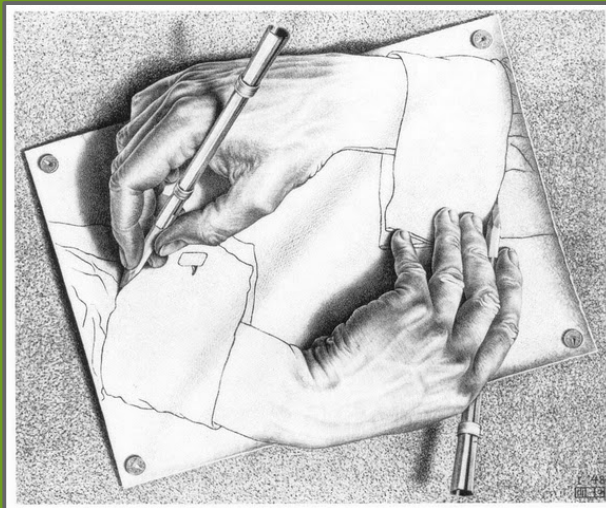


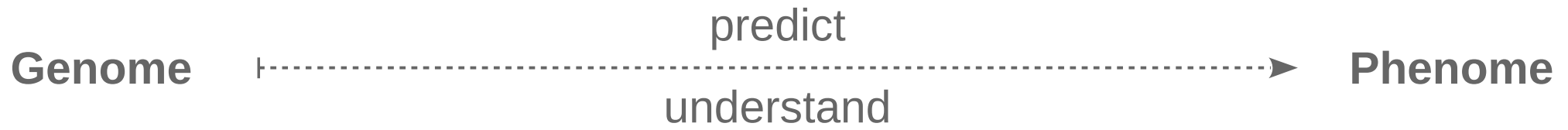


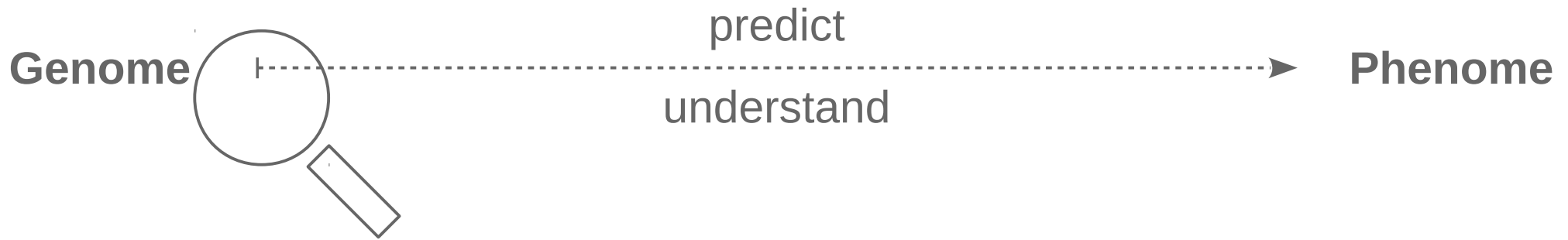
Functional annotation of livestock genomes

Chromatin structure and gene expression



Sylvain Foissac, INRA Toulouse, France
ASAS-ADSA Midwest Meeting, March 2019, Omaha, USA

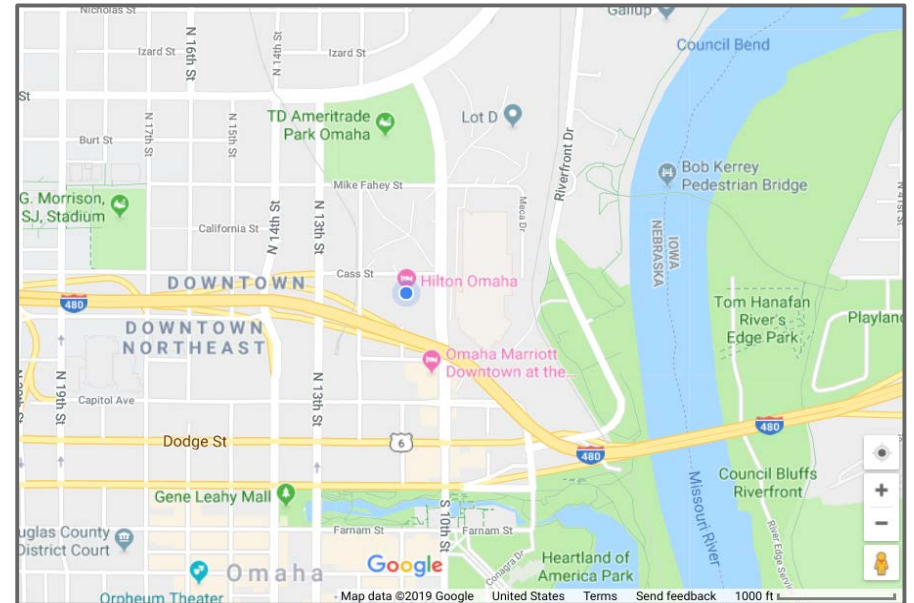


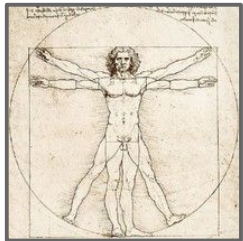
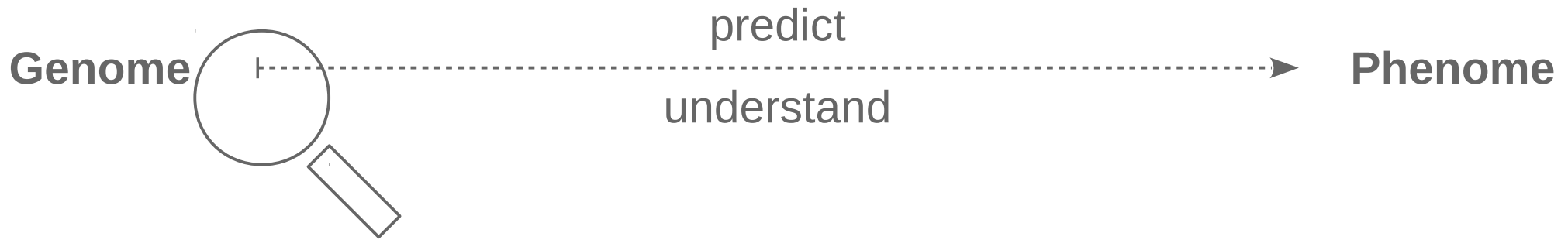


raw DNA sequence



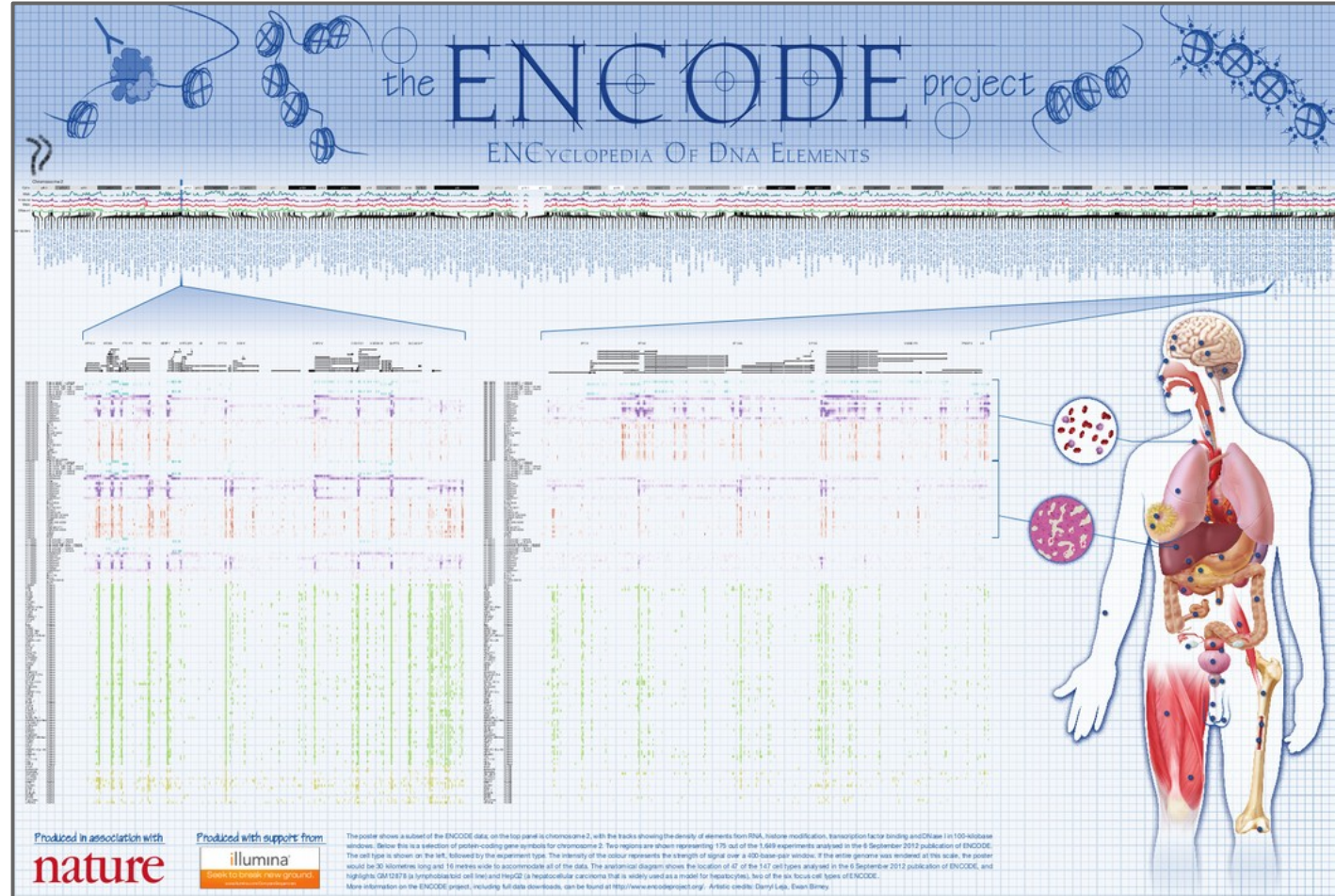
genomic annotation

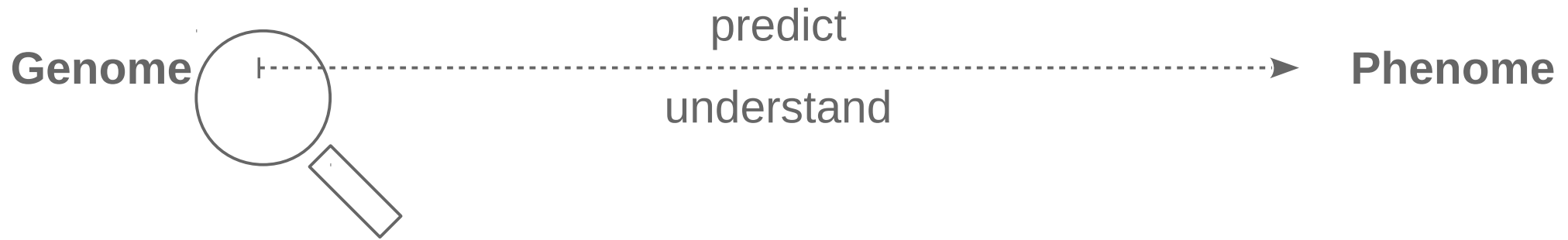




&

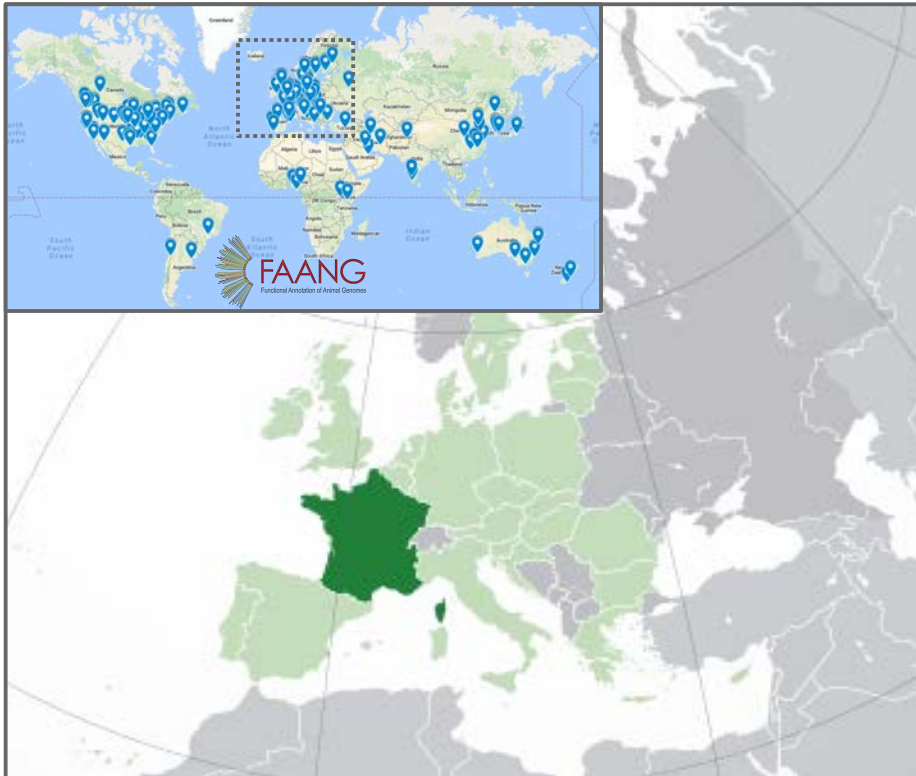
\Rightarrow





&

\Rightarrow ?



