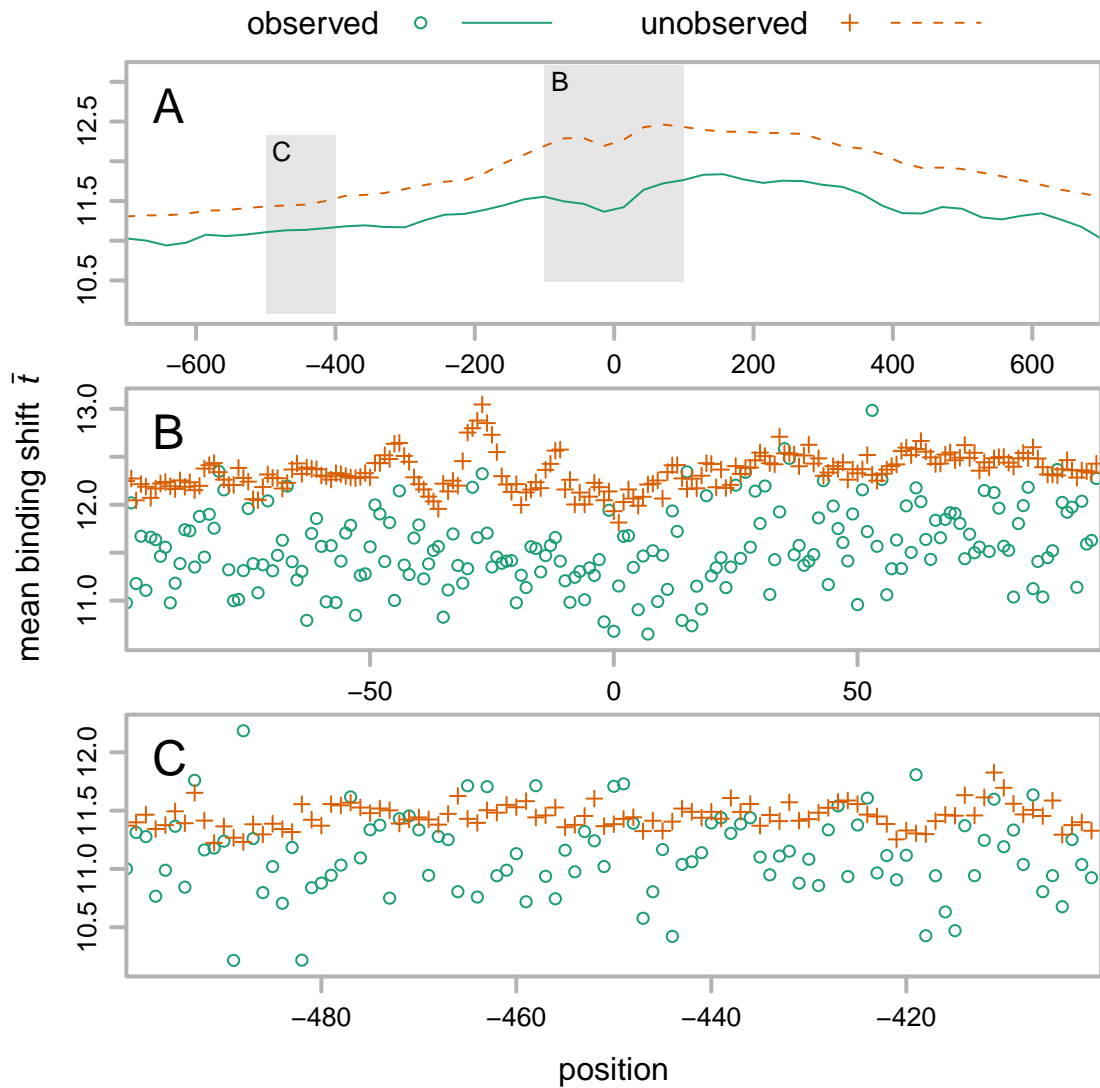
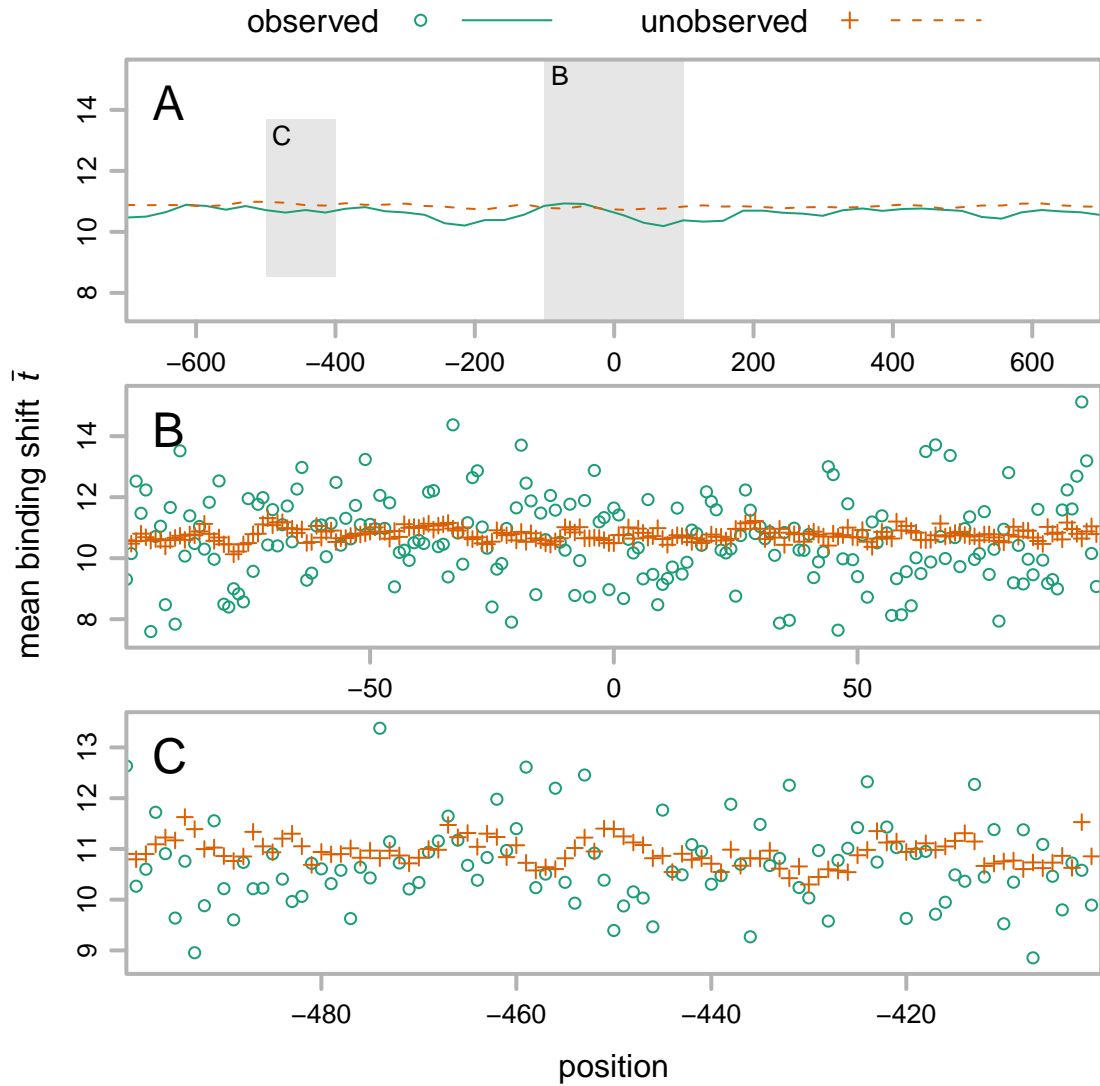


