

Bits and pieces about transcription start sites

Hedi Peterson

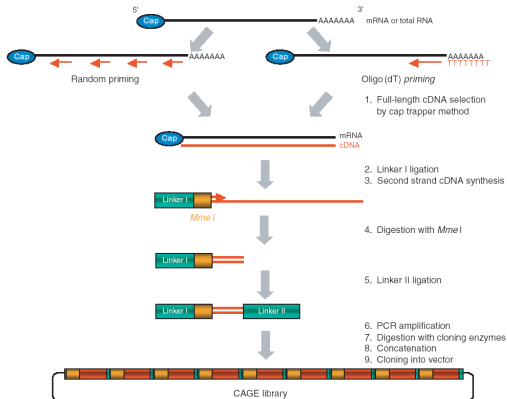
17.11.2006
JClub

Transcription start site

- Position on DNA where transcription starts - RNA polymerase (II) starts mRNA synthesis
- Transcription is triggered by transcription factors that bind to short DNA motifs upstream from the TSS
- Besides TFBS there are other common regulative elements - TATA-box, CCAAT-box, Gc-box and CpG islands

CAGE - what is it?

- **Cap Analysis Gene Expression**
- produces 20-21bp long tags



- produces 20-21bp long tags from 5' end of full-length cDNA
- tags will be mapped to corresponding genome
- overlapping tags from same strand form a tag cluster (TC)
- TC + surroundings form a core promoter

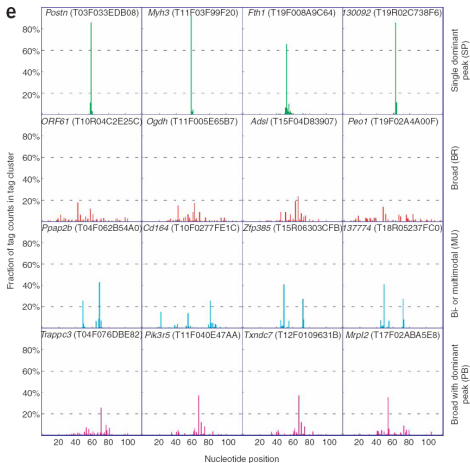
Human and mouse data

- 209(145) cDNA libraries for 23 tissues for mouse
- 43(41) libraries for human
- identified 729504 mouse and 665278 human potential TSSs (80-95% belong to TC)

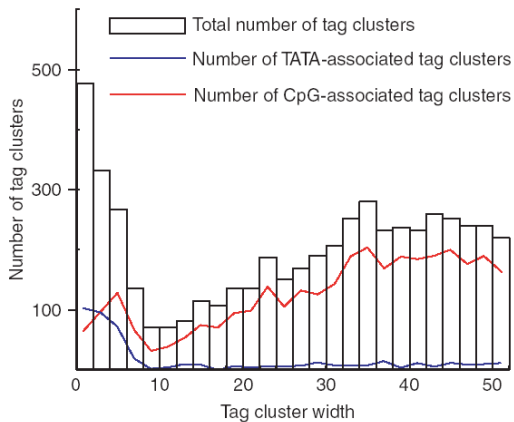
CAGE tag clusters

Four type of tag clusters
(based on 8000 mouse
and 6000 human TC, 100
tags per TC)

- single peak
- broad shape
- bi/multimodal
- broad with dominant peak

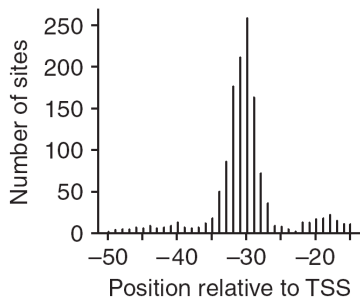


CAGE tag clusters



TATA-box

- strongly overrepresented in sharp TSS promoters
- highly conserved promoter region across species
- TATA-independent transc. init. occurs within a CpG island (90%)
- Sp1 recruits TBP when TATA-box is not present (in broad promoters)
- tissue specificity when PB promoters + strong TATA peak



- strongly associated with broad TSS regions
- more G nucleotides on "+" strand (could indicate promoter orientation)
- associated to house-keeping genes
- associated to bidirectional promoter activity
- central nervous system-specific promoters are especially CpG rich

