

# Profiling the landscape of transcription, chromatin accessibility and chromosome conformation of cattle, pig, chicken and goat genomes

*[FAANG pilot project “FR-AgENCODE”]*

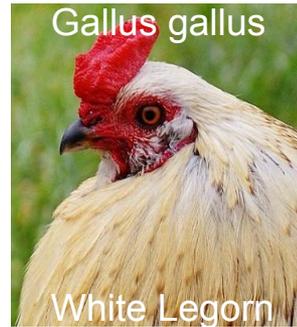
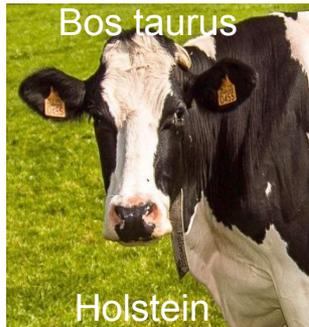
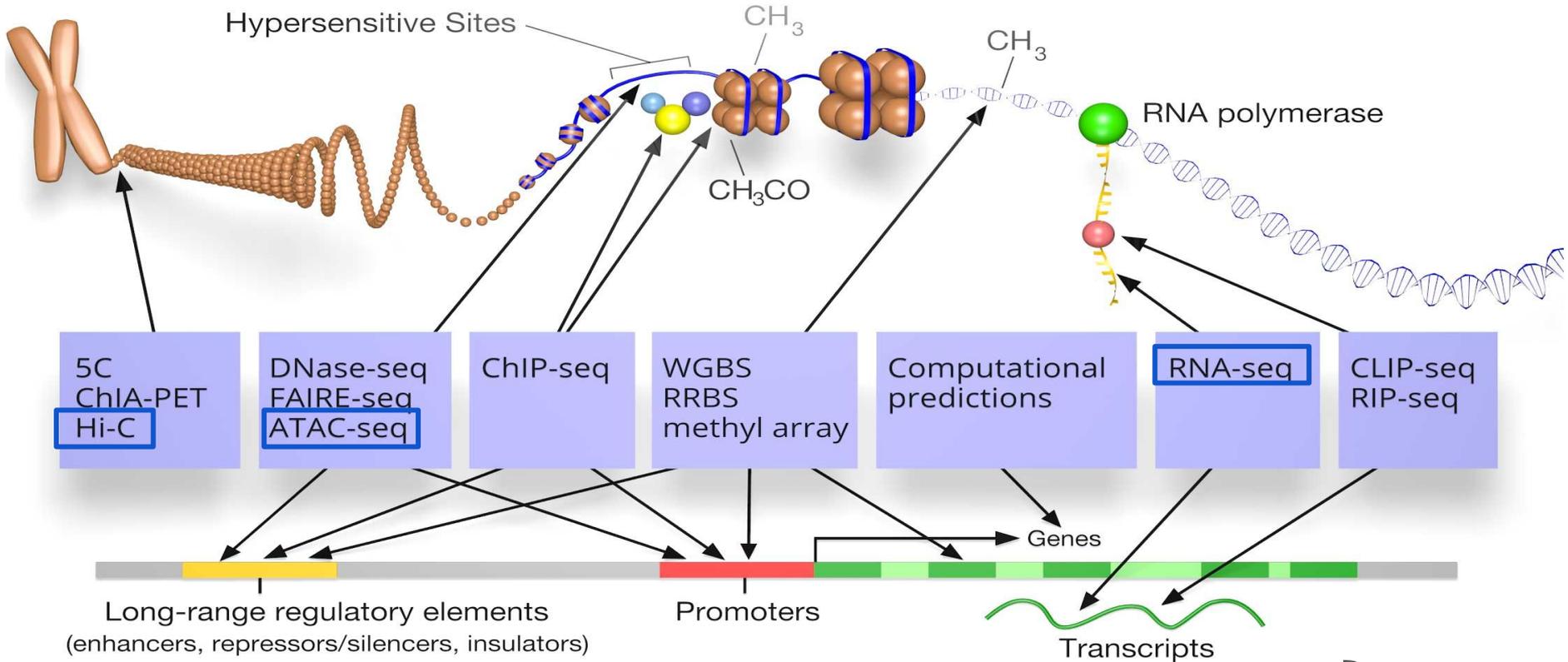
**Sarah Djebali**

INRA - Toulouse - France

Genetics, Physiology and Breeding Systems Laboratory



# FR-AgENCODE data



2 males  
2 females  
Liver  
CD4  
CD8

# FR-AgENCODER data generation progress

Species	ATAC-seq*			HiC (liver only)			Long RNA-seq			Small RNA-seq		
	Library preparation	Sequencing	Data processing									
Bos taurus			**						**			
Capra hircus			**			**			**			
Gallus gallus			**			**			**			
Sus scrofa			**			**			**			

\* No liver for bos\_taurus ATAC-seq / \*\* Integrative analysis on-going

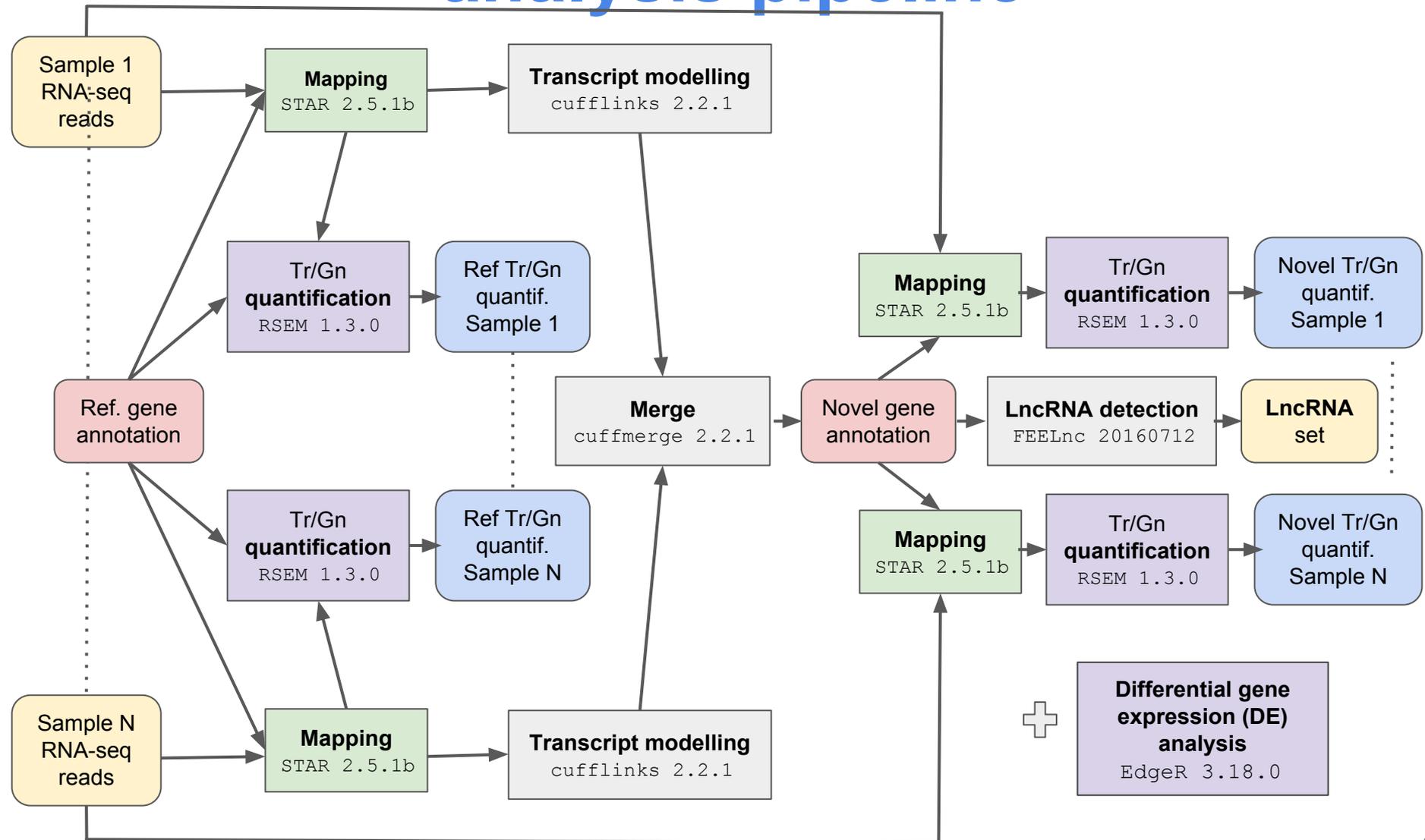
Our 16 animals are called cattle 1,2,3,4, goat 1,2,3,4, chicken 1,2,3,4 and pig 1,2,3,4:

- Animals 1 and 2 are **males**
- Animals 3 and 4 are **females**

# RNA-seq for expression and annotation of long coding and non-coding RNAs

- RNAs:
  - longer than 200bp
  - Poly-a-selection before cDNA synthesis
- Sequencing:
  - Directed
  - 2x150bp
  - 100 million read pairs per sample

# RNA-seq data processing and analysis pipeline



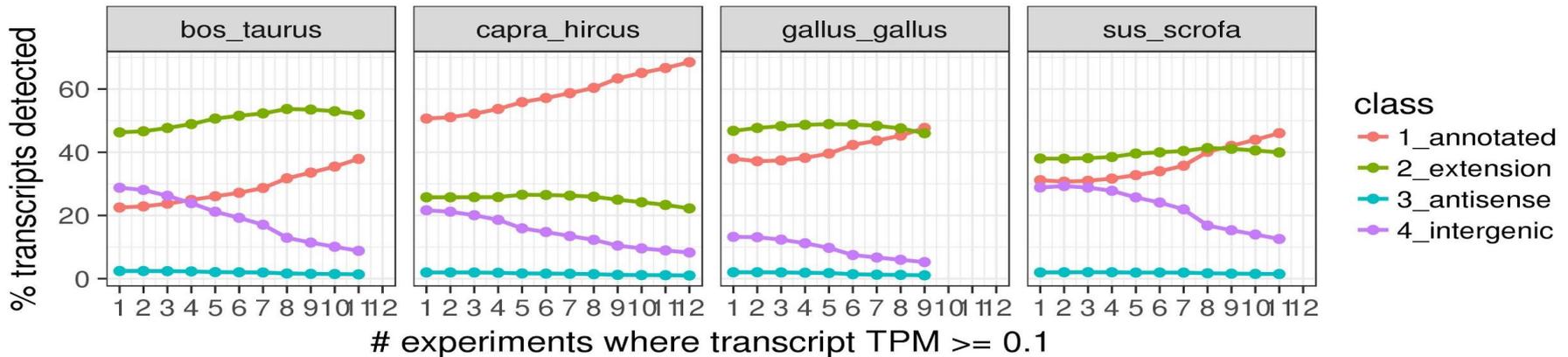
Nextflow implementation @ <https://github.com/skptic/LncRNA-Annotation-nf>

