

Ch 4 – L1.2

ENCODE transcriptome

An integrated encyclopedia of DNA elements in the human genome

The ENCODE Project Consortium*

This is the leading article that describes all the ENCODE project and gives a overall resumé of results obtained in the 2nd phase.

The human genome encodes the blueprint of life, but the function of the vast majority of its nearly three billion bases is unknown. The Encyclopedia of DNA Elements (ENCODE) project has systematically mapped regions of transcription, transcription factor association, chromatin structure and histone modification. These data enabled us to assign biochemical functions for 80% of the genome, in particular outside of the well-studied protein-coding regions. Many discovered candidate regulatory elements are physically associated with one another and with expressed genes, providing new insights into the mechanisms of gene regulation. The newly identified elements also show a statistical correspondence to sequence variants linked to human disease, and can thereby guide interpretation of this variation. Overall, the project provides new insights into the organization and regulation of our genes and genome, and is an expansive resource of functional annotations for biomedical research.

6 SEPTEMBER 2012 | VOL 489 | NATURE | 57

ENCODE official website: <https://www.encodeproject.org/>

ENCODE at the NHGRI: <http://www.genome.gov/encode/>

Nature ENCODE: <http://www.nature.com/encode/#/threads>

Landscape of transcription in human cells

Sarah Djebali^{1*}, Carrie A. Davis^{2*}, Angelika Merkel¹, Alex Dobin², Timo Lassmann³, Ali Mortazavi^{4,5}, Andrea Tanzer¹, Julien Lagarde¹, Wei Lin², Felix Schlesinger², Chenghai Xue², Georgi K. Marinov⁴, Jainab Khatun⁶, Brian A. Williams⁴, Chris Zaleski², Joel Rozowsky^{7,8}, Maik Röder¹, Felix Kokocinski⁹, Rehab F. Abdelhamid³, Tyler Alioto^{1,10}, Igor Antoshechkin⁴, Michael T. Baer², Nadav S. Bar¹¹, Philippe Batut², Kimberly Bell², Ian Bell¹², Sudipto Chakraborty², Xian Chen¹³, Jacqueline Chrast¹⁴, Joao Curado¹, Thomas Derrien¹, Jorg Drenkow², Erica Dumais¹², Jacqueline Dumais¹², Radha Duttagupta¹², Emilie Falconnet¹⁵, Meagan Fastuca², Kata Fejes-Toth², Pedro Ferreira¹, Sylvain Foissac¹², Melissa J. Fullwood¹⁶, Hui Gao¹², David Gonzalez¹, Assaf Gordon², Harsha Gunawardena¹³, Cedric Howald¹⁴, Sonali Jha², Rory Johnson¹, Philipp Kapranov^{12,17}, Brandon King⁴, Colin Kingswood^{1,10}, Oscar J. Luo¹⁶, Eddie Park⁵, Kimberly Persaud², Jonathan B. Preall², Paolo Ribeca^{1,10}, Brian Risk⁶, Daniel Robyr¹⁵, Michael Sammeth^{1,10}, Lorian Schaffer⁴, Lei-Hoon See², Atif Shahab¹⁶, Jorgen Skancke^{1,11}, Ana Maria Suzuki³, Hazuki Takahashi³, Hagen Tilgner^{1†}, Diane Trout⁴, Nathalie Walters¹⁴, Huaien Wang², John Wrobel⁶, Yanbao Yu¹³, Xiaoran Ruan¹⁶, Yoshihide Hayashizaki³, Jennifer Harrow⁹, Mark Gerstein^{7,8,18}, Tim Hubbard⁹, Alexandre Reymond¹⁴, Stylianos E. Antonarakis¹⁵, Gregory Hannon², Morgan C. Giddings^{6,13}, Yijun Ruan¹⁶, Barbara Wold⁴, Piero Carninci³, Roderic Guigó^{1,19} & Thomas R. Gingeras^{2,12}

Eukaryotic cells make many types of primary and processed RNAs that are found either in specific subcellular compartments or throughout the cells. A complete catalogue of these RNAs is not yet available and their characteristic subcellular localizations are also poorly understood. Because RNA represents the direct output of the genetic information encoded by genomes and a significant proportion of a cell's regulatory capabilities are focused on its synthesis, processing, transport, modification and translation, the generation of such a catalogue is crucial for understanding genome function. Here we report evidence that three-quarters of the human genome is capable of being transcribed, as well as observations about the range and levels of expression, localization, processing fates, regulatory regions and modifications of almost all currently annotated and thousands of previously unannotated RNAs. These observations, taken together, prompt a redefinition of the concept of a gene.

from Djebali et al., 2012

.....

Here we report identification and characterization of annotated and novel RNAs that are enriched in either of the two major cellular subcompartments (nucleus and cytosol) for all **15 cell lines studied**, and in three additional subnuclear compartments in one cell line.

In addition, we have sought to determine whether identified transcripts are modified at their 5' and 3' termini by the presence of a 7-methyl guanosine cap or polyadenylation, respectively.

These results considerably extend the current genome-wide annotated catalogue of long polyadenylated and small RNAs collected by the **GENCODE** annotation group.

ENCODE - Transcriptome

Djebali et al., 2012

RNA-Seq: identification of annotated and novel RNAs from either of the two major cellular subcompartments (nucleus and cytosol) for 15 cell lines.