- CCAAT- and GC-box associated with sharp peak TSS class
- secondary peak in 3'UTR with Inr GGG

Table 2 Overrepresentation and underrepresentation of transcriptional starting site sequence

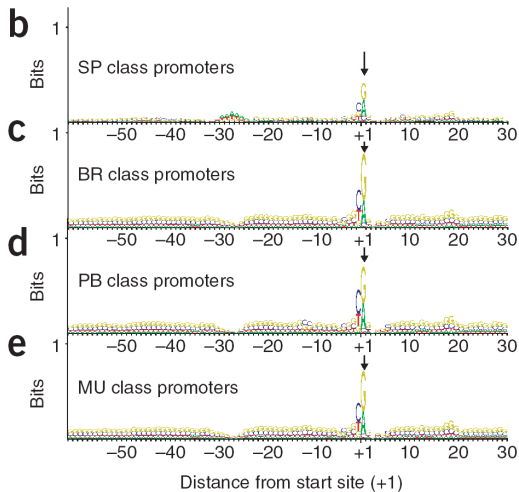
Overall analysis	SP	BR	PB	MU
TATA (all)	3.1×10^{-73}	1.9×10^{-16}	1.8×10^{-10}	2.4×10^{-9}
CCAAT (all)	0.04	0.42	0.37	0.49
GC (all)	1×10^{-4}	0.20	0.40	0.33
CpG (all)	1.0×10^{-137}	1.4×10^{-65}	8.7×10^{-6}	0.02
CpG promoters versus non-CpG promoters	SP	BR	PB	MU
TATA (no CpG)	2.6×10^{-77}	1.6×10^{-16}	2.8×10^{-16}	1.0×10^{-9}
CCAAT (no CpG)	6.8×10^{-23}	9.2×10^{-16}	0.11	0.42
GC (no CpG)	7.8×10^{-25}	5.9×10^{-18}	0.48	0.35
CpG (no TATA, CCAAT or GC)	4.8×10^{-45}	4.7×10^{-17}	3.4×10^{-5}	0.87

For each shape class, we determined whether a TATA box (within 50 bp) or a CCAAT, GC or CpG (within 200 bp) upstream of the start site of the clusters was present. *P* values were determined using the Fisher exact test (**Supplementary Note**). *P* values in boldface and italics indicate significant underrepresentation ($P < 0.01$); *P* values in boldface alone indicate significant overrepresentation ($P < 0.01$). In the lower part of the table, we separated pure CpG-island-overlapping promoters (without TATA, CCAAT and GC elements) and TATA, CCAAT and GC promoters (without CpG islands).

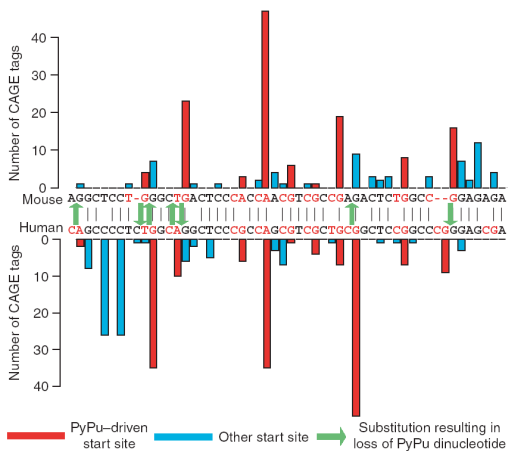
Initiator sequence (Inr element)

- previously known Inr element -> Py-Py-A-N-T/A-Py-Py
- Transcription preferentially starts (+1 pos) with a purine
- +1 pos preference for pyrimidine (A/G) (58.6%)
- CA for highly exp. transcripts
- AG, GG for rarely exp. transcripts

Inr elements for TC types



Pyrimidine/Purine substitutions



Promoter evolution

- Within 200bp upstream of TSS identity percentage declines
- TATA-containing promoters had a lower substitution rate than other 3 TC types
- human broad shape promoters have higher substitution rate than mouse promoter regions
- PyPu substitutions change the promoter expression
- Direct correlation between the global properties of promoters and TFBS density
- 3' UTR TSSs are at least some cases functional