# Many novel transcripts are found, which are globally less expressed than the annotated ones

Species	Genome / Gene annotation	Annotated Transcripts			Number of novel transcripts detected*	Number of novel lncRNAs
		Total number	Detected*			
			#	% of total		
Bos taurus	UMD 3.1 / Ensembl 84	<b>26,740</b>	16,100	<b>60.2</b>	<b>65,538</b>	<b>7,929</b>
Capra hircus	CHIR_ARS 1 / NCBI	<b>53,266</b>	34,442	<b>64.7</b>	<b>38,197</b>	<b>NA</b>
Gallus gallus	GalGal 5 / Ensembl 87	<b>38,118</b>	22,898	<b>60.1</b>	<b>41,116</b>	<b>3,264</b>
Sus scrofa	SScrofa 10.2 / Ensembl 84	<b>30,585</b>	18,746	<b>61.3</b>	<b>57,978</b>	<b>6,581</b>

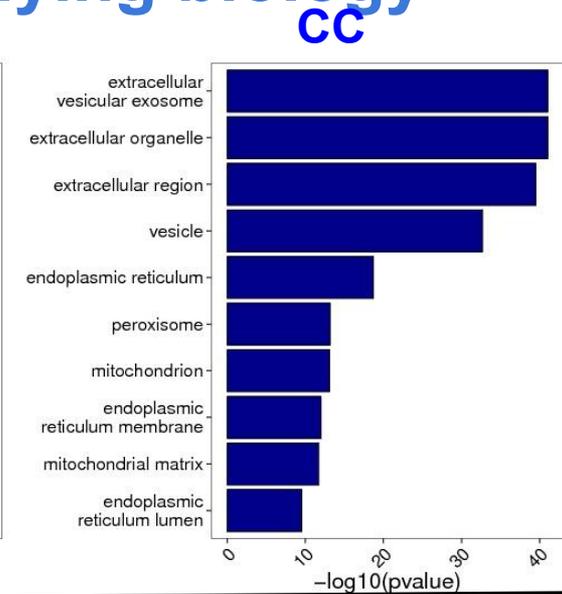
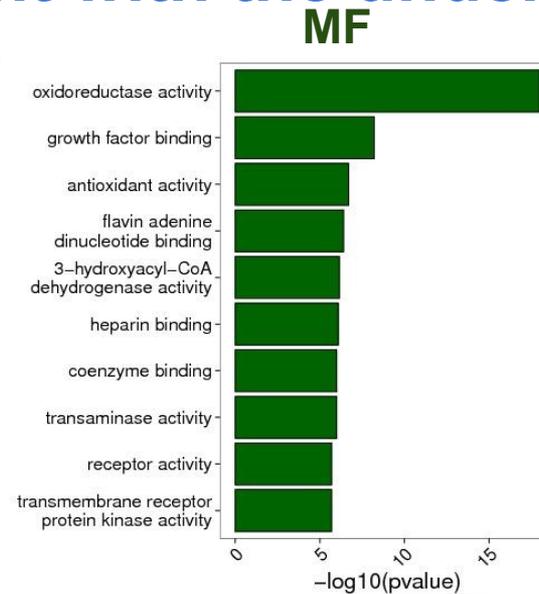
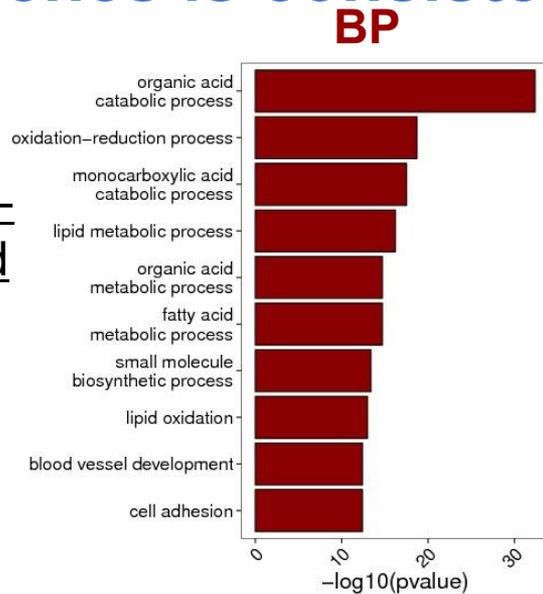
See Kevin Muret's talk just after this talk

\* with TPM  $\geq 0.1$  in  $\geq 2$  samples

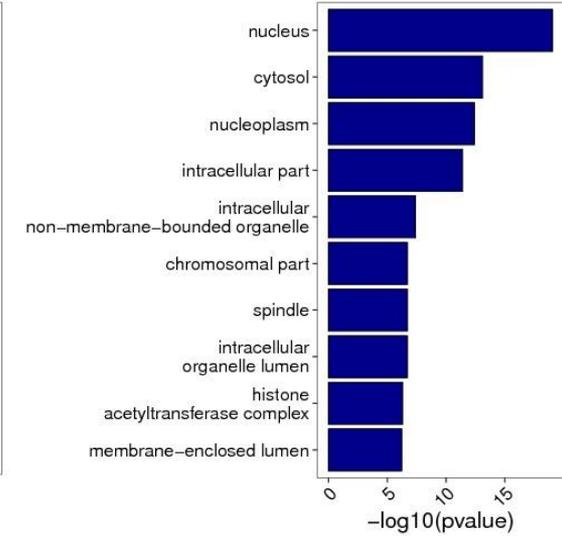
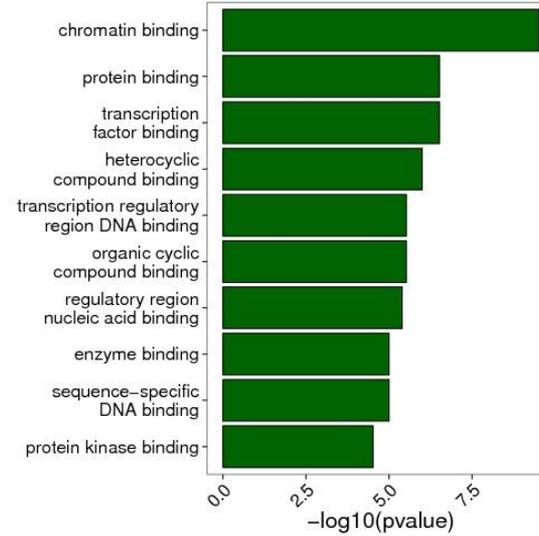
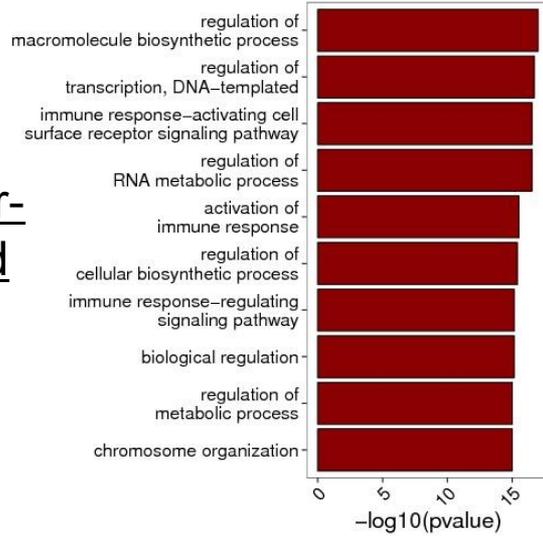


# Biological process (BP), molecular function (MF), cell compartment (CC) GO term enrichment on DE genes is consistent with the underlying biology

