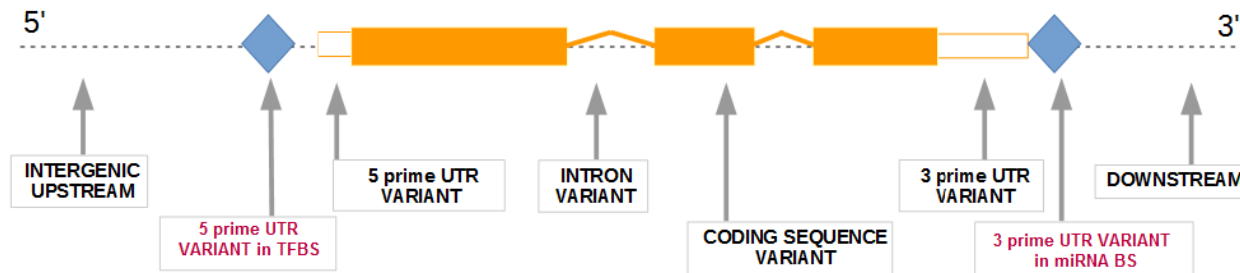
- **Individual whole genome sequencing**
 - Alignment on reference genome
 - Variant calling
 - Calculation of allelic frequencies
 - Variant annotation



Comprehensive list of variants in the genome with estimation of their frequency

Annotating candidate mutations

- Exhaustive list of variants (SNP and indels)
- Location of the variant / genes position in reference sequence



- Estimation of variant consequences on transcription and protein structure

Annotating candidate mutations

- **Tools**

e.g.



SnEff

Genetic variant annotation and effect prediction toolbox.



- **Use sequence ontology**

- i.e. set of terms and relationships used to describe the features and attributes of biological sequence
- Initially developed by the GO consortium



- **Integrate different information available on variants, e.g.**

- SIFT
- GERP scores

Annotating candidate mutations

- **Use sequence ontology**
 - i.e. set of terms and relationships used to describe the features and attributes of biological sequence
 - Initially developed by the GO consortium

[sequence_attribute](#)

[sequence_collection](#)

[sequence_comparison](#)

[sequence_feature](#)

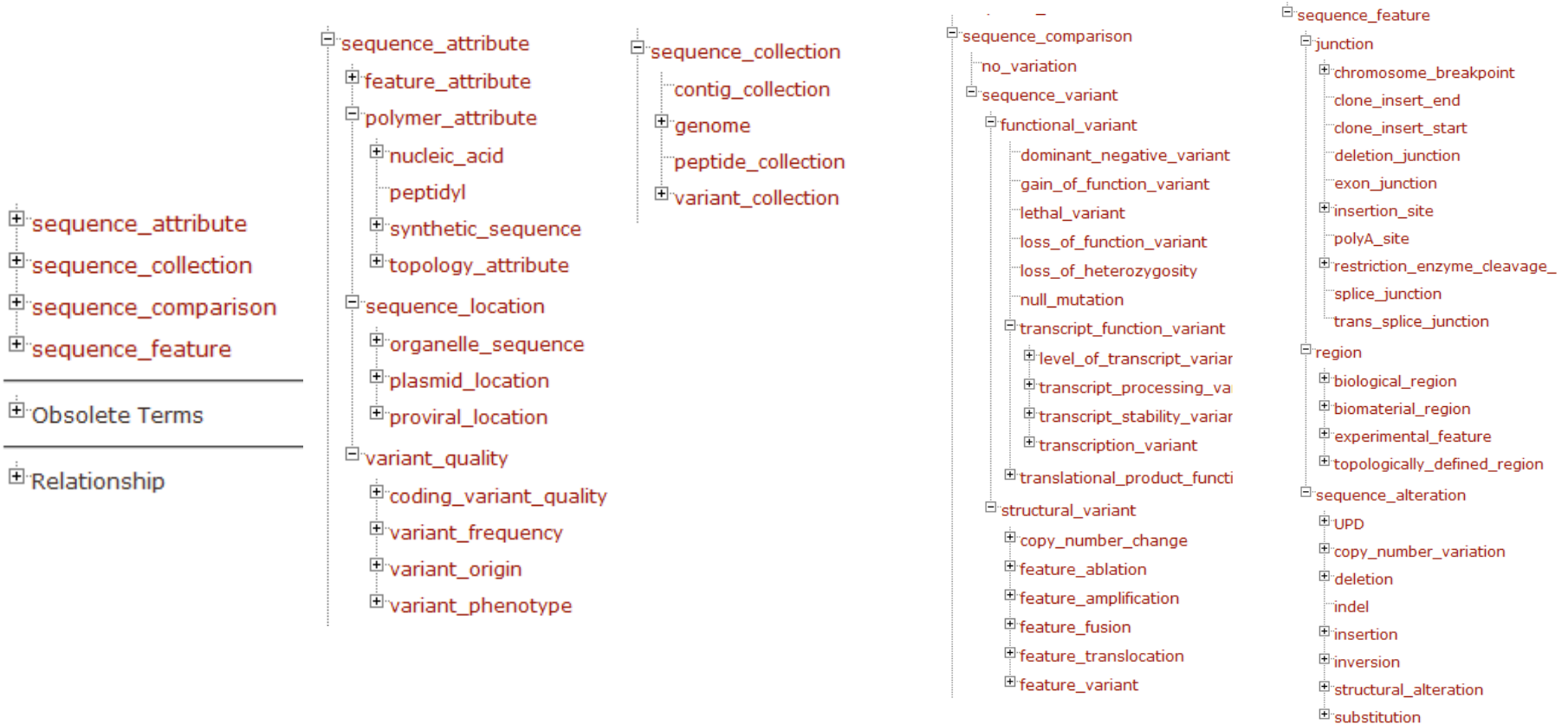
[Obsolete Terms](#)

[Relationship](#)



