

L4.2

# L4.2

## Agenda

- other transcripts from the catalogue
- what is a gene today ?
- RNA Polymerase pausing
- extensive AS of lncRNAs
- mechanisms of exon splicing
- *cis*-regulatory sequences of splicing and AS
- trans-acting factors for AS (RBP)
- tissue-specific AS regulators

# An RNA catalogue

It is worth spending few minutes on the Statistics to consider how many different types of long- and short-noncoding RNA have been catalogued



GENCODE


Data

Stats

Browser


Blog

DATA  
STATISTICS  
BROWSER

A close-up photograph of a young woman's face, split vertically down the middle. The left side shows her skin and features, while the right side is a solid white color, representing a mask or a specific data visualization.

HUMAN  
GENCODE 24 (09.12.15)

DATA  
STATISTICS  
BROWSER

A photograph of a small, grey mouse sitting on a white surface, looking directly at the camera.

MOUSE  
GENCODE M9 (10.03.16)


The GENCODE project produces high quality reference gene annotation and experimental validation for human and mouse genomes

## Announcements

- Read this [article](#) comparing GENCODE and RefSeq annotations
- NEW** Human GENCODE releases mapped back to GRCh37 now available [here](#)
- NEW** GENCODE M9, corresponding to Ensembl 84, has been released!
- Next human (25) and mouse (M10) releases scheduled for July 2016

## Tweets by @GenodeGenes

 GencodeGenes Retweeted 

 **APPRIS** @appris\_cnio  
{APPRIS} 2016\_03.v15 is out!! Principal Isoforms for multiple species in #Ensembl84 @GencodeGenes bit.ly/1pvjf5f 

 **GencodeGenes** @GencodeGenes  
Human GENCODE 24, released Aug 2015 

Embed

[View on Twitter](#)

<http://www.gencodegenes.org/>

## RNA Biotypes

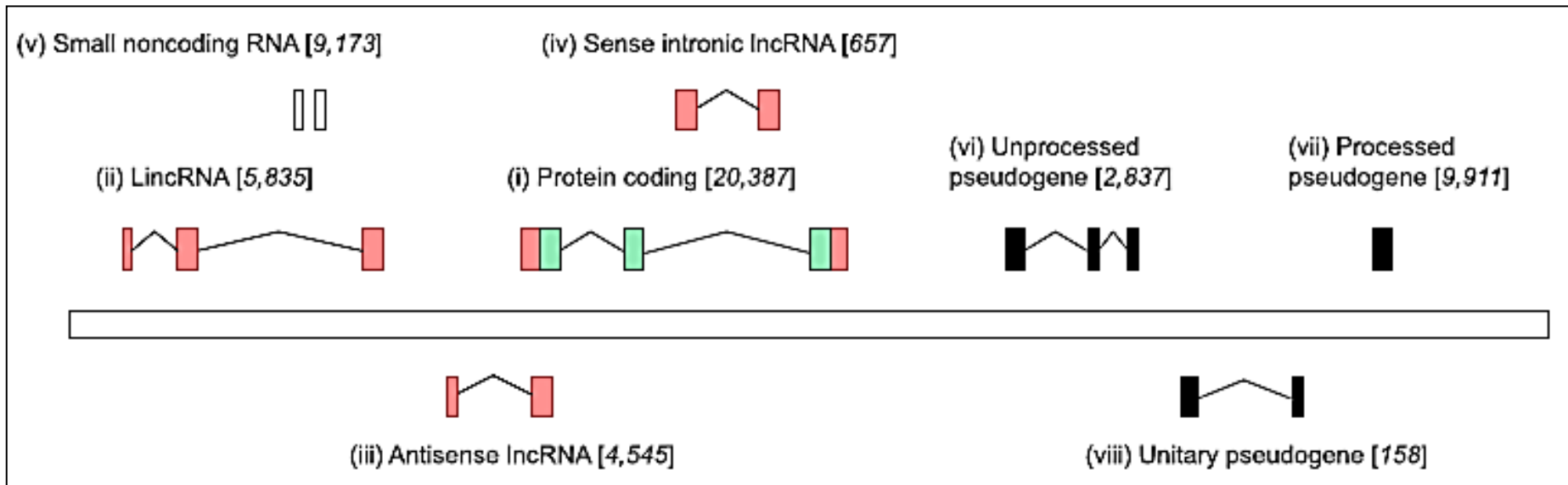


Figure 2. A summary of locus biotypes in GENCODE.

# Functional transcriptomics in the post-ENCODE era

Jonathan M. Mudge,<sup>1</sup> Adam Frankish, and Jennifer Harrow

*Department of Informatics, Wellcome Trust Sanger Institute, Hinxton CB10 1SA, United Kingdom*

The last decade has seen tremendous effort committed to the annotation of the human genome sequence, most notably perhaps in the form of the ENCODE project. One of the major findings of ENCODE, and other genome analysis projects, is that the human transcriptome is far larger and more complex than previously thought. This complexity manifests, for example, as alternative splicing within protein-coding genes, as well as in the discovery of thousands of long noncoding RNAs. It is also possible that significant numbers of human transcripts have not yet been described by annotation projects, while existing transcript models are frequently incomplete. The question as to what proportion of this complexity is truly functional remains open, however, and this ambiguity presents a serious challenge to genome scientists. In this article, we will discuss the current state of human transcriptome annotation, drawing on our experience gained in generating the GENCODE gene annotation set. We highlight the gaps in our knowledge of transcript functionality that remain, and consider the potential computational and experimental strategies that can be used to help close them. We propose that an understanding of the true overlap between transcriptional complexity and functionality will not be gained in the short term. However, significant steps toward obtaining this knowledge can now be taken by using an integrated strategy, combining all of the experimental resources at our disposal.

# «Classical» versus modern view of a gene

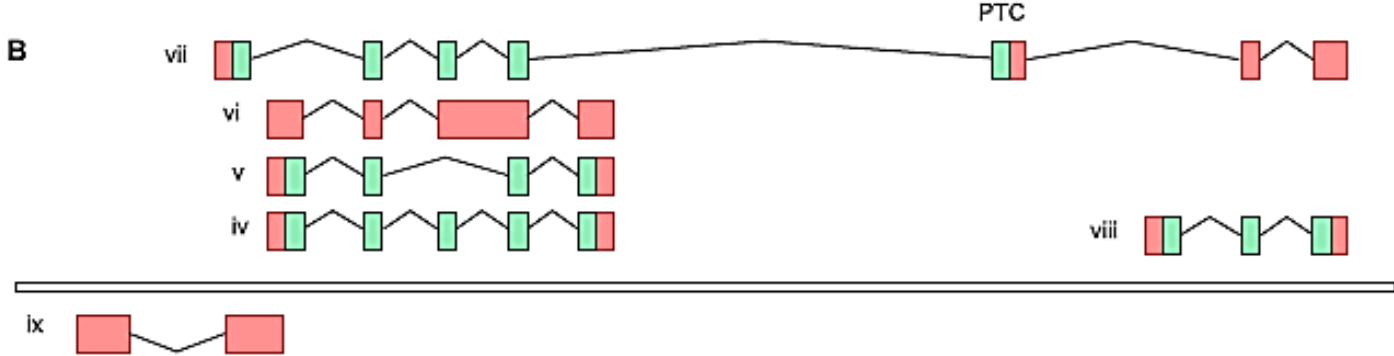
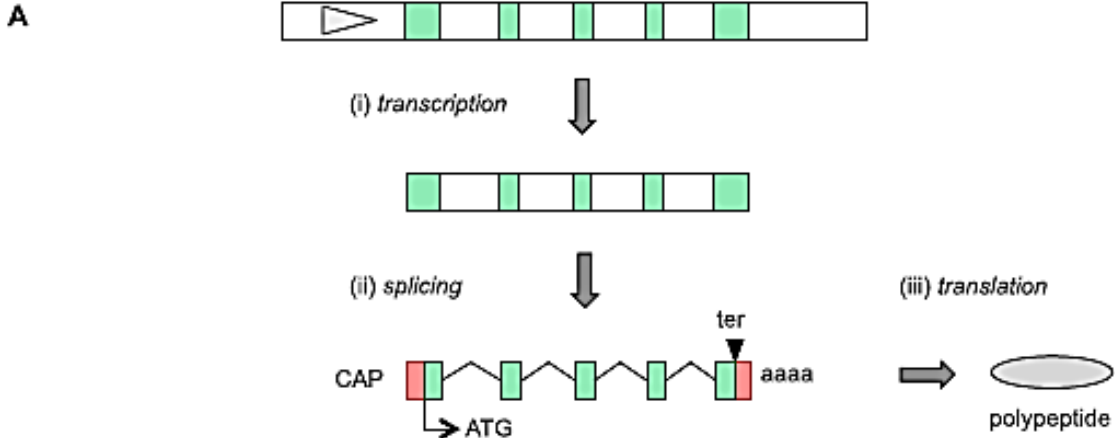
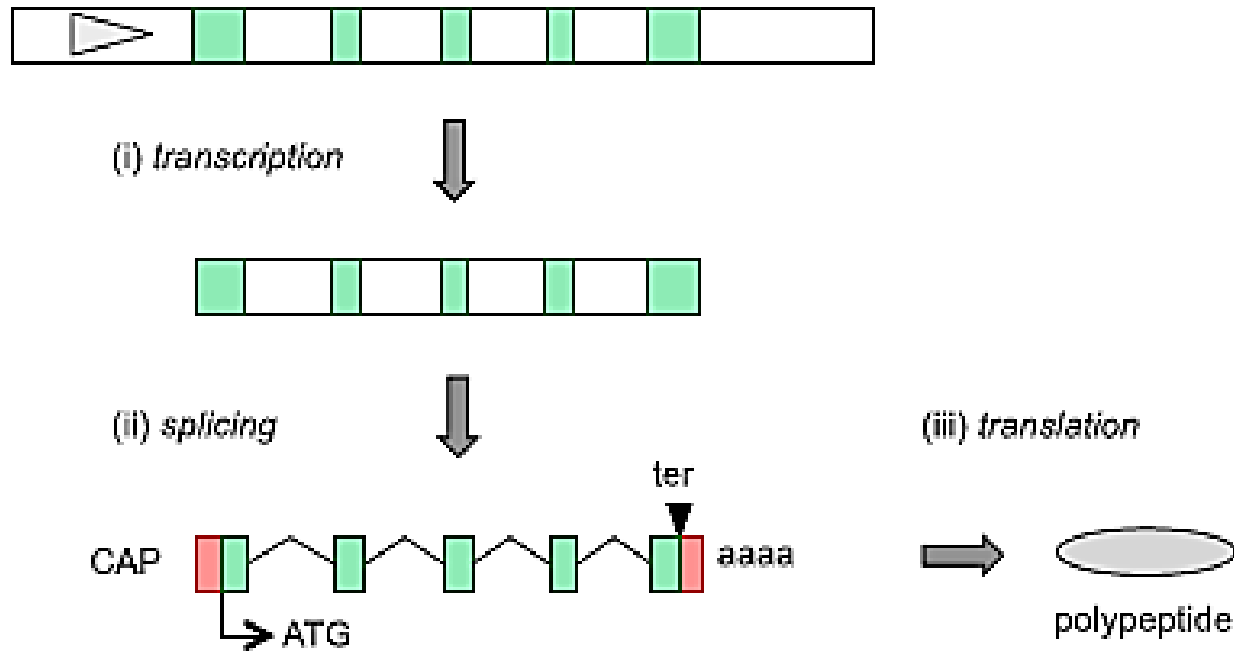
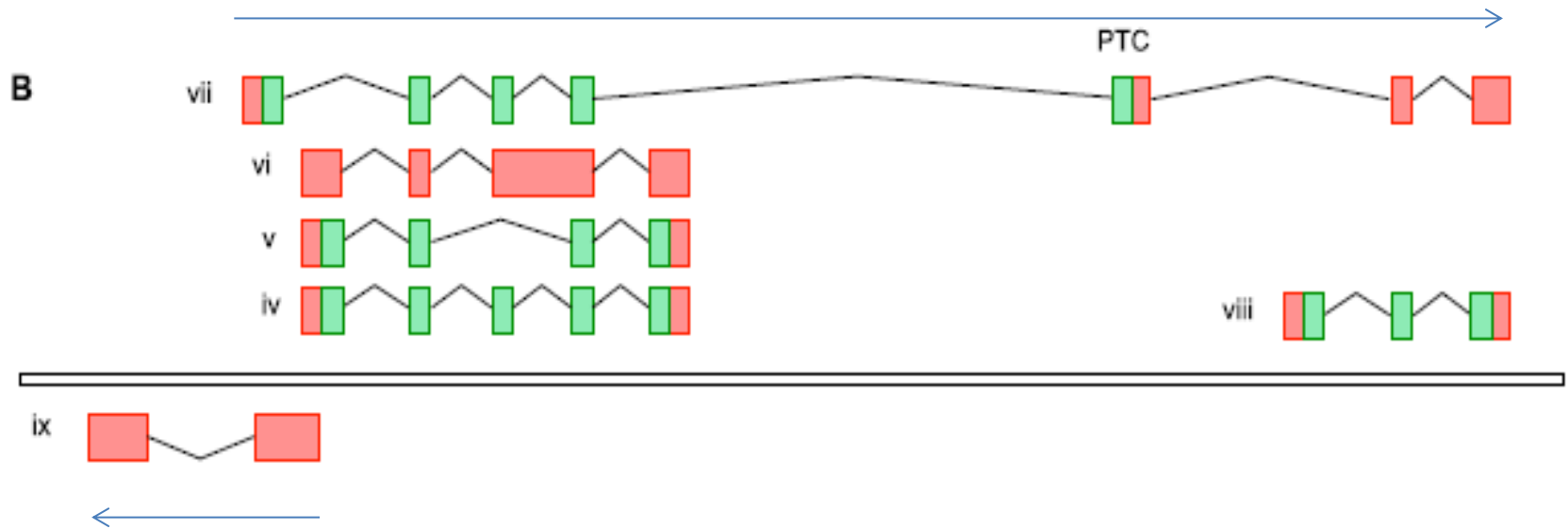


Figure 1. The evolving dogma of gene transcription.



(A) The **historical “central dogma”** of molecular biology. By this model, (i) transcription generates the primary transcript (exons in green, introns in white), with the initial interaction between the RNA polymerase complex and the genome being mediated by a promoter region (gray triangle). (ii) The introns of the primary transcript are removed by the spliceosome, and a mature mRNA is generated by 5' end capping (CAP) and polyadenylation (aaaa) (coding region [CDS] shown in green, untranslated 5' and 3' UTRs in red). (iii) The mRNA is translated into a polypeptide by the ribosome complex, with translation proceeding from the initiation codon (ATG) and ending at the termination codon (ter).