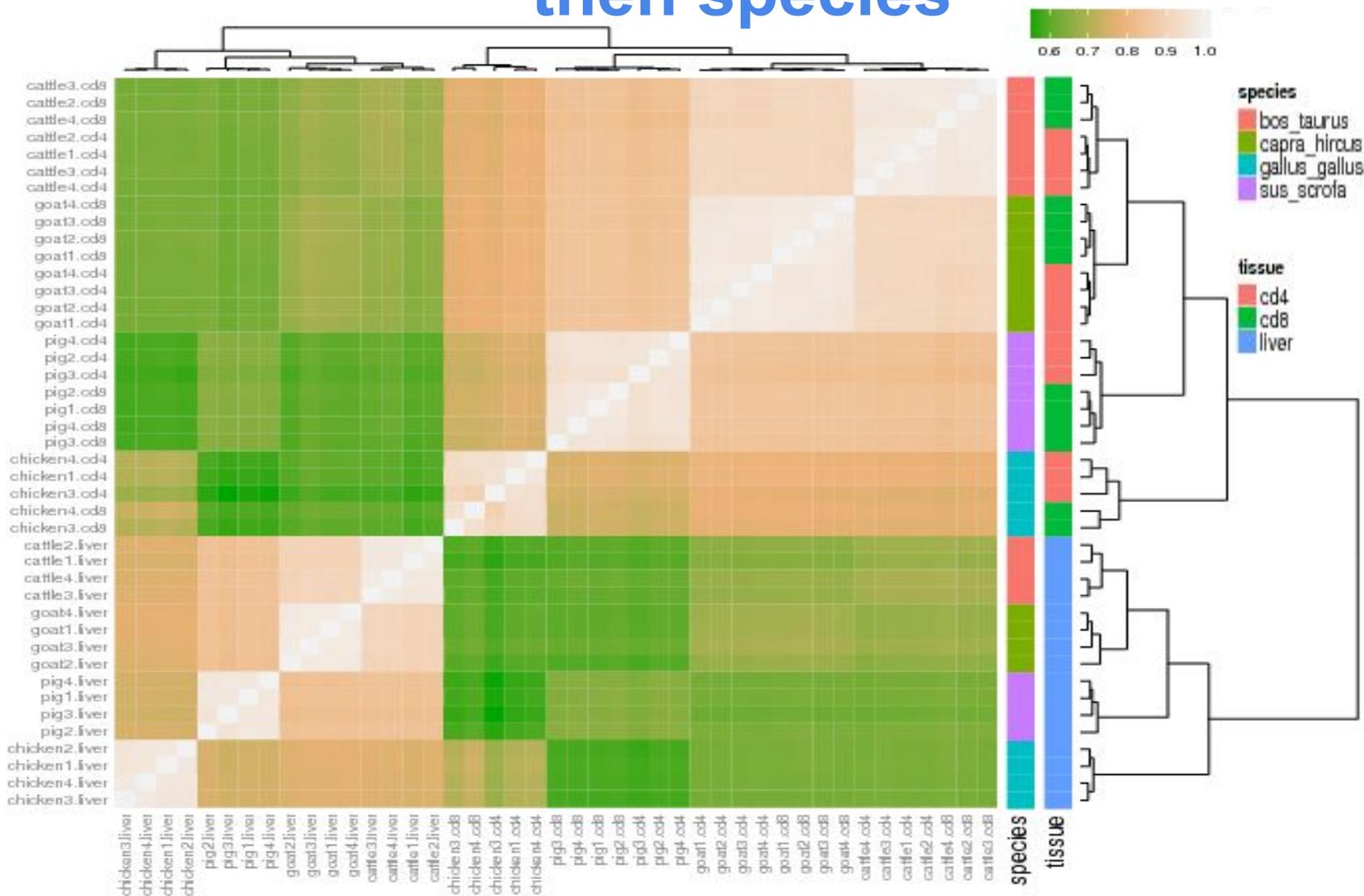
Liver over-expressed genes



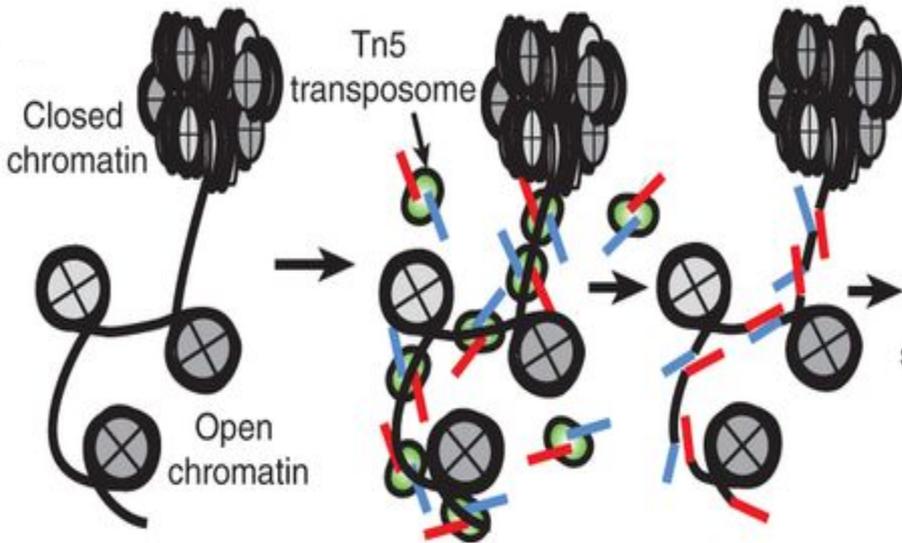
T cell over-expressed genes



# All species RNA-seq hierarchical clustering first separates liver from immune cells, and then species

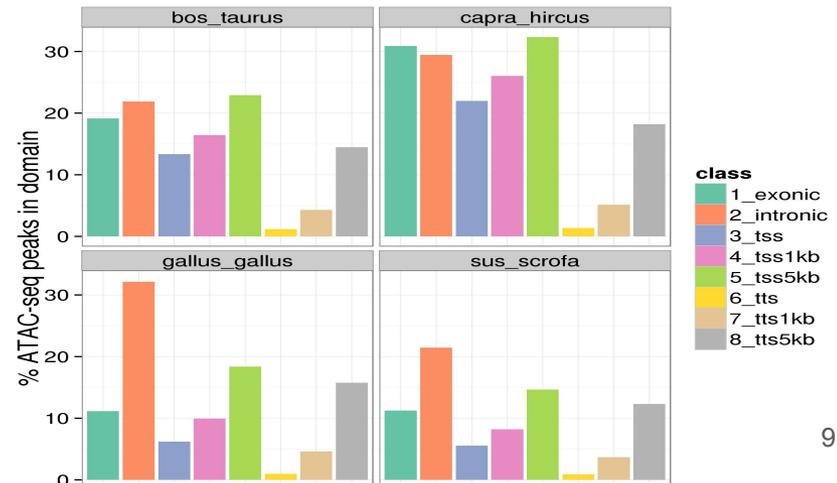
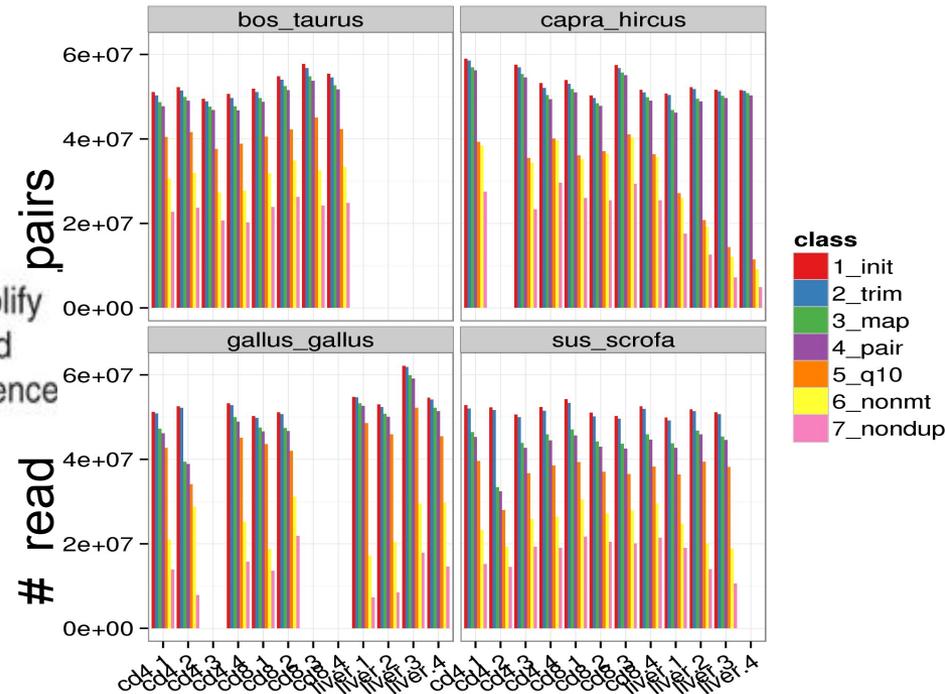