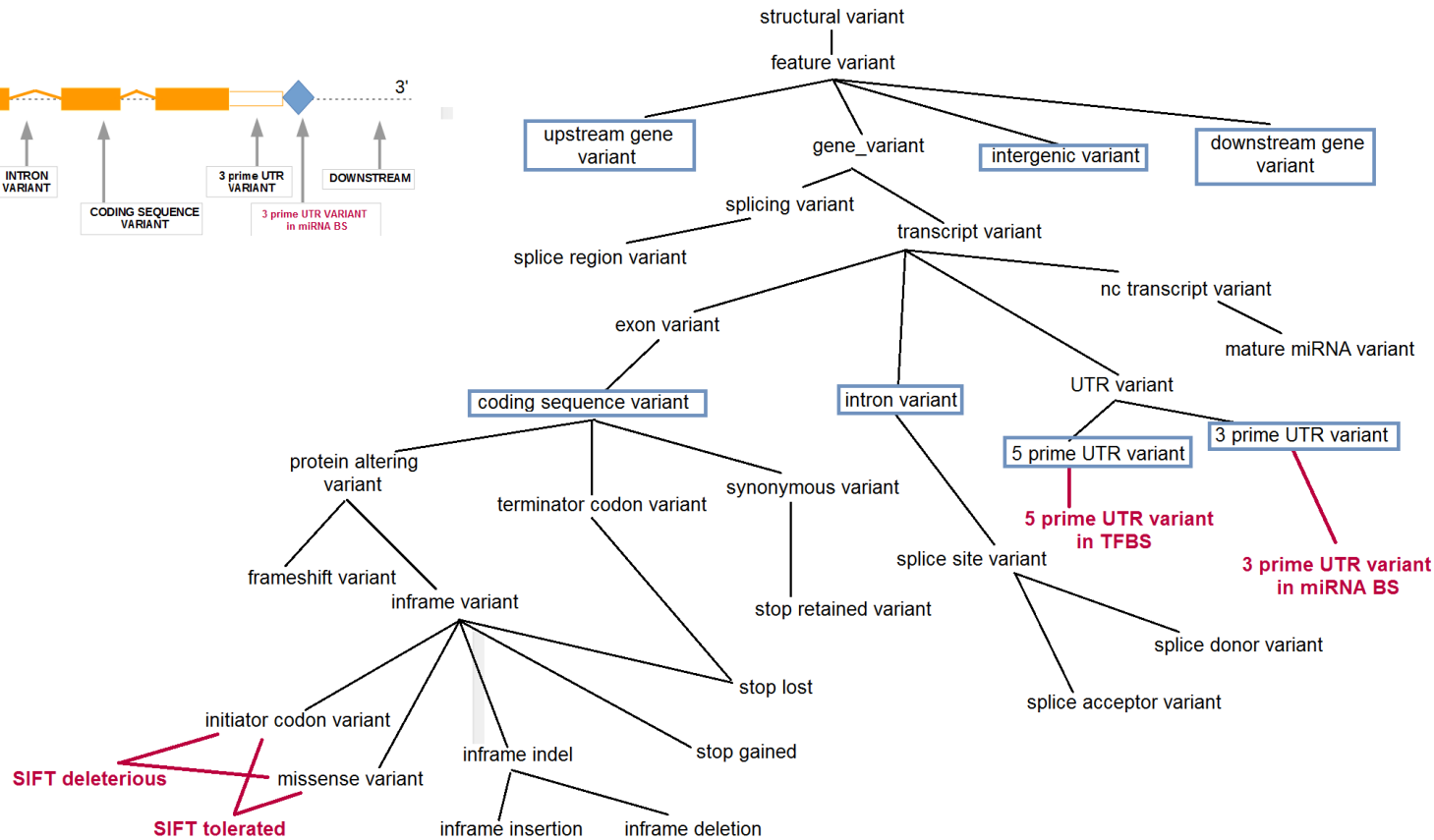
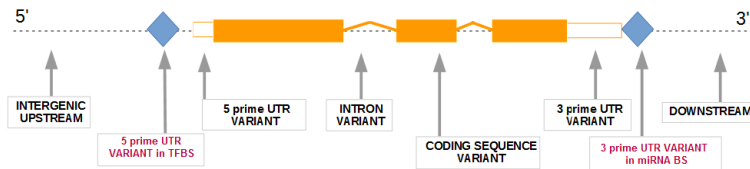
Annotating candidate mutations

- Sequence ontology



Annotating candidate mutations

- Sequence ontology in VEP (ensembl)



Annotating candidate mutations

- **Prediction of regulatory variant effect on transcription/expression**
 - e.g. variant within TFBS
 - Integrated in VEP for some species
 - Partial annotation in regulatory regions for other species
need other tools: e.g. Lasagna, mrSNP
- **Prediction of variant effect on protein structure:**
 - **SIFT:** applies on non-synonymous polymorphism and predicts whether an amino acid substitution affects protein function.
 - based on the degree of conservation of amino acid residues in sequence alignments derived from closely related sequences.
 - Tolerated/Deleterious (high, low confidence)

Annotating candidate mutations

- **Prediction of variant effect in the whole sequence**
 - **Genomic Evolutionary Rate Profiling (GERP) :**

a method for producing position-specific estimates of evolutionary constraint using maximum likelihood evolutionary rate estimation and discovering "constrained elements" that is indicative of a putative functional element.

 - Based on multiple sequence alignments and phylogenetic tree
 - Constraint intensity quantified in terms of a "rejected substitutions" (RS) score, i.e. the number of substitutions expected under neutrality minus the number of substitutions "observed" at the position.
 - Positive scores represent a substitution deficit and thus indicate that a site may be under evolutionary constraint. Negative scores indicate that a site is probably evolving neutrally



Variants prioritization

Different categories based on impact on protein level or protein structure, e.g.:

- Strong impact:
 - Non synonymous mutation deleterious (SIFT)
 - 5'UTR within TFBS
 - 3'UTR within miRNA binding site
 - In splice region
 - With positive GERP score

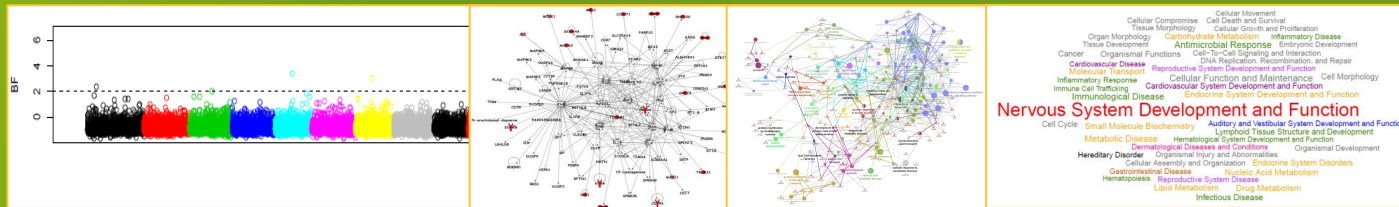
- Weakest impact
 - 5'UTR not in TFBS
 - 3'UTR not in miRNA binding site
 - Downstream
 - Upstream
 - Intronic
 - Non synonymous mutation tolerated



Candidate genes prioritization

- **Criteria**
 - Proximity to the position with the highest score (peak)
 - Number of variants in categories with strong effect (after variant annotation)
 - Allelic frequencies
 - mutation fixed (freq=1)
 - Mutation with high frequency
- **Functional annotation of these new list of candidate genes using systems biology tools**

2. Some examples of biological interpretation





Course outline

1. From regions under selection to biological interpretation

- Annotation of candidate genes using systems biology tools
 - Functional enrichment analysis: Gene ontology, pathway analysis
 - Gene networks analysis
- Inferring the main selective pressures
- Interpretation and story-telling
- Including phenotypes and environmental covariates
- Identifying candidate mutations and prioritizing candidate genes

2. Interpretation of selection footprints: some examples

- Phenotype-free approaches
 - Manual functional annotation: Senepol cattle breed
 - Functional annotation using systems biology tools: french dairy cattle breeds, West-African cattle breeds, European bison/cattle
- Association with covariates
 - Phenotypes: dairy traits in French cattle breeds

Manual functional annotation

The exemple of the Senepol cattle breed

- **Senepol: a European taurine breed with a small proportion of zebu ancestry.**
- **Living in tropical area** (Caribbean, St Croix island)
- **Identification of selection footprints**
 - 153 individuals genotyped on 47,365 SNPs
 - Tests based on the extent of haplotype homozygosity: iHS (Voight et al, 2007) and Rsb (Tang et al, 2007)
Package rehh, Gautier & Vitalis, 2012
 - 1Mb sliding windows (0.5Mb overlap)
 - Selection of regions with at least 2 SNPs exceeding the significant threshold (P<0.0001)



Flori et al, 2012, PLoS One

Manual functional annotation

The exemple of the Senepol cattle breed

- Only four regions under selection

Region	BTA	Position (Mb)	Peak position (Mb)	<i>iHS</i> _{SEN}	<i>Rsb</i> _{EUT/SEN}	<i>Rsb</i> _{ZEB/SEN}	Nb of significant SNPs	Gene closest to the maximum
#1	1	52.6–53.6	52.9	3.984	NS	3.984	2	No gene found
#2	1	2.4–3.4	3.2	NS	4.161	NS	6	TIAM1
#3	1	4.7–5.8	5.5	NS	4.538	NS	8	GRIK1
#4	20	38.6–39.6	39.5	4.055	4.576	4.055	3–6	RAI14