To see the EXPERIMENTAL GRID :

<http://genome.ucsc.edu/ENCODE/dataMatrix/encodeDataMatrixHuman.html>

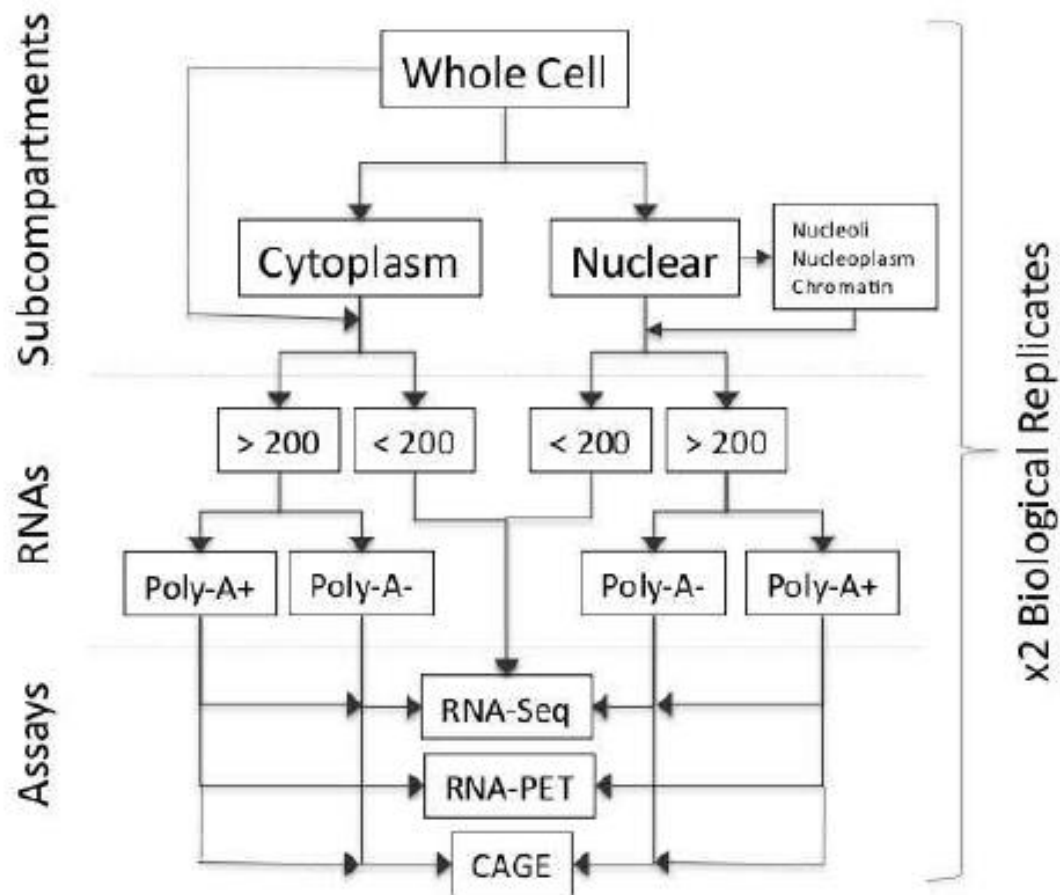
- 62.1% of genome covered by processed transcripts; 74.7% by unprocessed transcripts.
- Novel elements cover 78% of intronic nucleotides and 34% of intergenic sequences.
- Multiple isoforms per gene expressed simultaneously, with a plateau at 10-12 isoforms per gene per cell line.
- eRNA – transcripts starting from enhancers
- 6% of coding and noncoding overlap with small RNA (probably precursors)

Question: is this feature «conclusive» ?

RNA data set generation

We performed subcellular compartment fractionation (whole cell, nucleus and cytosol) before RNA isolation in 15 cell lines (Supplementary Table 1) to interrogate deeply the human transcriptome. For the K562 cell line, we also performed additional nuclear subfractionation into chromatin, nucleoplasm and nucleoli. The RNAs from each of these subcompartments were prepared in replica and were separated based on length into >200 nucleotides (long) and <200 nucleotides (short). Long RNAs were further fractionated into polyadenylated and non-polyadenylated transcripts. A number of complementary technologies were used to characterize these RNA fractions as to their sequence (RNA-seq), sites of initiation of transcription (cap-analysis of gene expression (CAGE)⁹) and sites of 5' and 3' transcript termini (paired end tags (PET)¹⁰; Supplementary Fig. 1). Sequence reads were