# ATAC-seq for open chromatin regions



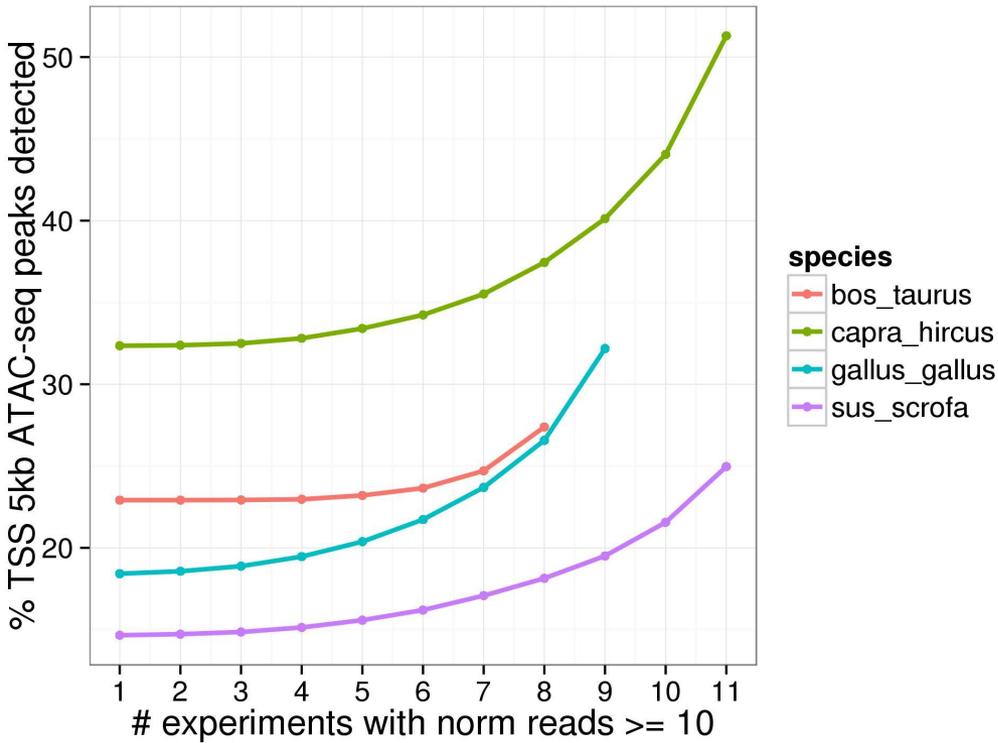
Buenrostro et al., Nature Methods, 2013

50 million read pairs per sample

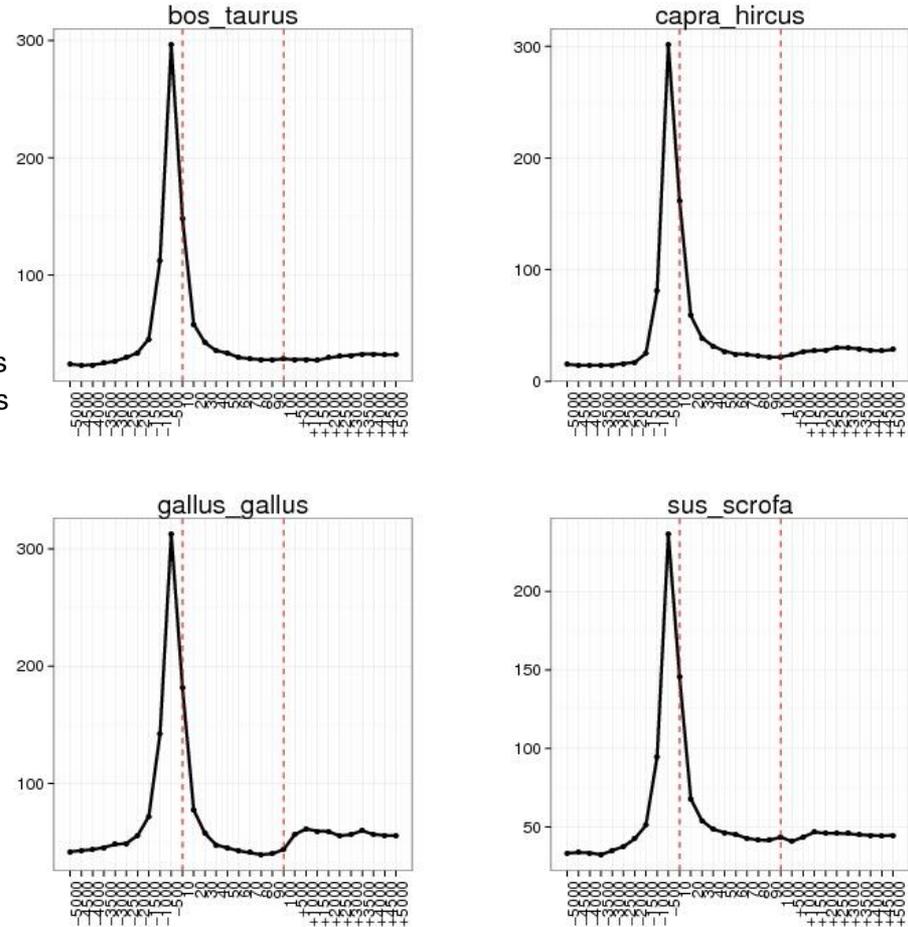


Species	Number of atac-seq peaks	Genome size (bp)	ATAC-seq peak coverage	
			# bp	% of genome
Bos taurus	<b>58,384</b>	2,670,422,299	42,288,741	<b>1.58</b>
Capra hircus	<b>46,901</b>	2,922,813,246	32,951,265	<b>1.13</b>
Gallus gallus	<b>116,893</b>	1,230,258,557	50,931,713	<b>4.14</b>
Sus scrofa	<b>120,914</b>	2,808,525,991	72,480,471	<b>2.58</b>

# ATAC-seq peaks at Transcription Start Sites (TSS); those peaks are more ubiquitous



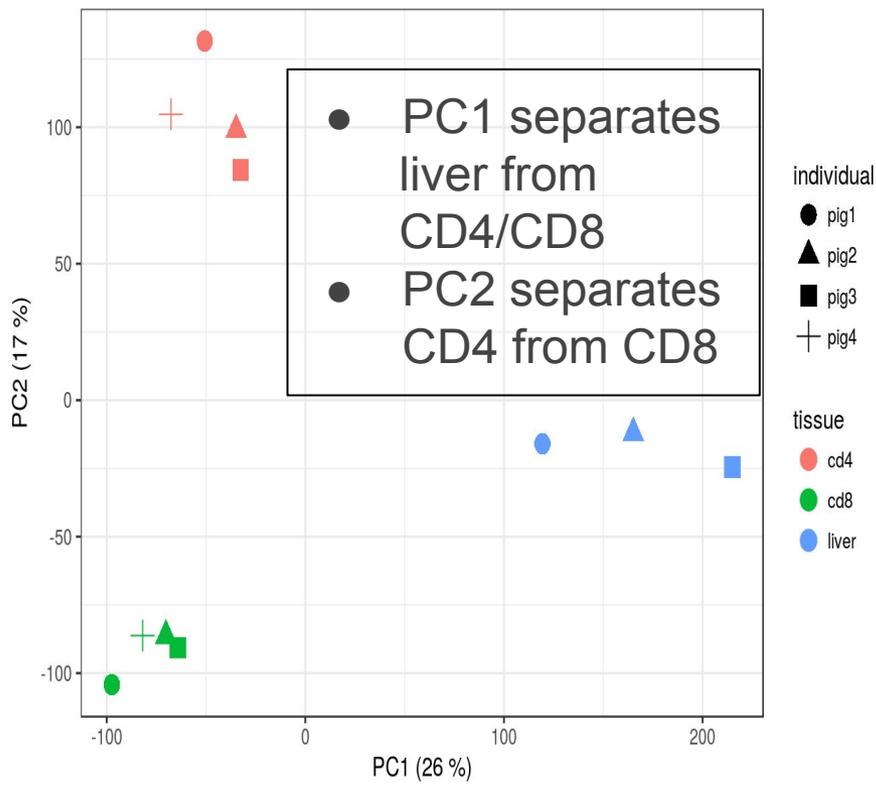
Open chromatin regions located close to TSSs are more ubiquitous than other open chromatin regions



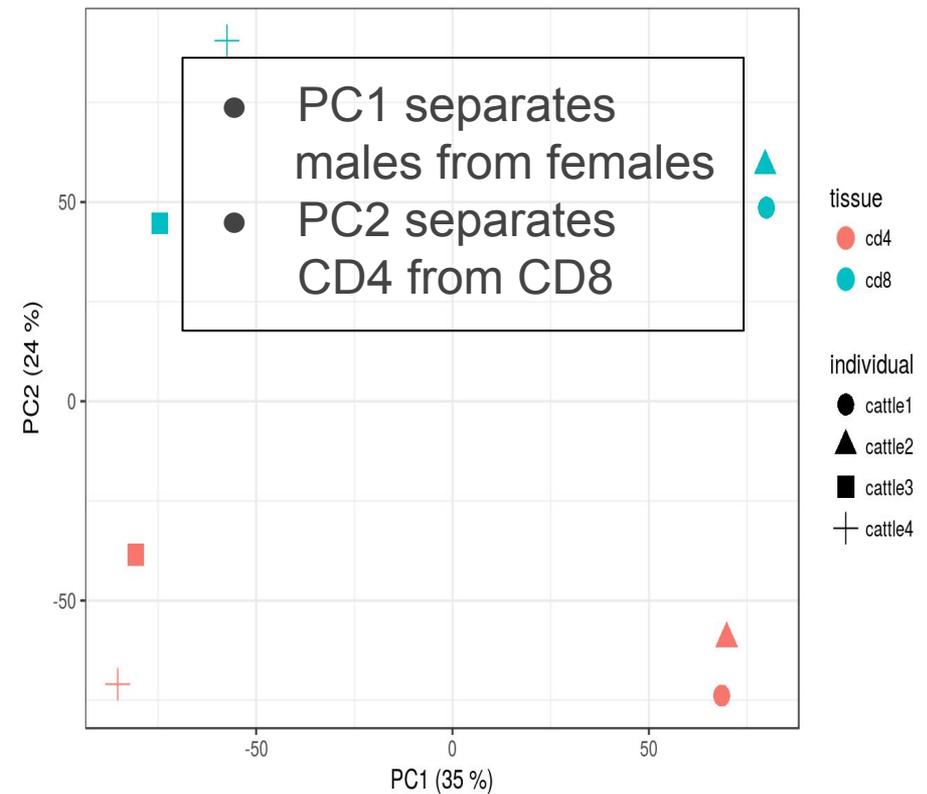
ATAC-seq coverage peaks at annotated TSS in all species<sup>10</sup>

# ATAC-seq PCA first separates liver from immune cells; in absence of liver, ATAC-seq PCA first separates males from females