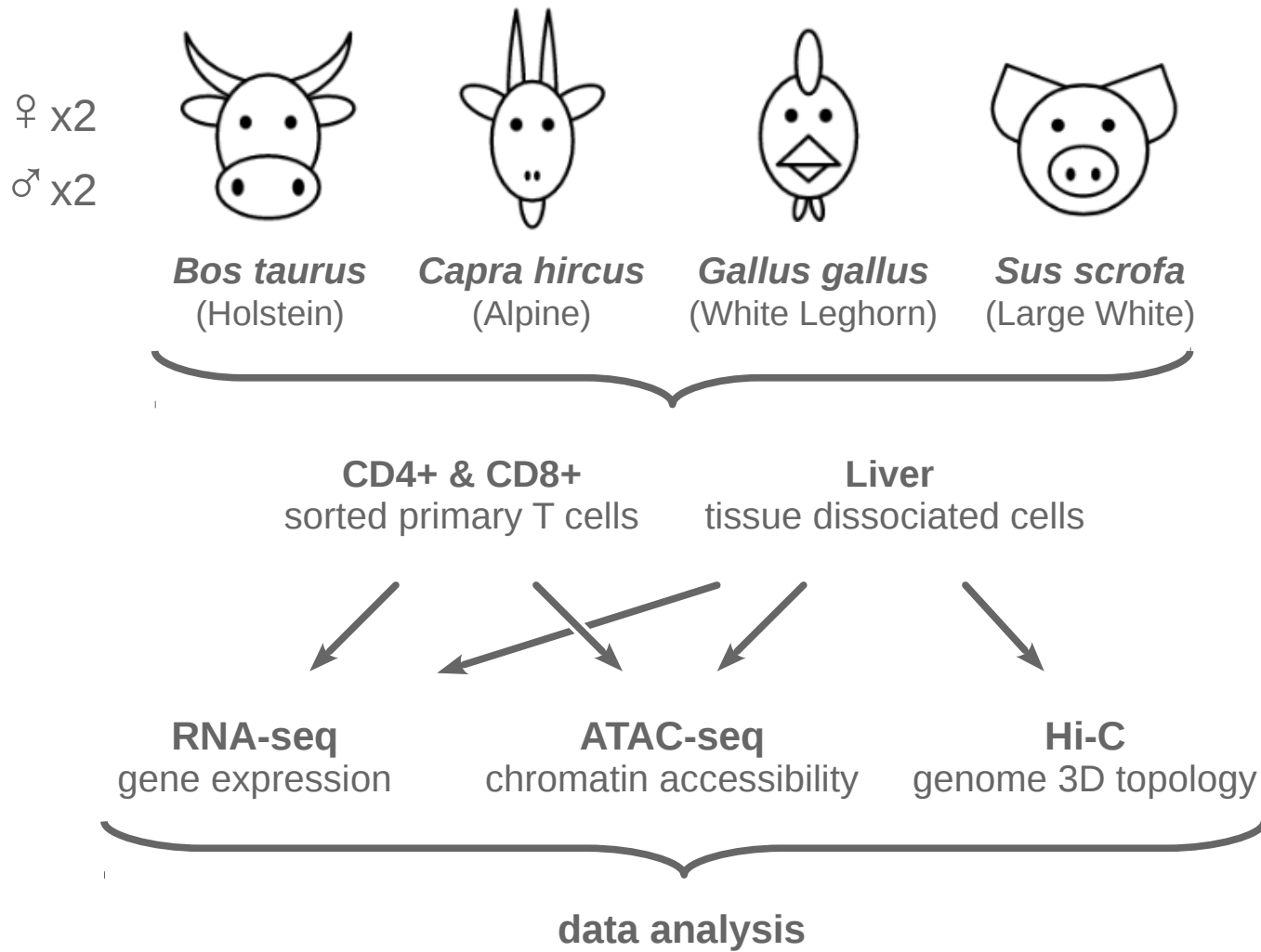
*Elisabetta
Giuffra*



*Sylvain
Foissac*



FR-AgENCODE: the french pilot project of FAANG from INRA



- ◆ 4 animals per species
- ◆ 40+ tissues per animal
Liver, CD4+ T cells, CD8+ T cells, sperm, plasma, heart, lung, skin, fat, duodenum, ileum, jejunum, cerebellum, frontal lobe, olfactory bulb, trigeminal ganglia, hypothalamus, pancreas, adrenals, kidney, muscle, bone, joints, spleen, lymphatic nodes, peyer's patches, ovary, oocytes, oviduct, uterus, mammary gland, acini, testis, seminal vesicle, etc.
- ◆ 2,000+ frozen samples
- ◆ Protocols and BioSamples entries at the EMBL-EBI
- ◆ Samples available at the INRA Bridge/CRB-anim biobank



Michèle
Tixier-
Biochard



Stéphane
Fabre

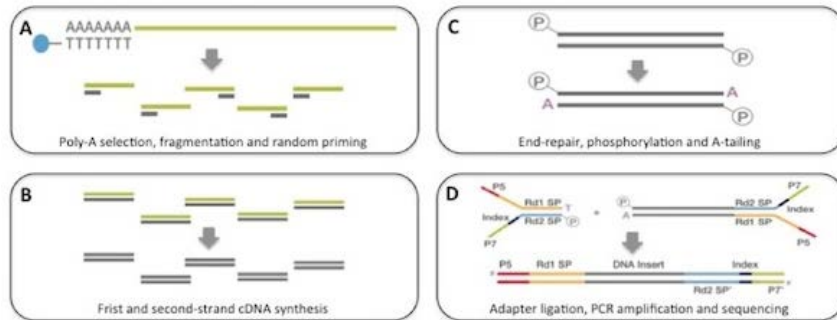
Index of <ftp://ftp.faang.ebi.ac.uk/ftp/protocols/samples/>

[Up to higher level directory](#)

Name	Size	Last Modified	
File: DEDJTR_SOP_CryofreezingTissue_20160317.pdf	69 KB	3/22/16	12:00:00 AM GMT+1
File: INRA_SOP_PBMC_purification_cattle_caprine_201...	60 KB	5/4/16	12:00:00 AM GMT+2
File: INRA_SOP_PBMC_seperation_swine_blood_2016...	89 KB	5/4/16	12:00:00 AM GMT+2
File: INRA_SOP_alveolar_macrophages_mammals_sam...	486 KB	7/29/16	12:00:00 AM GMT+2
File: INRA_SOP_chicken_splenocytes_sampling_20160...	131 KB	7/29/16	12:00:00 AM GMT+2
File: INRA_SOP_liver_spleen_mammarygland_forHiC_s...	327 KB	7/29/16	12:00:00 AM GMT+2
File: INRA_SOP_oocytes_granulosa_mammals_samplin...	246 KB	7/29/16	12:00:00 AM GMT+2

The screenshot shows the BioSamples website interface. At the top, there is a navigation bar with 'EMBL-EBI', 'Services', 'Research', 'Training', and 'About us'. The main header features the 'BioSamples' logo and a search bar containing the text '"FR-AgENCODE"'. Below the search bar, there are tabs for 'Home', 'Search', 'Submit', 'Documentation', and 'About'. The main content area displays 'Search results' for the query, showing 'Showing 1 to 10 of 68 results'. The first result is 'Submission GSB-721' with the identifier 'SAMEG317390' and an update date of '01-06-2018 19:05'. The submission description is 'Samples collected in frame of the Fr-AgENCODE project'. Below the description, there are several rows of 'has member' links, each with a unique identifier like 'SAMEA7981918'. The second result is 'Submission GSB-99' with the identifier 'SAMEG100045' and an update date of '23-02-2018 10:07'. Its description is 'FR-AgENCODE : Sample used for RNA-Seq Sequencing', and it also lists several 'has member' links with identifiers like 'SAMEA1088309'.

RNA-seq: transcriptome profiling

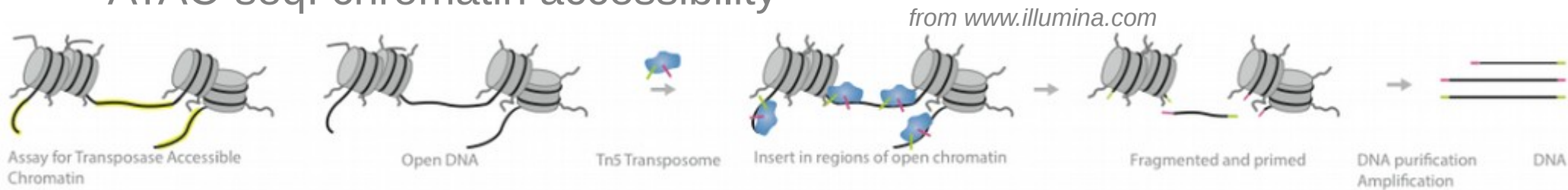


Diane Esquerré

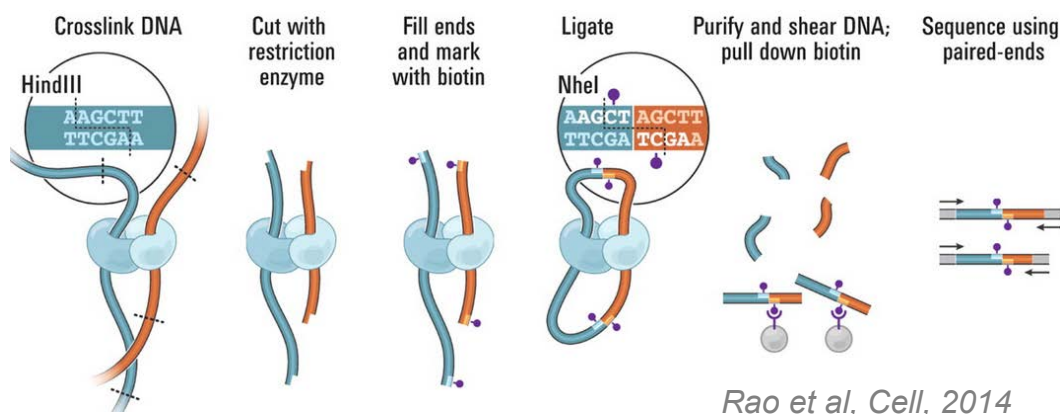


Hervé Acloque

ATAC-seq: chromatin accessibility



Hi-C: chromosome conformation



Index of <ftp://ftp.faang.ebi.ac.uk/ftp/protocols/assays/>

📁 [Up to higher level directory](#)

Name	Size	Last Modified	
File: DEDJTR_SOP_Sampl...	134 KB	3/22/16	12:00:00 AM GMT+1
File: INRA_SOP_ATAC-se...	398 KB	8/5/16	12:00:00 AM GMT+2
File: INRA_SOP_Hi-C_HA...	900 KB	12/8/16	12:00:00 AM GMT+1
File: INRA_SOP_RNA_extr...	216 KB	5/4/16	12:00:00 AM GMT+2
File: INRA_SOP_mRNA-se...	392 KB	11/24/17	12:00:00 AM GMT+1
File: ISU_SOP_ChIP_swini...	1688 KB	1/25/19	3:34:00 PM GMT+1
File: ISU_SOP_Fetaliver_...	279 KB	1/25/19	3:34:00 PM GMT+1
File: ISU_SOP_Macrophag...	88 KB	2/13/19	10:14:00 AM GMT+1
File: ROSLIN_SOP_Chrom...	458 KB	11/18/16	12:00:00 AM GMT+1
File: ROSLIN_SOP_Isolati...	74 KB	11/16/16	12:00:00 AM GMT+1
File: ROSLIN_SOP_Isolati...	59 KB	11/16/16	12:00:00 AM GMT+1
File: ROSLIN_SOP_Tissu...	328 KB	11/16/16	12:00:00 AM GMT+1

Completed experiments

RNA-seq

Cattle					Goat				
	cattle1	cattle2	cattle3	cattle4		goat1	goat2	goat3	goat4
cd4					cd4				
cd8	NA				cd8				
liver					liver				
Chicken					Pig				
	chicken1	chicken2	chicken3	chicken4		pig1	pig2	pig3	pig4
cd4	NA	NA			cd4	NA			
cd8	NA	NA	NA		cd8				
liver					liver				

RNA-seq:
~5B read pairs

ATAC-seq

Cattle					Goat				
	cattle1	cattle2	cattle3	cattle4		goat1	goat2	goat3	goat4
cd4					cd4		NA		
cd8					cd8				
liver	NA	NA	NA	NA	liver				
Chicken					Pig				
	chicken1	chicken2	chicken3	chicken4		pig1	pig2	pig3	pig4
cd4				NA	cd4				
cd8	NA		NA	NA	cd8				
liver					liver				NA

ATAC-seq:
~3B read pairs

HiC

Cattle					Goat				
	cattle1	cattle2	cattle3	cattle4		goat1	goat2	goat3	goat4
liver	NA	NA	NA	NA	liver				
Chicken					Pig				
	chicken1	chicken2	chicken3	chicken4		pig1	pig2	pig3	pig4
liver					liver				

