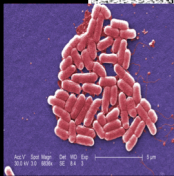
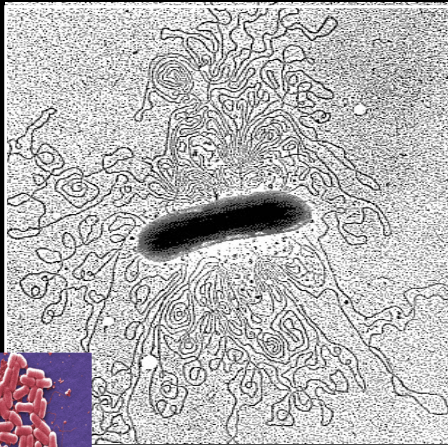


PSEUDOGENE DERIVED lncRNAs

Reason 1: The non-coding genome (r)evolution

E.coli



C. elegans

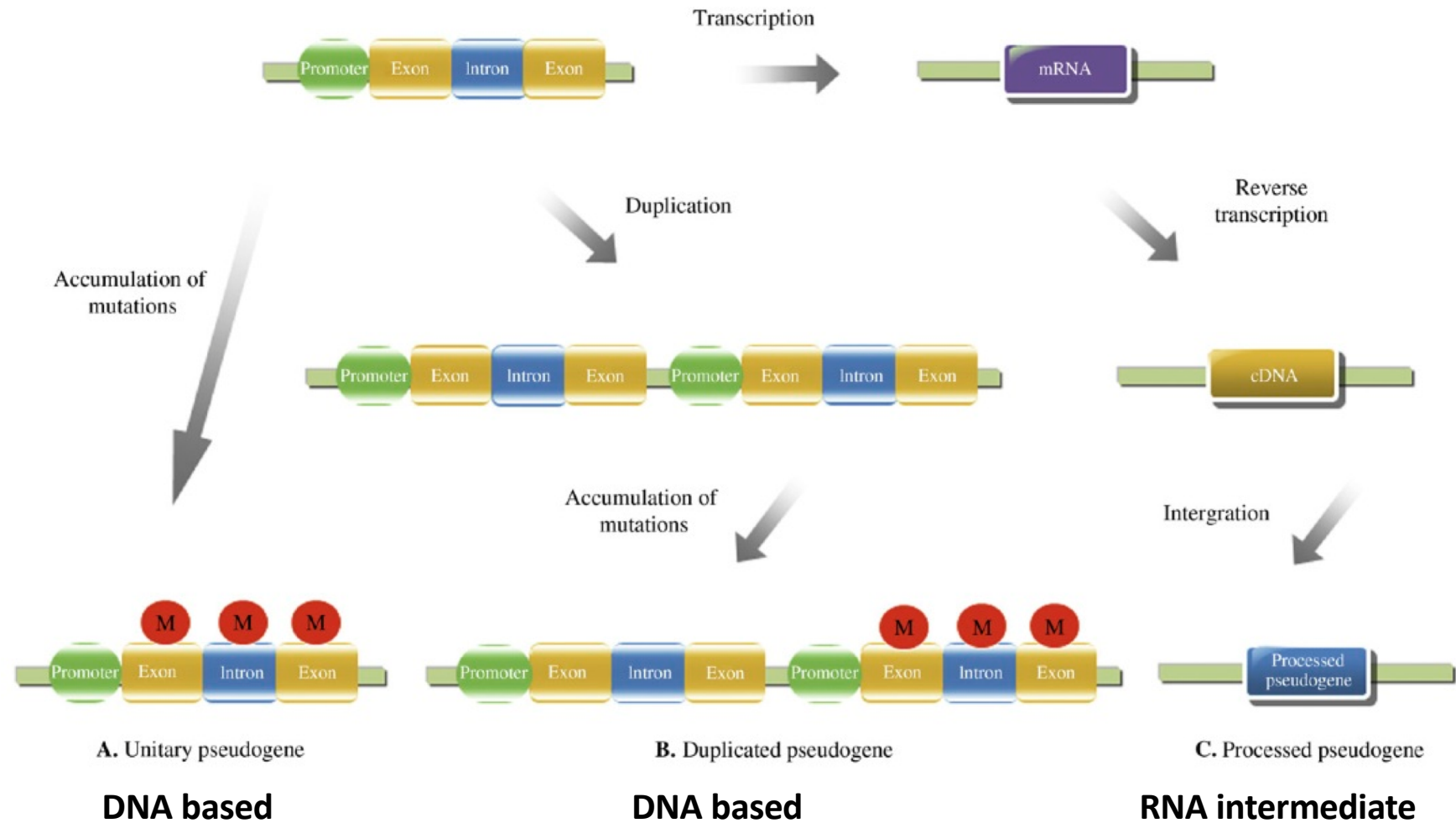


H. sapiens

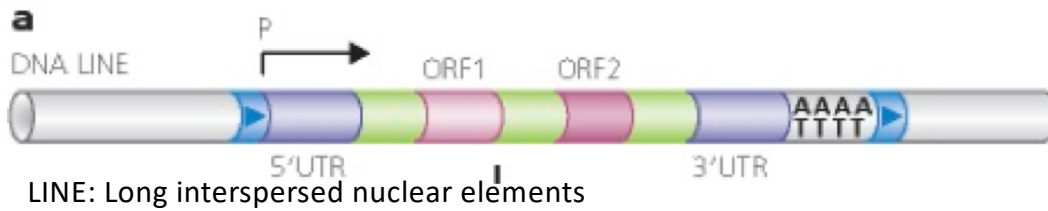


	Genome	5×10^6 bp	1×10^8 bp	3×10^9 bp
Chromosomes		1	6	23
Coding genes		6692	20541	21995
ncDNA		5%	60%	98%
non-coding RNA genes		15	23136	ca. 40000
miRNAs		0	224	4274
pseudogenes		21	1522	10616

Protein coding genes give rise to pseudogenes



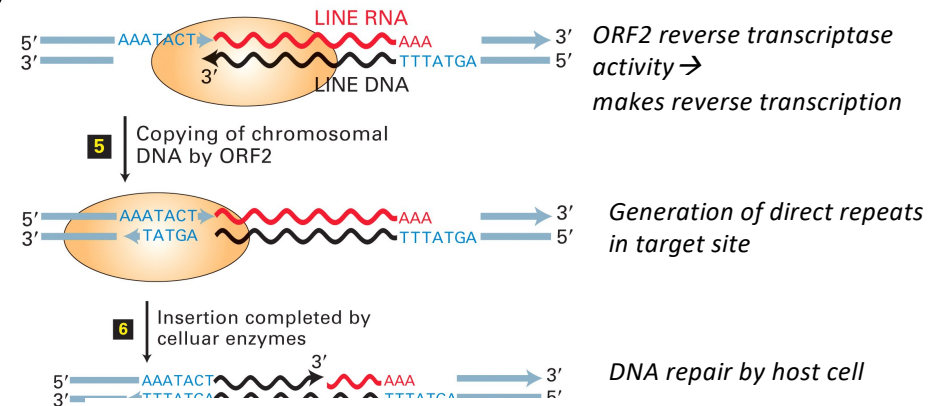
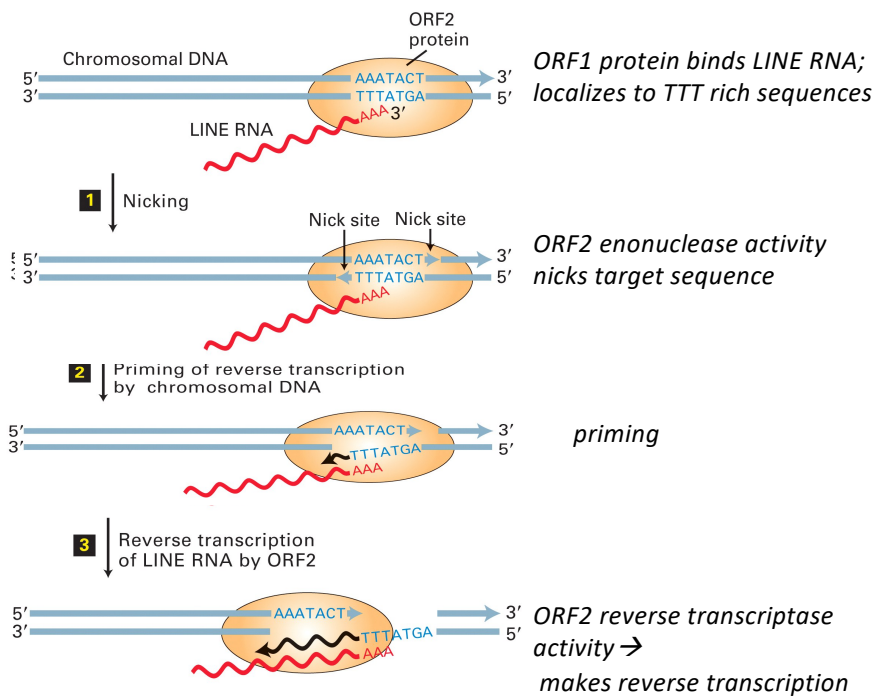
Transposition of Retrotransposons



