Lecture 5 - Promoters and transcription

contatto per studentessa cagliari

#### Analysis of transcriptome

-all kind of organisms from S.cerevisiae to human cells -an armamentary of methods, from EST to microarrays to RNA-Seq

have led to the concept that the genome is transcribed much more than previously thought

This phenomenon was named "pervasive transcription".

Most of "unusual" transcript are present at very low level and they are at least in part questioned. However, there are many RNA that have been characterized.

Conceptually, these RNA should be classified as:

- a) Related to genes (promoters, intronic, antisense)
- b) unrelated to genes (intergenic)
  - 1) stable
  - 2) unstable
    - I. coding
    - II. noncoding

Common problem: How to distinguish "transcripts" from RNA fragments or processed RNAs ?

in addition

RNA transcripts (even capped) may arise from unexpected phenomena, such as cleavage of longer transcript, *trans*-splicing followed by processing, editing, and even by RNA-dependent RNA polymerization (even though such a polymerase, found in lower eukaryotes (e.g *C. elegans*), has never been observed in Vertebrates).

Other functional annotations may help, such as:

> the presence of RNA polymerase(s)

the presence of other proteins necessary for transcription (and RNA processing)
histone modifications that detect the local functional status of the chromatin
DNA covalent modifications (e.g. CpG methylation)

≻....

Il problema della corretta mappatura delle T.U.

Definizione del 5' e 3' (CAGE, SAGE, RNA-Seq)

Definizione esatta degli esoni (Tiling arrays, RNA-Seq)

Definizione del promotore (CAGE, RNA-Seq + studies)

Promoters and transcription

- what we know from classical books

-How to define a TSS ?

- gene by gene methods
- genomic methods
- One study for the identification of promoters genome-wide
- Classification of non protein-coding RNAs (intragenic and intergenic transcripts)

# Mammalian RNA polymerase II core promoters: insights from genome-wide studies

## Review 1

Albin Sandelin\*<sup>‡</sup>, Piero Carninci<sup>‡</sup>, Boris Lenhard<sup>11</sup>, Jasmina Ponjavic<sup>¶</sup>, Yoshihide Hayashizaki<sup>‡</sup> and David A. Hume<sup>#</sup>

Abstract | The identification and characterization of mammalian core promoters and transcription start sites is a prerequisite to understanding how RNA polymerase II transcription is controlled. New experimental technologies have enabled genome-wide discovery and characterization of core promoters, revealing that most mammalian genes do not conform to the simple model in which a TATA box directs transcription from a single defined nucleotide position. In fact, most genes have multiple promoters, within which there are multiple start sites, and alternative promoter usage generates diversity and complexity in the mammalian transcriptome and proteome. Promoters can be described by their start site usage distribution, which is coupled to the occurrence of *cis*-regulatory elements, gene function and evolutionary constraints. A comprehensive survey of mammalian promoters is a major step towards describing and understanding transcriptional control networks.

424 JUNE 2007 VOLUME 8

© 2007 Nature Publishing Group

www.nature.com/reviews/genetics

How many important concepts of eukaryotic transcription are resumed in this figure ?





### table 26-1

#### Proteins Required for Transcription at the RNA Polymerase II Promoters of Eukaryotes

Transcription factor	Number of subunits	Subunit <i>M</i> ,	Functions
Initiation			
RNA polymerase II	12	10,000-220,000	Catalyzes RNA synthesis
TBP (TATA-binding protein)	1	38,000	Specifically recognizes the TATA box
TFIIA	3	12,000, 19,000, 35,000	Stabilizes binding of TFIIB and TBP to the promoter
TFIIB	1	35.000	Binds to TBP: recruits RNA polymerase-TFIIF complex
TFIID	12	15,000-250,000	Interacts with positive and negative regulatory proteins
TFIIE	2	34,000, 57,000	Recruits TFIIH; ATPase and helicase activities
TFIIF	2	30,000, 74,000	Binds tightly to RNA polymerase II; binds to TFIIB and prevents binding of RNA polymerase to nonspecific DNA sequences
TFIIH	12	35,000-89,000	Unwinds DNA at promoter; phosphorylates RNA polymerase; recruits nucleotide-excision repair
Elongation*			complex
ELL <sup>†</sup>	1	80,000	
P-TEFb	2	43,000, 124,000	
SII (TFIIS)	1	38,000	
Elongin (SIII)	3	15,000, 18,000, 110,000	

\*All elongation factors suppress the pausing or arrest of transcription by the RNA polymerase II – TFIIF complex.