TSS distances between human and mouse

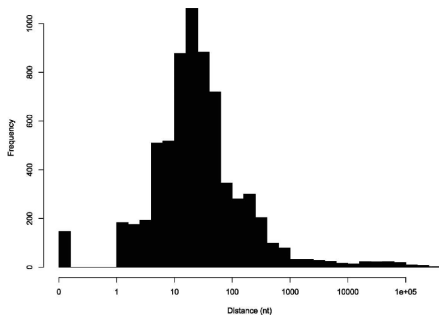
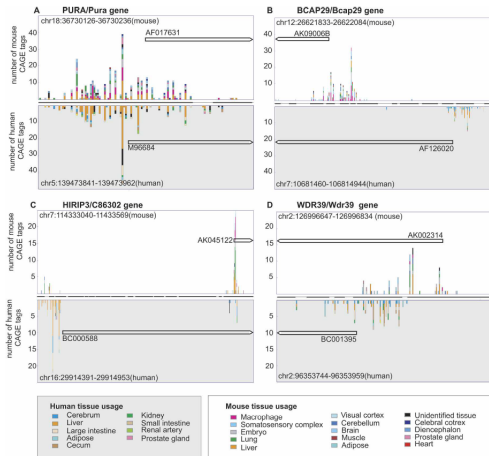


Figure 1. Histogram of distances between transcription start sites of homologous transcripts. The x-axis indicates the distance between the human TSS and the human position aligned to the mouse TSS.

TSS turnover

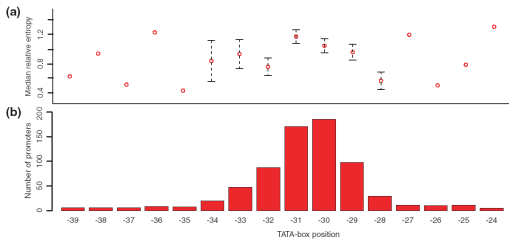


3'UTR promoters

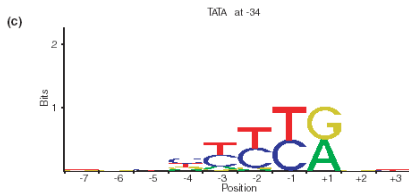
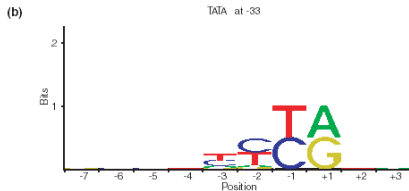
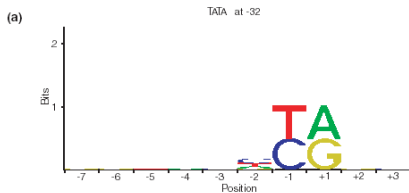
- Specific motif - 3 guanines just before TSSs.
- Highly conserved region +40..+90bp
- Proposed sense-antisense regulatory mechanism for downstream genes
- 3' TSS is rather tissue specific and prevalent in cerebellum and lung, reduced in embryo

- Tissue specificity is achieved by having TATA-box -32..-29bp upstream of TSS
- positions -31 and -30 are optimal
- TATA-box motifs closer than -28bp are generally non-functional
- TATA-box at position -34 has consensus motif TATATAA (other positions have consensus TATAAA), so -34 can also be -32 TATA-box

TATA-box distribution

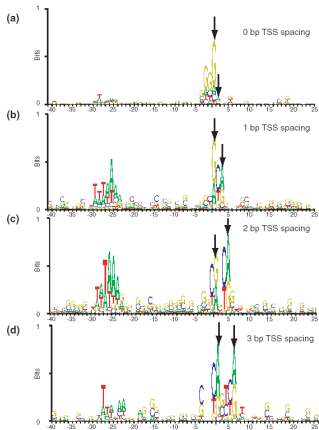


HMMs for upstream Inr



- Distance between two major TSS peaks are less than 5
- If the distance is 0 then specific Inr (GGG) is preferred

Twin TSSs



- 58% protein-coding transcriptional units use alternative promoters (previously known 20%)
- 93% transcriptional units, having at least 2 alternative promoters, have distinct methionine start codons
- two positions (-31,-30) for TATA-box are very specific for highly expressed genes