## **Supplementary Information**

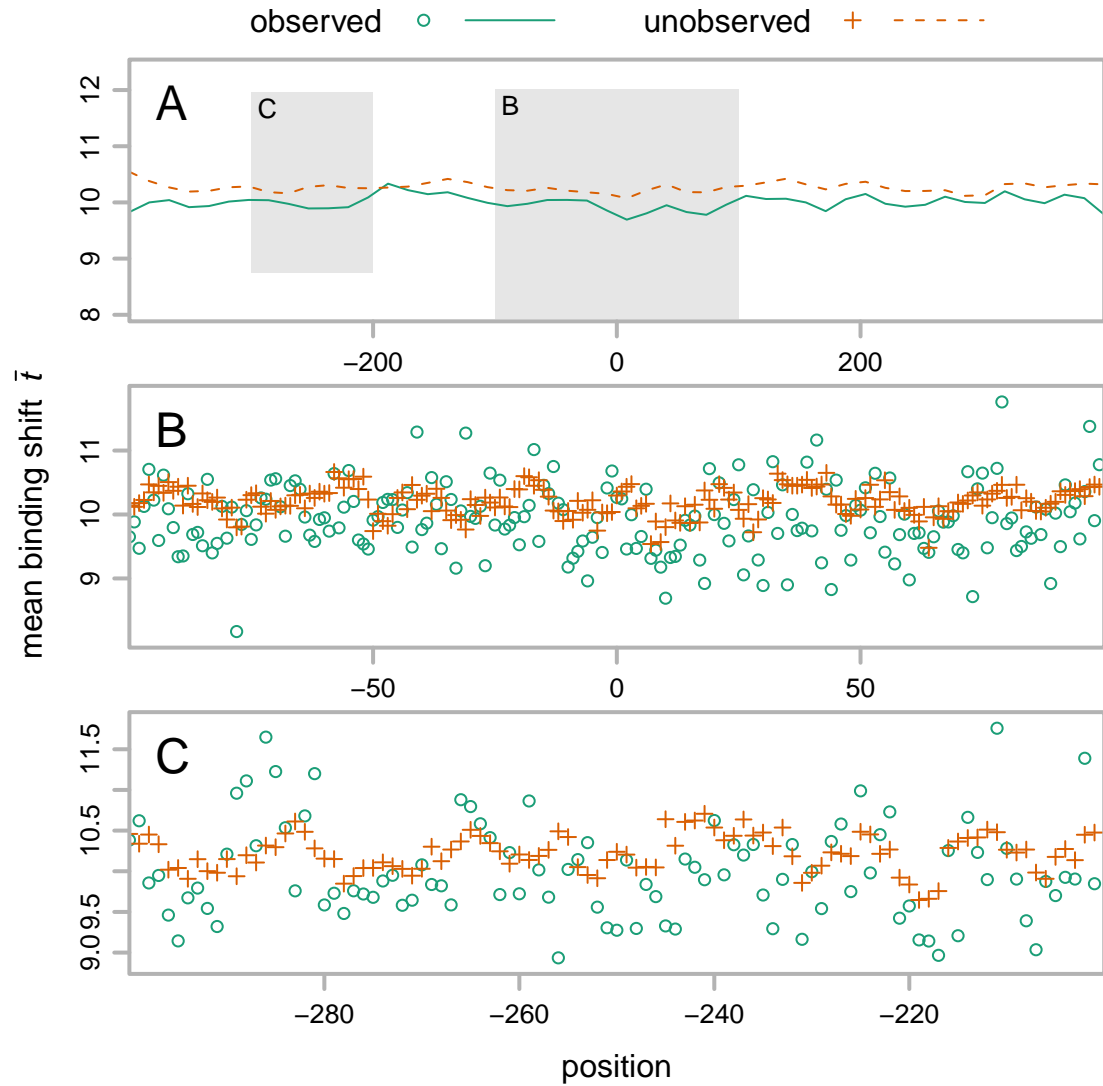
### **Supplementary Figures and Legends**