Hi-C
~2B read pairs

>80% of experiments completed

◆ Bioinformatics pipelines

<u>RNA-seq</u>	<u>ATAC-seq</u>	<u>HiC</u> (HiC-Pro)
<ul style="list-style-type: none"> - Read trimming (trimgalore) - Read mapping (STAR2) - Transcript modelling (Cufflinks2) - New gene annotation (Cuffmerge2) - Transcript/gene express^o quantification (RSEM) - LncRNA calling (FEELnc) 	<ul style="list-style-type: none"> - Read trimming (trimgalore) - Read mapping to genome (Bowtie2) - PCR duplicate removal (Samtools) - Mitochondrial read removal (Samtools) - Peak calling (MACS2) 	<ul style="list-style-type: none"> - Read trimming (cutadapt) - Read mapping to genome (Bowtie2) - Inconsistent pairs filtering (Samtools) - Contact matrix generation and normalization (ICE) - TAD calling (Armatus) - Visualization (Juicebox)

◆ Data integration and comparison

◆ Statistical analyses



Sarah
Djebali

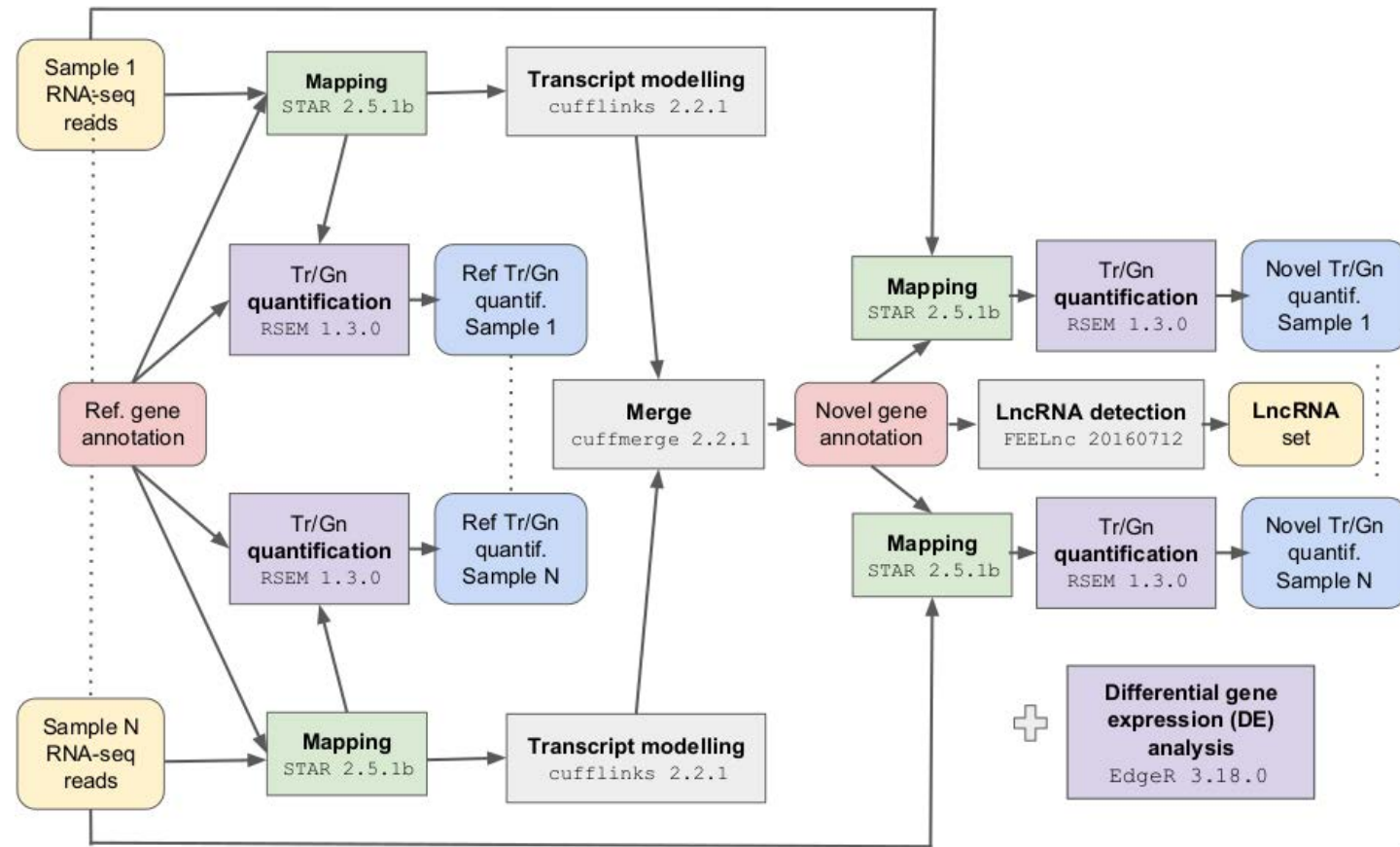


Nathalie
Vialaneix

◆ Bioinformatics pipelines

RNA-seq

- Read trimming (trimgalore)
- Read mapping (STAR2)
- Transcript modelling (Cufflinks2)
- New gene annotation (Cuffmerge2)
- Transcript/gene express^o quantification (RSEM)
- LncRNA calling (FEELnc)



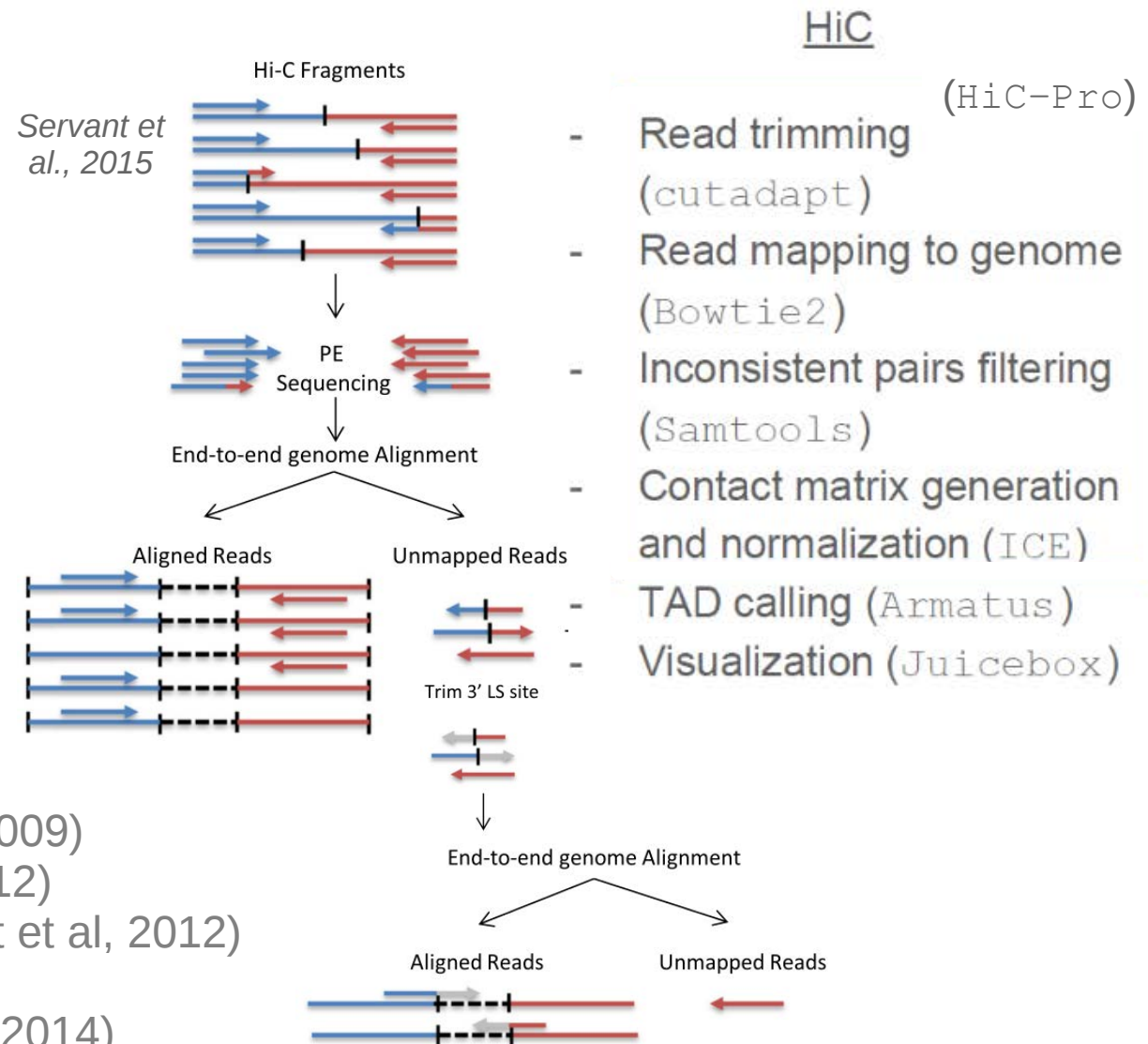
◆ Bioinformatics pipelines

Hi-C detailed workflow

- Trim reads (ligation site)
- **Map on reference genome**
- Discard inconsistent pairs
- **Count reads in pairs of genomic bins & generate contact matrix**
- Normalize contact matrix (non parametric, matrix balancing)
- Generate html report
- **Identify Topological Associated Domains, *cis* and *trans* interactions**
- **Identify A and B compartments**

Software

- HiC-Pro pipeline (Servant et al 2015)
- Bowtie2 mapping (Langmead et al, 2009)
- ICE normalization (Imakaev et al, 2012)
- HiTC display and A/B comp. (Servant et al, 2012)
- HiFive pipeline (Sauria et al, 2015)
- Armatus TAD finding (Filippova et al, 2014)
- Juicebox browser (Durand et al, 2016)

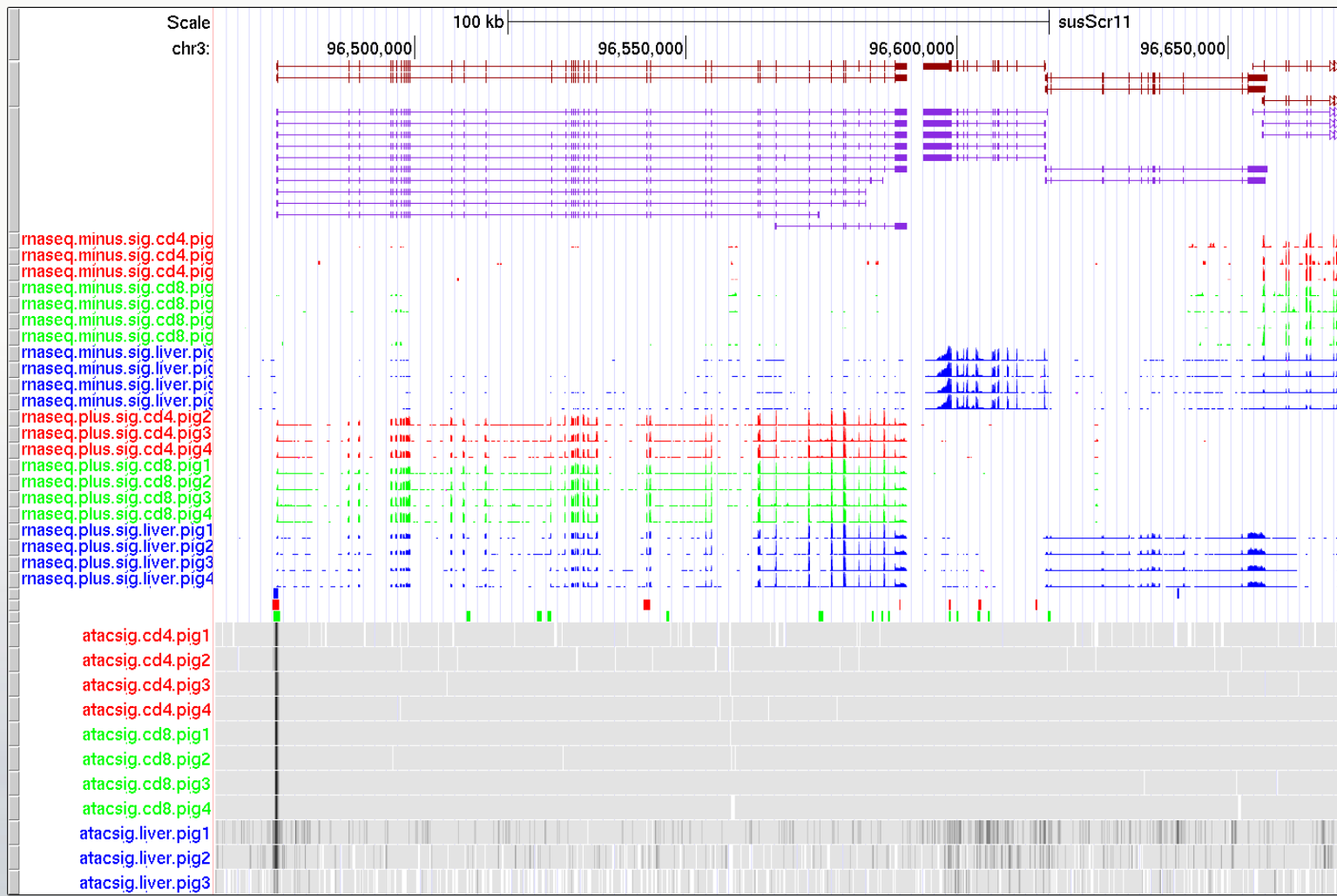


UCSC Genome Browser on Pig Feb. 2017 (Sscrofa11.1/susScr11) Assembly

move <<< << < > >> >>> zoom in 1.5x 3x 10x base zoom out 1.5x 3x 10x 100x

chr3:96,463,151-96,670,740 207,590 bp. enter position, gene symbol or search terms

chr3 3



move start < 2.0 > Click on a feature for details. Click or drag in the base position track to zoom in. Click side bars for track options. Drag side bars or labels up or down to reorder tracks. Drag tracks left or right to new position. Press "?" for keyboard shortcuts. move end < 2.0 >

track search default tracks default order hide all add custom tracks track hubs configure multi-region reverse resize refresh