BTA1 => *polled* locus

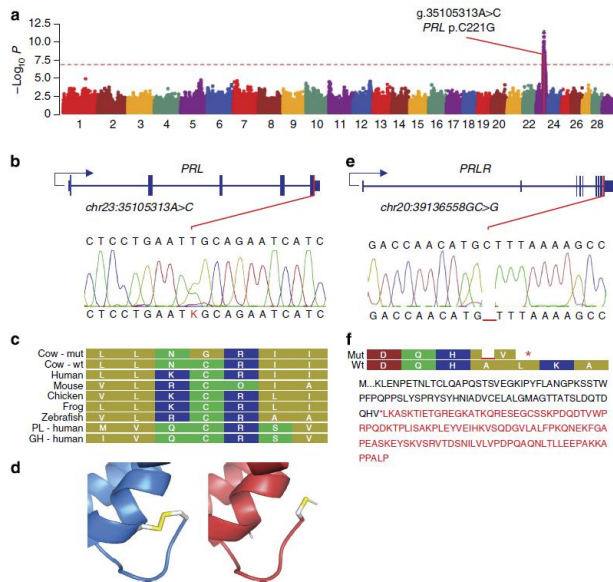
BTA20 => *slick* locus involved in a short and sleek hair coat and in thermotolerance/adaptation to tropical conditions

- Identificarion of candidate gene: boundaries located less than 25kb from the peak position

Manual functional annotation

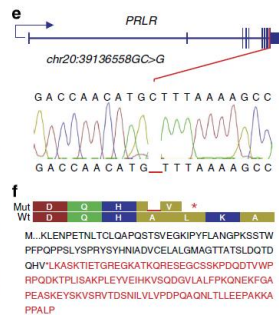
The exemple of the Senepol cattle breed

- Identification of the mutation responsible of the slick phenotype within the PRLR (Littlejohn et al, 2014)
- PRLR is located in the region #4 found under selection in Senepol



PRLP
p.Leu462

WT

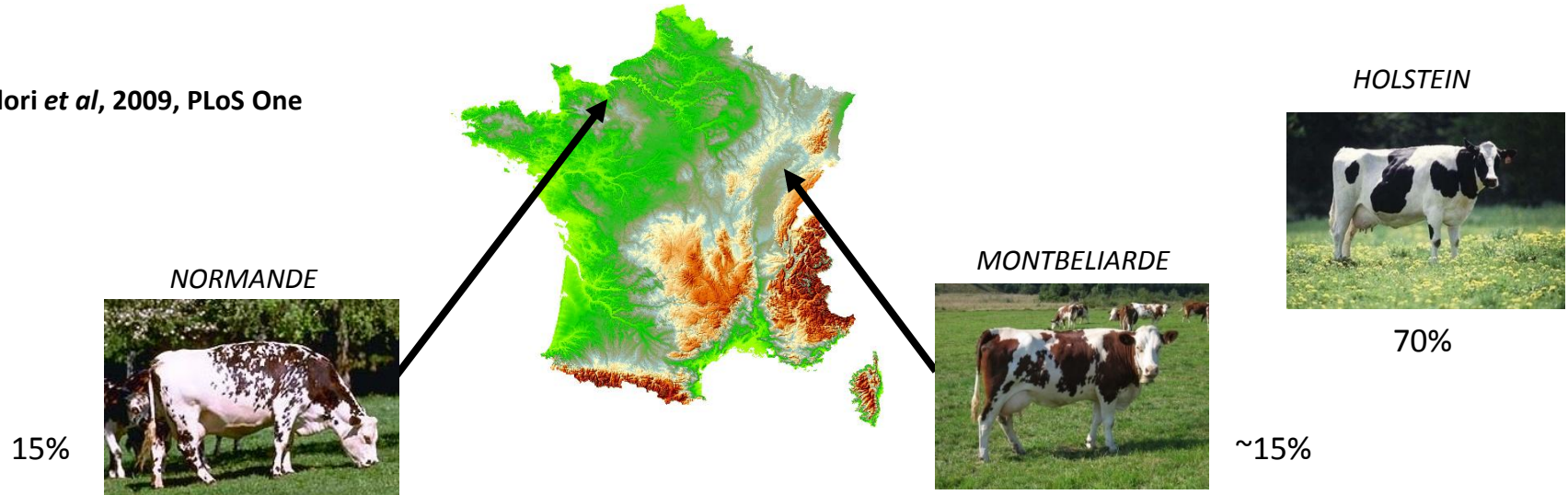


Annotation using systems biology tools

The example of the French dairy cattle breeds

- The three main French dairy cattle breeds => artificial selection

Flori *et al*, 2009, PLoS One



Breed	Herd-Book creation	Population size (2002)	Ne (généalogie)	Milk production (kg)	Content (°/oo fat/Protein)
MON	1872	1,799,200	34	7,441	38.8 / 32.5
NOR	1883	2,106,000	61	6,595	44.2 / 36.0
HOL	1922	11,535,378	42	8,628	40.9 / 31.6

- Increase in milk production but decline in reproductive performances

Annotation using systems biology tools

The example of the French dairy cattle breeds

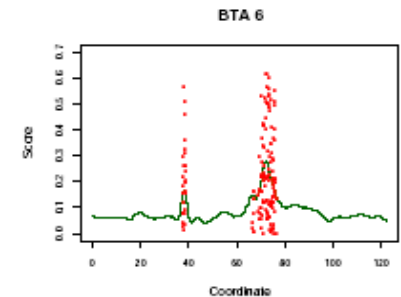
- **Identification of selection footprints**
 - 2,803 bulls genotyped on 42,486 SNPs
 - Differentiation (F_{st}) between breeds (Nicholson et al, 2002)
 - Empirical approach
 - combining information between closely related SNPs
 - Smoothed individual SNP F_{st} values over each chromosome for observed and simulated datasets
 - Calculation of local q-values
- **13 regions under selection** ($q_{value} < 0.05$)
- **Some regions contain genes with causal variants with a strong effect on milk production trait (GHR) or coloration (MC1R)**

Annotation using systems biology tools

The example of the French dairy cattle breeds

- Regions under selection

#	BTA	Start-End (peak position) in Mb	F_{ST} at the peak position (qvalue)	candidate gene	Breeds within which region is also significant
1	3	57.084–58.505 (58.343)	0.375 (0.0298)	CCCBL2	
2	4	78.833–80.43 (79.701)	0.667 (0.0298)	NUDCD3	
3	5	20.301–23.091 (21.02)	0.483 (0.0298)	na	NOR, HOL
4	5	97.803–100.826 (98.26)	0.557 (0.0298)	PIK3C2G	NOR, HOL
5	5	108.461–109.236 (109.182)	0.403 (0.0401)	CD163	
6	5	110.286–111.861 (111.552)	0.46 (0.0435)	ANO2	
7	6	37.433–38.756 (37.963)	0.566 (0.0298)	LAP3/LCORL	MON
8	6	66.599–66.935 (66.809)	0.165 (0.0435)	na	
9	6	68.938–76.32 (72.024)	0.616 (0)	PDGFRA	NOR
10	14	22.02–25.567 (22.634)	0.591 (0)	na	MON, NOR
11	18	12.987–14.058 (13.36)	0.632 (0)	MC1R	MON, HOL
12	20	31.964–33.757 (32.277)	0.523 (0.0298)	GHR	
13	26	22.137–23.191 (22.983)	0.509 (0.0298)	C10ORF76	



- Some regions contain genes with causal variants with a strong effect on milk production trait (GHR) or coloration (MC1R)
- Some regions contain loci involved in morphological traits

