

RESEARCH ARTICLE

Open Access



# Multi-species annotation of transcriptome and chromatin structure in domesticated animals

Sylvain Foissac<sup>1\*†</sup> , Sarah Djebali<sup>1†</sup>, Kylie Munyard<sup>2</sup>, Nathalie Vialaneix<sup>3</sup>, Andrea Rau<sup>4</sup>, Kevin Muret<sup>5</sup>, Diane Esquerré<sup>1,6</sup>, Matthias Zytnicki<sup>3</sup>, Thomas Derrien<sup>7</sup>, Philippe Bardou<sup>1</sup>, Fany Blanc<sup>4</sup>, Cédric Cabau<sup>1</sup>, Elisa Crisci<sup>4,10</sup>, Sophie Dhorne-Pollet<sup>4</sup>, Françoise Drouet<sup>8</sup>, Thomas Faraut<sup>1</sup>, Ignacio Gonzalez<sup>3</sup>, Adeline Goubil<sup>4</sup>, Sonia Lacroix-Lamandé<sup>8</sup>, Fabrice Laurent<sup>8</sup>, Sylvain Marthey<sup>4</sup>, Maria Marti-Marimon<sup>1</sup>, Raphaëlle Momal-Leisenring<sup>4</sup>, Florence Mompert<sup>1</sup>, Pascale Quéré<sup>8</sup>, David Robelin<sup>1</sup>, Magali San Cristobal<sup>1</sup>, Gwenola Tosser-Klopp<sup>1</sup>, Silvia Vincent-Naulleau<sup>9</sup>, Stéphane Fabre<sup>1</sup>, Marie-Hélène Pinard-Van der Laan<sup>4</sup>, Christophe Klopp<sup>3</sup>, Michèle Tixier-Boichard<sup>4</sup>, Hervé Acloque<sup>1,4</sup>, Sandrine Lagarrigue<sup>5</sup> and Elisabetta Giuffra<sup>4\*</sup>

## Abstract

**Background:** Comparative genomics studies are central in identifying the coding and non-coding elements associated with complex traits, and the functional annotation of genomes is a critical step to decipher the genotype-to-phenotype relationships in livestock animals. As part of the Functional Annotation of Animal Genomes (FAANG) action, the FR-AgENCODE project aimed to create reference functional maps of domesticated animals by profiling the landscape of transcription (RNA-seq), chromatin accessibility (ATAC-seq) and conformation (Hi-C) in species representing ruminants (cattle, goat), monogastrics (pig) and birds (chicken), using three target samples related to metabolism (liver) and immunity (CD4+ and CD8+ T cells).

**Results:** RNA-seq assays considerably extended the available catalog of annotated transcripts and identified differentially expressed genes with unknown function, including new syntenic lncRNAs. ATAC-seq highlighted an enrichment for transcription factor binding sites in differentially accessible regions of the chromatin. Comparative analyses revealed a core set of conserved regulatory regions across species. Topologically associating domains (TADs) and epigenetic A/B compartments annotated from Hi-C data were consistent with RNA-seq and ATAC-seq data. Multi-species comparisons showed that conserved TAD boundaries had stronger insulation properties than species-specific ones and that the genomic distribution of orthologous genes in A/B compartments was significantly conserved across species.

**Conclusions:** We report the first multi-species and multi-assay genome annotation results obtained by a FAANG project. Beyond the generation of reference annotations and the confirmation of previous findings on model animals, the integrative analysis of data from multiple assays and species sheds a new light on the multi-scale selective

(Continued on next page)

\*Correspondence: [sylvain.foissac@inra.fr](mailto:sylvain.foissac@inra.fr); [elisabetta.giuffra@inra.fr](mailto:elisabetta.giuffra@inra.fr)

†Sylvain Foissac and Sarah Djebali contributed equally to this work.

<sup>1</sup>GenPhySE, Université de Toulouse, INRA, INPT, ENVT, Chemin de Borde Rouge, F-31326 Castanet-Tolosan Cedex, France

<sup>4</sup>GABI, AgroParisTech, INRA, Université Paris Saclay, F-78350 Jouy-en-Josas, France

Full list of author information is available at the end of the article



(Continued from previous page)

pressure shaping genome organization from birds to mammals. Overall, these results emphasize the value of FAANG for research on domesticated animals and reinforces the importance of future meta-analyses of the reference datasets being generated by this community on different species.

**Keywords:** Functional annotation, Livestock, RNA-seq, ATAC-seq, Hi-C

## Background

Most complex trait-associated loci lie outside protein-coding regions, and comparative genomics studies have shown that the majority of mammalian-conserved and recently adapted regions consist of non-coding elements [1–3]. This evidence prompted the first large-scale efforts into genome annotation for human and model organisms [4–6]. The genome-wide annotation maps generated by these projects helped to shed light on the main features of genome activity. For example, chromatin conformation or transcription factor occupancy at regulatory elements can often be directly tied to the biology of the specific cell or tissue under study [3, 7, 8]. Moreover, although a subset of core regulatory systems are largely conserved across humans and mice, the underlying regulatory systems often diverge substantially [9–11], implying that understanding the phenotypes of interest requires organism-specific information for any specific physiological phase, tissue and cell.

The Functional Annotation of Animal Genomes (FAANG) initiative [12] aims to support and coordinate the community in the endeavor of creating reference functional maps of the genomes of domesticated animals across different species, tissues, and developmental stages, with an initial focus on farm and companion animals [13–16]. FAANG carries out activities to standardize assay protocols and analysis pipelines, to coordinate and facilitate data sharing. The FAANG Data Coordination Center provides an infrastructure for genotype-to-phenotype data [17, 18]. Substantial efforts are being dedicated to farm animal species, as deciphering the genotype-to-phenotype relationships underlying complex traits such as production efficiency and disease resistance is a prerequisite for exploiting the full potential of livestock [13, 16].

Here we report the main results of a pilot project (FR-AgENCODE [19]) launched at the beginning of the FAANG initiative. The broad aim was to generate standardized FAANG reference datasets from four livestock species (cattle, goat, chicken, and pig) through the adaptation and optimization of molecular assays and analysis pipelines. We first collected a panel of samples from more than 40 tissues from two males and two females of four species: *Bos taurus* (cattle, Holstein breed), *Capra hircus* (goat, Alpine breed), *Gallus gallus* (chicken, White Leghorn breed), and *Sus scrofa* (pig, Large White

breed), generating a total of 4115 corresponding entries registered at the EMBL-EBI BioSamples database (see “Methods” section). For molecular characterization, three tissues were chosen to represent a “hub” organ (liver) and two broad immune cell populations (CD4+ and CD8+ T cells). This allowed the acquisition of a partial representation of energy metabolism and immunity functions, as well as the optimization of the protocols for experimental assays for both tissue-dissociated and primary sorted cells. In addition to the transcriptome, we analyzed chromatin accessibility by the assay for transposase-accessible chromatin using sequencing (ATAC-seq, [20]), and we characterized the three-dimensional (3D) genome architecture by coupling proximity-based ligation with massively parallel sequencing (Hi-C, [21]) (Fig. 1). Using this combination of tissues/assays, we assessed the expression of a large set of coding and non-coding transcripts in the four species, evaluated their patterns of differential expression in light of chromatin accessibility at promoter regions, and characterized active and inactive topological domains of these genomes. The integrative analysis showed a global consistency across all data, emphasizing the value of a coordinated action to improve the genomic annotation of livestock species, and revealed multiple layers of evolutionary conservation from birds to mammals.

## Results and discussion

### High-depth RNA-seq assays provide gene expression profiles in liver and immune cells from cattle, goat, chicken, and pig

For each animal (two males, two females) of the four species, we used RNA-seq to profile the transcriptome of liver, CD4+ and CD8+ T cells (see “Methods” section, Fig. 1 and Additional file 1: Table S1). We prepared stranded libraries from polyA+ selected RNAs longer than 200 bp, and we sequenced them on an Illumina HiSeq3000 (see “Methods” section). Between 250M (chicken) and 515M (goat) read pairs were obtained per tissue on average, of which 94% (chicken) to 98% (pig) mapped to their respective genome using STAR [22, 23] (see “Methods” section, Additional file 1: Figure S1 and Tables S2–S4). As an initial quality control step, we processed the mapped reads with RSEM [24] to estimate the expression levels of all genes and transcripts from the Ensembl reference annotation (hereafter called “reference” genes/transcripts/annotation) (Additional file 1:

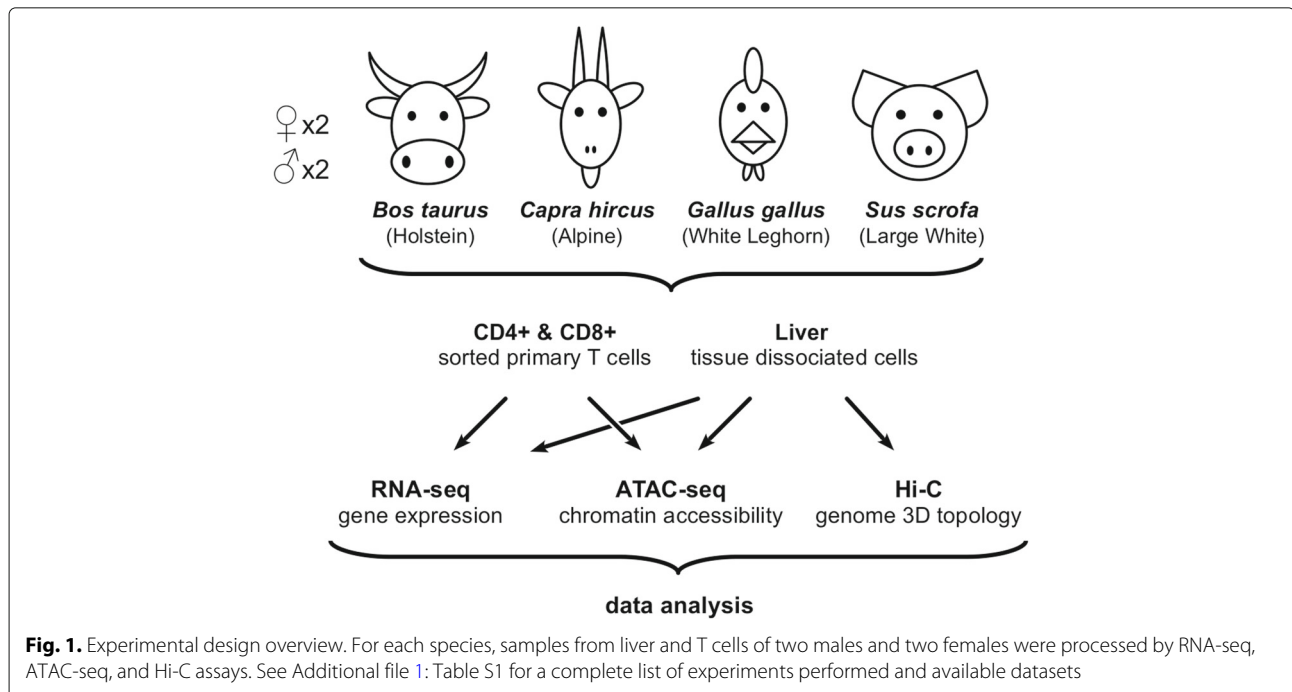


Table S2 and Figure S1; Additional file 2). As expected, a large majority of the reads (from 62% in cattle to 72% in goat) fell in annotated exons of the reference genes (Additional file 1: Figure S2). In spite of the specialized scope of this initial study limited to liver and immune cells, a large share of all reference genes were detected (from 58% in chicken to 65% in goat), even considering only transcript isoforms with an expression level higher than 0.1 TPM in at least two samples (see “Methods” section and Table 1).

For each species, we explored the similarity among samples using the expression profiles of the reference genes. Principal component analysis (PCA) revealed quite consistent patterns across species, where the first principal component (explaining 84 to 91% of the variability among samples) clearly separated samples according to their tissue of origin (liver vs. T cells). A more moderate yet systematic separation was observed between CD4+ and CD8+ T cells on the second principal component (Additional file 1: Figure S3). The consistency of these patterns across species supports the reliability of our RNA-seq data.

To compare the expression pattern of reference genes between species, we first checked that the male-to-female expression ratio was globally uniform genome-wide with the exception of the sex chromosomes in chicken (Additional file 1: Figure S4). Dosage compensation (leading for instance to X chromosome inactivation in mammals) can indeed be observed in all species except chicken, which is consistent with previous reports on dosage compensation [25]. Next, we hierarchically clustered our samples

using the expression of 9461 genes found to be orthologous across the four species (Fig. 2, “Methods” section and Additional file 3). Regardless of the species, liver and T cell samples clustered separately. Interestingly, T cell samples clustered first by species and then by subtypes (i.e., CD4+ versus CD8+). This suggests a strong specialization of the immune system during speciation, although the specific clustering pattern of CD4+ and CD8+ samples might also be driven by a small subset of genes whose expression varies largely across species and little across cell subtypes [26]. These results also depend on the set of orthologous genes available in the reference annotation. For most species and tissues, samples also clustered by sex, possibly due in part to the physiological state of females in lactation or laying eggs. These results highlight a high conservation of the liver gene expression program across vertebrates and that global transcriptome comparisons across several tissues and species can result in samples clustering by either factor, as shown in other studies [26–28].

#### Most reference genes are differentially expressed between liver and T cells

To provide functional evidence supporting our RNA-seq data, we performed a differential gene expression analysis across tissues per species for each gene in the reference annotation. Gene read counts provided by RSEM [24] were TMM-normalized [29] (see “Methods” section). Taking into account the specificities of our experimental design, in which samples from different tissues come from the same animal, we fitted generalized linear

**Table 1** Reference and FR-AgENCODE detected transcripts. This table provides the total number of reference transcripts for each species, the number and percent of those that were detected by RNA-seq (TPM  $\geq$  0.1 in at least 2 samples), the total number of FR-AgENCODE transcripts, and the subsets of them that were mRNAs (known and novel) and lncRNAs (known and novel). Overall, the transcript repertoire is increased by about 100% in most of the species. As these results naturally depend on the input data, details about the genome assemblies and reference annotations that were used for this study are listed in Additional file 1: Table S2

Species	Reference transcripts			FR-AgENCODE transcripts				
	All	Expressed		#	mRNAs		lncRNAs	
		#	% of total		Known	Novel	Known	Novel
Cattle	26,740	16,100	60.2	84,971	11,576	48,225	13	22,711
Goat	53,266	34,442	64.7	78,091	26,973	31,854	2247	11,617
Chicken	38,118	22,180	58.2	57,817	14,765	32,802	1314	6797
Pig	49,448	29,786	60.2	77,540	23,701	40,020	327	12,284

models (GLM) to identify genes with differential expression in either all pairwise comparisons between liver, CD4+ and CD8+ (model 1), or liver versus T cells globally (model 2).

As expected, the liver to T cell comparison yielded the largest number of differentially expressed genes, and relatively few genes were found to be differentially expressed between the two T cell populations (CD4+ and CD8+, see Additional file 1: Table S5 and Additional file 4). Strikingly, most genes showed significantly different expression between liver and immune cells in each species (from 7000 genes in chicken to 10,500 genes in goat), reflecting the difference between the physiological functions of these highly specialized cell types, in line with findings from the GTEx project [30]. Accordingly, Gene Ontology (GO) analysis provided results in line with the role of liver in metabolism and of T cells in immunity for all species (Additional file 1: Figure S5–8 and “Methods” section).