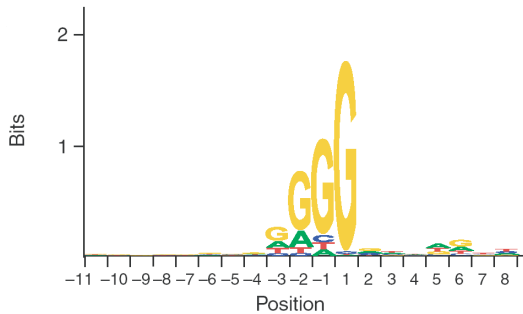
- Frith *et al* "**Evolutionary turnover of mammalian transcription start sites**", Genome Research, apr 2006
- Carninci *et al* "**Genome-wide analysis of mammalian promoter architecture and evolution**", Nature Genetics, jun 2006
- Ponjavic *et al* "**Transcriptional and structural impact of TATA-initiation site spacing in mammalian core promoters**", Genome Biology, aug 2006

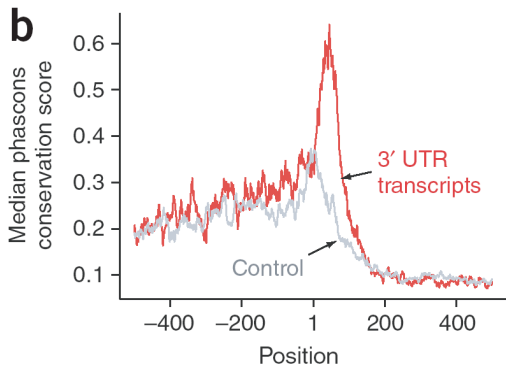
Questions?

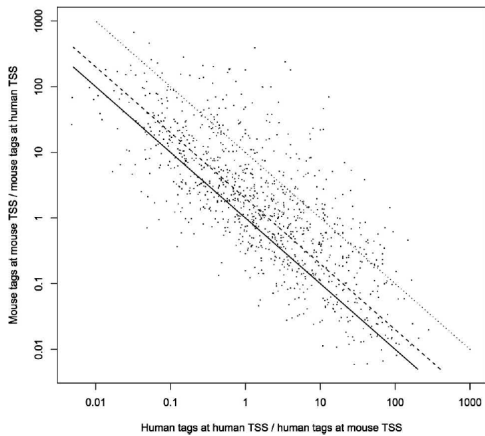
- Any questions?