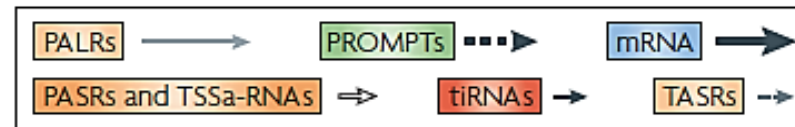
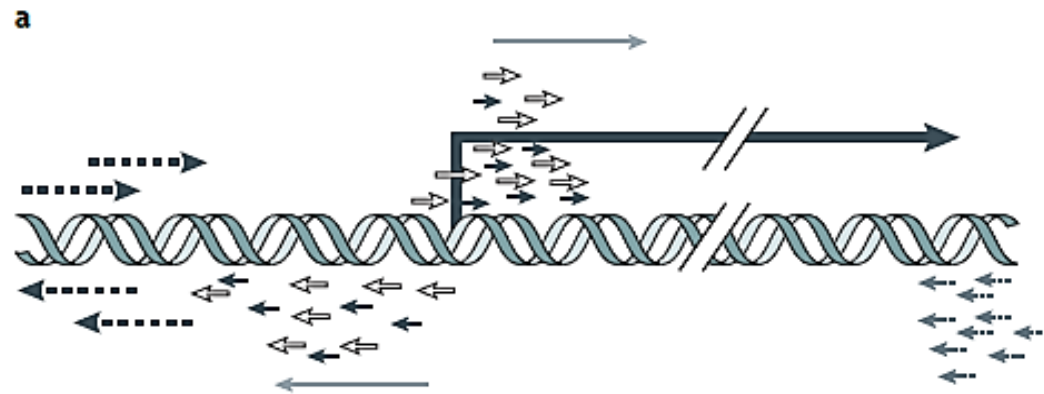


(B) An **updated model** reflecting a modern view of transcriptional complexity. Here, the same gene (iv) undergoes alternative splicing (AS), for example an exon skipping event that does not change the frame of the CDS (v); this event thus has the potential to generate an alternative protein isoform. However, products of AS cannot be assumed to be functional; this gene has generated a retained intron transcript (vi), perhaps due to the failure of the spliceosome to remove this intron. Further complexity comes from a read-through transcription event (vii), whereby a transcript is generated that also includes exons from a neighboring protein-coding locus (viii). In this example, the read-through transcript has an alternative first exon compared with the upstream gene that contains a potential alternative ATG codon, although the presence of a subsequent premature termination codon (PTC) prior to two splice junctions indicates that this transcript is likely subjected to the nonsense mediated decay (NMD) degradation pathway. Finally, model ix is a transcript that is antisense to the upstream gene; both loci are potentially generated under the control of a bidirectional promoter.

in addition....

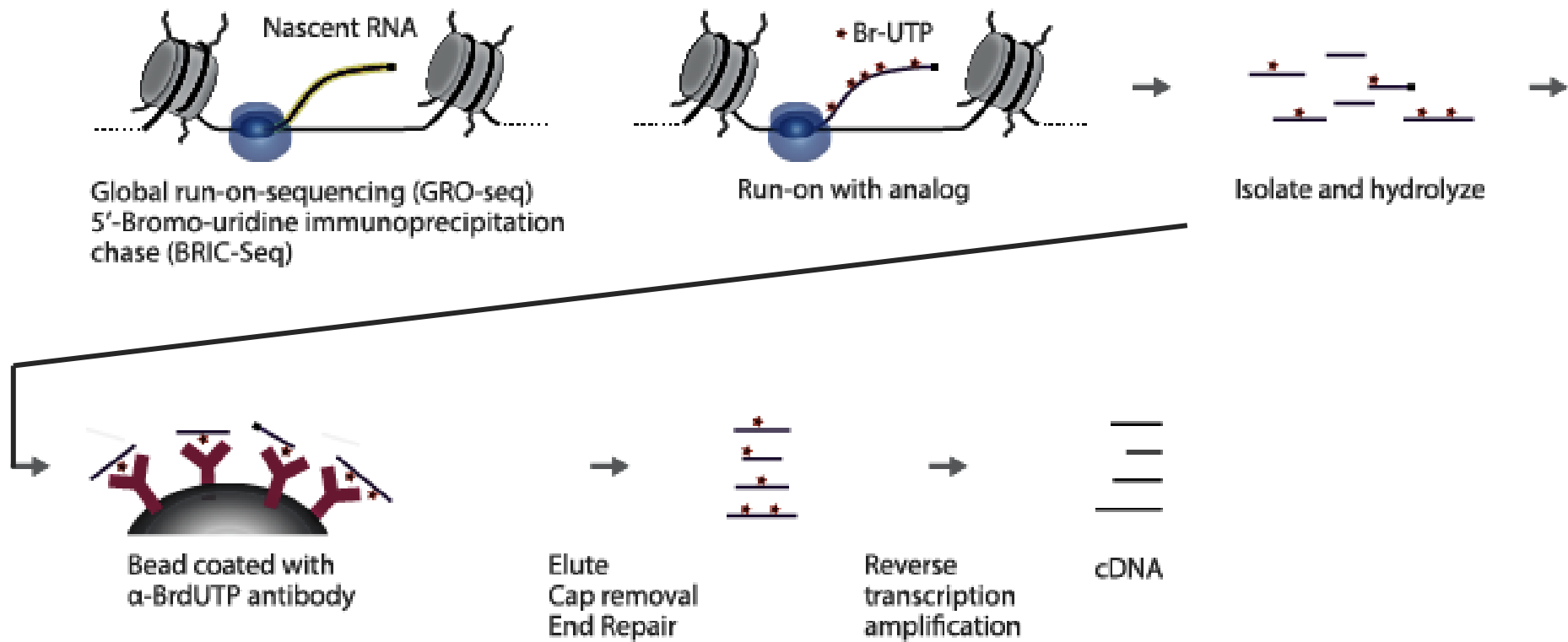
Using different technologies, including tiling microarrays, GRO-seq, CAGE, Sage and others, unstable short RNAs were also observed close to promoters

Unstable small RNA  
accompanying gene  
transcription



Short name of RNA classes	Full name of RNA classes	
PALRs	Promoter-associated long RNAs	Hundreds nt long RNAs spanning regions on proximal promoters to the first exon
PASRs	Promoter-associated short RNAs	20–70 nt long RNAs spanning regions around core promoters
TASRs	Termini-associated short RNAs	20–70 nt long RNAs spanning regions around transcription termination sites
PROMPTs	Promoter upstream transcripts	Unstable transcripts mapping 0.5–2 kb upstream the transcription starting sites
TSSa-RNAs	Transcription start sites antisense RNAs	RNAs, generally short and non-coding, generated from bidirectional activity of mammalian RNA Polymerase II
NRO-RNAs	Nuclear run-on assay derived RNAs	Short RNA detected by nuclear run-on assays, mapping 20 to 50 downstream to transcriptions starting sites of mRNAs
RE RNAs	Retrotransposon-derived RNAs	A heterogeneous class of RNAs which starting sites overlap retrotransposon elements
tiRNAs	Tiny transcription initiation RNAs	RNAs about 18 nt long, positioned about 20 bp after the transcription starting sites of highly expressed mRNAs

# GRO-Seq



from the Illumina website

As secondary product, GRO-seq localize active Polymerase



promoter-proximal pausing

# Promoter-proximal pausing of RNA polymerase II: emerging roles in metazoans

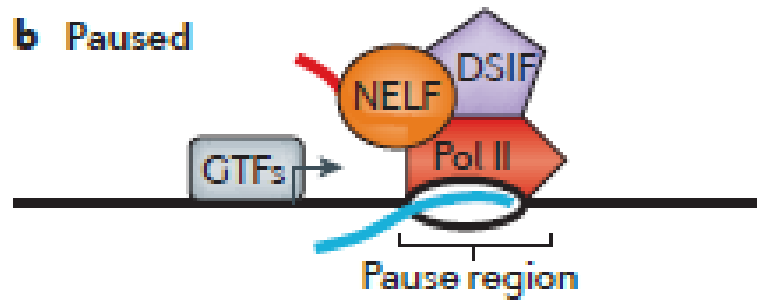
*Karen Adelman<sup>1</sup> and John T. Lis<sup>2</sup>*

**Abstract** | Recent years have witnessed a sea change in our understanding of transcription regulation: whereas traditional models focused solely on the events that brought RNA polymerase II (Pol II) to a gene promoter to initiate RNA synthesis, emerging evidence points to the pausing of Pol II during early elongation as a widespread regulatory mechanism in higher eukaryotes. Current data indicate that pausing is particularly enriched at genes in signal-responsive pathways. Here the evidence for pausing of Pol II from recent high-throughput studies will be discussed, as well as the potential interconnected functions of promoter-proximally paused Pol II.

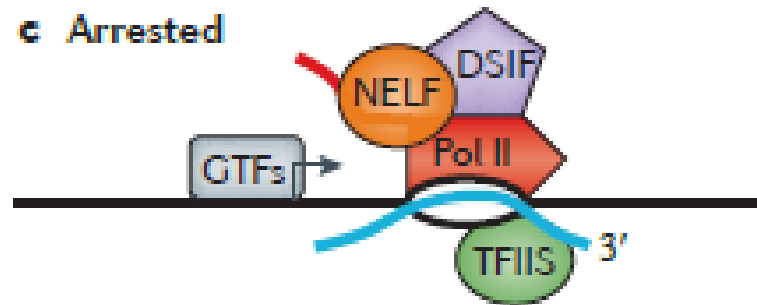
**a Pre-initiation complex**



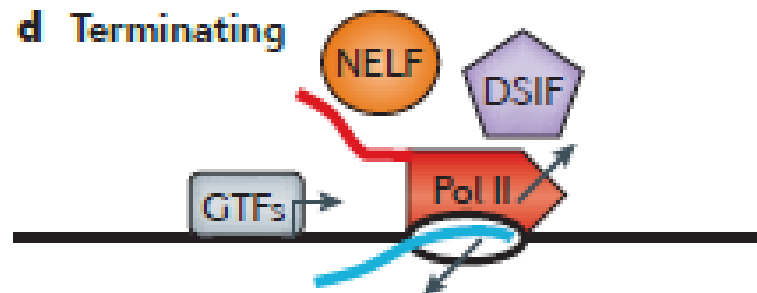
**b Paused**



**c Arrested**



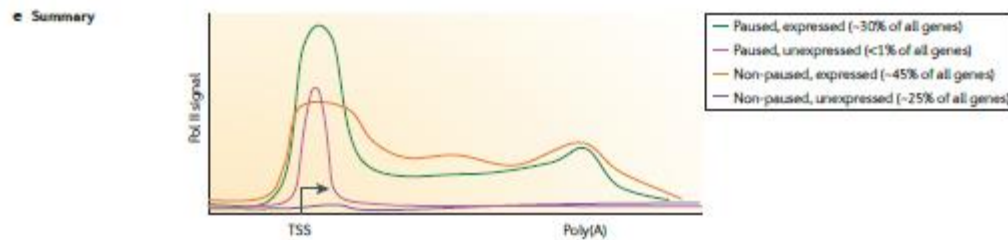
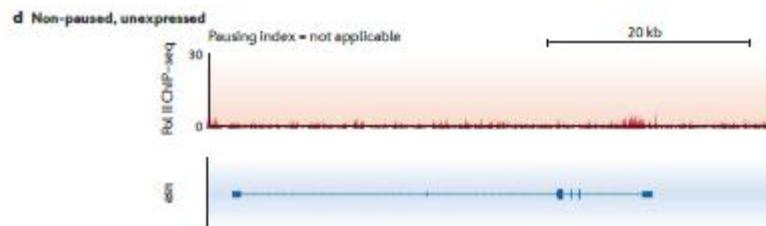
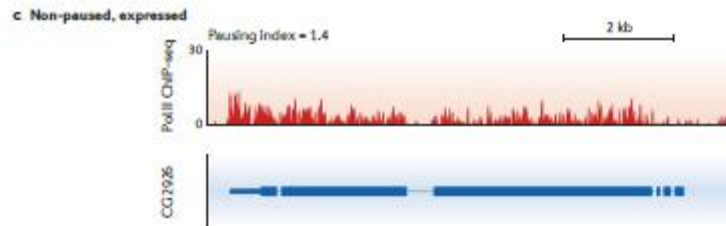
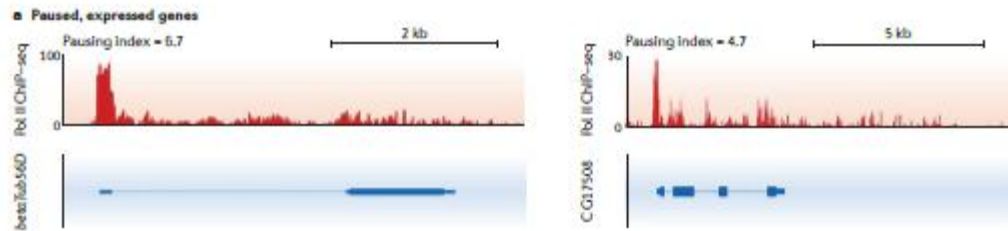
**d Terminating**