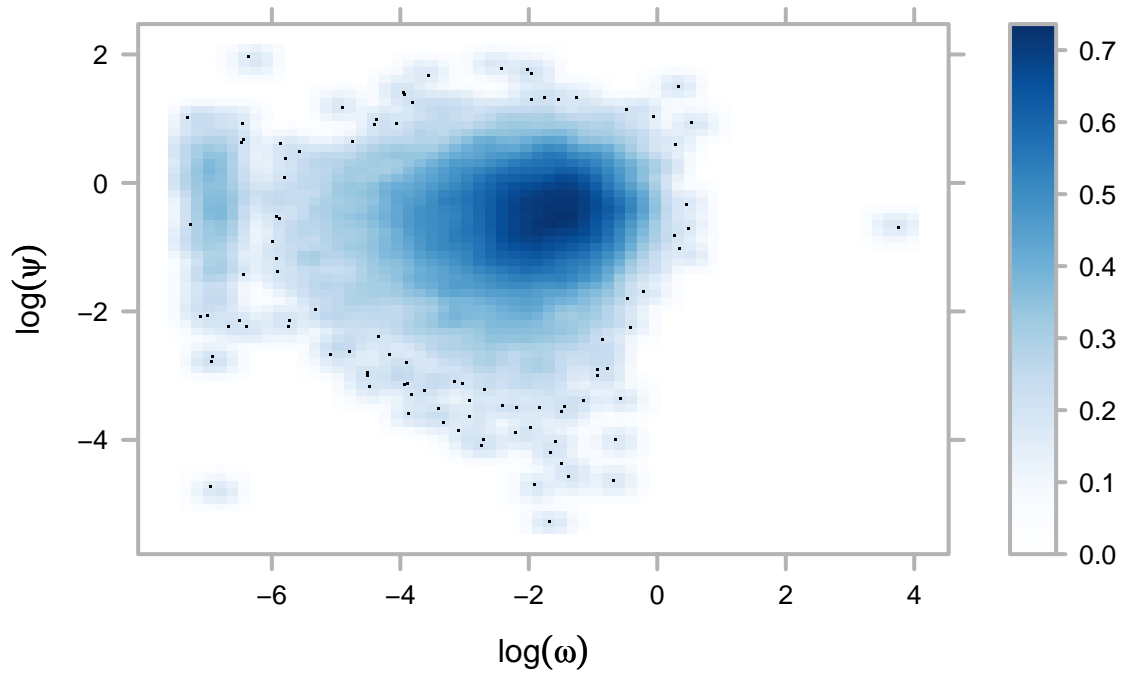
Supplementary Figure 1: Aggregation plot of the binding shifts of 17,194 mouse genes, averaged within two groups: one where the simulated mutation was observed in rat (green circles and solid line), and one where it was unobserved (orange crosses and dashed line). (A) Local regressions for  $\pm 700$  bp around the TSS, estimated with the loess (Cleveland and Devlin 1988) function in R (R Development Core Team 2007), with second-degree polynomials and  $\alpha = 0.1$ . Shaded regions in this plot are magnified as separate panels beneath to show mean binding shifts at individual positions (B) proximal to and (C) more distal from the TSS.



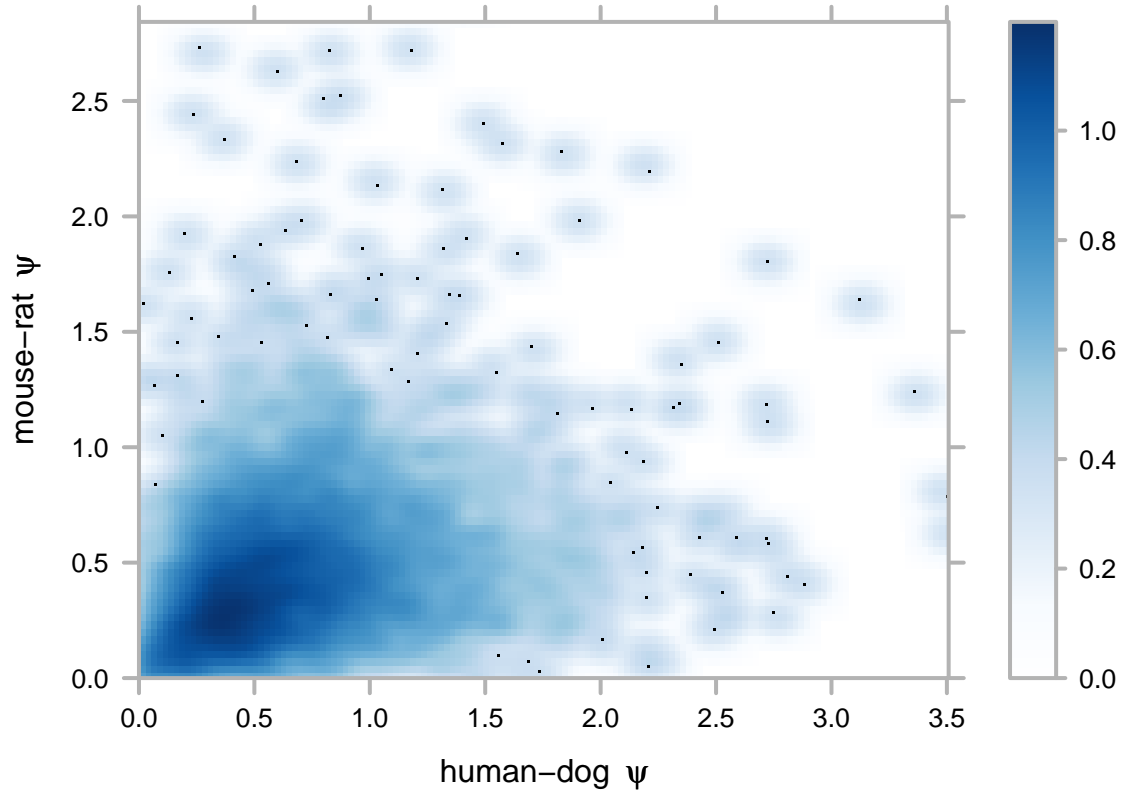
Supplementary Figure 2: Aggregation plot of the binding shifts of 496 human regions with enhancer activity validated in transgenic mice (Pennacchio et al. 2006), averaged within two groups: one where the simulated mutation was observed in mouse (green circles and solid line), and one where it was unobserved (orange crosses and dashed line). (A) Local regressions for  $\pm 700$  bp around the midpoint of the enhancer region, estimated with the loess (Cleveland and Devlin 1988) function in R (R Development Core Team 2007), with second-degree polynomials and  $\alpha = 0.1$ . Shaded regions in this plot are magnified as separate panels beneath to show mean binding shifts at individual positions (B) proximal to and (C) more distal from the midpoint of the enhancer region.