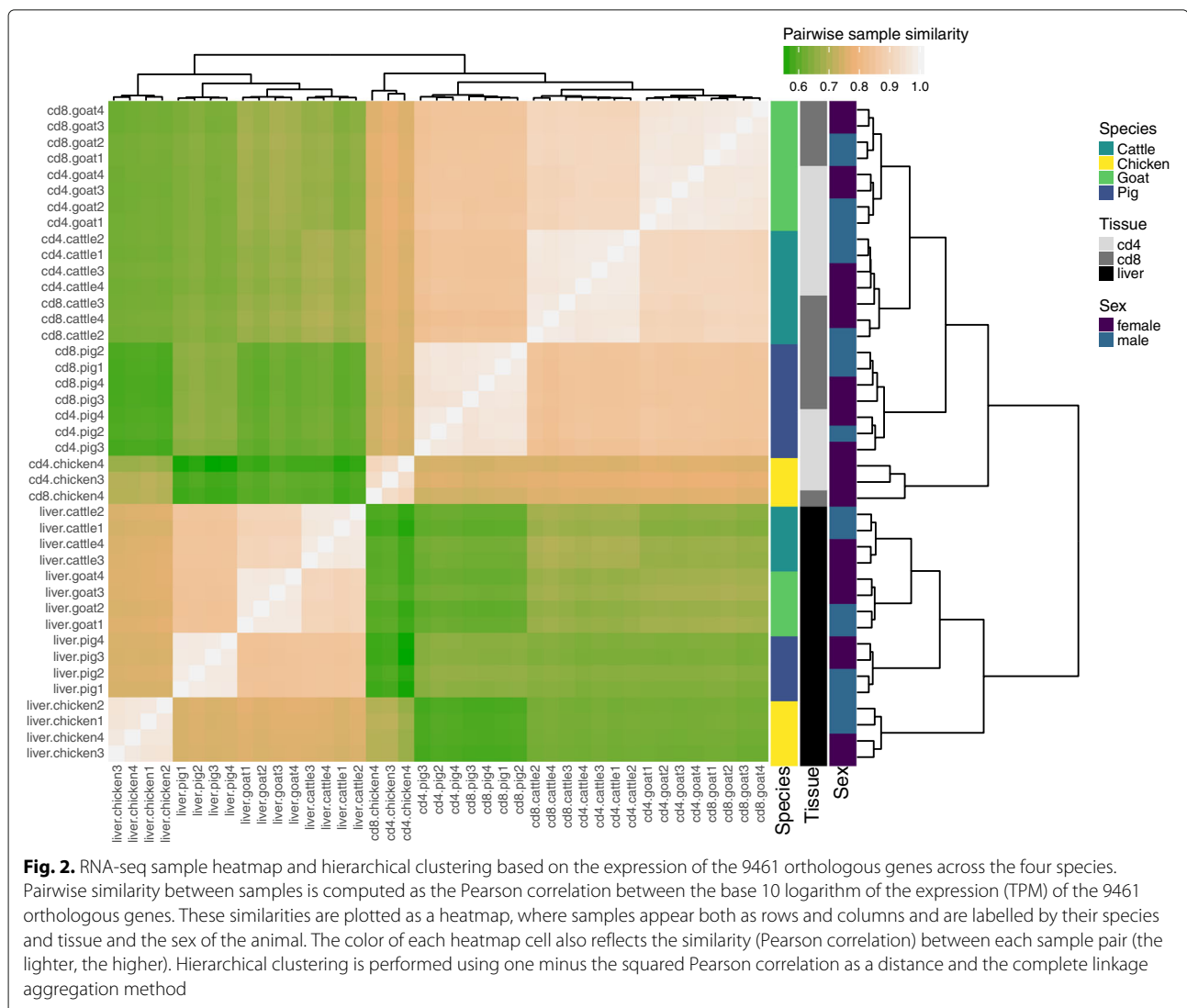
In accordance with the results of the hierarchical clustering (Fig. 2), most orthologous genes found to be differentially expressed between CD4+ and CD8+ T cells within species showed high variability of expression levels across species (not shown). This variability is likely caused by the natural heterogeneity in the relative proportions of T cell subtypes among the different species, as already reported between mammals [31, 32]. Nevertheless, 39 orthologous genes could consistently differentiate CD4+ and CD8+ in the four species, including mammals and chicken, which is significantly more than expected by chance ( $p$  value  $< 10^{-3}$ , permutation test). Among those, 10 and 29 genes showed significant overexpression in CD4+ and CD8+ cells respectively (Additional file 1: Table S6). We searched for the human orthologs of these genes in the baseline expression dataset of human CD4+ and CD8+  $\alpha\beta$  T cell subsets generated by the Blueprint Epigenome Project [33, 34] and considered their relative enrichment in each cell subset. With one exception (*ACVRL1*), all genes were

found to be expressed in human CD4+ and/or CD8+  $\alpha\beta$  T cells and 25 of them showed a relative enrichment in CD4+ (or CD8+) human cells consistent with our data across the four species. Out of these 25 genes, six and eight genes, respectively, could be associated with CD4+ and CD8+ T cell differentiation, activation, and function according to the literature (Additional file 1: Table S6).

#### Analysis of new transcripts improves and extends gene structure annotation

In order to test if our data could improve the reference gene annotation for each species, we used STAR and Cufflinks to identify all transcripts present in our samples and predict their exon-intron structures. We then quantified their expression in each sample using STAR and RSEM (see “Methods” section and Additional file 1: Figure S1) and only retained the transcripts and corresponding genes expressed in at least two samples with TPM  $\geq$  0.1. We identified between 58,000 and 85,000 transcripts depending on the species (Table 1, Additional file 5), hereafter called “FR-AgENCODE transcripts”.

To characterize these FR-AgENCODE transcripts with respect to the reference transcripts, we grouped them into four positional classes (see “Methods” section): (1) *known*: a FR-AgENCODE transcript whose exon-intron structure is strictly identical to a reference transcript (i.e., exact same exons and introns); (2) *extension*: same as (1), but the FR-AgENCODE transcript extends a reference transcript by at least 1 bp on at least one side; (3) *alternative*: a FR-AgENCODE transcript that shares at least one intron with a reference transcript but does not belong to the previous categories (only multi-exonic transcripts can be in this class); and (4) *novel*: a FR-AgENCODE transcript that does not belong to any of the above categories. We found that most FR-AgENCODE transcripts (between 37% for goat and 49% for chicken)



were of the alternative class, therefore enriching the reference annotation with new splice variants of known genes (Additional file 1: Table S7). The proportion of completely novel transcripts was relatively high for cattle, which is likely due to the incompleteness of the UMD3.1 version of the Ensembl annotation used at the time of the study (Table 1, Additional file 1: Figure S2, S9 and Table S7).

In order to identify interesting new transcripts involved in immunity and metabolism, we first selected the novel FR-AgENCODING coding transcripts that unambiguously project to a single human coding gene. We identified 93 (cattle), 52 (goat), 74 (chicken), and 26 (pig) genes, of which 12 are common to at least 2 livestock species (see “Methods” section, Additional file 1: Table S8, Figure S10A and Additional file 6). Gene set enrichment analyses on these gene lists confirmed their relevance for

T cell biology (Additional file 1: Figure S10B) and the added value of the FR-AgENCODING novel transcripts in terms of annotation of complex but important loci like TRBV and TRAV (Additional file 1: Figure S10C).

In addition, we performed a differential gene expression analysis similar to the one done on reference genes (see above and “Methods” section). Results were globally similar, with more than 88% of correspondence between the differentially expressed genes from the reference and the FR-AgENCODING annotation (Additional file 1: Figure S11, Tables S5 and S9; Additional file 7). Among the latter, between 202 (chicken) and 1032 (goat) genes were coding (at least one coding transcript predicted by FEELnc—see below) and did not overlap any reference gene on the same strand. This highlights the potential to identify novel interesting candidates for further functional characterization.



### Identification, classification, and comparative analysis of lncRNAs

Since deep RNA-seq libraries allow the detection of weakly expressed transcripts [35], we sought to identify the proportion of long non-coding RNAs (lncRNAs) among the FR-AgENCODING transcripts. Using the FEELnc classifier [36] trained on the reference annotation (see “Methods” section), we identified from 7502 (chicken) to 22,724 (cattle) lncRNA transcripts, among which a large majority were not previously reported (Additional file 1: Tables S10–11; Additional file 8). The high number of lncRNAs found in cattle is likely due in part to the incomplete genome annotation and genome assembly used at the time of the study (Additional file 1: Table S2). Consistent with previous reports in several species including human [37], dog [36], and chicken [38], predicted lncRNA genes had lower expression levels than reference protein-coding genes (Additional file 1: Figure S12). The structural features of these predicted lncRNA transcripts were consistent between the four species: lncRNAs are spliced but with fewer exons (1.5 vs. 10) and higher median exon length (660 vs. 130 bp) compared to mRNAs (Additional file 1: Figure S12). lncRNAs are also smaller than mRNAs (1800 vs. 3600 bp). Notably, the lower number of exons and consequent smaller size of lncRNAs compared to mRNAs could also be due to the weaker expression of lncRNAs, which makes their structure more difficult to identify [39].

In addition to the coding/non-coding classification, FEELnc can also categorize lncRNAs as intergenic or intragenic based on their genomic positions with respect to a provided set of reference genes (usually protein coding), and considering their transcription orientation with respect to these reference genes. This analysis revealed an overwhelming majority of intergenic lncRNA genes over intragenic ones (Additional file 1: Table S10), which is consistent with results obtained in human [37] and in chicken [38].

We and others previously showed a sharp decrease in lncRNA sequence conservation with increasing phylogenetic distance [37, 38, 40], in particular between chicken and human that diverged 300M years ago. We therefore analyzed lncRNA conservation between the four livestock species using a synteny approach based on the orthology of protein-coding genes surrounding the lncRNA and not on the lncRNA sequence conservation itself [38] (see “Methods” section). We found 73 such conserved, or syntenic, lncRNAs across cattle, goat, and pig, 19 across cattle, chicken, and pig, and 6 across all four species (Additional file 8). All were expressed in these species and located in the same orientation with respect to the flanking orthologous genes. An example of such a conserved lncRNA, hitherto unknown in our four species, is provided in Fig. 3. In human, this lncRNA is called

*CREMos* for “*CREM* opposite sense” since it is in a divergent position with respect to the neighboring *CREM* protein-coding gene. Interestingly, synteny is conserved across species from fishes to mammals and the *CREMos* lncRNA is overexpressed in T cells while the *CREM* protein-coding gene is overexpressed in liver in goat, cattle and chicken (Fig. 3). Additional examples of syntenic lncRNAs are provided in Additional file 8, and the ones found to be conserved between the 4 species are represented in Additional file 1: Figure S13.

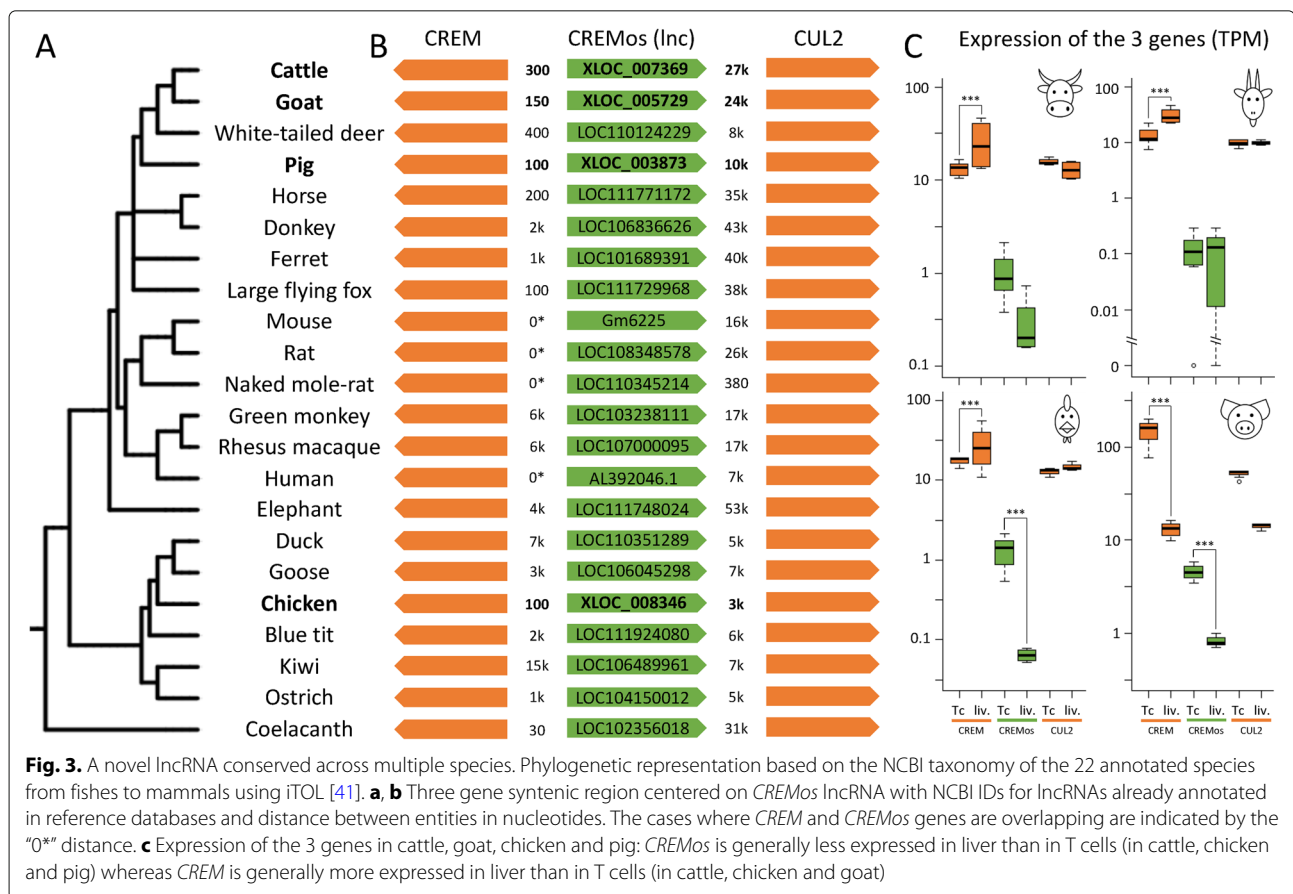
### Landscape of chromatin accessibility in cattle, goat, chicken, and pig

We used ATAC-seq to profile the accessible chromatin of liver, CD4+ and CD8+ T cells in animals from the four species. Between 480M (chicken) and 950M (pig) ATAC-seq fragments were sequenced per species, and were processed by a standard pipeline (“Methods” section and Additional file 1: Figure S14). Peaks were called in each tissue separately (see “Methods” section), resulting in between 26,000 (pig, liver) and 111,000 (pig, cd8) peaks per tissue (Additional file 1: Table S12). Those peaks were further merged into a unique set of peaks per species, resulting in between 75,000 (goat) and 149,000 (pig) peaks (Additional file 1: Table S12; Additional file 9), covering 1 to 5% of the genome. The average peak size was around 600 bp for all species, except for chicken where it was less than 500 bp. Merging tissue peaks did not result in much wider peaks (Additional file 1: Figure S15).

In comparison to the reference annotation, about 10–15% of the peaks lie at 5 kb or less from a Transcription Start Site (TSS) and can be considered to be promoter peaks. The precise distribution of these promoter peaks showed a clear higher accumulation at the TSS for all species (Fig. 4), supporting the quality of both the annotation and our data. Importantly, this signal was also observed around the TSS of novel FR-AgENCODING transcripts (i.e., those not from the known class; Additional file 1: Figure S16).

The vast majority of the peaks, however, were either intronic or intergenic (Fig. 4, Additional file 1: Figure S17; Additional file 9), similar to GWAS variants associated with human diseases [3]. In particular, from 38% (goat) to 55% (cattle) of the peaks were located at least 5 kb away from any reference gene (Additional file 1: Figure S17), indicating that ATAC-seq can detect both proximal and distal regulatory regions.