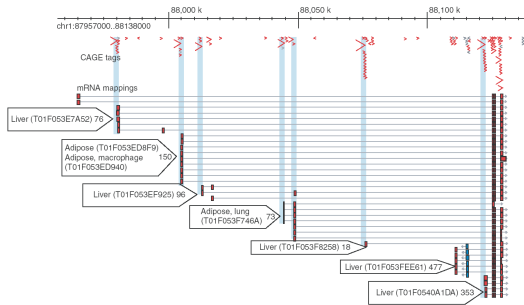


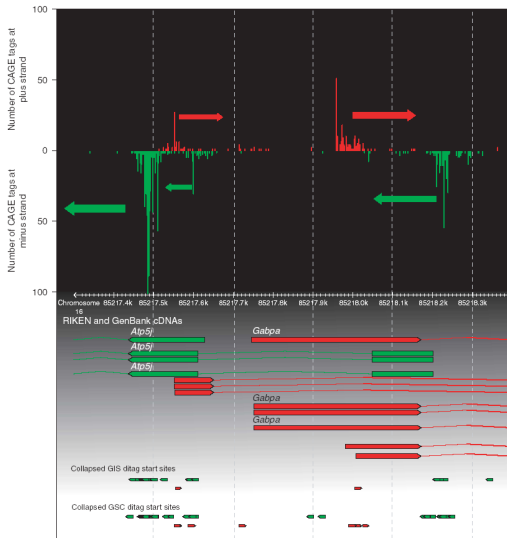


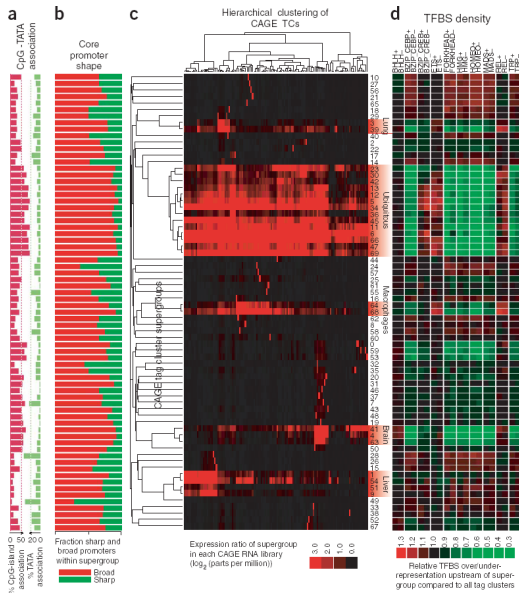


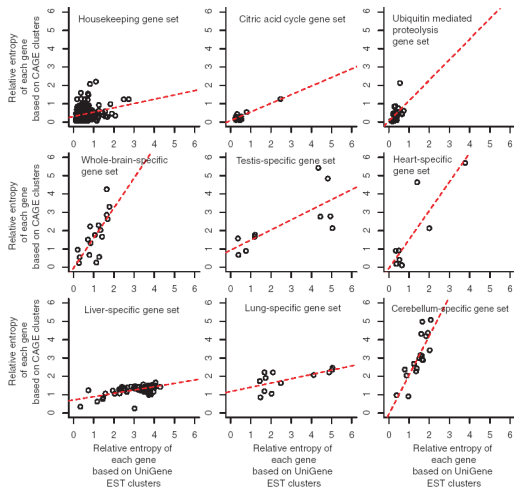


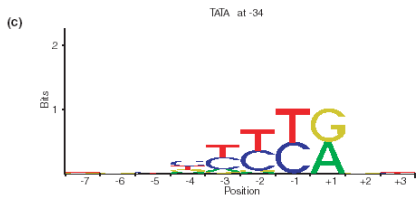
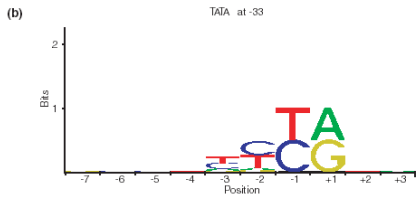
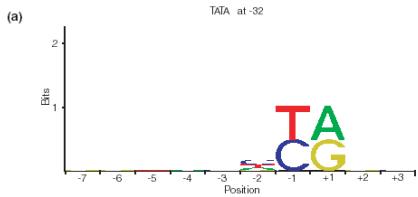


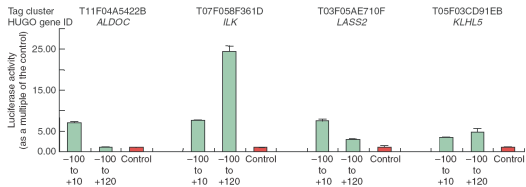


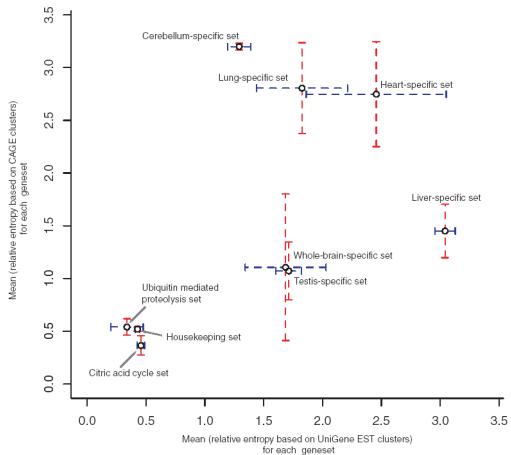












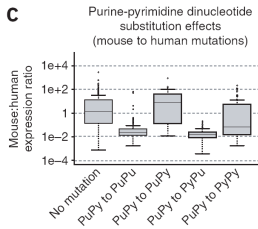
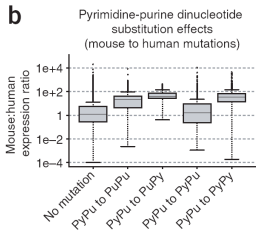
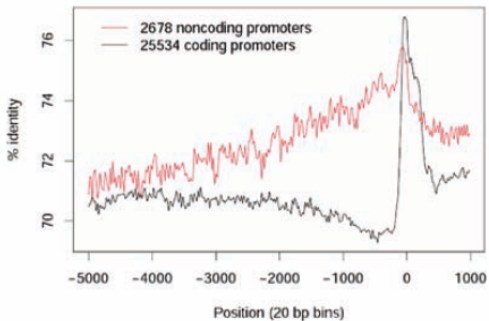