<sup>†</sup>The name is derived from the term *e*leven-nineteen *l*ysine-rich *l*eukemia. The gene for the factor ELL is the site of chromosomal recombination events frequently associated with the cancerous condition known as acute myeloid leukemia.



The "textbook" promoter

#### "consensus" sequences are found by computing the frequencies of bases at any position in several genes. Example:

gene 1:	a	t	t	g	t	С	t	а		
gene 2:	a	t	а	g	t	g	а	g		
gene 3:	a	t	t	С	t	g	g	С		
gene 4:	a	t	t	g	t	С	С	g		
gene 5:	a	t	t	С	t	g	а	t		
gene 6:	a	t	a	g	g	С	t	С		
etc:	-	-	-	-	-	-	-	-		
frequency a	100	0	35		0	0	0	25	20	
t	0	100	65		0	90	0	25	20	
С	0	0	0		35	0	50	20	30	
a	0	0	0		65	10	50	30	30	
Consensus:	A	т	W	I	S	т	S	N	N	

Symbol	Name
А	Adenine
С	Cytosine
Н	Thymine
G	Guanine
U	Uracil
W	A or T
R	A or G
K	G or T
Y	C or T
S	C or G
Μ	A or C
В	C, G or T
Н	A, C, or T
D	A, G, or T
V	A, C, or G
N	A, C, G, or T