Since active enhancers are expected to be enriched in chromatin regions that are both accessible and tagged with specific histone modification marks, we compared our ATAC-seq peaks to histone ChIP-seq peaks from another functional genomics study [42]. In that study, two histone modification marks (H3K4me3 and H3K27ac) were profiled in the genome of 20 mammals including

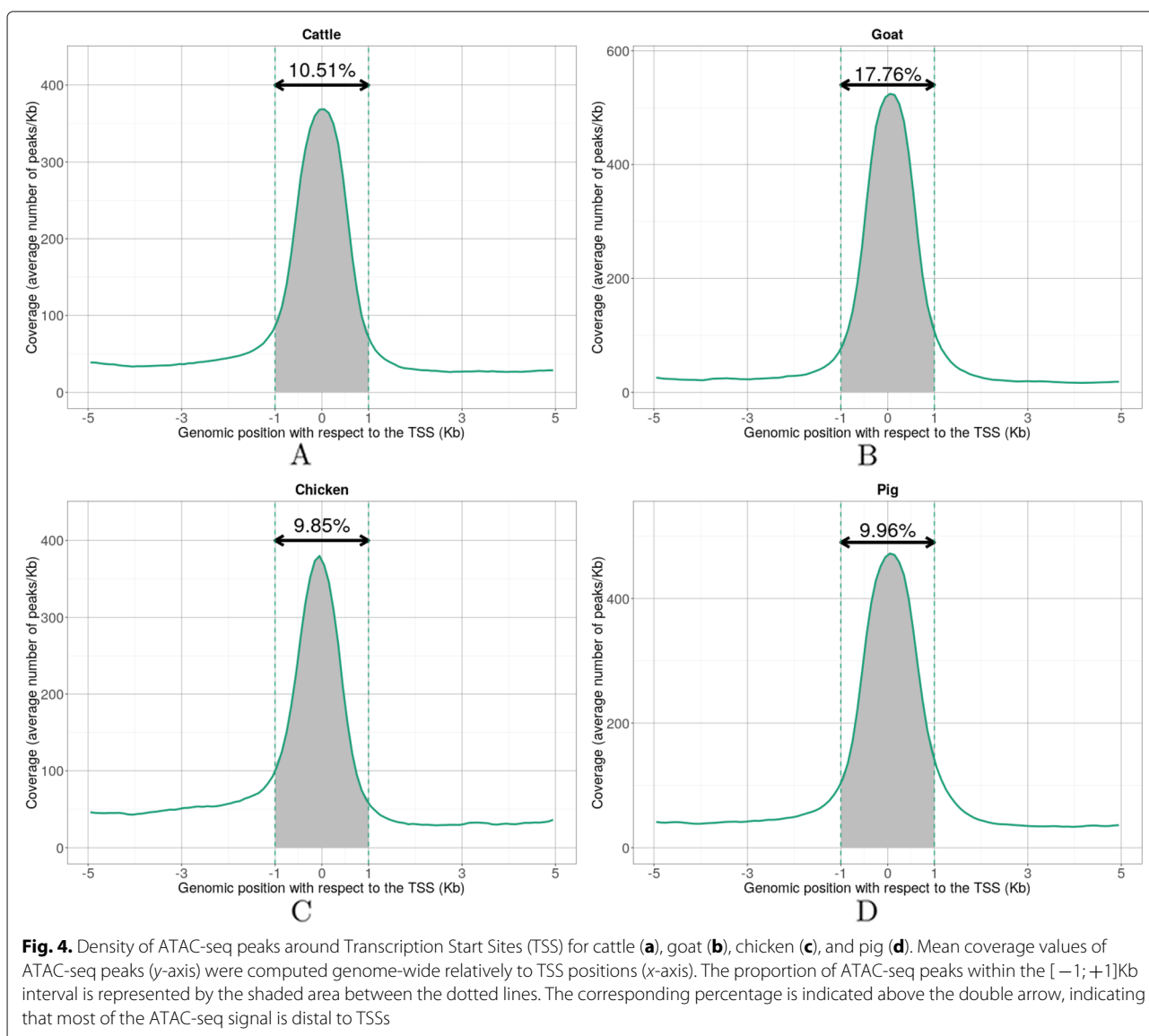


pig, for which we have ATAC-seq data in the same tissue (liver). This comparison showed that 6773 out of the 9632 H3K4me3 peaks (70.3%) and 8821 out of the 33,930 H3K27ac peaks (28.3%) overlapped our liver ATAC-seq peaks. These numbers were significantly higher than expected by chance as measured by shuffling peak positions ( $p$  value  $< 10^{-3}$ , permutation tests). Moreover, this subset of overlapping peaks have significantly higher  $q$ -value signal scores than their non-overlapping counterparts ( $p$  value  $< 2.2 \times 10^{-16}$ , Wilcoxon tests), which confirms the existence of a common signal between the datasets.

To further characterize functional regulatory sites in our samples, we compared chromatin accessibility between liver and T cells. The ATAC-seq peaks of each species were quantified in each sample and resulting read counts were normalized using a loess correction (see "Methods" section). A differential analysis similar to the one used for RNA-seq genes was then performed on normalized counts (see "Methods" section). We identified from 4800 (goat) to 13,600 (chicken) differentially accessible (DA) peaks between T cells and liver (Additional file 1: Table S13; Additional file 10). To test for the presence of

regulatory signals in these regions, we computed the density of transcription factor binding sites (TFBS) in ATAC-seq peaks genome-wide. Interestingly, TFBS density was significantly higher in DA ATAC-seq peaks compared to non-DA ATAC-seq peaks (Model 2;  $p$  value  $< 7.1 \times 10^{-4}$  for goat and  $p$  value  $< 2.2 \times 10^{-16}$  for chicken and pig, Wilcoxon tests, see "Methods" section). This enrichment was also observed for distal ATAC-seq peaks, at least 5 kb away from promoters (not shown), and suggests that differentially accessible peaks are more likely to have a regulatory role than globally accessible peaks.

**Promoter accessibility is associated with both positive and negative regulation of gene expression**  
 Accessible promoters are commonly associated with gene activation [43, 44]. Given the specific distribution of the ATAC-seq signal, we initially focused on proximal chromatin peaks (i.e., at 1 kb or less from a gene TSS) and used them to assign a promoter accessibility value to each gene. Using normalized read counts (see "Methods" section), we investigated the correlation between ATAC-seq and RNA-seq data either across all genes within each sample, or across all samples for each gene.

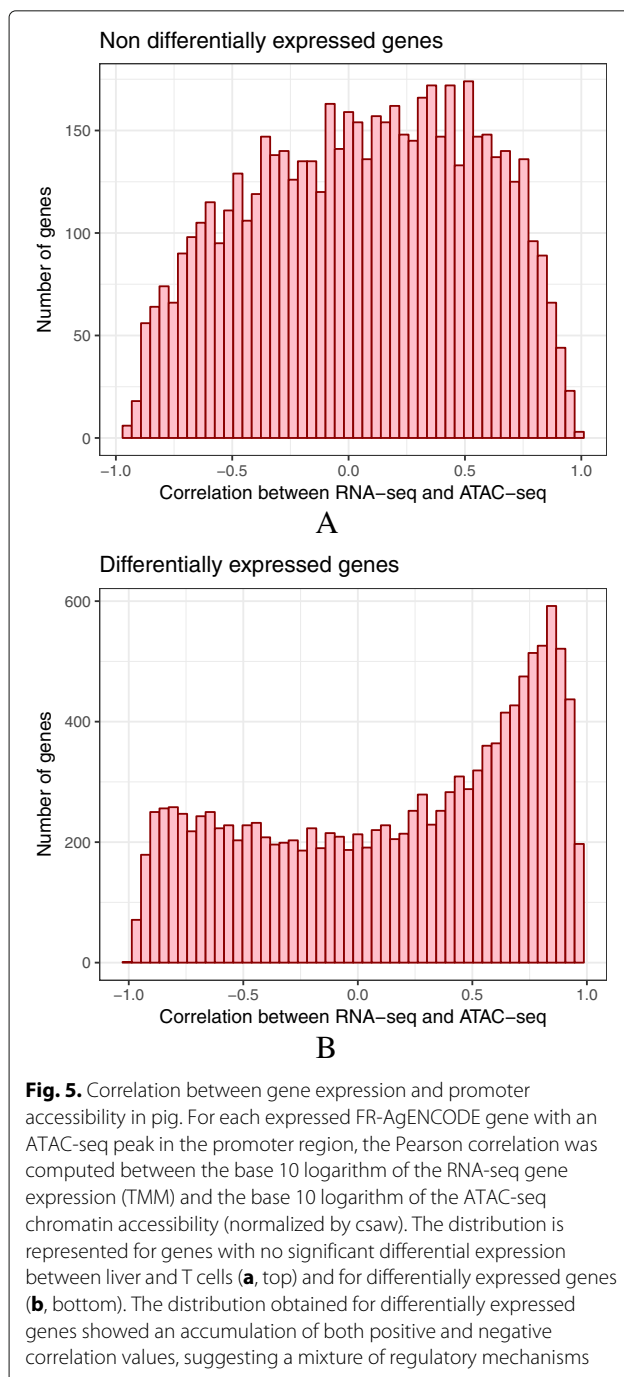


Within each sample, genes with highly accessible promoters showed higher expression values globally (Additional file 1: Figure S18), as already reported in mouse and human [45]. For pig and goat, the number of available samples further allowed us to compute for each gene the correlation between promoter accessibility and gene expression across all samples (Additional file 1: Figure S19 and “Methods” section). Interestingly, while the correlation value distribution appeared to be unimodal for non-differentially expressed genes, it was bimodal for differentially expressed genes, with an accumulation of both positive and negative correlation values (Fig. 5 and Additional file 1: Figure S20). This pattern supports the existence of different types of molecular mechanisms (i.e., both positive and negative) involved in gene expression regulation.

#### Comparative genomics reveals a core set of conserved chromatin accessible sites

We then investigated the evolution of chromatin accessibility genome-wide by performing a comparative analysis of all (proximal and distal) conserved ATAC-seq peaks across species. We identified conserved peaks by aligning all the sequences that corresponded to peaks from each species (both proximal and distal) to the human genome (see “Methods” section). Most peaks could be mapped globally, with an expected strong difference between mammals (72–80% of the peaks) and chicken (12% of the peaks). After keeping the best sequence hits, merging them on the human genome and retaining the unambiguous ones (see “Methods” section), we obtained a set of 212,021 human projections of livestock accessible chromatin regions, that we call human hits.





A large majority of the human hits (about 86%) originated from a single livestock species, which is consistent with previous reports about the fast evolution of regulatory elements and the species-specific feature of many enhancers [5, 42]. Nevertheless, the remaining 28,292 human hits (14%) had a conserved accessibility across two or more livestock species (“Methods” section and Additional file 11). As they share both sequence information and experimental evidence between several species, we

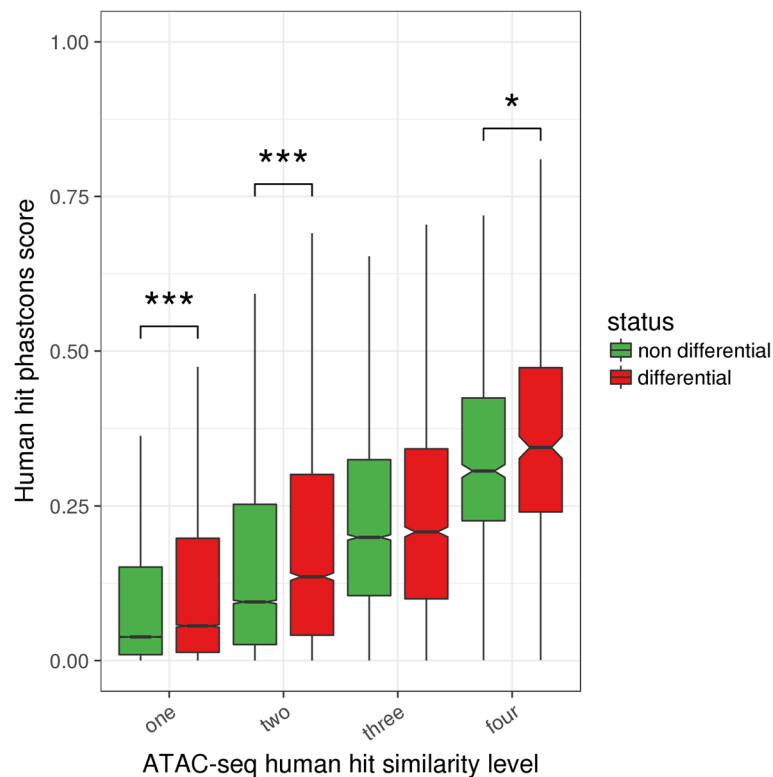
refer to those human hits as “conserved peaks” and to the number of species sharing the peak as their “similarity level”. Among these human hits, 1083 had a similarity level of 4, i.e., were shared by all 4 livestock species. Human hits from a single species were assigned a similarity level of 1. As previously done with the orthologous genes using RNA-seq data, we performed a hierarchical clustering of the samples based on the normalized accessibility values of the peaks with a similarity level of 4 (Additional file 1: Figure S21). Contrary to what was observed from the expression data, samples here mostly clustered according to species first, with the chicken as a clear outlier. However, for the two phylogenetically closest species (goat and cattle), we observed that all T cells clustered together, separately from liver. This suggests a stronger divergence and specialization of the regulatory mechanisms compared to the gene expression programs.

In addition, shuffling the peak positions within each species did not drastically change the mapping efficiency on the human genome overall but resulted in a much lower proportion of orthologous peaks (from 14 to 3% human hits with a similarity level > 1, see “Methods” section). Also, the overlap on the human genome between all the 212,021 human hits and ENCODE DNase I hypersensitive sites from liver and T cell samples [46] was three to four times higher than with the random set of shuffled peaks (25–39% per species vs. 7–9%).

Lastly, human hits that were identified as differentially accessible between liver and T cells in at least one species had higher PhastCons conservation scores on the human genome than the non differential peaks of the same similarity level (Fig. 6). This difference was significant for three out of the four similarity levels ( $p$  values < 0.01 overall, Wilcoxon tests), supporting a selective pressure on functionally active regulatory regions. Remarkably, this contrast was even stronger after discarding human hits close to a TSS in any of the species (Additional file 1: Figure S22,  $p$  values <  $10^{-6}$  overall, Wilcoxon tests, Additional file 11), in line with a specific conservation of distal regulatory elements beyond the promoter regions. Altogether, these results highlight a core set of conserved regulatory regions from birds to mammals that include both proximal and distal sites.

### Comprehensive maps of topological domains and genomic compartments in goat, chicken, and pig

In order to profile the structural organization of the genome in the nucleus, we performed in situ Hi-C on liver cells from the two male and the two female samples of pig, goat, and chicken. The in situ Hi-C protocol was applied as previously described [47] with slight modifications (see FAANG protocols online and “Methods” section). Reads were processed using a bioinformatics pipeline based on HiC-Pro [48] (“Methods” section). From 83 to 91% of



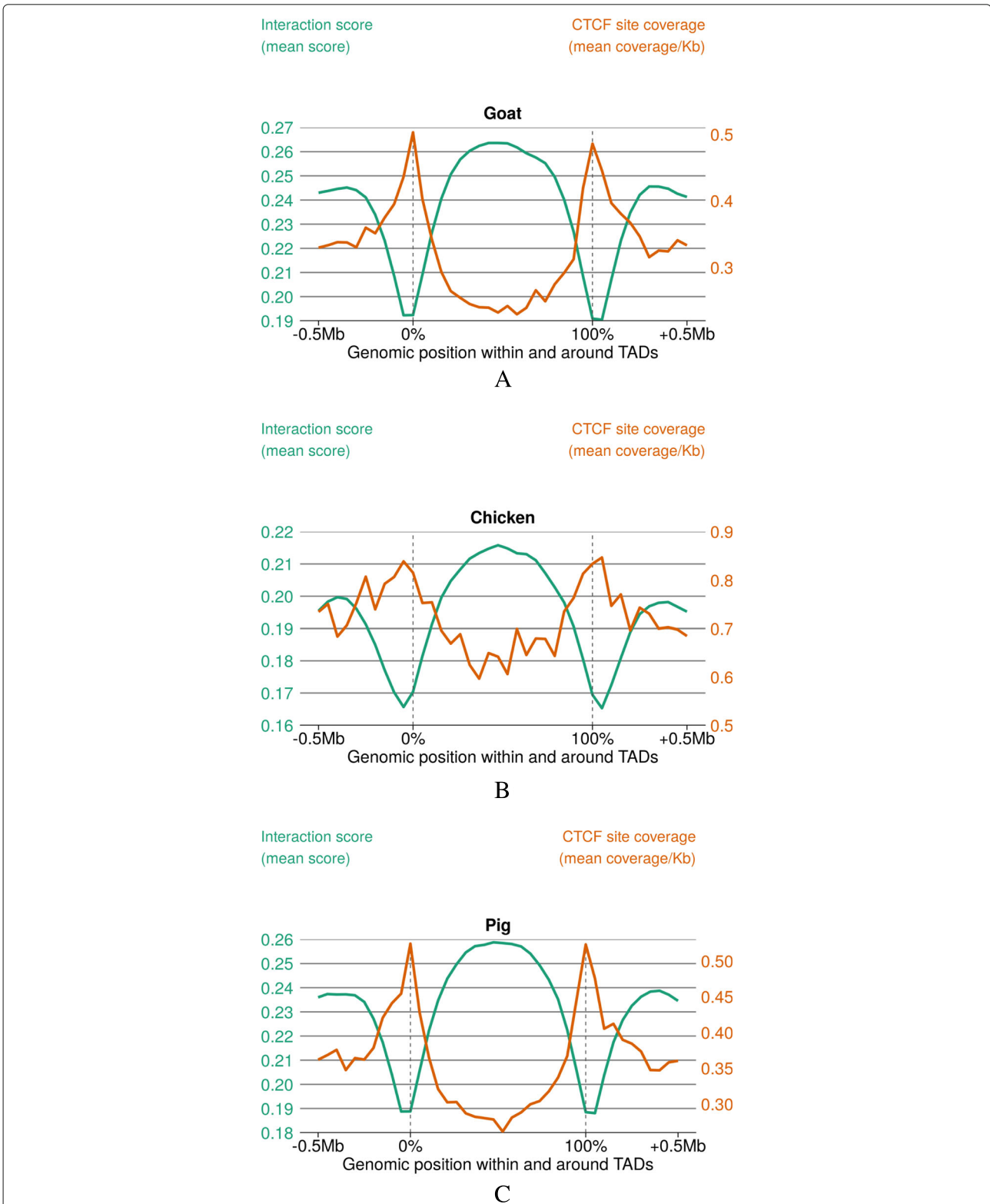
**Fig. 6.** Relationship between chromatin accessibility conservation and differential accessibility Phastcons scores of ATAC-seq human hits were plotted after dividing the human hits according to both their similarity level (between 1 and 4, x-axis) and their differential accessibility (DA) status (DA in at least one species or DA in none of the 4 species, boxplot color). Although the phastcons score obviously increases with the similarity level, it is clear that, for a given similarity level, the phastcons score is higher for DA human hits than for non DA human hits (all similarity levels except 3,  $p$  values < 0.01 overall, Wilcoxon tests) (number of elements in the boxplots from left to right: 163509, 21578, 16329, 4437, 6231, 6231, 2241, 878, 417)

the reads could be genomically mapped depending on the sample, and after filtering out all inconsistent mapping configurations we obtained a total of 182, 262, and 290M valid read pairs in goat, chicken, and pig respectively (Additional file 1: Table S14 and Figure S23). These sequencing depths allowed us to build interaction matrices (or Hi-C contact heatmaps) at the resolution of 40 and 500 kb in order to detect Topologically Associating Domains (TADs) and A/B compartments respectively (Additional file 1: Figure S24).