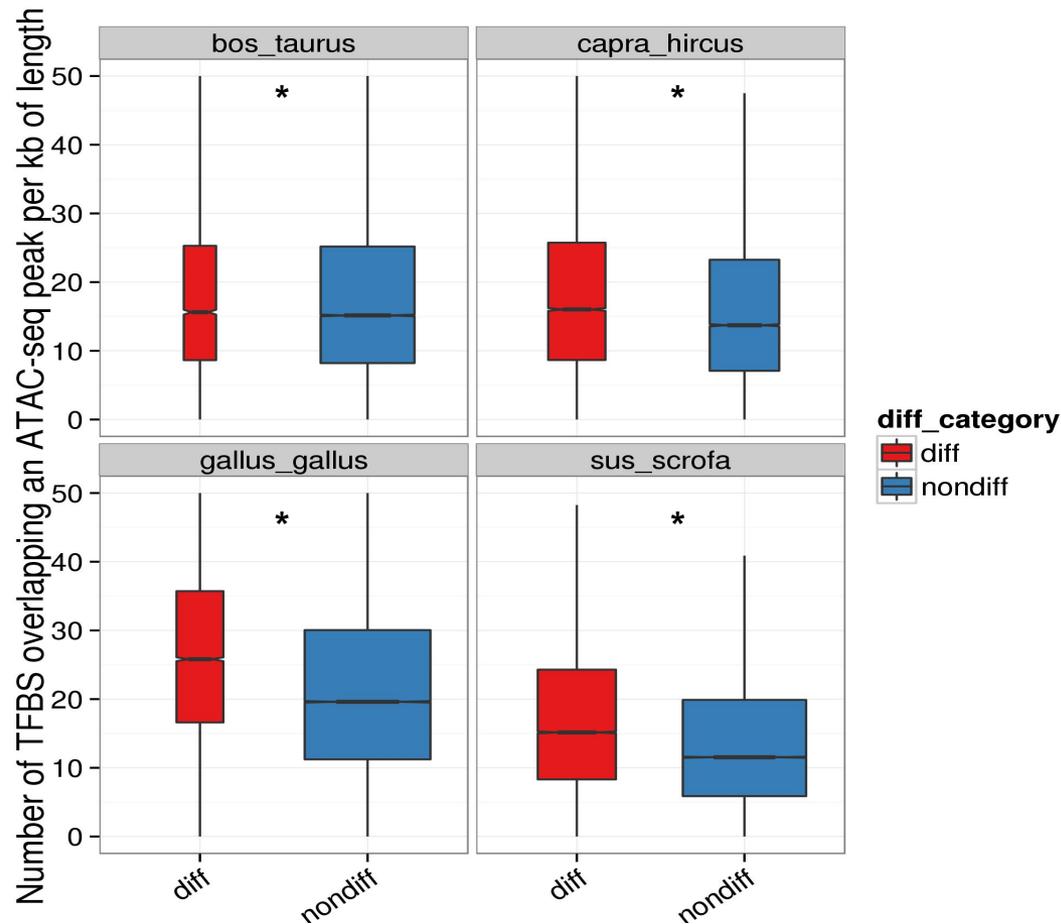
Sus scrofa



Bos taurus



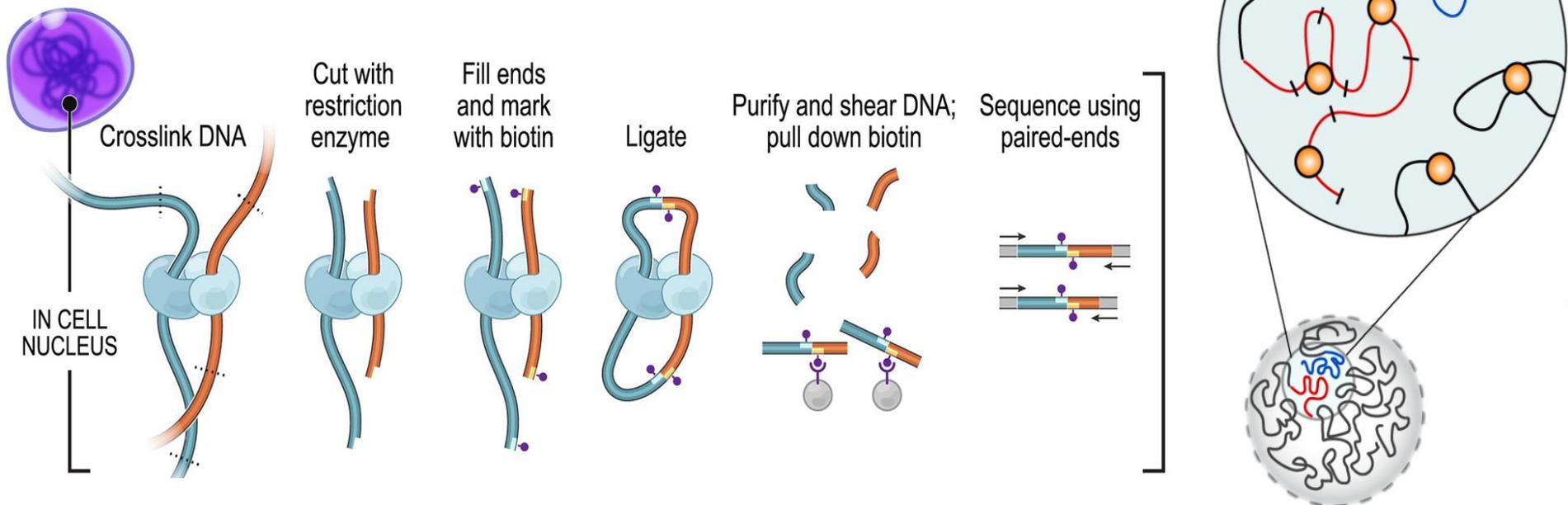
# Differential ATAC-seq peaks are more likely to be regulatory



Between tissue **differential** ATAC-seq peaks have a higher **TFBS density** than non differential ATAC-seq peaks (Wilcoxon test,  $p\text{-value} < 10^{-15}$ )

→ Differential ATAC-seq peaks are more likely to have a **regulatory** role

# HiC for 3D genomic structure



Rao et al, Cell, 2014

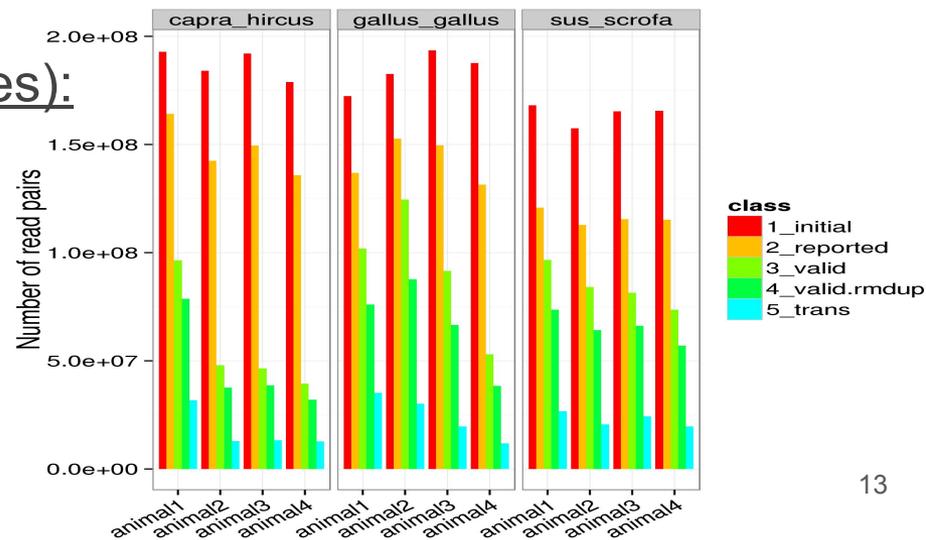
Belton et al, Methods, 2012

- HiC data (liver of 4 animals of 3 species):

- 180 million read pairs per sample

- HiC data analysis pipeline:

- HiC-Pro: Read Mapping/Filtering/Normalization
- Armatus: TAD calling
- HiTC: A/B compartment calling

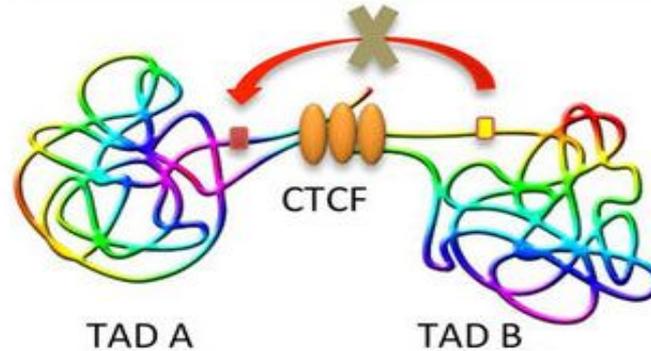


# Predicted CTCF binding sites peak at Topologically Associating Domain (TAD) boundaries

Hi-C contact matrix

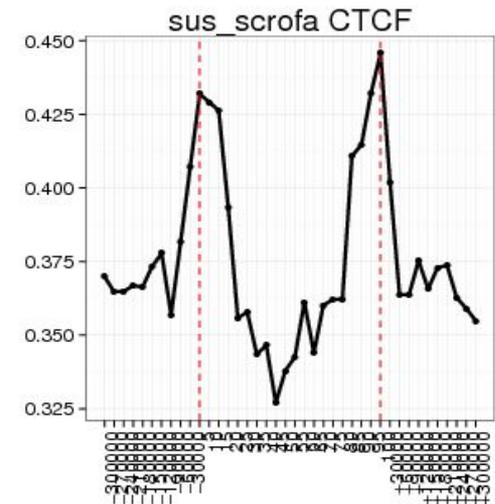
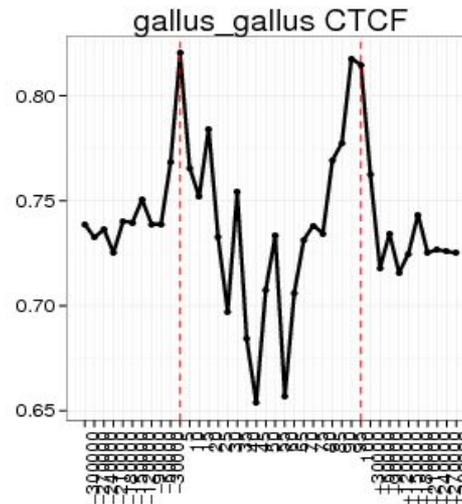
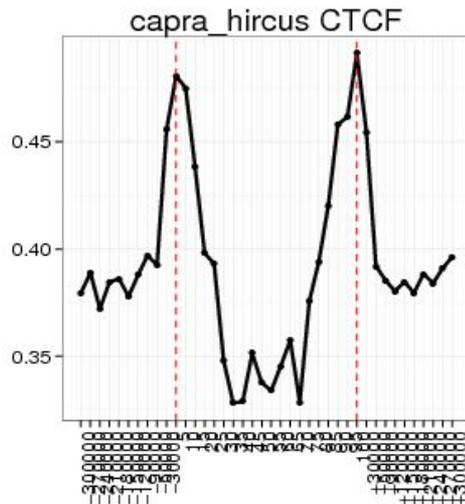


