

RESEARCH ARTICLE

Open Access



Genome variation and conserved regulation identify genomic regions responsible for strain specific phenotypes in rat

David Martín-Gálvez¹, Denis Dunoyer de Segonzac¹, Man Chun John Ma^{2,3,4}, Anne E. Kwitek^{2,3}, David Thybert^{1,5*} and Paul Flicek^{1*} 

Abstract

Background: The genomes of laboratory rat strains are characterised by a mosaic haplotype structure caused by their unique breeding history. These mosaic haplotypes have been recently mapped by extensive sequencing of key strains. Comparison of genomic variation between two closely related rat strains with different phenotypes has been proposed as an effective strategy for the discovery of candidate strain-specific regions involved in phenotypic differences. We developed a method to prioritise strain-specific haplotypes by integrating genomic variation and genomic regulatory data predicted to be involved in specific phenotypes. Specifically, we aimed to identify genomic regions associated with Metabolic Syndrome (MetS), a disorder of energy utilization and storage affecting several organ systems.

Results: We compared two Lyon rat strains, Lyon Hypertensive (LH) which is susceptible to MetS, and Lyon Low pressure (LL), which is susceptible to obesity as an intermediate MetS phenotype, with a third strain (Lyon Normotensive, LN) that is resistant to both MetS and obesity. Applying a novel metric, we ranked the identified strain-specific haplotypes using evolutionary conservation of the occupancy three liver-specific transcription factors (HNF4A, CEBPA, and FOXA1) in five rodents including rat. Consideration of regulatory information effectively identified regions with liver-associated genes and rat orthologues of human GWAS variants related to obesity and metabolic traits. We attempted to find possible causative variants and compared them with the candidate genes proposed by previous studies. In strain-specific regions with conserved regulation, we found a significant enrichment for published evidence to obesity—one of the metabolic symptoms shown by the Lyon strains—amongst the genes assigned to promoters with strain-specific variation.

Conclusions: Our results show that the use of functional regulatory conservation is a potentially effective approach to select strain-specific genomic regions associated with phenotypic differences among Lyon rats and could be extended to other systems.

Keywords: Metabolic syndrome, Genome regulation, Evolution

* Correspondence: david.thybert@earlham.ac.uk; flicek@ebi.ac.uk

¹European Molecular Biology Laboratory, European Bioinformatics Institute, Wellcome Genome Campus, Hinxton, Cambridge CB10 1SD, UK
Full list of author information is available at the end of the article



Background

Phenotypic diversity is ultimately driven by genetic differences. The connections between DNA sequence and observed phenotypes are often difficult to determine and may be confounded by non-genetic causes including environmental effects. Regardless, it is increasingly clear that differences in transcriptional regulation are an important factor explaining phenotypic diversity [1–6]. This is especially true between closely-related species [2, 7–9]. Accordingly, a number of efforts have been made to combine transcriptional regulatory data with genome variation to select candidate genomic regions involved in producing phenotypic characteristics of interest [10–13].

The rat is a key animal model for biomedical research [14–16]. More than 600 laboratory rat strains have been created over the last century in order to study specific traits including those which are more informative in rat than in other model species, such as behaviour and neurodegenerative diseases, cardiovascular diseases and metabolic disorders [17–19]. One focus over the last decade has been the identification of genes and other genomic loci associated with these strain-specific traits [13, 20]. Despite the great number of quantitative trait loci (QTL) identified in rat models using a number of techniques [21], only a small number of causative genes have been determined for complex traits or diseases [13, 22, 23].

Most genomic variants in an individual are expected to be neutral, and therefore have no impact on reproduction or survival [24, 25]. In the case of laboratory rats, the existing variation among strains (e.g. [26]) is the sum of the ancestral variation among individuals used in the process of strain development and the novel variation that originated and accumulated in the genome during the establishment and maintenance of the strains. Like humans [27] and laboratory mice [28], genetic variation among rat strains is not randomly distributed across the genome; instead it is organised in haplotype blocks [29–31], which are caused by meiotic crossover of the shared ancestral variation. Comparison of these haplotype blocks among rat strains with different phenotypes has proven to be a powerful strategy for genetic mapping of complex traits and diseases [29, 31]. For example, Atanur and colleagues analysed the genomes of 27 rat strains, and found that haplotype blocks with variants that are unique to a single strain were positively selected in the initial phenotype-driven derivation of these strains, and thus variants associated with strain-specific phenotypes are predicted to be in these regions [30]. However, the genomic extent of such regions and the number of sequence variants found within them are nearly always too large for an effective determination of candidate loci influencing the phenotype of interest (see e.g. [32]).

Regulatory activity such as active promoters, enhancers and transcription factor binding sites (TFBS) can be effectively mapped genome-wide with current techniques such as

chromatin immunoprecipitation followed by high-throughput sequencing (ChIP-seq) [33]. Previous studies have suggested that both the number and conservation level of transcription factor binding sites in a given region affect the level of gene expression [3, 34–37]. Since tissue characteristics are directed to a large extent by the activity of tissue-specific transcription factors, the location of these regulatory elements might be useful when selecting haplotype blocks associated with specific phenotypes or diseases.

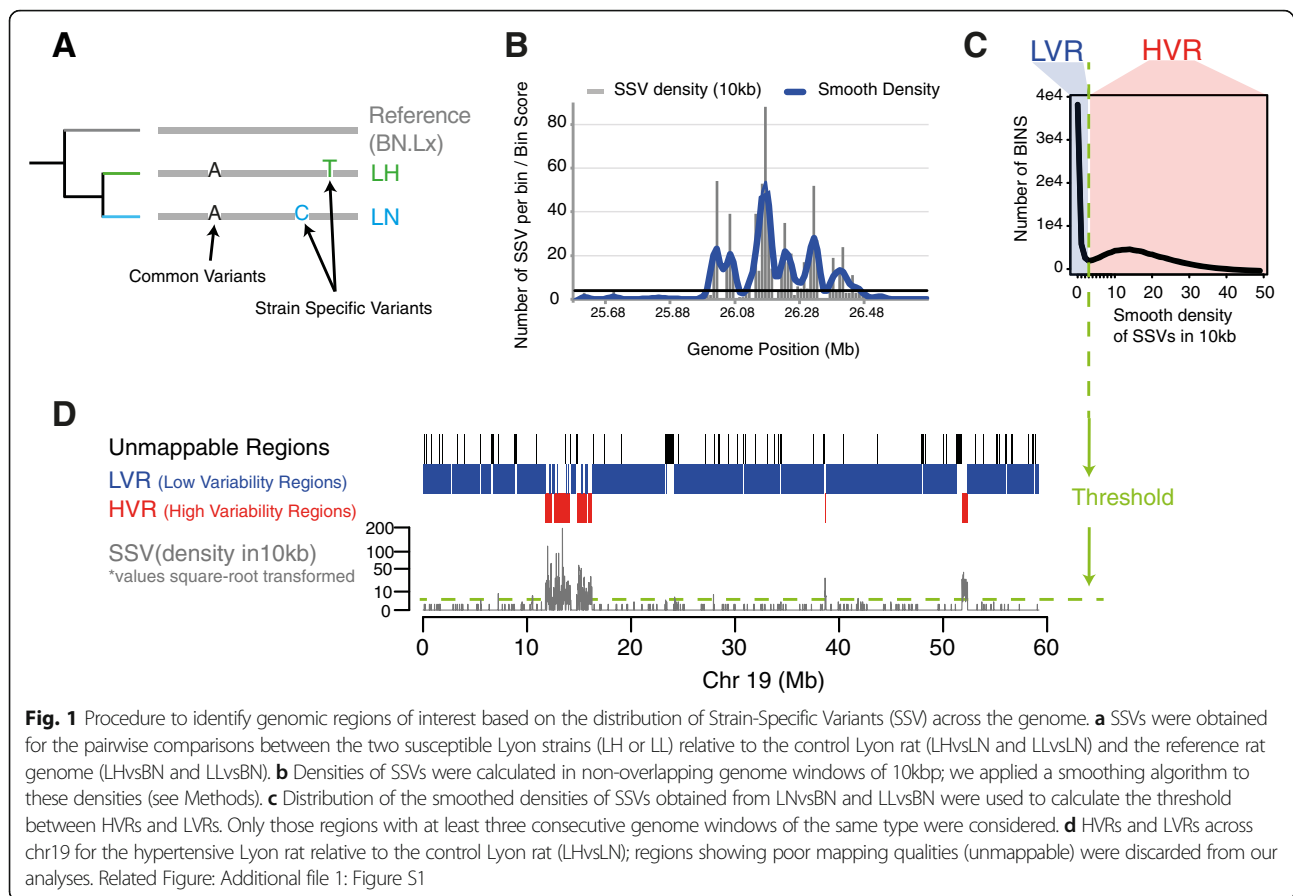
In this study, we characterise the haplotype blocks holding strain-specific genome variation among three closely related strains of the Lyon rat. Although Lyon rats were initially established as a model of hypertension [38], several additional symptoms related to metabolic syndrome (MetS), such as obesity, dyslipidaemia and susceptibility to insulin resistance have been found in the Lyon Hypertensive (LH/Mav) strain [39–41]. Only obesity is observed in the Lyon Low pressure (LL/Mav) and all MetS related phenotypes are absent in the Lyon Normotensive (LN/Mav) strain [39–42]. Since both liver and kidney are involved in MetS [43], we generated RNA-seq expression data from liver of LL rats and from kidney of all three strains and integrated these with relevant existing data including the level of regulatory conservation for three liver-specific transcription factors (CEBPA, FOXA1 and HNF4A [9]) between rat and five related mouse species and strains. We show that the level of functional regulatory conservation can help select strain-specific haplotype blocks putatively associated with phenotypic differences among Lyon rats.

Results

85% of strain-specific variation among Lyon rat strains is concentrated in less than 9% of the genome

To define haplotype blocks, we partitioned the rat genome into 10 kb windows and calculated the number of strain-specific variants (SSVs) in each window relative to the reference rat genome assembly (see Methods). We observed a bimodal distribution in the number of SSVs and used this distribution to define the resulting haplotype blocks as having either a high density of SSVs (High Variability Region, HVR) or a low density of SSVs (Low Variability Region, LVR) (see Methods, Fig. 1 and Additional file 1: Figure S1).

The distribution of SSVs across the genome was similar for the two pairwise comparisons of Lyon rats susceptible to MetS and obesity (LH and LL) and the Lyon Normotensive (LN) that is resistant (i.e. LHvsLN and LLvsLN, see Fig. 2a, Additional file 1: Figures S2A and S3A). In both cases, the vast majority of strain-specific variants were concentrated in HVRs (LHvsLN: 84.96% and LLvsLN: 85.09%), and these regions only covered a small part of the genome (LHvsLN: 8.55% and LLvsLN: 7.10%) (Fig. 2b, Additional file 1: Figures S2B and S3B, Table 1). These regions were partly



overlapping: 42.6% of LHvsLN HVRs overlap with LLvsLN HVRs, while 51.2% of LLvsLN HVRs overlap with LHvsLN. SSV overlaps have similar fractions (Fig. 2c). The fraction of the genome that we identify as highly strain-specific is similar to that obtained previously for these and other rat strains (see [30, 31], Methods and Additional file 1: Figure S2B). HVRs characterise a substantial reduction in the portion of the genome that is most likely to be involved in MetS phenotypes and therefore form the primary focus of our subsequent analysis.