CC	CAAT	CRE						ATG/KOZAK						
47	690	CAATGGGA	11.3	50	484	TGACGTCA	18.4	52	9601	CCA	AGATG	7.5		
47	601	CAATCAGC	13.1	49	282	ATGACGTC	8.5	52	617!	GCA	AGATG	13.7		
46	708	CAATCAGA	14.4	50	503	CTGACGTC	9.3	51	543	GCG	CCATG	9.3		
45	310	CCAATCGG	8.1	48	635	GTGACGTC	GTGACGTC 13.5 53 6881 GCA		CCATG	9.9				
46	871	CCAATCCC	8.0	50	313	GTGACGCA	7.4	52	1152	CAG	CCATG	11.1		
48	620	CCRATCAC	11.7	49	345	AGTGACGT	9.4	53	1005	CA	CCATGG	8.6		
47	1041	CCRRTCAC	23 6	49	294	COTCACOC	8.0	52	426	CO	CCATCC	9.0		
47	206	CONNECCO	13 3	40	280	COTCACOT	10.2	60	931	00	CCATCC	9.4		
42	300	COMMICOC	21 1	49	200	COTCACCT	2 1	52	10011	00	ACATOC	20 6		
47	000	CCAATOGG	34.1	50	3/9	GGIGAGGI	11 4	61	1010	0.0	ADATOG	0.0		
4/	0.00	GUCAATCA	44.5	40	204	DOCTOROGI	10.2	51	1010		AAATGG	16.0		
40	361	GCCAATAG	9.2	49	241	ACGTGACG	10.3	21	1048	5	AAATGGC	10.0		
48	357	GCCAATCG	12.4	49	472	ACGTGACC	8.4	52	12021	P	AGATGGC	36.9		
46	578	GCCAATGA	17.0	49	332	AAGTGACG	23.9	52	8811		AGATGGCG	40.2		
47	775	GCCAATGG	26.8	50	769	GAAGTGAC	10.4	51	654		ACATGGCG	13.5		
46	553	GCCAATCC	9.7					51	436		AAATGGCG	23.4		
47	537	TCCAATCA	7.0	TA	TA			51	414		AATGGCGG	12.9		
47	220	ACCAATCG	14.7	49	486	CCTATAAA 9.		52	1026!		GATGGCGG	27.2		
47	469	ACCAATGG	17.8	49	571	GCTATAAA	7.1	52	920		CATGGCGG	18.4		
46	583	ACCAATCA	17.4	49	496	CTATAAAG	10.1	54	291	CATGGCGT		11.1		
47	384	GACCAATG	9.8	49	809	TTATAAAG	10.9	51	583	ATGGCGCC		23.6		
47	400	GACCAATC	19.2	49	8611	TATABAAG	11.4	52	1125!	ATGGCGGC		27.7		
47	893	AGCCAATC	19.3	49	417	TATATAAG	9.7	52	619		ATGGCGGG	8.2		
46	748	AGCCAATG	13.8	10	5421	TATABACC	28.0	52	4681		ATGGCGGA	16.0		
47	680	GGCCAATG	11.7	42	9601	202322200	12.1	52	966		ATGGCTGC	15.8		
48	658	GGCCAATC	24.0	49	620	MINAAAGG	0.0	54	7911		CCAGGTAA	7.1		
47	547	CROCCART	10.2	49	020	TANAAGGC	3.9	56	3071		GCAGGTA	8.2		
47	324	CCACCANT	8 8	3995	12107			51	443	1	GCAGTCT	R 1		
47	403	CCCCCANT	12.4	NR	(F-1			55	16391	G	TGAGTG	7.5		
10	403	OGGCCAA1	10.0	50	1240	TGCGCCTG	11.9	6.2	0401	TO	TCACT	7 9		
48	1009	CGGCCAAT	10.9	50	2300	GCGCCTGC	12.3	ುತ	040:	190	STGAGI	1.2		
47	1039	CAGCCAAT	31.4	50	1767	CGCCTGCG	11.6	6.9	1414	~	CACOMO	7 4		
47	1/4	TCAGCCAA	10.5	50	463	CGCATGCG	15.5	52	20071	Gr	CONCERNICE OF	0.1		
47	1036	GCAGCCAA	7.0	50	2154	GCCTGCGC	7.8	53	366/1		CIGCIGCI	9.1		
				48	1205	GCGTGCGC	7.4	23	30101		TGCIGUIG	0.0		
SF	21			50	1041	CCTGCGCA	12.9							
48	1332	GCCACGCC	15.7	50	903	ACTGCGCC	8.0							
48	8136	GCCCCGCC	25.2	50	572	TGCGCATG	8.5							
48	3078	CGCCCCTC	7.3	49	386	CGCGCATG	11.1							
48	5248	CGCCCCGC	13.7	50	1179	CCCCATCC	18 5	Highlighted Sequences						
48	3141	CGCCCCCT	7.4	-00	44/2	000001100	10.00	CCAAT						
48	7055	CCGCCCCC	18.1					442-49 10 0 PRCCAATER			1170			
47	2106	CCGCCCAC	8.1					10.74	6-42 4	0.0	RECENTION	4470		
4.9	5783	CCCCCTCC	7.0	EI	S				12.2.0					
47	5204	cccccccc	15.5	49	1546	AGGAAGTG	7.6	SF	1					
40	3600	CCCCCCCC	12 6	49	1324	GGAAGTGA	16.2	*4	4-50	8.8	CCCCGCCC	3424		
40	10767	00000000	20 2	49	923	GGAAGTGC	11.9	*4	4-50	8.3	SCCCCCCCC	2687		
90	10/0/	NORCOLL	15 4	50	1892	GGAAGTGG	7.5	4	4-50	8.7	CCCGCCCC	2257		
48	1170	ACGCCCCC	15.4	49	284	CGGAAGTA	23.1							
48	829	ACGCCCCG	1.9	50	484	CGGAAGCA	13.B	Ch	us1					
48	1639	CACGCCCC	13.9	50	426	CGGAAGTC	24.8	+4	2-51	0 2	TOTOTOT	211		
48	2890	CCCGCCCT	8.9	51	402	CGGAAGTT	8.0		2-21	0.4	TOTOGOGA	211		
47	2334	CCCGCCCA	10.8	50	991	COGAAGTG	29.5							
48	2462	TCCGCCCC	8.4	51	356	CCCABATC	7.9	05	SF					
48	4767	CCCGCCTC	18.8	49	567	CCCAACCT	R /	*4	6-50	7.3	CCACGTGA	123		
48	3366	CTCCGCCC	11.8	49	004	COGRAGOI	10.7	*4	7-51	7.6	TCACGTGA	89		
48	11029	CCCCGCCC	31.3	00	1150	CGGRAGCG	19.2							
48	3190	CCCCGCCT	12.5	49	1150	CCGGAAGC	20.9	CF	RE					
49	918	TTCCGCCC	17.8	50	1030	CCGGAAGT	31.9		5-50	9 E	TCACCRCA	190		
48	2673	GCTCCGCC	7.2	51	459	CCGGAAAC	13.1	- 4	5-50	5.5	TGACGICA	125		
49	1213	CTTCCGCC	7.9	50	600	ACCGGAAG	40.6 *45-51 5.1 TGAT		TOATGTCA	140				
49	4947	GGCCCCGC	7.1	50	1096	GCCGGAAG	23.2	2 46-50 7.1 TTGCGTCA		TTGCGTCA	48			
47	5130	CCTCCCTC	8 1	49	1224	CCCGGAAG	20.1							
40	7005	CCCCTCCC	7 4	51	603	ACCCGGAA	7.8	TA	ATA					
40	1303	00001000	1.1.4	50	382	CACCGGAA	12.9	*4	8-49!	7.7	TATAAAD	472		
				49	401	GACCGGAA	7.4	4	8-491	2.4	TATATAD	349		