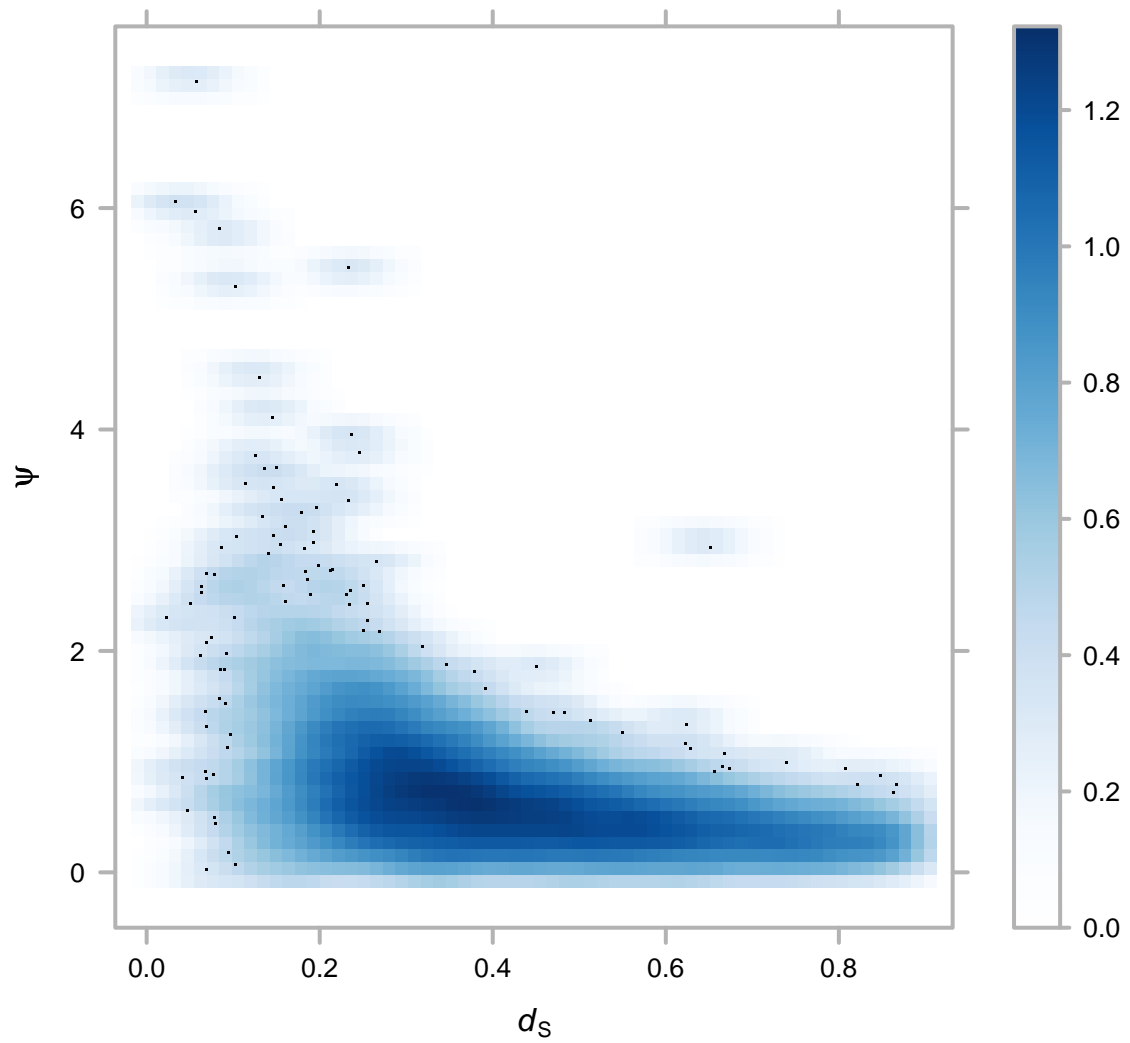
Supplementary Figure 3: Aggregation plot of the binding shifts of 790 human ancestral repeat regions (Paten et al. 2008) with length of at least 1000, averaged within two groups: one where the simulated mutation was observed in mouse (green circles and solid line), and one where it was unobserved (orange crosses and dashed line). (A) Local regressions for  $\pm 400$  bp around the midpoint of the ancestral repeats, estimated with the loess (Cleveland and Devlin 1988) function in R (R Development Core Team 2007), with second-degree polynomials and  $\alpha = 0.1$ . Shaded regions in this plot are magnified as separate panels beneath to show mean binding shifts at individual positions (B) proximal to and (C) more distal from the midpoint of the ancestral repeats.