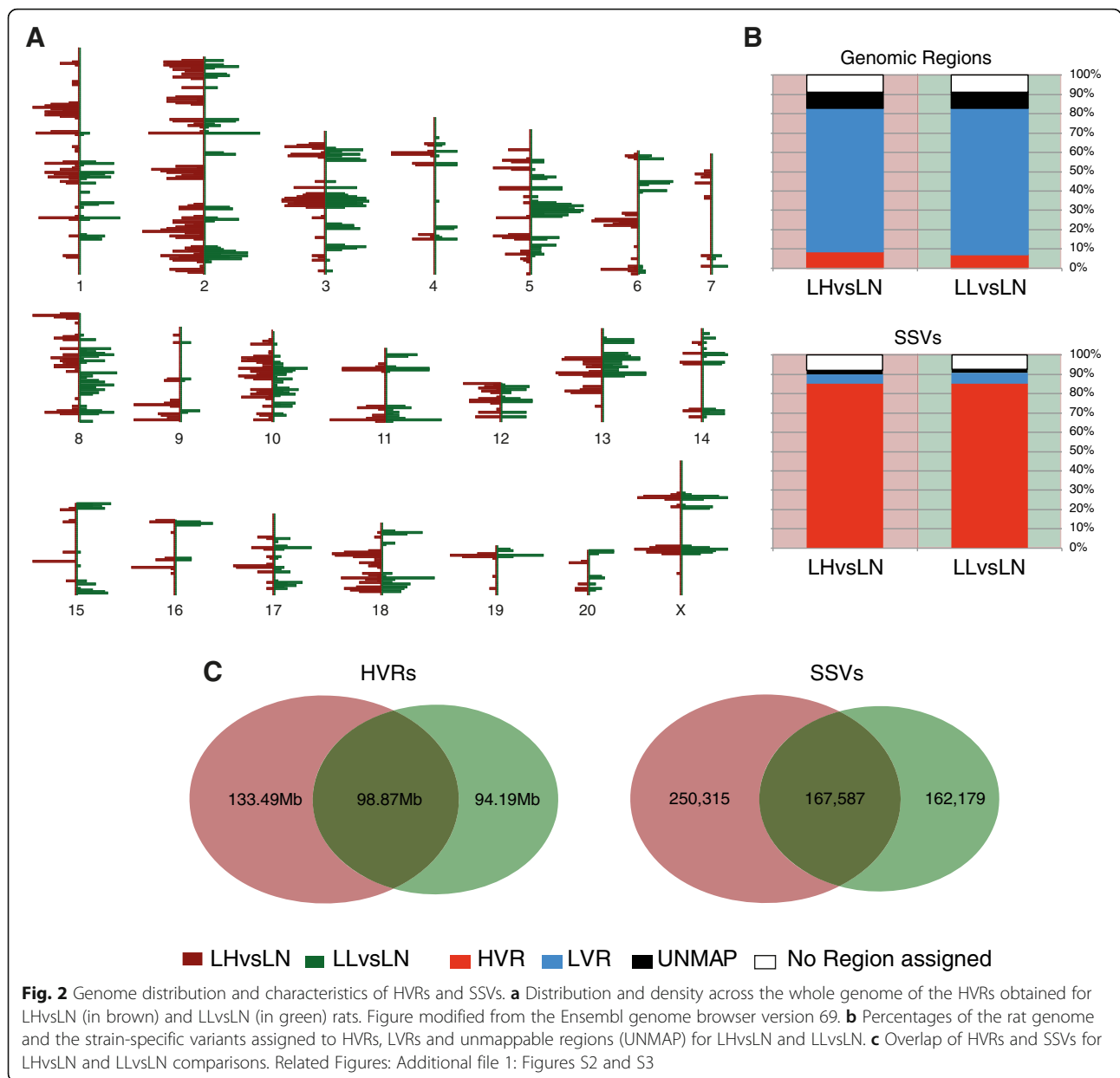
Evidence for the functionality of high variability regions in the Lyon rats associated with MetS

We then sought to determine if the HVRs preferentially contain features that could explain the phenotypic differences among Lyon rats by comparing them to other regions of the genome (see Methods and Fig. 3). Specifically, we tested whether there is a significant enrichment in HVRs of the following elements: i) annotated genes, ii) genes associated with metabolic-related traits, iii) genes differently expressed among Lyon rats, iv) occupancy in rat of three liver-specific transcription factors, and v) regions orthologous to human variants associated by GWAS to obesity and metabolic traits.

The number of Ensembl genes [44] that overlap at least one HVR was marginally greater than expected by chance and this overlap was significant for the LHvsLN comparison ($p < 0.05$), but not for LLvsLN ($p > 0.05$) (Fig. 3b). Additionally, but only in the case of LLvsLN, there was a significant enrichment of genes associated with Type 1 diabetes mellitus (KEGG PATHWAY database, gene count: 25, $p < 10^{-9}$, see results for 'All HVRs' in Additional file 1: Figures S4 and S5 (DAVID web services v6.7 [45, 46]). Gene enrichment in the HVRs was more significant when considering only the genes whose expression in either liver or kidney differed between LH and LN strains, and between LL and LN strains (see Methods and Additional file 1: Table S1). (Fig. 3c and Additional file 1: Table S2).

We next considered whether the HVRs were enriched for either the occupancy of three specific transcription factors (HNF4A, CEBPA, and FOXA1) or the 418 rat orthologues of Human GWAS variants associated with obesity and metabolic traits. In both cases we did not observe a significant enrichment (Fig. 3d–e, Additional file 1: Figure S6 and Table S3).

In summary, the observation that both annotated and differentially expressed genes are enriched in HVRs supports the hypothesis that HVRs harbour functional



regions that could be responsible for phenotypic differences observed among the Lyon rats. However, given the overall genomic span of identified HVRs and the large number of SSVs both in coding and non-coding regions in the HVRs (see Table 1), these analyses on their own are inadequate to suggest either the causative genes or

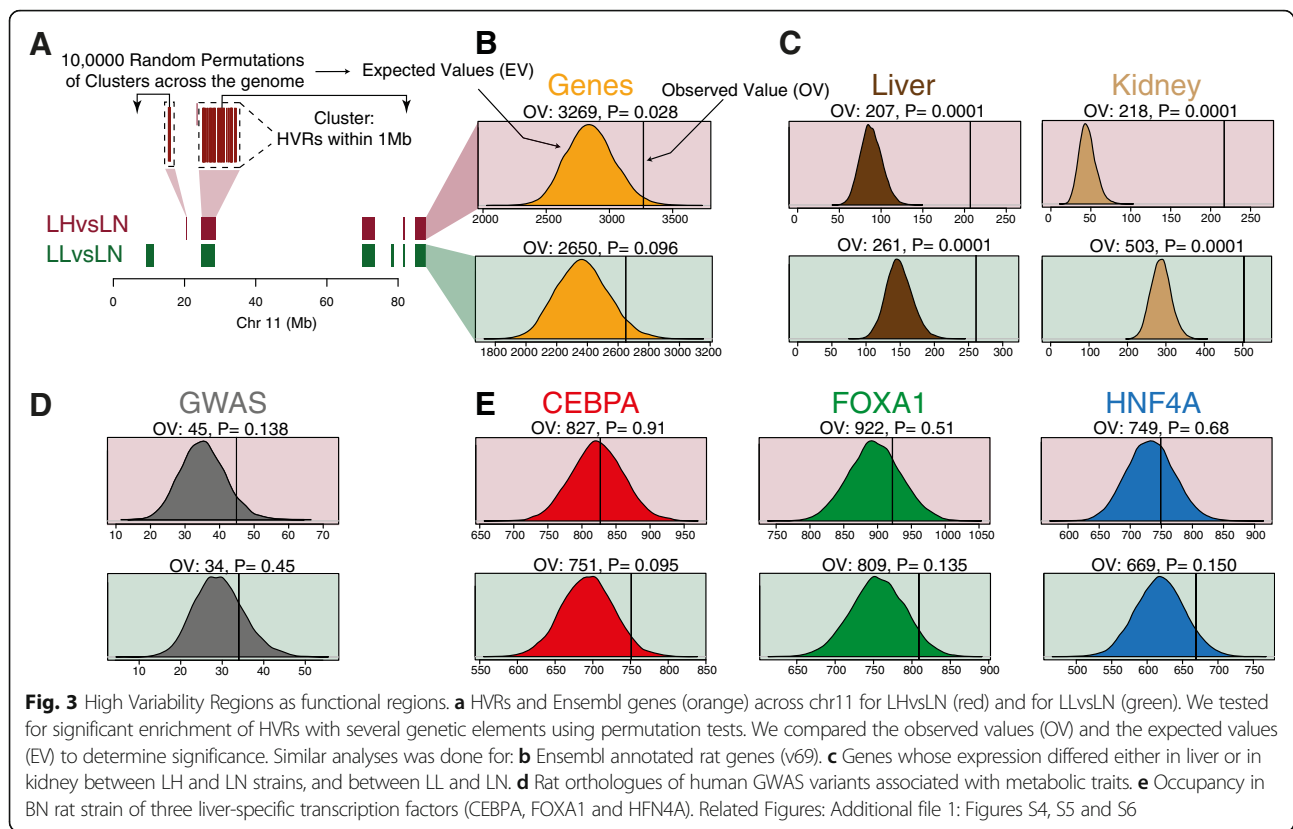
the causative variants influencing MetS or obesity across the whole genome.

Liver-specific regulation data can prioritise regions of strain-specific variation in Lyon rats associated with MetS

We next integrated strain-specific genomic variation with available genomic regulatory data from tissue relevant to MetS in order to prioritise the HVRs using the occupancy and level of conservation of the three liver-specific transcription factors. We created subsets of HVRs characterised by occupancy of the factors and the level of conservation among mice and rats using a factor-specific Conservation Enrichment score (CE_f , see Methods). Briefly, CE_f is the fraction of transcription factor binding events in a 10 kb

Table 1 Descriptive statistics for the SSVs and HVRs obtained for the susceptible Lyon rats (LH and LL) relative to the resistant Lyon strain (LN)

Strains	Total SSVs	SSVs in HVRs	HVR count	Total HVRs size
LHvsLN	413,068	351,459 (85.09%)	2319	232.36 Mb (8.55%)
LLvsLN	324,932	276,053 (84.96%)	1941	193.06 Mb (7.10%)



window that are conserved between rat and mouse for each transcription factor. The score was determined independently for each of the three factors ($f = \text{CEBPA}$, FOXA1 , or HNF4A). Thus, for each transcription factor and for both the LHvsLN and LLvsLN comparisons, we created seven subsets of HVRs: 'All HVRs' including those without any binding event; 'HVR w/TFBS' with at least one TFBS regardless of conservation; and five subsets containing HVRs with CE_f greater than 0, 0.2, 0.4, 0.6 and 0.8, respectively. We then reassessed the evidence for functionality of these HVR subsets in a similar way to that done with the whole set of HVRs as above.

Enrichment of Ensembl rat genes in HVR subsets corresponded with the occupancy and level of conservation of the three liver-specific transcription factors (Fig. 4a). With the exception of the subset of LLvsLN with all HVRs, all tests in the HVR subsets were statistically significant ($p < 0.05$). For LHvsLN and for the three factors, maximum significance possible ($p < 10^{-3}$) was obtained for the subset of HVRs with at least one TFBS (HVR w/TFBS), and for the subsets with $\text{CE}_f > 0.0$ and $\text{CE}_f > 0.2$. In the case of LLvsLN and for the three transcription factors, the maximum significance was obtained in the subset of HVRs with at least one TFBS, and HVR subset with $\text{CE}_f > 0.0$ (i.e. HVR subsets with the darkest colour in Fig. 4a).