FitzGerald PC et al. Clustering of DNA Sequences in Human Promoters. <u>Genome Res. 2004 August; 14(8): 1562–</u> <u>1574.</u>



9

One hundred fifty-six DNA sequences are grouped into related sequences and arranged by their peak position relative to the TSS. From the left the table contains the most abundant bin, the number of times the sequence occurs in the distribution, the 8-mer sequence, and finally the P value (see text). The end of the table contains consensus sequences. Here the leftmost numbers are the bins defining the peak, followed by the clustering factor (CF), the consensus sequence, and finally the number of occurrences of the sequence in the bins that comprise the peak. Exclamation point (!) denotes sequences that are at least threefold more abundant in the maximum bin on the DNA strand presented in the table than on the opposite strand. The asterisk (\*) denotes sequences used in Tables 2 and 3. IUPAC letters used to represent degenerate bases are R (G,A), W (A,T), Y (T,C), K (G,T), V (G, C, A), D (G,A,T), and N (A,T,G,C).

The "Initiator" element (INR) is ill-defined and much less conserved than the TATA-box

The INR was demonstrated by bioinformatics <u>but also</u> using direct interaction studies using isolated recombinant subunits of the TFIID protein, the so-called TAFs (TBP -associated factors).

In vitro transcripiton experiments using recombinant TBP, TFIIB and individual TAFs showed that TAF150+TAF250 (drosophila) can correctly initiate transcription from TATA-less promoters.

Question: do TAF150/TAF250 recognize INR? How is this element composed?



The **DPE** (downstream promoter element) was demonstrated by bioinformatic and functional analysis. TAFII40 and TAFII60 in reconstitution experiments were enough to sustain transcription initiation in TATA-less, DPE-containing promoters.



FIG. 4. The DPE appears to be present in many Drosophila promoters.

(A) The frequency of occurrence of the DPE appears to be comparable to that of the TATA box in *Drosophila* core promoters. A *Drosophila* core promoter database was created by aligning sequences of 205 *Drosophila* core promoters with accurately determined transcription start sites. The number of promoters that appear to possess a TATA box only, a DPE only, both elements, or neither element is shown. TATA boxes were defined as sequences with at least a 5 out of 6 match with the TATAAA sequence upstream of 220 relative to the transcription start site. DPE motifs were defined as sequences with at least a 5 out of 6 match with the start site. The *Drosophila* core promoter database is available at the website <a href="http://www-biology.ucsd.edu/labs/Kadonaga/DCPD.html">http://www-biology.ucsd.edu/labs/Kadonaga/DCPD.html</a>.

from: Kutach & Kadonaga 2000, Mol Cell Biol 20: 4754-64.



FIG. 7. A model of two distinct interactions of TFIID with TATA- versus DPE-driven core promoters. The model is discussed in the text. TAFs, TBP-associated factors.

from: Kutach & Kadonaga 2000, Mol Cell Biol 20: 4754-64.

A further promoter element, the **BRE** (TFIIB-response element) was unraveled essentially by structural analysis



Architecture of the human ternary TFIIBc-TBPc-MLP (Major late promoter) complex. (A) Ribbons and space-filling representation of one asymmetric unit. (**B**) An isolated ternary complex as viewed down the pseudo-2-fold axis of the hTBPc-TATA-box interaction. (**C**) Species-specific differences in the TBP-TFIIB interface. Residues of hTBPc are colored in red, and those of A.thaliana (aTBP2) in vellow. Residues of hTFIIBc are shown in light green, those of the previously determined structure in blue.