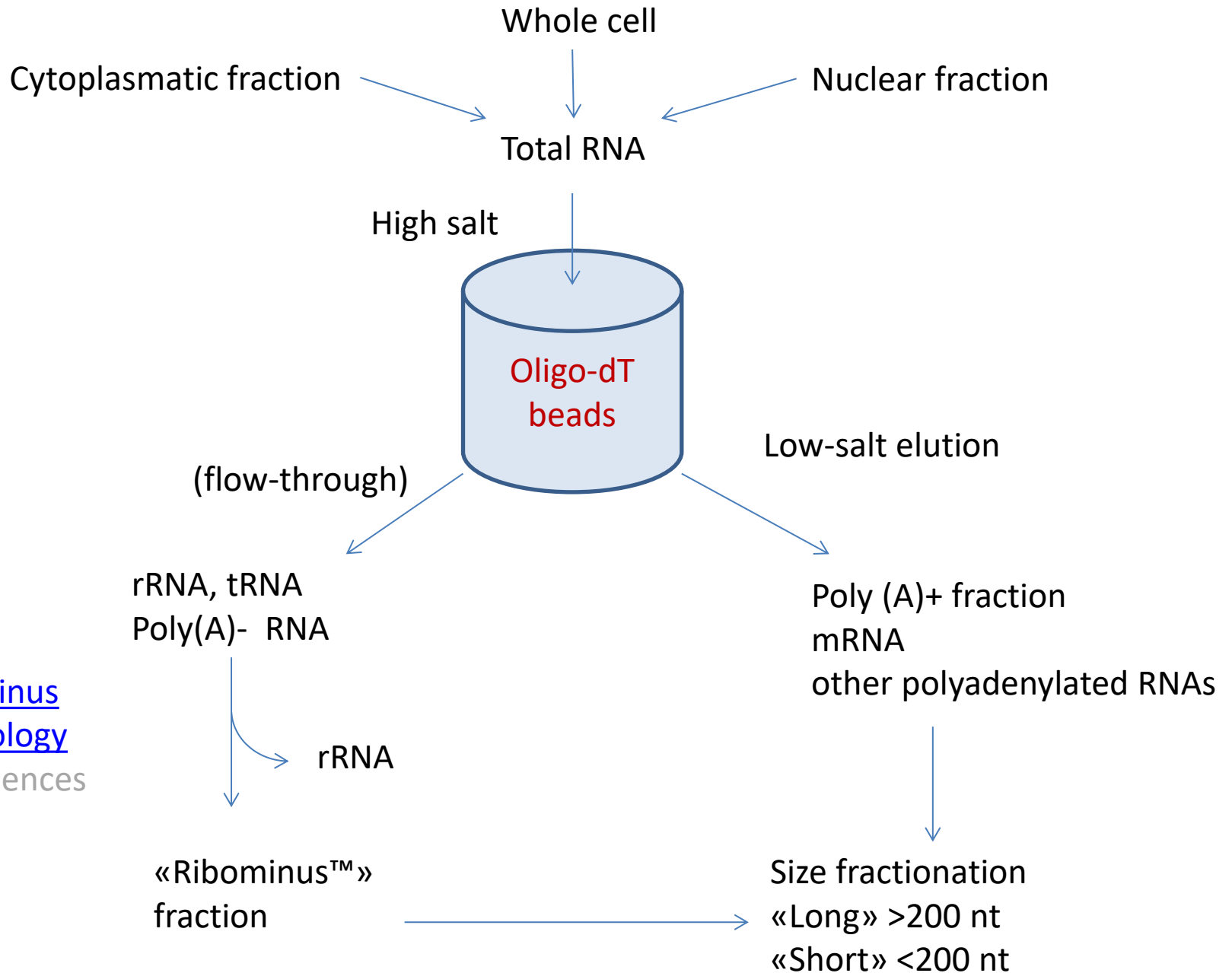
RNA-PET is a paired-end tag (PET) sequencing method for full-length mRNA analysis. RNA-PET captures and sequences the 5'- and 3'-end tags of full-length cDNA fragments of all expressed genes in a biological sample



From Supplementary 2

Supplementary Figure S1

Sample Flowchart. The ENCODE transcriptome data are obtained from several cell lines which have been cultured in replicates. They were either left intact (whole cell) and/or fractionated into cytoplasm and nucleus prior to RNA isolation. Total RNA was then isolated and partitioned into RNA \geq 200bp (long) and $<$ 200bp (short). The long RNA was further partitioned over an oligo-dT column into polyA+ and polyA- fractions. The K562 cell line also underwent additional fractionation into nucleoli, nucleoplasm and chromatin, but no further partition into polyA+ and polyA- was done. RNA-seq was conducted on polyA+, polyA- and total (K562) RNA samples. CAGE was conducted primarily on polyA+ and total RNA but also on some polyA- samples. RNA-PET was conducted on PolyA+ samples only (not shown here are RNA-seq experiments performed at CalTech on polyA+ whole cell RNA extracts).