3D model

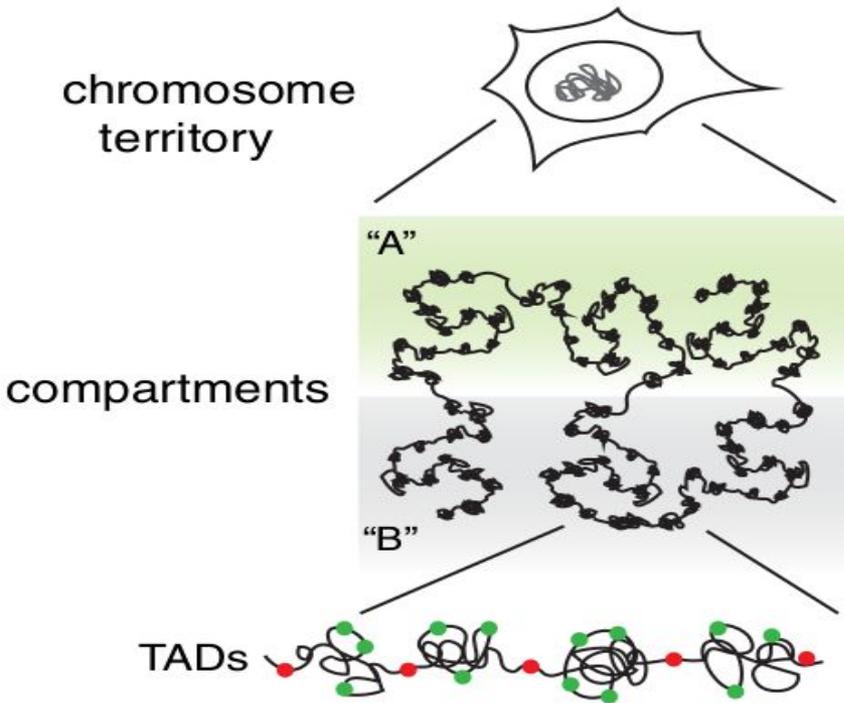


*Li et al, Scientific Reports, 2016*

AVG coverage

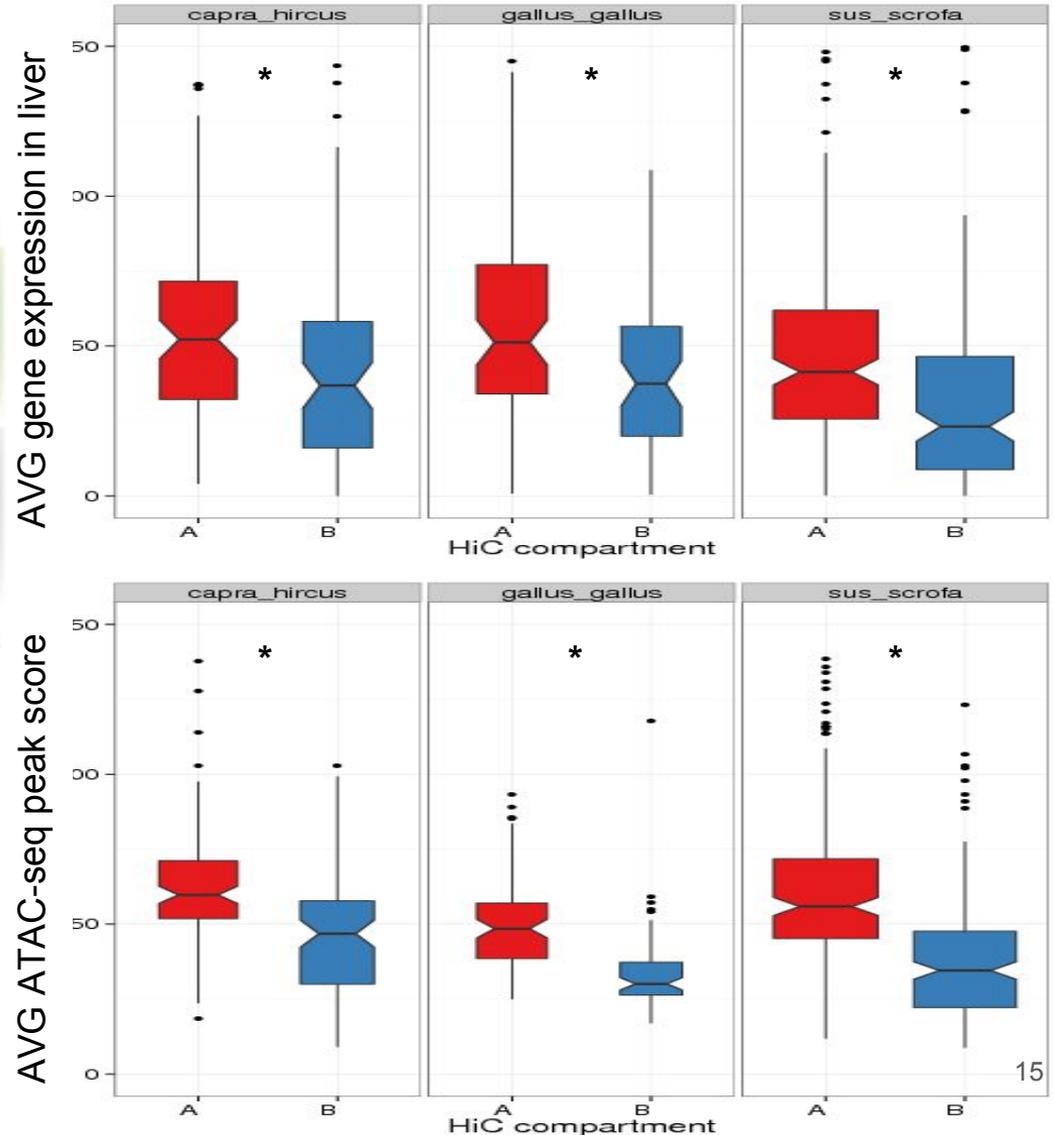


# There is global consistency between RNA-seq, ATAC-seq and HiC data



**A compartments: open, expressed**  
**B compartments: closed, repressed**

See Sylvain Foissac's talk at the *Pig Genetics and Genomics* session tomorrow Tuesday (talk # 144)



- **Management:**

- Elisabetta Giuffra
- Sylvain Foissac
- Sandrine Lagarrigue
- Marie-Hélène Pinard

- **Sampling:**

- Michèle Tixier-Boichard
- Stéphane Fabre
- Gwenola Tosser-Klopp
- Pascale Queré
- Fany Blanc
- Fabrice Laurent

- **Assays:**

- Hervé Acloque
- Diane Esquerre
- Sophie Pollet
- Adeline Goubil
- Florence Mompert
- Françoise Drouet
- Silvia Vincent-Nailleau

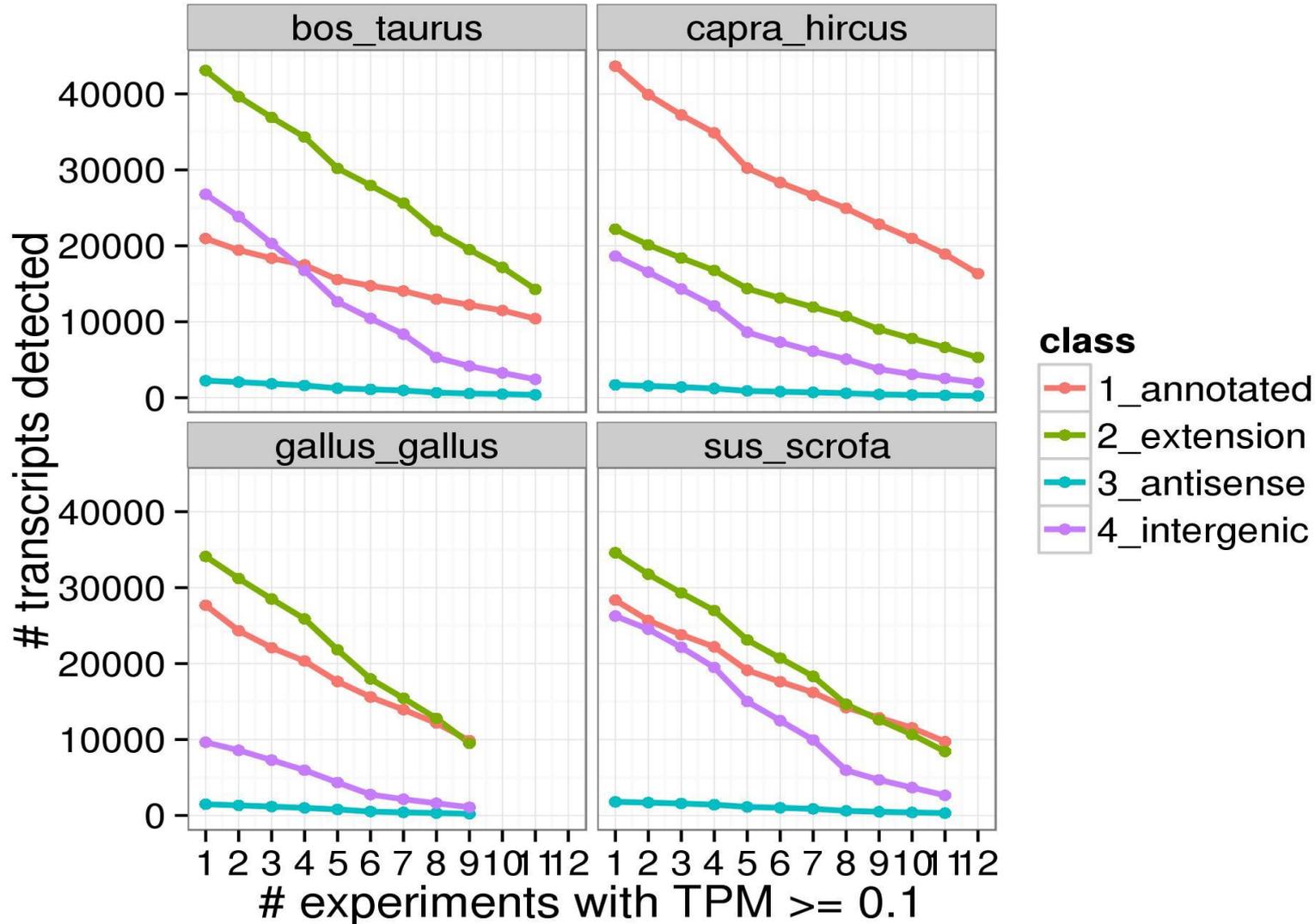
- **Data analysis:**

- Kylie Munyard
- Cédric Cabau
- Nathalie Villa-Vialaneix
- Matthias Zytnicki
- Kévin Muret
- Andrea Rau
- Thomas Derrien
- Christine Gaspin
- Christophe Klopp
- Ignacio Gonzalez
- David Robelin
- Magali San Cristobal
- Maria Marti
- Sylvain Marthey
- Philippe Bardou

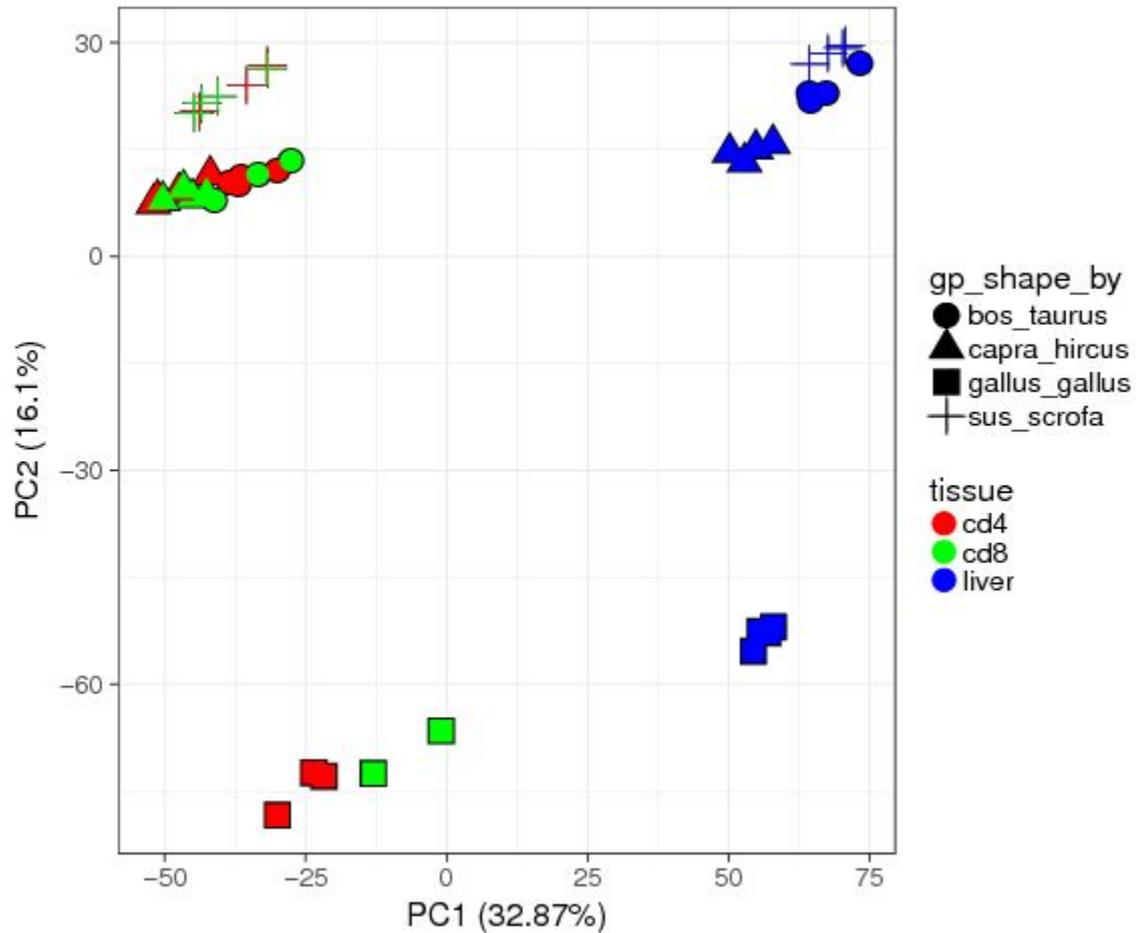
Thanks for your attention!