Tsai FTF & Siegler PB. EMBO J., 19: 25-36, 2000. Structural basis of preinitiation complex assembly on human Pol II promoters



**Fig. 3.** Protein-DNA contacts that specify the orientation of the hTFIIBc-hTBPc-MLP complex. Complete schematic illustrating all protein-DNA interactions in the ternary complex. Pink: TBP. Green: TFIIB. Arrows indicate the location of the BRE. An oval indicates an interaction between the promotor and the protein side chain, and a

promoter and the protein side chain, and a square an interaction with the protein main chain.

Amino acid residues that are in contact with the major groove are shown in upper-case letters, and those in contact with the minor groove in lower-case letters. Hydrogen bonds are represented by dotted lines.

from: Tsai FTF & Siegler PB. EMBO J., 19: 25-36, 2000

#### **Regulatory factors**

The Human Genome contains genes encoding 1700-1900 transcription factors.



Transcription factors are composed of DNA-binding domains (DBA) plus trans-activating domains (TA)

These domains are independent from each other

DBD is sequence-specific



The simple DNA binding per se does not influence transcription

DBDs and TA act independently as demonstrated by domain-switching experiments



Transcription factors have multiple roles:

- 1. Protein-protein interaction with the basal complex (PIC) through Mediator
- 2. Recruitment of histone modifying enzymes
- 3. Recruitment of Chromatin Remodeling Complexes

Distant regulatory regions contact basal promoters by looping



cis-regulatory elements or modules



Quindi, quando un gene deve essere attivato, la sequenza di eventi è:

il riconoscimento del gene da parte di uno o più fattori di trascrizione regolatori, seguito da associazione di enzimi che rimodellano la cromatina localmente (HAT+ATPremodeling complexes), che permette il legame di ulteriori fattori regolatori, del Mediatore e del PIC, ed infine la trascrizione del gene Let's go back to the promoter story....

This kind of information ("**textbook promoter**") was obtained historically by studying a limited number of promoters with well-defined TSS, clear promoter activity and defined regulatory elements.

Several promoters, less defined and more difficult to study, were left apart, typical example the CpG – type promoters.

Genomic studies have partially changed our knowledge of promoters.

Studies oriented to define the **TSS** genome-wide, such as CAGE and 5' - SAGE, were especially instructive.

These studies demonstrated, first, that the "textbook promoter" is present at no more that 10-20% of mammalian genes (17% in human and mouse), which represent a group of inducible, tissue-specific genes.

Remaining transcription units have different structures, more often relying on CpG islands.



#### RACE= rapid amplification of cDNA ends







**Figure 3** Characterization of a novel testis transcript using tiling arrays. **a**, An EVG discovered in the analysis of chromosome 22 (Fig. 2e) was localized to a 10-kb region at one end of the insert of BAC clone AL031587. Both strands of this 113-kb genomic interval were tiled with 60-mer probes at 10-bp steps. The tiling array was hybridized with RNA from human testis. **b**, Hybridization signals corresponding to tiling probes from this region were filtered and plotted as log<sub>10</sub> values of the normalized signal strengths. Of the

six Genscan predicted exons in this region, two (exons 3 and 6) were at variance with the hybridization data. **c**, Detailed views of tiling data showing one correctly predicted exon and one incorrectly predicted exon. **d**, Typically, tiling data narrow the search window for an intron/exon boundary to 20–30-bp. The exact splice junction is then identified using consensus sequences (GT-AG rule) and ORF information. The exact splice junction can also be determined by sequencing RT–PCR products.

# Direct isolation and identification of promoters in the human genome

Tae Hoon Kim,<sup>1</sup> Leah O. Barrera,<sup>1</sup> Chunxu Qu,<sup>1</sup> Sara Van Calcar,<sup>1</sup> Nathan D. Trinklein,<sup>4</sup> Sara J. Cooper,<sup>4</sup> Rosa M. Luna,<sup>2</sup> Christopher K. Glass,<sup>2</sup> Michael G. Rosenfeld,<sup>3</sup> Richard M. Myers,<sup>4</sup> and Bing Ren<sup>1,2,5</sup>

<sup>1</sup>Ludwig Institute for Cancer Research, <sup>2</sup>Department of Cellular and Molecular Medicine, and <sup>3</sup>Howard Hughes Medical Institu University of California, San Diego, La Jolla, California 92093, USA; <sup>4</sup>Department of Genetics, Stanford University School of Medicine, Stanford, California 94305, USA