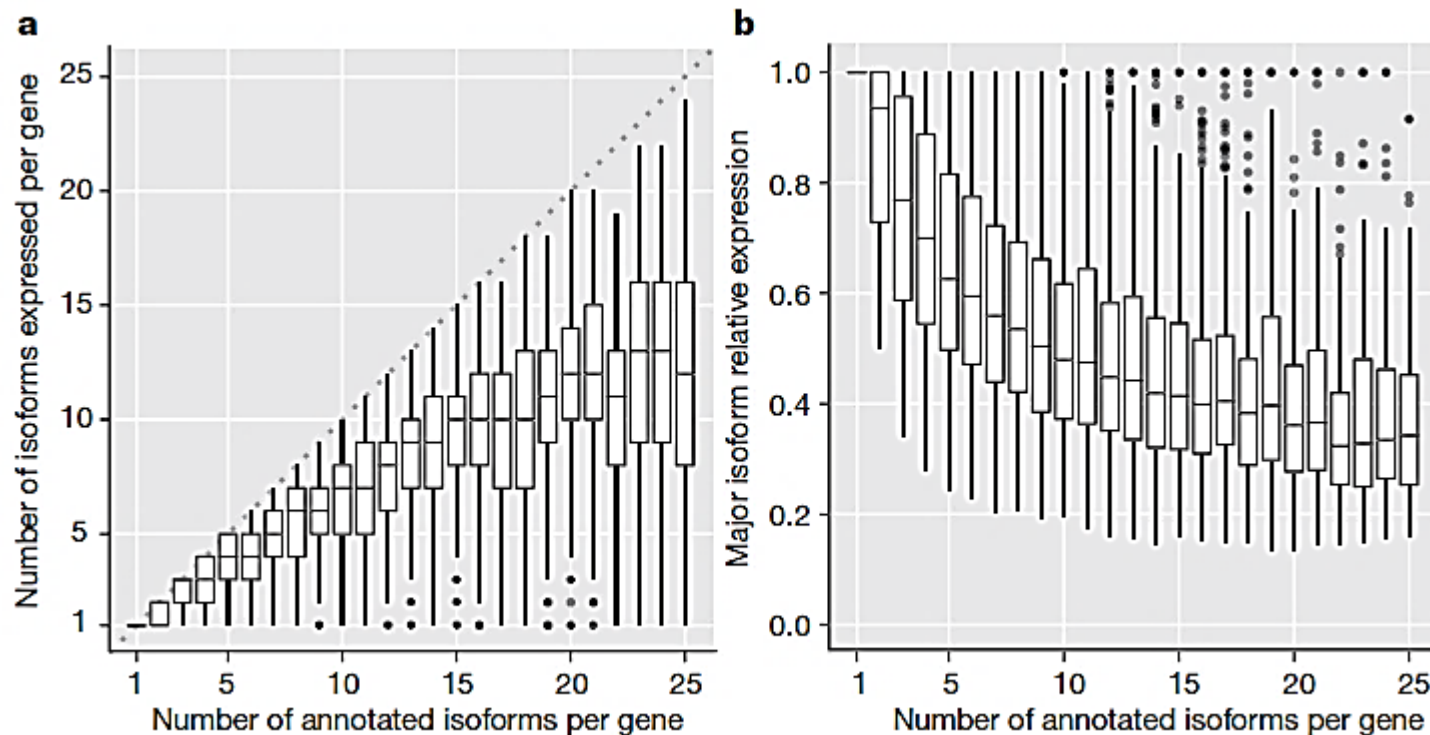


Such a deep, complete and detailed study in many cell lines allowed the discovery of:

- Many novel transcripts from known protein-coding genes (*)
- some undescribed potentially protein coding genes
- a large number of long noncoding RNAs from unsuspected genomic regions (intergenic)
- a large number of intragenic sense and antisense RNAs
- a large number of short and middle-short noncoding RNAs

(*) several transcript isoforms derived from alternative exon splicing, alternative TSS usage, alternative poly(A) site usage

Djebali 2012, Figure 4 – Transcript Isoforms



- Number of expressed isoforms per gene per cell line. A plateau is evident between 10 and 12
- Relative expression of the most abundant isoform per gene per cell line.

See an example at: <https://www.ncbi.nlm.nih.gov/gene/7157>

Alternative transcription initiation and termination.

a total of 128,021 TSSs were detected across all cell lines (97,778 previously annotated ; 30,243 were novel intergenic/antisense TSSs).

CAGE tags.... identified a total of 82,783 nonredundant TSSs

48% of the CAGE-identified TSSs located within 500 base pairs (bp) of an annotated RNA-seq-detected GENCODE TSS,
additional 3% within 500 bp of a novel TSS

A large number of novel transcripts were classified as lncRNA

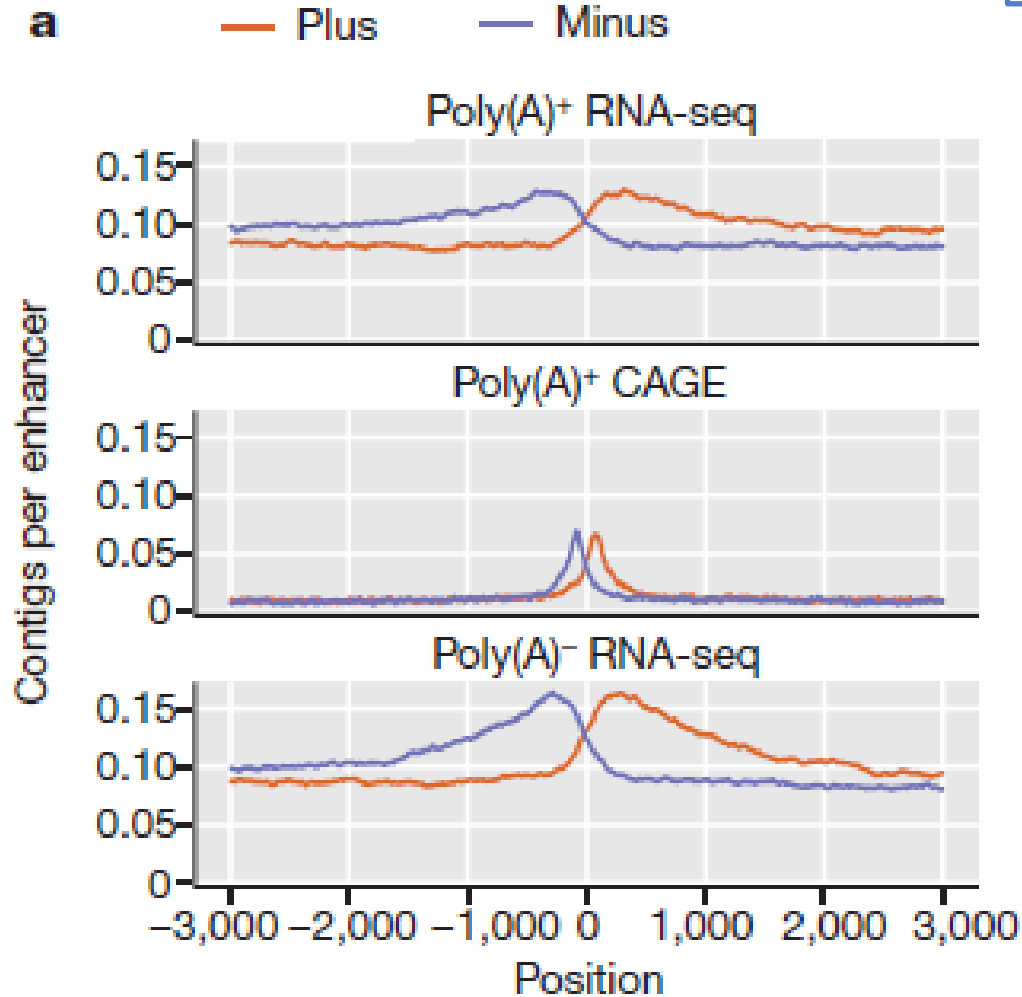
long noncoding RNA

How is the classification «*noncoding*» attributed ?