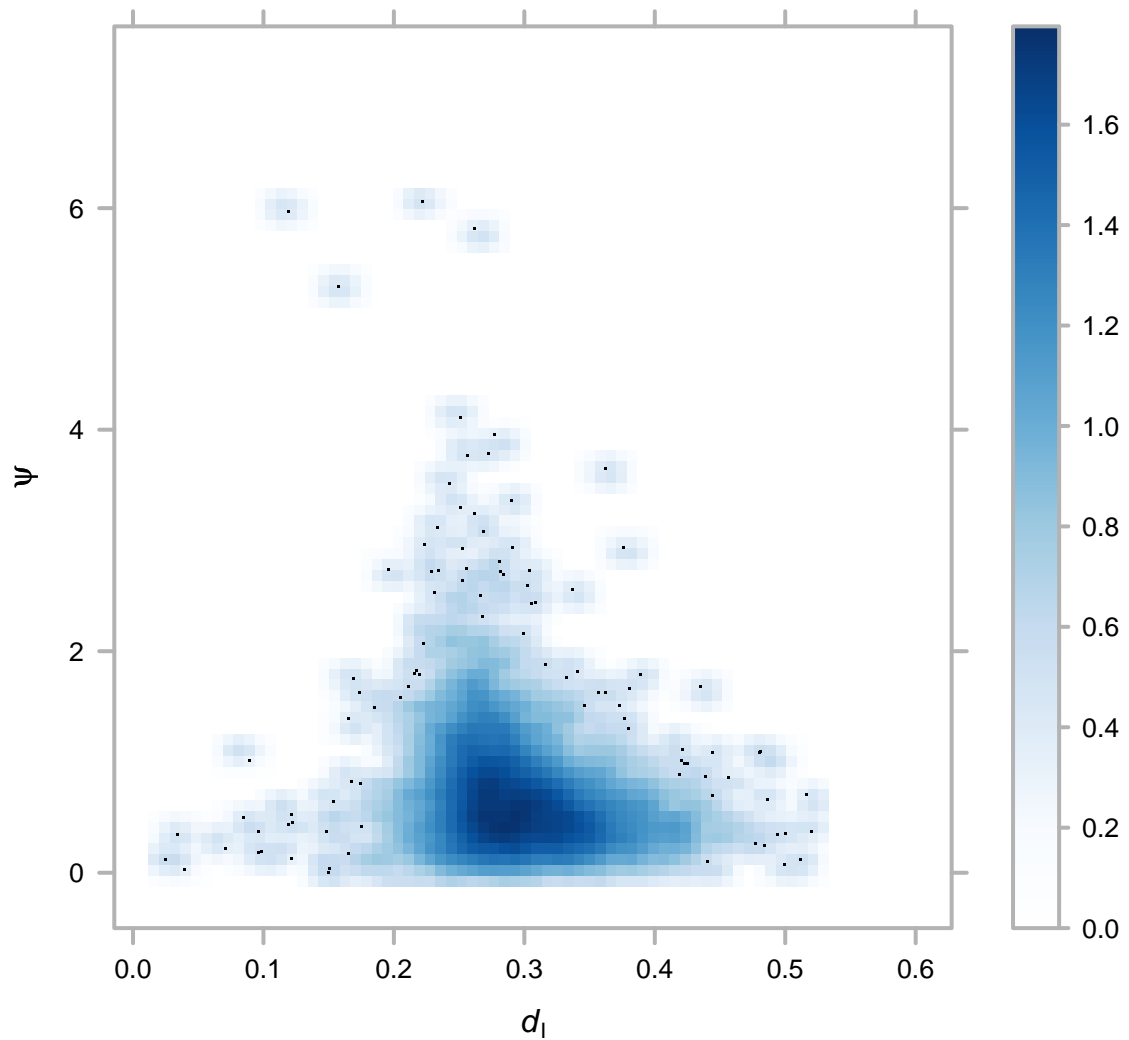
Supplementary Figure 4: Smoothed log-log scatter plot of  $\psi$  against  $\omega$  for 9060 human-dog transcripts, excluding transcripts where  $d_T$ ,  $d_N$ , or  $d_S$  is 0. The plot space is divided into a number of cells, and the color of each cell indicates the fourth root of a two-dimensional kernel density estimate of the number of data points in that cell, as performed by the *geneplotter* package of Bioconductor (Gentleman et al. 2004). The color key to the right of the plot indicates how various colors correspond to the transformed density estimate values.



Supplementary Figure 5: Scatter plot of mouse-rat  $\psi$  against human-dog  $\psi$  for 4547 1:1:1:1 orthologous transcripts. The plot space is divided into a number of cells, and the color of each cell indicates the fourth root of a two-dimensional kernel density estimate of the number of data points in that cell, as performed by the `geneplotter` package of Bioconductor (Gentleman et al. 2004). Only the bottom 99.9 percent of mouse-rat and human-dog  $\psi$  values are displayed, to eliminate the distorting effect of outliers. The color key to the right of the plot indicates how various colors correspond to the transformed density estimate values.