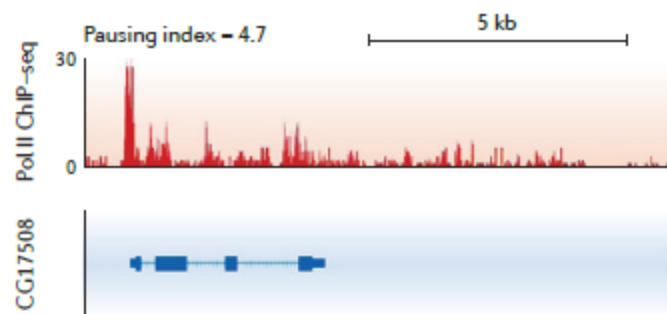
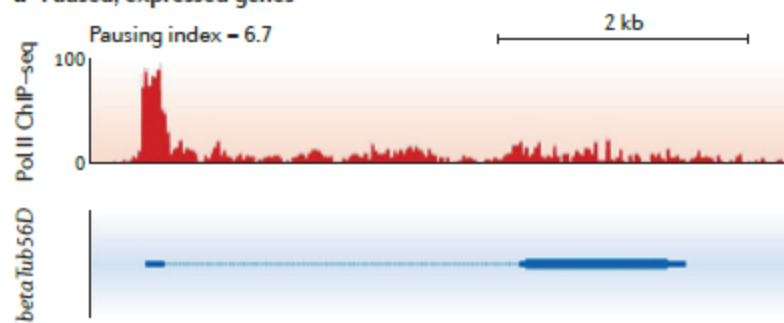
● Poised

f Stalled

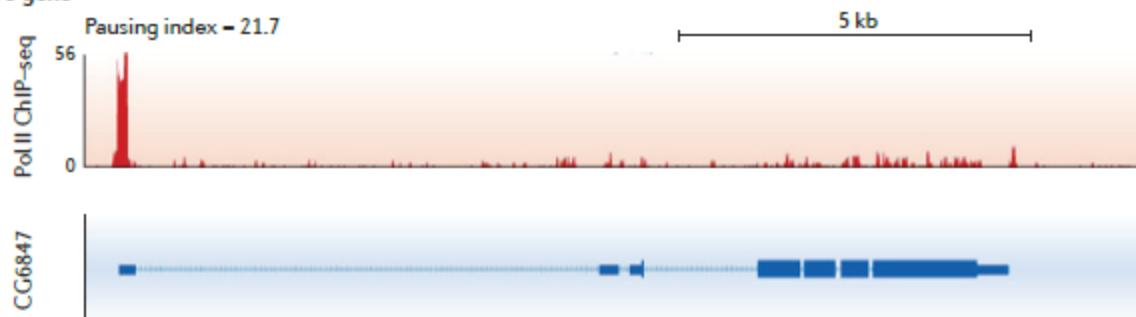




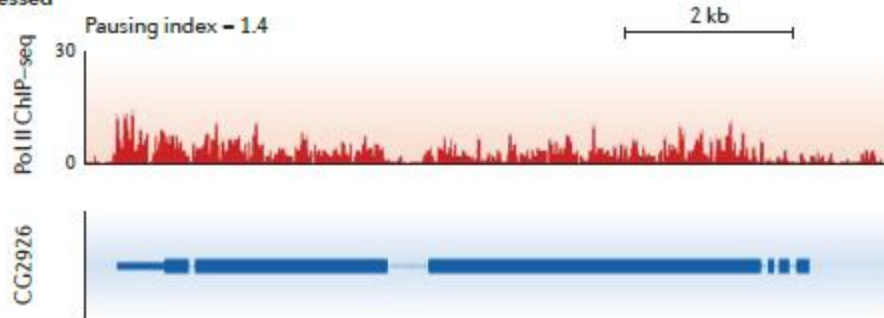
**a Paused, expressed genes**



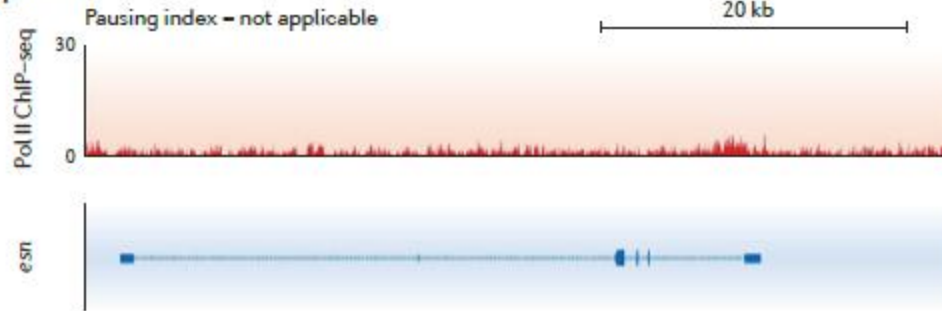
**b Paused but inactive gene**



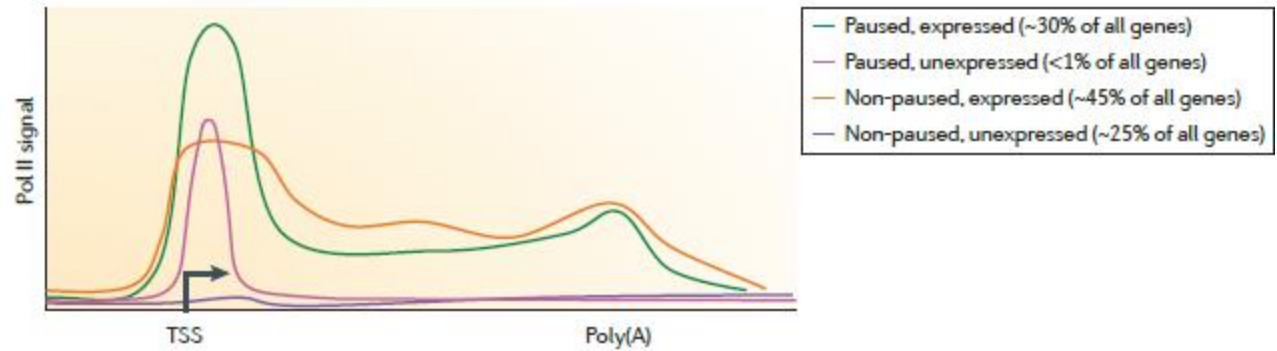
**c Non-paused, expressed**

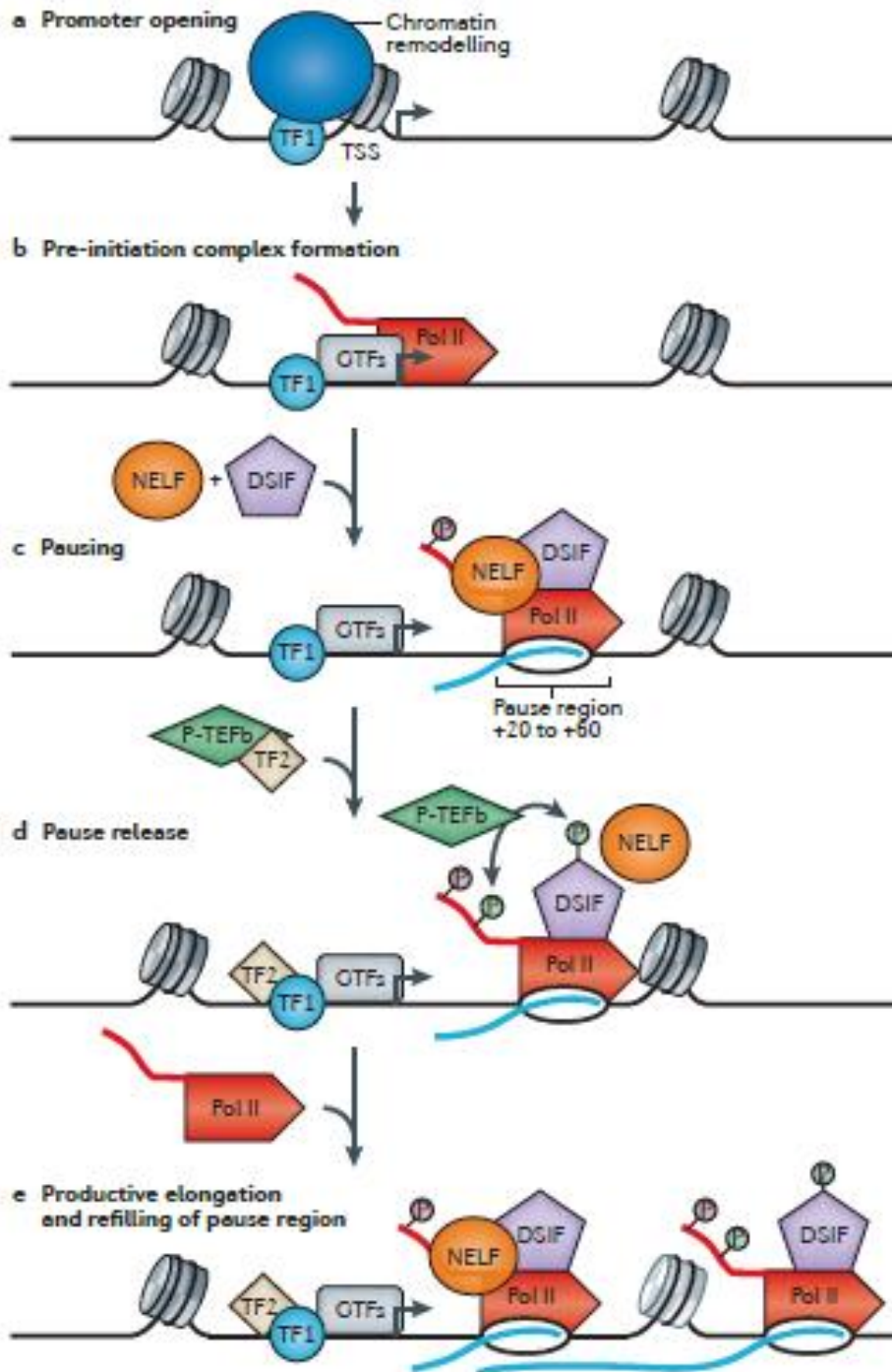


**d Non-paused, unexpressed**



**• Summary**

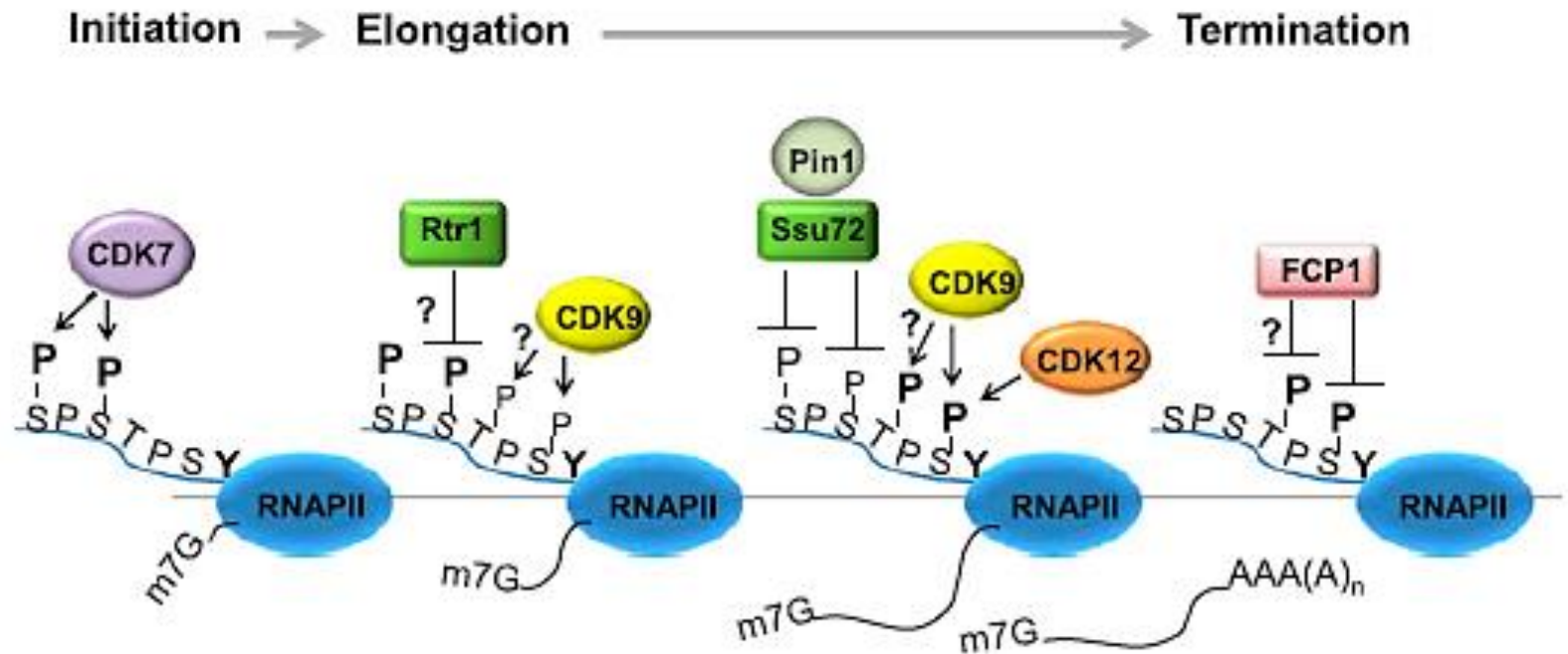




CTD phosphorylation in Ser-2 leads to SETD2 association. SETD2 is the lysine Methyl Transferase specific for H3K36 methylation

Phosphorylation of DSIF leads to NELF dissociation, while DSIF is converted to elongation factor

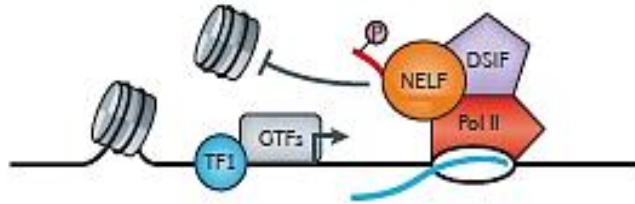
## RNA Polymerase cycle and CTD phosphorylation



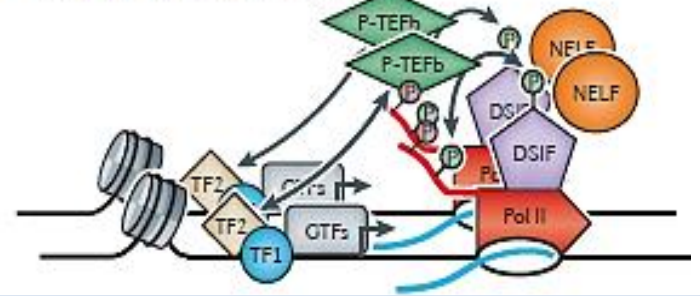
Dynamic modification of the CTD during the transcription cycle.

# Possible functions of PolII pausing

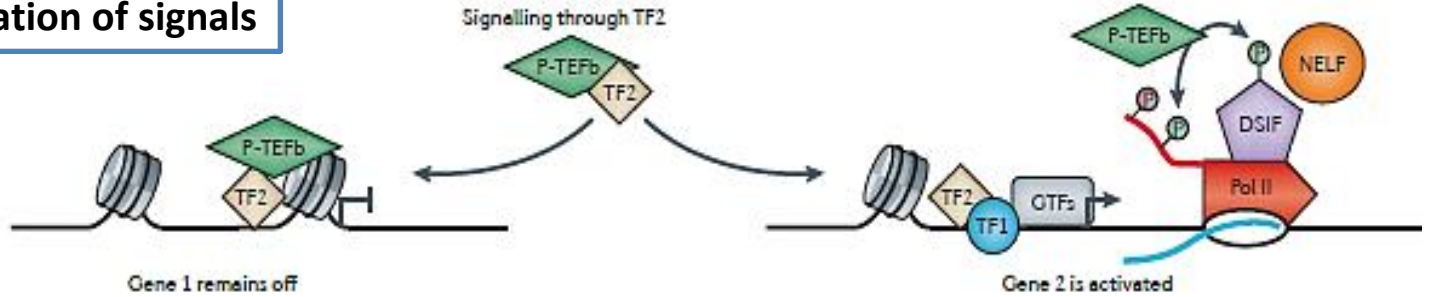
## establishing permissive chromatin



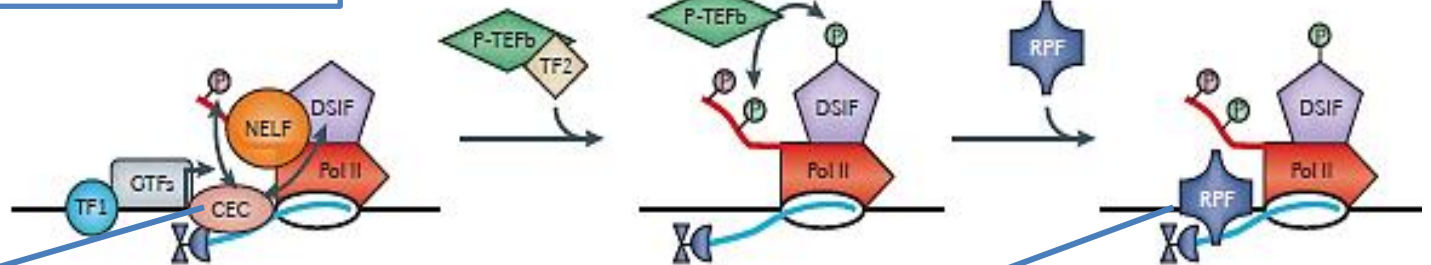
## synchronous rapid activation



## integration of signals



## checkpoint in early elongation ?



Capping Enzyme Complex

RNA Processing Factors





## Human

### Statistics about the current GENCODE Release (version 30)

The statistics derive from the [gtf file](#) that contains only the annotation of the main chromosomes.