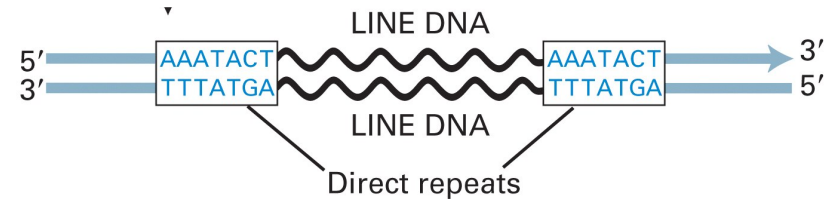
LINE elements (L1,L2,L3)
 (21% of genome; 800.000 copies)

ORF1: RNA binding protein

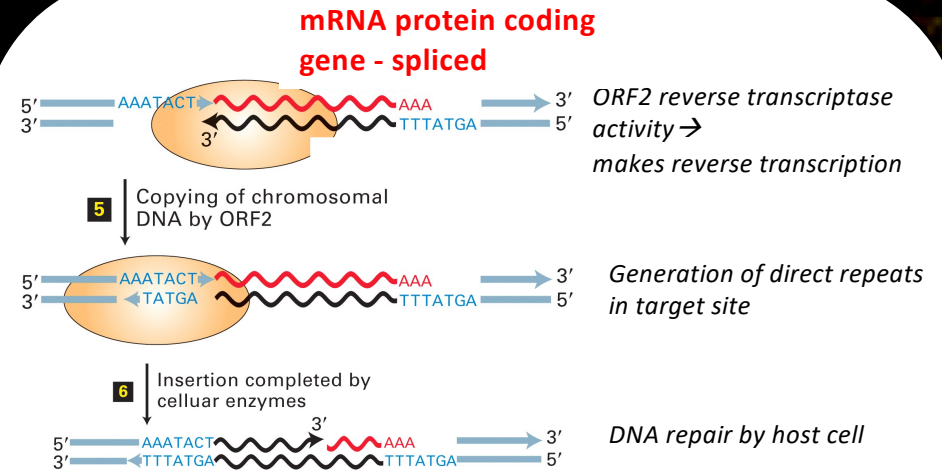
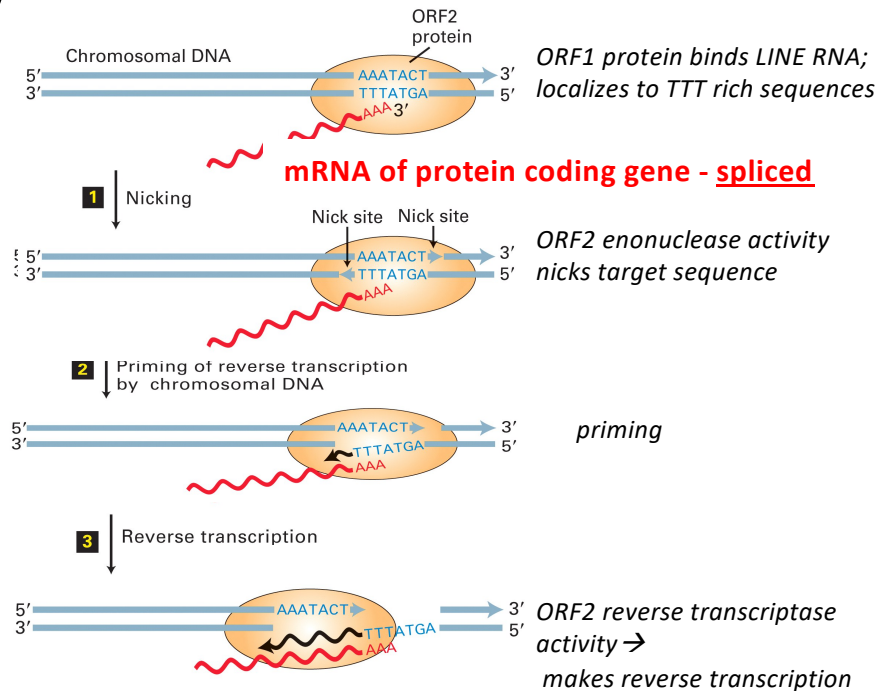
ORF2: Endonuclease, Reverse transcriptase



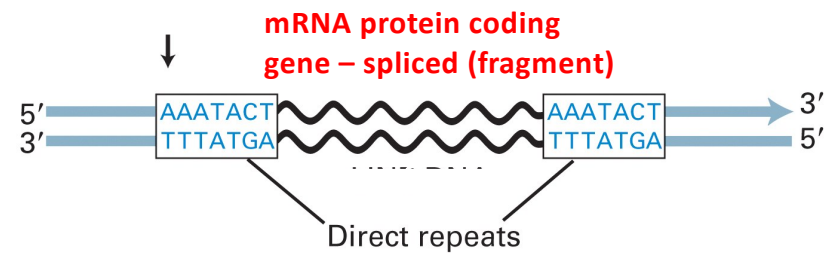
FINAL PRODUCT



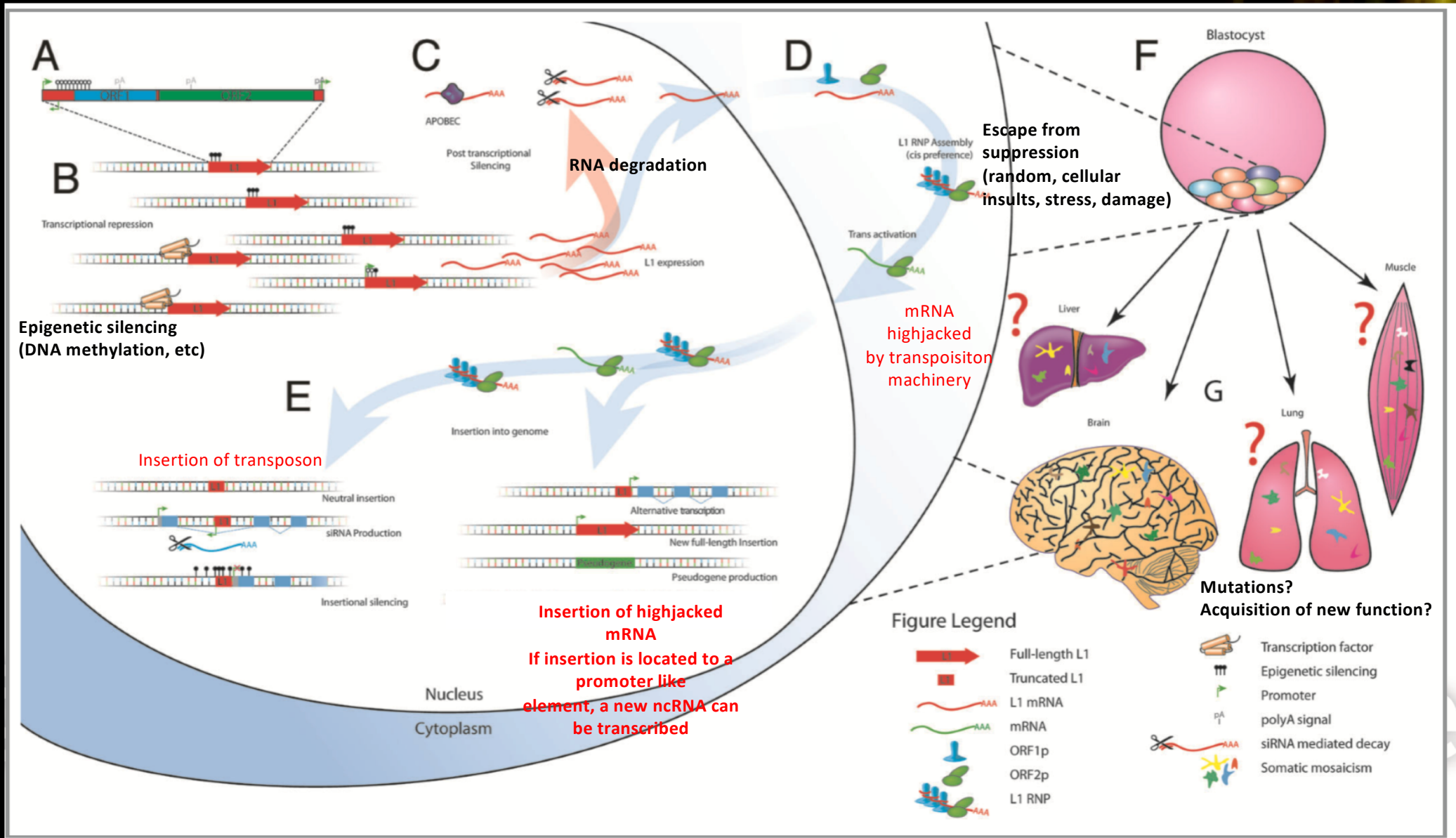
Retro-transposition machinery hijacks endogenous mRNAs