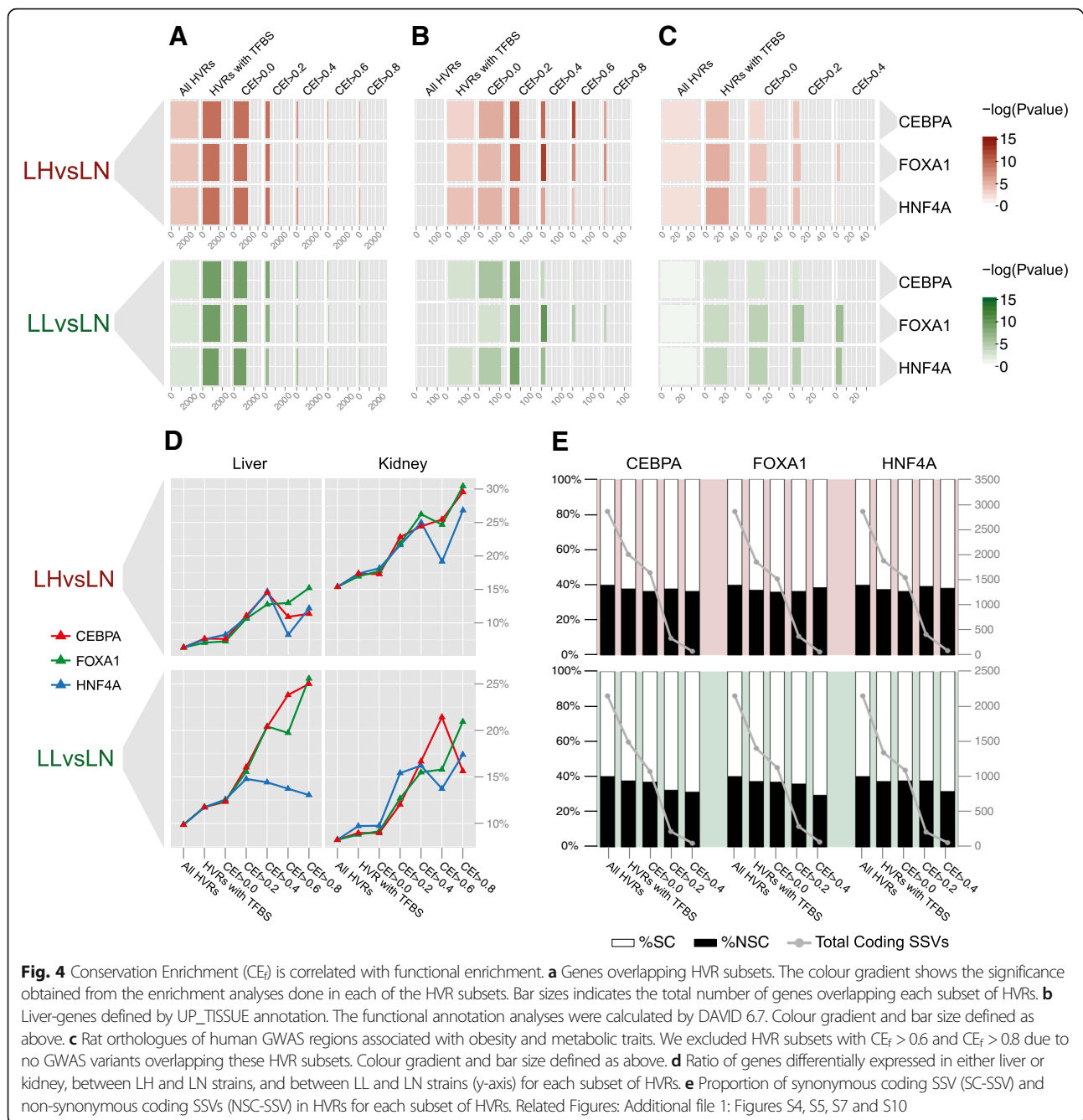
As above, we analysed the functional annotation enrichment in the HVR subsets using DAVID [45, 46]

with KEGG PATHWAY [47, 48] and UP TISSUE [49] databases (see Methods). In all cases (both LHvsLN and LLvsLN for the three liver-specific transcription factors) the term 'liver' from the UP TISSUE database, had the greatest accumulated significance across HVR subsets (see Additional file 1: Figures S4A and S5A). These results indicate that HVRs selected according to information from liver-specific regulation data are enriched in genes associated with liver function. Importantly, but as expected, this association was not evident without using genomic regulation data to select HVRs (see Fig. 4b).

In the case of LHvsLN, the greatest enrichments in genes associated with liver function were obtained for the $\text{CE}_{\text{CEBPA}} > 0.6$ ($p < 10^{-5}$), $\text{CE}_{\text{FOXA1}} > 0.4$ ($p < 10^{-5}$) and $\text{CE}_{\text{HNF4A}} > 0.2$ ($p < 10^{-3}$) subsets (HVR subsets with the darkest colour in Fig. 4b). For LLvsLN, the greatest enrichments were obtained for $\text{CE}_{\text{CEBPA}} > 0.2$ ($p < 10^{-3}$), $\text{CE}_{\text{FOXA1}} > 0.4$ ($p < 10^{-4}$) and $\text{CE}_{\text{HNF4A}} > 0.2$ ($p < 10^{-3}$) (Fig. 4b).

For the analyses using KEGG_PATHWAY database, we did not find a consistent increase in significance associated with an increase in CE_f , although we did identify some functional terms that were statistically significant (see Additional file 1: Figures S4B and S5B).

Additionally, we looked for an enrichment of putative GWAS positive regions in HVR subsets by determining the orthologous location in rat of NHGRI-EBI GWAS Catalog SNPs associated with obesity and metabolic-related in



humans [50] (specific terms listed in Additional file 1: Table S3). The use of the regulatory information from liver-specific transcription factors identified significant enrichments of GWAS variants in relevant subsets of HVRs. For example, we found significant enrichments (i.e. $p < 0.05$) for LHvsLN in the subsets of HVRs w/ TFBS and $CE_f > 0.2$ for the three liver-specific factors, in $CE_f > 0.0$ for FOXA1 and HNF4A and in $CE_f > 0.4$ for FOXA1 (Fig. 4c, Additional file 1: Table S4). Regarding LLvsLN, we found significant enrichments for in HVRs w/TFBS, $CE_f > 0.0$, $CE_f > 0.2$, $CE_f > 0.4$ for both FOXA1 and HNF4A (Fig. 4c, Additional file 1: Table S4).

Genes differentially expressed between LH and LN strains in either liver or kidney (Additional file 1: Table S1) are significantly enriched for all subsets of HVRs (Additional file 1: Figure S7), regardless of which transcription factor is considered. The same is true for genes differentially expressed between the LL and LN strains. To compare the differences in enrichment between subsets, we computed the fraction of all Ensembl rat genes that are differentially expressed for each HVR. In all cases, the fraction of differentially expressed genes was positively correlated with CE_f (Fig. 4d). For example, for LHvsLN

and data from liver, the fraction increased from 6% ('all HVRs' subset) to 15% ($CE_{FOXA1} > 0.8$), while for kidney, it increased from 15% ('all HVRs' subset) to 30% ($CE_{CEBPA} > 0.8$; $CE_{FOXA1} > 0.8$). A similar pattern was observed for LLvsLN (Fig. 4d).

We hypothesised that there may be a correlation between selection pressures leading to regulatory conservation as measured by CE_f and changes to the sequence of protein coding genes within the same sets of HVRs. The ratio of non-synonymous coding SSVs (NSC-SSVs) to synonymous coding SSVs (SC-SSVs) was therefore compared across HVR subsets (see Methods). Although we find relatively little difference in the ratio of the non-synonymous changes, especially for the case of the LHvsLN comparison, in the LLvsLN comparison, non-synonymous changes do appear to be depleted when HVRs have higher regulatory conservation (i.e. higher CE_f) (Fig. 4e). This may be the effect of simultaneous selection on both protein coding genes and regulatory networks for a subset of regions in the LL genome.

In summary, the use of CE_f (i.e. the conservation level between rat and mouse in the occupancy of three liver-specific transcription factors) is effective for selecting candidate regions involved in phenotypic differences between Lyon rats. In most cases, we observed an increase in statistical significance as a function of CE_f . Moreover, the consideration of regulatory information was required to identify HVRs significantly enriched for genes associated with metabolic related-trait genes and enriched for human GWAS variants related to obesity and metabolic traits.

Integrating results from the three liver-specific transcription factors to prioritise the strain-specific variation in Lyon rats associated with MetS

Given the observed stability of combinatorially bound transcription factors [9] and connection of these regions to human disease [51], we assessed if the number of liver-specific transcription factors used to estimate the conservation level could more efficiently prioritise candidate HVRs. For this purpose, we used the HVR subsets with $CE_f > 0$ (i.e. all HVRs with at least one conserved TFBS).

We observed that conservation of more than one type of factor in a given HVR was common: 41% (LHvsLN) and 43% (LLvsLN) of the HVR $CE_f > 0$ regions had conserved peaks for all three of the liver-specific transcription factors (Fig. 5a). We then partitioned the HVRs with conserved peaks by the diversity of factors that were conserved in the given HVR. Specifically, 'HVR 1TF' includes HVRs with one or more conserved TFBS from at least one factor; while 'HVR 2TF' and 'HVR 3TF' refer to HVRs with conserved TFBS from at least two or all three factors (i.e. HVR 3TF is a strict subset of HVR 2TF, which is strict subset of HVR 1TF). We then assessed these HVRs subsets to determine if an

increased diversity of conserved liver-specific transcription factors is an effective method to prioritise HVRs.

Although the presence of genes is significantly enriched in each of these HVR subsets, there are no differences among HVR 1TF, HVR 2TF and HVR 3TF: in all cases enrichment significances were equal to $p = 10^{-4}$ (Fig. 5b). Similar results were obtained when considering all liver-associated genes in HVRs (see Methods) for LLvsLN (HVR TF1: gene count = 126, $p < 0.05$; HVR TF2: gene count = 105, $p < 0.05$; HVR TF3: gene count = 87, $p < 10^{-2}$) and LHvsLN (HVR TF1: gene count = 153, $p < 0.05$; HVR TF2: gene count = 131, $p < 10^{-2}$; HVR TF3: gene count = 102, $p < 10^{-2}$) (Fig. 5c).

HVRs with an increased diversity of conserved peaks were generally significantly enriched (permutation tests, $p < 0.05$) for orthologous regions of human GWAS SNPs except for the case of the HVR 3TF subset with the LHvsLN SSVs (Fig. 5d and Additional file 1: Table S5).

The significance of enrichments of genes differentially expressed in either liver or kidney was $p < 10^{-3}$ for the subsets with at least one conserved peak for one, two and three liver-specific transcription factors, respectively, for both LHvsLN and LLvsLN strain comparisons. Considering the ratios of differentially expressed genes, we observed that they kept relatively constant across HVR subsets as the number of liver-specific transcription factors with conserved peaks increased (Fig. 5e).

These results suggest that knowledge of which TFBSs are conserved and whether a given region of the genome has conserved TFBSs from multiple factors may be effective in some situations at prioritising regions with strain specific variation involved with tissue specific functions. For example, we did not observe enrichments in HVRs associated with liver-expressed genes and the orthologous rat regions associated with human GWAS without using the conservation level as measured by CE_f (see above and Fig. 4).

Analysing the genes obtained from the selected high variability regions

To gain insight into genes from the prioritised HVRs that may be important for MetS or obesity, we focused on the most conserved and regulatorily complex HVR subset, i.e. the set containing at least one conserved peak for all three liver-specific factors (the HVR 3TF subset, see above).

We performed two analyses based on possible functional mechanisms underpinning the phenotypic differences among Lyon rats. First, we characterised those genes with non-synonymous coding strain-specific variants (NSC-SSVs; see Methods) overlapping the selected HVRs. Such variation would result in changes to the amino acid sequence that may be responsible for functional changes in the resulting proteins. Second, we