We identified from  $\approx 650$  (chicken) to 2000 (pig) TADs of variable sizes ( $\approx 1$  Mb on average, Additional file 1: Table S15; Additional file 12), with a 73–89% genome-wide coverage. To validate these domains predicted by Juicer [49] (see “Methods” section), we computed three metrics along the genome: the Directionality Index (DI), to quantify the degree of upstream or downstream interaction bias for any genomic region [50], the local interaction score to represent the insulation profile along the genome [51, 52], and the density of in silico predicted CTCF binding sites, expected to be prevalent at TAD boundaries [53, 54]. For each species, we observed

that the distribution of these three metrics was consistent with previous reports on model organisms (Fig. 7 and Additional file 1: Figure S25), supporting the relevance of our topological annotation. Similar results were obtained using another TAD finding tool called Armatus [55], although predicted domains were smaller (150 to 220 kb on average) and consequently more abundant (Additional file 1: Figure S25).

At a higher organizational level, we identified “active” (A) and “inactive” (B) epigenetic compartments as defined by [21] (see “Methods” section and Additional file 1: Figure S24). We obtained from  $\approx 580$  to 700 compartments per genome with a mean size between 1.6 Mb (chicken) and 3.4 Mb (goat) and covering between 71.9% (goat) and 88.6% (pig) of the genome (see Additional file 1: Table S15; Additional file 12). We also observed high consistency of the compartments between replicates (same compartment for 80% of the loci in all 4 animals, see Additional file 1: Figure S26). In model organisms, A compartments represent genomic regions enriched for open chromatin and transcription compared to B compartments [50]. By using RNA-seq and ATAC-seq data



**Fig. 7.** CTCF motif density and local interaction score within and around TADs. Local interaction score across any position measured from Hi-C matrices and represented on the y-axis (left). The mean density of predicted CTCF binding sites is also shown on the y-axis (right). Mean interaction score and CTCF density are plotted relative to the positions of Hi-C-derived Topologically Associating Domains. Dotted lines indicate TAD boundaries. Absolute scale is used on the x-axis up to 0.5 Mb on each side of the TADs while relative positions are used inside the domains (from 0 to 100% of the TAD length)

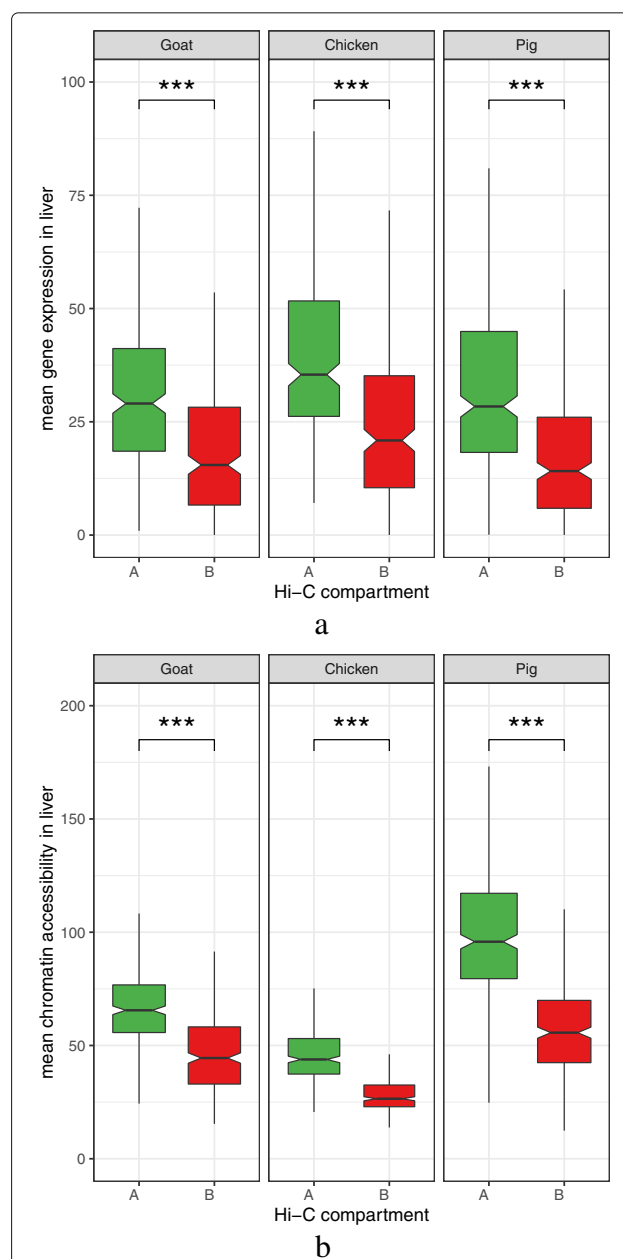
from the same liver samples as those for which we had Hi-C data, we observed that, as expected, both the average gene expression and the average chromatin accessibility were significantly higher in A than in B compartments (Fig. 8,  $p$  value  $< 2.2 \times 10^{-16}$  for each comparison, Wilcoxon tests), emphasizing the biological consistency of our results across all molecular assays and species.

#### Genome structure comparison reveals a multi-scale selective pressure on topological features across evolution

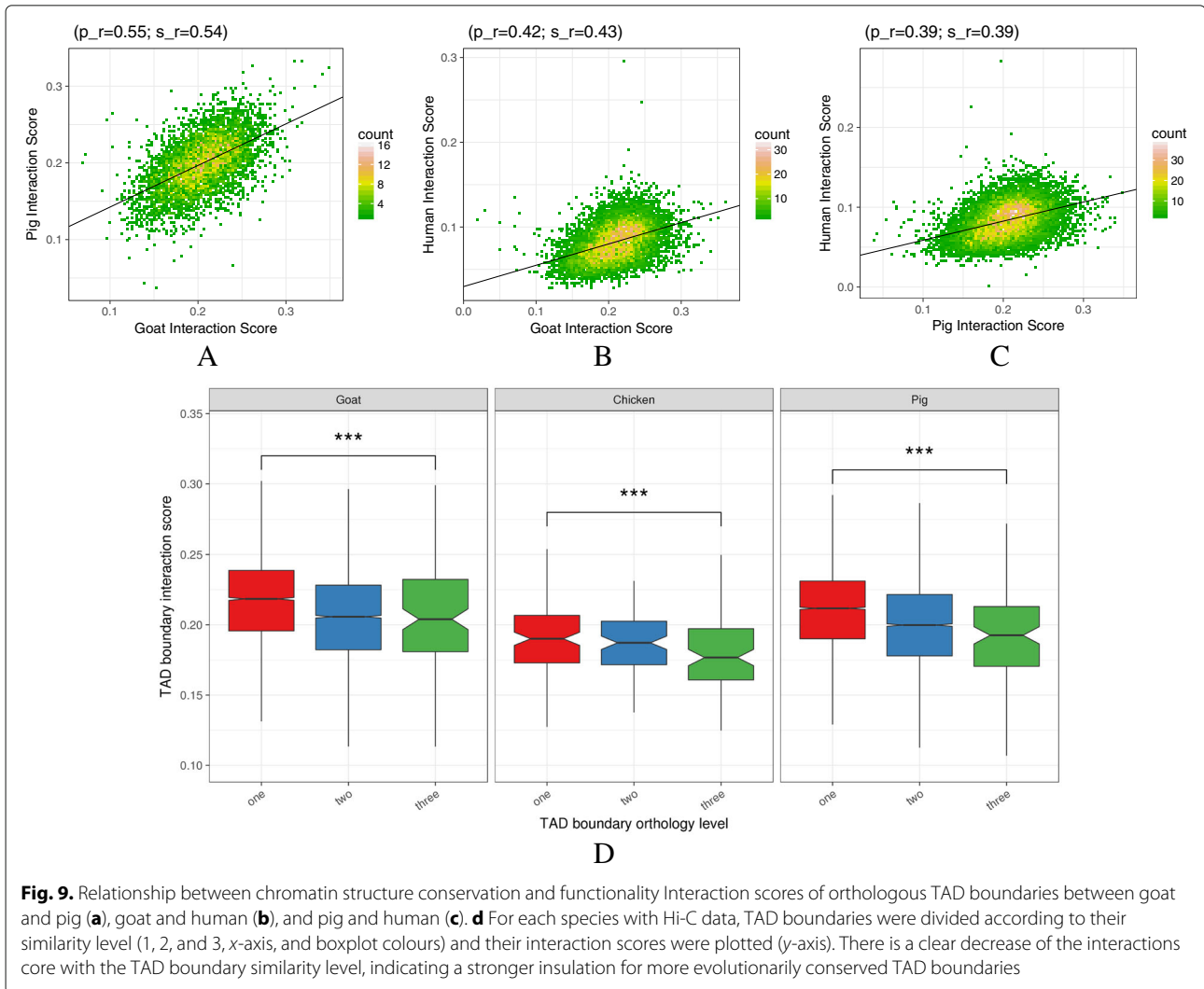
It has been shown that the general organization in TADs tends to be conserved across species [54, 56] and that the presence of specific TAD boundaries can be crucial for biological functions like development [57]. In line with these reports, we wondered if TAD boundaries might play a fine grain regulatory role beyond a binary model of simple absence/presence. Under this assumption, we hypothesized that the insulating capacity of conserved TAD boundaries could be under selective pressure. We therefore assessed the link between their insulation potential and their evolutionary conservation.

As previously done with ATAC-seq peaks (see “Methods” section and above) we first mapped all the TAD boundaries from each species to the human genome to identify the orthologous ones. Pairwise comparisons of their local interaction scores showed a clear correlation between our species (Fig. 9a). Since the interaction score here reflects the proportion of cis-contacts across a TAD boundary, such a correlation supports a conservation of the insulation strength between adjacent TADs. Strikingly, similar correlations were obtained between each of our mammalian species and human (GM12878 cell line, see “Methods” section and Fig. 9b, c) [47]. Beside confirming a general conservation of the TAD structure throughout evolution, these results emphasize the quantitative nature of this activity, in line with previous findings [54, 58, 59]. Moreover, a conserved insulation level at TAD boundaries suggests various degrees of functional impact and a fine control of their regulatory role, involving complex molecular mechanisms.

To further characterize this link between conservation and TAD strength, we assigned to each boundary a similarity level depending on the number of livestock species with a common hit on the human genome, as we did for ATAC-seq peaks (see above and Additional file 13). In all 3 species, we observed that the higher the similarity level of a TAD boundary, the lower its interaction score (Fig. 9d). These results revealed that TAD boundaries under stronger selective pressure had higher insulation activities and, expectedly, a more important role in terms of genome architecture and regulatory function. These conclusions complement previous findings about genome structure conservation across various evolutionary distances [58, 60].



**Fig. 8.** Gene expression (a) and chromatin accessibility (b) in A and B topological compartments. For the three species with Hi-C-derived A and B compartments, the distribution of the RNA-seq gene expression values (normalized read counts, top panel) and ATAC-seq chromatin accessibility values (normalized read counts, bottom panel) is shown per compartment type. A “active” compartments. B “repressed” compartments. As Hi-C data was only available from liver, only RNA-seq and ATAC-seq values from the same samples were considered. The significant and systematic difference of gene expression and chromatin accessibility values between A and B compartments ( $p$  values  $< 2.2 \times 10^{-16}$  overall, Wilcoxon tests) confirms a general consistency between RNA-seq, ATAC-seq and Hi-C data across species



**Fig. 9.** Relationship between chromatin structure conservation and functionality Interaction scores of orthologous TAD boundaries between goat and pig (a), goat and human (b), and pig and human (c). d For each species with Hi-C data, TAD boundaries were divided according to their similarity level (1, 2, and 3, x-axis, and boxplot colours) and their interaction scores were plotted (y-axis). There is a clear decrease of the interactions core with the TAD boundary similarity level, indicating a stronger insulation for more evolutionarily conserved TAD boundaries

Unlike TADs, chromosomal A/B compartments have often been reported as highly variable between tissues or developmental stages, involving dynamic mechanisms of epigenetic control [61–63]. Here, we postulated that despite its plasticity, the structural organization in genome compartments for a given tissue could also be under selective pressure across a large phylogenetic spectrum, as shown in closely related species [58]. As active compartments are known to be gene-rich we first confirmed that, although both compartment types roughly have the same size, most genes were found in A compartments in each species. In addition, we observed that the general proportion of genes in A compartments was remarkably stable across species (66.9%, 69.7%, and 70.1% of all genes in chicken, goat, and pig respectively). The 5728 orthologous genes with a predicted compartment in all three species were also found to be preferentially localized in active compartments, with slightly higher proportions than for all genes in general (69.5%, 75.9%, and

76.4% for chicken, goat, and pig respectively), probably due to the fact that conserved genes tend to have higher expression levels.

Since all these orthologous genes were assigned a compartment type (i.e., A or B label) in each species separately, we tested whether any significant conservation of compartment type across species could be detected. Among the 5728 orthologous genes, 3583 had the same compartment type in all species, which was 49% more than expected by chance assuming independence between species. This cross-species conservation was observed for both A and B compartments ( $p$  value  $< 2.2 \times 10^{-17}$  for both,  $\chi^2$  goodness-of-fit test), suggesting that such conservation was not restricted to regions of higher gene expression.

Altogether, results from the cross-species comparisons of ATAC-seq peaks, TAD boundaries and A/B compartments reveal a general conservation of the genome structure at different organizational levels from birds to



mammals, and shed a new light on the complex interplay between genome structure and function.

## Conclusion

We report the first multi-species and multi-assay genome annotation results obtained by a FAANG project. The main outcomes were consistent with our expectations and provide new evolutionary insights about regulatory and structural aspects of the genome:

- Despite only three different tissues being used, a majority of the reference transcripts could be detected. Moreover, the newly identified transcripts considerably enrich the reference annotations of these species.
- Differential analyses of gene expression in liver and T cells yielded results consistent with known metabolism and immunity functions and identified novel interesting candidates for functional analyses, including conserved syntenic lncRNAs.
- ATAC-seq data allowed an abundance of potential regulatory regions to be mapped, and, upon integration with RNA-seq data, suggested complex mechanisms of gene expression regulation. Comparative genomics analyses revealed evolutionary conservation both for proximal and distal regulators.
- Hi-C experiments provided the first set of genome-wide 3D interaction maps of the same tissue from three livestock species. Beyond the chromosome topology annotation, the analysis showed high consistency with gene expression and chromatin accessibility. Multi-species analyses revealed a global selective pressure on organizational features of the genome structure at different scales, beyond the TAD level.

Therefore, the FR-AgENCODE group has delivered a strong proof of concept of a successful collaborative approach at a national scale to apply FAANG guidelines to various experimental procedures and animal models. This notably includes the set up of a combination of sequencing assays on primary cells and tissue-dissociated cells, as well as a large collection of documented tissue samples available for further projects. It also confirmed, in line with several studies in model species [4–6, 8] the value of combining molecular assays on the same samples to simultaneously identify the transcriptomes and investigate underlying regulatory mechanisms.

In the context of the global domesticated animal genome annotation effort, lessons learned from this pilot project confirm conclusions drawn by the FAANG community regarding the challenges to be addressed in the future [13]. Furthermore, the mosaic nature of a global annotation effort that gathers contributions from various

partners worldwide emphasizes the challenge of translating recent advances from the field of data science into efficient methods for the integrative analysis of 'omics data and the importance of future meta-analyses of several datasets [16].