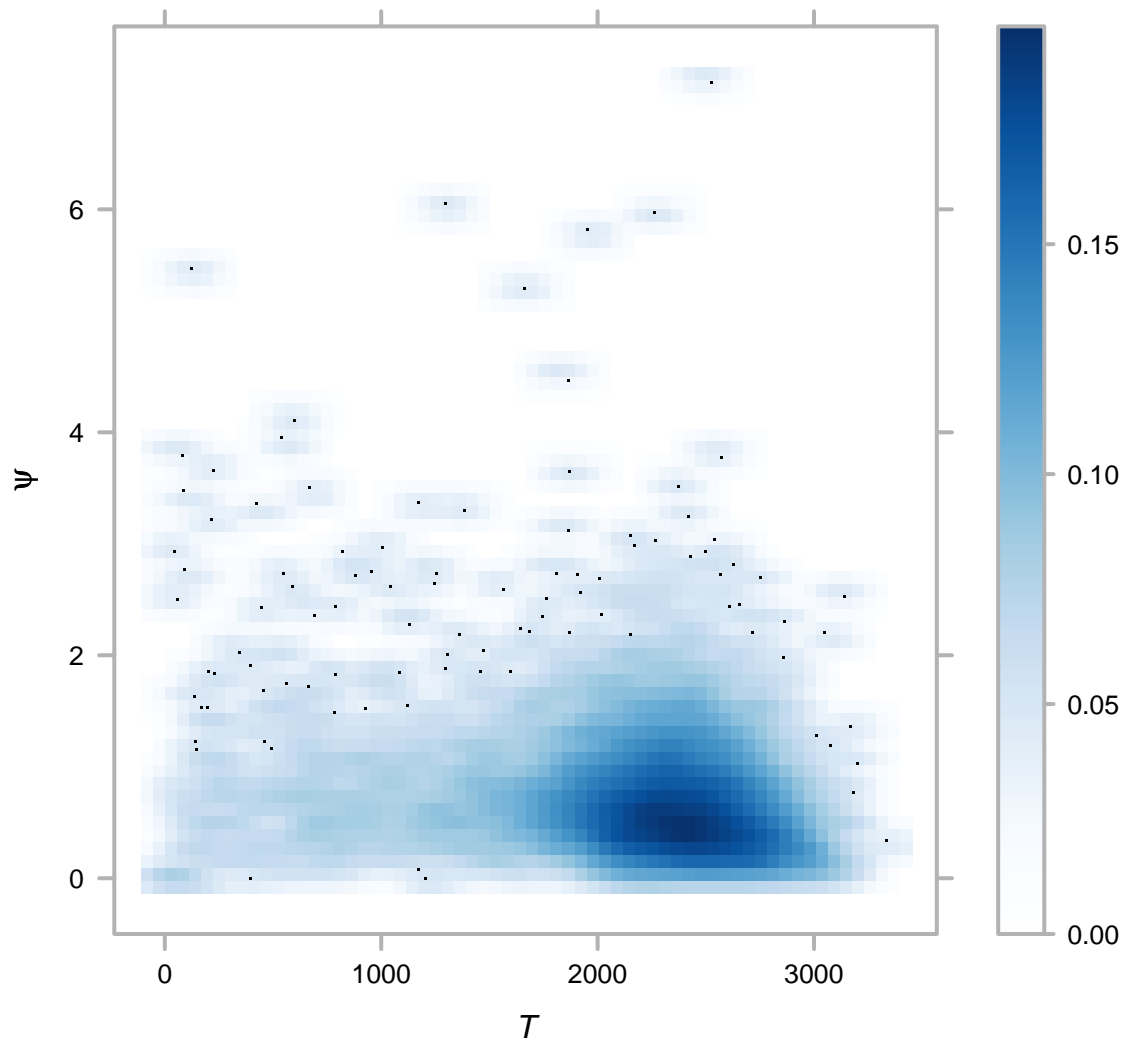


Supplementary Figure 6: Smoothed scatter plot of  $\psi$  against  $d_S$  for 17,600 human-dog transcripts. The plot space is divided into a number of cells, and the color of each cell indicates the fourth root of a two-dimensional kernel density estimate of the number of data points in that cell, as performed by the `genepLOTter` package of Bioconductor (Gentleman et al. 2004). The color key to the right of the plot indicates how various colors correspond to the transformed density estimate values.



Supplementary Figure 7: Smoothed scatter plot of  $\psi$  against  $d_I$  for 17,600 human-dog transcripts. The plot space is divided into a number of cells, and the color of each cell indicates the fourth root of a two-dimensional kernel density estimate of the number of data points in that cell, as performed by the `geneplotter` package of Bioconductor (Gentleman et al. 2004). The color key to the right of the plot indicates how various colors correspond to the transformed density estimate values.





Supplementary Figure 8: Smoothed scatter plot of  $\psi$  against  $T$  for 17,600 human-dog transcripts. The plot space is divided into a number of cells, and the color of each cell indicates the fourth root of a two-dimensional kernel density estimate of the number of data points in that cell, as performed by the `geneflotter` package of Bioconductor (Gentleman et al. 2004). The color key to the right of the plot indicates how various colors correspond to the transformed density estimate values.

## Supplementary Tables

Supplementary Table 1: Identifiers used by JASPAR CORE for the 89 TFs used in the model, along with the HGNC symbol for the TF's gene or its human ortholog.

JASPAR identifier	HGNC symbol
Ar	<i>AR</i>
Arnt	<i>ARNT</i>
Arnt-Ahr	<i>AHR</i>
Bapx1	<i>NKX3-2</i>
c-ETS	<i>ETS1</i>
cEBP	<i>CPEB1</i>
Chop-cEBP	<i>DDIT3</i>
CREB1	<i>CREB1</i>
deltaEF1	<i>ZEB1</i>
E2F1	<i>E2F1</i>
ELK1	<i>ELK1</i>
ELK4	<i>ELK4</i>
En1	<i>EN1</i>
ESR1	<i>ESR1</i>
Evi1	<i>EVI1</i>
Fos	<i>FOS</i>
Foxa2	<i>FOXA2</i>
FOXC1	<i>FOXC1</i>
FOXD1	<i>FOXD1</i>
Foxd3	<i>FOXD3</i>
FOXF2	<i>FOXF2</i>
FOXI1	<i>FOXI1</i>
FOXL1	<i>FOXL1</i>
Foxq1	<i>FOXQ1</i>
GABPA	<i>GABPA</i>
Gata1	<i>GATA1</i>
GATA2	<i>GATA2</i>
GATA3	<i>GATA3</i>
Gfi	<i>GFI1</i>
HAND1-TCF3	<i>TCF3</i>
HLF	<i>HLF</i>
HNF4	<i>HNF4A</i>
IRF1	<i>IRF1</i>
IRF2	<i>IRF2</i>
Klf4	<i>KLF4</i>
MafB	<i>MAFB</i>
MAX	<i>MAX</i>
MEF2A	<i>MEF2A</i>
Myb	<i>MYB</i>

Supplementary Table 1: (continued)

JASPAR identifier	HGNC symbol
MYC-MAX	<i>MYC</i>
Mycn	<i>MYCN</i>
Myf	<i>MYOD1</i>
NF-kappaB	<i>NFKB1</i>
NFIL3	<i>NFIL3</i>
NFKB1	<i>NFKB1</i>
NHLH1	<i>NHLH1</i>
Nkx2-5	<i>NKX2-5</i>
NR1H2-RXR	<i>NR1H2</i>
NR2F1	<i>NR2F1</i>
NR3C1	<i>NR3C1</i>
Pax2	<i>PAX2</i>
Pax4	<i>PAX4</i>
Pax5	<i>PAX5</i>
Pax6	<i>PAX6</i>
Pbx	<i>PBX1</i>
PPARG	<i>PPARG</i>
PPARG-RXRA	<i>RXRA</i>
Prrx2	<i>PRRX2</i>
REL	<i>REL</i>
RELA	<i>RELA</i>
Roaz	<i>ZNF423</i>
RORA	<i>RORA</i>
RORA1	<i>RORA</i>
RREB1	<i>RREB1</i>
RUNX1	<i>RUNX1</i>
RUSH1-alfa	<i>GATA4</i>
RXR-VDR	<i>VDR</i>
Sox17	<i>SOX17</i>
Sox5	<i>SOX5</i>
SOX9	<i>SOX9</i>
SP1	<i>SP1</i>
SPI1	<i>SPI1</i>
SPIB	<i>SPIB</i>
Spz1	<i>SPZ1</i>
SRF	<i>SRF</i>
SRY	<i>SRY</i>
Staf	<i>ZNF143</i>
T	<i>T</i>
TAL1-TCF3	<i>TCF3</i>
TBP	<i>TBP</i>
TCF1	<i>HNF1A</i>