Use drop-down controls below and press refresh to alter tracks displayed

- ◆ most of all the known transcripts are detected in liver and T cells
- ◆ reference annotation can be extended by a factor 2 to 3
- ◆ most of the new transcripts are alternative isoforms of coding genes

Reference and FR-AgENCODE transcripts detected

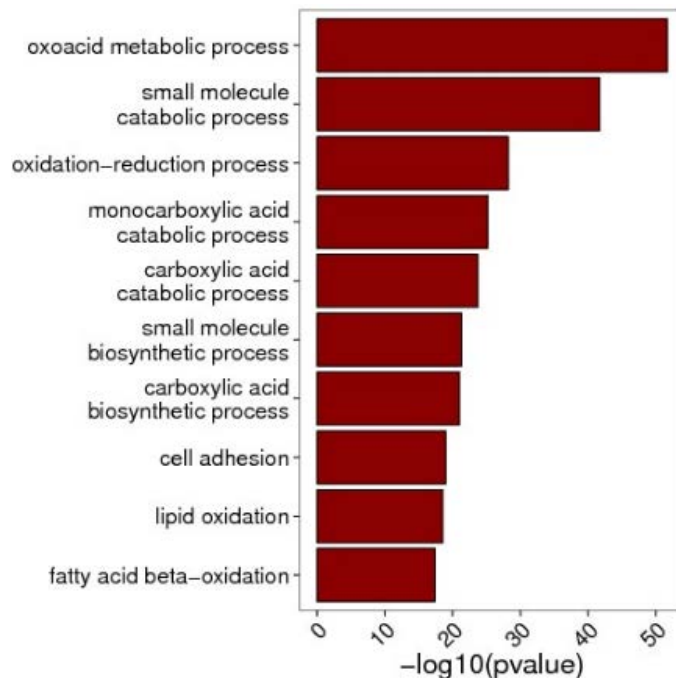
Species	Reference transcripts			FR-AgENCODE transcripts		
	All	Expressed		#	mRNAs	lncRNAs
		#	% of total			
Cattle	26,740	16,100	60.2	84,971	59,801	22,724
Goat	53,266	34,442	64.7	78,091	64,962	13,864
Chicken	38,118	22,180	58.2	57,817	47,567	7,502
Pig	49,448	29,786	60.2	77,540	63,721	12,587

- ◆ most of all the known transcripts are detected in liver and T cells
- ◆ reference annotation can be extended by a factor 2 to 3
- ◆ most of the new transcripts are alternative isoforms of coding genes
- ◆ differential expression between liver and T cells:
 - ◆ most of the genes are Differentially Expressed (DE)

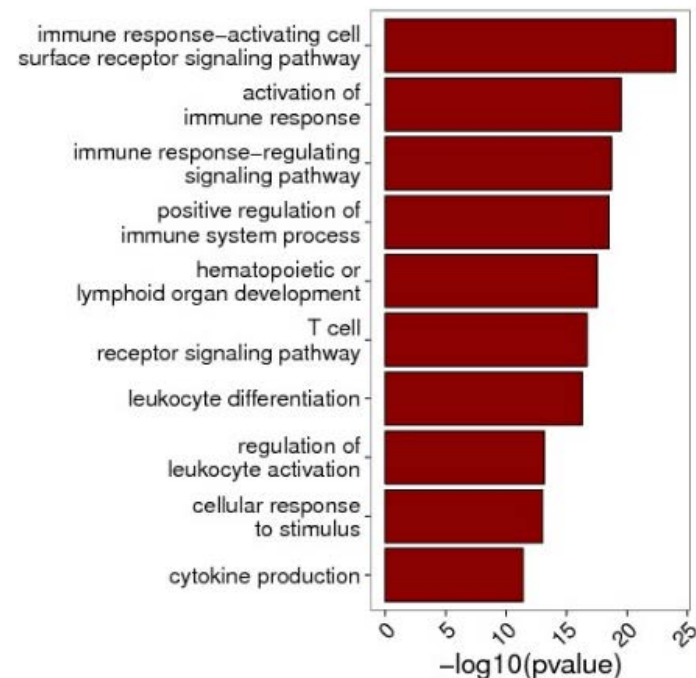
		Number of DE reference genes			
		Cattle	Goat	Chicken	Pig
liver > T cells		4,992	6,188	4,307	5,666
T cells > liver		3,943	4,384	2,640	3,772

- ◆ most of all the known transcripts are detected in liver and T cells
- ◆ reference annotation can be extended by a factor 2 to 3
- ◆ most of the new transcripts are alternative isoforms of coding genes
- ◆ differential expression between liver and T cells:
 - ◆ most of the genes are Differentially Expressed (DE)
 - ◆ DE genes have consistent GO functions