Transcriptional regulatory elements play essential roles in gene expression during animal development and cellular response to environmental signals, but our knowledge of these regions in the human genome is limited despite the availability of the complete genome sequence. Promoters mark the start of every transcript and are an important class of regulatory elements. A large, complex protein structure known as the pre-initiation complex (PIC) is assembled on all active promoters, and the presence of these proteins distinguishes promoters from other sequences in the genome. Using components of the PIC as tags, we isolated promoters directly from human cells as protein–DNA complexes and identified the resulting DNA sequences using genomic tiling microarrays. Our experiments in four human cell lines uncovered 252 PIC-binding sites in 44 semirandomly selected human genomic regions comprising 1% (30 megabase pairs) of the human genome. Nearly 72% of the identified fragments overlap or immediately flank 5' ends of known cDNA sequences, while the remainder is found in other genomic regions that likely harbor putative promoters of unannotated transcripts. Indeed, molecular analysis of the RNA isolated from one cell line uncovered transcripts initiated from over half of the putative promoter fragments, and transient transfection assays revealed promoter activity for a significant proportion of fragments when they were fused to a luciferase reporter gene. These results demonstrate the specificity of a genome-wide analysis method for mapping transcriptional regulatory elements and also indicate that a small, yet significant number of human genes remains to be discovered.

[Supplemental material is available online at www.genome.org.]

Other methods to identify genome-wide promoters have been used.

One strategy is to identify all the sequences that are bound by RNA Polymerase II and by basal transcription factors, such as TAFs, by using ChIP-on-chip.

As a part of the ENCODE project, for example, 1%-human genome coverage tiling arrays were used to identify sequences bound by RNA PolII and TAF1 (see TAF nomenclature).





Figure 1. Direct isolation and identification of promoters in the human genome. (*A*) A schematic of GWLA for mapping RNAP and TAF1-binding sites. Growing cells are cross-linked with formaldehyde, and their nuclei are isolated , the chromatin extracted and sonicated.

The resulting chromatin (protein–DNA) complexes are incubated with either anti-RNAP (-RNAP) or anti-TAF1 (-TAF1) antibody. The immunoprecipitated DNA is subjected to ligation-mediated PCR, labeled with Cy5 dye, and competitively hybridized to the ENCODE array with the Cy3-labeled unenriched chromatin (Li et al. 2003).





(*B*) A typical detailed view of the TAF1- and RNAP-binding data (within the ENCODE region, ENr231). Negative logarithmic *P* values of enrichment by RNAP or TAF1 ChIP for each probe are plotted in green. Relative enrichment (R) values by RNAP or TAF1 ChIP for each probe fragment are plotted in red on the inverted axis. (*C*) A representative view of an entire ENCODE locus (ENr231) with annotated RNAP- and TAF1-binding states that have *P* values <0.0001, noted in red blocks (detailed view in *B* is highlighted in gray).

This is a modern way to show information on a genome scheme, and is used today by most genomic browsers

**Genome annotation** taken by multiple studies, is aligned to genome maps, thus giving a representational picture of "what is known" on each section of the genome in a continuous fashion.

### Genome annotation

http://genome.ucsc.edu/cgi-bin/hgGateway





Annotazione bioinformatica di un segmento di 15 kb del genoma umano, contenente un gene, utilizzando il browser Genotator.

Tutto il genoma viene continuamente ri-annotato considerando le nuove evidenze sperimentali, se liberamente disponibili.



### Figure 2. Summary of identified PIC-binding sites matched to transcripts

in IMR90. Histograms plotting relative distance of the RNAP-, TFIID-, or PIC-binding sites to the nearest 5 ends of fulllength GenBank RefSeq transcripts. (A) Distribution (number of RNAP-binding sites, y-axis)

of relative distances (in Kbp, *x-axis) of the RNAP-binding sites to the* nearest 5 ends of full length mRNA. The first and last bars are counts for those RNAP-binding sites that are >2.5 Kbp upstream (<) or downstream (>) from the nearest 5 end. (*B*) *Distribution (number of TFIID-binding* sites, *y-axis) of relative distances (in Kbp, x-axis) of the TFIID-binding sites* to the nearest 5 ends of full-length mRNA. The first and last bars are counts for those TFIID-binding sites that are >2.5 Kbp upstream (<) or downstream (>) from the nearest 5 end. (*C*) *Distribution (number of* PIC-binding sites, *y-axis) of relative distances (in Kbp, x-axis) of the PICbinding* sites to the nearest 5 ends of full-length mRNA. The first and last bars are counts for those PIC-binding sites that are >2.5 Kbp upstream (<) or downstream (>) from the nearest 5 end. When a sequence is identified as a putative promoter, the functional validation consists in the classical **reporter assay** 

It this assay, a minigene is constructed fusing the putative promoter sequence to a **reporter** gene, i.e. a gene whose product can be easily measured.