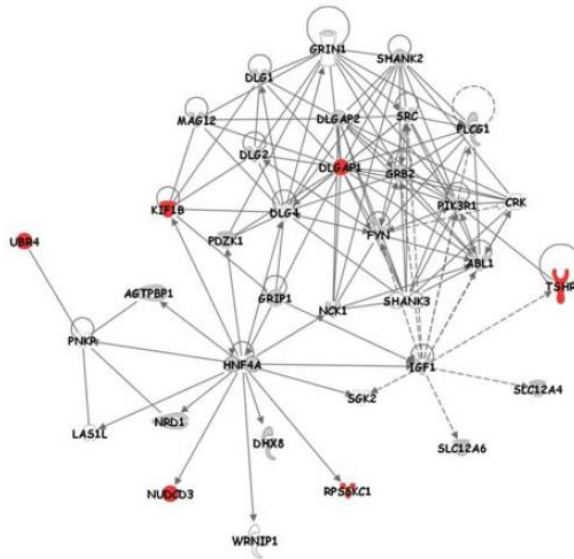
Annotation using systems biology tools

The example of the French dairy cattle breeds

Holstein



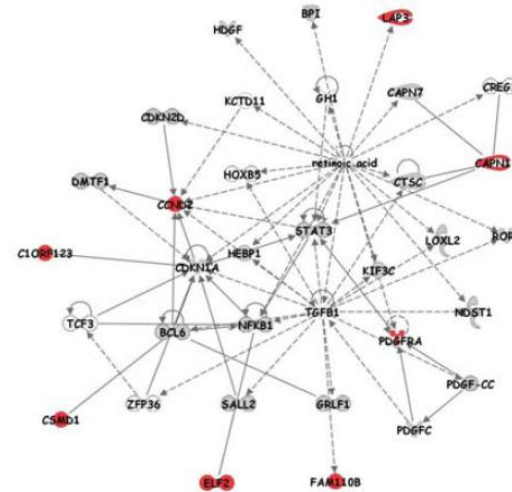
7 eligible genes/8



Montbeliarde



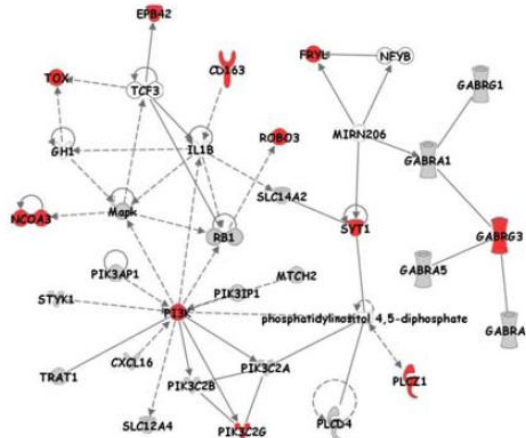
13 eligible genes /16



Normande



14 eligible genes/19

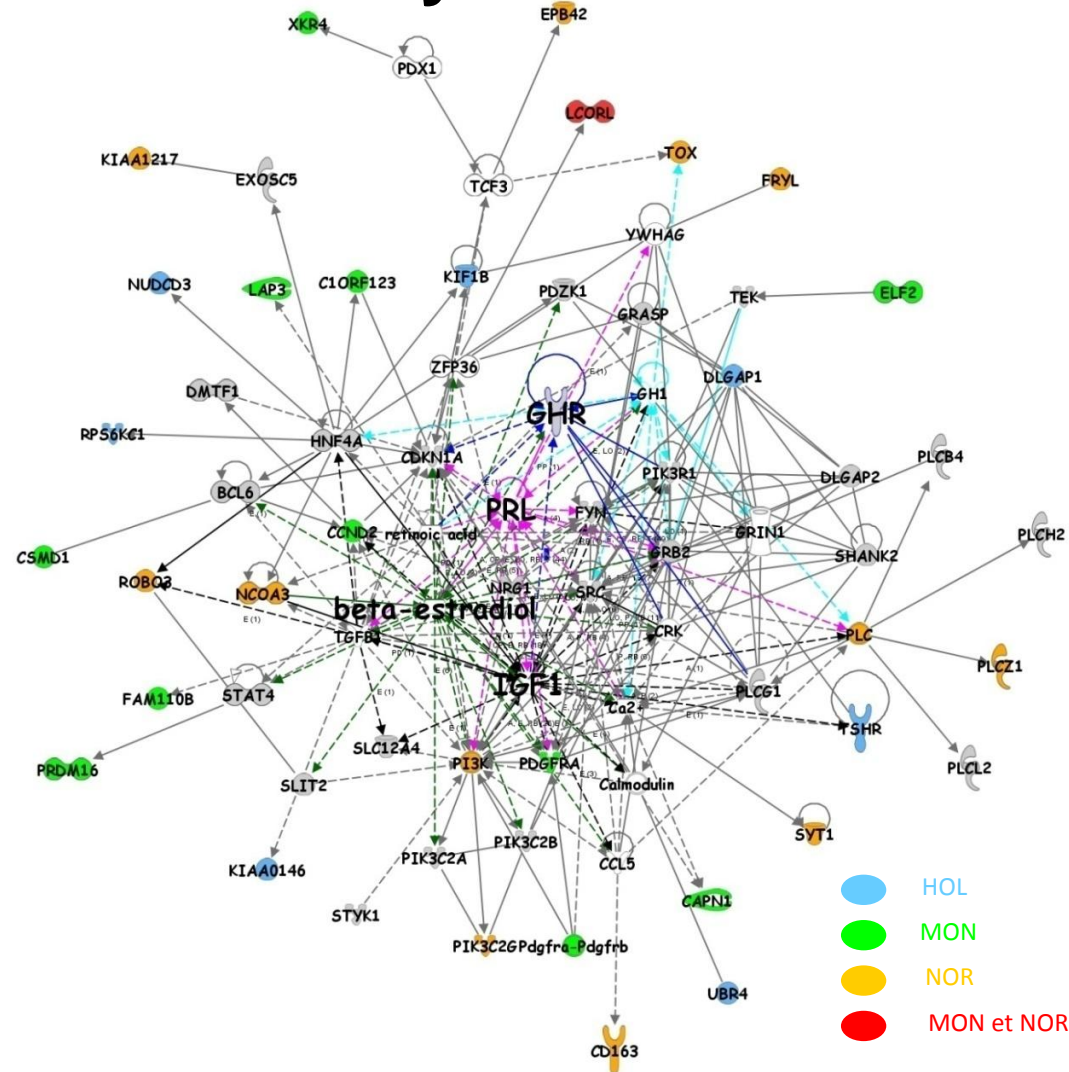


- Significant networks contain different genes under selection in the three breeds
- But genes under selection in the each breed are involved in the same biological pathway
- Genome Plasticity of genome response to the same selective pressure

Annotation using systems biology tools

The example of the French dairy cattle breeds

- **Global gene network**
 - contains genes under selection in at least one breed
- **Central role of somatotropic and gonadotropic axes in response to artificial selection**
- **Illustrates the antagonism between milk production and reproduction**





Course outline

1. From regions under selection to biological interpretation

- Annotation of candidate genes using systems biology tools
 - Functional enrichment analysis: Gene ontology, pathway analysis
 - Gene networks analysis
- Inferring the main selective pressures
- Interpretation and story-telling
- Including phenotypes and environmental covariates
- Identifying candidate mutations and prioritizing candidate genes

2. Interpretation of selection footprints: some examples

- Phenotype-free approaches
 - Manual functional annotation: Senepol cattle breed
 - Functional annotation using systems biology tools: French dairy cattle breeds, West-African cattle breeds, European bison/cattle
- Association with covariates
 - Phenotypes: dairy traits in French cattle breeds

Annotation using systems biology tools

The exemple of the West-African cattle breeds

- West-African cattle: models of adaptation to tropical conditions



Gautier et al, 2009, BMC Genomics

Exemple of the West-African cattle breeds : annotation using systems biology tools