Table 1. Evolution of PyPu initiator dinucleotides at CAGE tag start sites in promoters with and without TSS turnover

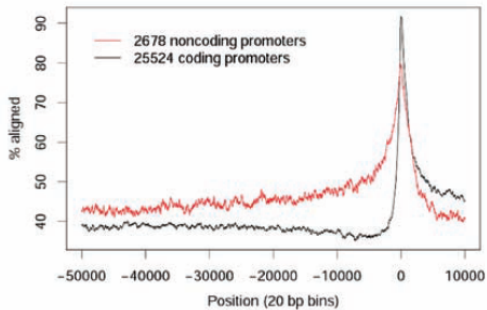
Type of change	PyPu→PuPy	PyPu→PuPu	PyPu→PyPy	PyPu→PyPu
Cases in turnover promoters (%)	26 1.95	194 14.56	189 14.19	923 69.29
Cases in reference promoters (%)	284 1.22	2544 10.92	2431 10.46	18,029 77.42
Turnover rate/reference rate	1.60	1.33	1.36	0.90
Significance of difference ^a	3.03×10^{-2}	7.76×10^{-5}	3.12×10^{-5}	3.04×10^{-11}

^aFisher's exact two-tail test.

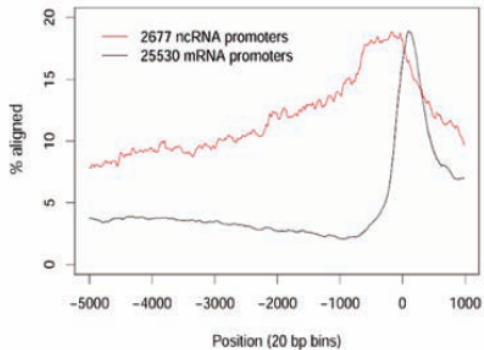
Sequence conservation of mouse promoters vs. human

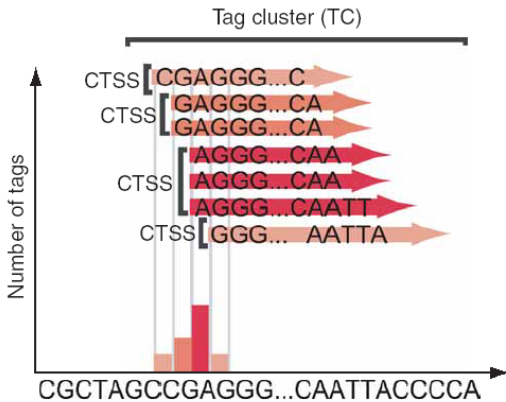


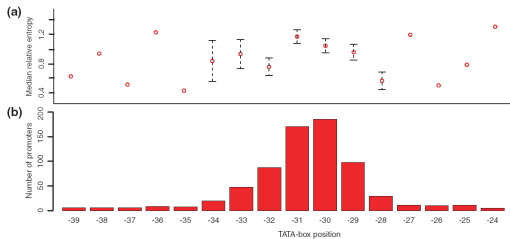
C Sequence conservation of mouse promoters vs. human