Supplementary Table 1: (continued)

JASPAR identifier	HGNC symbol
TCF11-MafG	<i>NFE2L1</i>
TEAD	<i>TEAD1</i>
TFAP2A	<i>TFAP2A</i>
TP53	<i>TP53</i>
USF1	<i>USF1</i>
YY1	<i>YY1</i>
ZNF42_1-4	<i>MZF1</i>
ZNF42_5-13	<i>MZF1</i>

GO term	low $T$		not in JASPAR PWMs	
	$p$	$q$	$p$	$q$
immune response	$4 \times 10^{-22}$	$< 1 \times 10^{-4}$	0.5	1
antigen processing and presentation	$3 \times 10^{-17}$	$< 1 \times 10^{-4}$	0.4	1
defense response	$1 \times 10^{-16}$	$< 1 \times 10^{-4}$	0.9	1
immune system process	$4 \times 10^{-15}$	$< 1 \times 10^{-4}$	1	1
sensory perception of smell	$6 \times 10^{-15}$	$< 1 \times 10^{-4}$	0.3	1
sensory perception of chemical stimulus	$4 \times 10^{-14}$	$< 1 \times 10^{-4}$	0.3	1
response to stimulus	$6 \times 10^{-13}$	$< 1 \times 10^{-4}$	1	1
antigen processing and presentation of peptide or polysaccharide antigen via MHC class II	$7 \times 10^{-10}$	$< 1 \times 10^{-4}$	0.8	1
inflammatory response	$2 \times 10^{-9}$	$< 1 \times 10^{-4}$	1	1
RNA processing	$6 \times 10^{-8}$	$< 1 \times 10^{-4}$	0.2	1
response to other organism	$5 \times 10^{-7}$	$2 \times 10^{-4}$	0.5	1
response to wounding	$7 \times 10^{-7}$	$1 \times 10^{-4}$	0.8	1
multi-organism process	$9 \times 10^{-7}$	$1 \times 10^{-4}$	0.3	1

Supplementary Table 2: GO biological process terms enriched in genes with low- $T$  TSS-flanking regions (Wilcoxon rank sum test;  $p < 10^{-6}$ ). The global analysis is significant (Kolmogorov-Smirnov-like test;  $p = 6 \times 10^{-3}$ ). For each term the significance of its enrichment in genes not in the JASPAR PWM set is also shown.

GO term	low $T$		not in JASPAR PWMs	
	$p$	$q$	$p$	$q$
olfactory receptor activity	$5 \times 10^{-17}$	$< 1 \times 10^{-4}$	0.2	1
chemokine receptor binding	$7 \times 10^{-8}$	$< 1 \times 10^{-4}$	0.8	1
MHC class II receptor activity	$1 \times 10^{-7}$	$< 1 \times 10^{-4}$	1	1
chemokine activity	$3 \times 10^{-7}$	$< 1 \times 10^{-4}$	0.8	1
G-protein-coupled receptor binding	$7 \times 10^{-7}$	$< 1 \times 10^{-4}$	0.7	1

Supplementary Table 3: GO molecular function terms enriched in genes with low- $T$  TSS-flanking regions (Wilcoxon rank sum test;  $p < 10^{-6}$ ). The global analysis is significant (Kolmogorov-Smirnov-like test;  $p < 10^{-4}$ ). For each term the significance of its enrichment in genes not in the JASPAR PWM set is also shown.

GO term	low $T$		not in JASPAR PWMs	
	$p$	$q$	$p$	$q$
MHC protein complex	$1 \times 10^{-17}$	$< 1 \times 10^{-4}$	0.5	1
MHC class II protein complex	$3 \times 10^{-13}$	$< 1 \times 10^{-4}$	0.7	1
ribonucleoprotein complex	$3 \times 10^{-9}$	$< 1 \times 10^{-4}$	0.07	0.6

Supplementary Table 4: GO cellular component terms enriched in genes with low- $T$  TSS-flanking regions (Wilcoxon rank sum test;  $p < 10^{-6}$ ). The global analysis is significant (Kolmogorov-Smirnov-like test;  $p = 3 \times 10^{-3}$ ). For each term the significance of its enrichment in genes not in the JASPAR PWM set is also shown.

GO term	high $T$		in JASPAR PWMs	
	$p$	$q$	$p$	$q$
post-translational protein modification	$5 \times 10^{-15}$	$< 1 \times 10^{-4}$	1	1
protein modification process	$9 \times 10^{-14}$	$< 1 \times 10^{-4}$	1	1
biopolymer modification	$4 \times 10^{-13}$	$< 1 \times 10^{-4}$	1	1
protein amino acid phosphorylation	$5 \times 10^{-12}$	$< 1 \times 10^{-4}$	1	1
phosphorus metabolic process	$4 \times 10^{-10}$	$< 1 \times 10^{-4}$	1	1
phosphate metabolic process	$4 \times 10^{-10}$	$< 1 \times 10^{-4}$	1	1
intracellular signaling cascade	$3 \times 10^{-9}$	$< 1 \times 10^{-4}$	0.8	1
phosphorylation	$8 \times 10^{-9}$	$< 1 \times 10^{-4}$	1	1
cell communication	$2 \times 10^{-8}$	$< 1 \times 10^{-4}$	0.6	1
small GTPase mediated signal transduction	$8 \times 10^{-8}$	$< 1 \times 10^{-4}$	1	1
organ morphogenesis	$1 \times 10^{-7}$	$< 1 \times 10^{-4}$	$6 \times 10^{-14}$	$< 1 \times 10^{-4}$
signal transduction	$1 \times 10^{-7}$	$< 1 \times 10^{-4}$	0.8	1
localization	$2 \times 10^{-7}$	$< 1 \times 10^{-4}$	0.9	1
anatomical structure morphogenesis	$2 \times 10^{-7}$	$< 1 \times 10^{-4}$	$3 \times 10^{-16}$	$< 1 \times 10^{-4}$
developmental process	$2 \times 10^{-7}$	$< 1 \times 10^{-4}$	$3 \times 10^{-22}$	$< 1 \times 10^{-4}$
biological regulation	$3 \times 10^{-7}$	$< 1 \times 10^{-4}$	$1 \times 10^{-26}$	$< 1 \times 10^{-4}$
endocytosis	$4 \times 10^{-7}$	$< 1 \times 10^{-4}$	0.6	1
membrane invagination	$4 \times 10^{-7}$	$< 1 \times 10^{-4}$	0.6	1
regulation of biological process	$4 \times 10^{-7}$	$< 1 \times 10^{-4}$	$6 \times 10^{-29}$	$< 1 \times 10^{-4}$
regulation of cellular process	$7 \times 10^{-7}$	$< 1 \times 10^{-4}$	$1 \times 10^{-29}$	$< 1 \times 10^{-4}$
anatomical structure development	$7 \times 10^{-7}$	$< 1 \times 10^{-4}$	$4 \times 10^{-19}$	$< 1 \times 10^{-4}$