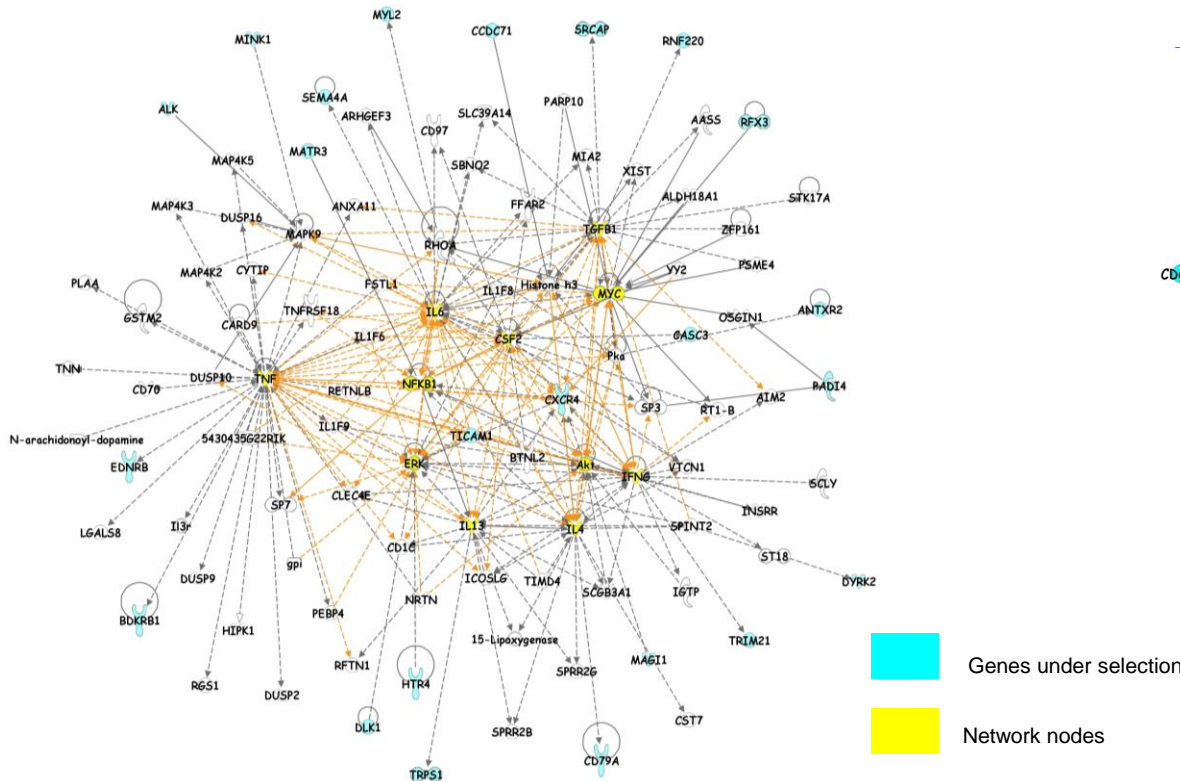
- **Identification of selection footprints**
 - 342 individuals genotyped on 36320 SNPs
 - Differentiation (F_{st} , Bayesian model) between 9 West-African populations
 - Decision rule to identify non-neutral loci based on BF, expressed in deciban unit ($dB_i = 10 \log_{10}(BF_i)$)
 - Smoothed individual SNP BF values over each chromosome
 - Permutations to estimate local p-values
 - Correction of the local p-values by computing q-value
- **53 regions under selection** (at the 5% local FDR)

Exemple of the West-African cattle breeds : annotation using systems biology tools

- **SNP annotation using Transmap Refseq**
 - 46,598 Refseq ID anchored to the bovine genome assembly
 - A SNP is considered representative of a gene if localized within gene boundaries extended by 15kb upstream and downstream
 - Annotation of the Refseq using IPA => 7,177 different genes
- **Identification of 42 candidate genes**
- **Functional analysis using IPA**
- **Network analysis using IPA and functional annotation of each significant network**
 - Each network contains at most 35 molecules

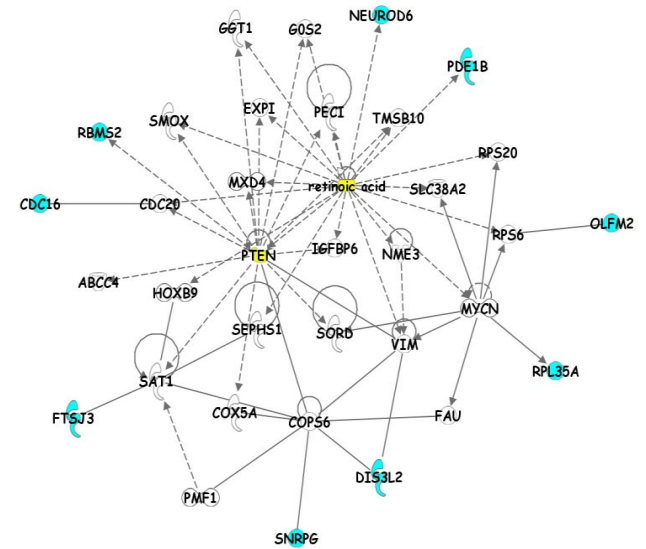
Exemple of the West-African cattle breeds : annotation using systems biology tools

A. Network N (obtained merging N1, N2 et N4) 22 genes under selection



=> Innate and adaptive immune response

B. Network N3 9 genes under selection

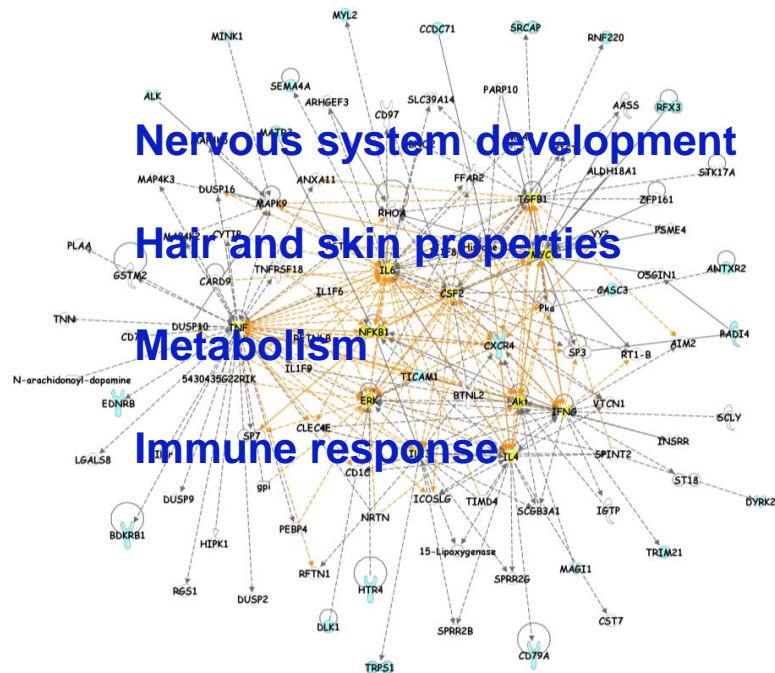




=> Nervous system diseases and disorders

INGENUITY
PATHWAY ANALYSIS

Exemple of the West-African cattle breeds : annotation using systems biology tools

Gautier *et al*, 2009, BMC Genomics
Flori *et al*, 2009, PLoS One
Flori *et al*, 2014, Mol. Ecol



 Candidate genes
 Principal Nodes



Selective pressures

Climatic conditions

Drought and food shortage

**Breeders' choices (coat color,
horn)**

Pathogens



Course outline

1. From regions under selection to biological interpretation

- Annotation of candidate genes using systems biology tools
 - Functional enrichment analysis: Gene ontology, pathway analysis
 - Gene networks analysis
- Inferring the main selective pressures
- Interpretation and story-telling
- Including phenotypes and environmental covariates
- Identifying candidate mutations and prioritizing candidate genes

2. Interpretation of selection footprints: some examples

- Phenotype-free approaches
 - Manual functional annotation: Senepol cattle breed
 - Functional annotation using systems biology tools: French dairy cattle breeds, West-African cattle breeds, European bison/cattle
- Association with covariates
 - Phenotypes: dairy traits in French cattle breeds

Annotation using systems biology tools

The example of the European bison

- **Wisent: the largest European herbivore, emblematic of the continent wildlife**
- **Identification of selection footprints between the European bison and cattle**