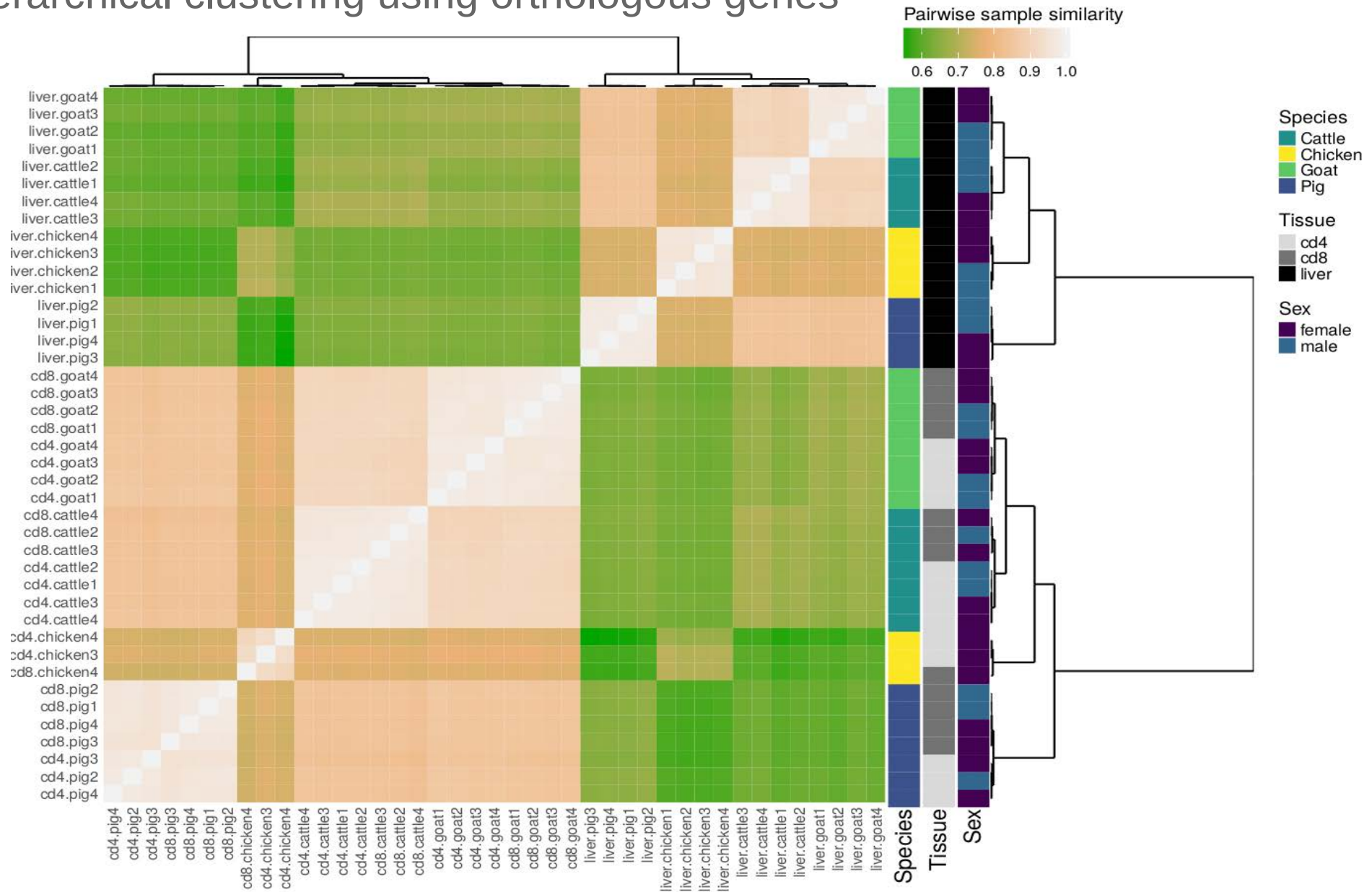
liver > T cells



T cells > liver

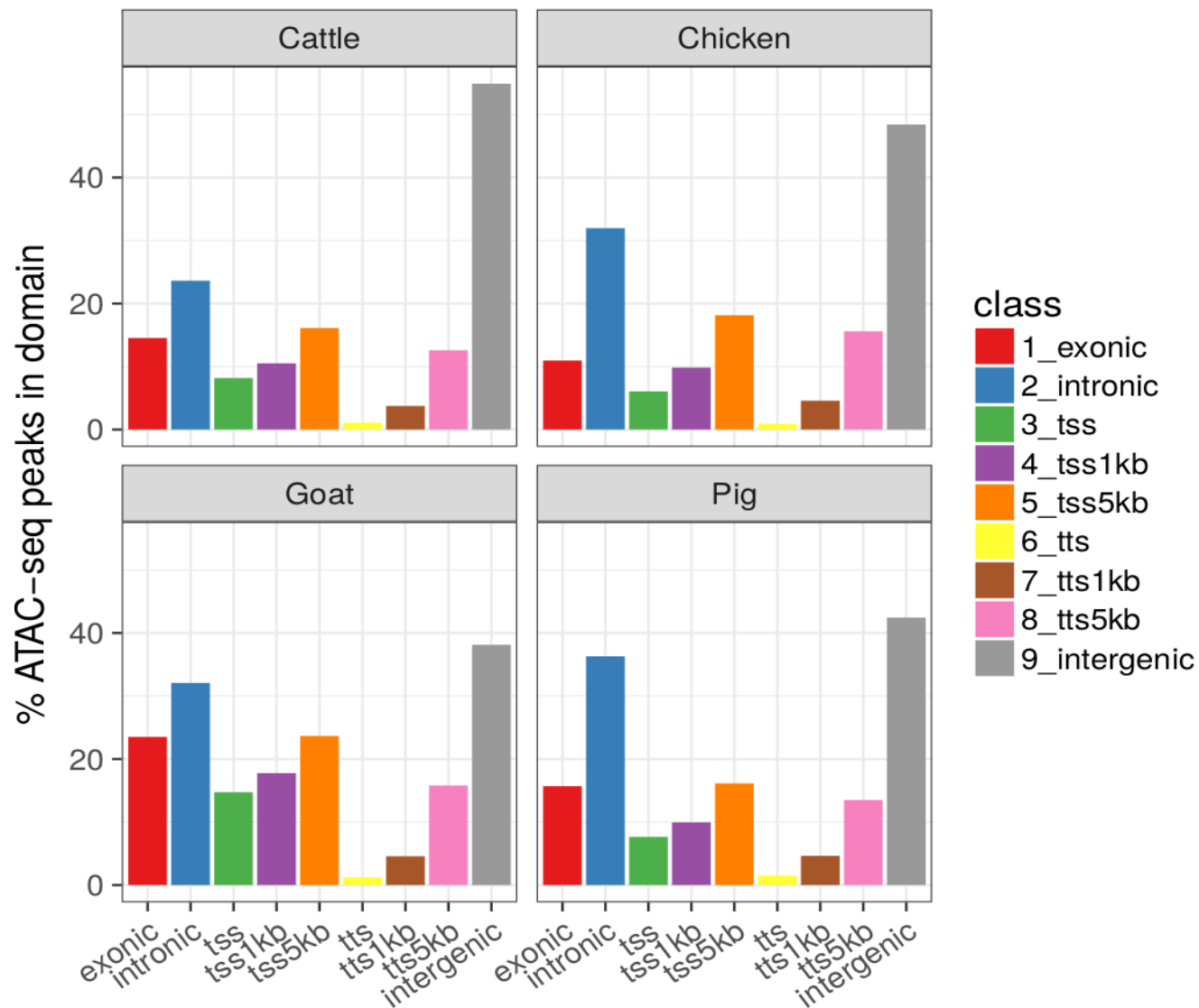


Hierarchical clustering using orthologous genes

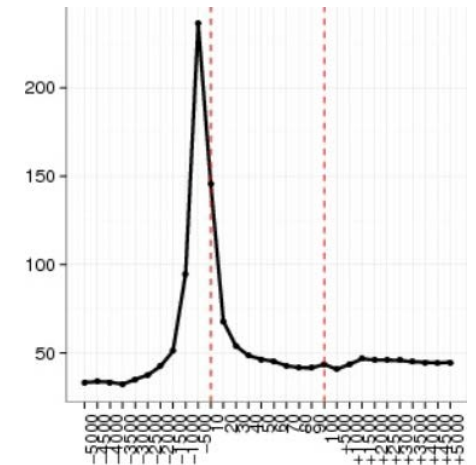
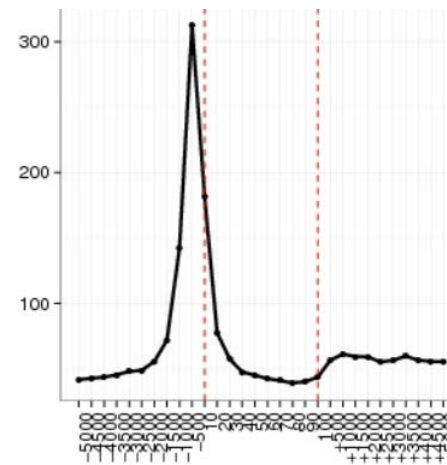
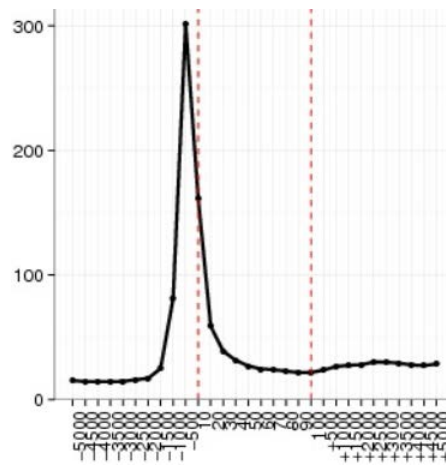
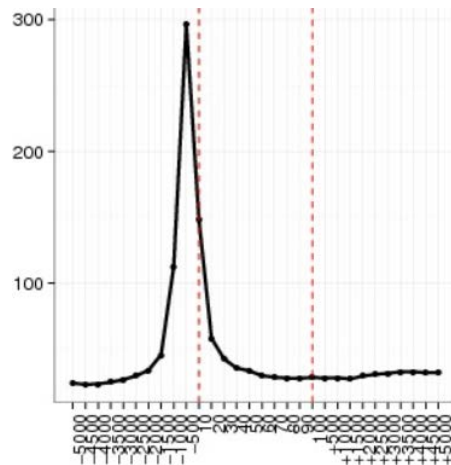


Samples cluster first by tissue, then by species

- ◆ 75,000-150,000 accessibility sites by species (~2-4% of the genome)
- ◆ Most of them are intergenic & intronic



- ◆ 75,000-150,000 accessibility sites by species (~2-4% of the genome)
- ◆ Most of them are intergenic & intronic
- ◆ Promoter accessibility: max within 1Kb upstream of gene starts

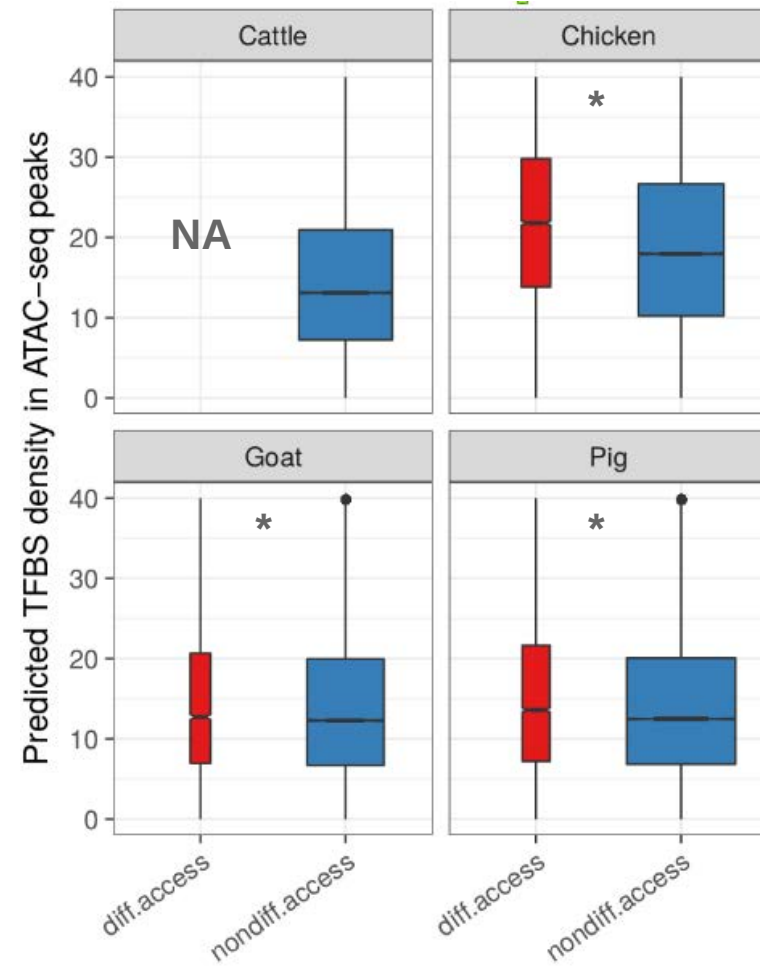


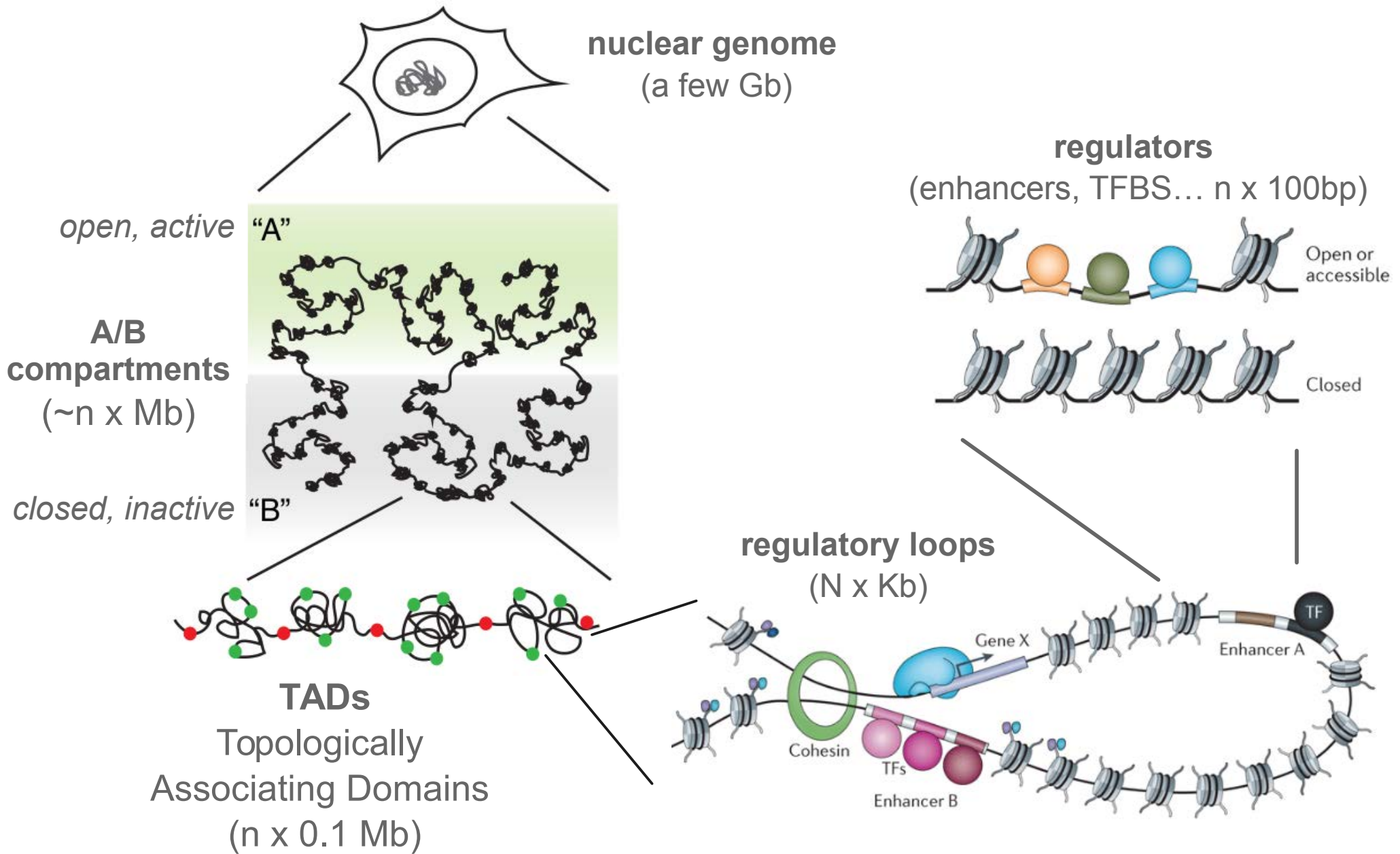
Mean ATAC-seq score around and within genes

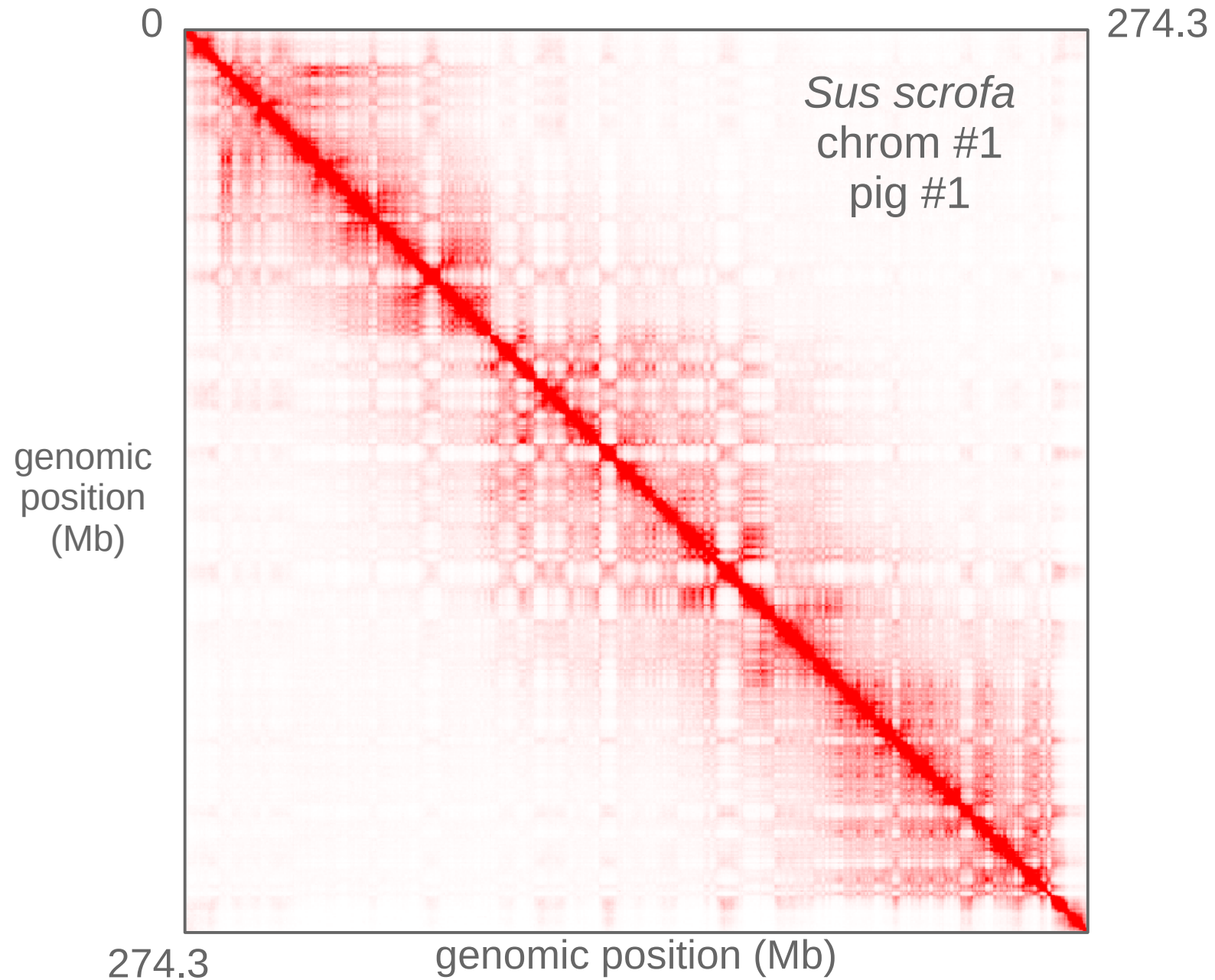
- ◆ 75,000-150,000 accessibility sites by species (~2-4% of the genome)
- ◆ Most of them are intergenic & intronic
- ◆ Promoter accessibility: max within 1Kb upstream of gene starts
- ◆ Comparative analysis between liver and T cells
 - ◆ 5,000 to 13,000 sites are Differentially Accessible (DA) by species