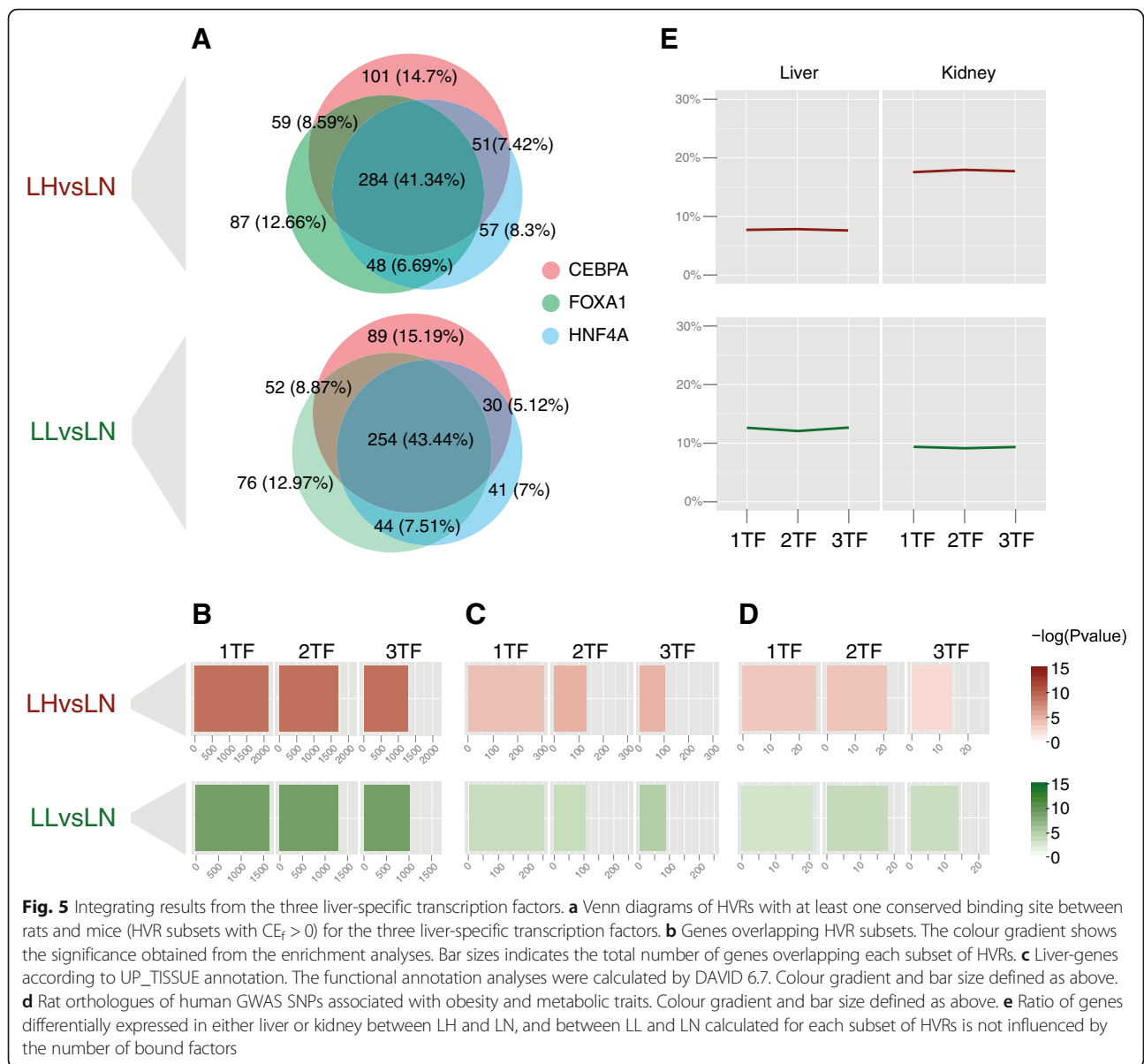


Fig. 5 Integrating results from the three liver-specific transcription factors. **a** Venn diagrams of HVRs with at least one conserved binding site between rats and mice (HVR subsets with $CE_T > 0$) for the three liver-specific transcription factors. **b** Genes overlapping HVR subsets. The colour gradient shows the significance obtained from the enrichment analyses. Bar sizes indicates the total number of genes overlapping each subset of HVRs. **c** Liver-genes according to UP_TISSUE annotation. The functional annotation analyses were calculated by DAVID 6.7. Colour gradient and bar size defined as above. **d** Rat orthologues of human GWAS SNPs associated with obesity and metabolic traits. Colour gradient and bar size defined as above. **e** Ratio of genes differentially expressed in either liver or kidney between LH and LN, and between LL and LN calculated for each subset of HVRs is not influenced by the number of bound factors

identified those genes located near putative promoters in rat obtained from Villar et al. [4] (see Methods) and with SSVs overlapping the selected HVRs. We assumed these SSVs might affect the expression of the proximal genes. For these analyses, we used only those genes expressed in liver as measured by RNA-seq data (Fragments per Kilobase Million (FPKM) > 1, see Methods).

We categorised the selected genes based on whether they were i) liver-specific genes (according to the UniProt tissue database, see Methods), ii) differentially expressed in liver and/or kidney when comparing the susceptible Lyon strains with the control Lyon strain, (see Methods); iii) overlapping human GWAS variants associated with obesity and metabolic traits overlapping the gene body in the case of genes with NSC-SSVs or overlapping the

promoter in the cases of genes linked to promoters with SSVs (see Table 2 and Additional files 2 and 3).

We also analysed the genes associated by published evidence to three symptoms showed by the LH strain (insulin resistance, dyslipidaemias) and by the LH and LL strains (obesity) plus two symptoms not obviously present in these strains as control (heart disease and Alzheimers), using DisGeNET (v4.0 [52, 53]) (see Methods). For this analysis, we used the corresponding human orthologues of the selected rat genes because the DisGeNET data is mainly for human (see Methods). A total of 7542 and 7520 rat genes had human orthologues and were expressed in liver in the LH and LL strains, respectively. DisGeNET identifies a small number of these genes as associated with the metabolic syndrome

Table 2 Number of rat genes and human orthologues expressed in liver of the susceptible Lyon rats (LH or LL) and associated with coding or promoter strain-specific variation overlapping HVRs with at least one conserved peak for three liver-specific transcription factors (i.e. the HVR 3TF subset)

	LHvsLN		LLvsLN	
	Rat genes	Human orthologues	Rat genes	Human orthologues
Coding Variation				
Genes with NSC-SSVs	133	111 (1, 1, 14)	96	78 (0, 0, 8)
Genes with GWAS ^a and NSC-SSVs	3	2 (0, 0, 2)	3	3 (0, 0, 2)
Liver-genes ^b with NSC-SSVs	15	9 (0, 1, 4)	10	6 (0, 0, 1)
Dif-liver ^c genes with NSC-SSVs	14	12 (0, 0, 0)	17	16 (0, 0, 1)
Dif-kidney ^c genes with NSC-SSVs	32	25 (0, 0, 4)	9	7 (0, 0, 0)
Promoters				
Genes assigned to promoters with SSVs	232	206 (7, 3, 32)	189	164 (4, 4, 18)
Genes assigned to promoters with GWAS ^a and SSVs	1	1 (0, 0, 1)	1	1 (0, 0, 1)
Liver-genes ^b assigned to promoters with SSVs	35	30 (2, 1, 9)	26	26 (1, 1, 4)
Dif-liver ^c assigned to promoters with SSVs	38	33 (1, 1, 8)	26	26 (1, 1, 3)
Dif-kidney ^c assigned to promoters with SSVs	52	44 (0, 1, 7)	13	13 (0, 1, 1)

^aHuman GWAS variants associated with Obesity and Metabolic traits in the NHGRI-EBI GWAS catalogue

^bGenes expressed specifically in liver according to UniProt tissue database and accessed by using DAVID web services

^cGenes differently expressed in liver or in kidney when comparing the susceptible Lyon rats (LH or LL) with the control Lyon rat (LH)

Numbers of human orthologues having evidence of associations with insulin resistance, dyslipidaemias and obesity are shown between parentheses

phenotypes and, as expected, these gene sets are highly similar for the LH and LL strains with approximately 140 (1.9%), 110 (1.5%) and 800 (10.6%) genes associated with insulin resistance, dyslipidaemia and obesity, respectively in each strain.

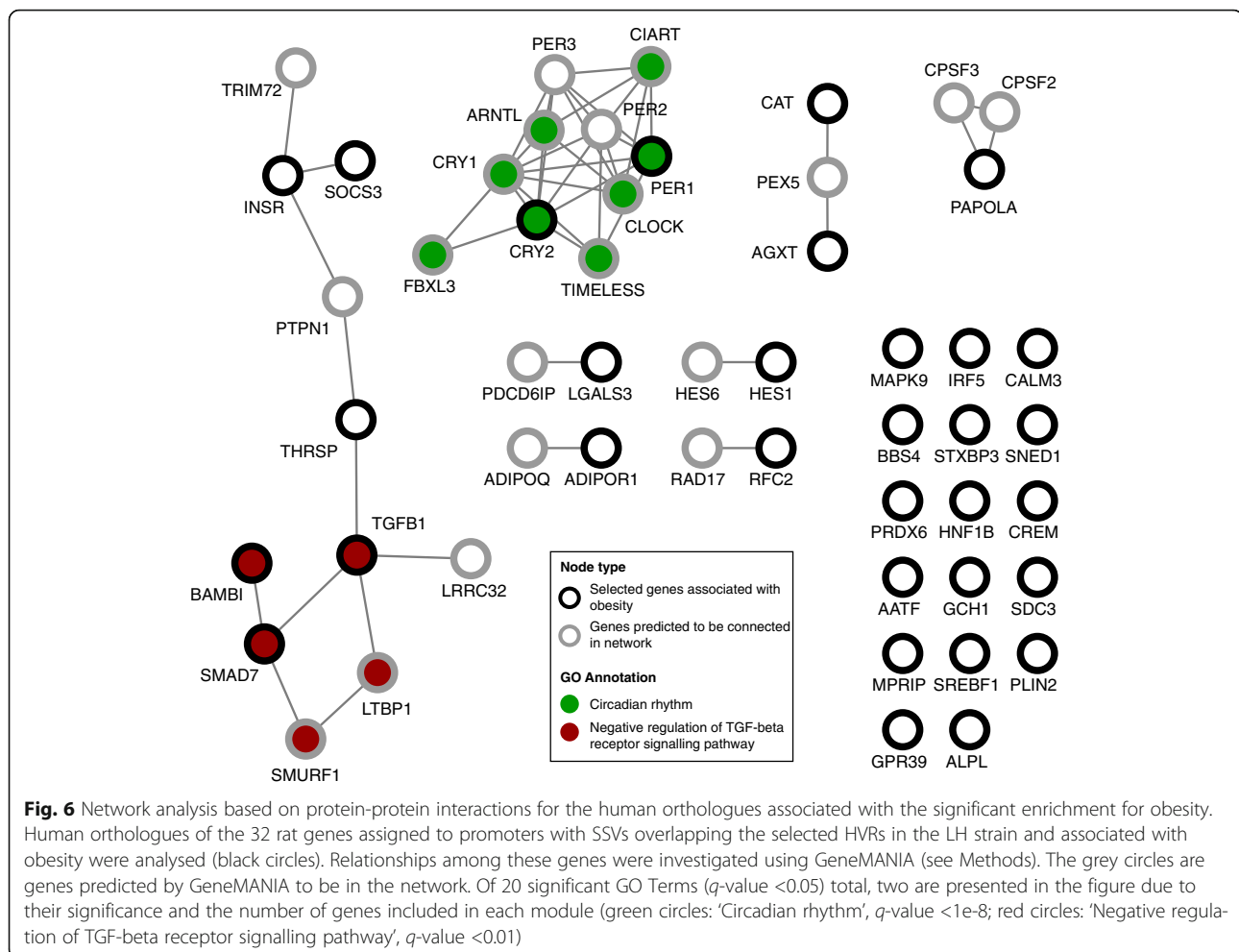
There were 173 protein-coding genes expressed in liver of LH or LL strains with at least one NSC-SSV in the selected HVRs. Of these 173 genes, 144 had identified human orthologues and can thus be compared with the DisGeNET data. This set of genes was not significantly enriched for published associations to obesity, insulin resistance or dyslipidaemias, (see Table 2, Additional file 2).

A larger number of genes were one-to-one associated with putative active promoters including 3865 and 3864 genes that were expressed in livers of LH and LL strain rats and had human orthologues, respectively. Of the 3865 genes with active promoters from the LH strain, a total of 85 (2.2%), 66 (1.7%) and 425 (11%) were associated with insulin resistance, dyslipidaemias and obesity, respectively. The numbers for LL are similar: 86 (2.2%), 65 (1.7%), 417 (10.8%) for the associations to insulin resistance, dyslipidaemias and obesity, respectively. Only a fraction of these promoters had SSVs in the selected HVRs: 206/3865 (5.3%) for LH and 164/3864 (4.2%) for LL. Thirty-two of the genes assigned to promoters with SSVs overlapping the selected HVRs in the LH strain were associated with obesity (15.5%, Fisher's exact test: $p < 0.05$, see Table 2, Additional file 1: Table S6 and Additional file 3). There were no significant enrichments for insulin resistance or dyslipidaemias in the LH strain or

any significant associations in the LL comparison (see Table 2 and Additional file 3).

We found no significant enrichments for the two symptoms used as control in either the comparison to genes with at least NSC-SSV in the selected HVRs (Fisher's exact tests: all p -values >0.08) or to genes one-to-one assigned to promoters with SSVs overlapping the selected HVRs (Fisher's exact tests: all p -values >0.05).

Of the set of 32 genes responsible for the significant enrichment for obesity in the LHvsLN comparison (Additional file 1: Figure S8), the gene with most published evidence of association with obesity was the insulin receptor gene *Insr* (ENSRNOG00000029986) (Additional file 1: Table S6); *Cat* (ENSG00000121691) was the human gene of that list assigned to the promoter with the greatest number of SSVs (58 SSVs) overlapping the HVRs. To assess interactions among the identified set of 32 genes and gain further insight into how they might connect to the Lyon rat phenotypes, we conducted a network analysis and a functional enrichment analyses using GeneMANIA (plugging for Cytospace v3.4.1 [54]). We observed two major modules (Fig. 6 and Additional file 1: Table S7) and significant enrichments for 20 GO terms (Additional file 1: Table S7). One module includes 10 genes interconnected with PER1 and CRY2 and which is responsible for the significant enrichment of GO terms related with circadian rhythm and regulatory binding regions in DNA (Fig. 6 and Additional file 1: Table S7). The other module includes 11 genes interconnected with INSR, SOCS3, THRSP, TGFBI, BAMBI and SMAD7. Genes in this module are responsible



for the significant enrichment of GO terms related to regulation of the transforming growth factor beta receptor (TGF-beta) signalling pathway and regulation of the transmembrane receptor protein serine/threonine kinase signalling pathway (Fig. 6 and Additional file 1: Table S7).