Gautier et al, under revision

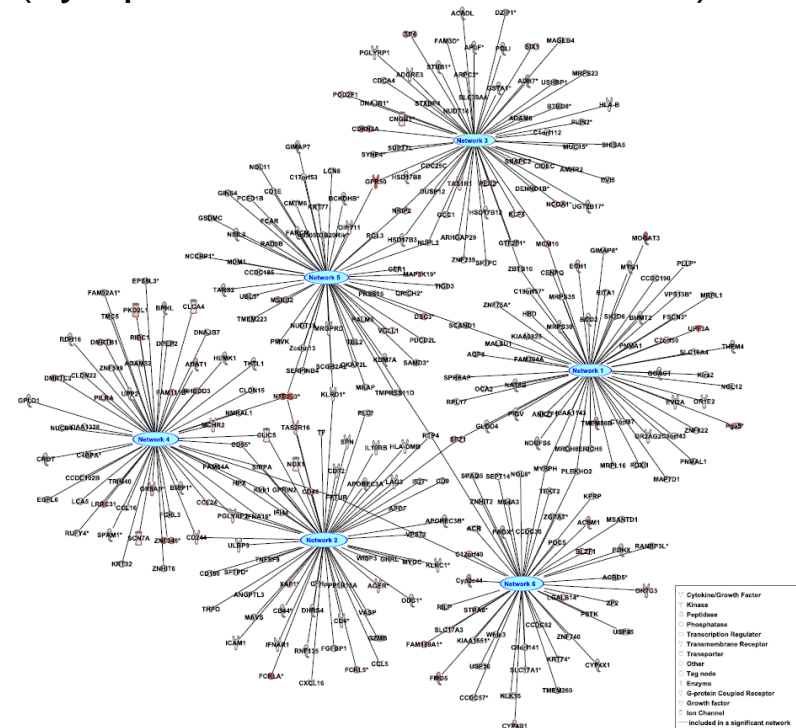
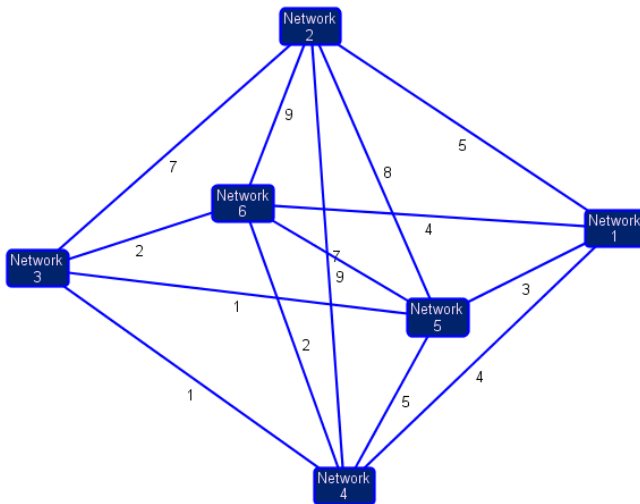


- Using individual whole-genome sequences (10X, 100nt paired-end)
- Calculating K_a/K_s (Kimura, 1983; *KaKs calculator*, Zhang et al, 2006):
 - the K_a/K_s ratio is an indicator of selective pressure acting on a protein-coding gene.
 - i.e. ratio of the number of non-synonymous substitutions per non-synonymous site (K_a) to the number of synonymous substitution per synonymous site (K_s), in a given period of time.
- Genes with a K_a/K_s ratio above 1 are evolving under positive selection

Annotation using systems biology tools

The example of the European bison

- **873 transcripts under selection** => 450 genes ready for functional and network analyses
- **70% of the genes participated to a global network** corresponding to six significant interconnected networks (by up to 12 common molecules)