- ◆ 75,000-150,000 accessibility sites by species (~2-4% of the genome)
- ◆ Most of them are intergenic & intronic
- ◆ Promoter accessibility: max within 1Kb upstream of gene starts
- ◆ Comparative analysis between liver and T cells
 - ◆ 5,000 to 13,000 sites are Differentially Accessible (DA) by species
 - ◆ more Transcription Factor Binding Sites in DA vs. non-DA sites

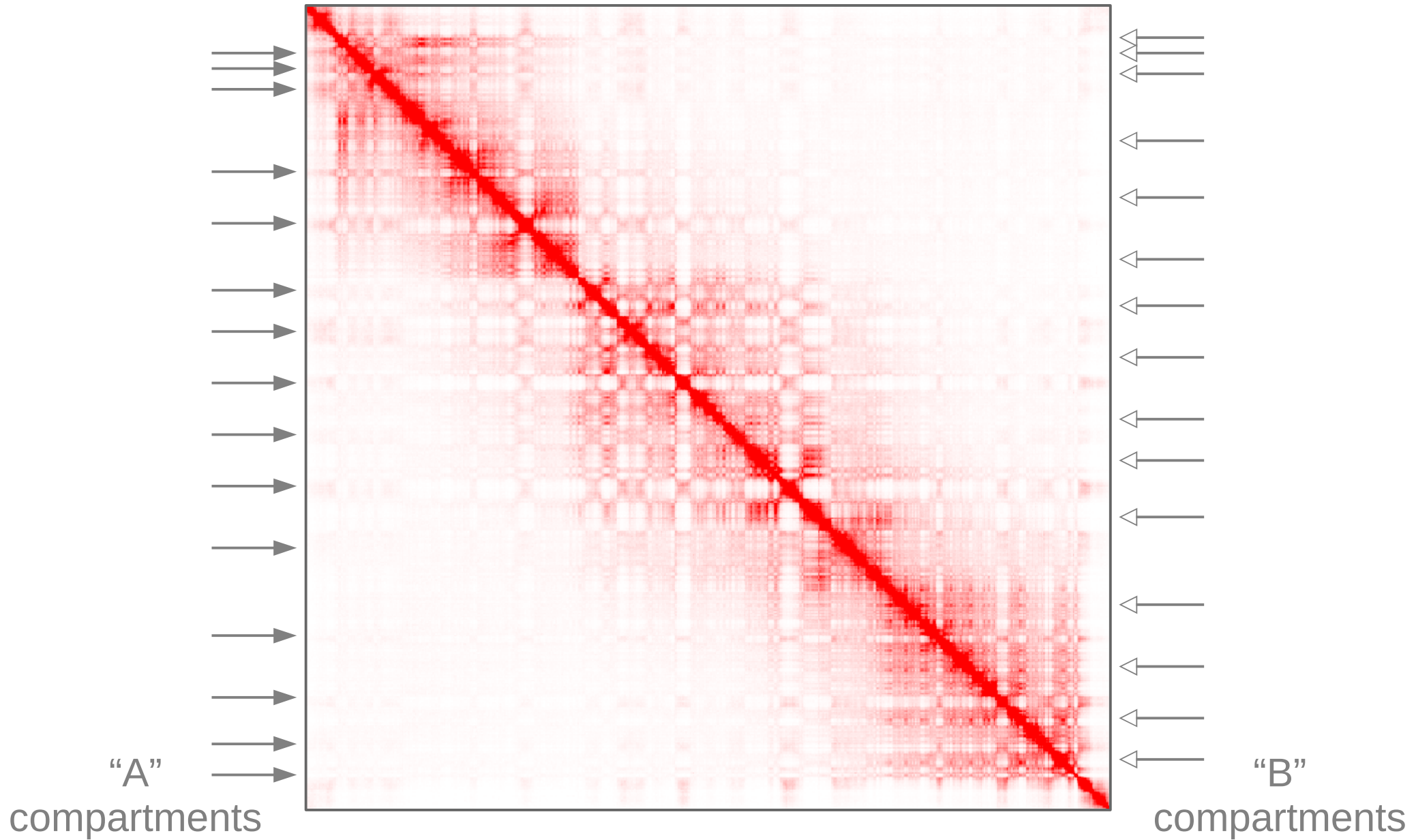
=> support for a regulatory role



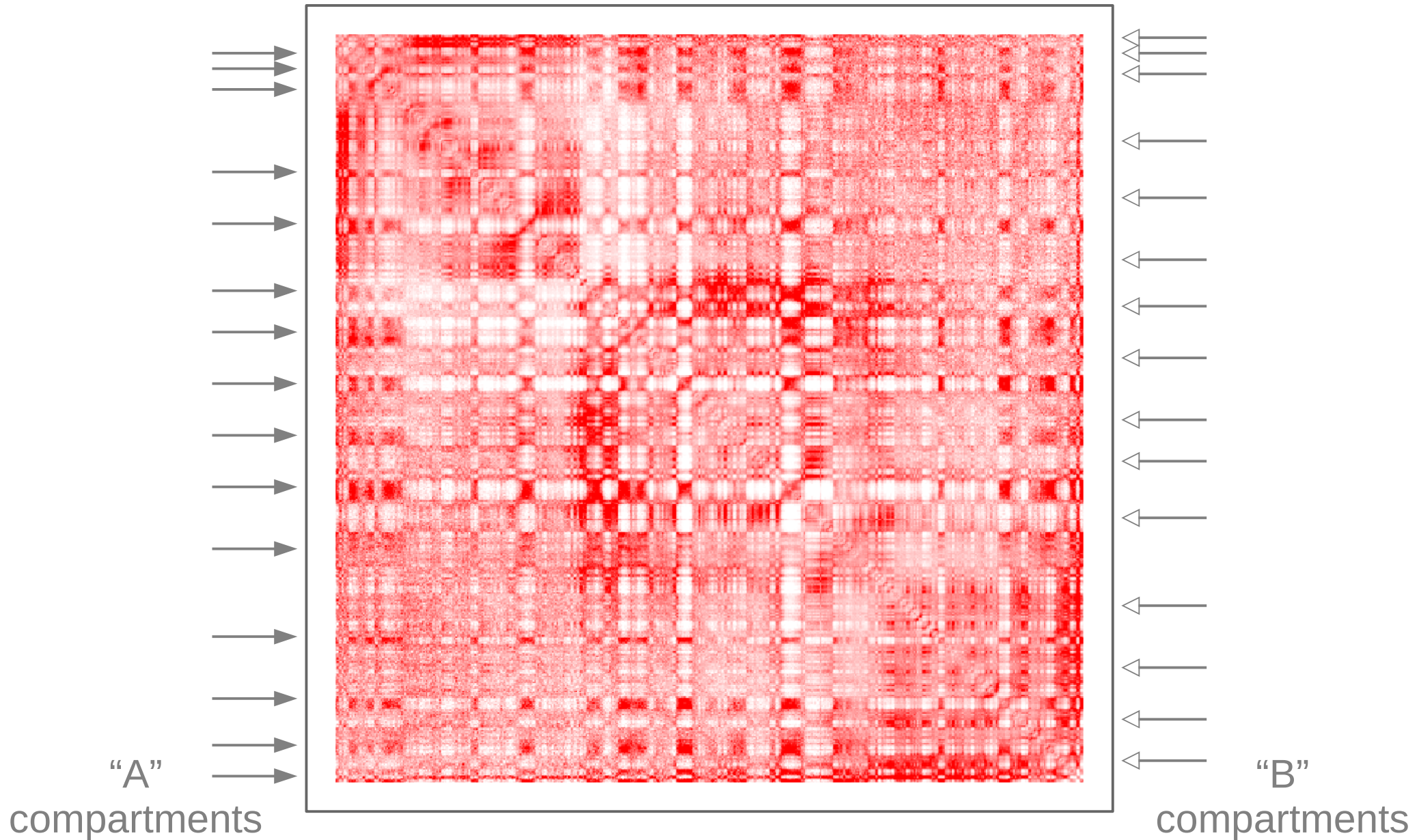




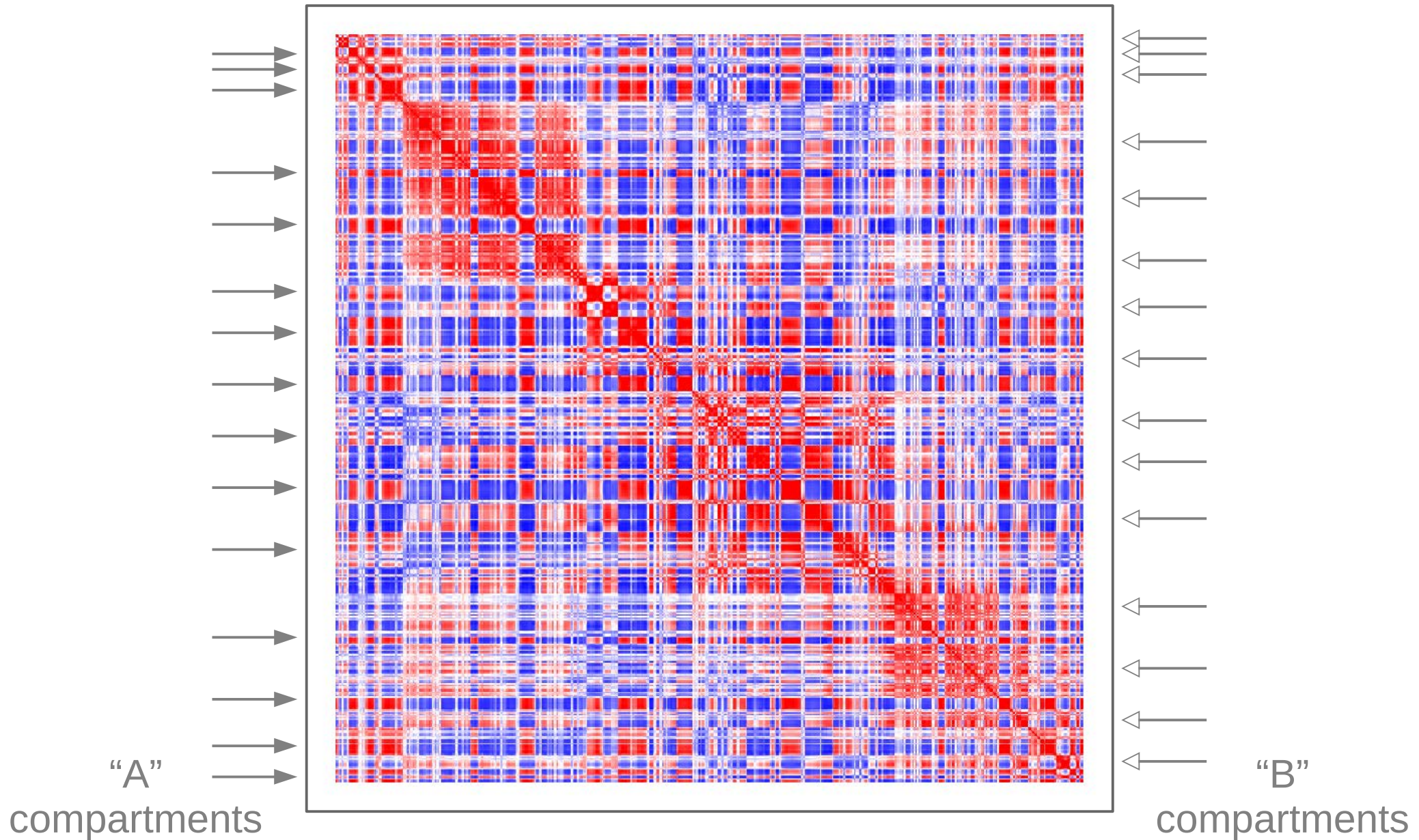
The A/B (epi)genome compartments



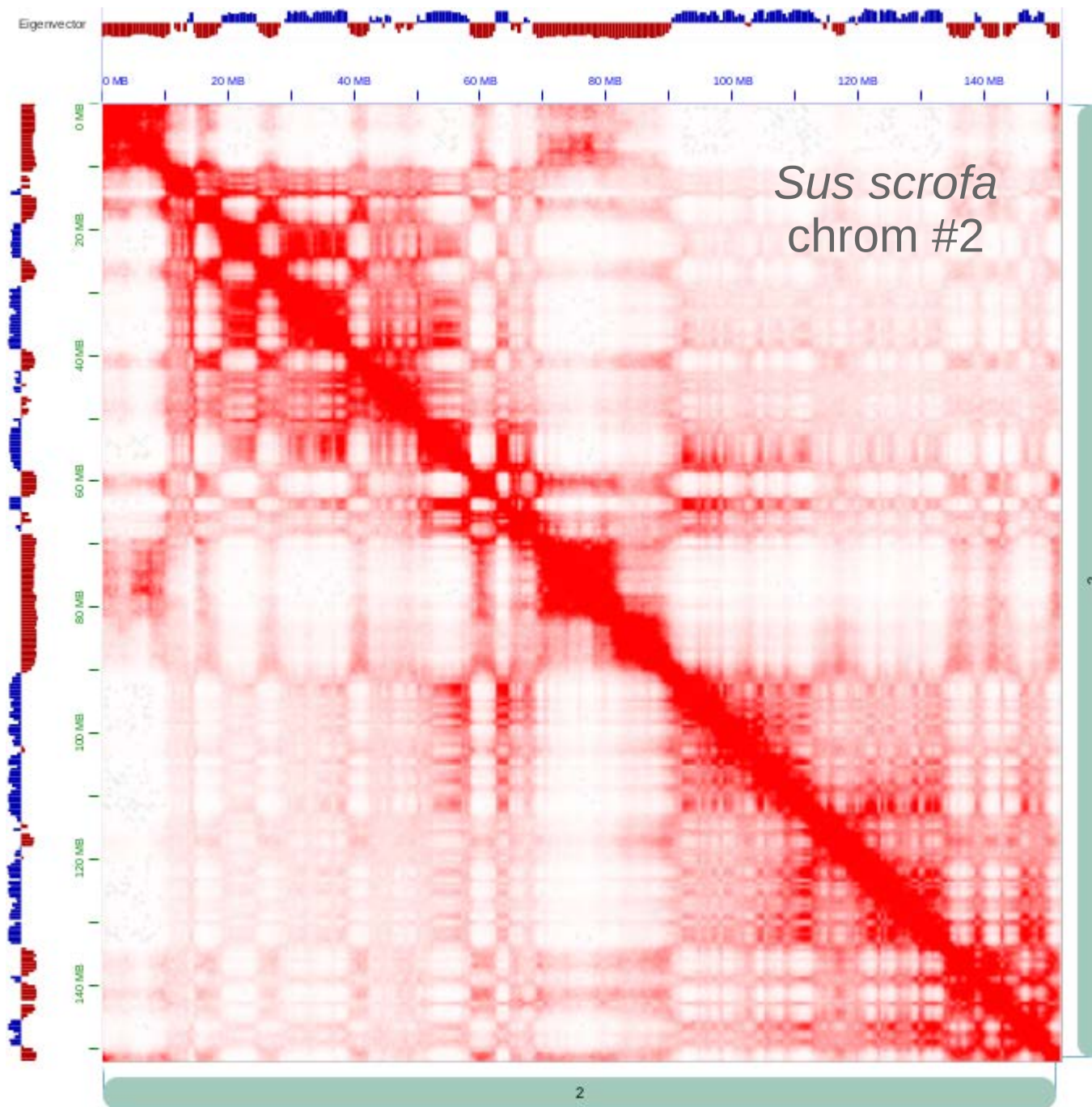
The A/B (epi)genome compartments



The A/B (epi)genome compartments

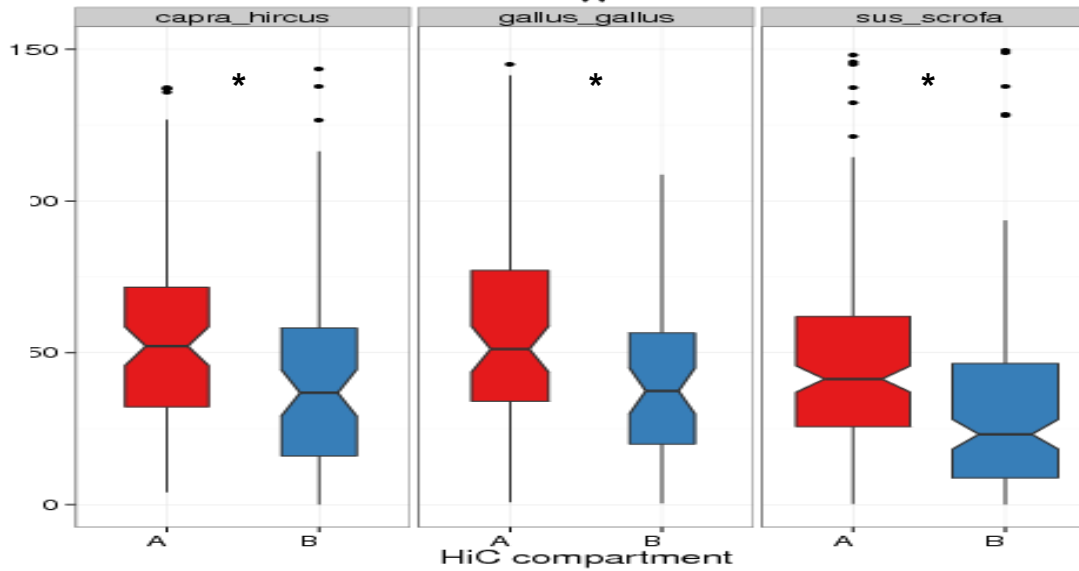


The A/B (epi)genome compartments

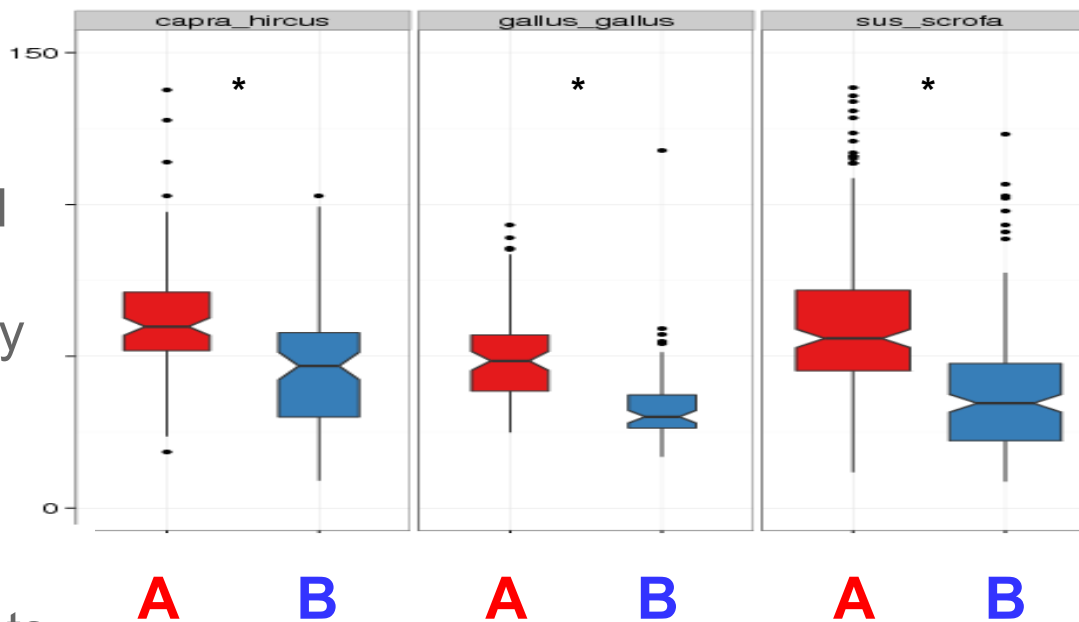




RNA-seq
gene
expression



ATAC-seq
chromatin
accessibility