Discussion

In this study, we have used the level of functional regulatory conservation between related species to prioritise genomic regions whose patterns of genome variation suggest that they are involved in phenotypic differences in a model of obesity and metabolic syndrome, the Lyon rat strains.

As a first step, we characterised haplotype blocks by density of strain-specific variants for the two comparisons between the susceptible Lyon strains with respect to the resistant Lyon strain (i.e. LHvsLN and LLvsLN). In agreement with similar analyses [30, 31], most of these variants were concentrated in a small part of the genomes, which we termed High Variability Regions (HVRs). Next, we classified the HVRs according to conserved occupancy between rat and mice for three

liver-specific transcription factors. Functional enrichment of selected HVRs was evident for those genetic elements where a significant enrichment was found in the whole HVR sets. Importantly, our approach revealed associations between HVRs with liver-genes and with rat orthologues of human GWAS linked to obesity and metabolic traits.

We also searched genes associated with genomic variation linked to two selected sets of HVRs, one from each strain comparison; namely, those sets with haplotype blocks having at least one conserved peak among rat and mice for each of the three liver-specific transcription factors (i.e. 'HVR 3TF' subset). In these two subsets, we determined those genes with non-synonymous strain-specific variants and genes assigned to promoters with strain-specific variation overlapping the selected haplotype blocks. We reported a list of these selected genes where we included additional information coming from functional analyses and supporting the association of these genes with human GWAS for obesity and metabolic traits and with traits in the susceptible Lyon strains (insulin resistance, dyslipidaemias and obesity) (Additional files 2 and 3). We found a

significant enrichment of liver-expressed genes associated with obesity that were assigned to promoters with strain-specific variation overlapping the selected haplotype block obtained from LHvsLN. Network analyses using these genes and based on protein-protein interactions identified modules implicated in circadian rhythms and in the TGF-beta and transmembrane receptor protein serine/threonine kinase signalling pathways.

Recent studies have described the crucial role of circadian rhythms in regulation of body weight and metabolic process in rodents and other mammals [55–57]. The expression of two genes with strain-specific variation, PER1 and CYR2, are part of the circadian pathway and are regulated in part by the binding of the CLOCK:ARNTL(BMAL1) heterodimer to their promoter regions. PER1 and CRY2 also feed-back to inhibit the CLOCK:ARNTL heterodimer, which itself regulates the transcription of other genes involved in lipid metabolism [56, 57]. Thus, the strain-specific changes in the promoter sequence of PER1 and CRY2 detected using our methodology are plausibly involved in the metabolic symptoms shown by LH strain. Further support for this interpretation is the observation that PER1 is downregulated in liver and kidney in LH and LL strains in comparison with LN strain, whereas CLOCK is upregulated in kidney in LL strain. A relationship also exists between obesity and the TGF-beta and transmembrane receptor protein serine/threonine kinase signalling pathways [58]. Specifically, SMAD7 is an inhibitor of the TGF-beta signalling pathway [58] and is also downregulated in liver in the LH strain. BAMBI cooperates with SMAD7 in the inhibition of the TGF-beta signalling pathway [59], and it is also downregulated in kidney in LH strain. Taken as a whole, the network analyses of the 32 genes associated with obesity demonstrates that at least a fraction of them are plausibly implicated in the obesity phenotype shown by the LH strain.

Ma et al. characterised the blocks with a high density of variants that are unique in the Lyon strains in order to fine-map QTLs for MetS previously identified in these rat strains [31]. As result, the candidate QTL were narrowed by 78%. By focusing their analyses to coding variants in the QTL on rat chromosome 17, they reduced the number of candidate genes to 27. We found that three of these genes had non-synonymous strain-specific variation overlapping the most stringent HVR 3TF subset (ENSRNOG00000014834, ENSRNOG00000022378, ENSRNOG00000027453, see Additional file 1: Table S8), however none of these genes were assigned to promoters with strain-specific variation overlapping HVR 3TF regions. More recently, Wang et al., reported 17 candidate genes involved in the phenotypic differences between LH and LN Lyon rats [41]. We found that only one of Wang et al.'s genes held non-synonymous variation overlapping an HVR 3TF region (ENSRNOG00000016109, Additional file 1: Table S8). In addition, none of the genes identified

by Wang et al., were assigned to promoters with strain-specific variation overlapping the HVR 3TF regions (Additional file 1: Table S8).

The gene *C17h6orf52* (ENSRNOG00000039379), which encodes a protein similar to chromosome 6 open reading frame 52 (C6ORF52), is the only one reported by both studies in the previous paragraph. It is suggested to be the most likely eQTL driver gene involved in phenotypic differences between LH and LN strains [41]. *C17h6orf52* is cis-regulated by an eQTL hotspot on chromosome 17 and is predicted to affect 100 of 278 trans-eQTL genes [41]. While *C17h6orf52* was not the single gene from Wang et al., linked to strain-specific variation overlapping the strict HVR 3TF subset, it is associated with the less restrictive HVRs 2TF subset obtained from LHvsLN comparison. Moreover, *C17h6orf52* was the only gene from the Ma et al. and Wang et al. lists with both non-synonymous coding and promoter assigned SSVs (Fig. 7). This result would suggest that *C17h6orf52* has been under positive selection during the phenotype-driven derivation of this strain and gives support to the predicted role of *C17h6orf52* affecting susceptibility in LH rats for the Metabolic syndrome reported by Ma et al. and Wang et al. The identification of *C17h6orf52* by our complementary method further supports its role in the phenotype and lends additional validation to our general approach.

Conclusions

Our results demonstrate both the potential and the limitations of using the level of functional regulatory conservation to prioritise genomic regions potentially associated with phenotypic differences among Lyon rats. This approach would be most easily extended to other systems with similar breeding histories including other rat strains and mice strains. Importantly, it is not needed to generate data from many individuals like QTL and eQTL approaches and allows the use information already available, as conservation in regulatory elements between rat and mice.

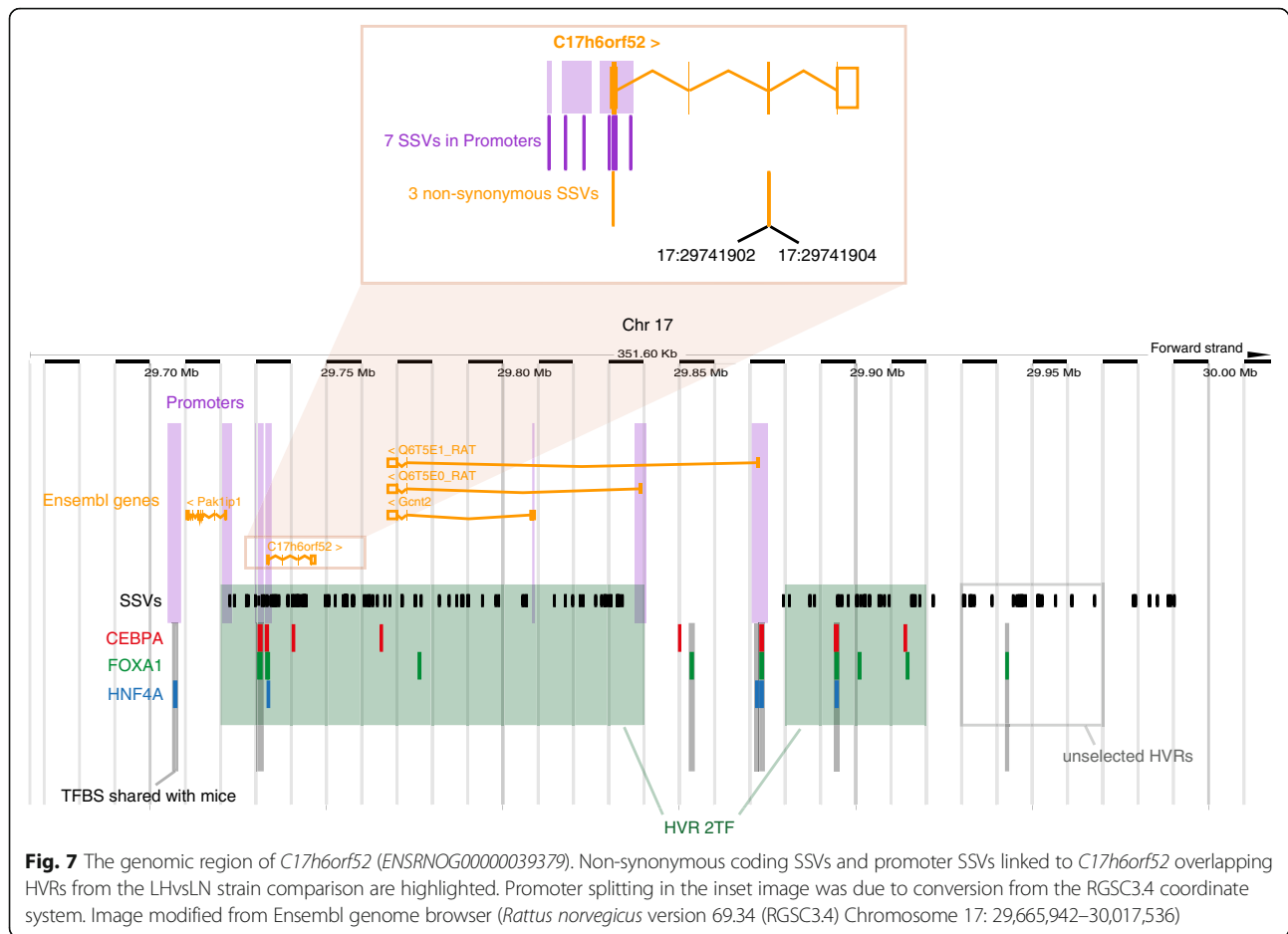
Methods

Determination of high variability regions, low variability regions and unmappable regions

Genomic sequences and single nucleotide variants

We used existing whole genome alignments (*ENA accession: ERP002160*) and single-nucleotide variants (available from the Rat Genome Database) of the three Lyon strains (LH, LL and LN) that were generated by Atanur et al. [30] in comparison to the BN reference genome (RGSC-3.4 [60]).

Strain-Specific Variant (SSV) We called a SSV for a given strain as a genomic position with an allele that is not present in the strain used as reference (Fig. 1a and Additional file 1: Figure S1A). Firstly, we obtained SSVs for Lyon strains compared to the BN reference genome



RGSC-3.4 and the resulting sets of SSVs are referred to as LHvsBN, LLvsBN and LNvsBN for the SSVs specific to the LH, LL and LN strains, respectively. These comparisons were used to calculate the threshold for the different types of genomic regions (see below and Fig. 1c and Additional file 1: Figure S1C). Secondly, we obtained SSVs for the two pairwise comparisons of Lyon rats that are susceptible to MetS and obesity phenotypes relative to the strain that is resistant (LHvsLN and LLvsLN). In these cases, we called a SSV as a genomic position with at least one allele that is not present in both LN and the reference BN genome. By doing this, we discard from LH and LL genomes the genetic variation shared with LN strain, which we assume are not associated with MetS (Fig. 2, Additional file 1: Figures S2 and S3). Furthermore, in a similar way as done for LN strain by comparing it to LH and LL strains (LNvsLH, LNvsLL); these comparisons were used as controls of our approach, because we expected to not find any association between LN strain SSVs and MetS (see Additional file 1: Figures S2 and S3).

In a previous study with the Lyon rats, Ma et al. [31] considered a SSV as any position that differed between the two strains that were being compared regardless of

whether the position was variable with respect to the reference BN genome (indicated as LH + LN, LN + LH in Additional file 1: Figure S2 and LL + LN and LN + LL in Additional file 1: Figure S3). In our case, we considered a SSV only if the allele both differed from the BN reference genome and was also in the strain that was used as query and not present in the strain used as control. For example, in Additional file 1: Figure S1, our approach considers only the G in the LNvsLN comparison, while Ma et al., would have included both the G and the C. This different criterion allowed us to remove from our LHvsLN and LLvsLN analyses those genomic regions specific to LN, and which are likely not associated with the phenotypic differences associated with MetS among Lyon rats (see Additional file 1: Figures S2 and S3). Other differences from methodology used by Ma et al., were i) we did not discard the roughly 5% of SSVs that were called heterozygous by Atanur et al. [30]; and ii) we did not use those genome regions with a low estimated accessibility (see below).