Annotation using systems biology tools

The example of the European bison

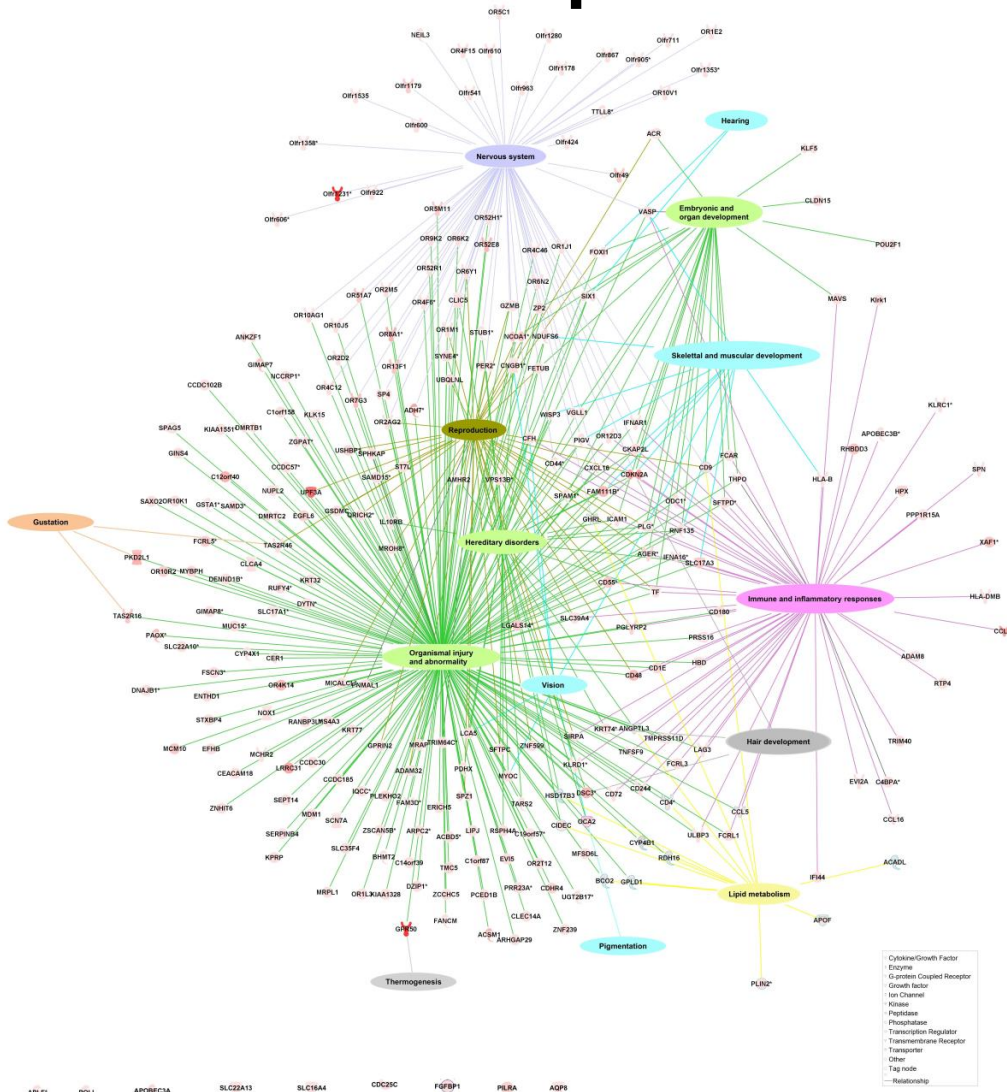
- **Functional analysis**

Main functional categories	Diseases and Biological Functions	p-value	Number of molecules
Physiological System Development and Function	Nervous System Development and Function	$1.26 \times 10^{-16} - 3.05 \times 10^{-02}$	54
	Hematological System Development and Function	$1.88 \times 10^{-04} - 3.05 \times 10^{-02}$	39
	Immune Cell Trafficking	$1.88 \times 10^{-04} - 3.05 \times 10^{-02}$	34
	Embryonic development	$2.07 \times 10^{-03} - 3.05 \times 10^{-02}$	15
	Organ development	$2.07 \times 10^{-03} - 3.05 \times 10^{-02}$	11
Diseases and Disorders	Inflammatory Response	$1.88 \times 10^{-04} - 3.05 \times 10^{-02}$	35
	Infectious Disease	$2.70 \times 10^{-04} - 3.05 \times 10^{-02}$	14
	Connective Tissue Disorders	$3.29 \times 10^{-04} - 3.05 \times 10^{-02}$	12
	Organismal Injury and Abnormalities	$3.29 \times 10^{-04} - 3.05 \times 10^{-02}$	208
Molecular and Cellular Functions	Skeletal and Muscular Disorders	$3.29 \times 10^{-04} - 3.05 \times 10^{-02}$	11
	Cell Death and Survival	$2.96 \times 10^{-05} - 3.05 \times 10^{-02}$	33
	Cellular Compromise	$2.96 \times 10^{-05} - 3.05 \times 10^{-02}$	24
	Cell Morphology	$1.10 \times 10^{-04} - 3.05 \times 10^{-02}$	22
	Cellular Assembly and Organization	$1.10 \times 10^{-04} - 3.05 \times 10^{-02}$	17
	Cell-To-Cell Signaling and Interaction	$1.17 \times 10^{-04} - 3.05 \times 10^{-02}$	94

— Most significant category: Nervous System Development and Function

Annotation using systems biology tools

The example of the European bison



- Many genes are related to several functions
=>pleiotropic role
- Some genes underlie obvious distinctive features between wisent and cattle

Annotation using systems biology tools

The example of the European bison

- **Some genes underlie obvious distinctive features between wisent and cattle**

Key functions

- Hair development and thermogenesis
- Olfactory and taste receptor genes
- Genes involved in immune response
- Genes involved in lipid metabolism and mammary gland development
- Genes involved in key processes of nervous system and in several by-product phenotypes of domestication

Presumed adaptation

- => to cold climatic conditions (wisent)
- => to food resources (Forest-habitat, wisent)
- => to different pathogen exposures (wildlife/domestication)
- => artificial selection in cattle (dairy traits)
- =>to domestication/wildlife

Annotation using systems biology tools

The example of the European bison

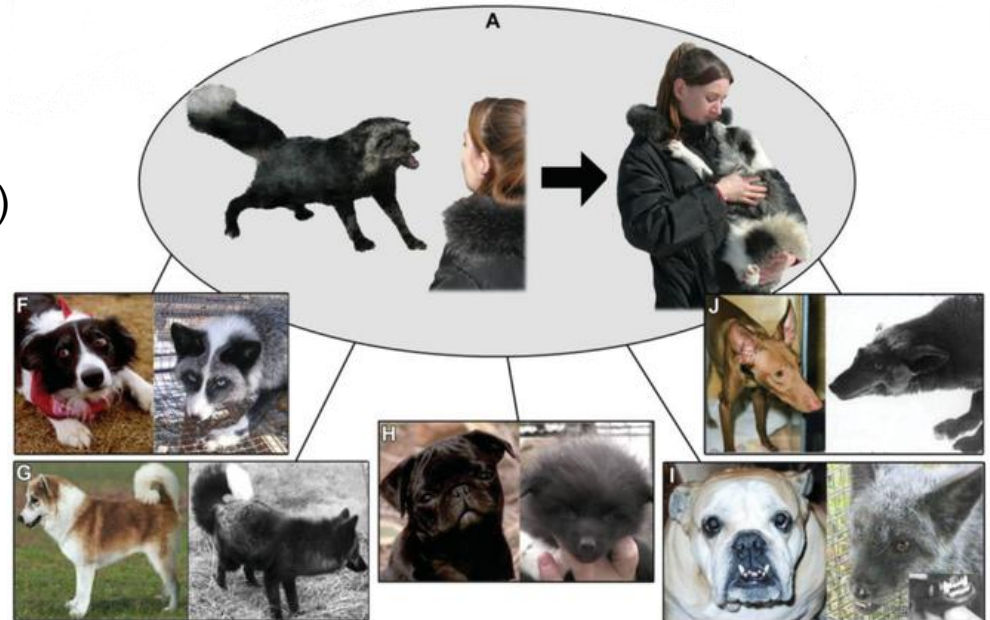
- Numerous genes involved in the domestication syndrom
- Main modified traits in domestication syndrom

Selection

- Docility
- Juvenile behavior
- Depigmentation (white patches)

By-products phenotypes of domestication

- Floppy and reduced ears
- Reduced muzzle and jaws
- Smaller teeth
- Shape of the tail (curly)
- More frequent oestrus

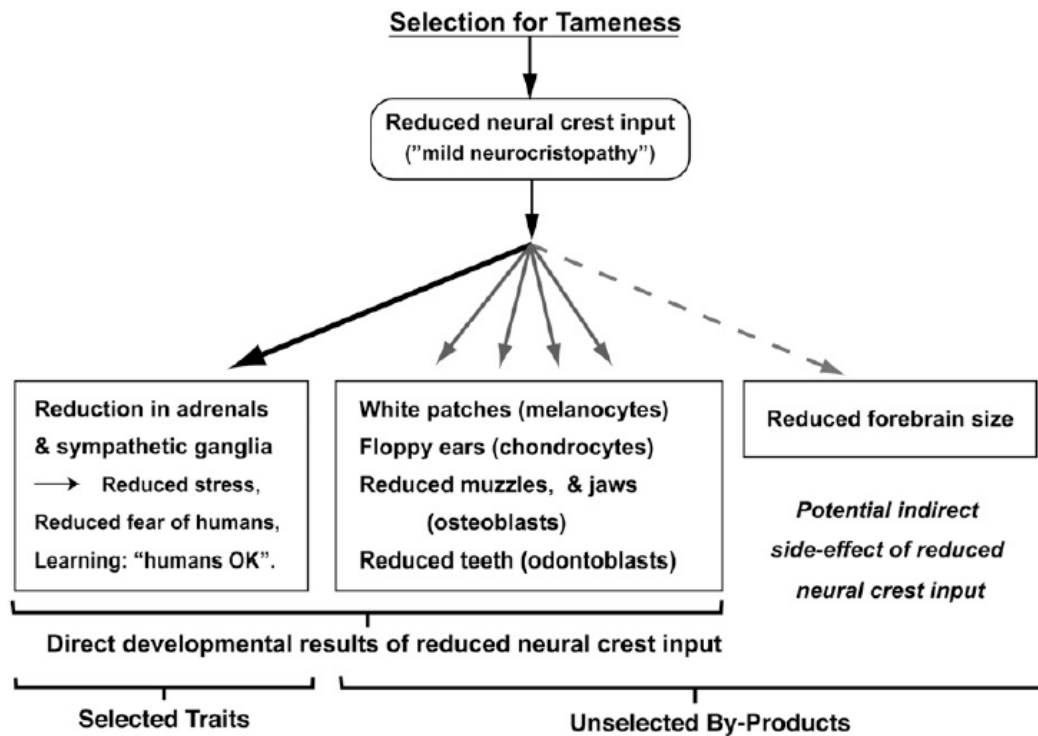


Belyaev and Trut , 1989; Trut et al, 2009

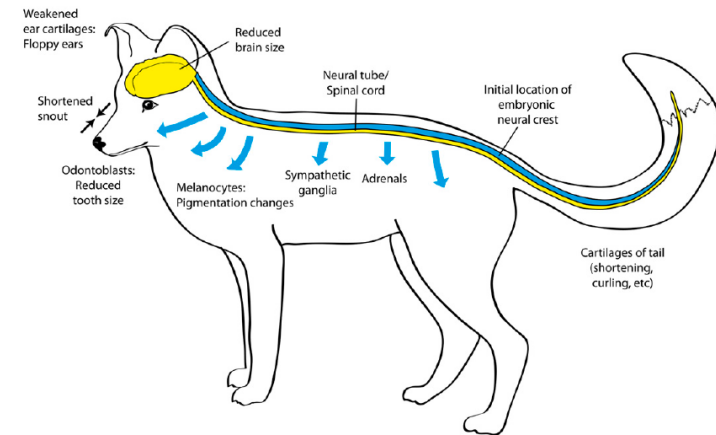
Annotation using systems biology tools

The example of the European bison

- Our analysis give an empirical support to the unified explanation of the domestication syndrome in mammals proposed by Wilkins et al, 2014