=>

Links
between
genome
structure
& function

A/B
compartments

- ◆ Gene distribution in A/B compartments

Goat: 7844 genes in A (**65.5%**) vs. 4133 in B (**34.5%**)

Chicken: 4571 genes in A (**64.0%**) vs. 2576 in B (**36.0%**)

Pig: 6737 genes in A (**63.4%**) vs. 3883 in B (**36.6%**)

◆ Gene distribution in A/B compartments

Goat: 7844 genes in A (**65.5%**) vs. 4133 in B (**34.5%**)

Chicken: 4571 genes in A (**64.0%**) vs. 2576 in B (**36.0%**)

Pig: 6737 genes in A (**63.4%**) vs. 3883 in B (**36.6%**)

◆ Focusing on orthologous genes: compartments across species

“A” in the 3 species

expected: 1529 genes (**26.7%**)

observed: 2972 genes (**51.9%**)

“B” in the 3 species

expected: 259 genes (**4.5%**)

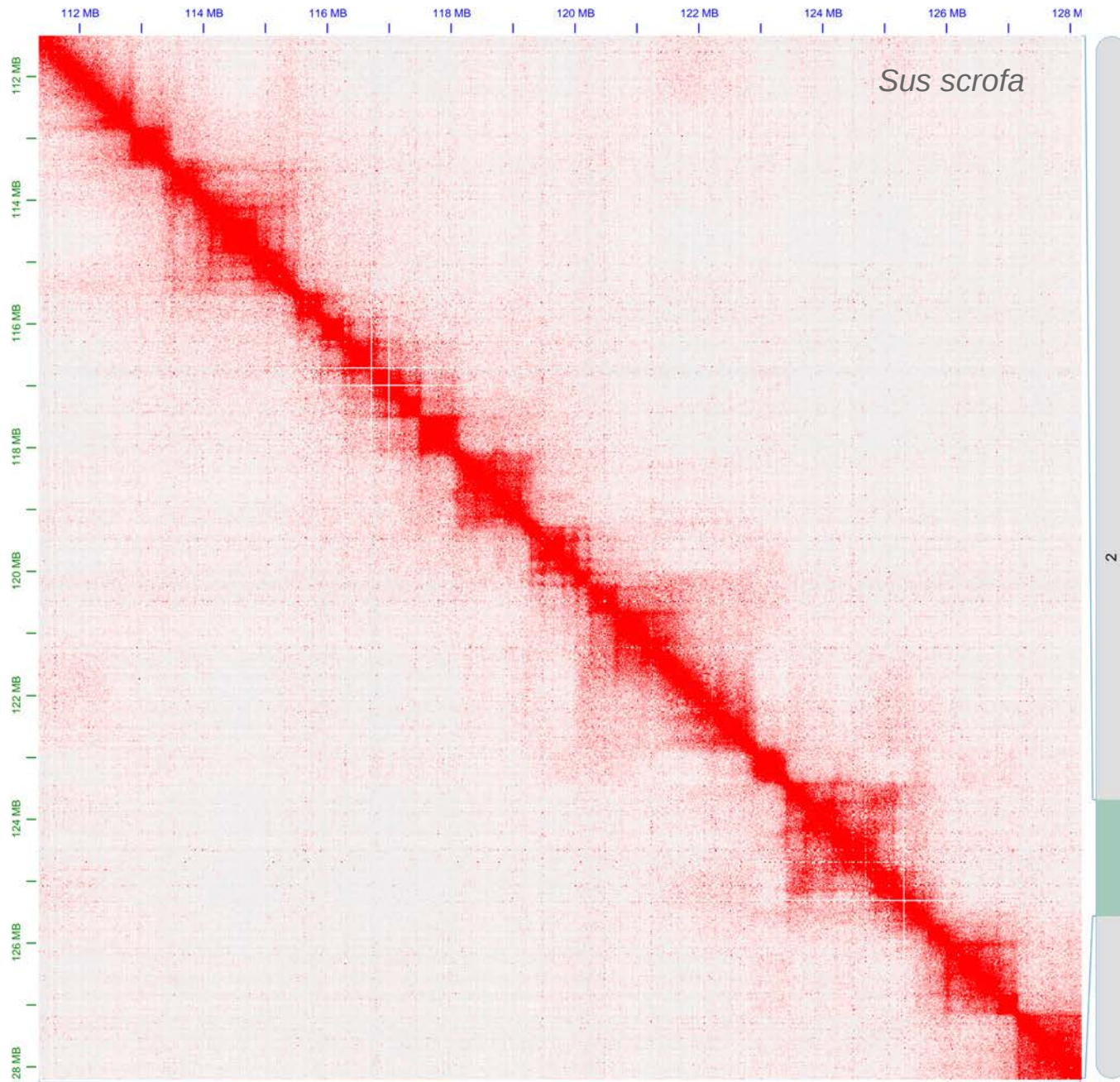
observed: 611 genes (**10.7%**)

(N=5728 orthologous genes with an assigned compartments in the 3 species)