For details about the calculation of these statistics please see the [README\\_stats.txt file](#).

#### General stats

Total No of Genes	58870	Total No of Transcripts	208621
Protein-coding genes	19986	Protein-coding transcripts	83688
<u>Long non-coding RNA genes</u>	<u>16193</u>	- full length protein-coding	57687
Small non-coding RNA genes	7576	- partial length protein-coding	26001
Pseudogenes	14706	Nonsense mediated decay transcripts	15550
- processed pseudogenes	10663	<u>Long non-coding RNA loci transcripts</u>	<u>30369</u>
- unprocessed pseudogenes	3525		
- unitary pseudogenes	221		
- polymorphic pseudogenes	42		
- pseudogenes	18	Total No of distinct translations	61870

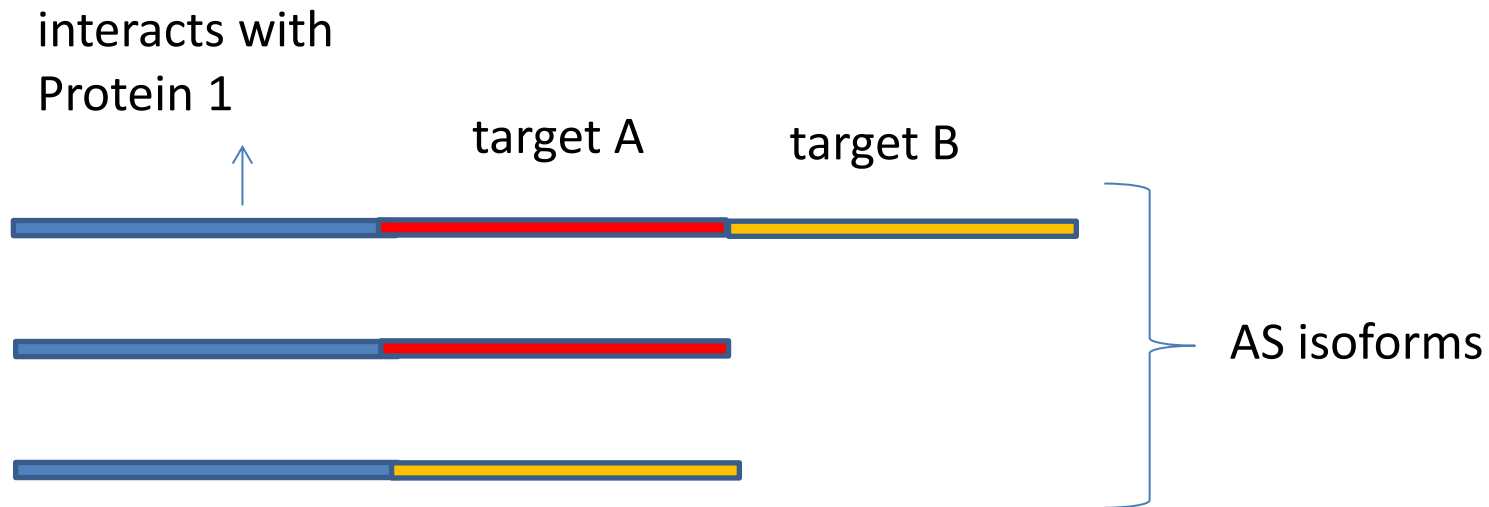


LncRNA undergo Alternative Splicing

They are capped and polyadenylated

What is the sense of making AS ?

Their function can be modulated by including/excluding certain parts.



Alternative Splicing of lncRNAs is guided by the same elements as protein-coding RNAs

However, while in protein-coding RNA alternative exons are few (average one-two on an average of 9 exons), lncRNA tend to have more alternatives.

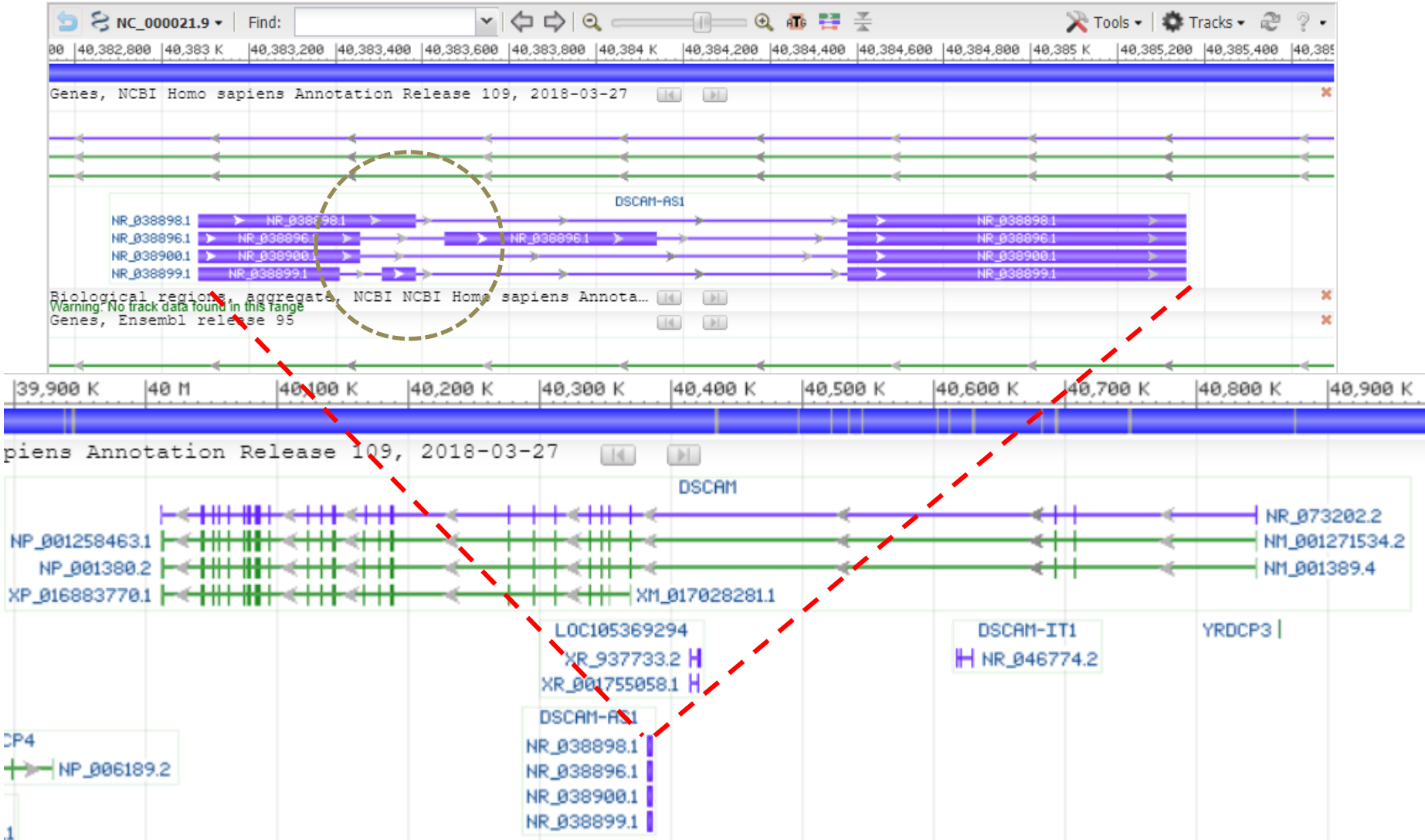
Note that lncRNAs do not have the constraint of the coding sequence.

Genomic regions, transcripts, and products

Go to [reference sequence details](#)

Genomic Sequence:

Go to nucleotide: [Graphics](#) [FASTA](#) [GenBank](#)



CP4  
NP\_006189.2

## Back to splicing

What do we know of splicing ?

When, where, how ?

Co-transcriptional or post-transcriptional ?

This is a long-standing, wearying discussion on this, lasting at least 25 years

let's see how ENCODE has approached it

## **The transcriptome of nuclear subcompartments**

For the K562 cell line, we also analysed RNA isolated from three subnuclear compartments (chromatin, nucleolus and nucleoplasm).

Almost half (18,330) of the GENCODE (v7) annotated genes detected for all 15 cell lines (35,494) were identified in the analysis of just these three nuclear subcompartments. In addition, there were as many novel unannotated genes found in K562 subcompartments as there were in all other data sets combined.

The interrogation of different subcellular RNA fractions provides snapshots of the status of the RNA population along the RNA processing pathway.

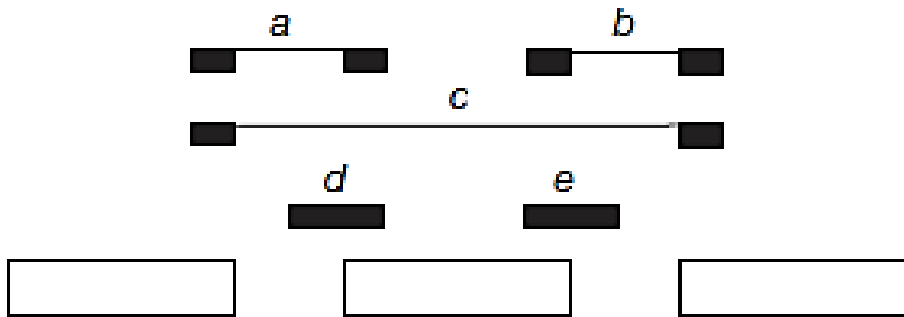
Thus, by analysing short and long RNAs in the different subcellular compartments, we confirm that splicing predominantly occurs during transcription.

By using RNA-seq to measure the degree of completion of splicing (Fig. 2a), we observed that around most exons, introns are already being spliced in chromatin-associated RNA—the fraction that includes RNAs in the process of being transcribed (Fig. 2b).

Concomitantly, we found strong enrichment specifically of spliceosomal small nuclear RNAs (snRNAs) in this RNA fraction

Co-transcriptional splicing provides an explanation for the increasing evidence connecting chromatin structure to splicing regulation, and we have observed that exons in the process of being spliced are enriched in a number of chromatin marks

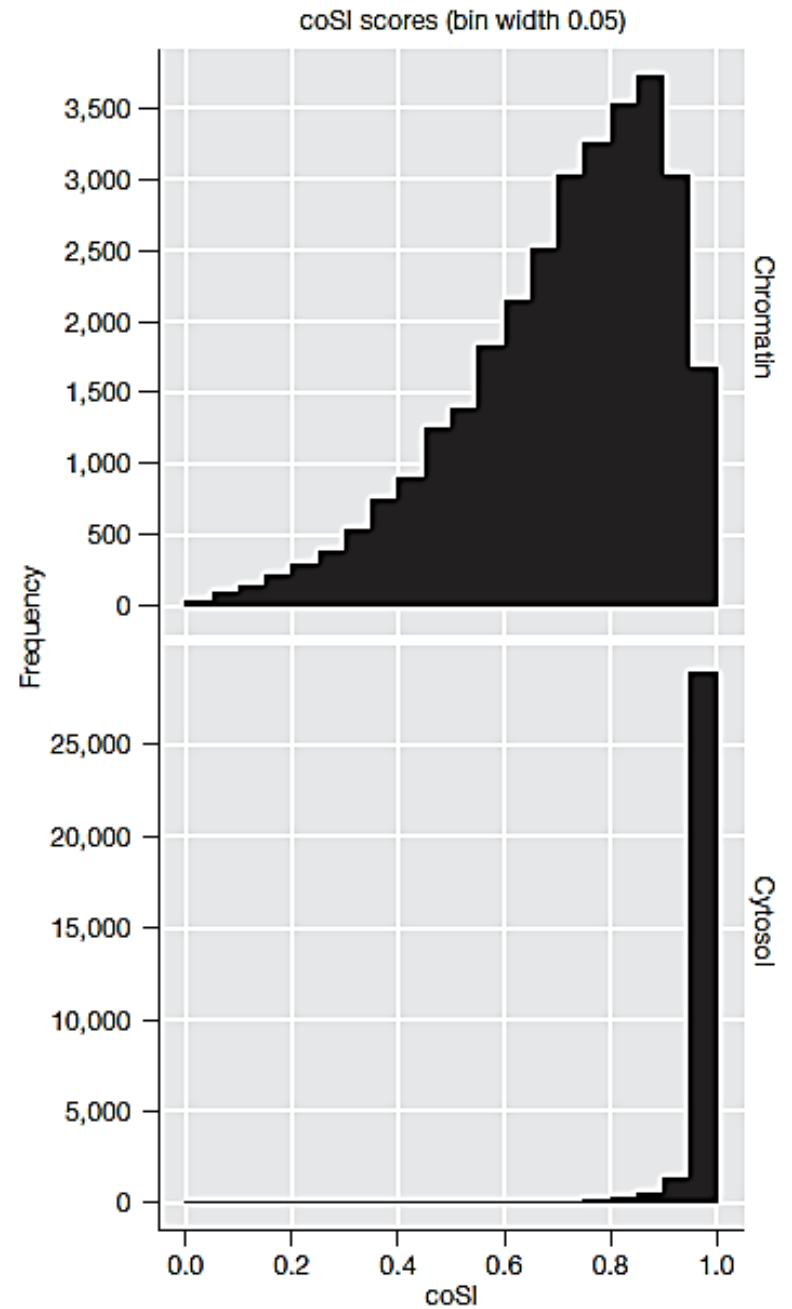
Figure 2  
Co-transcriptional splicing evaluation



coSI = the ratio between junction reads and exon-intron reads

Djebali et al., 20120

**b**



AS a matter of fact your Textbook, written by Alberto Kornblhitt, one of the major scientists in the field, bases discussion non mechanisms on splicing being mostly co-transcriptional.