Altogether, these annotation results will be useful for future studies aiming to determine which subsets of putative regulatory elements are conserved, or diverge, across animal genomes representing different phylogenetic taxa. This will be beneficial for devising efficient annotation strategies for the genomes of emerging domesticated species.

## Methods

### Animals, sampling, and tissue collections

#### Animals and breeds

Well-characterized breeds were chosen in order to obtain well-documented samples. Holstein is the most widely used breed for dairy cattle. For goats, the Alpine breed is one of the two most commonly used dairy breeds, and for pigs, the Large white breed is widely used as a dam line. For chickens, the White Leghorn breed was chosen as it provides the genetic basis for numerous experimental lines and is widely used for egg production.

Four animals were sampled for each species, two males and two females. They all had a known pedigree. Animals were sampled at an adult stage, so that they were sexually mature and had performance records, obtained in known environmental conditions. Females were either lactating or laying eggs.

All animals were fasted at least 12 h before slaughter. No chemicals were injected before slaughtering, animals were stunned and bled in a licensed slaughter facility at the INRA research center in Nouzilly.

#### Samples

Liver samples of 0.5 cm<sup>3</sup> were taken from the edge of the organ, avoiding proximity with the gallbladder and avoiding blood vessels. Time from slaughter to sampling varied from 5 min for chickens to 30 min for goats and pigs and 45 min for cattle. For the purpose of RNA-seq, samples were immediately snapfrozen in liquid nitrogen, stored in 2-ml cryotubes and temporarily kept in dry ice until final storage at –80°C.

For mammals, whole blood was sampled into EDTA tubes before slaughter; at least one sampling took place well before slaughter (at least 1 month) and another just before slaughter, in order to obtain at least 50 ml of whole blood for separation of lymphocytes (PBMC). PBMC were re-suspended in a medium containing 10% FCS, counted, conditioned with 10% DMSO and stored in liquid nitrogen prior to the sorting of specific cell types: CD3+CD4+ ("CD4") and CD3+CD8+ ("CD8").

For chicken, spleen was sampled after exsanguination. Spleen leucocytes were purified by density-gradient separation to remove nucleated erythrocytes contamination and stored in liquid nitrogen prior to CD4+ and CD8+ T cell sorting.

All protocols for liver sampling, PBMC separation, splenocyte purification, and T cell sorting can be found at <http://ftp.faaang.ebi.ac.uk/ftp/protocols/samples/>

### Experimental assays and protocols

All assays were performed according to FAANG guidelines and recommendations, available at <http://www.faaang.org>. All detailed protocols used for RNA extraction and libraries production for RNA-seq, ATAC-seq, and Hi-C are available at <http://ftp.faaang.ebi.ac.uk/ftp/protocols/assays/>.

#### RNA extraction

Cells and tissues were homogenized in TRIzol reagent (Thermo) using an ULTRA-TURRAX (IKA-Werke) and total RNAs were extracted from the aqueous phase. They were then treated with TURBO DNase (Ambion) to remove remaining genomic DNA and then processed to separate long and small RNAs using the mirVana miRNA Isolation kit. Small and long RNA quality was assessed using an Agilent 2100 Bioanalyzer and RNA 6000 nano kits (Agilent) and quantified on a Nanodrop spectrophotometer.

#### RNA-seq

Stranded mRNA libraries were prepared using the TruSeq Stranded mRNA Sample Prep Kit -V2 (Illumina) on 200 ng to 1  $\mu$ g of total long RNA with a RNA Integrity Number (RIN) over 8 following the manufacturer's instructions. Libraries were PCR amplified for 11 cycles and library quality was assessed using the High Sensitivity NGS Fragment Analysis Kit DNF-474 and the Fragment Analyser system (AATI). Libraries were loaded onto a High-seq 3000 (Illumina) to reach a minimum read numbers of 100M paired reads for each library.

#### Hi-C

In situ Hi-C libraries were made according to [47] with a few modifications. For all species, fresh liver biopsies were dissociated using Accutase, and each resulting cell suspension was filtered using a 70  $\mu$ m cell strainer. Cells were then fixed with 1% formaldehyde for 10 min at 37 °C and fixation was stopped by adding Glycine to a final concentration of 0.125M. After two washes with PBS, cells were pelleted and kept at -80°C for long term storage. Subsequently, cells were thawed on ice and 5 million cells were processed for each Hi-C library. Cell membranes were disrupted using a potter-Elvehjem PTFE pestle and nuclei were then permeabilized using 0.5% SDS with

digestion overnight with HindIII endonuclease. HindIII-cut restriction sites were then end-filled in the presence of biotin-dCTP using the Klenow large fragment and were religated overnight at 4 °C. Nucleus integrity was checked using DAPI labelling and fluorescence microscopy. Nuclei were then lysed and DNA was precipitated and purified using Agencourt AMPure XP beads (Beckman Coulter) and quantified using the Qubit fluorimetric quantification system (Thermo). Hi-C efficiency was controlled by PCR using specific primers for each species and, if this step was successful, DNA was used for library production. DNA was first treated with T4 DNA polymerase to remove unligated biotinylated ends and sheared by sonication using a M220 Covaris ultra-sonicator with the DNA 550 pb Snap-Cap microtube program (Program length: 45 s; Picpower 50; DutyF 20; Cycle 200; Temperature 20 °C).

Sheared DNA was then size-selected using magnetic beads, and biotinylated fragments were purified using M280 Streptavidin Dynabeads (Thermo) and reagents from the Nextera\_Mate\_Pair Sample preparation kit (Illumina). Purified biotinylated DNA was then processed using the TrueSeq nano DNA kit (Illumina) following the manufacturer's instructions. Libraries were amplified for 10 cycles and then purified using Agencourt AMPure XP beads. Library quality was assessed on a Fragment Analyser (AATI) and by endonuclease digestion using NheI endonuclease. Once validated, each library was sequenced on an Illumina Hi-Seq 3000 to reach a minimum number of 150M paired reads per library. Libraries from the cattle samples failed the Quality Control steps (proportion of mapped reads, number of valid interactions) and were not included in the analysis.

#### ATAC-seq

ATAC-seq libraries were prepared according to [20] with a few modifications. For liver, cells were dissociated from the fresh tissue to obtain a single cell suspension. Cells were counted and 50,000 cells were processed for each assay. Transposition reactions were performed using the Tn5 Transposase and TD reaction buffer from the Nextera DNA library preparation kit (Illumina) for 30 min at 37 °C. DNA was then purified using the Qiagen MinElute PCR purification kit. Libraries were first amplified for 5 cycles using custom-synthesized index primers and then a second amplification was performed. The appropriate number of additional PCR cycles was determined using real-time PCR, permitting the cessation of amplification prior to saturation. The additional number of cycles needed was determined by plotting the Rn versus Cycle and then selecting the cycle number corresponding to one-third of the maximum fluorescent intensity. After PCR amplification, libraries were purified using a Qiagen MinElute PCR purification kit followed by an additional clean-up and sizing step using AMPure XP beads (160  $\mu$ l

of bead stock solution was added to 100  $\mu\text{l}$  of DNA in EB buffer) following the manufacturer's instructions. Library quality was assessed on a BioAnalyser (Agilent) using Agilent High Sensitivity DNA kit (Agilent), and libraries were quantified using a Qubit Fluorometer (Thermo). Considering that the Hi-C protocol was not successful on the liver samples from cattle, ATAC-seq was not attempted on these samples either.

### Bioinformatics and data analysis

All software used in this project along with the corresponding versions are listed in Additional file 1: Table S3. The reference gene annotation was obtained from the Ensembl v90 release (pig: Scrofa11.1, chicken: GalGal5, cattle: UMD3.1, goat: ARS1). Since *Capra hircus* was not part of the Ensembl release, we used the NCBI CHIR\_ARS1 annotation (see Additional file 1: Table S2).

### RNA-seq

**RNA-seq pipeline** Prior to any processing, all RNA-seq reads were trimmed using cutadapt version 1.8.3. Reads were then mapped twice using STAR v2.5.1.b [22, 23]: first on the genome indexed with the reference gene annotation to quantify expression of reference transcripts, and secondly on the same genome indexed with the newly generated gene annotation (FR-AgENCODE transcripts) (see below and Additional file 1: Figure S1) [64]. The STAR `-quantMode TranscriptomeSAM` option was used in both cases in order to additionally generate a transcriptome alignment (bam) file. After read mapping and CIGAR-based softclip removal, each sample alignment file (bam file) was processed with Cufflinks 2.2.1 [65, 66] with the `max-intron-length` (100000) and `overlap-radius` (5) options, guided by the reference gene annotation (`-GTF-guide` option) ([64], Additional file 1: Figure S1). All cufflinks models were then merged into a single gene annotation using Cuffmerge 2.2.1 [65, 66] with the `-ref-gtf` option. The transcript and gene expressions on both the reference and the newly generated gene annotation were quantified as TPM (transcripts per million) using RSEM 1.3.0 [24] on the corresponding transcriptome alignment files ([64], Additional file 1: Figure S1). The newly generated transcripts were then processed with FEELnc version 0.1.0 [36] in order to classify them into "lncRNA", "mRNA" and "otherRNA" (Additional file 1: Figure S1, Tables S10–11, Figure S9). The newly generated transcripts with a TPM value of at least 0.1 in at least 2 samples were called FR-AgENCODE transcripts and kept as part of the new annotation. The 0.1 threshold was chosen knowing that the expression values of polyadenylated transcripts usually go from 0.01 to 10,000 [35] and that we wanted to simultaneously capture long non coding RNAs that are generally lowly expressed and reduce the risk of calling artefactual transcripts.

**PCA based on gene expression** Principal Component Analysis (PCA) was performed using the **mixOmics** R package [67] on the RNA-seq sample quantifications of each species. This was done using the expression (TPM) of two different sets of genes: reference genes with TPM 0.1 in at least two samples (Additional file 1: Figure S3) and FR-AgENCODE genes with TPM 0.1 in at least two samples (Additional file 1: Figure S11).

**Annotated gene orthologs** We used Ensembl Biomart [68] to define the set of orthologous genes across cattle, chicken and pig. Only "1 to 1" similarity relationships were kept (11,001 genes). Since goat was not part of the Ensembl annotation, goat gene IDs were added to this list using gene name as a correspondence term. The resulting 4-species orthologous set contained 9461 genes (Additional file 3).

**RNA-seq sample hierarchical clustering** Based on the expression of the 9461 orthologous genes in the 39 RNA-seq samples from the four species, the sample-by-sample correlation matrix was computed using the Pearson correlation of the  $\log_{10}$  gene TPM values (after adding a pseudocount of  $10^{-3}$ ). We then represented this sample by sample correlation matrix as a heatmap where the samples were also clustered using a complete linkage hierarchical clustering (Fig. 2).

**RNA-seq normalization and differential analysis** To perform the differential analysis of gene expression, we used the expected read counts provided by RSEM [24]. RNA-seq library size normalization factors were calculated using the weighted Trimmed Mean of M-values (TMM) approach of [69] as implemented in the R/Bioconductor package **edgeR** [29]. The same package was used to fit three different per-gene negative binomial (NB) generalized log-linear models [70].

- In **Model 1**, the expression of each gene was explained by a tissue effect; because all three tissues (liver, CD4, CD8) were collected from each animal, an animal effect was also included to account for these repeated measures:

$$\frac{\log \mu_{gi}}{s_i} = \beta_{g,\text{tissue}(i)} + \gamma_{g,\text{animal}(i)},$$

where  $\mu_{gi}$  represents the mean expression of gene  $g$  in sample  $i$ ,  $s_i$  the TMM normalization factor for sample  $i$ ,  $\text{tissue}(i) \in \{\text{liver, CD4, CD8}\}$  and  $\text{animal}(i) \in \{1, 2, 3, 4\}$  the tissue and animal corresponding to sample  $i$ , and  $\beta_{g,\text{tissue}(i)}$  and  $\gamma_{g,\text{animal}(i)}$  the fixed tissue and animal effects, respectively, of gene  $g$  in sample  $i$ . Hypothesis tests were performed to identify significantly differentially

expressed genes among each pair of tissues, e.g.,

$$\mathcal{H}_{0g} : \beta_{g,\text{liver}} = \beta_{g,\text{CD4}}.$$

- **Model 2** is identical to the previous model, where gene expression was modeled using both a tissue and an animal effect, with the exception that the CD4 and CD8 tissues were collapsed into a single group. In this model, the only hypothesis of interest is thus between the liver and global CD cell group:

$$\mathcal{H}_{0g} : \beta_{g,\text{liver}} = \beta_{g,\text{CD}}.$$

All hypothesis tests were performed using likelihood-ratio tests and were corrected for multiple testing with the Benjamini-Hochberg (FDR, [71]) procedure. Genes with an FDR smaller than 5% and an absolute log-fold change larger than 2 were declared differentially expressed.

**GO analysis of differentially expressed genes** For each species, GO term enrichment analysis was performed on the genes found to be over- or under-expressed in liver versus T cells. This analysis was done separately for each species (Additional file 1: Figure S5-S7) and subsequently for genes identified for all species (Additional file 1: Fig. S8), using the three following ontologies: biological process (BP), molecular function (MF) and cell compartment (CC), and using the **Gostat** R/Bioconductor package [72]) for only those genes with a human ortholog.

#### FR-AgENCODE transcript positional classification

The FR-AgENCODE transcript models were first classified according to their position with respect to reference transcripts:

**known** the FR-AgENCODE transcript is strictly identical to a reference transcript (same intron chain and same initial and terminal exons)

**extension** the FR-AgENCODE transcript extends a reference transcript (same intron chain but one of its two most extreme exons extends the reference transcript by at least one base pair)

**alternative** the FR-AgENCODE transcript corresponds to a new isoform (or variant) of a reference gene, i.e., the FR-AgENCODE transcript shares at least one intron with a reference transcript but does not belong to the above categories

**novel** the FR-AgENCODE transcript is in none of the above classes

#### FR-AgENCODE transcript coding classification

The FR-AgENCODE transcript models were also classified according to their coding potential. For this, the FEELnc program (release v0.1.0) was used to discriminate long non-coding RNAs from protein-coding RNAs. FEELnc includes three consecutive modules:

FEELnc<sub>filter</sub>, FEELnc<sub>codpot</sub> and FEELnc<sub>classifier</sub>. The first module, FEELnc<sub>filter</sub>, filters out non-lncRNA transcripts from the assembled models, such as transcripts smaller than 200 nucleotides or those with exons strandedly overlapping exons from the reference annotation. This module was used with default parameters except `-b transcript_biotype=protein_coding`, `pseudogene` to remove novel transcripts overlapping protein-coding and pseudogene exons from the reference. The FEELnc<sub>codpot</sub> module then calculates a coding potential score (CPS) for the remaining transcripts based on several predictors (such as multi k-mer frequencies and ORF coverage) incorporated into a random forest algorithm [73]. In order to increase the robustness of the final set of novel lncRNAs and mRNAs, the options `-mode=shuffle` and `-spethres=0.98,0.98` were set. Finally, the FEELnc<sub>classifier</sub> classifies the resulting lncRNAs according to their positions and transcriptional orientations with respect to the closest annotated reference transcripts (sense or antisense, genic or intergenic) in a 1Mb window (`-maxwindow=1000000`).

It is worth noting that between 83 and 2718 lncRNA transcripts were not classified because of their localization on the numerous unassembled contigs in livestock species with no annotated genes.

#### FR-AgENCODE gene conservation between species

In order to obtain gene orthology relationships, we projected FR-AgENCODE transcripts from the four livestock species to the human GRCh38 genome using the UCSC pslMap program (<https://github.com/ENCODE-DCC/kentUtils/tree/master/src/hg/utills/pslMap>, v302). More precisely, we used the UCSC `hg38To[species.assembly].over.chain.gz` chain file for each species (created in-house for goat following UCSC instructions) and retained only the best hit for each transcript (according to the pslMap score). We further required each FR-AgENCODE gene to project to a single human gene that did not strandedly overlap any other projected FR-AgENCODE gene.

**Syntenic conservation of lncRNAs** Briefly, a lncRNA was considered as “syntenically” conserved between two species if (1) the lncRNA was located between two orthologous protein-coding genes, (2) the lncRNA was the only one in each species between the two protein-coding genes, and (3) the relative gene order and orientation of the resulting triplet was identical between species. Using these criteria, we found six triplets shared between the four species, 73 triplets shared between cattle, goat, and pig, and 19 triplets shared between cattle, chicken, and pig.