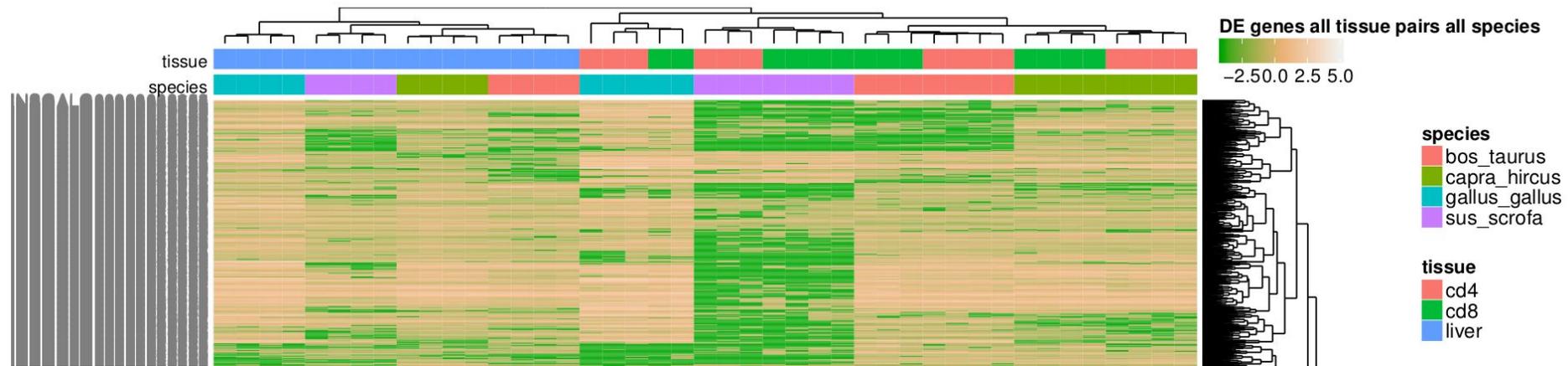
**Additional slides**

# Transcripts detected and their expression

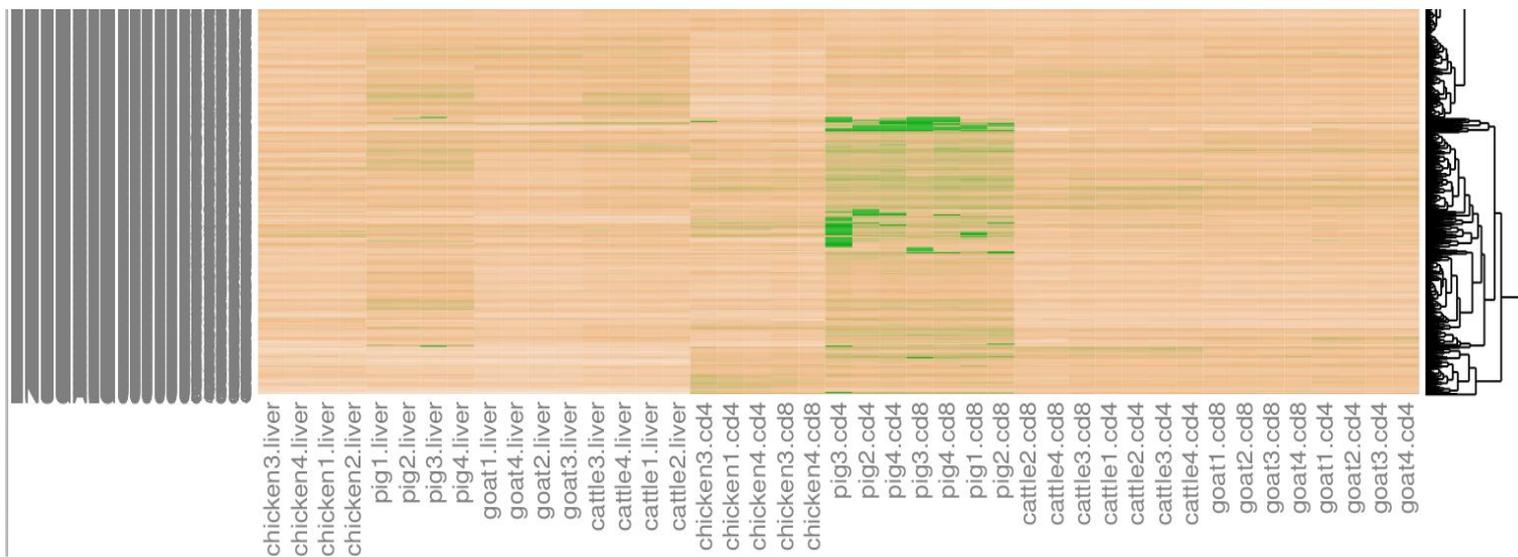


# All RNA-seq experiment PCA (~7300 orthologous genes)

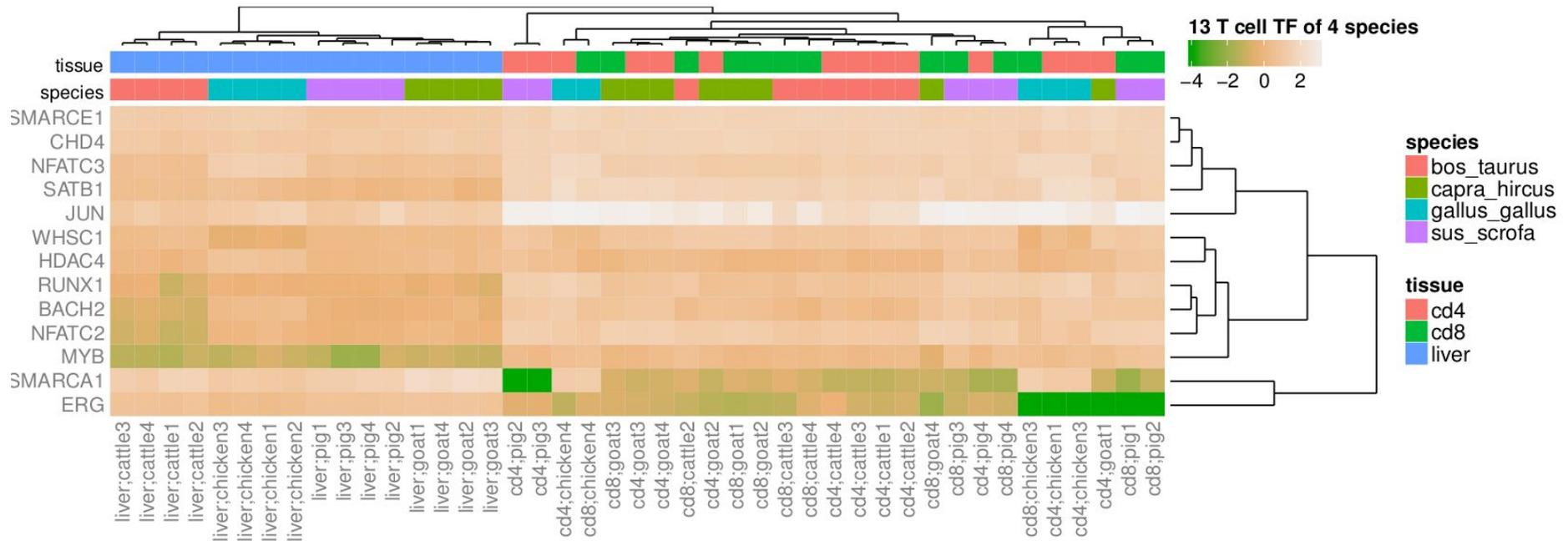




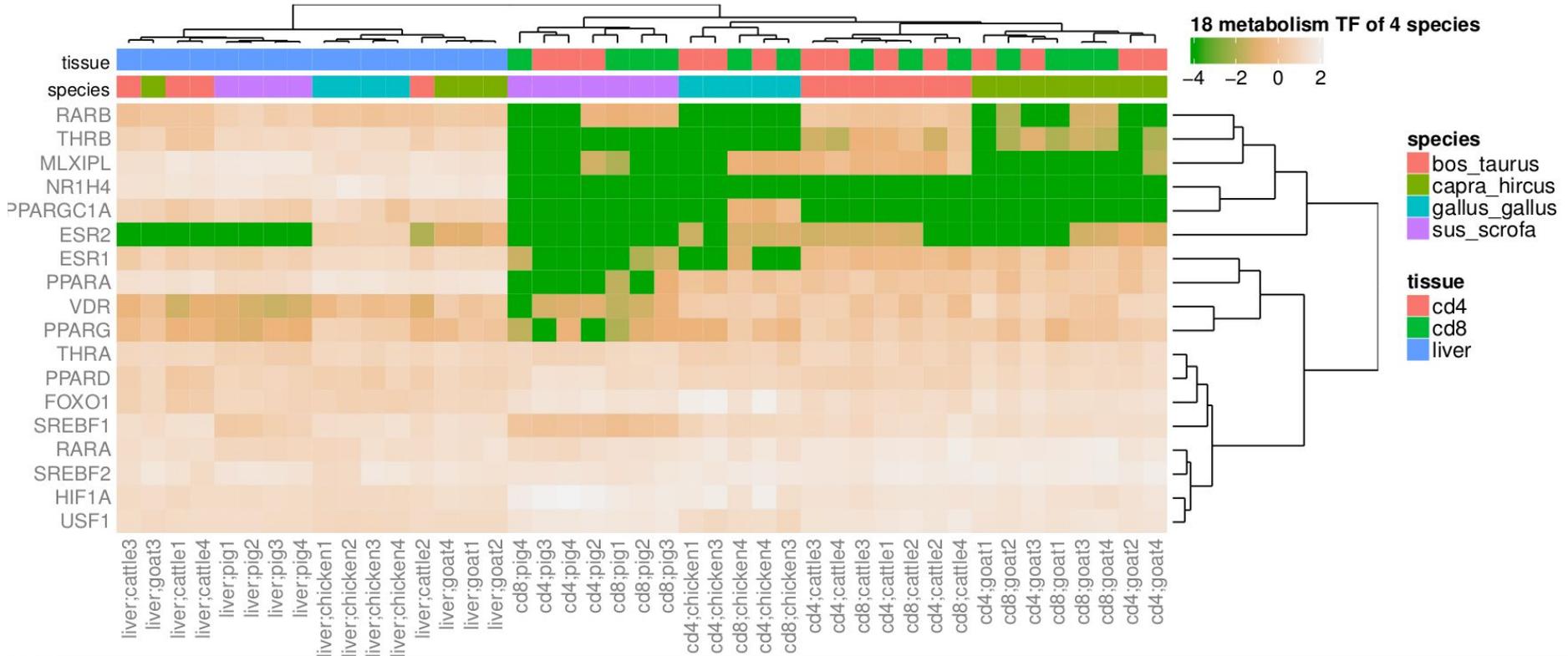
.....



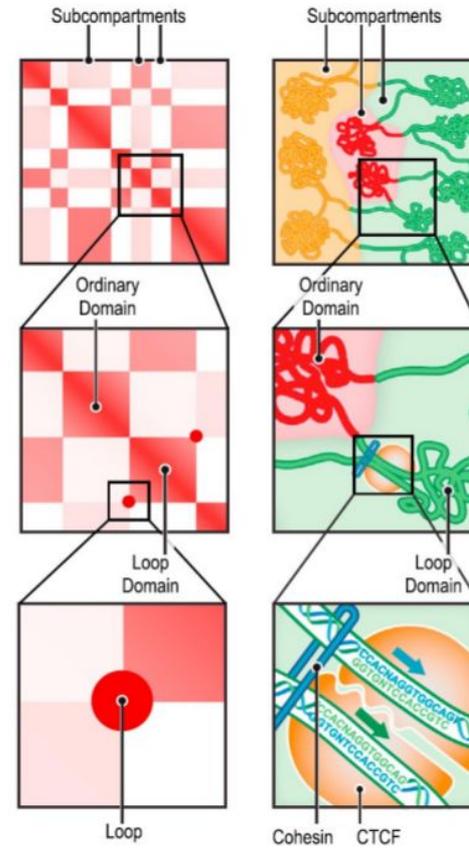
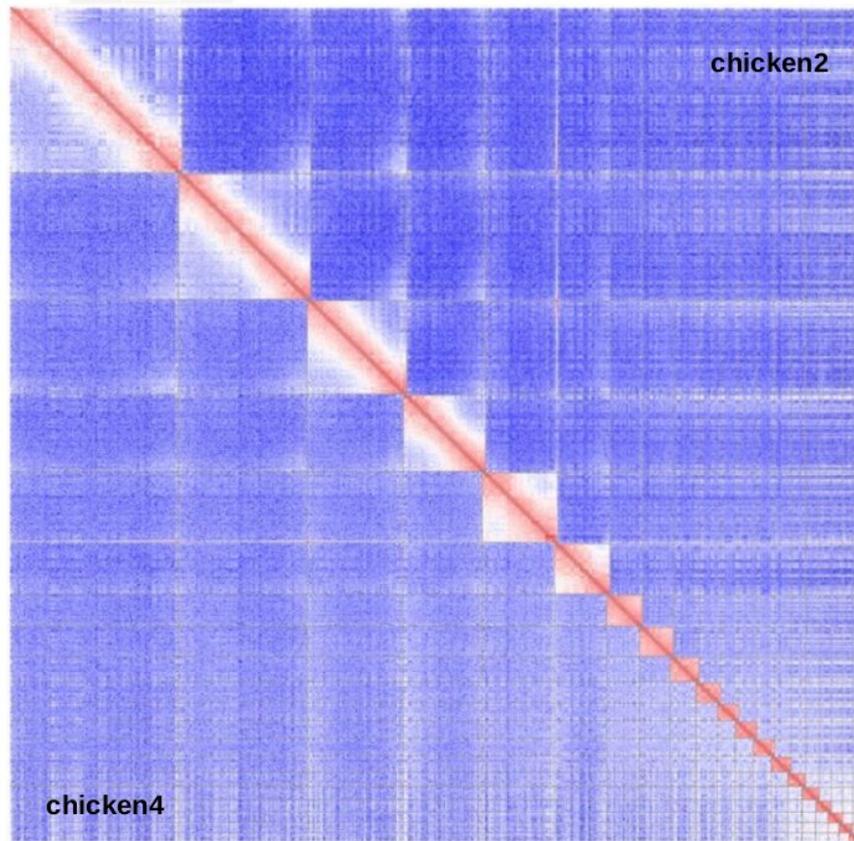
# Expression of the 13 T cell TF common to the 4 species (adding goat using gene name)



# Expression of the 18 metabolism TF common to the 4 species (adding goat using gene name)



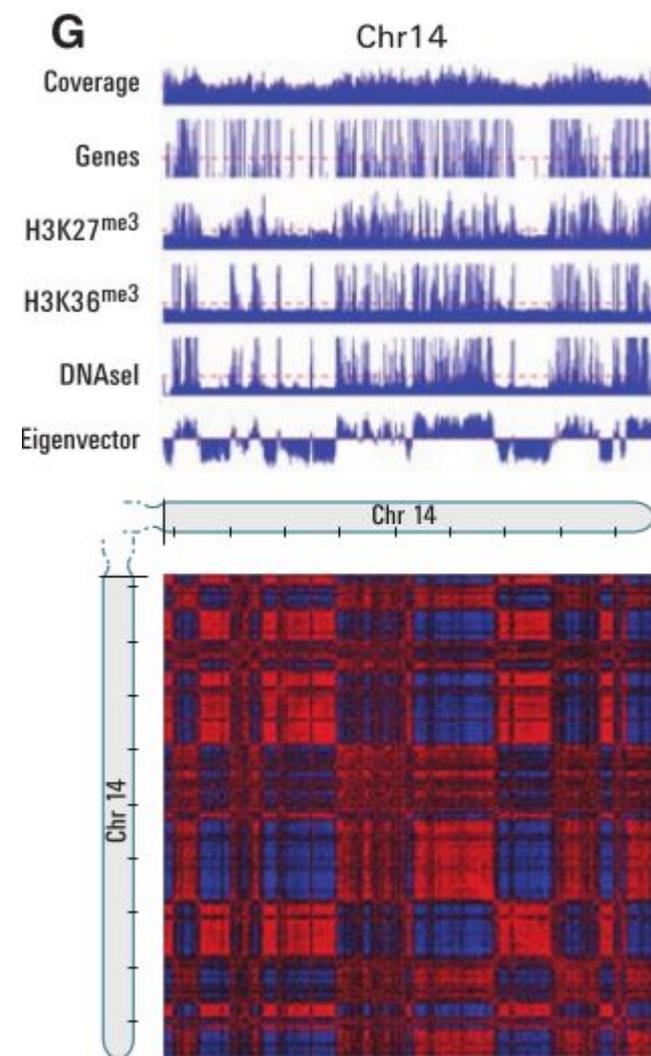
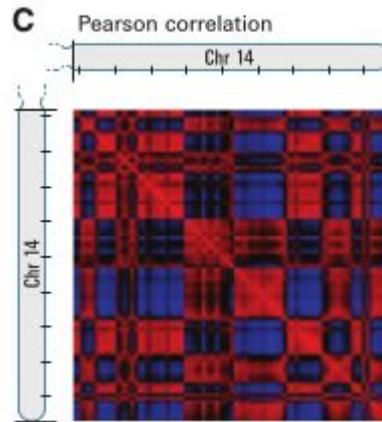
# HiC results: interaction matrices



Lieberman-aiden et al, 2009  
 Comprehensive Mapping of Long-Range Interactions Reveals  
 Folding Principles of the Human Genome

*The normalized matrix shows many large blocks of enriched and depleted interactions, generating a plaid pattern (Fig. 3B). If two loci (here 1-Mb regions) are nearby in space, we reasoned that they will share neighbors and have correlated interaction profiles. We therefore defined a correlation matrix  $C$  in which  $c_{ij}$  is the Pearson correlation between the  $i$ th row and  $j$ th column of  $M^*$ . This process dramatically sharpened the plaid pattern (Fig. 3C);*