This construct is then transfected into cultured cells and the product measured after a period of time necessary for transgnene import and expression.

Mammalian **reporter genes** are usually nonmammalian genes, such as

- •CAT cloramphenicol acetyltransferase
- •Luc firefly luciferase
- •GFP Jellyfish green fluorescent protein
- ullet eta-gal beta-galactosidase





The reporter gene assay for promoter and enhancers





**Figure 3.** Experimental validation of the putative promoters by reporter assays. (*A*) A schematic of reporter assay used to determine whether the identified putative promoter fragment can support transcription. Each putative promoter fragment was segmented into 750-bp fragments by PCR and cloned into the luciferase reporter construct, pGL3, in either forward or reverse orientations. The resulting reporter constructs were individually transfected into HT1080 cells and the resulting luciferase activity was measured.





**Putative Promoter Fragments** 

**Figure 3.** Experimental validation of the putative promoters by reporter assays. (*B*) Reporter activities of 17 putative promoter fragments are shown (remaining five of 22 were not tested). Relative reporter activity was determined by comparing the luciferase activity of the test fragment and the control genomic DNA fragments. Promoter fragments with significant reporter activity (exceeding three times the standard deviation of all control fragments) are highlighted in gray.











However, the most impressive data came from the FANTOM project carried out using CAGEX

The most extensive core promoter identification study undertaken so far used CAGE tags to identify 184,379 human and 177,349 mouse core promoters, many of which might contain a cluster of individual TSSs<sup>24</sup>. A previous analysis that involved full-length cDNA sequencing identified 30,964 human and 19,023 mouse promoters42. But even the most recent figures are likely to be a substantial underestimate. First, sequencing 50–100,000 tags in each library can reliably detect only those transcripts that are expressed at a level of at least 10 copies in each cell (as there are at least 400,000 mRNAs in an average mammalian cell43). Many transcripts are not present at this level, either because they are of low abundance in individual cells or are expressed in only a small subset of cells in the tissues that have been studied.

# Genome-wide analysis of mammalian promoter architecture and evolution

Piero Carninci<sup>1,2,21</sup>, Albin Sandelin<sup>1,3,21</sup>, Boris Lenhard<sup>1,3,20,21</sup>, Shintaro Katayama<sup>1</sup>, Kazuro Shimokawa<sup>1</sup>, Jasmina Ponjavic<sup>1,20</sup>, Colin A M Semple<sup>1,4</sup>, Martin S Taylor<sup>1,5</sup>, Pär G Engström<sup>3</sup>, Martin C Frith<sup>1,6</sup>, Alistair R R Forrest<sup>6</sup>, Wynand B Alkema<sup>3</sup>, Sin Lam Tan<sup>7</sup>, Charles Plessy<sup>2</sup>, Rimantas Kodzius<sup>1,2</sup>, Timothy Ravasi<sup>1,6,8</sup>, Takeya Kasukawa<sup>1,9</sup>, Shiro Fukuda<sup>1</sup>, Mutsumi Kanamori-Katayama<sup>1</sup>, Yayoi Kitazume<sup>1</sup>, Hideya Kawaji<sup>1,9</sup>, Chikatoshi Kai<sup>1</sup>, Mari Nakamura<sup>1</sup>, Hideaki Konno<sup>1</sup>, Kenji Nakano<sup>1,9</sup>, Salim Mottagui-Tabar<sup>3,20</sup>, Peter Arner<sup>10</sup>, Alessandra Chesi<sup>11</sup>, Stefano Gustincich<sup>11</sup>, Francesca Persichetti<sup>12</sup>, Harukazu Suzuki<sup>1</sup>, Sean M Grimmond<sup>6</sup>, Christine A Wells<sup>19</sup>, Valerio Orlando<sup>13</sup>, Claes Wahlestedt<sup>3,20</sup>, Edison T Liu<sup>14</sup>, Matthias Harbers<sup>15</sup>, Jun Kawai<sup>1,2</sup>, Vladimir B Bajic<sup>1,7,16</sup>, David A Hume<sup>1,6,21</sup> & Yoshihide Hayashizaki<sup>1,2,17,18</sup>