Note:

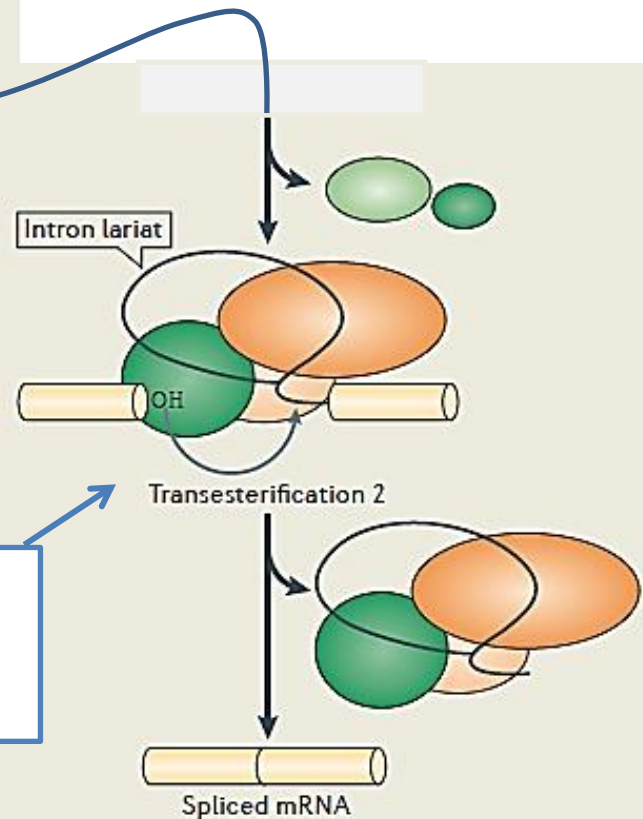
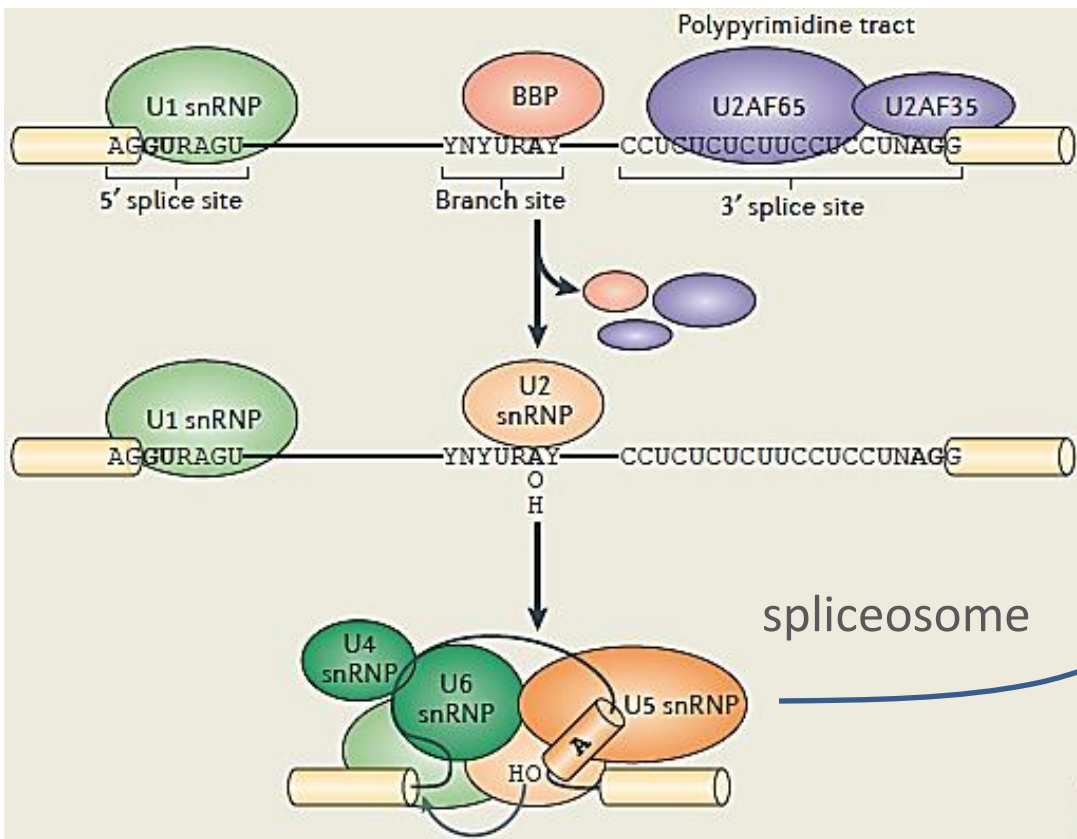
saying that splicing occurs co-transcriptionally does not mean that all catalytic events of splicing occur immediately as soon as the RNA progressively emerges from RNA PolIII.



# Alternative splicing: a pivotal step between eukaryotic transcription and translation

*Alberto R. Kornblihtt, Ignacio E. Schor, Mariano Alló, Gwendal Dujardin, Ezequiel Petrillo and Manuel J. Muñoz*

**Abstract** | Alternative splicing was discovered simultaneously with splicing over three decades ago. Since then, an enormous body of evidence has demonstrated the prevalence of alternative splicing in multicellular eukaryotes, its key roles in determining tissue- and species-specific differentiation patterns, the multiple post- and co-transcriptional regulatory mechanisms that control it, and its causal role in hereditary disease and cancer. The emerging evidence places alternative splicing in a central position in the flow of eukaryotic genetic information, between transcription and translation, in that it can respond not only to various signalling pathways that target the splicing machinery but also to transcription factors and chromatin structure.



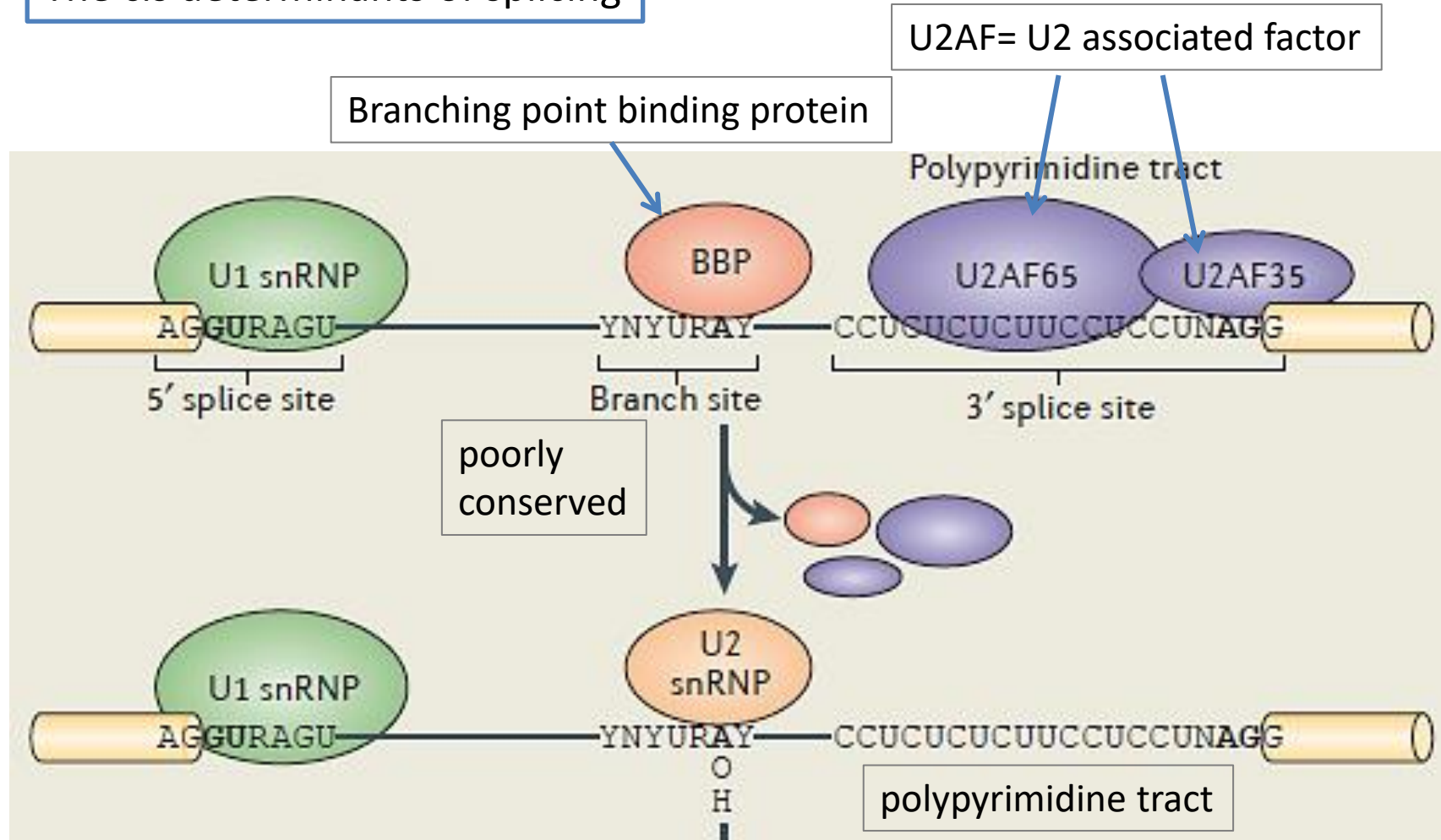
1st transesterification: the 2'-OH of the branching point «A» attacks the phosphodiester bond between exon and intron

2nd transesterification: the 3'-OH of the exon attacks the phosphodiester bond between intron and exon

**Table 1. Core Subunits of Human U snRNPs**

U snRNP	subunit gene	common subunit name(s) <sup>a</sup>	molecular mass (kDa) <sup>b</sup>	% U snRNP	recognizable domain/functional site <sup>c</sup>
U1 (248.1 kDa)	RNU1	U1 snRNA	53.5	21.6	
	SNRPB, -B2, -D1, -D2, -D3, -E, -F, -G	seven Sm proteins	94.3	38.0	Sm
	SNRNPA	U1-A	31.3	12.6	RRM
	SNRNP70	U1-70k	51.6	20.8	RRM; SR repeat
	SNRNPC	U1-C	17.4	7.0	Znf
U2 (987.4 kDa)	RNU2	U2 snRNA	61.2	6.2	
	SNRPB, -B2, -D1, -D2, -D3, -E, -F, -G	seven Sm proteins	94.3	9.6	Sm
	SNRPA1	U2A'	28.4	2.9	LRR
	SNRPB2	U2B''	25.4	2.6	RRM
	SF3A1	SF3a120	88.9	9.0	SWAP; UBQ domain
	SF3A2	SF3a66	49.3	5.0	Znf
	SF3A3	SF3a60	58.6	5.9	Znf; SAP
	SF3B1	SF3b155	145.8	14.8	HEAT repeat
	SF3B2	SF3b145	100.2	10.1	SAP
	SF3B3	SF3b130	135.5	13.7	DExH/D
	SF3B4	SF3b49	44.4	4.5	RRM
	SF3B5	SF3b10	10.1	1.0	
	SF3B14	SF3b14a; p14	14.6	1.5	RRM
	PHF5A	SF3b14b; Rds3	12.4	1.3	PHD-like
	DDX46	DDX46; hPrp5p	117.4	11.9	DExH/D; SR repeat
	SMNDC1	SPF30/SMNrp	26.7	2.7	Tudor domain
	U5 (1055.7 kDa)	RNU5	U5 snRNA	37.6	3.6
SNRPB, -B2, -D1, -D2, -D3, -E, -F, -G		seven Sm proteins	94.3	8.9	Sm
TXNL4A		U5-15K	16.9	1.6	TRX
SNRNP40		U5-40K	39.3	3.7	WD40
CD2BP2		U5-52K	37.6	3.6	GYF
DDX23		U5-100K; hPrp28	95.6	9.1	DExH/D; SR repeat
PRPF6		U5-102K; hPrp6	106.9	10.1	HAT/TPR repeats
EFTUD2		U5-116K; hSnu114	109.4	10.4	EF2-like fold; GTPase
SNRNP20		U5-200K; hBrr2	244.5	23.2	DExH/D
PRPF8		U5-220k; hPrp8	273.6	25.9	RNase H-fold; RRM; Jab1/MPN
U4/U6 (589.1 kDa)		RNU4	U4 snRNA	46.9	8.0
	RNU6	U6 snRNA	34.6	5.9	
	SNRPB, -B2, -D1, -D2, -D3, -E, -F, -G	seven Sm proteins	94.3	16.0	Sm
	LSM2, -3, -4, -5, -6, -7, -8	seven LSm proteins	78.9	13.4	Sm
	NHP2L1	15.5K	14.2	2.4	
	PPIH	U4/U6-20K; SnuCyp-20	19.2	3.3	cyclophilin-like
	PRPF31	U4/U6-61K; hPrp31	55.5	9.4	Nop
	PRPF4	U4/U6-60K; hPrp4	58.4	9.9	WD40
	PRPF3	U4/U6-90K; hPrp3	77.5	13.1	PWI
	SART3	p110; SART3; hPrp24	109.6	18.6	HAT repeats; RRM

# The cis determinants of splicing



This is GT-AG introns (by far the most frequent)

A secondary type exist (xx-xx), requiring U11 and U12 snRNP

<b>Code</b>	<b>Represents</b>	<b>Complement</b>
A	Adenine	T
G	Guanine	C
C	Cytosine	G
T	Thymine	A
Y	Pyrimidine (C or T)	R
R	Purine (A or G)	Y
W	weak (A or T)	W
S	strong (G or C)	S
K	keto (T or G)	M
M	amino (C or A)	K
D	A, G, T (not C)	H
V	A, C, G (not T)	B
H	A, C, T (not G)	D
B	C, G, T (not A)	V
X/N	any base	X/N
-	Gap	-

DNA base code

## Alternative splicing

Sequences at the borders of exon-intron and within the intron are similar but can vary.

Splice sites can be **strong** or **weak** depending on how far their sequences diverge from the consensus sequence.

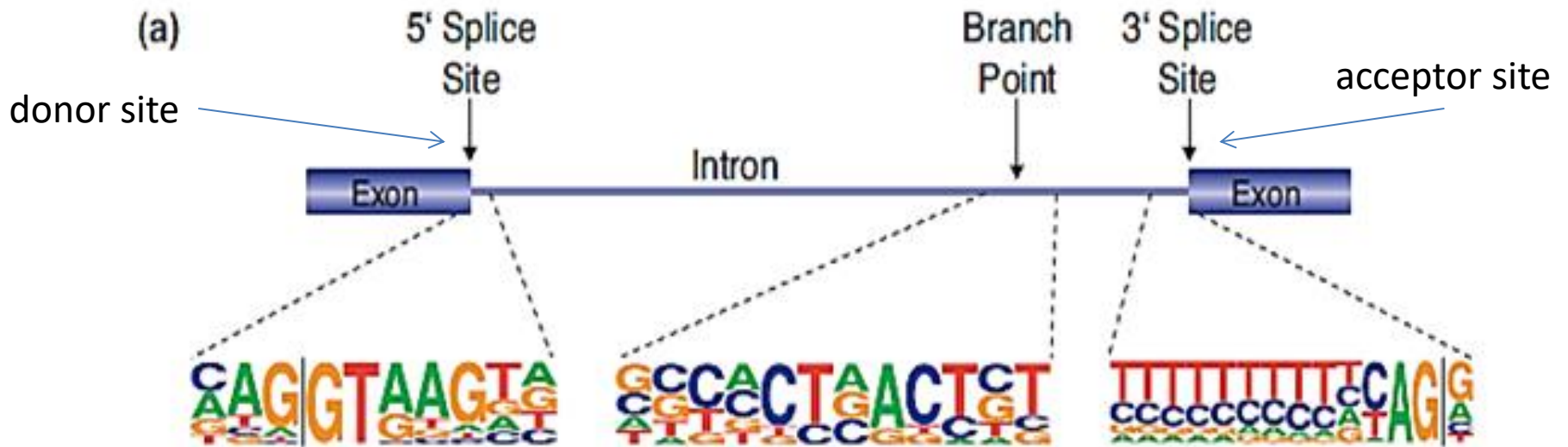
This determines their affinities for cognate splicing factors

In general, strong splice sites lead to constitutive splicing and full usage of the site

*Pay attention to this concept, it will ground the discussion on PolII elongation rate as determinant of AS*



The first chance to obtain regulation derives from how exons are recognized



Variations of these sequences can give «stronger» and «weaker» splicing sites

Indeed, **in addition** to the sequences directly regulating binding of spliceosome components, in both Exons and Introns sequence motifs exist that regulate the use of splice sites.