*The plaid pattern suggests that each chromosome can be decomposed into two sets of loci (arbitrarily labeled A and B) such that contacts within each set are enriched and contacts between sets are depleted. We partitioned each chromosome in this way by using principal component analysis. For all but two chromosomes, the first principal component (PC) clearly corresponded to the plaid pattern (positive values defining one set, negative values the other) (fig. S1). For chromosomes 4 and 5, the first PC corresponded to the two chromosome arms, but the second PC corresponded to the plaid pattern. The entries of the PC vector reflected the sharp transitions from compartment to compartment observed within the plaid heatmaps.*



# Method

- Get contact matrix (raw counts)
- Extract intra-chromosomal sub-matrices (& regenerate corresponding indices)
- For each chrom independently:
  - ICE-normalize counts (matrix balancing)
  - Normalize by expected counts (scale by the distance factor)
  - Run a PCA of the bins using these counts => “**direct**” method
  - Generate pearson correlation matrix
  - Run a PCA on the bins using the correlations => “**corr**” method
  - Run a PCA on the bins using the correlations with HitC package (that does additional filtering) => “**hitc**” method
  - Extract from the PCA of each method the 3 first PCs
  - Choose the PC by comparing with the ICE-normalized counts on the diagonal (2 ways: PC vs. count correlation or t.test PC sign vs. counts)
- Segmentation by merging adjacent bins with same PC sign (+/-) and assign A/B compartment using the sign of the PC vs. count correlation.

# HiC read density peaks at distal ATAC-seq peaks but is depleted at TSS ATAC-seq peaks