**ATAC-seq**

**ATAC-seq data analysis pipeline** ATAC-seq reads were trimmed with trimgalore 0.4.0 using the `-stringency 3`, `-q 20`, `-paired` and `-nextera` options (Additional file 1: Table S3). The trimmed reads were then mapped to the genome using bowtie 2 2.3.3.1 with the `-X 2000` and `-S` options [74]. The resulting sam file was then converted to a bam file with samtools 1.3.1, and this bam file was sorted and indexed with samtools 1.3.1 [75]. The reads for which the mate was also mapped and with a MAPQ  $\geq 10$  were retained using samtools 1.3.1 (`-F 12` and `-q 10` options, [75]), and finally only the fragments where both reads had a MAPQ  $\geq 10$  and which were on the same chromosome were retained.

Mitochondrial reads were then filtered out, as well as duplicate reads (with picard tools, MarkDuplicates subtool). The highest proportion of filtering was due to the MAPQ 10 and PCR duplicate filters (Additional file 1: Figure S14). The peaks were called using MACS2 2.1.1.20160309 [76] in each tissue separately using all the mapped reads from the given tissue (`-t` option) and with the `-nomodel`, `-f BAMPE` and `-keep-dup all` options. To get a single set of peaks per species, the tissue peaks were merged using mergeBed version 2.26.0 [77]. These peaks were also quantified in each sample by simply counting the number of mapped reads overlapping the peak.

ATAC-seq peaks were also classified with respect to the reference gene annotation using these eight genomic domains and allowing a peak to be in several genomic domains:

- exonic** the ATAC-seq peak overlaps an annotated exon by at least one bp
- intronic** the ATAC-seq peak is totally included in an annotated intron
- tss** the ATAC-seq peak includes an annotated TSS
- tss1Kb** the ATAC-seq peak overlaps an annotated TSS extended 1 kb both 5' and 3'
- tss5Kb** the ATAC-seq peak overlaps an annotated TSS extended 5 kb both 5' and 3'
- tts** the ATAC-seq peak includes an annotated TTS
- tts1Kb** the ATAC-seq peak overlaps an annotated TTS extended 1 kb both 5' and 3'
- tts5Kb** the ATAC-seq peak overlaps an annotated TTS extended 5 kb both 5' and 3'
- intergenic** the ATAC-seq peak does not overlap any gene extended by 5 kb on each side

**ATAC-seq differential analysis: normalization and model** Contrary to RNA-seq counts, ATAC-seq counts exhibited trended biases visible in log ratio-mean average (MA) plots between pairwise samples after normalization using the TMM approach, suggesting that an alternative

normalization strategy was needed. In particular, trended biases are problematic as they can potentially inflate variance estimates or log fold-changes for some peaks. To address this issue, a fast loess approach [78] implemented in the `normOffsets` function of the R/Bioconductor package `csaw` [79] was used to correct differences in log-counts vs log-average counts observed between pairs of samples.

As for RNA-seq, we used two different differential models: Model 1 for tissue pair comparisons, Model 2 for T cell versus liver comparisons (see corresponding “RNA-seq” section for more details).

**ATAC-seq peak TFBS density** In order to identify Transcription Factor Binding Sites (TFBS) genome-wide, we used the FIMO [80] software (Additional file 1: Table S3) to look for genomic occurrences of the 519 TFs catalogued and defined in the Vertebrate 2016 JASPAR database [81]. We then intersected these occurrences with the ATAC-seq peaks of each species and computed the TFBS density in differential vs non differential ATAC-seq peaks. Among the predicted TFBSs, those obtained from the CTCF motif were used to profile the resulting density with respect to Topological Associating Domains from Hi-C data (Fig. 7, Additional file 1: Figure S25).

**Comparison between ATAC-seq peaks and ChIP-seq histone mark peaks** Pig liver H3K4me3 and H3K27ac ChIP-seq peaks from the Villar et al. study [42] were downloaded from ArrayExpress (experiment number E-MTAB-2633). As these peaks were provided on the 10.2 pig genome assembly, they were first lifted over to the 11.1 pig genome assembly using the UCSC liftover program (<https://genome.sph.umich.edu/wiki/LiftOver>). About 86.7% (9632 out of 11,114) of the H3K4me3 peaks and 91.8% (31,161 out of 33,930) of the H3K27ac peaks could be lifted over to the 11.1 genome assembly. The median peak size was 1944 bp for H3K4me3 and 2786 bp for H3K27ac, and the peak size distribution was very similar for the initial 10.2 and the lifted over 11.1 peaks. As for genome coverage, the H3K4me3 and H3K27ac peaks covered 0.9% and 4.7% of the 11.1 pig genome, respectively. In comparison, there were 25,885 pig liver ATAC-seq peaks with a median size of 360 bp and covering 0.5% of the pig genome. Consistent with what was expected from the two histone marks, the vast majority (94.9%) of the H3K4me3 peaks (known to be associated to promoter regions) overlapped (bedtools intersect program) with the H3K27ac peaks (known to be associated to both promoter and enhancer regions), and about 30% of the H3K27ac peaks overlapped with the H3K4me3 peaks. Comparing our pig liver ATAC-seq peaks to the histone mark peaks, we found that 27.1% (7012 out of 25,885) and 36.4% (9410 out of 25,885) of our pig liver ATAC-seq peaks overlapped with



the H3K4me3 and H3K27ac peaks, respectively. Reciprocally, 70.3% (6773 out of 9632) and 28.3% (8821 out of 31,161) of the H3K4me3 and H3K27ac peaks respectively overlapped with our pig liver ATAC-seq peaks.

To assess if these numbers were higher than expected by chance, we shuffled (bedtools shuffle program) the 25,885 pig liver ATAC-seq peaks 1000 times on the pig genome and recomputed their intersection with the two sets of histone mark peaks (H3K4me3 and H3K27ac). After doing so, we never obtained percentages of H3K4me3 and H3K27ac peaks, respectively, overlapping the shuffled ATAC-seq peaks that were equal or higher than the ones obtained with the real ATAC-seq peaks. This means that indeed, 70.3% and 28.3% of the histone mark peaks overlapping our ATAC-seq peaks are percentages that are significantly higher than expected by chance ( $p$  value  $< 10^{-3}$ ).

We also compared the ATAC-seq, H3K4me3 and H3K27ac peak scores (fold enrichment against random Poisson distribution with local lambda for ATAC-seq peaks and fold-enrichment over background for ChIP-seq peaks) of the common peaks versus the other peaks. In doing so, we found that common peaks had significantly higher scores than non common peaks (median 94 versus 32,  $p$  value  $< 2.2 \times 10^{-16}$  for ATAC-seq peaks, median 57 versus 22,  $p$  value  $< 2.2 \times 10^{-16}$  for H3K4me3 peaks and median 32 versus 12,  $p$  value  $< 2.2 \times 10^{-16}$  for H3K27ac peaks, Wilcoxon tests), highlighting a common signal between the two techniques.

### Chromatin accessibility conservation across species

In order to investigate the conservation of chromatin accessibility across our 4 livestock species, we used the human GRCh38 genome as a reference. After indexing the softmasked GRCh38 genome (main chromosomes) using lastdb (last version 956, -uMAM4 and -cR11 options, <http://last.cbrc.jp/>), we used the lastal program followed by the last-split program (-m1 and -no-split options) (last version 956, <http://last.cbrc.jp/>) to project the 104,985 cattle, 74,805 goat, 119,894 chicken, and 149,333 pig ATAC-seq peaks onto the human genome. In doing so and consistent with the phylogenetic distance between our species and human, we were able to project 72.6% (76,253) cattle, 73.7% (55,113) goat, 12.3% (14,792) chicken, and 80.1% (119,680) pig peaks to the human genome. The percentage of bases of the initial peaks that could be aligned was around 40% for mammals and 14% for chicken. Then, for each peak that could be projected onto the human genome, we retained its best hit (as provided by lastal) and then merged all these best hits (i.e., from the 4 species) on the human genome (using bedtools merge). A total of 215,620 human regions were obtained, from which we kept the 212,021 that came from a maximum of 1 peak

from each species. Those 212,021 regions were called human hits.

Based on the 1083 four-species orthologous peaks in the 38 ATAC-seq samples, the sample-by-sample correlation matrix was computed using the Pearson correlation of the  $\log_{10}$  normalized ATAC-seq values (after adding a pseudo-count of  $10^{-3}$  to the values). We then represented this sample-by-sample correlation matrix as a heatmap where the samples were also clustered using a complete linkage hierarchical clustering (Additional file 1: Figure S21). Chicken ATAC-seq samples clustered completely separately from mammal ATAC-seq samples. T cell samples from cattle and goat were also closer to each other than to liver samples.

To shuffle the 104,985 cattle, 74,805 goat, 119,894 chicken, and 149,333 pig ATAC-seq peaks, we used the bedtools shuffle program on their respective genomes and projected these shuffled peaks to the human genome as was done for the real peaks.

We also compared the human hits to the combined set of 519,616 ENCODE human DNase I peaks from two CD4+, two CD8+ and one "right lobe of liver" samples (experiment accessions ENCSR683QJJ, ENCSR167JFX, ENCSR020LUD, ENCSR316UDN, and ENCSR555QAY from the encode portal <https://www.encodeproject.org/>, by merging the peaks from the 5 samples into a single set of peaks using bedtools merge). We found that 23.1% (48,893 out of 212,021) of the human hits obtained from the real ATAC-seq peaks overlapped human DNase I peaks, whereas only 8.5% (21,159 out of 249,943) of the human hits obtained from shuffled ATAC-seq peaks overlapped human DNase I peaks. This further supports the biological signal present in these data.

Finally we used the phastcons measure of vertebrate sequence conservation obtained from the multiple alignment of 100 vertebrate species genomes including human (hg38.phastCons100way.bw bigwig file from the UCSC web site <https://genome.ucsc.edu/>). For each human hit, we computed its phastcons score using the bigWigAverageOverBed utility from UCSC (<https://github.com/ENCODE-DCC/kentUtils>).

### Hi-C

**Hi-C data analysis pipeline** Our Hi-C analysis pipeline includes HiC-Pro v2.9.0 [82] (Additional file 1: Table S3) for the read cleaning, trimming, mapping (this part is internally delegated to Bowtie 2 v2.3.3.1), matrix construction, and matrix balancing ICE normalization [83]. HiC-Pro parameters: BOWTIE2\_GLOBAL\_OPTIONS = -very-sensitive -L 30 -score-min L, -1, -0.1 -end-to-end -reorder, BOWTIE2\_LOCAL\_OPTIONS = -very-sensitive -L 20 -score-min L, -0.6, -0.2 -end-to-end -reorder, LIGATION\_SITE = AAGCTAGCTT,

MIN\_INSERT\_SIZE = 20, MAX\_INSERT\_SIZE=1000, GET\_ALL\_INTERACTION\_CLASSES = 1, GET\_PROCESS\_SAM = 1, RM\_SINGLETON = 1, RM\_MULTI = 0, RM\_DUP = 1, MAX\_ITER = 100, FILTER\_LOW\_COUNT\_PERC = 0.02, FILTER\_HIGH\_COUNT\_PERC = 0, EPS = 0.1. TAD finding was performed using two methods: arrowhead from the Juicer tool V1.5.3 [49] at the 10 kb resolution using the matrix balancing normalization (arrowhead -r 10000 -k KR option), and with Armatus V2.1 [55] with default parameters and gamma=0.5 (Additional file 1: Table S3). Graphical visualization of the matrices was produced with the HiTC R/Bioconductor package v1.18.1 [48] (Additional file 1: Table S3). Export to JuiceBox [84] was done through Juicer Tools V1.5.3 (Additional file 1: Table S3). These tools were called through a pipeline implemented in Python. Because of the high number of unassembled scaffolds (e.g., for goat) and/or micro-chromosomes (e.g., for chicken) in our reference genomes, only the longest 25 chromosomes were considered for TAD and A/B compartment calling. For these processes, each chromosome was considered separately.

The Directionality Index (DI) was computed using the original definition introduced by [50] to indicate the upstream vs. downstream interaction bias of each genomic region. Interaction matrices of each chromosome were merged across replicates and the score was computed for each bin of 40 kb. CTCF sites were predicted along the genomes by running FIMO with the JASPAR TFBS catalogue (see “ATAC-seq peak TFBS density” section).

A/B compartments were obtained using the method described in [21] as illustrated in Additional file 1: Figure S24: first, ICE-normalized counts,  $K_{ij}$ , were corrected for a distance effect with:

$$\widehat{K}_{ij} = \frac{K_{ij} - \overline{K}^d}{\sigma^d},$$

in which  $\widehat{K}_{ij}$  is the distance-corrected count for the bins  $i$  and  $j$ ,  $\overline{K}^d$  is the average count over all pairs of bins at distance  $d = d(i, j)$  and  $\sigma^d$  is the standard deviation of the counts over all pairs of bins at distance  $d$ . Within-chromosome Pearson correlation matrices were then computed for all pairs of bins based on their distance-corrected counts and a PCA was performed on this matrix. The overall process was performed similarly to the method implemented in the R/Bioconductor package HiTC [48]. Boundaries between A and B compartments were identified according to the sign of the first PC (eigenvector). Since PCAs had to be performed on each chromosome separately, the average counts on the diagonal of the normalized matrix were used to identify which PC sign (+/-) should be assigned to A and B compartments for

each chromosome. This allowed a homogenous assignment across chromosomes to be obtained, without relying on the reference annotation. In line with what was originally observed in humans, where the first PC was the best criterion for separating A from B compartments (apart from a few exceptions like chromosome 14 for instance [21]), we also observed a good agreement between the plaid patterns of the normalized correlation matrices and the sign of the first PC (Additional file 1: Figure S24).

To estimate the robustness of A/B compartment calling, the method was tested on each replicate separately (four animals). Since the HiTC filtering method can discard a few bins in some matrices, resulting in missing A/B labels, the proportion of bins with no conflicting labels across replicates was computed among the bins that had at least two informative replicates (Additional file 1: Figure S26).

**Chromatin structure conservation across species** To get insight into chromatin structure conservation across species, similar to what was done with chromatin accessibility data (see above), we projected the 11,711 goat, 6866 chicken, and 14,130 pig 40 kb TAD boundaries to the human GRCh38 genome using lastal followed by last-split (-m1 and -no-split options, last version 956, <http://last.cbrc.jp/>), using the same indexed GRCh38 softmasked genome as was used for ATAC-seq, see above). For this analysis we considered TADs from Armatus because of the high number of boundaries that were identified by the method. As expected from their length, TAD boundary projections were highly fragmented (median 16, 2, and 19 blocks per projection representing 3%, 0.6%, and 3% of the initial segment, for the best hit of goat, chicken, and pig, respectively). In order to recover conserved segments, we chained the alignments using a python in-house script (program available on demand, used with stranded mode, coverage=0.4, score=3000, and length\_cutoff=5000). Doing so, we managed to project 90% of the mammalian and 5% of the chicken TAD boundaries onto the human genome. Similar to what was done for ATAC-seq, for each projected TAD boundary, its best hit (according to the chaining score) was retained. The median length of those best hits represented 95% and 78% of the initial query size for mammals and chicken respectively. Merging these best hits on the human genome (using bedtools merge), we obtained 16,870 human regions with a median length of 44.6 kb (similar to the initial TAD boundary size of 40 kb). Out of those, 16,468 were considered non ambiguous (i.e., not coming from several TAD boundaries from the same species) and were retained for further analyses. As was found for the ATAC-seq peaks, the majority (65.6%) of the hits were single species (similarity level 1), a substantial percentage of them (34%) were 2 species hits (similarity

level 2), and seventy one of them (0.4%) were 3 species hits (similarity level 3).

To estimate the structural impact of each TAD boundary, we used the local interaction score as used by [51] and [52] and sometimes referred to as “interaction ratio” or “insulation profile”. Within a sliding window of 500 kb along the genome (step=10 kb), the insulation score ratio is defined as the proportion of read pairs that span across the middle of the window. The score ratio is reported at the middle position of the window and represents the local density of the chromatin contacts around this point. This proportion is expected to be maximal in regions with many local interactions (typically TADs) and minimal over insulators (typically TAD boundaries). Intuitively, a TAD boundary with a low interaction score (which indicates strong insulation properties) has a good capacity to prevent interactions that cross it while a TAD with a relatively high interaction score has a “weak” insulation strength. Here, only valid interactions (“valid pairs”) in *in cis* (inter-chromosomal contacts were not considered in the ratio) were considered after applying all HiC-Pro QC and filters. Computing a ratio among all read pairs that have both reads within the sliding window reduces the impact of potential biases (read coverage, restriction site density, GC content, etc.). Consequently, the interaction profiles from the 4 replicates along the genome of each species were highly similar (not shown), allowing to merge them in order to assign each TAD boundary a single score per species. For orthologous TAD boundaries, the scores from different species could be used to compute pairwise correlations. Human data were obtained from <http://aidenlab.org/data.html> [47] for the GM12878 cell line (ENCODE batch 1, HIC048.hic file from <https://bcm.app.box.com/v/aidenlab/file/95512487145>). The .hic file was parsed by the Juicer tool (“dump” mode with options “observed KR”) to compute the corresponding interaction score as described above. The LiftOver tool was used to convert the genomic positions of the human TAD boundaries (version hg19 vs. hg38) before comparing the interaction scores with livestock species.