Named after their location and effect:

ESE: exonic splicing enhancers

ESS: exonic splicing silencers

ISE: intronic splicing enhancers

ISS: intronic splicing silencers

They are therefore ***cis-elements*** for splicing regulation



**Trans-acting factors** of splicing are proteins binding to these elements.

They belong to the general class of RNA Binding Proteins (**RBP**)

Three categories:

1. SR proteins and SR-like
2. hnRNP
3. tissue-specific and context-specific factors

## **SR proteins = splicing regulators**

### Domains:

The most typical domain is an alternating Arginine-Serine-rich domain, called “RS domain”: it is a protein-protein interaction domain.

### Regulation:

SR are phosphorylated at Ser by several kinases → regulates interaction with each other and with other proteins.

### Other interactants:

SR proteins also interact with the CAP-binding protein and with poly-A binding proteins.

### Binding sites:

Mostly at Exons, sometimes also to ISE (intronic splicing enhancers)

### Activity:

Mostly activatory toward the most proximal exon.

Exon definition.

Canonical SR proteins

Name*	Domains	Binding sequence	Target genes
<i>Canonical SR proteins</i>			
SRp20 (SFRS3)	RRM and RS	GCUCCUCUUC	SRP20, CALCA and INSR
SC35 (SFRS2)	RRM and RS	UGCUGUU	ACHE and GRIA1–GRIA4
ASF/SF2 (SFRS1)	RRM, RRMH and RS	RGAAGAAC	HIPK3, CAMK2D, HIV RNAs and GRIA1–GRIA4
SRp40 (SFRS5)	RRM, RRMH and RS	AGGAGAAGGGA	HIPK3, PRKCB and FN1
SRp55 (SFRS6)	RRM, RRMH and RS	GGCAGCACCUG	TNNT2 and CD44
SRp75 (SFRS4)	RRM, RRMH and RS	GAAGGA	FN1, E1A and CD45
9G8 (SFRS7)	RRM, zinc finger and RS	(GAC) <sub>n</sub>	TAU, GNRH and SFRS7
SRp30c (SFRS9)	RRM, RRMH and RS	CUGGAUU	BCL2L1, TAU and HNRNPA1
SRp38 (FUSIP1)	RRM and RS	AAAGACAAA	GRIA2 and TRD
<i>Other SR proteins</i>			
SRp54	RRM and RS	ND	TAU
SRp46 (SFRS2B)	RRM and RS	ND	NA
RNPS1	RRM and Ser-rich	ND	TRA2B
SRp35	RRM and RS	ND	NA
SRp86 (SRp508 and SFRS12)	RRM and RS	ND	NA
TRA2α	RRM and two Arg-rich	GAAARGARR	dsx
TRA2β	RRM and two RS	(GAA) <sub>n</sub>	SMN1, CD44 and TAU
RBM5	RRM and RS	ND	CD95
CAPER (RBM39)	RRM and RS	ND	VEGF

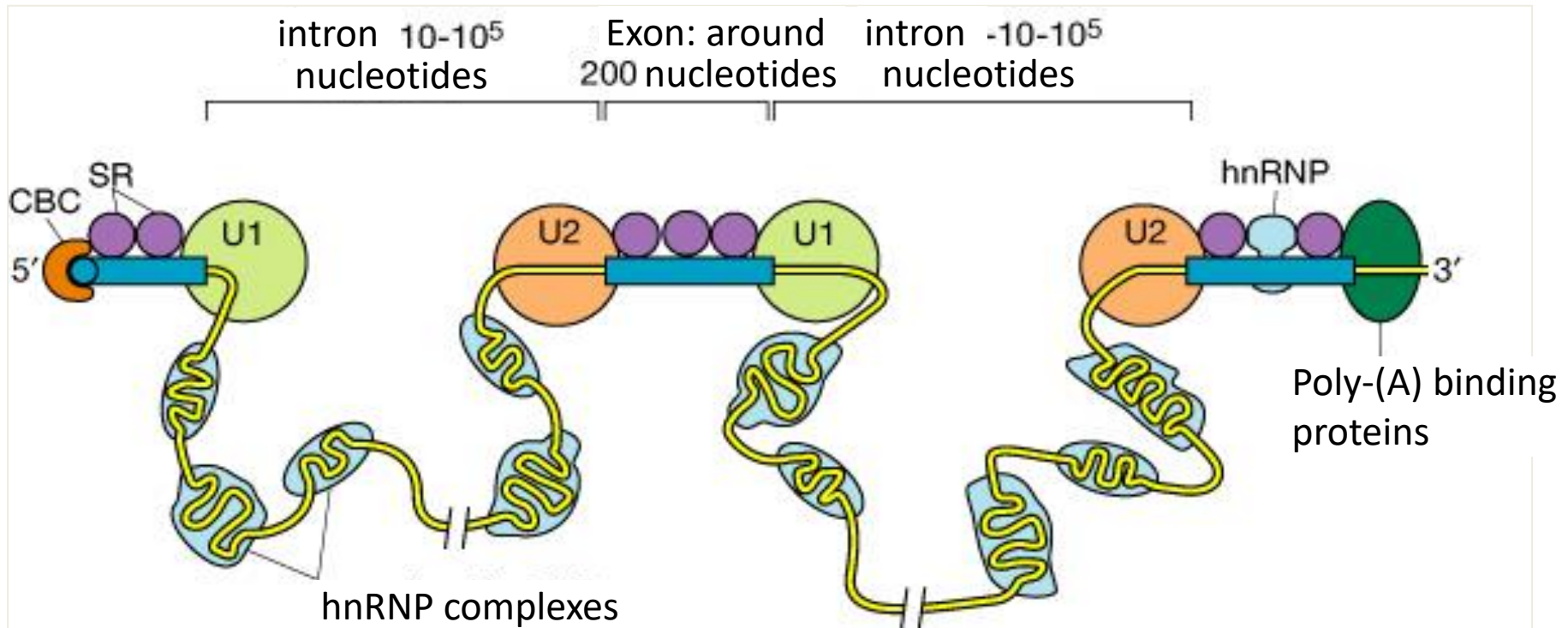
SR-like  
(other protein containing an RS domain)

## **hnRNP proteins** (heterogeneous nuclear Ribo Nucleic Protein)

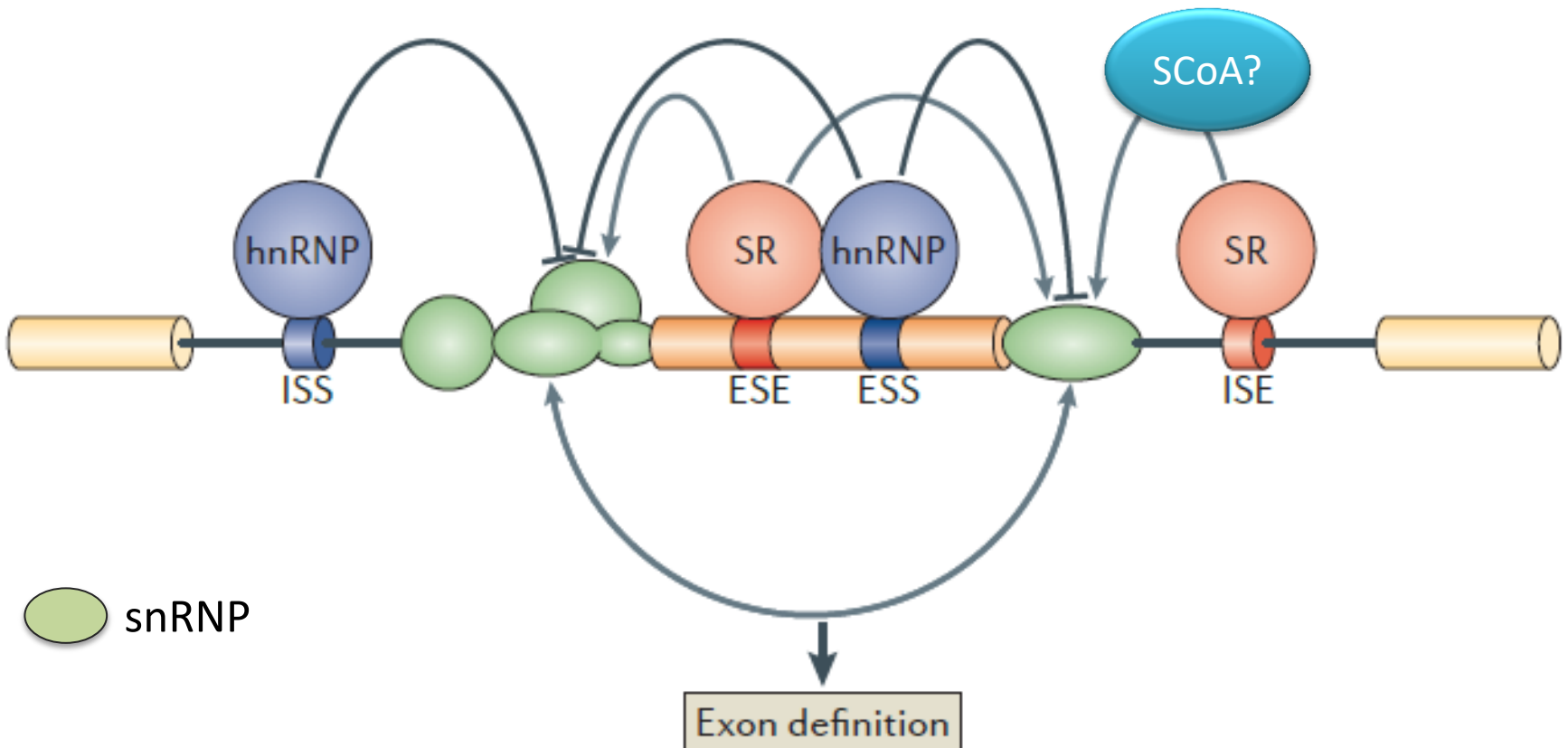
- Many different families
- Usually bind intronic sites
- Intron definition
- Several other roles have been ascribed to individual members, e.g. cytoplasmic localization.

**Table 1 | Ribonucleoproteins that are involved in pre-mRNA splicing**

Name	Other names	Domains*	Binding sequences	Target genes
hnRNP A1	NA	RRM, RGG and G	UAGGGA/U	SMN2 and RAS
hnRNP A2	NA	RRM, RGG and G	(UUAGGG) <sub>n</sub>	HIV <i>tat</i> and IKBKAP
hnRNP B1				
hnRNP C1	AUF1	RRM	U rich	APP
hnRNP C2				
hnRNP F	NA	RRM, RGG and GY	GGGA and G rich	PLP, SRC and BCL2L2
hnRNP G	NA	RRM and SRGY	CC(A/C) and AAGU	SMN2 and TMP1
hnRNP H	DSEF1	RRM, RGG, GYR and GY	GGGA and G rich	PLP, HIV <i>tat</i> and BCL2L1
hnRNP H'				
hnRNP I	PTB	RRM	UCUU and CUCUCU	PTB, nPTB, SRC, CD95, TINT2, CALCA and GRIN3B
hnRNP L	NA	RRM	C and A rich	NOS and CD45
hnRNP LL	SRRF	RRM	C and A rich	CD45
hnRNP M	NA	RRM and GY	ND	FGFR2
hnRNP Q	NA	RRM and RGG	ND	SMN2



The effect of SR and hnRNP binding is either to stabilize or destabilize the interaction of basal splicing factors (snRNPs) with the splicing sites  
This action can be direct protein-protein contact, or mediated by splicing co-activators.



Albeit SR protein can be regulated by signalling pathways (*e.g. by phosphorylation, and several examples are given in your Kornblitt Textbook*), as well as some hnRNP proteins, the fact that most SRs and most hnRNPs have ubiquitous expression suggests that additional tissue-specific factors should be involved in tissue-specific AS.

Hence, additional cis-elements could be present in regulated pre-mRNAs.