Smooth density of SSVs For downstream analyses, we used a weighted sliding window approach—triangular

smoothing—to calculate the number of SSVs in non-overlapping 10 kb genome windows (Fig. 1b and Additional file 1: Figure S1B). This method smoothes differences among windows that were caused by the genome compartmentalization. For a given window in the genome at position x , we calculated the smoothed density of SSVs as the following floating mean with weights:

$$\text{Smooth Density of SSVs} = \sum_{i=-k}^k SSV_{x+i} \times (k-|i|) / \sum_{i=-k}^k k-|i|$$

where SSV_{x+i} is the number of Strain Specific Variants in the window with position $x+i$, and k is the number of neighbouring windows up and downstream used for smoothing. We use $k=3$ in our analyses empirically (data not shown) because this value gives a clear distinction between two types of genomic regions (see below and Fig. 1c and Additional file 1: Figure S1C).

Genomic regions

High Variability Regions and Low Variability Regions

The smoothed density of SSVs in genome windows between two rat strains shows a bimodal distribution (Fig. 1c and Additional file 1: Figure S1C). The left peak in the bimodal distribution contains regions of the genome identical by descent, with low a density of SSVs (Low Variability Region, LVR). The right peak contains regions of the genome that are divergent between the two strains with a high density of SSVs. A distinct valley separates the two peaks, which we used as a threshold to differentiate HVRs and LVRs. We calculated this threshold for the three comparisons between the Lyon strain rats and the reference rat genome (RGSC-3.4). In all three cases the threshold obtained was three (Additional file 1: Figure S1C); that is, windows with a smoothed SSV density greater than three variants in 10 kb were classified as HVRs, and windows whose smoothed SSV density was less than or equal to three were classified as LVRs. Only those regions with at least three consecutive genome windows of the same type were considered for further analyses (Fig. 1d).

Unmappable regions We performed two analyses on the BAM files to estimate the parameters to characterise the non-accessible genome regions of the Lyon strains. Firstly, we obtained the distribution of mapping qualities (i.e. $-10 \log_{10} \text{Pr}(\text{mapping position is wrong})$, <http://samtools.github.io/hts-specs/>) by using `QualityScoreDistribution.jar` from Picard tools (v1.81 (1299), <http://broadinstitute.github.io/picard/>) with the option `VALIDATION_STRINGENCY = LENIENT` (Additional file 1: Figure S9A). Secondly, we

calculate genome coverage per base by using *genomeCoverageBed* form *Bedtools* (v2.17.0 [61]) with default parameters (Additional file 1: Figure S9B). According to results obtained from the later analysis, we considered a region as unmappable when at least three consecutive windows with an average mapping quality less than or equal to 30 and/or with an average coverage greater than or equal to 100 (Figs. 1d, 2, Additional file 1: Figures S2B and S3B).

Animals

LH/MRrrcAek, LN/MRrrcAek, and LL/MRrrcAek rats were bred and maintained in an approved animal facility at the University of Iowa on a 12-h light-dark cycle and provided food and water ad libitum. Male offspring were used in this study. The rats were phenotyped and tissues collected as previously described [41]. Briefly, at 3 weeks of age the rats were weaned onto normal chow (Teklad 7913 – Harlan Teklad NIH-31 irradiated, 18% protein, 6% fat). At 15 weeks of age they were switched to a 4% NaCl diet (Teklad 7913 modified with 4% NaCl) until they were humanely euthanized with CO_2 at 18 weeks of age after an overnight fast. Tissues were collected and stored in RNAlater (Life Technologies, Grand Island, N.Y.) at -80°C for subsequent RNA extraction.

Gene expression

RNA-seq data

RNA was isolated from liver and kidney tissue using standard TRIzol methods [62]. RNA quality was measured (BioAnalyzer 2100, Agilent Technologies, Santa Clara, CA, USA), using an RIN threshold of 7. Libraries were prepared using TruSeq RNA Sample Preparation Kits v2 (Illumina, San Diego, CA) according to manufacturer's instructions. RNA sequencing was performed on an Illumina HiSeq 2000, with paired-end, 50 bp cycles, at the Iowa Institute of Human Genetics – Genomics Division. Six samples were multiplexed per lane, yielding approximately 30 million reads per sample. All data consisted of six biological replicates for LH and LL liver and five biological replicates for LL liver and LH, LL, and LN kidney. Sequence data from LH and LN liver was previously described (Wang et al. [41]; GSE50027). Remaining sequence data created for this study has been deposited in the ArrayExpress database at EMBL-EBI (www.ebi.ac.uk/arrayexpress) under accession number E-MTAB-5939.

We analysed the read quality using FASTQC software (v 0.10.1 <http://www.bioinformatics.babraham.ac.uk/projects/fastqc/>). Reads were trimmed by using Trimmomatic (v0.32 [63]) if the Phred score of any base was below 25 (LEADING:25 TRAILING:25). We used reads with at least 36 bases (MINLEN:36) and only those paired reads that remained after trimming.

Gene expression analyses

We estimated differential gene expression between LH and LN, and between LL and LN. TopHat (v2.0.13 [64]) was used to map reads to the rat reference genome (RGSC-3.4). Read alignments with more than two mismatches were discarded (*-read-mismatches 2*). We also used the option *-no-novel-juncs* to look for reads across junctions already annotated. The Cufflinks package (v2.2.1 [65]) was used to assemble transcripts separately for each replicate. We used *cuffmerge* to merge the transcript assemblies from replicates to be analysed, and finally we used *cuffdiff* to find differently expressed genes and with Benjamini-Hochberg corrected False Discovery Rates (FDR) of 0.05. Thus, we obtained the differently expressed genes either in liver or in kidney between LH and LN, and between LL and LN (see Additional file 1: Table S1). We also used the FPKM (Fragments per kilobase of exon per million of fragments mapped) values obtained from these analyses to get the list of genes that were expressed in the livers of the LH and LL strains. We considered a gene expressed in liver if its FPKM was greater than 1.0.

Regulation data: liver-specific transcription factors

We used the liver ChIP-seq datasets generated by Stefflova et al. [9] for BN rat strain (ArrayExpress accession: E-MTAB-1414) and for five mouse species/strains (*Mus musculus* (strains C57BL/6 J and A), *Mus caroli*, *Mus castaneus* and *Mus spretus*, ArrayExpress accession: E-MTAB-1414). The dataset comprised two biological replicates for each species/strain and for three liver-specific transcription factors (CEBPA, HNF4A and FOXA1). Reads were aligned using BWA [66] with default parameters. Peak locations were called by SWEMBL (<https://github.com/stevenwilder/SWEMBL>). Final peak sets contained peaks present in both biological replicates.

Conservation of occupancy of three liver-specific transcription factors between rat and mouse strains

We compared peaks generated from ChIP-seq datasets among the five mouse species and the rat. We used only the genomic regions present in the BLAST-Z alignment between mouse and rat available in Ensembl (v59 [44]) and using the NCBI37 mouse genome as references for comparing datasets from the different species. We considered as conserved peaks between rat and mouse, the overlapping peaks between rat and at least one mouse species/strain. Coordinates of conserved peaks were converted to the rat genome reference (RGSC-3.4). For each HVR and liver-transcription factor, we calculated its Conservation Enrichment score (CE_f) as the number of conserved peaks in 10 kb for a given transcription factor.

Permutation tests

Permutation tests were used to find significant enrichments in HVRs. We clustered HVRs within an empirically defined distance of 1 Mb because HVRs have a non-uniform distribution across the genome (see Fig. 3a) and we assume nearby HVRs are regulatorily non-independent. Clusters were randomly permuted across the whole genome by using the command *shuffle* from BEDTools (v2.22.0 [61]) and accessed from *pybedtools* (v0.6.9 [67]); the relative coordinates of HVRs inside of clusters were maintained. We estimated the distribution of expected values by calculating either the total number or average of genetic elements overlapping the set of HVRs inside of the shuffled cluster for each permutation. We performed 10,000 permutations in each test. Significance of the enrichment of the genetic element in HVRs was obtained by calculating the two-tailed *p*-value according to this formula:

$$\text{two-tailed } P\text{value} = \frac{\text{card}(|EV - \overline{EV}| \geq |OV - \overline{EV}|) + 1}{\text{total of permutations} + 1}$$

where *OV* is the value obtained from the observed HVRs, and *EV* is the expected value calculated from each of the 10,000 sets of permuted HVRs (see Fig. 3a). The minimal *p*-value possible with 10,000 permutations is 1×10^{-4} .

Functional analyses of HVRs

Ensembl genes overlapping HVRs

We used the set of Ensembl genes from Ensembl (v69 [44]) for the BN reference genome RGSC-3.4. We used a permutation test (see above) to determine if genes overlapped HVRs more often than expected by chance. We calculated the number of genes overlapping at least one HVR in the observed permuted sets (Fig. 3b).

Differentially expressed genes overlapping HVRs

We tested if genes that were differentially expressed between LH and LN and between LL and LN overlapped HVRs more often than expected by chance (Fig. 3c). We used permutation tests for these analyses (see above). We used the list of genes differentially expressed that were obtained from RNA-seq data (Additional file 1: Table S1). We calculated the number of these genes that overlapped at least one HVR in the observed and permuted sets of HVRs (Fig. 3c and Additional file 1: Table S2).

Gene-annotation enrichment analysis of HVRs

We analysed if there was a functional enrichment associated with metabolic or obesity phenotypes for the genes overlapping at least one HVR that were obtained from LHvsLN ('All HVRs' in Additional file 1: Figure S4) and LLvsLN ('All HVRs' in Additional file 1: Figure S5). We

tested for this enrichment by using DAVID web services v6.7 (python client [45, 46]) for KEGG PATHWAY [47, 48] and UP TISSUE (Uniprot Consortium [49]) databases (release/download date: Sep 2009, <https://david.ncicrf.gov/content.jsp?file=update.html>). We used DAVID v6.7 (Sep 2009) for our analyses rather than DAVID v6.8 (October 2016), because most of the data used in our study (gene annotation, SNPs and occupancies of liver-transcription factors) are based on the RGSCv3.4 assembly, which is also that used by DAVID v6.7. DAVID v6.8 uses the Rnor 6.0 assembly and differences in the gene sets and/or gene nomenclature between these two rat assemblies create inconsistencies that affect the accuracy of our results (data not shown). In addition, although the KEGG PATHWAY resource was updated in DAVID 6.8, the UP TISSUE dataset, which we used to report the expected association between the term liver and the level of functional regulatory conservation (see Results section), was not updated in DAVID v6.8 (in both versions UP TISSUE is dated Sep 2009). We recognise that, in their recent paper, Wadi et al. [68] showed that the use of out-dated gene annotation prevents the identification of all significant terms in enrichment analyses. However, in our case, even when using DAVID v6.7, we found significant results and the expected correlation between gene enrichment and the level of functional regulatory conservation. Thus, the DAVID supporting database that we use are largely the same between v6.7 and v6.8, it is more important for us to be consistent on the assembly and gene set for our analysis.

Liver-specific transcription factor overlapping HVRs

We also tested if the number of peaks in rat overlapping HVRs (Fig. 3e and Additional file 1: Figure S6) and the average CE_f (Additional file 1: Figure S10) observed for each one of the three liver-specific transcription factors was significantly greater than that expected by chance. We used permutation tests (see above) for these analyses. For the observed values, we used either the total number of peaks overlapping the HVRs or the average CE_f for a given transcription factor. For the expected values, we calculated the two latter values for each one of the 10,000 permuted sets of HVRs (see above).

Human GWAS variants associated with metabolic traits overlapping HVRs

We obtained from the NHGRI-EBI GWAS catalogue [50] the list of SNPs associated with obesity and metabolic-related traits in humans (search terms used in Additional file 1: Table S3). SNPs coordinates were converted from the GRCh38 human assembly to the rat RGSC-3.4 assembly using mapping from GRCh38 to Rno6.0 and then from Rno6.0 to Rno5.0 and RGSC-3.4. All conversions used the Ensembl Perl API and the Ensembl assembly converter software (v87 [44]). As with other genetic elements

analysed, we then used permutation tests to determine if there was a significant enrichment of rat orthologous positions for these GWAS variants overlapping HVRs. For the observed value, we used the total number of GWAS variants overlapping HVRs. The expected value was calculated as the number of GWAS variants overlapping each one of the 10,000 sets of permuted sets of HVRs.