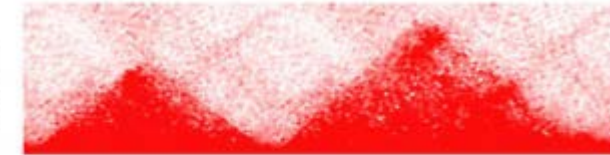
**conservation of
genome compartments**

=>

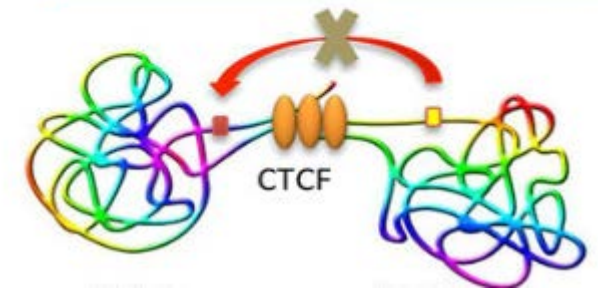
**evidence of
functional role**



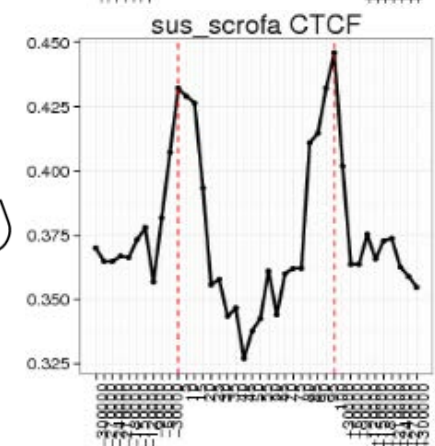
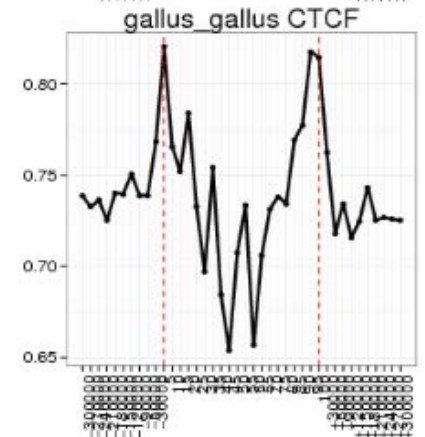
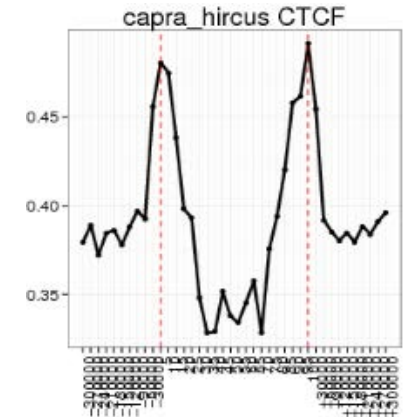
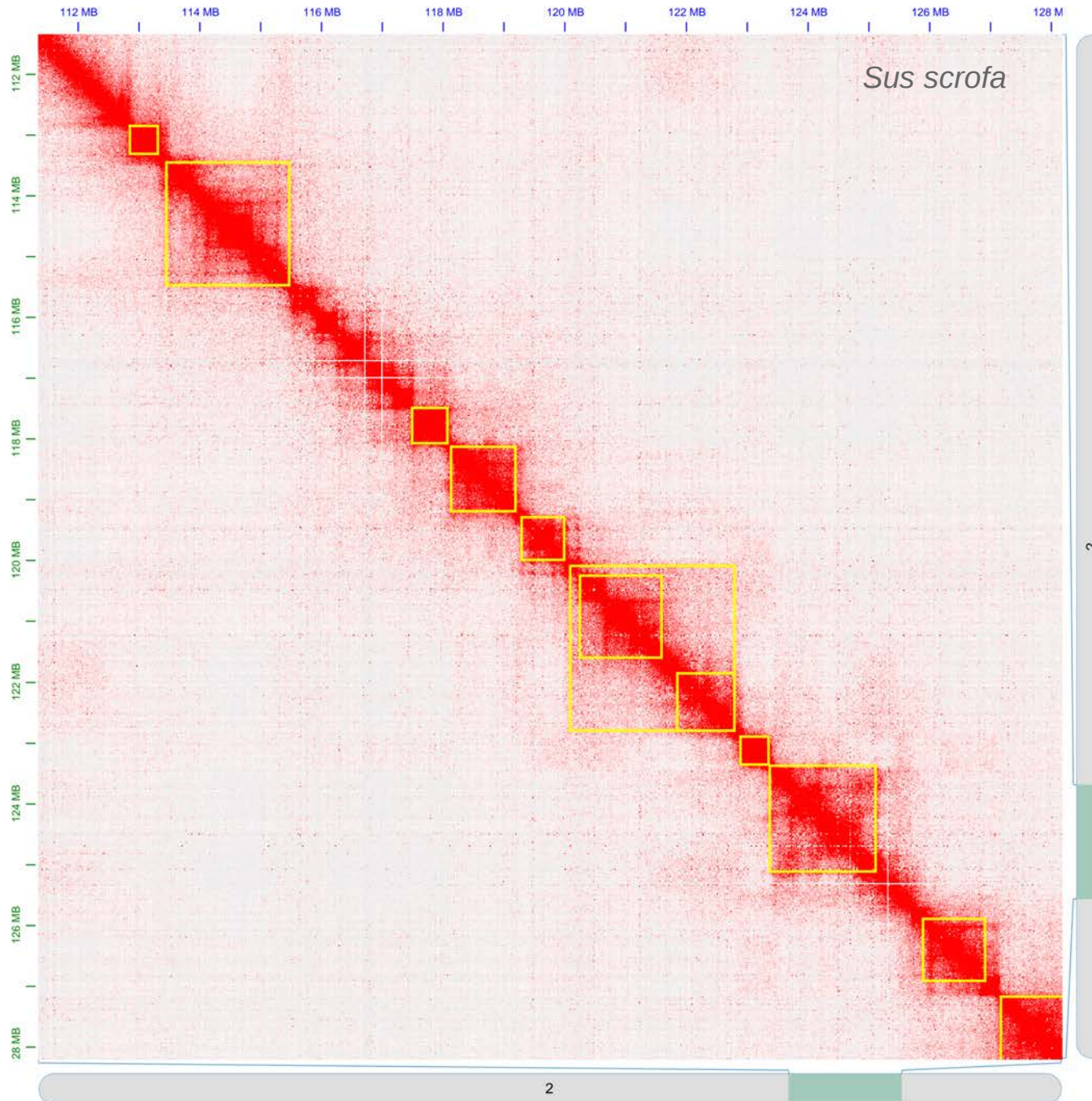
Hi-C contact matrix



3D model



Li et al, 2016, Scientific Reports



- ◆ Identify “orthologous” TAD boundaries (between 2 or 3 species):

N = 16,468 non ambiguous hits

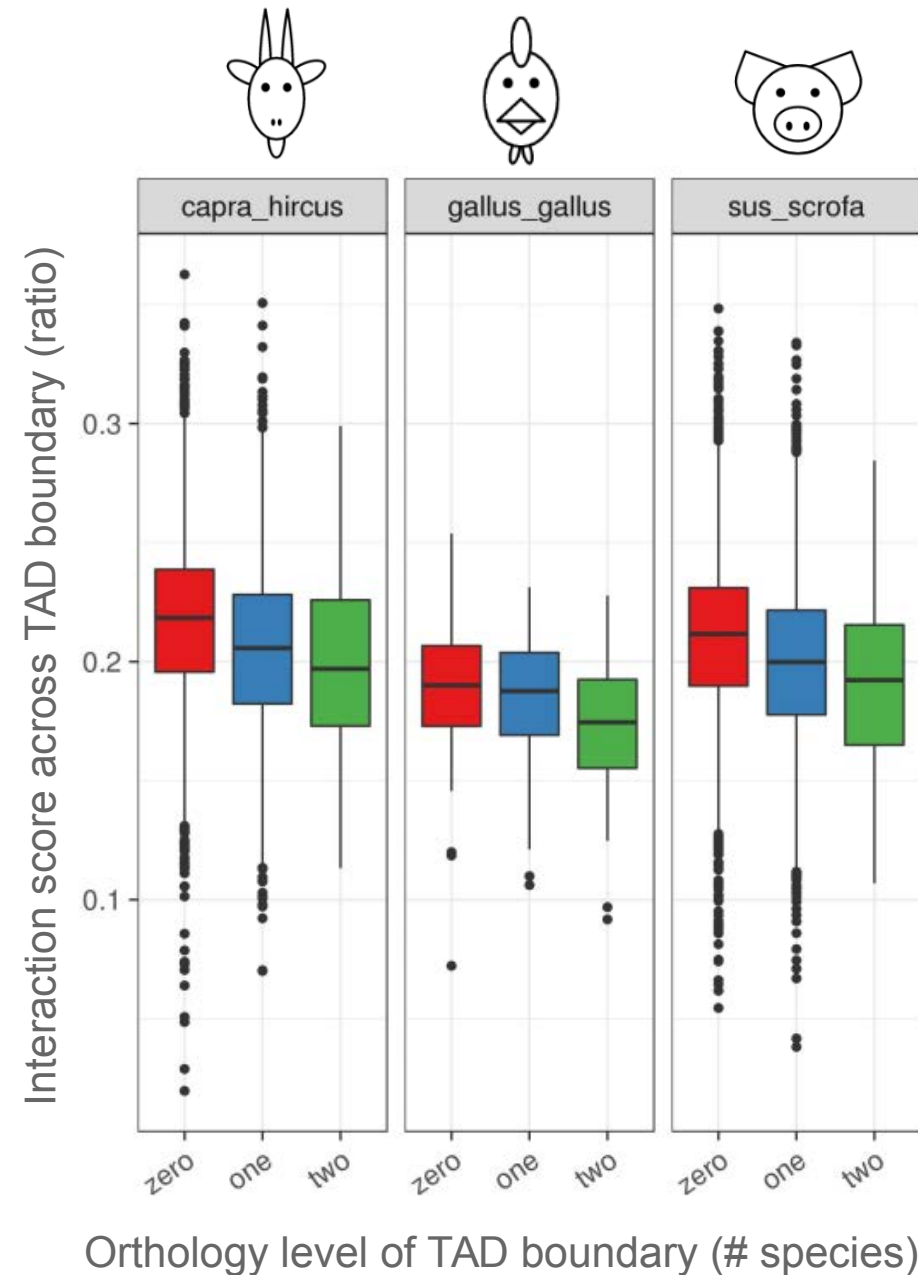
- ◆ 10,805 from 1 species => level zero
 - ◆ 5592 from 2 species => level one
 - ◆ 71 from 3 species => level two
-
- ◆ Compute interaction score across TAD boundaries of each level (the lower the score, the stronger the insulation)

- Identify “orthologous” TAD boundaries (between 2 or 3 species):

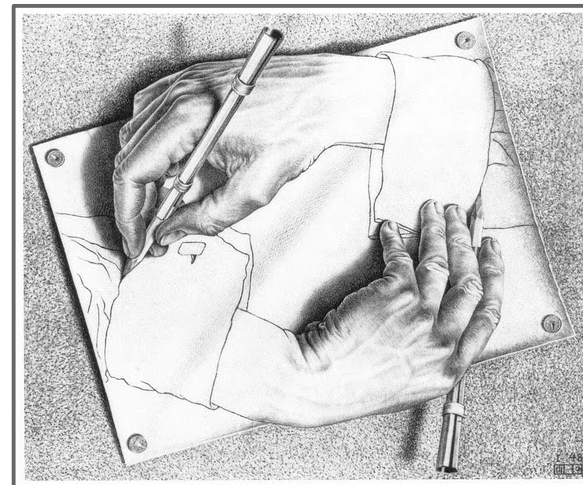
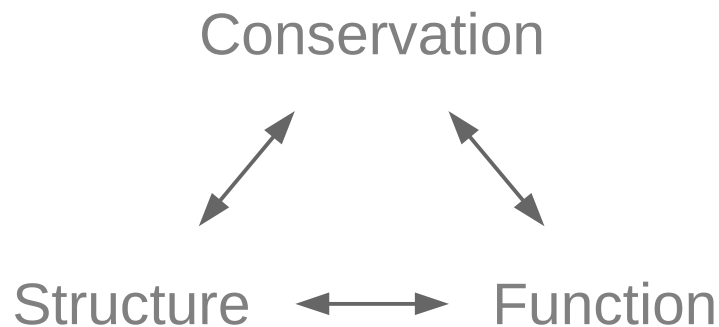
N = 16,468 non ambiguous hits

- 10,805 from 1 species => level zero
 - 5592 from 2 species => level one
 - 71 from 3 species => level two
- Compute interaction score across TAD boundaries of each level (the lower the score, the stronger the insulation)

conserved TAD boundaries have stronger insulations => evidence of selective pressure



- ◆ The FR-AgENCODE project contributes to the FAANG action
- ◆ Substantial extension of the transcriptional map
- ◆ Identification of potential regulatory sites with accessible chromatin
- ◆ Integrative analysis across assays and/or across species:
a powerful approach to investigate gene expression
- ◆ Chromatin conformation is under selective pressure at various organizational levels: accessibility sites, TADs & compartments



Management

Elisabetta Giuffra (coordination)
 Sylvain Foissac (coordination)
 Sandrine Lagarrigue
 Marie-Hélène Pinard-Van der Laan



Sampling and assays

Hervé Acloque
 Cécile Berri
 Fany Blanc
 Sophie Dhorne-Pollet
 Françoise Drouet
 Diane Esquerre
 Stéphane Fabre
 Joël Gautron
 Adeline Goubil
 Sonia Lacroix-Lamandé
 Fabrice Laurent
 Florence Mompert
 Pascale Queré
 Michèle Tixier-Boichard
 Gwenola Tosser-Klopp
 Silvia Vincent-Naulleau

[...]



Data analysis

Philippe Bardou
 Cédric Cabau
 Elisa Crisci
 Thomas Derrien
 Sarah Djebali-Quelen
 Sylvain Foissac
 Christine Gaspin
 Ignacio Gonzalez
 Christophe Klopp
 Sandrine Lagarrigue
 Sylvain Marthey
 Maria Marti-Marimon
 Raphaelle Momal-Leisenring
 Kylie Munyard
 Kévin Muret
 Andrea Rau
 David Robelin
 Magali San Cristobal
 Nathalie Vialaneix
 Matthias Zytnicki



+ Hi-C team @ Toulouse

Hervé Acloque
 Martine Bouissou-Matet
 Sarah Djebali
 Yvette Lahbib
 Laurence Liaubet
 Maria Marti
 Florence Mompert
 Pierre Neuvial
 David Robelin
 Nathalie Vialaneix
 Alain Vignal
 Matthias Zytnicki

+ INRA Platforms & Facilities

@BRIDGE biorepository
 GeT-PlaGe sequencing
 GenoToul bioinformatics
 GenoToul biostatistics
 Experimental & animal facilities UE Le Pin, Bourges, GenESI, PEAT, PAO ; UR BOA, PRC CIRE

SelGen INRA metaprogram