FINAL PRODUCT: PROCESSED PSEUDOGENE



Retrotransposons can change genetic context



PSEUDOGENE BIOTYPES

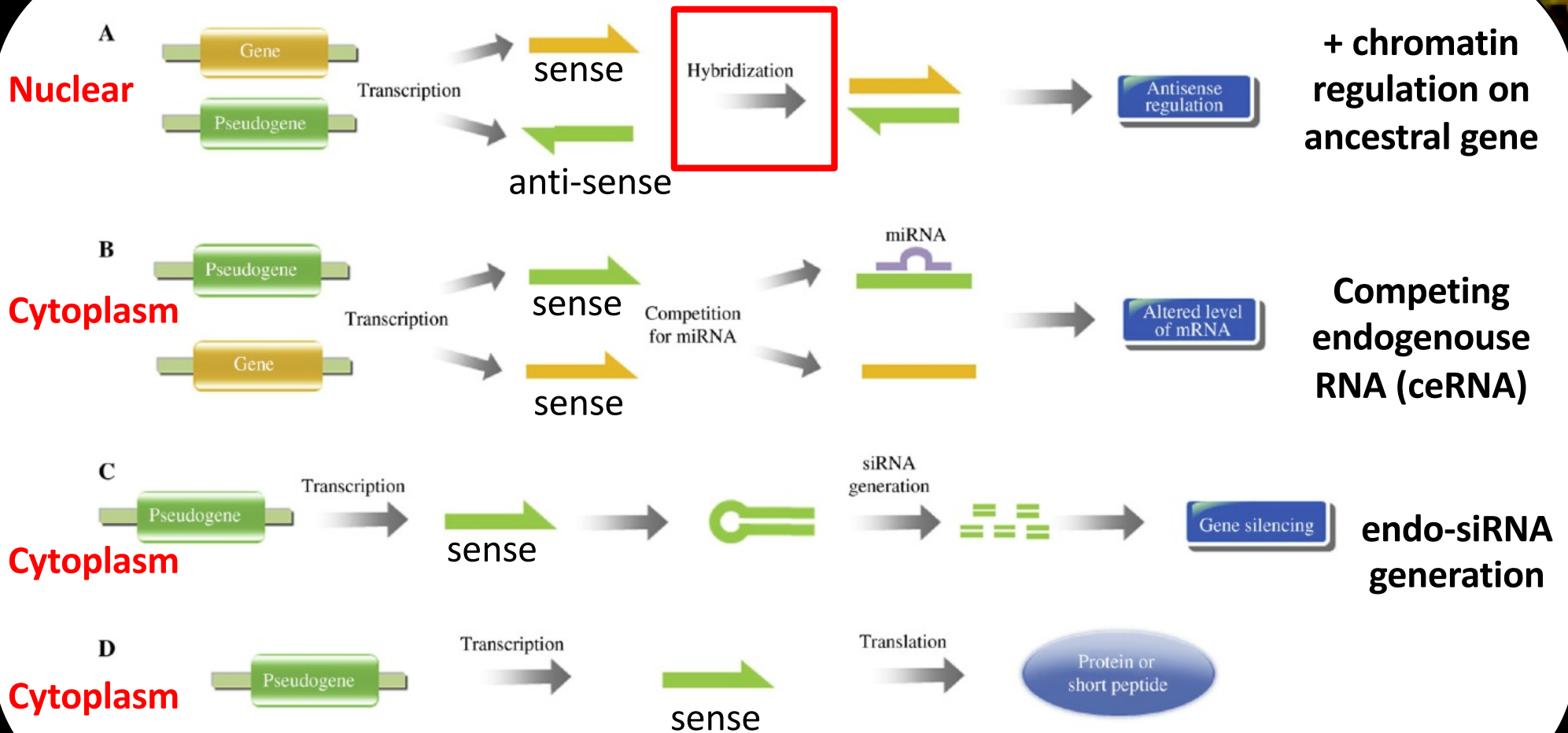
Table 2 Pseudogene biotypes

Biotype	Definition
Processed pseudogene	Pseudogene created via retrotransposition of the mRNA of a functional protein-coding parent gene followed by accumulation of disabling mutations
Duplicated pseudogene	Pseudogene created via genomic duplication of a functional protein-coding parent gene followed by accumulation of disabling mutations
Unitary pseudogene	Pseudogene for which the ortholog in a reference species (mouse) is coding but the human locus has accumulated fixed disabling mutations
Polymorphic pseudogene	Locus known to be coding in some individuals but with disabling mutations in the reference genome
IG pseudogene	Immunoglobulin gene segment with disabling mutations
TR pseudogene	T-cell receptor gene segment with disabling mutations

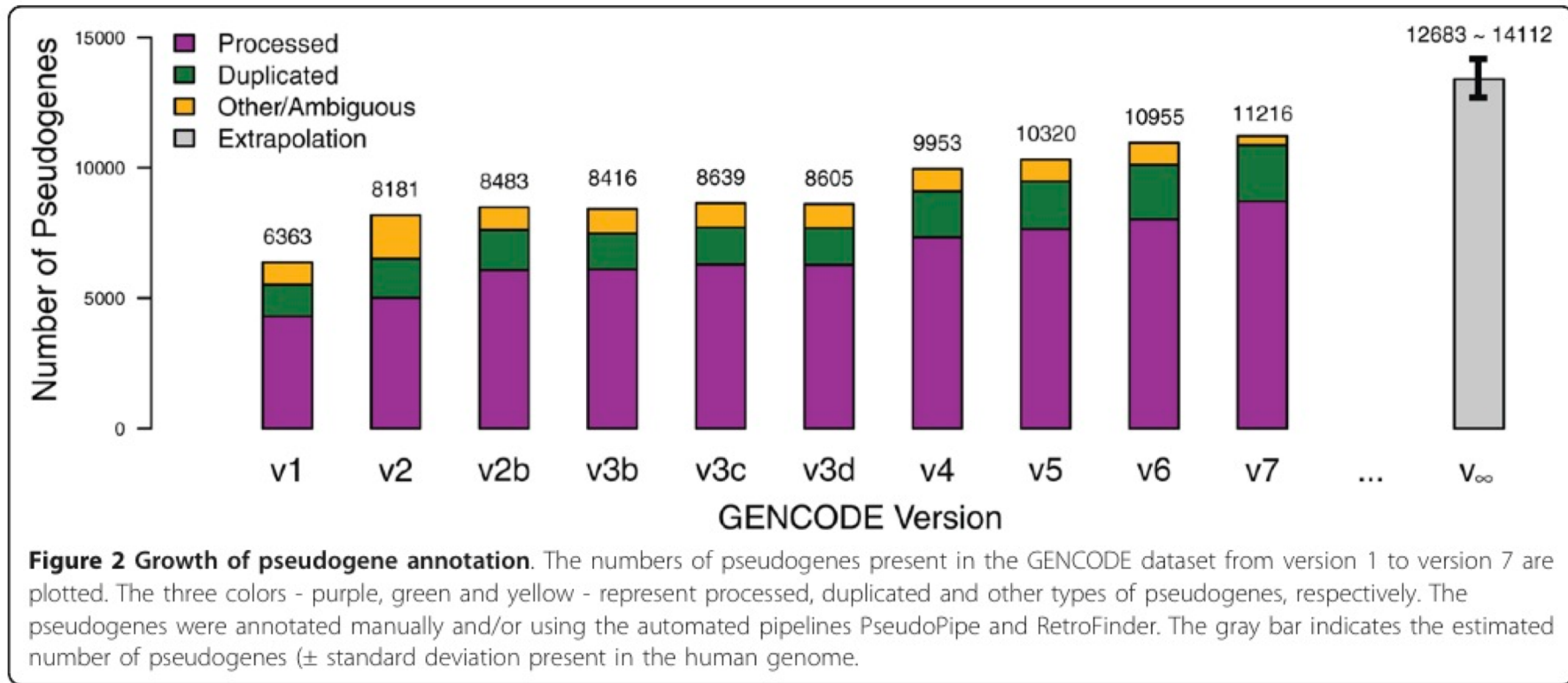
Duplicated/Unitary pseudogenes: can bring regulatory sequences, often spliced

Processed pseudogenes: hitch hike on regulatory elements dispersed throughout the genome; expression depends on the vicinity of regulatory elements

Pseudogene derived RNAs can acquire new functions



PSEUDOGENE BIOTYPES



*The majority of pseudogenes are processed pseudogenes:
Burst of retro-transposition events in recent phase of evolution*

Total No of Genes	60498
Protein-coding genes	19797
Long non-coding RNA genes	15931
Small non-coding RNA genes	9882
Pseudogenes	14477
- processed pseudogenes:	10727
- unprocessed pseudogenes:	3271
- unitary pseudogenes:	172
- polymorphic pseudogenes:	59

GENOMICS STRATEGIES TO IDENTIFY AND CLASSIFY PSEUDOGENES

Table 3 Fields for pseudogene features in the psiDR annotation file **Pseudogene decoration resource**