Central role of the neural crest





Course outline

1. From regions under selection to biological interpretation

- Annotation of candidate genes using systems biology tools
 - Functional enrichment analysis: Gene ontology, pathway analysis
 - Gene networks analysis
- Inferring the main selective pressures
- Interpretation and story-telling
- Including phenotypes and environmental covariates
- Identifying candidate mutations and prioritizing candidate genes

2. Interpretation of selection footprints: some examples

- Phenotype-free approaches
 - Manual functional annotation: Senepol cattle breed
 - Functional annotation using systems biology tools: French dairy cattle breeds, West-African cattle breeds, European bison/cattle
- Association with covariates
 - Phenotypes: dairy traits in French cattle breeds

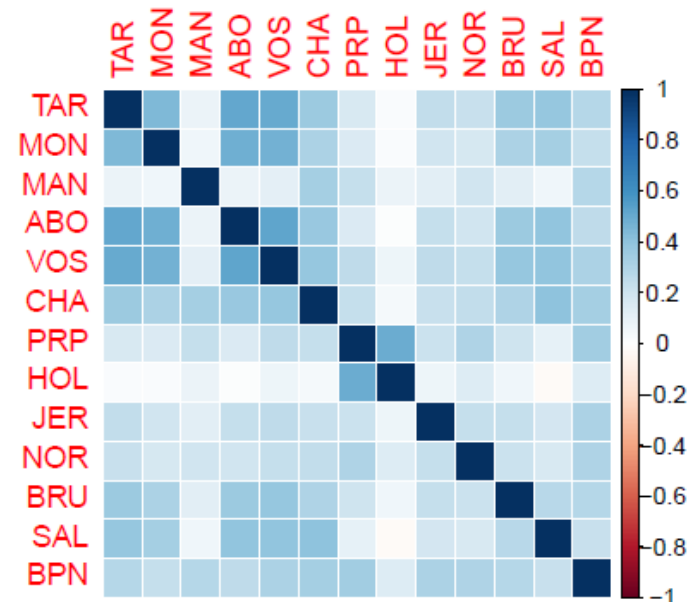
Association with covariates

The example of the French cattle breeds

Gautier and Flori, *in prep*

- Identification of regions under selection associated with dairy traits
 - Dairy traits: average milk production, lactation length, fat and protein content
 - 13 breeds
 - 50K SNP chip
 - Baypass (Gautier, 2015):
 - Auxiliary covariate model

Correlation plot based on Ω

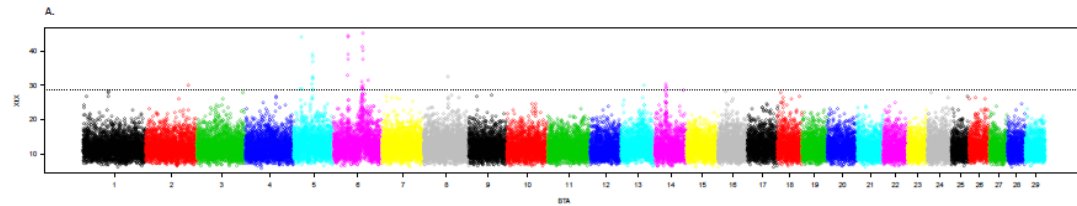


Association with covariates

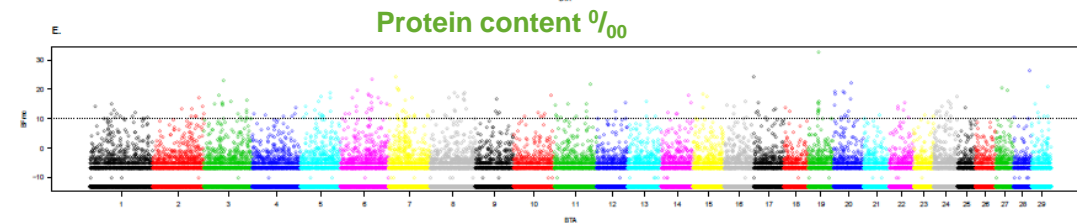
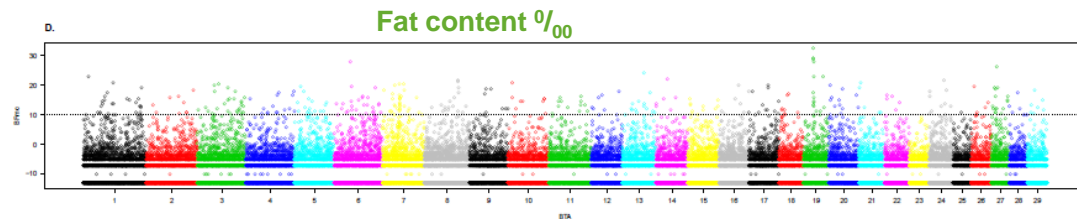
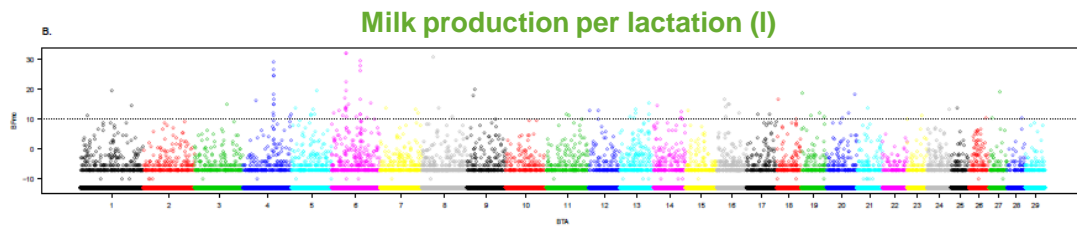
The example of the French cattle breeds

- Identification of regions under selection associated with dairy traits

- $XtX > 28.6$ (0.01% PODs)



- $BF_{mc} > 10 \Rightarrow$ strong evidence
 - $BF_{mc} > 20 \Rightarrow$ decisive evidence



Association with covariates

The example of the French cattle breeds

- **Identification of regions under selection associated with dairy traits**
 - Using Refgene file from UMD3.1 assembly on ucsc (14627 Refseq ID)
 - Annotation of SNPs using the Refgene file: a SNP is representative of a gene if it is located within gene boundaries +/- 15kb
=>94492 single ID
 - $BF_{mc} > 20$ =>decisive evidence
 - » Milk production: NUDC3, LAP3, LCORL
 - » Fat content: 24 RefSeq
 - » Protein content 5 Refseq
 - Choice of a less stringent criteria would be more informative (e.g. $BF_{mc} > 10$)



Conclusion

- **Functional and network analysis might be powerful to find a biological interpretation of footprints of selection**
- **Taking special care to avoid storytelling**
 - Criteria to define a candidate region and candidate genes (e.g. significant threshold, nb of significant SNP/gene or region, distance from peak)
 - » If too stringent criteria => loss of information
 - Integrating other information
 - » Sequencing data=>exhaustive list of variants that can be annotated and allelic frequencies
 - » QTL information
 - » Using association methods with population environmental covariates or phenotypes



Some references

- **Charitou et al, GSE, 2016.** Using biological network to integrate, visualize and analyse genomic data
- **Khatri et al, PLoS Computational biology, 2012.** Ten years of pathway analysis: current approaches and outstanding challenges
- **Pavlidis et al, MBE, 2012.** A critical assessment of storytelling: Gene ontology categories and the importance of validating genomic scans.