Supplementary Table 5: GO biological process terms enriched in genes with high- $T$  TSS-flanking regions (Wilcoxon rank sum test;  $p < 10^{-6}$ ). The global analysis is significant (Kolmogorov-Smirnov-like test;  $p < 10^{-4}$ ). For each term the significance of its enrichment in genes in the JASPAR PWM set is also shown.



GO term	high $T$		in JASPAR PWMs	
	$p$	$q$	$p$	$q$
phosphotransferase activity, alcohol group as acceptor	$1 \times 10^{-14}$	$< 1 \times 10^{-4}$	1	1
kinase activity	$3 \times 10^{-14}$	$< 1 \times 10^{-4}$	1	1
transferase activity, transferring phosphorus-containing groups	$8 \times 10^{-13}$	$< 1 \times 10^{-4}$	1	1
protein kinase activity	$2 \times 10^{-12}$	$< 1 \times 10^{-4}$	1	1
protein serine/threonine kinase activity	$2 \times 10^{-12}$	$< 1 \times 10^{-4}$	1	1
protein tyrosine kinase activity	$3 \times 10^{-11}$	$< 1 \times 10^{-4}$	1	1
transferase activity	$2 \times 10^{-10}$	$< 1 \times 10^{-4}$	1	1
protein binding	$3 \times 10^{-8}$	$< 1 \times 10^{-4}$	$8 \times 10^{-4}$	0.01
magnesium ion binding	$3 \times 10^{-7}$	$< 1 \times 10^{-4}$	1	1

Supplementary Table 6: GO molecular function terms enriched in genes with high- $T$  TSS-flanking regions (Wilcoxon rank sum test;  $p < 10^{-6}$ ). The global analysis is significant (Kolmogorov-Smirnov-like test;  $p < 10^{-4}$ ). For each term the significance of its enrichment in genes in the JASPAR PWM set is also shown.

GO term	high $T$		in JASPAR PWMs	
	$p$	$q$	$p$	$q$
cytoplasm	$4 \times 10^{-12}$	$< 1 \times 10^{-4}$	1	1

Supplementary Table 7: GO cellular component terms enriched in genes with high- $T$  TSS-flanking regions (Wilcoxon rank sum test;  $p < 10^{-6}$ ). The global analysis is significant (Kolmogorov-Smirnov-like test;  $p < 10^{-4}$ ). For each term the significance of its enrichment in genes in the JASPAR PWM set is also shown.

GO term	low $\psi$		high number of gaps	
	$p$	$q$	$p$	$q$
extracellular region	$2 \times 10^{-17}$	$< 1 \times 10^{-4}$	0.4	1
extracellular region part	$5 \times 10^{-12}$	$< 1 \times 10^{-4}$	0.3	1
extracellular space	$6 \times 10^{-8}$	$< 1 \times 10^{-4}$	0.2	1
plasma membrane	$3 \times 10^{-7}$	$< 1 \times 10^{-4}$	1	1
plasma membrane part	$4 \times 10^{-7}$	$< 1 \times 10^{-4}$	1	1
extracellular matrix	$5 \times 10^{-7}$	$< 1 \times 10^{-4}$	0.6	1
proteinaceous extracellular matrix	$7 \times 10^{-7}$	$< 1 \times 10^{-4}$	0.6	1

Supplementary Table 8: GO cellular component terms enriched in genes with low- $\psi$  TSS-flanking regions (Wilcoxon rank sum test;  $p < 10^{-6}$ ). The global analysis is significant (Kolmogorov-Smirnov-like test;  $p = 2 \times 10^{-3}$ ). For each term the significance of its enrichment in genes with a high number of gaps is also shown.

GO term	high $\psi$		high number of gaps	
	$p$	$q$	$p$	$q$
cytoplasm	$< 2 \times 10^{-308}$	$< 1 \times 10^{-4}$	0.3	1
cytoplasmic part	$8 \times 10^{-16}$	$< 1 \times 10^{-4}$	0.6	1
mitochondrion	$5 \times 10^{-15}$	$< 1 \times 10^{-4}$	0.3	1
intracellular	$2 \times 10^{-12}$	$< 1 \times 10^{-4}$	1	1
intracellular part	$2 \times 10^{-11}$	$< 1 \times 10^{-4}$	1	1
organelle membrane	$1 \times 10^{-8}$	$< 1 \times 10^{-4}$	1	1
cell part	$6 \times 10^{-8}$	$< 1 \times 10^{-4}$	0.8	1
cell	$8 \times 10^{-8}$	$< 1 \times 10^{-4}$	0.8	1
ribonucleoprotein complex	$2 \times 10^{-7}$	$< 1 \times 10^{-4}$	1	1
organelle part	$9 \times 10^{-7}$	$< 1 \times 10^{-4}$	1	1

Supplementary Table 9: GO cellular component terms enriched in genes with high- $\psi$  TSS-flanking regions (Wilcoxon rank sum test;  $p < 10^{-6}$ ). The global analysis is significant (Kolmogorov-Smirnov-like test;  $p < 10^{-4}$ ). For each term the significance of its enrichment in genes with a high number of gaps is also shown.