Field	Explanation	psiDR value
Transcript ID	Pseudogene ID from GENCODE annotation. Used for cross-referencing	
Parent	Protein ID, Gene ID, chromosome, start, end and strand. Detailed in section ' <i>Parents of pseudogenes</i> '	
Sequence similarity	The percentage of pseudogene sequence preserved from parent	
Transcription	Evidence for pseudogene transcription and validation results. May be tagged as EST, BodyMap, RT-PCR or None, which represent pseudogene expression evidence from corresponding data sources. Multiple tags are separated by commas. Detailed in section ' <i>Transcription of pseudogenes</i> '	1, transcription; 0, otherwise
DNaseI hypersensitivity	A categorical result indicating whether the pseudogene has easily accessible chromatin, predicted by a model integrating DNaseI hypersensitivity values within 4 kb genomic regions upstream and downstream of the 5' end of pseudogenes. Detailed in section ' <i>Chromatin signatures of pseudogenes</i> '	1, has Dnase hypersensitivity in upstream; 0, otherwise
Chromatin state	Whether a pseudogene maintains an active chromatin state, as predicted by a model using Segway segmentation. Detailed in section ' <i>Chromatin signatures of pseudogenes</i> '	1, active chromatin; 0, otherwise
Active Pol2* binding	Whether Pol2 binds to the upstream region of a pseudogene. Detailed in section ' <i>Upstream regulatory elements</i> '	1, active binding site; 0, otherwise
Active promoter region	Whether there are active promoter regions in the upstream of pseudogenes. Detailed in section ' <i>Upstream regulatory elements</i> '	1, active binding site; 0, otherwise
Conservation	Conservation of pseudogenes is derived from the divergence between human, chimp and mouse DNA sequences. Detailed in section ' <i>Evolutionary constraint on pseudogenes</i> '	1, conserved; 0, otherwise

*Pol2, RNA polymerase II.

- **Parent gene/ancestral gene = functional gene with greatest sequence similarity**
- **Ancestral gene can be identified for ca. 90% of pseudogenes**
- **10% of pseudogenes are highly degraded and is derived from a parent gene with highly similar paralogs**
- **Or parent gene contains a commonly found functional domain**
- **NOTE: most parental genes have only 1 pseudogene**
- **NOTE: some parental genes – mainly housekeeping genes - have MANY pseudogenes:**
 - **Robosomal protein L21: 143 pseudogenes**
 - **Gapdh: 68 pseudogenes**

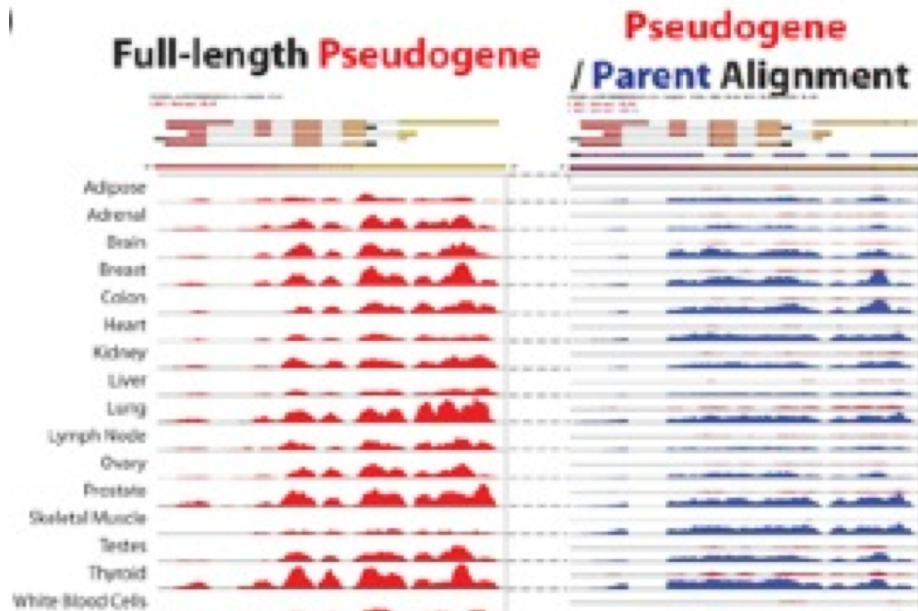
Features of transcribed pseudogenes

Problem: precise analysis of RNA-seq/array data: high sequence similarity pseudogene – parental gene

2012: ca 9000 pseudogenes: 873 are transcribed according to STRINGENT psiDR parameters (real number is higher)

tissue specific expression

transcription of pseudogene



transcription of pseudogene and parental gene

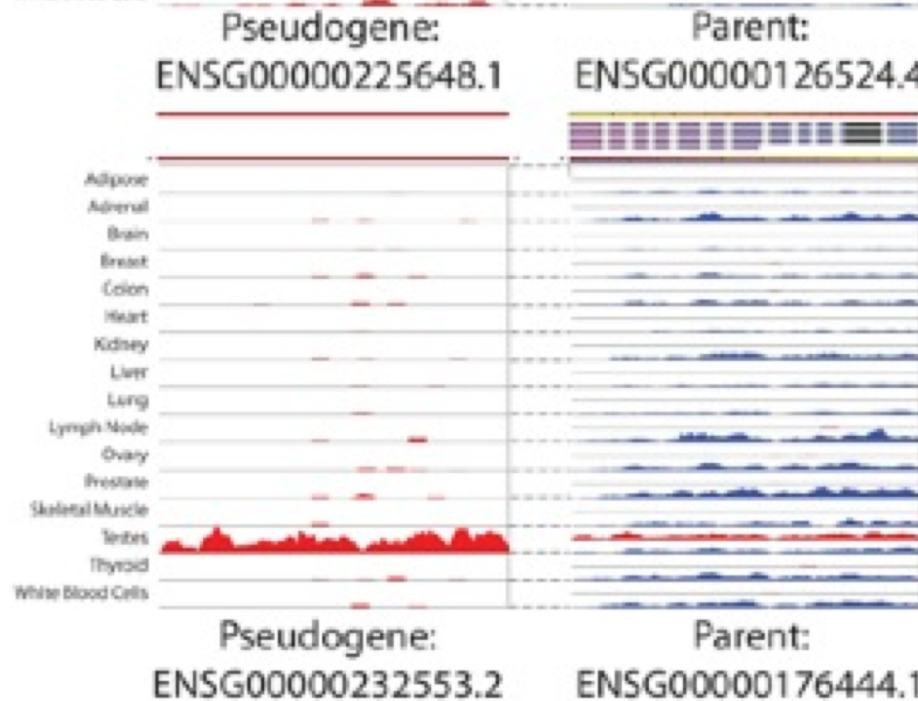
differential expression parental/pseudogene

Pseudogene expression levels are LOWER than coding gene expression

Pseudogenes are expressed in a different manner compared to parental mRNAs (different tissues)