<sup>7</sup> Mammalian promoters can be separated into two classes, conserved TATA box-enriched promoters, which initiate at a welldefined site, and more plastic, broad and evolvable CpG-rich promoters. We have sequenced tags corresponding to several hundred thousand transcription start sites (TSSs) in the mouse and human genomes, allowing precise analysis of the sequence architecture and evolution of distinct promoter classes. Different tissues and families of genes differentially use distinct types of promoters. Our tagging methods allow quantitative analysis of promoter usage in different tissues and show that differentially regulated alternative TSSs are a common feature in protein-coding genes and commonly generate alternative N termini. Among the TSSs, we identified new start sites associated with the majority of exons and with 3' UTRs. These data permit genome-scale identification of tissue-specific promoters and analysis of the *cis*-acting elements associated with them.

VOLUME 38 | NUMBER 6 | JUNE 2006 NATURE GENETICS

# 146 mouse cDNA libraries41 human cDNA libraries



626



Figure 1 Definition and characteristics of CAGE tag clusters. (a) Tag clusters are produced by grouping overlapping tags on the same strand. Hence, tag clusters are defined by a start and end position, a count of tags and a distribution of these counts. Unique tag starts within the tag cluster form CAGE tag starting sites (CTSSs).

Students you must go back to previous lectures to find out how a CAGE library is generated and what a CAGE tag is.





Figure 1. (c) Association of tag cluster **width** (minimal length of the sequence fragment containing >80% of all tags in the cluster) with TATA boxes and CpG islands for tag clusters with >100 tags.



Figure 1. (e) Arrays of representative tag clusters for different shape classes. Histograms indicate the fraction of tags in the tag cluster mapping into each position in a 120-bp window centered on the tag cluster. The single peak (SP) class is characterized by a sharp peak, indicative of a single, well-defined TSS. The broad (BR) shape indicate multiple, weakly defined TSSs. The bimodal/multimodal (MU) shape class implies multiple welldefined TSSs within one cluster. Combination of a welldefined TSS surrounded by weaker TSSs results in a broad with dominant peak shape (PB). HUGO gene names or transcriptional unit identifiers for cognate genes and tag cluster identifiers are shown above each tag cluster.





Figure 2. TATA-box and TSS spacing definition and consensus. (a) Accurate distribution of the spacing between TATA-box promoter and initiation sites.



Legend to the previous slide

Figure 2. TATA-box and TSS spacing definition and consensus. (b–e) Sequence logos for promoter sequences aligned at the TSSs constructed by counting each tag and its flanking region as one sequence, divided by promoter shape class. The y axis shows the information content (measured in bits), reviewed in ref. 15. In all cases, there is a clear preference for a pyrimidine-purine initiation site at -1,+1. A TATA-like motif is visible around the -30 position in the SP class promoters (b). In the BR class promoters, as most of those promoters are overlapped by CpG islands, the entire region is GC-rich; there is anisotropy of nucleotide content: there are more guanine than cytosine nucleotides in the plus strand upstream of the TSS (c). The logos of PB (d) and MU (e) class promoters look similar to this, indicating that these two ambiguous two categories are more likely to share the common initiation mechanism with BR promoters than with the SP ones. The PB class has a certain proportion of mixed cases, with both a CpG island and a TATA-box.



Figure 4 Pyrimidine-purine dinucleotides drive expression.

(a) A detailed view of the core promoter of the mouse Ptprn gene (TC 73140) and corresponding human region illustrates the usage of pyrimidine-purine dinucleotides as dominant start sites and the expression changes resulting from mutations at these positions.

Ba Mouse Syn1 promoter (TC id: T0XR0125C59F)



Promoter structures are conserved



Human PURA promoter (TC id: T05F085033E0)

#### Bb Mouse Pura promoter (TC id: T18F0230753D)

What have we learnt from genome-wide expression analysis ?

-more protein-coding genes than known

-most of genes produce more RNA transcripts than expected (derived both from nonclassical sense or antisense transcription and by alternative RNA processing)

-unexpected variety of noncoding RNA transcripts of various size (from very short to very long, from intergenic to intragenic), whose functions are mostly unknown (with the exception of miRNA and few long ncRNA)

-tissue- , cell- and stage-specific expression of both protein-coding and noncoding RNA widely confirmed

Most transcript use RNA Polymerase II