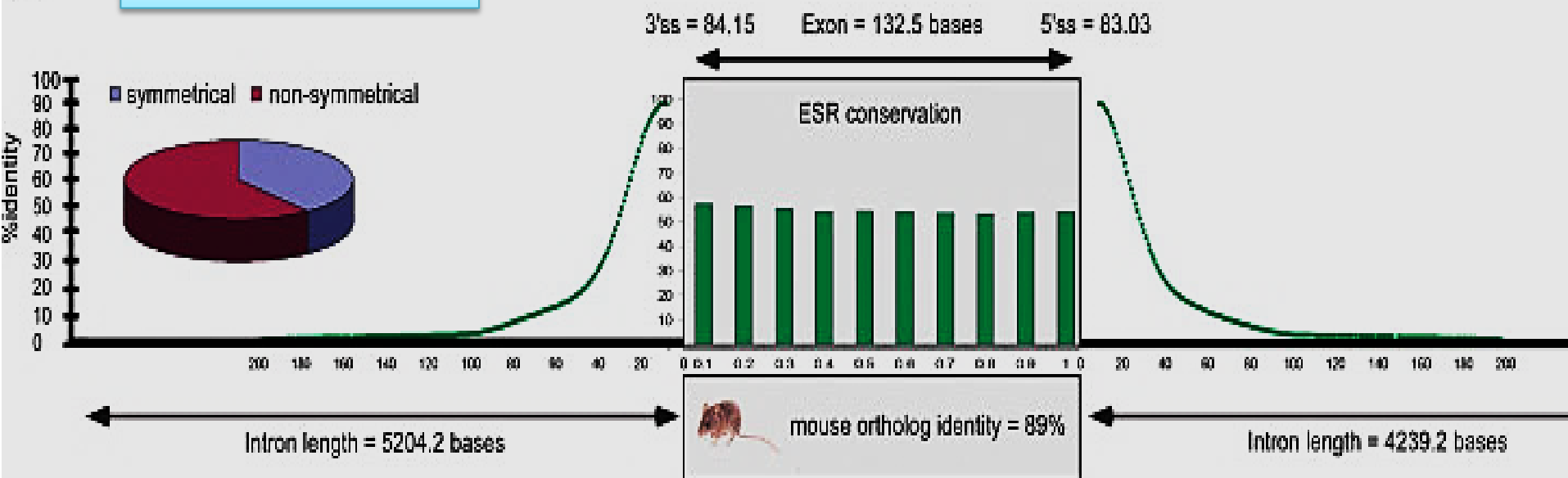
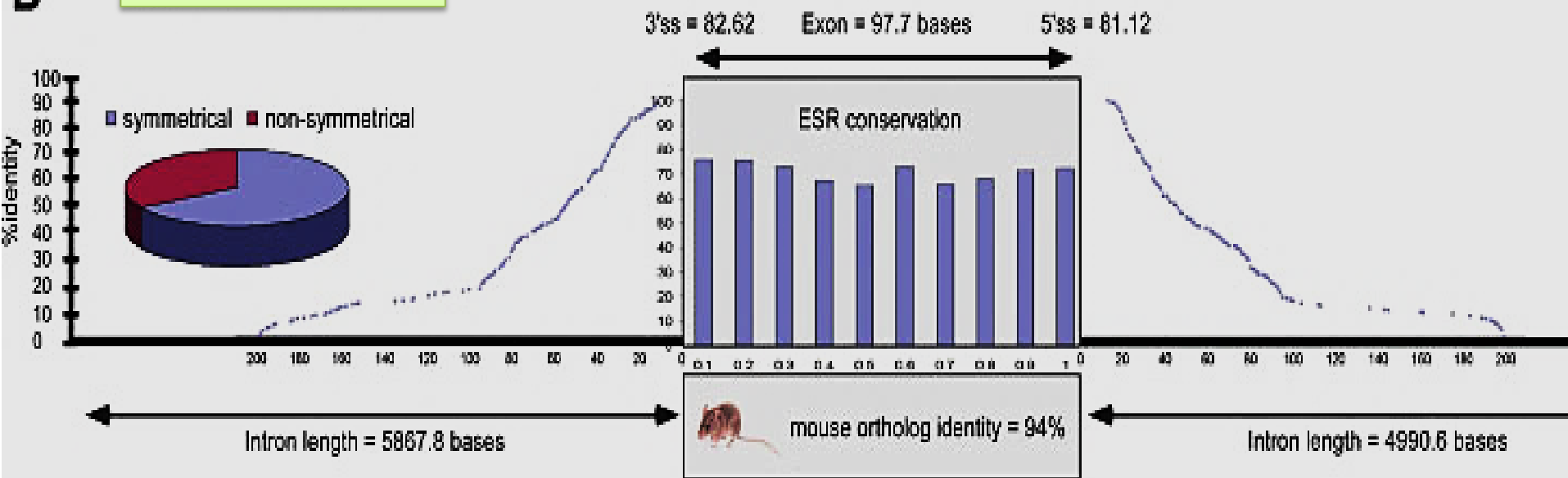
One way to answer this question is to compare sequences around **alternative** exons to those of **constitutive** exons.



Some Authors used comparative genomics to obtain insight. Human and mouse transcriptomes are quite well characterized, making it possible to classify exons as constitutive or alternative based on real expression data (microarrays, RNA-seq).

Sequences were then compared. Exons were normalized in length and flanking introns were explored within 200 pb.

From Kim et al., BioEssays 30:38–47

**A****Constitutive - strong****B****Alternative - weak**

## 2.1 Features of **Alternative** versus **Constitutive** Exons

### More conserved

Exons that are alternatively spliced in both human and mouse are **more conserved** than constitutive exons

Conservation is higher toward exon edges and extends farther in introns:  
**cis-regulatory sequences ?**

### Weak splice sites

Cassette exons<sup>(1)</sup> have **weak splice sites**, compared to the strong ones in constitutive exons

### Shorter

Alternative cassette exons are also **shorter** and are flanked by longer introns than constitutively spliced ones.

### Symmetric

The percentage of **symmetrical** exon is definitely higher in alternative exons (symmetrical means “divisible-by-three” number of base pairs)

<sup>(1)</sup>*Cassette exons: refers to exon skipping*

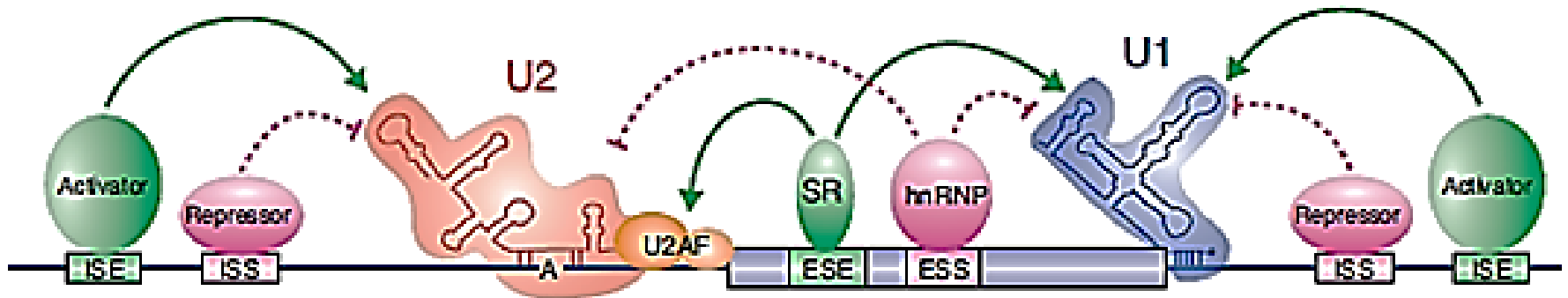
*From Kim et al., 2007, BioEssays 30:38–47*

All this points to the existence of additional AS regulators

# Regulatory

- Tissue-specific splicing
- Regulated splicing
- Epigenetic establishment of splicing patterns

- 1) Tissue-specific splicing factors
- 2) Signal transduction regulated factors
- 3) Chromatin effects on splicing choice



*Major SR proteins and hnRNP can hardly explain tissue-specific splicing*

How to identify tissue-specific AS regulators

## Appendix: The search for tissue-specific splicing factors

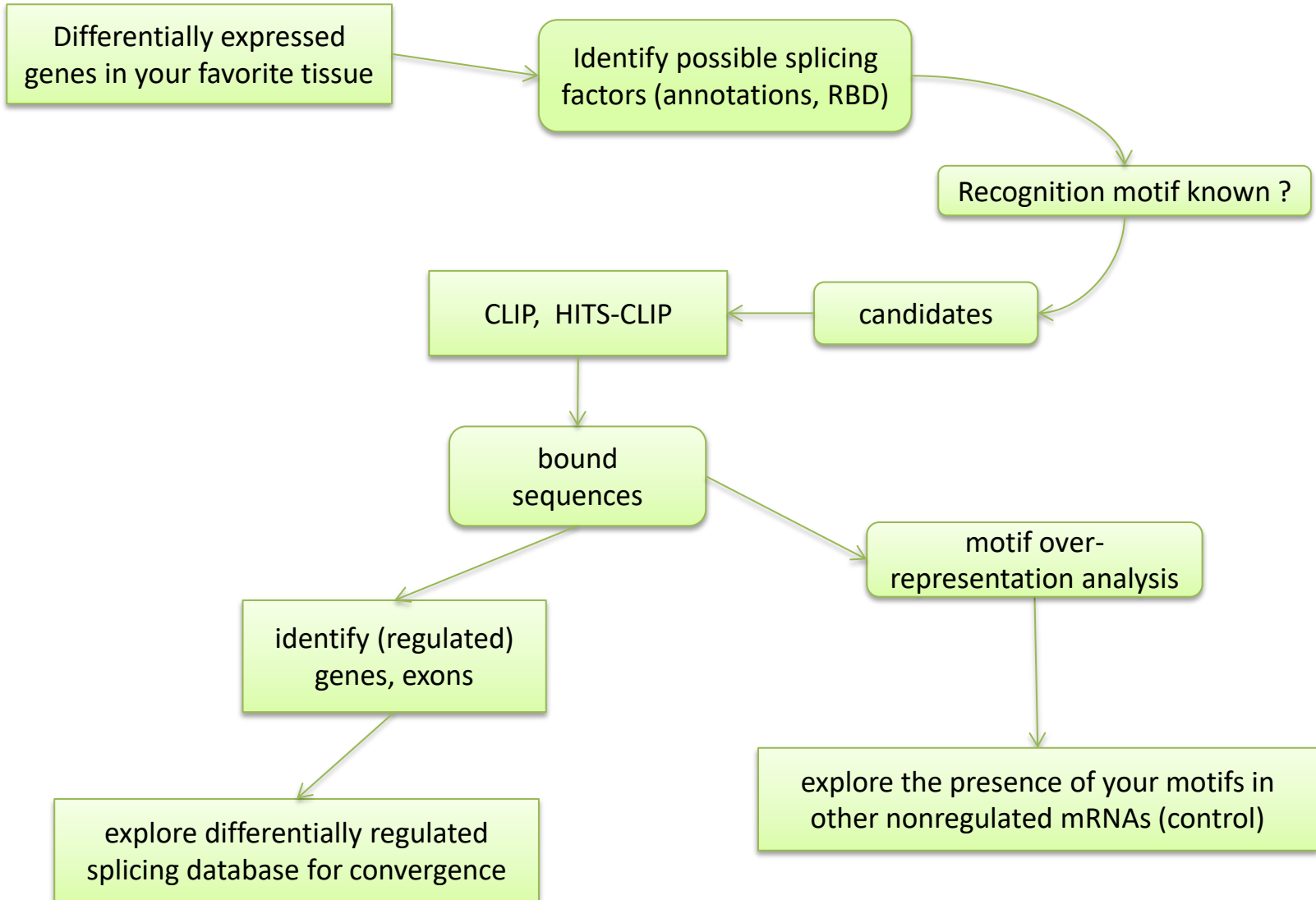


Table 2 | **Tissue-specific alternative splicing factors**

Name	Other names	Binding domain	Binding motif	Tissue expression	Target genes
nPTB	brPTB and PTBP2	RRM	CUCUCU	Neurons, myoblasts and testes	<i>BIN1</i> , <i>GLYRA2</i> , <i>ATP2B1</i> , <i>MEF2</i> , <i>NASP</i> , <i>SPAG9</i> and <i>SRC</i>
NOVA1	NA	KH	YCA Y	Neurons of the hindbrain and spinal cord	<i>GABRG2</i> , <i>GLYRA2</i> and <i>NOVA1</i>
NOVA2	NA	KH	YCA Y	Neurons of the cortex, hippocampus and dorsal spinal cord	<i>KCNJ</i> , <i>APLP2</i> , <i>GPHN</i> , <i>JNK2</i> , <i>NEO</i> , <i>GRIN1</i> and <i>PLCB4</i>
FOX1	A2BP1	RRM	(U)GCAUG	Muscle, heart and neurons	<i>ACTN</i> , <i>EWSR1</i> , <i>FGFR2</i> , <i>FN1</i> and <i>SRC</i>
FOX2	RBM9	RRM	(U)GCAUG	Muscle, heart and neurons	<i>EWS</i> , <i>FGFR2</i> , <i>FN1</i> and <i>SRC</i>
RBM35a	ESRP1	RRM	GU rich	Epithelial cells	<i>FGFR2</i> , <i>CD44</i> , <i>CTNND1</i> and <i>ENAH</i>
RBM35b	ESRP2	RRM	GU rich	Epithelial cells	<i>FGFR2</i> , <i>CD44</i> , <i>CTNND1</i> and <i>ENAH</i>
TIA1	mTIA1	RRM	U rich	Brain, spleen and testes	<i>MYPT1</i> , <i>CD95</i> , <i>CALCA</i> , <i>FGFR2</i> , <i>TIAR</i> , <i>IL8</i> , <i>VEGF</i> , <i>NF1</i> and <i>COL2A1</i>
TIAR	TIAL1 and mTIAR	RRM	U rich	Brain, spleen, lung, liver and testes	<i>TIA1</i> , <i>CALCA</i> , <i>TIAR</i> , <i>NF1</i> and <i>CD95</i>
SLM2	KHDRBS3 and TSTAR	KH	UAAA	Brain, tests and heart	<i>CD44</i> and <i>VEGFA</i>
Quaking	QK and QKL	KH	ACUAAY[...]UAAY	Brain	<i>MAG</i> and <i>PLP</i>
HUB	HUC, HUD and ELAV2	RRM	AU rich	Neurons	<i>CALCA</i> , <i>CD95</i> and <i>NF1</i>



MBNL	NA	CCCH zinc finger domain	YGCU(U/G)Y	Muscles, uterus and ovaries	<i>TNTT2</i> , <i>INSR</i> , <i>CLCN1</i> and <i>TNNT3</i>
CELF1	BRUNOL2	RRM	U and G rich	Brain	<i>TNTT2</i> and <i>INSR</i>
ETR3	CELF2 and BRUNOL3	RRM	U and G rich	Heart, skeletal muscle and brain	<i>TNTT2</i> , <i>TAU</i> and <i>COX2</i>
CELF4	BRUNOL4	RRM	U and G rich	Muscle	<i>MTMR1</i> and <i>TNTT2</i>
CELF5	BRUNOL5 and NAPOR	RRM	U and G rich	Heart, skeletal muscle and brain	<i>ACTN</i> , <i>TNTT2</i> and <i>GRIN1</i>
CELF6	BRUNOL6	RRM	U and G rich	Kidney, brain and testes	<i>TNTT2</i>

*A2BP1*, ataxin 2-binding protein 1; *ACTN*,  $\alpha$ -actinin; *APLP2*, amyloid- $\beta$  precursor-like protein 2; *ATP2B1*, ATPase, Ca<sup>2+</sup> transporting, plasma membrane 1; *BIN1*, bridging integrator 1; *CALCA*, calcitonin-related polypeptide- $\alpha$ ; *CELF*, CUGBP- and ETR3-like factor; *CLCN1*, chloride channel 1; *COL2A1*, collagen, type II,  $\alpha$ 1; *COX2*, cytochrome c oxidase II; *CTNND1*, catenin  $\delta$ 1; *EWSR1*, Ewing sarcoma breakpoint region 1; *FGFR2*, fibroblast growth factor receptor 2; *FN1*, fibronectin 1; *GABRG2*, GABA A receptor,  $\gamma$ 2; *GLYRA2*, glycine receptor,  $\alpha$ 2 subunit; *GPHN*, gephyrin; *GRIN1*, glutamate receptor, ionotropic, NMDA 3B; *IL8*, interleukin-8; *INSR*, insulin receptor; *JNK2*, Jun N-terminal kinase 2; *KCNJ*, potassium inwardly-rectifying channel, subfamily; *KHDRBS3*, KH domain-containing, RNA-binding, signal transduction-associated protein 3; *MAG*, myelin associated glycoprotein; *MBNL*, muscleblind; *MEF2*, myocyte enhancing factor 2; *MTMR1*, myotubularin-related protein 1; *NASP*, nuclear autoantigenic sperm protein; *NEO*, neogenin; *NF1*, neurofibromin 1; *NOVA*, neuro-oncological ventral antigen; *PLCB4*, phospholipase C  $\beta$ 4; *PLP*, proteolipid protein; *PTB*, polypyrimidine-tract binding protein; *RBM*, RNA-binding protein; *RRM*, RNA recognition motif; *SLM2*, SAM68-like mammalian protein 2; *SPAG9*, sperm associated antigen 9; *TIA1*, T cell-restricted intracellular antigen 1; *TIAR*, TIA1-related protein; *TNTT2*, troponin T type 2; *VEGF*, vascular endothelial growth factor.