tissue specific expression

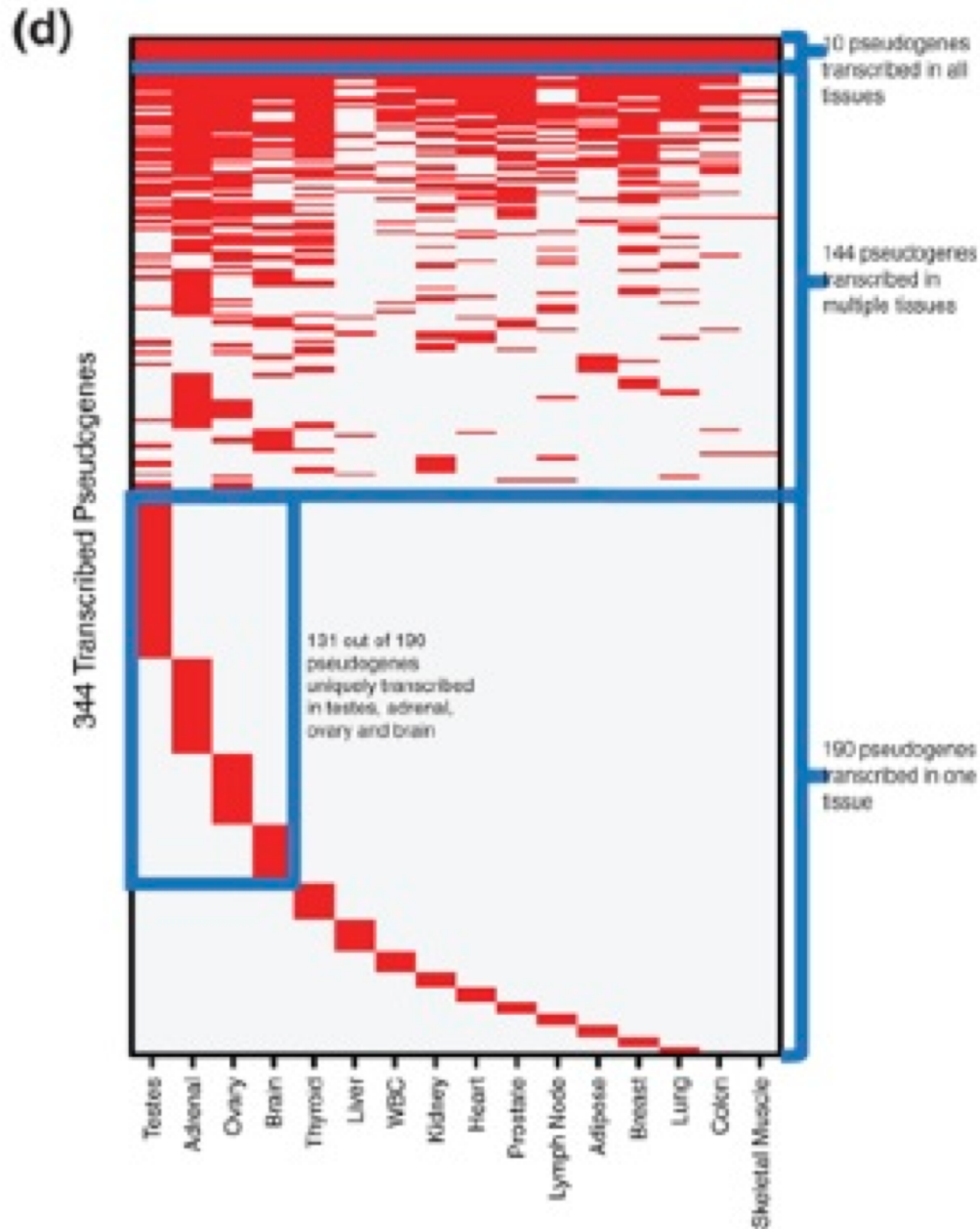
transcription of pseudogene



transcription of pseudogene and parental gene

differential expression parental/pseudogene

The majority of pseudogenes show tissue specific expression



Categories:

- Expressed in all tissues
(10 out of 344 tested pseudogenes)
- 144/344 pseudogenes expressed in more than 1 tissue
- 190/344 pseudogenes exclusively expressed in 1 tissue

duplicate/processed pseudogenes have specific regulatory elements!!

Evolutionary constraint on pseudogenes in different species

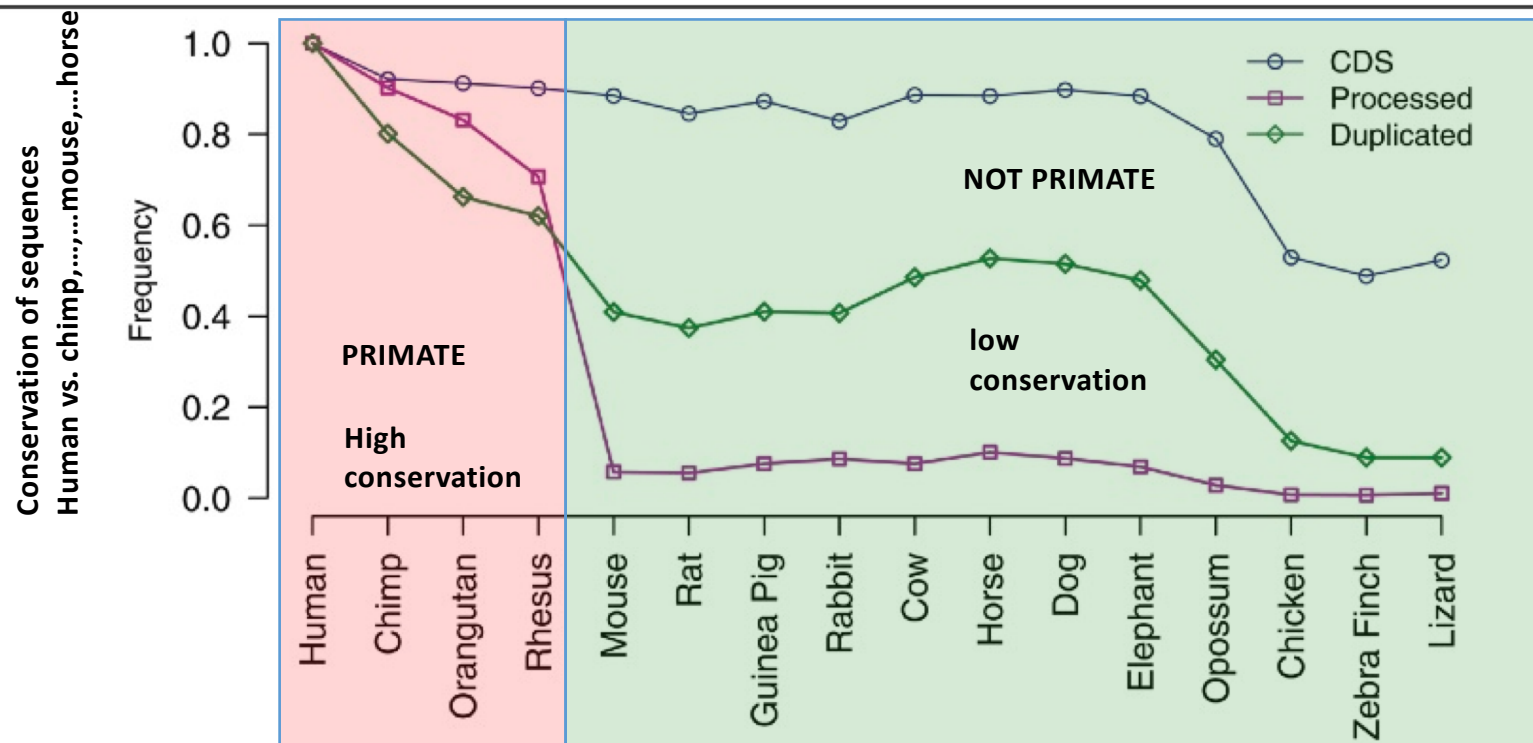


Figure 6 Preservation of human coding sequences, processed pseudogenes and duplicated pseudogenes. Sequences orthologous to human genomic regions from different species were studied. The sequence preservation rate was calculated as the percentage of sequences aligned to human sequence from each species. The calculation was based on a MultiZ multiple genome sequence alignment.

dogenes. While the preservation of duplicated pseudogenes decreases gradually with the increase of evolutionary distance of the species from human, the preservation of processed pseudogenes exhibits an abrupt decrease from macaque to mouse and remains low within the species more divergent than mouse.

These results are in agreement with previous findings showing that most processed pseudogenes in humans and mice are lineage-specific, arising from distinct retrotransposition bursts happening in the two organisms after they diverged [13,41].

Selective constraint in »inside« specific pseudogene lncRNAs

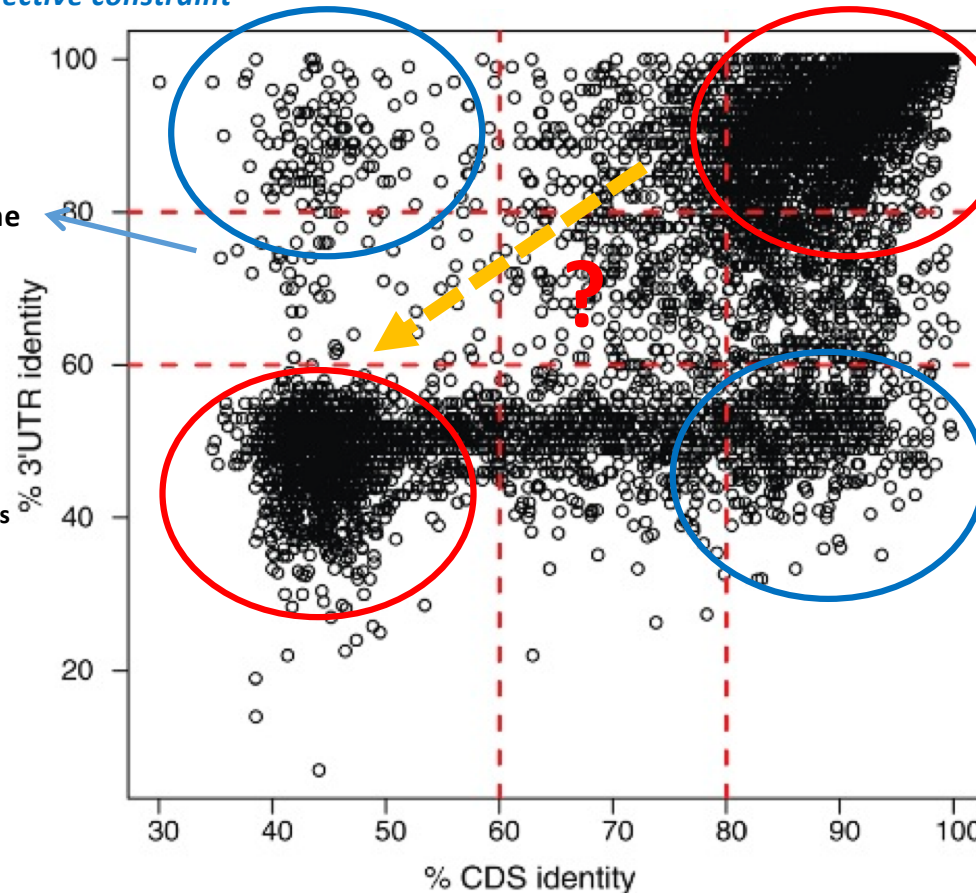
Sequence identity between parental and pseudogenes with focus on coding sequence (CDS) and 3'UTRs of ancestral mRNAs

high 3'UTR conservation
low CDS conservation
= high selective constraint

Fraction of pseudogene lncRNAs show selective conservation of 3'UTR

single pseudogene

Fraction of pseudogene lncRNAs show overall loss of conservation



High CDS conservation
High 3'UTR conservation
= high selective constraint

Most pseudogenes have similar sequence identity between CDS and UTR (high: >80%; low: <60%)

low 3'UTR conservation
high CDS conservation
= high selective constraint

Fraction of pseudogene lncRNAs show selective conservation of CDS

Mutations were rejected by natural selection non-randomly. Certain regions in the sequence may be under higher evolutionary constraint than the others.