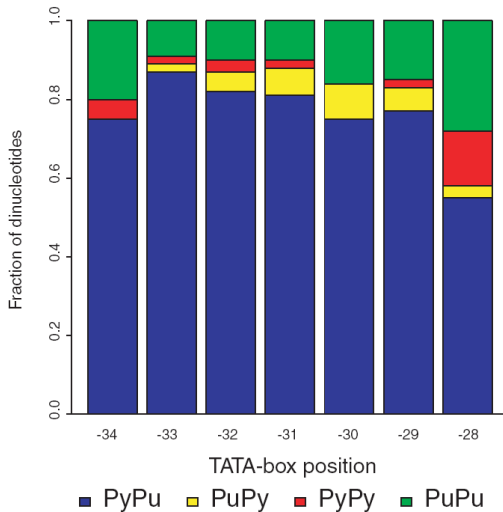


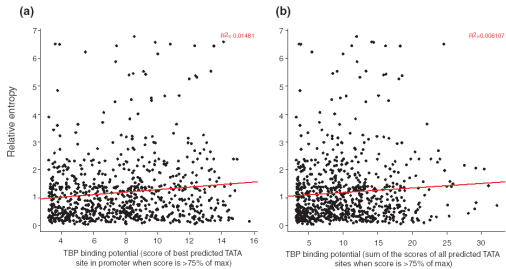
Sequence conservation of mouse promoters vs. chicken

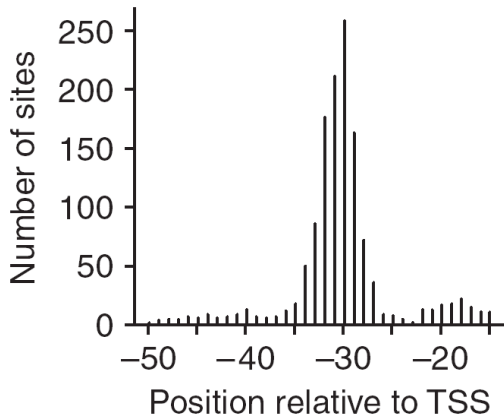


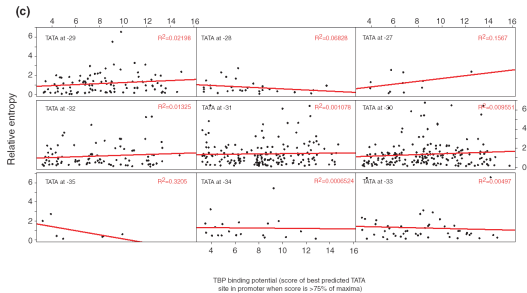




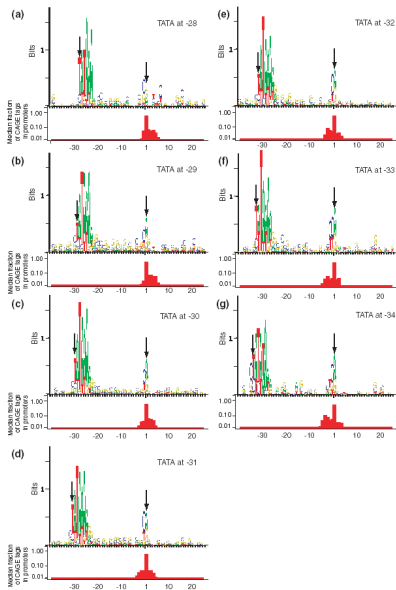


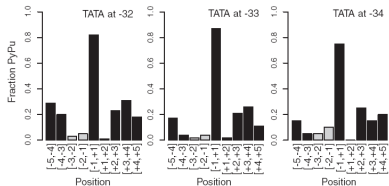
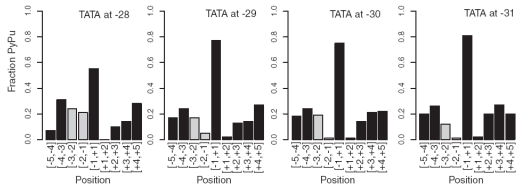