- ORF search in all the possible frames
- Short ORFs evaluated on «codon usage»
- Proteomic database interrogated
- Association with ribosomes (*poly-ribosome purification and RNA-seq*)

eRNA

Enhancer attribution by means of PTMs+TF CHIP-Seq data



Expression level - quantity

Transcripts range in a 6-order magnitude (poly A+)(10^{-2} to 10^4 rpkm) or 5 orders of magnitude (poly A-) (10^{-2} to 10^3 rpkm)

Assuming that 1–4 r.p.k.m. approximates to 1 copy per cell (*Montazavi et al., 2008*):

- one quarter of protein-coding RNAs *and*
- 80% of long noncoding RNAs (lncRNA)

Are expressed at 1 or <1 molecules per cell

i.e. the majority of lncRNAs are expressed at a very low level

Novel lncRNAs discovered here contains also a class showing *rpkm* from 10^{-4} to 10^{-1} : << extremely low expression >>

Question: what does it mean «less than one molecule per cell ?

Expression level by class

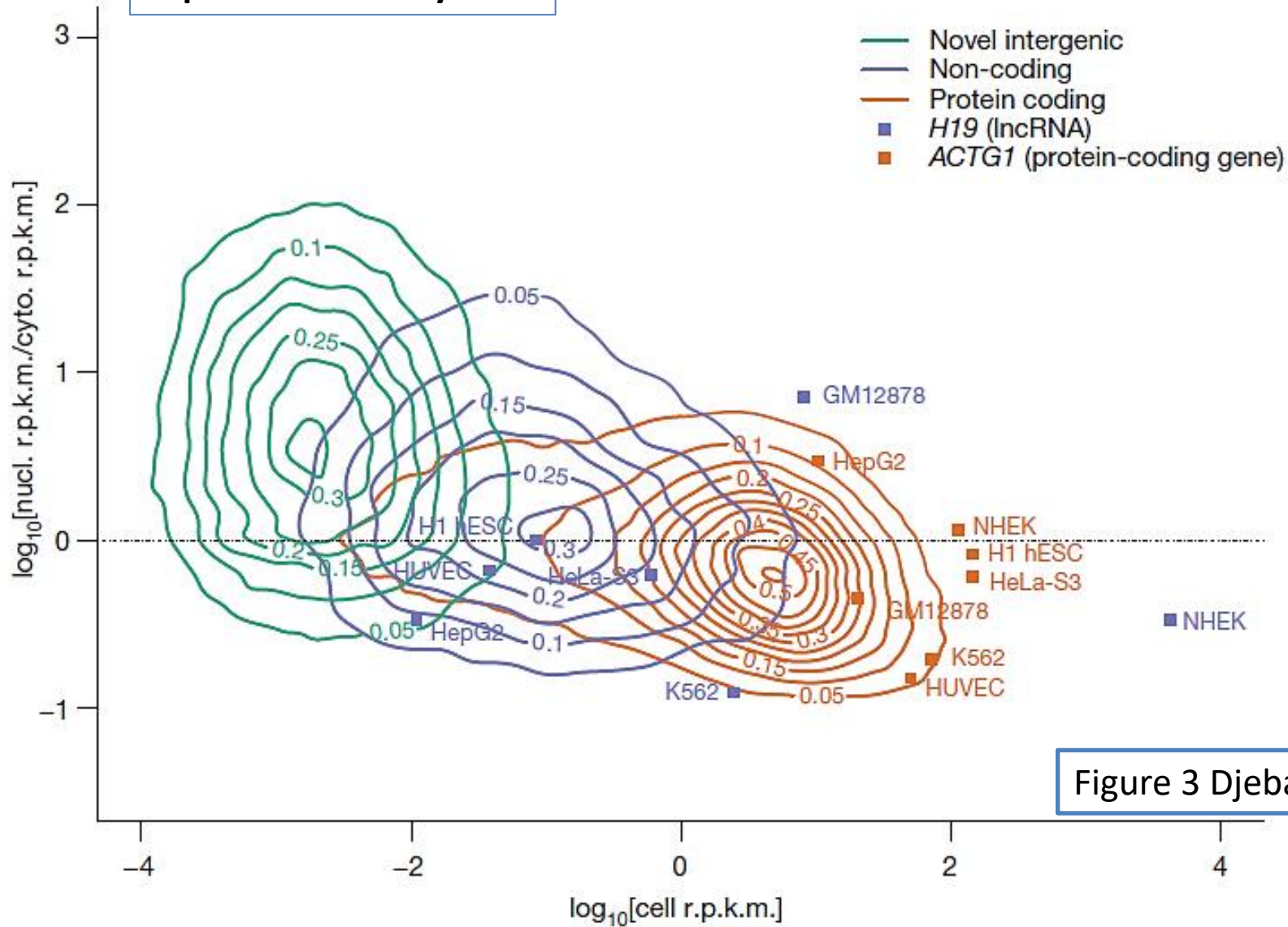


Figure 3 Djebali et al. 2012

Protein coding transcripts are the only class that is enriched in the cytoplasm

Review

The Dimensions, Dynamics, and Relevance of the Mammalian Noncoding Transcriptome

Ira W. Deveson,^{1,2} Simon A. Hardwick,^{1,3} Tim R. Mercer,^{1,3}
and John S. Mattick^{1,2,3,*}

464 Trends in Genetics, July 2017, Vol. 33, No. 7 <http://dx.doi.org/10.1016/j.tig.2017.04.004>
© 2017 Elsevier Ltd. All rights reserved.