The number and proportion of genes (all or only the orthologous ones) in each compartment type was computed using bedtools map (-distinct option on the gene ID field). Orthologous genes were taken from Ensembl as previously described. Under the independence assumption of compartment assignment between species, the expected proportion of orthologous genes with “triple A” (resp. with “triple B”) assignments between species is equal to the product of the observed frequencies for A (resp. for B) compartments in the three species. The observed frequencies of “triple A” and “triple B” assignments in orthologous genes was compared to this expected proportion using a  $\chi^2$  goodness-of-fit test.

### Multi-assay integration

**ATAC-seq vs. RNA-seq correlation: intra- and inter-sample analysis** For each ATAC-seq peak that overlapped a promoter region (1 kb upstream of the TSS, as suggested in Fig. 4) its less-normalized read count value (see differential analysis) was associated with the TMM-normalized expression of the corresponding gene from the reference annotation. Intra- and inter-sample correlations were then investigated: within each sample, genes were ranked according to their expressions and the distribution of the corresponding ATAC-seq values was computed for each quartile (Additional file 1: Figure S18). Across samples, the Pearson correlation coefficient was computed for each gene using only the samples for which both the ATAC-seq and the RNA-seq normalized values were available (e.g.,  $n = 10$  for pig, Additional file 1: Figure S19–20). Similar results were obtained with Spearman correlations (not shown).

**Chromatin accessibility and gene expression in A/B compartments** To compute the general chromatin accessibility in A and B compartments, we first computed the average of the normalized read count values across all liver samples for each ATAC-seq peak. For each compartment, the mean value of all contained peaks was then reported and the resulting distributions for all A and B compartments were reported (Fig. 8).

The same approach was used to assess the general expression of genes in A and B compartments, using the average of the normalized expression values from the liver samples. Difference between A and B distributions was tested for statistical significance using a Wilcoxon test.

### Supplementary information

Supplementary information accompanies this paper at <https://doi.org/10.1186/s12915-019-0726-5>.

Additional files are available in the Additional file section and on the FR-AgENCODE website [www.fragencode.org](http://www.fragencode.org).

**Additional file 1:** Additional file 1: Supplementary figures (S1-S26) and tables (S1-S15).

**Additional file 2:** Reference genes and transcripts (structure, expression) of the 4 species. Archive content:

- `bos_taurus.gtf`
- `bos_taurus.refgn.tpm.tsv`
- `capra_hircus.gtf`
- `capra_hircus.refgn.tpm.tsv`
- `gallus_gallus.gtf`
- `gallus_gallus.refgn.tpm.tsv`
- `sus_scrofa.gtf`
- `sus_scrofa.refgn.tpm.tsv`

**Additional file 3:** Orthologs between the 4 livestock species. We used Biomart to retrieve the 1 to 1 orthology relationships between chicken, pig and cattle and added goat via gene name. The human gene id is given for reference.

**Additional file 4:** Reference DE genes (all combinations): the archive contains four folders, one for each species (*bos\_taurus*, *capra\_hircus*, *gallus\_gallus*, *sus\_scrofa*). Each folder contains itself two subfolders, one for each model: *diffcounts.nominsum* (Model 1) and *diffcounts.cdvs\_liver* (Model 2). Results of Model 1 are given in:

- *refgenes.counts.min2tpm0.1.normcounts.diff.readme.idx*
- *refgenes.counts.min2tpm0.1.normcounts.diff.cd4.cd8.bed*
- *refgenes.counts.min2tpm0.1.normcounts.diff.cd4.liver.bed*
- *refgenes.counts.min2tpm0.1.normcounts.diff.cd8.liver.bed*

Results of Model 2 are given in:

- *refgenes.counts.min2tpm0.1.normcounts.diff.readme.idx*
- *refgenes.counts.min2tpm0.1.normcounts.diff.cd.liver.bed*

All bed files contain the coordinates and id of the genes found to be differentially expressed between the two conditions. The file also contains the normalized read counts of those genes in the different samples as well as the adjusted pvalue, logFC and normLogFC (see *readme.idx* file for more details).

**Additional file 5:** FR-AgENCODE genes and transcripts (structure, expression, positional and coding classes).

- *bos\_taurus\_cuff\_tpm0.1\_2sample\_complete.gff*
- *bos\_taurus\_cuff\_tpm0.1\_2sample\_trid\_4posclasses\_3codingclasses\_booleans.tsv*
- *bos\_taurus\_frag.gnid.posclasslist.codclasslist.tsv*
- *bos\_taurus\_fraggn.tpm.tsv*
- *capra\_hircus\_cuff\_tpm0.1\_2sample\_complete.gff*
- *capra\_hircus\_cuff\_tpm0.1\_2sample\_trid\_4posclasses\_3codingclasses\_booleans.tsv*
- *capra\_hircus\_frag.gnid.posclasslist.codclasslist.tsv*
- *capra\_hircus\_fraggn.tpm.tsv*
- *gallus\_gallus\_cuff\_tpm0.1\_2sample\_complete.gff*
- *gallus\_gallus\_cuff\_tpm0.1\_2sample\_trid\_4posclasses\_3codingclasses\_booleans.tsv*
- *gallus\_gallus\_frag.gnid.posclasslist.codclasslist.tsv*
- *gallus\_gallus\_fraggn.tpm.tsv*
- *sus\_scrofa\_cuff\_tpm0.1\_2sample\_complete.gff*
- *sus\_scrofa\_cuff\_tpm0.1\_2sample\_trid\_4posclasses\_3codingclasses\_booleans.tsv*
- *sus\_scrofa\_frag.gnid.posclasslist.codclasslist.tsv*
- *sus\_scrofa\_fraggn.tpm.tsv*

**Additional file 6:** Four livestock species FR-AgENCODE gene orthology.

**Additional file 7:** FR-AgENCODE DE genes (all combinations). The archive has the same structure than *de.refgn.tar.gz* with names starting with *cuffgenes* instead of *refgenes*.

**Additional file 8:** lncRNAs (information from FEELnc, orthology, structure, etc). Archive content:

- *bos\_taurus.lncrna.TPM0.1in2samples.classif.tsv*
- *capra\_hircus.lncrna.TPM0.1in2samples.classif.tsv*
- *ConservedLncRNABySynteny\_73\_19\_6.xlsx*
- *gallus\_gallus.lncrna.TPM0.1in2samples.classif.tsv*
- *sus\_scrofa.lncrna.TPM0.1in2samples.classif.tsv*

**Additional file 9:** ATAC-seq peaks (coordinates, quantification, positional classification): the archive contains four folders, one for each species (*bos\_taurus*, *capra\_hircus*, *gallus\_gallus*, *sus\_scrofa*). Each folder contains the following six files:

- *mergedpeaks\_allinfo\_gn\_frag.tsv*
- *mergedpeaks\_allinfo\_tr\_frag.tsv*
- *mergedpeaks\_allinfo\_tr\_ref.tsv*
- *mergedpeaks\_allinfo\_gn\_ref.tsv*
- *mergedpeaks\_peaknb.allexp.readnb.bed.readme.idx*
- *mergedpeaks\_peaknb.allexp.readnb.bed*

**Additional file 10:** DA ATAC-seq peaks (all combinations). The archive has the same structure as *de.refgn.tar.gz* with names starting with *mergedpeaks\_peaknb.allexp.readnb* instead of *refgenes.counts.min2tpm0.1*.

**Additional file 11:** Four livestock species ATAC-seq peak orthology.

**Additional file 12:** Hi-C TADs and A/B compartments: the archive contains three folders, one for each species (*capra\_hircus*, *gallus\_gallus*, *sus\_scrofa*). Each folder contains the following two files:

- *compartments.bed*
- *mat.40000.longest25chr.tad.consensus.bed*

**Additional file 13:** Three livestock species TAD boundary orthology.

## Abbreviations

DE: Differentially expressed FAANG: Functional annotation of animal genomes lncRNA: Long non-coding RNA mRNA: Messenger RNA PE: Paired-end PolyA: PolyAdenylation RT-PCR: Reverse transcriptase polymerase chain reaction TAD: Topological associating domain TF: Transcription factor TFBS: Transcription factor binding site TPM: Transcript per million TSS: Transcription start site TTS: Transcription termination site

## Acknowledgements

We would like to thank the FAANG community for the general support and in particular A. Archibald (Roslin Institute, UK), L. Clarke, P. Harrison (EMBL-EBI, UK), and M. Groenen (WUR, Netherlands) for their collaboration in the organization of the project, and B. Rosen (ARS, USDA) for providing information about the goat sexual chromosomes. We wish to thank all field operators at INRA experimental animal facilities, units, and platforms in France for the access and the assistance in animal handling and sampling, including the following: Y. Gallard, UE Le Pin (Gouffern en Auge), F. Bouvier and T. Fassier, UE Bourges (Osmoy), S. Ferchaud, UE GenESI (Magneraud/Rouillé), Y. Baumard, UE PEAT (Nouzilly), C. Berri and J. Gautron, UR BOA (Nouzilly), G. Gomot, J.P. Dubois and A. Arnould, UR PRC CIRE (Nouzilly), E. Guettier, D. Capo and J. Savoie, UE PAO (Nouzilly).

We are grateful to A. Breschi (Stanford, USA) and J. Lagarde (CRG, Spain) for sharing scripts and to N. Servant (Institut Curie, France) for assistance on HiC-PRO. Additional acknowledgements go to C. Donnadiou and O. Bouchez from the Get-Plage sequencing platform (INRA Toulouse) and to C. Gaspin and her staff at the GenoToul bioinformatics platform (INRA Toulouse), especially D. Laborie and M.S. Trotard for IT support.

## Authors' contributions

MTB and Sfa (coordination), AG, EG, FB, FD, FL, GTK, HA, MTB, PQ, SDP, Sfa, SLL, SVN (sampling and cell sorting) contributed to animal and sampling. DE and HA (coordination), DE, SDP (RNA-seq libraries), AG, EG, KMun (ATAC-seq libraries), FM, HA, and MM (Hi-C libraries) contributed to molecular assays. CK,



SD, Sfo (coordination), CC, SD (RNA-seq), KMun, SD, Sfo (ATAC-seq), KMur, SL, TD (lncRNAs), AR, NV, RML (differential analyses), DR, IG, MM, MSC, MZ, NV, Sfo (Hi-C pipeline), EC, EG, SD, SL, TD (functional analyses), SD, Sfo (metadata), AR, NV, SD, Sfo, TF (integrative analyses), HA, SD, Sfo, SM, and PB (data submission) contributed to bioinformatics and data analysis. AR, CK, EG, HA, KMun, KMur, MTB, MZ, NV, SD, Sfo, SL, and TD contributed to the manuscript writing. EG, MHP, Sfo, and SL are the management committee. EG and Sfo are responsible for project coordination. All authors have read and approved the manuscript.

### Funding

This study has been supported by the INRA "SelGen metaprogramme", grant "FR-AgENCODE: A French pilot project to enrich the annotation of livestock genomes" (2015–2017). S. Djebali, A. Rau and E. Crisci are supported by the AgreenSkills+ fellowship program with funding from the EU's Seventh Framework Program under grant agreement FP7-609398. E. Crisci was also supported by the Animal Health and Welfare ERA-Net (anihwa) - project KILLeuPRRSV. Additional financial support for tissue biobanking was provided by the CRB-Anim infrastructure project, ANR-11-INBS-0003, funded by the French National Research Agency in the context of the "Investing for the Future" program.

### Availability of data and materials

Experimental protocols for tissue sampling and molecular assays are publicly available at [18]. Sample records are available at the BioSamples database [85] under submission codes GSB-99 and GSB-721 [86]. Biological material from the 16 animals (tissue samples and aliquots) are stored at the INRA CRB-Anim BioBanking facility and are available on request. Raw sequences and metadata are available at the EMBL-EBI's European Nucleotide Archive ENA and at the FAANG Data Coordination Center using accessions PRJEB27455 (RNA-seq), PRJEB27111 (ATAC-seq) and PRJEB27364 (Hi-C) [18]. Additional data and results (including annotation files, list of differentially expressed genes with normalized expression values, annotated ATAC-seq peaks with raw read counts, differentially accessible peaks with normalized read counts, Hi-C matrices, TADs and A/B compartments) are available at [19]. All data generated or analysed during this study are included in this published article and its supplementary information files.

### Ethics approval and consent to participate

All animal handling and sampling was realized in conformity with the French legislation on animal experimentation. For mammals, blood samples were collected in the context of the following approvals (APAFIS/project#): 334-2015031615255004\_v4 and 333-2015031613482601\_v4 (pigs), 3066-201511301610897\_v2 (cattle), 03936.02 and 8613-2017012013585646\_v4 (goats). Chicken immune cells were obtained from spleen sampled after slaughter (no need for animal experiment authorization).

### Consent for publication

Not applicable.

### Competing interests

The authors declare that they have no competing interests.

### Author details

<sup>1</sup>GenPhySE, Université de Toulouse, INRA, INPT, ENVT, Chemin de Borde Rouge, F-31326 Castanet-Tolosan Cedex, France. <sup>2</sup>Curtin University, School of Pharmacy & Biomedical Sciences, CHIRI Biosciences, 24105 Perth, Australia. <sup>3</sup>MIAT, Université de Toulouse, INRA, Chemin de Borde Rouge, F-31326 Castanet-Tolosan Cedex, France. <sup>4</sup>GABI, AgroParisTech, INRA, Université Paris Saclay, F-78350 Jouy-en-Josas, France. <sup>5</sup>PEGASE, Agrocampus-Ouest, INRA, F-35590 Saint-Gilles Cedex, France. <sup>6</sup>INRA, US1426, GeT-PlaGe, Genotoul, Chemin de Borde Rouge, F-31326 Castanet-Tolosan Cedex, France. <sup>7</sup>UMR6290 IGDR, CNRS, Université Rennes 1, Rennes, France. <sup>8</sup>UMR1282 ISP, INRA, F-37380 Nouzilly, France. <sup>9</sup>IRCM, CEA, Université Paris Saclay, F-78350 Jouy-en-Josas, France. <sup>10</sup>Department of Population Health and Pathobiology, College of Veterinary Medicine, North Carolina State University, NC 27607 Raleigh, USA.