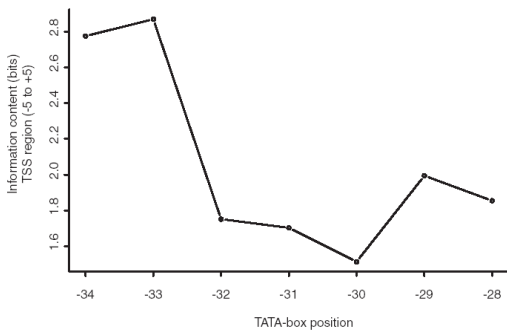




tataseqlogo

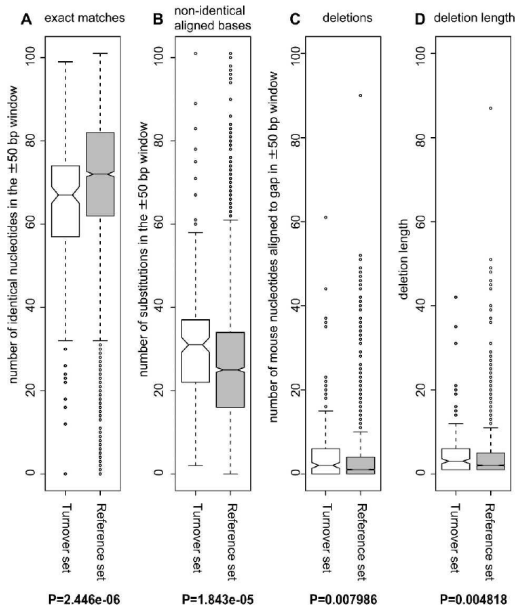






Correlation of tissue specificities measured by relative entropy in CAGE and UniGene EST clusters

Gene set	EST versus CAGE Spearman rank correlation coefficient	Spearman rank correlation P value	Number of genes
Whole brain specific	216	1.10×10^{-3}	17
Testis specific	48	9.68×10^{-2}	9
Heart specific	40	1.48×10^{-2}	10
Liver specific	20,898	1.32×10^{-6}	66
Lung specific	92	1.81×10^{-2}	12
Cerebellum specific	186	$<2.20 \times 10^{-16}$	20
Citric acid cycle	318	2.90×10^{-1}	14
Ubiquitin-mediated proteolysis pathway	886	5.94×10^{-3}	23
Housekeeping genes	2,208,352	8.54×10^{-6}	263
All sets combined	5,269,164	$<2.20 \times 10^{-16}$	434



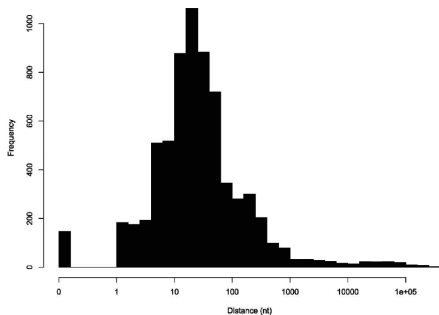
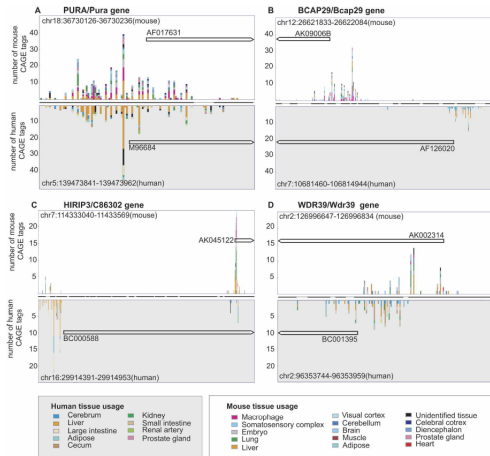
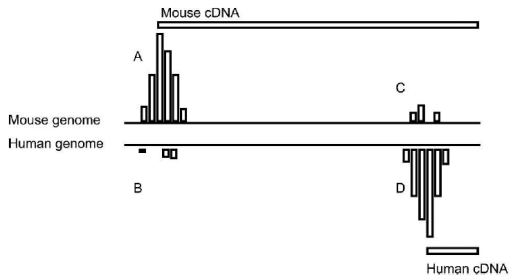
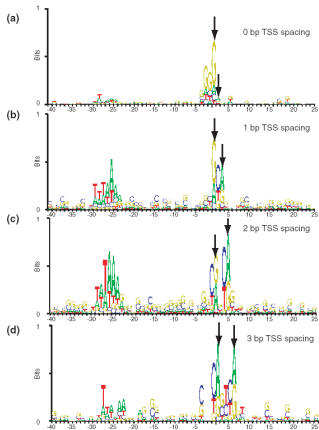


Figure 1. Histogram of distances between transcription start sites of homologous transcripts. The x-axis indicates the distance between the human TSS and the human position aligned to the mouse TSS.









Numbers





TSS evolution/turnover

