Selection of HVRs according to CE_f

For downstream analyses, we created seven subsets of HVRs according to the occupancy for the three liver-specific transcription factors and their CE_f for each on each one of the three liver-specific transcription factors (CE_{CEBPA} , CE_{FOXA1} and CE_{HNF4A}): all HVRs; HVRs with at least one peak (HVR w/TFBS); HVRs with CE_f greater than 0 (i.e. HVRs with at least one conserved peak), and HVRs with CE_f greater than 0.2, 0.4, 0.6, 0.8 respectively (Additional file 1: Tables S9 and S10 show sizes and number of SSVs of HVR subsets). We analysed each one of the subsets of HVRs in a similar way to that used for the full set of HVRs as described above. Then, we compared the results obtained in each analysis across the subsets of HVRs (Fig. 4). Specifically, we analysed the enrichment in HVRs for i) Ensembl genes, ii) gene annotation from DAVID (UP TISSUE and KEGG PATHWAY databases), iii) differentially expressed genes in liver and kidney (Additional file 1: Table S1) and iv) rat orthologues of human GWAS variants associated with obesity and metabolic-related traits (Additional file 1: Tables S3 and S4). Additionally, we also tested if the proportion of non-synonymous coding SSVs (NSC-SSVs) and synonymous coding SSVs (SC-SSVs) in HVRs differed between the subsets of HVRs. For this, we estimated the effect of SSVs in HVRs by using the Ensembl Variant Effect Predictor (VEP) tool (standalone perl script v2.7 associated with Ensembl v69 [69]). We considered in the analyses those NSC-SSVs whose most severe effect was '*missense_variant*', '*stop_gained*' or '*stop_lost*'.

Selection of HVRs by the number of liver-specific transcription factors with conserved peaks

Three subsets of HVRs were created according to how many liver-specific transcription factors had conserved peaks: the '*HVR 1TF*' subset included HVRs with conserved peaks for at least one liver-specific factor, the '*HVR 2TF*' subset had HVRs with conserved peaks for at least two factors, and the '*HVR 3TF*' subset had HVRs with conserved peaks for all three liver-specific factors (Fig. 5a and Additional file 1: Table S11). We compared the functionality among these HVRs subsets to test the importance of the number of transcription factors used to define the conservation level (Fig. 5). We analysed each one of these three subsets of HVRs in a similar way as used for the full set of HVRs and for the HVRs

subsets created with different conservation levels as described in the previous section. Specifically, we compared the enrichment in HVRs among 'HVR 1TF', 'HVR 2TF' and 'HVR 3TF' subsets for i) Ensembl genes, ii) gene annotation from DAVID (liver term of UP_TISSUE database), iii) differentially expressed genes in liver and kidney and iv) rat orthologues of human GWAS variants associated with obesity and metabolic-related traits.

Analyses of SSVs of the selected subsets of HVRs

For these analyses, we selected the subset of HVRs that had at least one conserved peak between rat and mouse strains/species for all three of the liver-specific transcription factors (i.e. 'HVR 3TF' subset) as they show enrichment for most of the functional elements and because of the observed stability of combinatorially bound transcription factors [9]. We limited our analysis to the genes that were both expressed in liver of LH or LL (FPKM >1) and associated with coding or non-coding strain-specific variation.

Non-synonymous coding SSVs (NSC-SSVs) in the selected subsets of HVRs

We assessed the effect of the SSVs on the protein by using the VEP tool (standalone perl script v2.7 [69]). We considered in the analyses those SSVs classified as non-synonymous variants and whose most severe effect was 'missense_variant', 'stop_gained' or 'stop_lost'.

SSVs of the selected subsets of HVRs sited in promoters

Positions of putative promoters in Rat were obtained from Villar et al. [4]. These authors characterised promoters and enhancers by using modifications to histone 3 lysine 27 (H3K27ac) and histone 3 lysine 4 (H3K4me3). Active promoters are marked by H3K4me3 and H3K27ac, while active enhancers are regions marked by H3K27ac [4]. Coordinates were converted from the Rnor5.0 assembly to the RGSC-3.4 assembly using the Ensembl assembly converter software (v80 [44]). Genes were assigned to promoters if the gene's transcription start site (TSS) overlapped or was within 5 kb downstream of the promoter. Only one-to-one gene-promoter assignments were used for our analysis.

Association between metabolic diseases and genes

Genes associated with the three metabolic-related symptoms showed by LH and LL strains (i.e. insulin resistance, dyslipidaemias and obesity) were obtained from DisGeNET (v4.0 [52, 53]). DisGeNET is a platform integrating information on associations between genes and human diseases from public data sources and literature. We analysed those genes expressed in liver and with either the selected NSC-SSV or assigned to selected promoters with SSVs. DisGeNET analysis used the

human orthologous genes of the selected rat genes with homology determined by the Ensembl Perl API (v69 [44]). Only the human orthologous genes with rat homology annotated as 'one2one' or 'apparently one2one' were used. From DisGeNET, we searched for disease gene associations using relevant Unified Medical Language System Concept Unique Identifiers (UMLS® CUIs, insulin resistance: C0021655, dyslipidaemias: C0242339 and obesity: C0028754). We also included two additional diseases not shown by the susceptible Lyon strains as controls for our analyses (heart diseases: C0018799 and Alzheimers: C0002395).

For each of the five diseases, we compared, using Fisher's exact test, the counts of rat genes expressed in liver with NSC-SSVs overlapping the selected subsets of HVRs and human orthologues associated with that disease with the total number of rat genes expressed in liver and human orthologues associated with the disease. A similar comparison was done for genes assigned to promoters with SSVs overlapping the selected HVRs. In this case, the total number of human orthologues of rat genes expressed in liver was limited to those that were one-to-one assigned to promoters.

Network analysis

A network analysis was done for the 32 human orthologues associated with the significant enrichment for obesity (Fig. 6). This analysis was done using Cytoscape 3.5.1 (v3.5.1 [70]) and GeneMANIA (plugging for Cytoscape v3.4.1 [54]) with default parameters, which involves the use of information from six type of sources ('Co-expression', 'Physical interactions', 'Predicted', 'Co-localization', 'Pathway' and 'Genetic Interactions'). For downstream analyses, we focused on the subnetwork obtained from the protein interaction databases ('Physical interactions') because it highlighted two modules in the resulting network with 10 and 11 genes, respectively. The functional enrichment analysis to identify which Gene Ontology terms were significantly enriched in the network was also done with GeneMANIA (Additional file 1: Table S7).

Additional files

Additional file 1: Supplementary figures and tables. (DOCX 724 kb)

Additional file 2: All relevant characteristics of genes that are overlapping HVRs with conserved TFBS from all three factors and have non-synonymous coding SSVs in at least one strain comparison. (XLSX 56 kb)

Additional file 3: All relevant characteristics of genes that are overlapping HVRs with conserved TFBS from all three factors and have SSVs in promoters that can be one-to-one associated with genes in at least one strain comparison. (XLSX 70 kb)

Abbreviations

(N)SC-SSV: (Non-)synonymous coding strain-specific variant; CE_i: Conservation Enrichment score; FPKM: Fragments per kilobase of exon per million of fragments mapped; HVR {1,2,3}TF: HVRs with one or more conserved TFBS

from at least {one, two or three} factor(s); HVR: High Variability Region; LH: Lyon Hypertensive; LL: Lyon Low pressure; LN: Lyon Normotensive; LVR: Low Variability Region; MetS: Metabolic syndrome; QTL: Quantitative trait loci; SSV: Strain-specific variant; TFBS: Transcription factor binding site

Acknowledgments

We thank John Marioni for helpful discussions and Duncan Odom for a critical reading of the manuscript. RNA-seq data were obtained at the Genomics Division of the Iowa Institute of Human Genetics, which is supported, in part, by the University of Iowa Carver College of Medicine.

Funding

The research leading to these results has received funding from the European Community's Seventh Framework Programme (FP7/2007–2013) under grant agreement N° HEALTH-F4-2010-241504 (EURATRANS), and from NIH R01 HL089895 (AEK) and R21 DK089417 (AEK). We acknowledge additional support from the Wellcome Trust (WT108749/Z/15/Z) and the European Molecular Biology Laboratory.

Availability of data and materials

Datasets generated or used during the current study are available for download from public repositories, these can be located by the accession numbers provided throughout the text and repeated here for completeness: ERP002160, E-MTAB-5939, E-MTAB-1414. Remaining data and materials are contained within the manuscript, described in the Additional files section, or available from the corresponding authors on reasonable request.

Authors' contributions

DMG, DT and PF conceived and designed the study. DDdS determined high variability regions, low variability regions and unmappable regions. DMG conducted the analysis of the gene expression and regulation data, selected HVRs and did the functional analysis of them. DMG did association analysis between diseases and genes and did the network analysis. AEK designed the RNA-seq component of the study. MCJM and AEK performed tissue collection, sample preparation and RNA-seq data generation. DMG, AEK, DT and PF wrote the manuscript and created figures with significant comments and critical assessment by the other authors. All authors read and approved the final manuscript.

Ethics approval

All animal protocols needed for these studies were approved by the Institutional Animal Care and Use Committee (IACUC) at the University of Iowa.

Consent for publication

Not applicable.

Competing interests

PF is a member of the scientific advisory boards of Fabric Genomics, Inc., and Eagle Genomics, Ltd.

Publisher's Note

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

Author details

¹European Molecular Biology Laboratory, European Bioinformatics Institute, Wellcome Genome Campus, Hinxton, Cambridge CB10 1SD, UK.

²Department of Pharmacology, University of Iowa, Iowa City, IA, USA. ³Iowa Institute of Human Genetics, University of Iowa, Iowa City, IA, USA. ⁴Present address: MD Anderson Cancer Center, University of Texas, Houston, TX, USA. ⁵Present address: Earlham Institute, Norwich research Park, Norwich NR4 7UH, UK.

Received: 7 August 2017 Accepted: 27 November 2017

Published online: 22 December 2017

References

- Wittkopp PJ, Kalay G. Cis-regulatory elements: molecular mechanisms and evolutionary processes underlying divergence. *Nat Rev Genet.* 2012; 13(1):59–69.