The proliferation and evolution of RNA-Seq, including the advent of methods for targeted, single-molecule, and single-cell sequencing, continues to enlarge our understanding of **transcriptional diversity**.

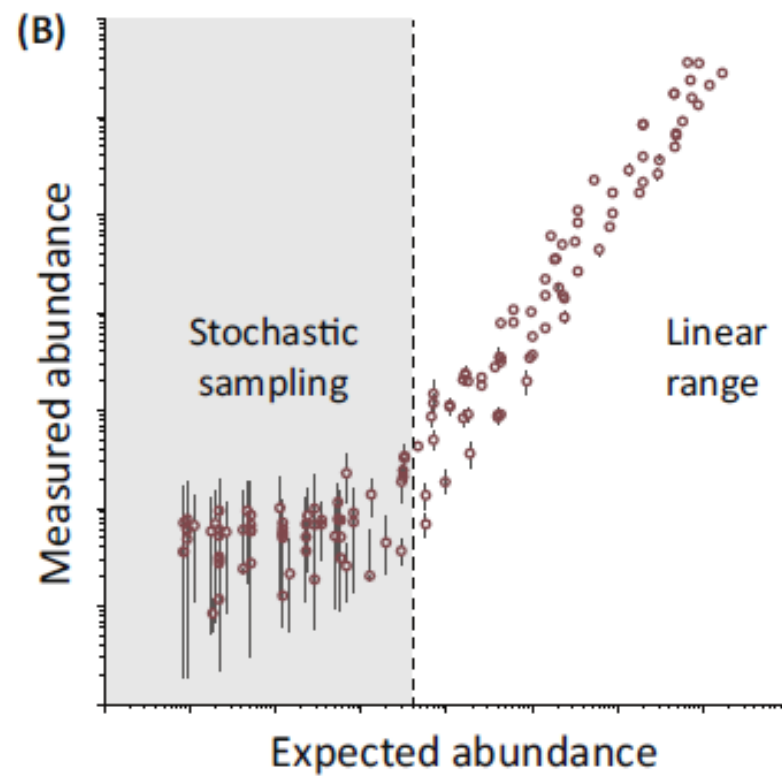
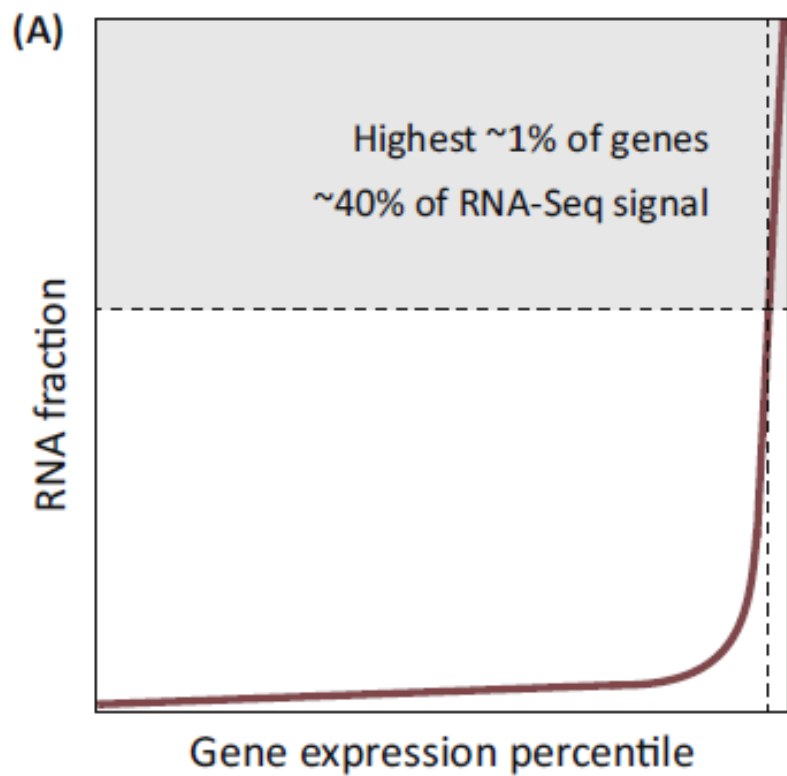
LncRNA are on average expressed at a lower range than protein-coding genes

To explore and understand less expressed lncRNAs, the Targeted RNA-seq method was developed. In practice, rare transcripts are selected using appropriate primers so that the sequencing library is enriched.

LncRNA databases vary greatly in number, and this is due to the criteria assumed to accept a lncRNA.

GENECODE is the most conservative,

MiTranscriptome lists 58,648 lncRNAs compared to 21,313 protein-coding



Expression of lncRNAs is highly tissue-specific

ENCODE: 50% of the features were seen only in one cell line.