Inconsistency implies that mutations were rejected by natural selection non-randomly. Certain regions in the sequence may be under higher evolutionary constraint than the others. 998 pseudogenes show a high (>80%) sequence identity to parent CDS and simultaneously poor (<60%) sequence identity to the 3' UTR, and 36 pseudogenes with high (>80%) sequence identity to the parent 3' UTR and small (<60%) sequence identity to CDS.

Chromatin at transcriptional start site of transcribed pseudogenes is similar to coding genes

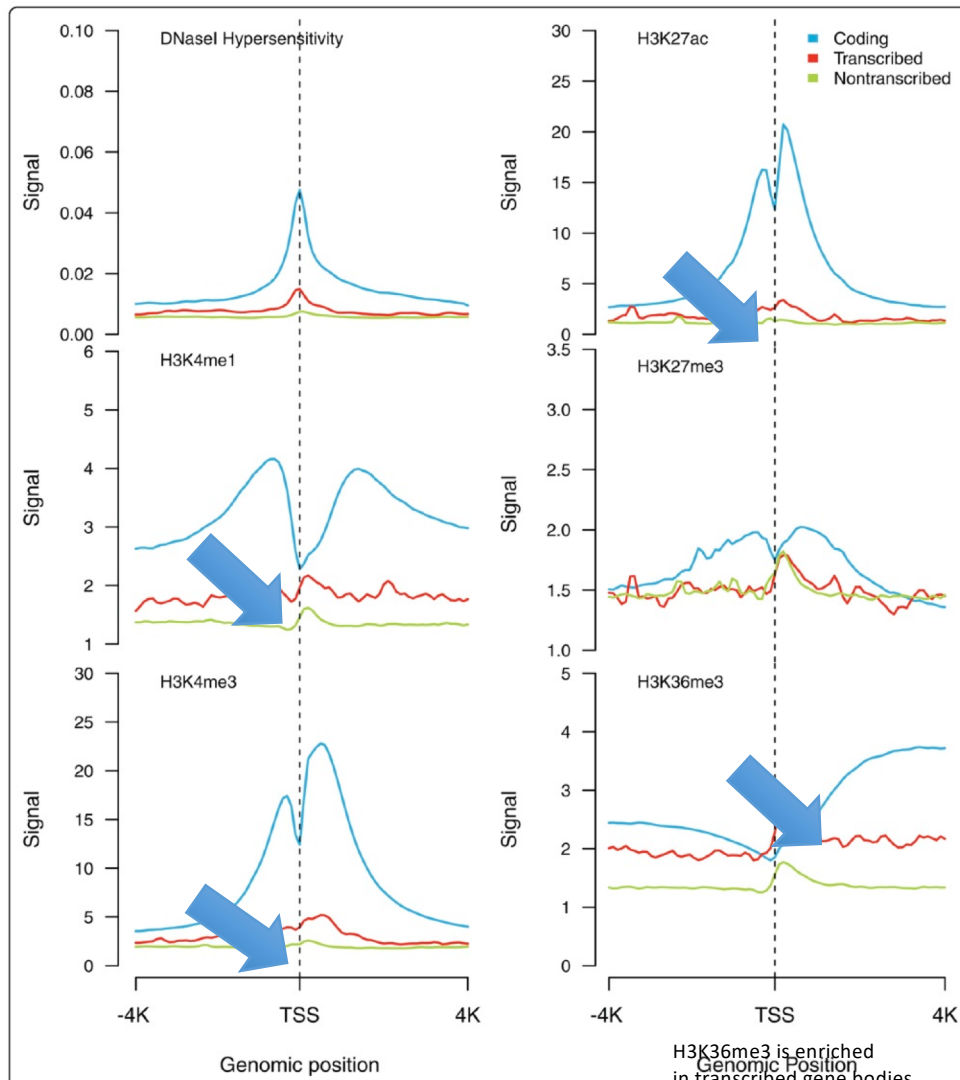
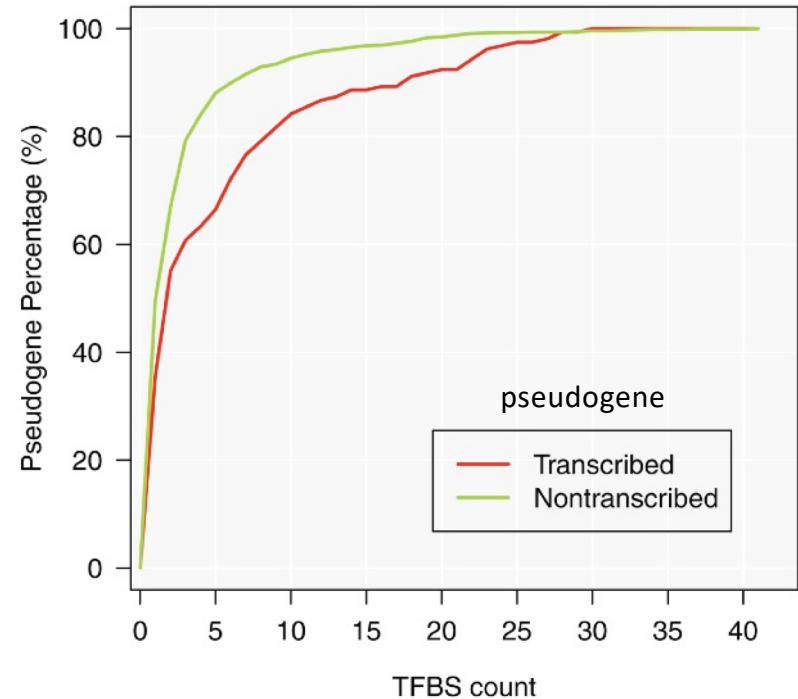


Figure 8 Chromatin signatures: DNaseI hypersensitivity and histone modification. Average chromatin accessibility profiles and various histone modifications surrounding the TSS for coding genes, transcribed pseudogenes, and non-transcribed pseudogenes. The coding gene histone modification profiles around the TSS follow known patterns - for example, enrichment of H3K4me1 around 1 kb upstream of the TSS and the H3K4me3 peaks close to the TSS [63]. Transcribed pseudogenes also show stronger H3K4 signals than non-transcribed pseudogenes. H3K27me3, a marker commonly associated with gene repression [64], showed depletion around the TSS for the coding gene and a distinctive peak in the same region for the pseudogenes. H3K36me3 also shows a similar pattern as H3K27me3 at TSSs, which may relate to nucleosome depletion.



Frequency of transcription factor binding sites enriched in transcribed Pseudogenes vs non-transcribed pseudogenes

Transcribed pseudogenes resemble coding genes; however: Peaks are not as clear defined = average chromatin marks are less concentrated:

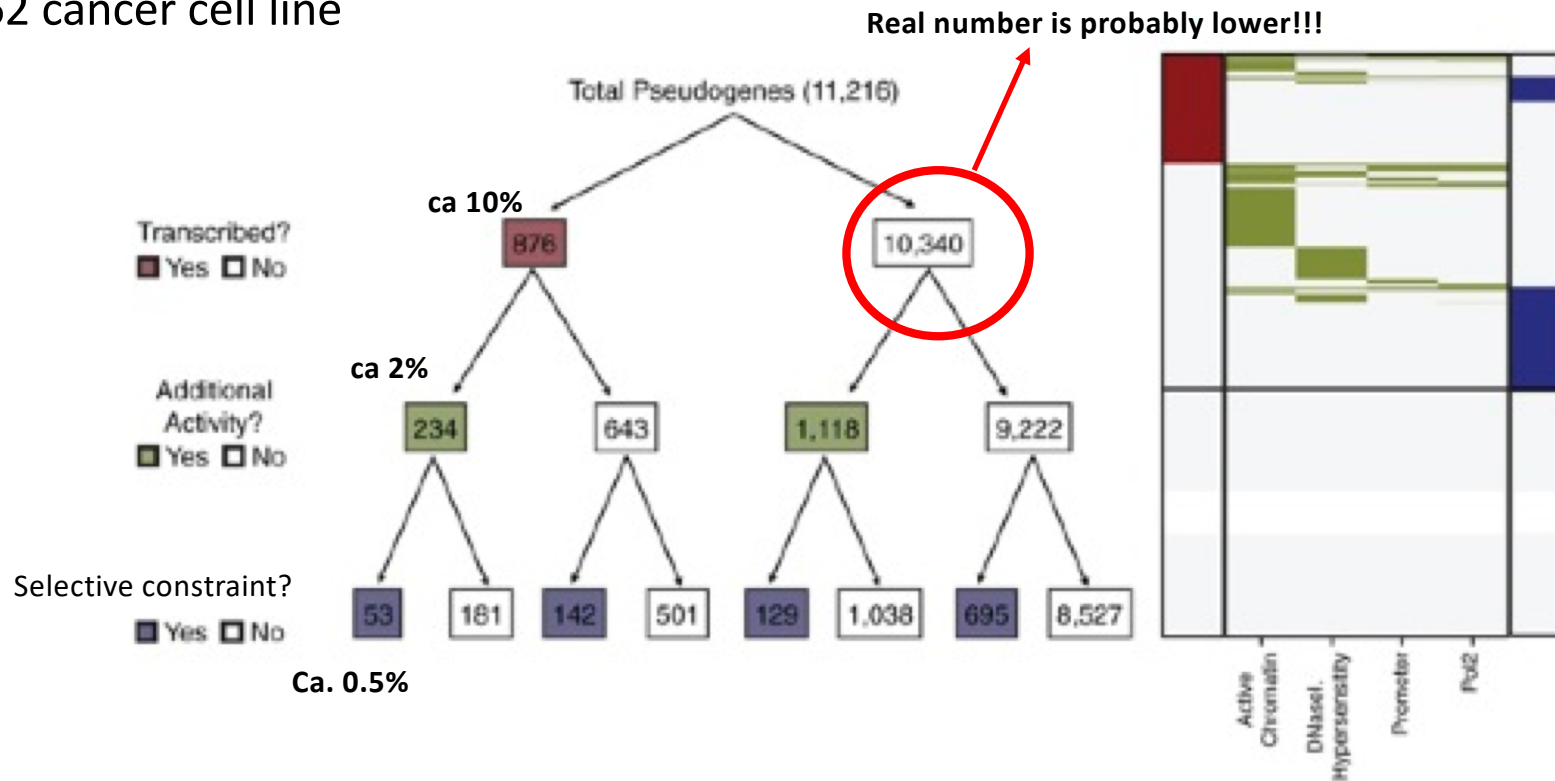
Reason:

→ lower expression

→ expressed pseudogenes do not show marks in an uniform manner

Pseudogenes are a diversified group of genetic elements

K562 cancer cell line



→ few pseudogenes show consistently active signals across all biological features that describe gene activity

→ many pseudogenes show little or no activity

Figure 12 Summary of pseudogene annotation and case studies. (a) A heatmap showing the annotation for transcribed pseudogenes including active chromatin segmentation, DNaseI hypersensitivity, active promoter, active Pol2, and conserved sequences. Raw data were from the K562 cell line. (b) A transcribed duplicated pseudogene (Ensembl gene ID: ENST00000434500.1; genomic location, chr7: 65216129-65228323)

Pseudogenes are a diversified group of genetic elements

(b)

Transcribed With Additional Activity

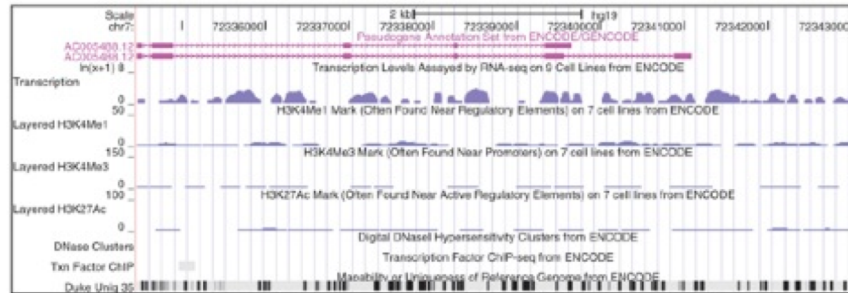


Transcribed
DNase hypersensitive sites
Histonemarks
Transcription factor

Pseudogene
under selective constraint
→ maintained

(c)

Transcribed Only

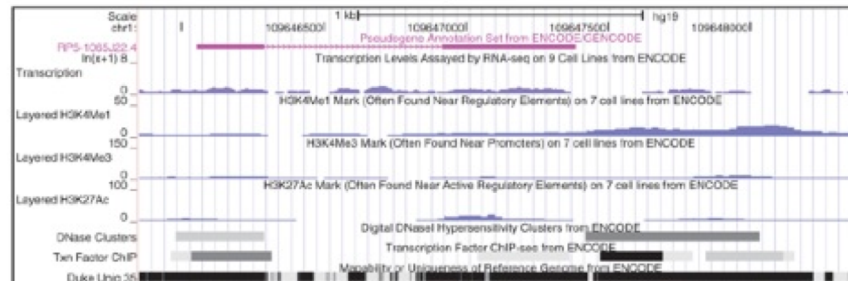


Transcribed
DNase hypersensitive sites
Histonemarks
Transcription factor

Pseudogenes
under low selective constraints
→ This stage also involves
acquisition of new splice sites –
resembles a stage of testing new
mutations for evolutionary
advantage. Result:
A. dying pseudogene or
B. acquisition of critical feature
leading to the resurrection to
become a functional pseudogene

(d)

Partially Active



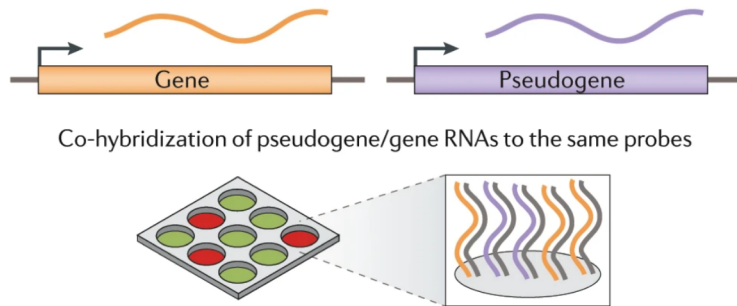
Transcribed
DNase hypersensitive sites
Histonemarks
Transcription factor

Figure 12 Summary of pseudogene annotation and case studies. (a) A heatmap showing the annotation for transcribed pseudogenes including active chromatin segmentation, DNase hypersensitivity, active promoter, active Pol2, and conserved sequences. Raw data were from the K562 cell line. (b) A transcribed duplicated pseudogene (Ensembl gene ID: ENST00000434500.1; genomic location, chr7: 65216129-65228323) showing consistent active chromatin accessibility, histone marks, and TFBSs in its upstream sequences. (c) A transcribed processed pseudogene (Ensembl gene ID: ENST00000355920.3; genomic location, chr7: 72333321-72339656) with no active chromatin features or conserved sequences. (d) A non-transcribed duplicated pseudogene showing partial activity patterns (Ensembl gene ID: ENST00000429752.2; genomic location, chr1: 109646053-109647388). (e) Examples of partially active pseudogenes. E1 and E2 are examples of duplicated pseudogenes. E1 shows UG71A2P

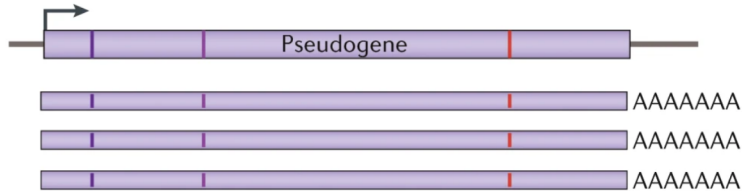
In light of these examples, we believe that the partial activity patterns are reflective of the pseudogene evolutionary process, where a pseudogene may be in the process of either resurrection as a ncRNA or gradually losing its functionality. Understanding why pseudogenes show partial activity may shed light on pseudogene evolution and function.