- Pai AA, Gilad Y. Comparative studies of gene regulatory mechanisms. *Curr Opin Genet Dev.* 2014;29:68–74.
- Villar D, Flicek P, Odom DT. Evolution of transcription factor binding in metazoans - mechanisms and functional implications. *Nat Rev Genet.* 2014; 15(4):221–33.
- Villar D, Berthelot C, Aldridge S, Rayner TF, Lukk M, Pignatelli M, Park TJ, Deaville R, Erichsen JT, Jasinska AJ, et al. Enhancer evolution across 20 mammalian species. *Cell.* 2015;160(3):554–66.
- Lowdon RF, Jang HS, Wang T. Evolution of epigenetic regulation in vertebrate genomes. *Trends Genet.* 2016;32(5):269–83.
- Mack KL, Nachman MW. Gene regulation and speciation. *Trends Genet.* 2017;33(1):68–80.
- Shibata Y, Sheffield NC, Fedrigo O, Babbitt CC, Wortham M, Tewari AK, London D, Song L, Lee B-K, Iyer VR. Extensive evolutionary changes in regulatory element activity during human origins are associated with altered gene expression and positive selection. *PLoS Genet.* 2012;8(6):e1002789.
- Romero IG, Ruvinsky I, Gilad Y. Comparative studies of gene expression and the evolution of gene regulation. *Nat Rev Genet.* 2012;13(7):505–16.
- Stefflova K, Thybert D, Wilson MD, Streeter I, Aleksic J, Karagianni P, Brazma A, Adams DJ, Talianidis I, Marioni JC, et al. Cooperativity and rapid evolution of cobound transcription factors in closely related mammals. *Cell.* 2013;154:530–40.
- Ward LD, Kellis M. Interpreting noncoding genetic variation in complex traits and human disease. *Nat Biotechnol.* 2012;30(11):1095–106.
- Nica AC, Dermitzakis ET. Expression quantitative trait loci: present and future. *Phil Trans R Soc B.* 2013;368(1620):20120362.
- Lowe WL, Reddy TE. Genomic approaches for understanding the genetics of complex disease. *Genome Res.* 2015;25(10):1432–41.
- Moreno-Moral A, Petretto E. From integrative genomics to systems genetics in the rat to link genotypes to phenotypes. *Dis Model Mech.* 2016;9(10):1097–110.
- Jacob HJ. The rat: a model used in biomedical research. *Methods Mol Biol.* 2010;597:1–11.
- Lindsey JR, Baker HJ. Chapter 1 - historical foundations. In: Franklin MASHWL, editor. *The laboratory rat.* 2nd ed. Burlington: Academic Press; 2006. p. 1–52.
- Aitman T, Dhillon P, Geurts AM. A RAtional choice for translational research? *Dis Model Mech.* 2016;9(10):1069–1072.
- Voigt B. Rat strain repositories. *Methods Mol Biol.* 2010;597:323–31.
- Mashimo T, Serikawa T. Rat resources in biomedical research. *Curr Pharm Biotechnol.* 2009;10(2):214–20.
- Yau AC, Holmdahl R. Rheumatoid arthritis: identifying and characterising polymorphisms using rat models. *Dis Model Mech.* 2016;9(10):1111–23.
- Dwinell MR, Lazar J, Geurts AM. The emerging role for rat models in gene discovery. *Mamm Genome.* 2011;22(7–8):466–75.
- Shimoyama M, De Pons J, Hayman GT, Laulederkind SJ, Liu W, Nigam R, Petri V, Smith JR, Tutaj M, Wang S-J. The rat genome database 2015: genomic, phenotypic and environmental variations and disease. *Nucleic Acids Res.* 2014; doi:10.1093/nar/gku1026.
- Aitman T, Petretto E, Behmoaras J. Genetic mapping and positional cloning. *Methods Mol Biol.* 2010;597:13–32.
- Baud A, Hermesen R, Guryev V, Stridh P, Graham D, McBride MW, Foroud T, Calderari S, Diez M, Ockinger J, et al. Combined sequence-based and genetic mapping analysis of complex traits in outbred rats. *Nat Genet.* 2013;45:767–75.
- Kimura M. Evolutionary rate at the molecular level. *Nature.* 1968;217(5129):624–6.
- King JL, Jukes TH. Non-Darwinian evolution. *Science.* 1969;164(3881):788–98.
- Hermesen R, de Ligt J, Spee W, Blokzijl F, Schafer S, Adami E, Boymans S, Flink S, van Bostel R, van der Weide RH, et al. Genomic landscape of rat strain and substrain variation. *BMC Genomics.* 2015;16:357.
- Consortium GP. A global reference for human genetic variation. *Nature.* 2015;526(7571):68–74.
- Adams DJ, Doran AG, Lilue J, Keane TM. The mouse genomes project: a repository of inbred laboratory mouse strain genomes. *Mamm Genome.* 2015;26(9–10):403–12.
- Saar K, Beck A, Bihoreau MT, Birney E, Brocklebank D, Chen Y, Cuppen E, Demonchy S, Dopazo J, Flicek P, et al. SNP and haplotype mapping for genetic analysis in the rat. *Nat Genet.* 2008;40(5):560–6.
- Atanur SS, Diaz AG, Maratou K, Sarkis A, Rotival M, Game L, Tschannen MR, Kaisaki PJ, Otto GW, Ma MC, et al. Genome sequencing reveals loci under artificial selection that underlie disease phenotypes in the laboratory rat. *Cell.* 2013;154(3):691–703.

31. Ma MCJ, Atanur S, Aitman T, Kwitek A. Genomic structure of nucleotide diversity among Lyon rat models of metabolic syndrome. *BMC Genomics*. 2014;15(1):197.
32. Cuppen E. Haplotype-based genetics in mice and rats. *Trends Genet*. 2005; 21(6):318–22.
33. Encode Project Consortium, Bernstein BE, Birney E, Dunham I, Green ED, Gunter C, Snyder M. An integrated encyclopedia of DNA elements in the human genome. *Nature*. 2012;489(7414):57–74.
34. Berman BP, Nibu Y, Pfeiffer BD, Tomancak P, Celniker SE, Levine M, Rubin GM, Eisen MB. Exploiting transcription factor binding site clustering to identify cis-regulatory modules involved in pattern formation in the drosophila genome. *Proc Natl Acad Sci U S A*. 2002;99(2):757–62.
35. Wong ES, Thybert D, Schmitt BM, Stefflova K, Odom DT, Flicek P. Decoupling of evolutionary changes in transcription factor binding and gene expression in mammals. *Genome Res*. 2015;25(2):167–78.
36. Pennacchio LA, Rubin EM. Genomic strategies to identify mammalian regulatory sequences. *Nat Rev Genet*. 2001;2(2):100–9.
37. Cheng Y, Ma Z, Kim B-H, Wu W, Cayting P, Boyle AP, Sundaram V, Xing X, Dogan N, Li J. Principles of regulatory information conservation between mouse and human. *Nature*. 2014;515(7527):371–5.
38. Dupont J, Dupont JC, Froment A, Milon H, Vincent M. Selection of three strains of rats with spontaneously different levels of blood pressure. *Biomedicine*. 1973;19(1):36–41.
39. Vincent M, Boussairi EH, Cartier R, Lo M, Sassolas A, Cerutti C, Barres C, Gustin MP, Cuisinaud G, Samani NJ, et al. High blood pressure and metabolic disorders are associated in the Lyon hypertensive rat. *J Hypertens*. 1993;11(11):1179–85.
40. Sassolas A, Vincent M, Benzoni D, Sassard J. Plasma lipids in genetically hypertensive rats of the Lyon strain. *J Cardiovasc Pharmacol*. 1981;3(5):1008–14.
41. Wang J, Ma MCJ, Mennie AK, Pettus JM, Xu Y, Lin L, Traxler MG, Jakoubek J, Atanur SS, Aitman TJ, et al. Systems biology with high-throughput sequencing reveals genetic mechanisms underlying the metabolic syndrome in the Lyon hypertensive rat. *Circ Cardiovasc Genet*. 2015;8(2):316–26.
42. Bilusic M, Bataillard A, Tschannen MR, Gao L, Barreto NE, Vincent M, Wang T, Jacob HJ, Sassard J, Kwitek AE. Mapping the genetic determinants of hypertension, metabolic diseases, and related phenotypes in the Lyon hypertensive rat. *Hypertension*. 2004;44(5):695–701.
43. Kaur J. A comprehensive review on metabolic syndrome. *Cardiol Res Pract*. 2014;2014:943162.
44. Yates A, Akanni W, Amode MR, Barrell D, Billis K, Carvalho-Silva D, Cummins C, Clapham P, Fitzgerald S, Gil L. Ensembl 2016. *Nucleic Acids Res*. 2016; 44(D1):D710–6.
45. Jiao X, Sherman BT, Huang DW, Stephens R, Baseler MW, Lane HC, Lempicki RA. DAVID-WIS: a stateful web service to facilitate gene/protein list analysis. *Bioinformatics*. 2012;28(13):1805–6.
46. Huang d W, Sherman BT, Lempicki RA. Bioinformatics enrichment tools: paths toward the comprehensive functional analysis of large gene lists. *Nucleic Acids Res*. 2009;37(1):1–13.
47. Kanehisa M, Goto S, Sato Y, Kawashima M, Furumichi M, Tanabe M. Data, information, knowledge and principle: back to metabolism in KEGG. *Nucleic Acids Res*. 2014;42(Database issue):D199–205.
48. Kanehisa M, Goto S. KEGG: kyoto encyclopedia of genes and genomes. *Nucleic Acids Res*. 2000;28(1):27–30.
49. UniProt Consortium. UniProt: a hub for protein information. *Nucleic Acids Res*. 2015;43(Database issue):D204–12.
50. Welter D, MacArthur J, Morales J, Burdett T, Hall P, Junkins H, Klemm A, Flicek P, Manolio T, Hindorf L, et al. The NHGRI GWAS catalog, a curated resource of SNP-trait associations. *Nucleic Acids Res*. 2014;42(Database issue):D1001–6.
51. Ballester B, Medina-Rivera A, Schmidt D, Gonzalez-Porta M, Carlucci M, Chen XT, Chessman K, Faure AJ, Funnell APW, Goncalves A, et al. Multi-species, multi-transcription factor binding highlights conserved control of tissue-specific biological pathways. *elife*. 2014;3:e02626.
52. Piñero J, Queralt-Rosinach N, Bravo À, Deu-Pons J, Bauer-Mehren A, Baron M, Sanz F, Furlong LI. DisGeNET: a discovery platform for the dynamical exploration of human diseases and their genes. *Database (Oxford)*. 2015; bav028.
53. Piñero J, Bravo À, Queralt-Rosinach N, Gutiérrez-Sacristán A, Deu-Pons J, Centeno E, García-García J, Sanz F, Furlong LI. DisGeNET: a comprehensive platform integrating information on human disease-associated genes and variants. *Nucleic Acids Res*. 2016; doi:10.1093/nar/gkw943.
54. Montojo J, Zuberi K, Rodriguez H, Kazi F, Wright G, Donaldson SL, Morris Q, Bader GD. GeneMANIA Cytoscape plugin: fast gene function predictions on the desktop. *Bioinformatics*. 2010;26(22):2927–8.
55. Summa KC, Turek FW. Chronobiology and obesity: interactions between circadian rhythms and energy regulation. *Adv Nutr*. 2014;5(3):312S–9S.
56. Froy O. Circadian rhythms and obesity in mammals. *ISRN Obes*. 2012; 437198.
57. Maury E, Ramsey KM, Bass J. Circadian rhythms and metabolic syndrome. *Circ Res*. 2010;106(3):447–62.
58. Aihara K-I, Ikeda Y, Yagi S, Akaike M, Matsumoto T. Transforming growth factor-β1 as a common target molecule for development of cardiovascular diseases, renal insufficiency and metabolic syndrome. *Cardiol Res Pract*. 2011;175381.
59. Yan X, Lin Z, Chen F, Zhao X, Chen H, Ning Y, Chen Y-G. Human BAMB1 cooperates with Smad7 to inhibit transforming growth factor-β signaling. *J Biol Chem*. 2009;284(44):30097–104.
60. Gibbs RA, Weinstock GM, Metzker ML, Muzny DM, Sodergren EJ, Scherer S, Scott G, Steffen D, Worley KC, Burch PE, et al. Genome sequence of the Brown Norway rat yields insights into mammalian evolution. *Nature*. 2004; 428(6982):493–521.
61. Quinlan AR, Hall IM. BEDTools: a flexible suite of utilities for comparing genomic features. *Bioinformatics*. 2010;26:841–2.
62. Chomczynski P, Sacchi N. Single-step method of RNA isolation by acid guanidinium thiocyanate-phenol-chloroform extraction. *Anal Biochem*. 1987; 162(1):156–9.
63. Bolger AM, Lohse M, Usadel B. Trimmomatic: a flexible trimmer for Illumina sequence data. *Bioinformatics*. 2014;30(15):2114–20.
64. Trapnell C, Pachter L, Salzberg SL. TopHat: discovering splice junctions with RNA-Seq. *Bioinformatics*. 2009;25(9):1105–11.
65. Trapnell C, Williams BA, Pertea G, Mortazavi A, Kwan G, van Baren MJ, Salzberg SL, Wold BJ, Pachter L. Transcript assembly and quantification by RNA-Seq reveals unannotated transcripts and isoform switching during cell differentiation. *Nat Biotechnol*. 2010;28(5):511–5.
66. Li H, Durbin R. Fast and accurate short read alignment with burrows-wheeler transform. *Bioinformatics*. 2009;25:1754–60.
67. Dale RK, Pedersen BS, Quinlan AR. Pybedtools: a flexible python library for manipulating genomic datasets and annotations. *Bioinformatics*. 2011; 27(24):3423–4.
68. Wadi L, Meyer M, Weiser J, Stein LD, Reimand J. Impact of outdated gene annotations on pathway enrichment analysis. *Nat Methods*. 2016;13(9):705–6.
69. McLaren W, Gil L, Hunt SE, Riat HS, Ritchie GR, Thormann A, Flicek P, Cunningham F. The ensembl variant effect predictor. *Genome Biol*. 2016; 17(1):122.
70. Shannon P, Markiel A, Ozier O, Baliga NS, Wang JT, Ramage D, Amin N, Schwikowski B, Ideker T. Cytoscape: a software environment for integrated models of biomolecular interaction networks. *Genome Res*. 2003;13(11):2498–504.

Submit your next manuscript to BioMed Central and we will help you at every step:

- We accept pre-submission inquiries
- Our selector tool helps you to find the most relevant journal
- We provide round the clock customer support
- Convenient online submission
- Thorough peer review
- Inclusion in PubMed and all major indexing services
- Maximum visibility for your research

Submit your manuscript at
www.biomedcentral.com/submit