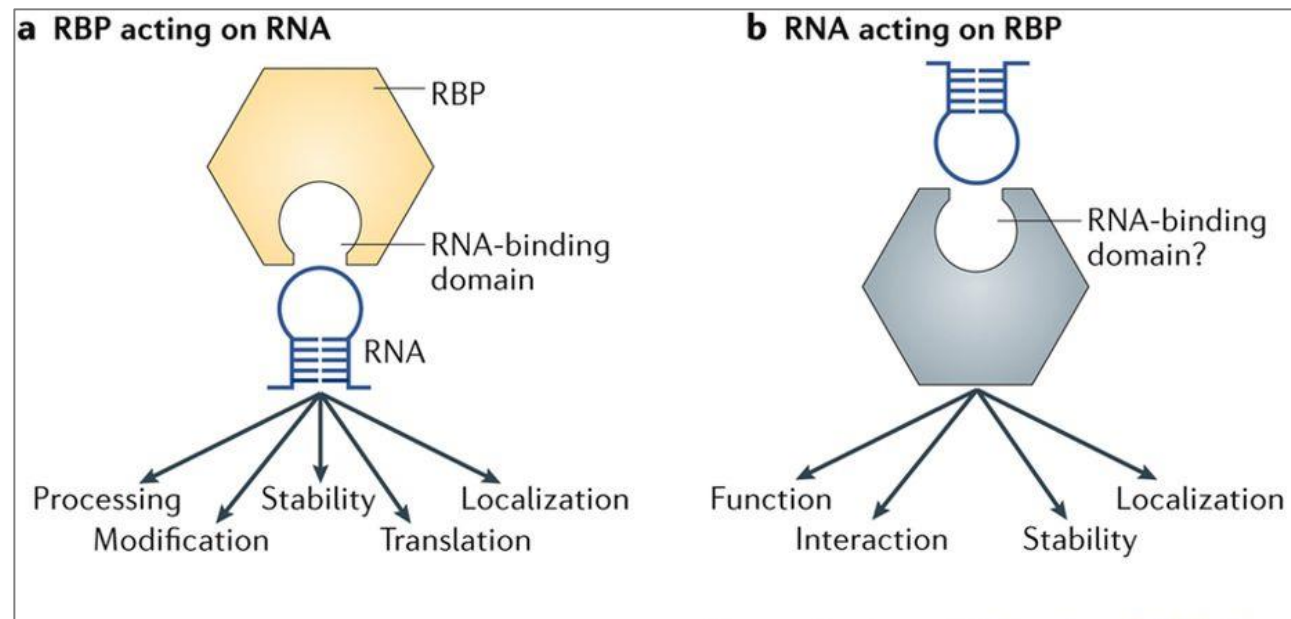
RNA Binding Proteins are a growing class ...

Generic name: RBP (RNA binding Proteins) GO category: RNA-Binding

Recent studies used RIC (RNA Interactome Capture) identified an exceptional number of RBPs (860 from HeLa and 791 from HEK293).

Many of these do not carry any of the known domains:

- RRM - RNA Recognition Motif
- KH
- DEAD box helicase
- Zn-fingers motifs

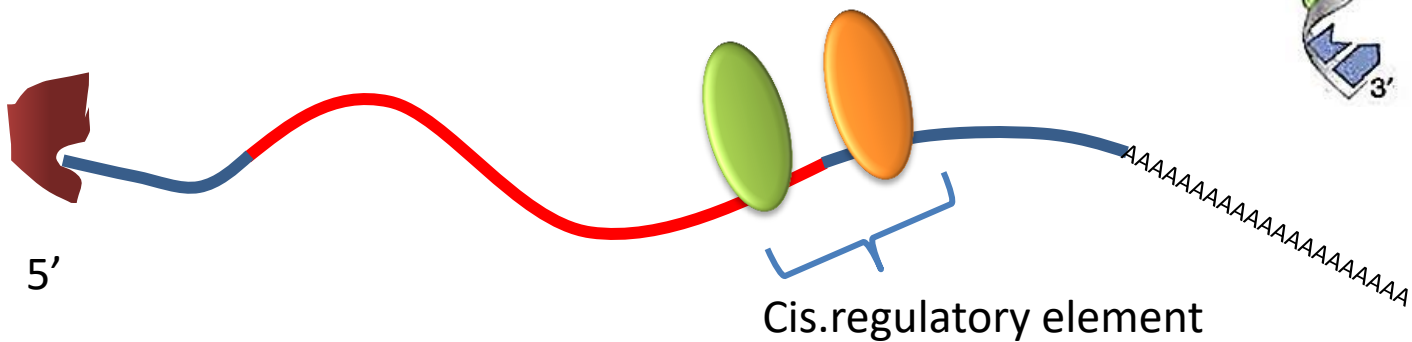
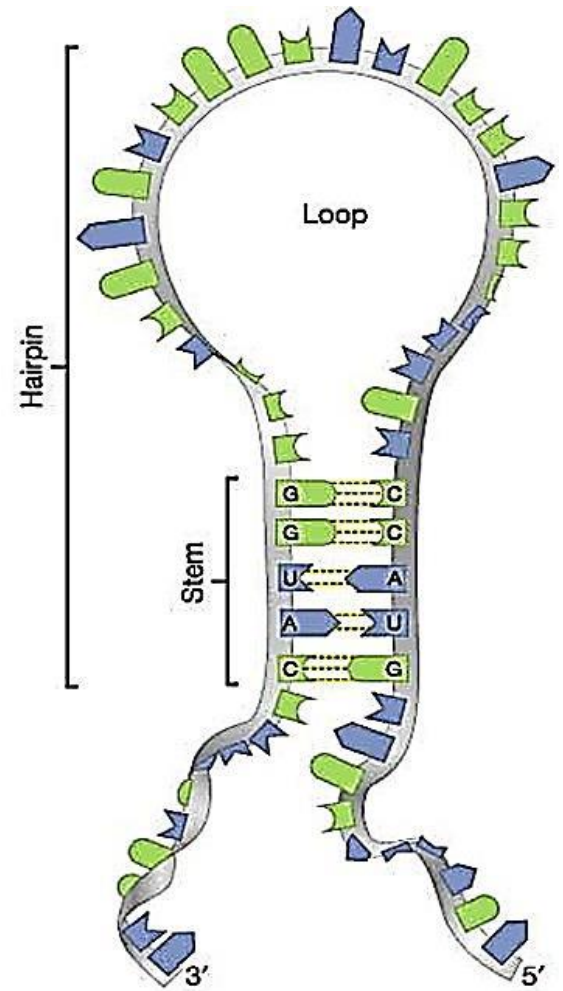


## Identification of RNA binding proteins motifs

Specificity of RNA binding: both «sequence» and «structure» elements

Problems in predicting regulatory motifs:

- Localization (intron length)
- Sometimes dispersed elements
- Sometimes the structural component prevails upon pure sequence



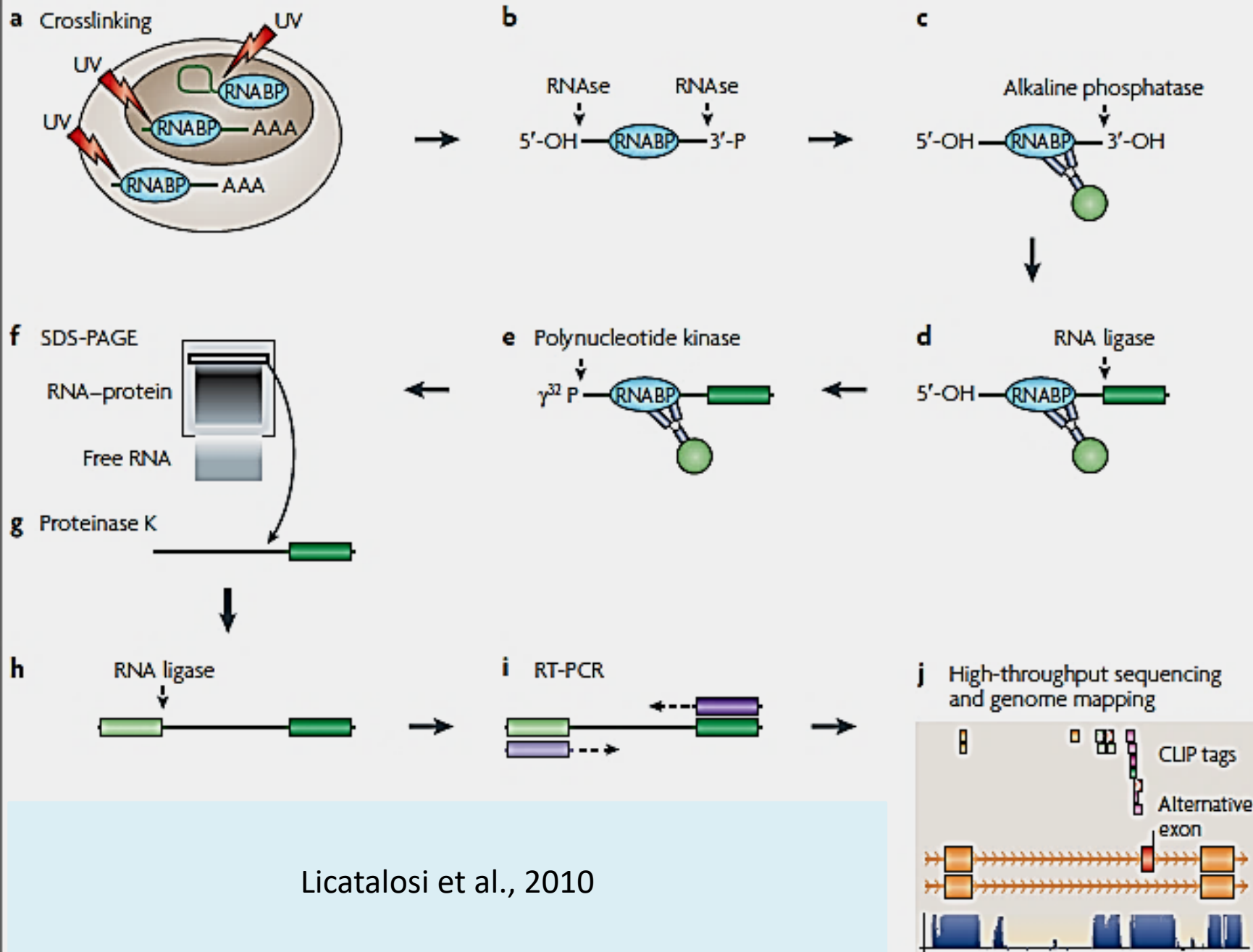
RRM: Sex-lethal <http://www.ebi.ac.uk/thornton-srv/databases/cgi-bin/pdbsum/GetPage.pl?pdbcode=1b7f>

RRM: PTB

<http://www.ebi.ac.uk/thornton-srv/databases/cgi-bin/pdbsum/GetPage.pl?pdbcode=2adc>

KH: NOVA1 <http://www.ebi.ac.uk/thornton-srv/databases/cgi-bin/pdbsum/GetPage.pl?pdbcode=1dt4>

## Box 2 | CLIP and HITS-CLIP methods



Example:

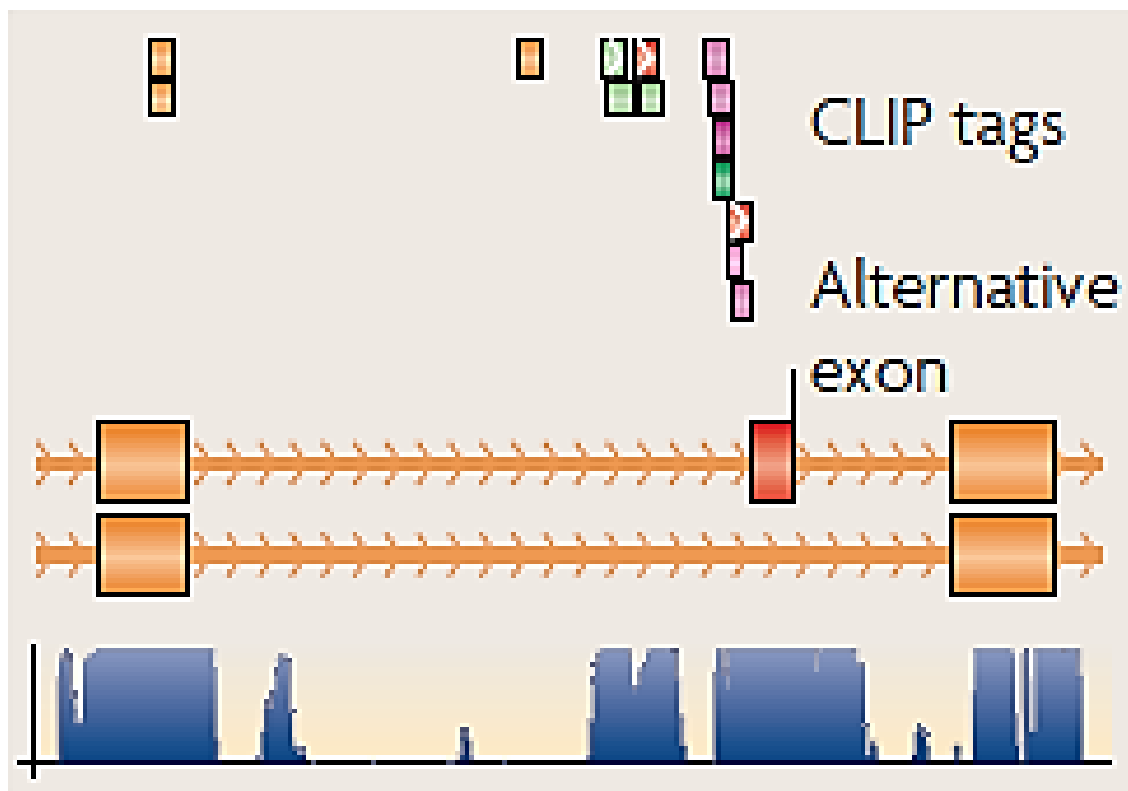
Individual Splicing Factors  
HITS-CLIP profiles mapped  
to genome and compared  
to RNA-Seq profiles.

## High-throughput sequencing and genome mapping

CLIP-seq reads  
mapping

Reference

RNA-seq reads



# The mechanisms of alternative splicing regulation