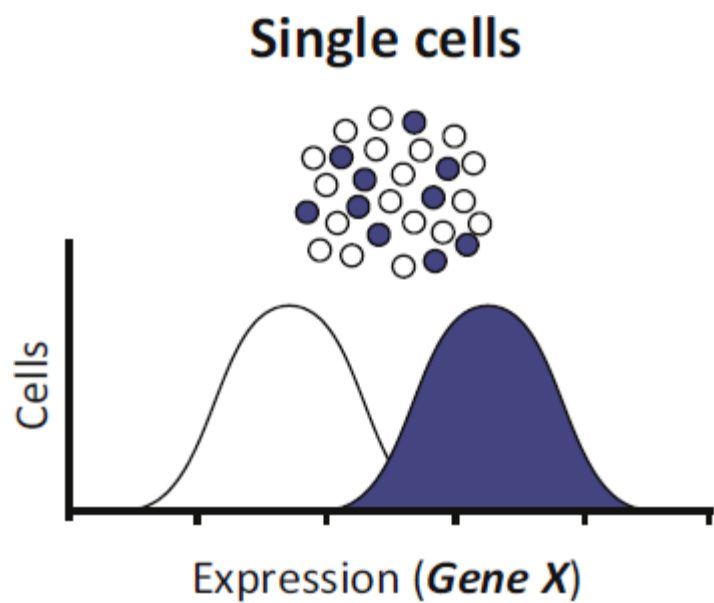
By FISH analysis: expression highly limited to cell types (in brain)

Single-cell RNA-Seq : lncRNAs expressed only in some cells.

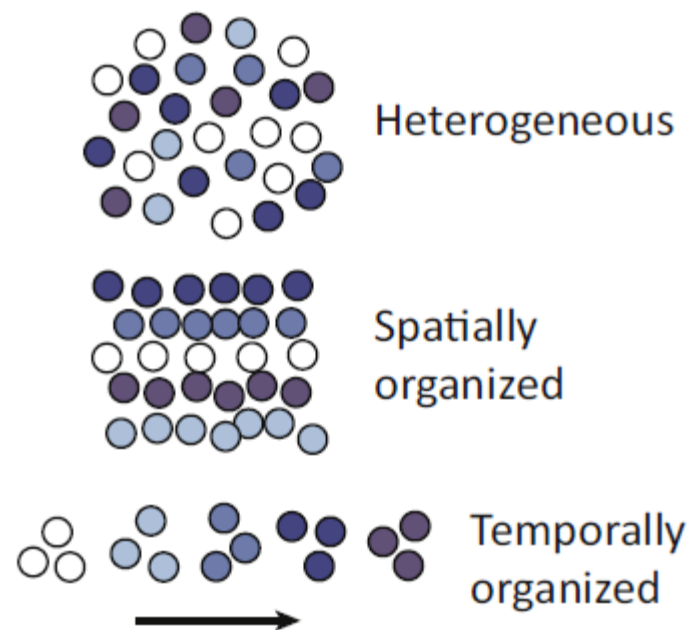
Thus, average expression is low, but single-cell expression can be high

Cell-subtype determinants?

Diversity not seen in cell cultures, possibly since much more homogeneous.

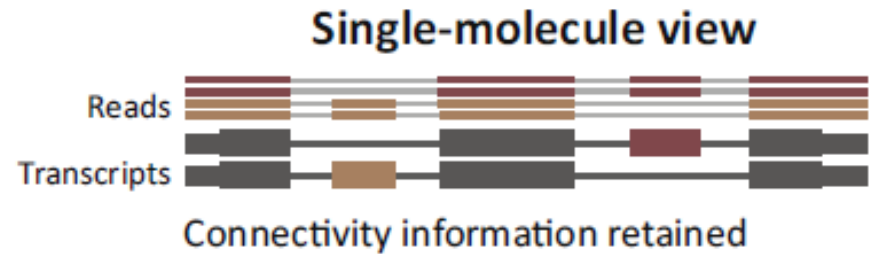
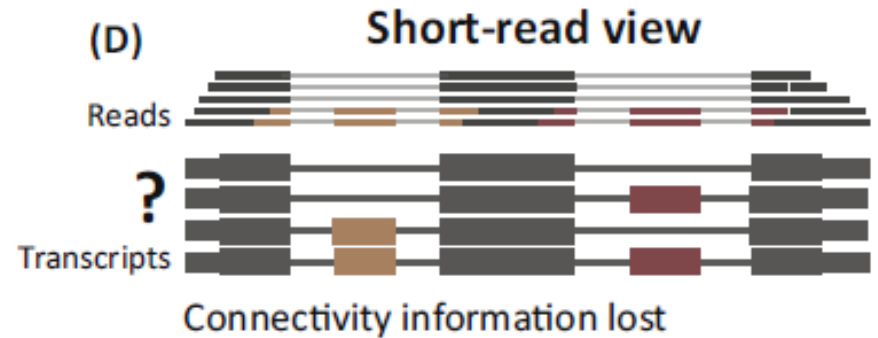
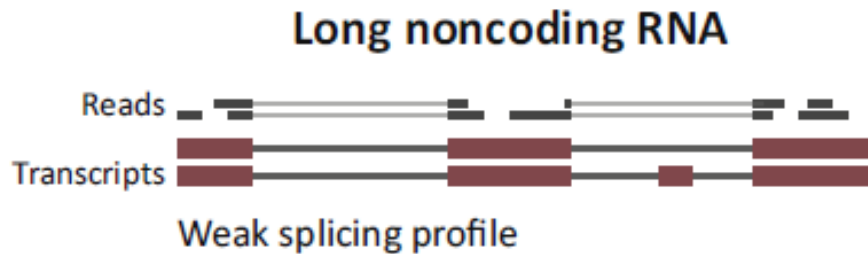
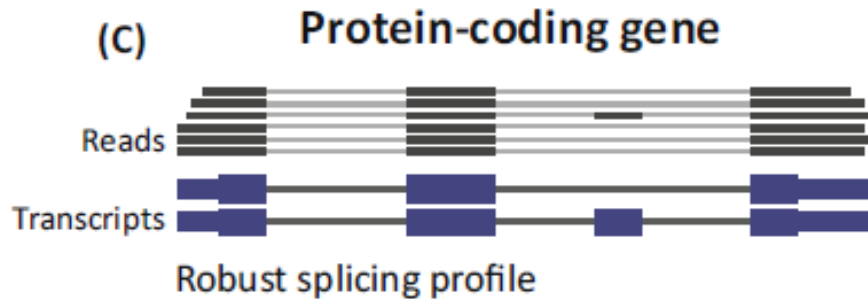