Received: 21 October 2019 Accepted: 19 November 2019

Published online: 30 December 2019

### References

- Siepel A, Bejerano G, Pedersen JS, Hinrichs AS, Hou M, Rosenbloom K, et al. Evolutionarily conserved elements in vertebrate, insect, worm, and yeast genomes. *Genome Res.* 2005;15:1034–50.
- Hindorff LA, Sethupathy P, Junkins HA, Ramos EM, Mehta JP, Collins FS, et al. Potential etiologic and functional implications of genome-wide association loci for human diseases and traits. *Proc Natl Acad Sci.* 2009;106:9362–7.
- Maurano MT, Humbert R, Rynes E, Thurman RE, Haugen E, Wang H, et al. Systematic localization of common disease-associated variation in regulatory DNA. *Science.* 2012;337:1190–95.
- The ENCODE Project Consortium. An integrated encyclopedia of DNA elements in the human genome. *Nature.* 2012;489:57–74.
- Yue F, Cheng Y, Breschi A, Vierstra J, Wu W, Ryba T, et al. A comparative encyclopedia of DNA elements in the mouse genome. *Nature.* 2014;515:355–64.
- Gerstein MB, Rozowsky J, Yan KK, Wang D, Cheng C, Brown JB, et al. Comparative analysis of the transcriptome across distant species. *Nature.* 2014;512:445–8.
- Vernot B, Stergachis AB, Maurano MT, Vierstra J, Neph S, Thurman RE, et al. Personal and population genomics of human regulatory variation. *Genome Res.* 2012;22:1689–97.
- Kundaje A, Meuleman W, Ernst J, Bilenyk M, Yen A, Heravi-Moussavi A, et al. Integrative analysis of 111 reference human epigenomes. *Nature.* 2015;518:317–30.
- Stergachis AB, Neph S, Sandstrom R, Haugen E, Reynolds AP, Zhang M, et al. Conservation of trans-acting circuitry during mammalian regulatory evolution. *Nature.* 2014;515:365–70.
- Cheng Y, Ma Z, Kim BH, Wu W, Cayting P, Boyle AP, et al. Principles of regulatory information conservation between mouse and human. *Nature.* 2014;515:371–5.
- Schmidt D, Wilson MD, Ballester B, Schwalie PC, Brown GD, Marshall A, et al. Five-vertebrate ChIP-seq reveals the evolutionary dynamics of transcription factor binding. *Science.* 2010;328:1036–40.
- The FAANG Consortium. Functional Annotation of Animal Genomes (FAANG): a coordinated international action to accelerate genome to phenome. <http://www.faaang.org>. Accessed 13 Nov 2019.
- Andersson L, Archibald AL, Bottema CD, Brauning R, Burgess SC, Burt DW, et al. Coordinated international action to accelerate genome-to-phenome with FAANG, the Functional Annotation of Animal Genomes project. *Genome Biol.* 2015;16:57.
- Tuggle CK, Giuffra E, White SN, Clarke L, Zhou H, Ross PJ, et al. GO-FAANG meeting: a gathering on functional annotation of animal genomes. *Anim Genet.* 2016;47:528–33.
- Kern C, Wang Y, Chitwood J, Korf I, Delany M, Cheng H, et al. Genome-wide identification of tissue-specific long non-coding RNA in three farm animal species. *BMC Genomics.* 2018;19:684.
- Giuffra E, Tuggle CK, FAANG Consortium. Functional annotation of animal genomes (FAANG): current achievements and roadmap. *Ann Rev Anim Biosci.* 2019;7:65–88.
- Harrison P, Fan J, Richardson D, Clarke L, Zerbino D, Cochrane G, et al. FAANG, establishing metadata standards, validation and best practices for the farmed and companion animal community. *Anim Genet.* 2018;49:520–6.
- The FAANG Consortium. The FAANG Data Coordination Center. <https://data.faaang.org>. Accessed 11 Nov 2019.
- The FR-AgENCODE group. FR-AgENCODE: a FAANG pilot project for the annotation of livestock genomes. <http://www.fragencode.org>. Accessed 11 Nov 2019.
- Buenrostro JD, Giresi PG, Zaba LC, Chang HY, Greenleaf WJ. Transposition of native chromatin for fast and sensitive epigenomic profiling of open chromatin, DNA-binding proteins and nucleosome position. *Nat Methods.* 2013;10:1213–8.
- Lieberman-Aiden E, Van Berkum NL, Williams L, Imakaev M, Ragoczy T, Telling A, et al. Comprehensive mapping of long-range interactions reveals folding principles of the human genome. *Science.* 2009;326:289–93.
- Dobin A, Davis CA, Schlesinger F, Drenkow J, Zaleski C, Jha S, et al. STAR: ultrafast universal RNA-seq aligner. *Bioinformatics.* 2013;29:15–21.
- Dobin A, Gingeras TR. Mapping RNA-seq reads with STAR. *Curr Protocol Bioinform.* 2015;51:11–4.



24. Li B, Dewey CN. RSEM: accurate transcript quantification from RNA-Seq data with or without a reference genome. *BMC Bioinformatics*. 2011;12:323.
25. Mank JE. Sex chromosome dosage compensation: definitely not for everyone. *Trends Genet*. 2013;29:677–83.
26. Breschi A, Djebali S, Gillis J, Pervouchine DD, Dobin A, Davis CA, et al. Gene-specific patterns of expression variation across organs and species. *Genome Biol*. 2016;17:151.
27. Lin S, Lin Y, Nery JR, Urich MA, Breschi A, Davis CA, et al. Comparison of the transcriptional landscapes between human and mouse tissues. *Proc Natl Acad Sci*. 2014;111:17224–9.
28. Sudmant PH, Alexis MS, Burge CB. Meta-analysis of RNA-seq expression data across species, tissues and studies. *Genome Biol*. 2016;17:16.
29. Robinson MD, McCarthy DJ, Smyth GK. edgeR: a Bioconductor package for differential expression analysis of digital gene expression data. *Bioinformatics*. 2010;26:139–40.
30. Melé M, Ferreira PG, Reverter F, DeLuca DS, Monlong J, Sammeth M, et al. The human transcriptome across tissues and individuals. *Science*. 2015;348:660–5.
31. Gerner W, Käser T, Saalmüller A. Porcine T lymphocytes and NK cells – an update. *Dev Comp Immunol*. 2009;33:310–20.
32. Guzman E, Hope J, Taylor G, Smith AL, Cubillos-Zapata C, Charleston B. Bovine  $\gamma\delta$  T cells are a major regulatory T cell subset. *J Immunol*. 2014;193:208–22.
33. Kapushesky M, Adamusiak T, Burdett T, Culhane A, Farne A, Filippov A, et al. Gene Expression Atlas update – a value-added database of microarray and sequencing-based functional genomics experiments. *Nucleic Acids Res*. 2011;40:D1077–81.
34. Petryszak R, Burdett T, Fiorelli B, Fonseca NA, Gonzalez-Porta M, Hastings E, et al. Expression Atlas update – a database of gene and transcript expression from microarray-and sequencing-based functional genomics experiments. *Nucleic Acids Res*. 2014;42:D926–32.
35. Djebali S, Davis CA, Merkel A, Dobin A, Lassmann T, Mortazavi A, et al. Landscape of transcription in human cells. *Nature*. 2012;489:101.
36. Wucher V, Legeai F, Hédan B, Rizk G, Lagoutte L, Leeb T, et al. FEELnc: a tool for long non-coding RNA annotation and its application to the dog transcriptome. *Nucleic Acids Res*. 2017;45:e57.
37. Derrien T, Johnson R, Bussotti G, Tanzer A, Djebali S, Tilgner H, et al. The GENCODE v7 catalog of human long noncoding RNAs: analysis of their gene structure, evolution, and expression. *Genome Res*. 2012;22:1775–89.
38. Muret K, Klopp C, Wucher V, Esquerré D, Legeai F, Lecerf F, et al. Long noncoding RNA repertoire in chicken liver and adipose tissue. *Genet Sel Evol*. 2017;49:6.
39. Lagarde J, Usczyńska-Ratajczak B, Santoyo-Lopez J, Gonzalez JM, Tapanari E, Mudge JM, et al. Extension of human lncRNA transcripts by RACE coupled with long-read high-throughput sequencing (RACE-Seq). *Nat Commun*. 2016;7:12339.
40. Hezroni H, Koppstein D, Schwartz MG, Avrutin A, Bartel DP, Ulitsky I. Principles of long noncoding RNA evolution derived from direct comparison of transcriptomes in 17 species. *Cell Rep*. 2015;11:1110–22.
41. Letunic I, Bork P. Interactive tree of life (iTOL) v3: an online tool for the display and annotation of phylogenetic and other trees. *Nucleic Acids Res*. 2016;44:W242–5.
42. Villar D, Berthelot C, Aldridge S, Rayner TF, Lukk M, Pignatelli M, et al. Enhancer evolution across 20 mammalian species. *Cell*. 2015;160:554–66. <https://doi.org/10.1016/j.cell.2015.01.006>.
43. Degner JF, Pai AA, Pique-Regi R, Veyrieras JB, Gaffney DJ, Pickrell JK, et al. DNase I sensitivity QTLs are a major determinant of human expression variation. *Nature*. 2012;482:390.
44. Qu K, Zaba LC, Giresi PG, Li R, Longmire M, Kim YH, et al. Individuality and variation of personal regulomes in primary human T cells. *Cell Syst*. 2015;1:51–61. <https://doi.org/10.1016/j.cels.2015.06.003>.
45. Scott-Brown JP, López-Moyado IF, Trifari S, Wong V, Chavez L, Rao A, et al. Dynamic changes in chromatin accessibility occur in CD8+ T cells responding to viral infection. *Immunity*. 2016;45:1327–40.
46. Thurman RE, Rynes E, Humbert R, Vierstra J, Maurano MT, Haugen E, et al. The accessible chromatin landscape of the human genome. *Nature*. 2012;489:75.
47. Rao S, Huntley M, Durand N, Stamenova E, Bochkov I, Robinson J, et al. A 3D map of the human genome at kilobase resolution reveals principles of chromatin looping. *Cell*. 2014;159:1665–80.
48. Servant N, Lajoie BR, Nora EP, Giorgetti L, Chen CJ, Heard E, et al. HiTC: exploration of high-throughput ‘C’ experiments. *Bioinformatics*. 2012;28:2843–4.
49. Durand NC, Shamim MS, Machol I, Rao SS, Huntley MH, Lander ES, et al. Juicer provides a one-click system for analyzing loop-resolution Hi-C experiments. *Cell Syst*. 2016;3:95–8.
50. Dixon JR, Selvaraj S, Yue F, Kim A, Li Y, Shen Y, et al. Topological domains in mammalian genomes identified by analysis of chromatin interactions. *Nature*. 2012;485:376–80.
51. Gong Y, Lazaris C, Sakellaropoulos T, Lozano A, Kambadur P, Ntziachristos P, et al. Stratification of TAD boundaries reveals preferential insulation of super-enhancers by strong boundaries. *Nat Commun*. 2018;9:542.
52. Crane E, Bian Q, McCord RP, Lajoie BR, Wheeler BS, Ralston EJ, et al. Condensin-driven remodelling of X chromosome topology during dosage compensation. *Nature*. 2015;523:240.
53. Sofueva S, Yaffe E, Chan WC, Georgopoulou D, Rudan MV, Mira-Bontenbal H, et al. Cohesin-mediated interactions organize chromosomal domain architecture. *EMBO J*. 2013;32:3119–29.
54. Rudan MV, Barrington C, Henderson S, Ernst C, Odom DT, Tanay A, et al. Comparative Hi-C reveals that CTCF underlies evolution of chromosomal domain architecture. *Cell Rep*. 2015;10:1297–309.
55. Filippova D, Patro R, Duggal G, Kingsford C. Identification of alternative topological domains in chromatin. *Algorithm Mol Bio*. 2014;9:14.
56. Dixon JR, Gorkin DU, Ren B. Chromatin domains: the unit of chromosome organization. *Mol Cell*. 2016;62:668–80.
57. Lupiáñez DG, Kraft K, Heinrich V, Krawitz P, Brancati F, Klopocki E, et al. Disruptions of topological chromatin domains cause pathogenic rewiring of gene-enhancer interactions. *Cell*. 2015;161:1012–25.
58. Yang Y, Zhang Y, Ren B, Dixon JR, Ma J. Comparing 3D genome organization in multiple species using Phylo-HMRF. *Cell Syst*. 2019. <https://doi.org/10.1101/552505>.
59. Fishman V, Battulin N, Nuriddinov M, Maslova A, Zlotina A, Strunov A, et al. 3D organization of chicken genome demonstrates evolutionary conservation of topologically associated domains and highlights unique architecture of erythrocytes’ chromatin. *Nucleic Acids Res*. 2018;47:648–65.
60. Harmston N, Ing-Simmons E, Tan G, Pery M, Merkschlager M, Lenhard B. Topologically associating domains are ancient features that coincide with Metazoan clusters of extreme noncoding conservation. *Nat Commun*. 2017;8:441.
61. Dixon JR, Jung I, Selvaraj S, Shen Y, Antosiewicz-Bourget JE, Lee AY, et al. Chromatin architecture reorganization during stem cell differentiation. *Nature*. 2015;518:331.
62. Doynova MD, Markworth JF, Cameron-Smith D, Vickers MH, O’Sullivan JM. Linkages between changes in the 3D organization of the genome and transcription during myotube differentiation in vitro. *Skelet Muscle*. 2017;7:5.
63. Schmitt AD, Hu M, Jung I, Xu Z, Qiu Y, Tan CL, et al. A compendium of chromatin contact maps reveals spatially active regions in the human genome. *Cell Rep*. 2016;17:2042–59.
64. Djebali S, Wucher V, Foissac S, Hitte C, Corre E, Derrien T. Bioinformatics pipeline for transcriptome sequencing analysis. In: U Ørom, Enhancer RNAs, volume 1468. New York: Humana Press; 2017. p. 201–219.
65. Trapnell C, Williams BA, Pertea G, Mortazavi A, Kwan G, Van Baren MJ, et al. Transcript assembly and quantification by RNA-Seq reveals unannotated transcripts and isoform switching during cell differentiation. *Nat Biotechnol*. 2010;28:511–5.
66. Roberts A, Pimentel H, Trapnell C, Pachter L. Identification of novel transcripts in annotated genomes using RNA-Seq. *Bioinformatics*. 2011;27:2325–9.
67. Rohart F, Gautier B, Singh A, Lê Cao KA. mixOmics: an R package for omics feature selection and multiple data integration. *PLoS Comput Biol*. 2017;13:005752.
68. Kinsella RJ, Kähäri A, Haider S, Zamora J, Proctor G, Spudich G, et al. Ensembl BioMarts: a hub for data retrieval across taxonomic space. *Database*. 2011;2011:bar030.
69. Robinson MD, Oshlack A. A scaling normalization method for differential expression analysis of RNA-seq data. *Genome Biol*. 2010;11:R25.
70. McCarthy DJ, Chen Y, Smyth GK. Differential expression analysis of multifactor RNA-Seq experiments with respect to biological variation. *Nucleic Acids Res*. 2012;40:4288–97.

71. Benjamini Y, Hochberg Y. Controlling the false discovery rate: a practical and powerful approach to multiple testing. *J R Stat Soc Ser B (Methodol)*. 57:289–300.
72. Falcon S, Gentleman R. Using GOstats to test gene lists for GO term association. *Bioinformatics*. 2006;23:257–8.
73. Breiman L. Random forests. *Mach Learn*. 2001;45:5–32.
74. Langmead B, Salzberg SL. Fast gapped-read alignment with Bowtie 2. *Nat Methods*. 2012;9:357–9.
75. Li H, Handsaker B, Wysoker A, Fennell T, Ruan J, Homer N, et al. The sequence alignment/map format and SAMtools. *Bioinformatics*. 2009;25:2078–9.
76. Feng J, Liu T, Qin B, Zhang Y, Liu XS. Identifying ChIP-seq enrichment using MACS. *Nat Protocol*. 2012;7:1728–40.
77. Quinlan AR. BEDTools: the Swiss-army tool for genome feature analysis. *Curr Protocol Bioinform*. 2014;47:11–2.
78. Ballman KV, Grill DE, Oberg AL, Therneau TM. Faster cyclic loess: normalizing RNA arrays via linear models. *Bioinformatics*. 2004;20:2778–86.
79. Lun AT, Smyth GK. csaw: a Bioconductor package for differential binding analysis of ChIP-seq data using sliding windows. *Nucleic Acids Res*. 2015;44:e45.
80. Grant CE, Bailey TL, Noble WS. FIMO: scanning for occurrences of a given motif. *Bioinformatics*. 2011;27:1017–8.
81. Mathelier A, Fornes O, Arenillas DJ, Chen Cy, Denay G, Lee J, et al. JASPAR 2016: a major expansion and update of the open-access database of transcription factor binding profiles. *Nucleic Acids Res*. 2015;44:D110–5.
82. Servant N, Varoquaux N, Lajoie BR, Viara E, Chen CJ, Vert JP, et al. HiC-Pro: an optimized and flexible pipeline for Hi-C data processing. *Genome Biol*. 2015;16:259.
83. Imakaev M, Fudenberg G, McCord RP, Naumova N, Goloborodko A, Lajoie BR, et al. Iterative correction of Hi-C data reveals hallmarks of chromosome organization. *Nat Methods*. 2012;9:999–1003.
84. Durand NC, Robinson JT, Shamim MS, Machol I, Mesirov JP, Lander ES, et al. Juicebox provides a visualization system for Hi-C contact maps with unlimited zoom. *Cell Syst*. 2016;3:99–101.
85. Barrett T, Clark K, Gevorgyan R, Gorelenkov V, Gribov E, Karsch-Mizrachi I, et al. BioProject and BioSample databases at NCBI: facilitating capture and organization of metadata. *Nucleic Acids Res*. 2012;40(Database issue):D57–63. <https://doi.org/10.1093/nar/gkr1163>.
86. EMBL-EBI. BioSamples. <https://www.ebi.ac.uk/biosamples>. Accessed 13 Nov 2019.

## Publisher's Note

Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.

Ready to submit your research? Choose BMC and benefit from:

- fast, convenient online submission
- thorough peer review by experienced researchers in your field
- rapid publication on acceptance
- support for research data, including large and complex data types
- gold Open Access which fosters wider collaboration and increased citations
- maximum visibility for your research: over 100M website views per year

At BMC, research is always in progress.

Learn more [biomedcentral.com/submissions](https://biomedcentral.com/submissions)