Challenges in studying pseudogene lncRNAs

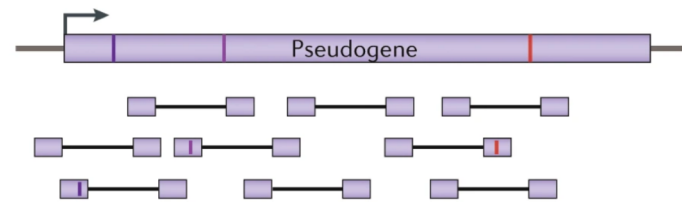
a Hybridization to DNA microarrays



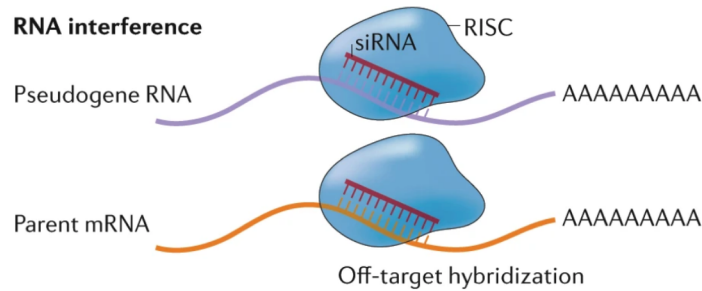
c Long-read cDNA sequencing



b Short-read cDNA sequencing



d RNA interference



a | Microarray analysis cannot distinguish between parent gene and pseudogene expression due to the co-hybridization of the two similar transcripts to the same oligonucleotide probes. **b** | Short-read cDNA sequencing is unable to confidently distinguish many pseudogene RNAs from their parent mRNAs due to insufficient nucleotide differences per read. **c** | Long-read cDNA sequencing allows accurate quantification of pseudogene RNAs due to a higher number of specific differences per read. **d** | RNA interference is poorly suited to analysis of pseudogenes due to off-target hybridization of small interfering RNAs (siRNAs) to the parent gene. RISC, RNA-induced silencing complex.

Expression

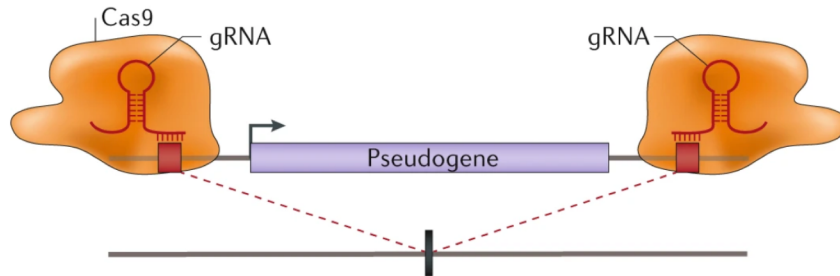
Validate with gene/pseudogene specific RT-PCR oligos

Knock-down with specific siRNAs and validation by specific RT-PCR

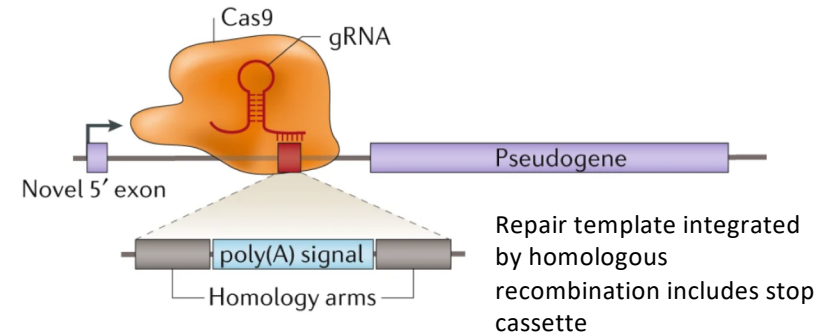
Challenges in studying pseudogene lncRNAs

specific alteration of pseudogenes by CRISPR technology

c CRISPR-mediated deletion

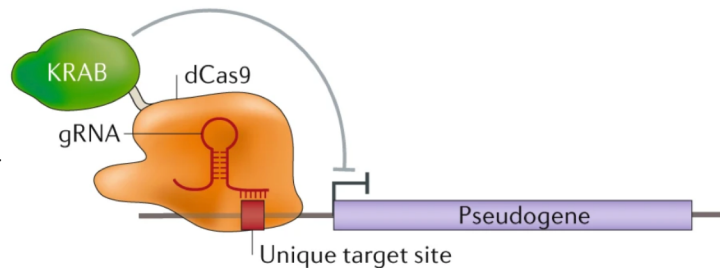


d Integration of transcriptional terminator



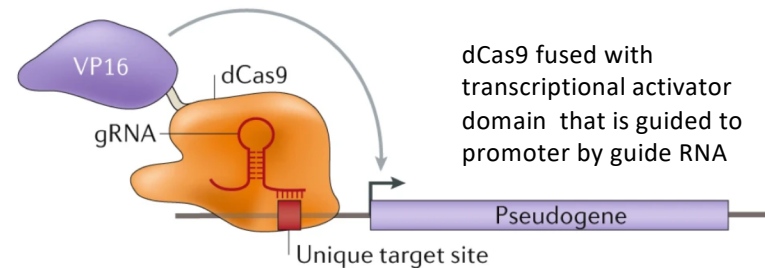
e CRISPR interference

dCas9 fused with transcriptional repressor domain that is guided to promoter by guide RNA



f CRISPR activation

dCas9 fused with transcriptional activator domain that is guided to promoter by guide RNA



- c)** CRISPR–Cas9 genome engineering allows deletion of pseudogenes by targeting unique flanking sequences. **d** | If pseudogenes have novel 5' exons, transcriptional terminators can be introduced by homologous recombination to deplete pseudogene transcripts. **e** | CRISPR-based transcriptional interference (CRISPRi) enables depletion of pseudogene transcription by targeting a dCas9–KRAB fusion protein to unique sequences upstream of the transcriptional start site. **f** | CRISPR-based transcriptional activation (CRISPRa) enables activation of pseudogene transcription by targeting a dCas9–VP16 fusion protein to unique sequences upstream of the transcription start site. dCas9, catalytically inactive Cas9; gRNA, guide RNA; PAM, protospacer-adjacent motif; poly(A), polyadenylation.

Evidence for functional relevance of pseudogene encoded lncRNAs

Pseudogene	Organism	Parent gene function	Biological impact of pseudogene	Regulatory mechanism	Refs
<i>PGK2</i>	Human	Phosphoglycerate kinase	Testis-specific enzyme that catalyses the conversion of 1,3-bisphosphoglycerate to 3-phosphoglycerate during glycolysis	Protein-based	32,33
<i>POU5F1B</i>	Human	Pou-domain transcription factor	Putative transcription factor that promotes tumour growth; amplified in gastric cancer	Protein-based	36
<i>NANOGP8</i>	Human	Homeodomain transcription	Putative transcription factor that promotes cell proliferation	Protein-based	126
<i>ΨCX43 (GJA1P1)</i>	Human	Gap junction protein	Putative gap junction protein that inhibits cell growth	Protein-based	127
<i>NOTCH2NL</i>	Human	Transmembrane receptor	Activates NOTCH signalling by sequestering the inhibitory ligand DELTA; expands cortical progenitor population	Protein-based	37,38
<i>SRGAP2C</i>	Human	Slit-Robo Rho GTPase activating protein	Dimerizes with SRGAP2, inhibiting its function	Protein-based	39,40
<i>NOS pseudogene</i>	<i>Lymnaea stagnalis</i>	Nitric oxide synthase	Antisense RNA prevents the translation of <i>NOS</i> by forming an RNA-RNA hybrid	RNA-based	41
<i>PTENP1</i>	Human	Phosphatase that converts PtdIns(4,5)P ₂ to PtdIns(3,4,5)P ₃ , inhibiting PI3K-AKT signalling	Increases PTEN expression by sequestering microRNAs; can act as a tumour suppressor	RNA-based	45
<i>BRAFP1</i>	Human	Serine/threonine protein kinase	Increases BRAFP1 expression by sequestering microRNAs; can act as an oncogene	RNA-based	46
<i>HMGA1-p</i>	Human	High-mobility-group chromatin protein	Inhibits HMGA1 expression by competing for the RNA-stabilizing protein αCP1	RNA-based	128
<i>Lethe</i>	Mouse	Ribosomal protein subunit S15A	Directly binds to and inhibits NF-κB, modulating inflammatory responses	RNA-based	44
<i>RNA5SP141</i>	Human	5S ribosomal RNA	Binds to RIG-I during herpesvirus infection, inducing interferon expression	RNA-based	28
<i>OCT4pg5-as</i>	Mouse	Pou-domain transcription factor	Suppresses OCT4 expression by increasing EZH2 occupancy at the <i>OCT4</i> promoter	RNA-based	129
<i>HBBP1</i>	Human	β-globin	Facilitates switching of fetal to adult globin expression by regulating contacts with the locus control region	DNA-based	50
Immunoglobulin pseudogenes	Chicken	Immunoglobulin segments	Generate immunoglobulin diversity by gene conversion	DNA-based	130,131
<i>PRSS3P2</i>	Human	Cationic trypsinogen	Causes hereditary pancreatitis by gene conversion with <i>PRSS3</i>	DNA-based	55
<i>CYP21A2P</i>	Human	21-Hydroxylase, a cytochrome P450 enzyme	Causes adrenal hyperplasia by gene conversion with <i>CYP21A2</i>	DNA-based	56
<i>CYP2A7</i>	Human	Cytochrome P450 enzyme	Increases <i>CYP2A6</i> mRNA stability due to a 3' UTR polymorphism formed by gene conversion	DNA-based	132

PI3K, phosphoinositide 3-kinase; PtdIns(4,5)P₂, phosphatidylinositol 4,5-bisphosphate; PtdIns(3,4,5)P₃, phosphatidylinositol 3,4,5-trisphosphate; PTEN, phosphatase and tensin homologue; UTR, untranslated region.