(B)



Most lncRNA have the same structure as protein-coding: Exons & Introns

Short-reads sequencing make it difficult to discriminate among transcripts



Definite improvements are expected for single-molecule, long-read sequencing technologies

Oxford Nanopore: <https://vimeo.com/211385238>

Pacific Bio:

<https://www.pacb.com/smrt-science/smrt-sequencing/>

Functional characterization: only few lncRNAs

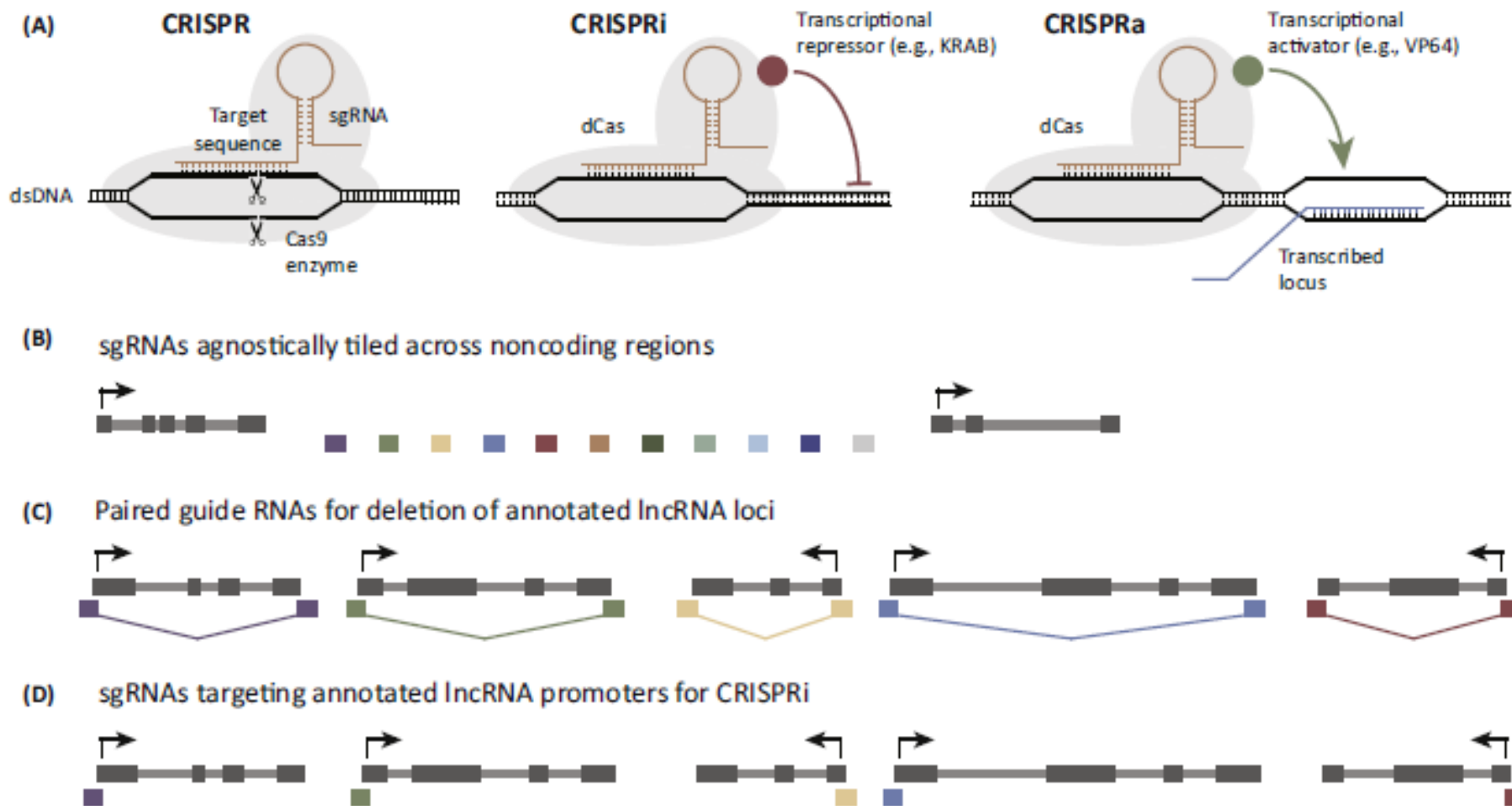
CRISPR screening are made today

Many lncRNA KO or KD → lethal phenotype

Many lncRNAs participate in epigenomic regulation

(examples from monoallelic expression lesson, interacting with PRC2)

- HOTAIR binds both PRC2 and LSD1 (KDM) (repressor scaffolds)
- NEAT1 in paraspeckles
- Xist in X-chr inactivation
- RNA-a at regulatory sequences (Enhancers?) – bind to Mediator for looping
- Other lncRNAs in imprinted regions (local repressor)



from Deveson et al. (TextBook)