



Comparative Analysis of Non-Coding Transcription using ENCODE and modENCODE data

ENCODE Users Meetings 2016

Joel Rozowsky

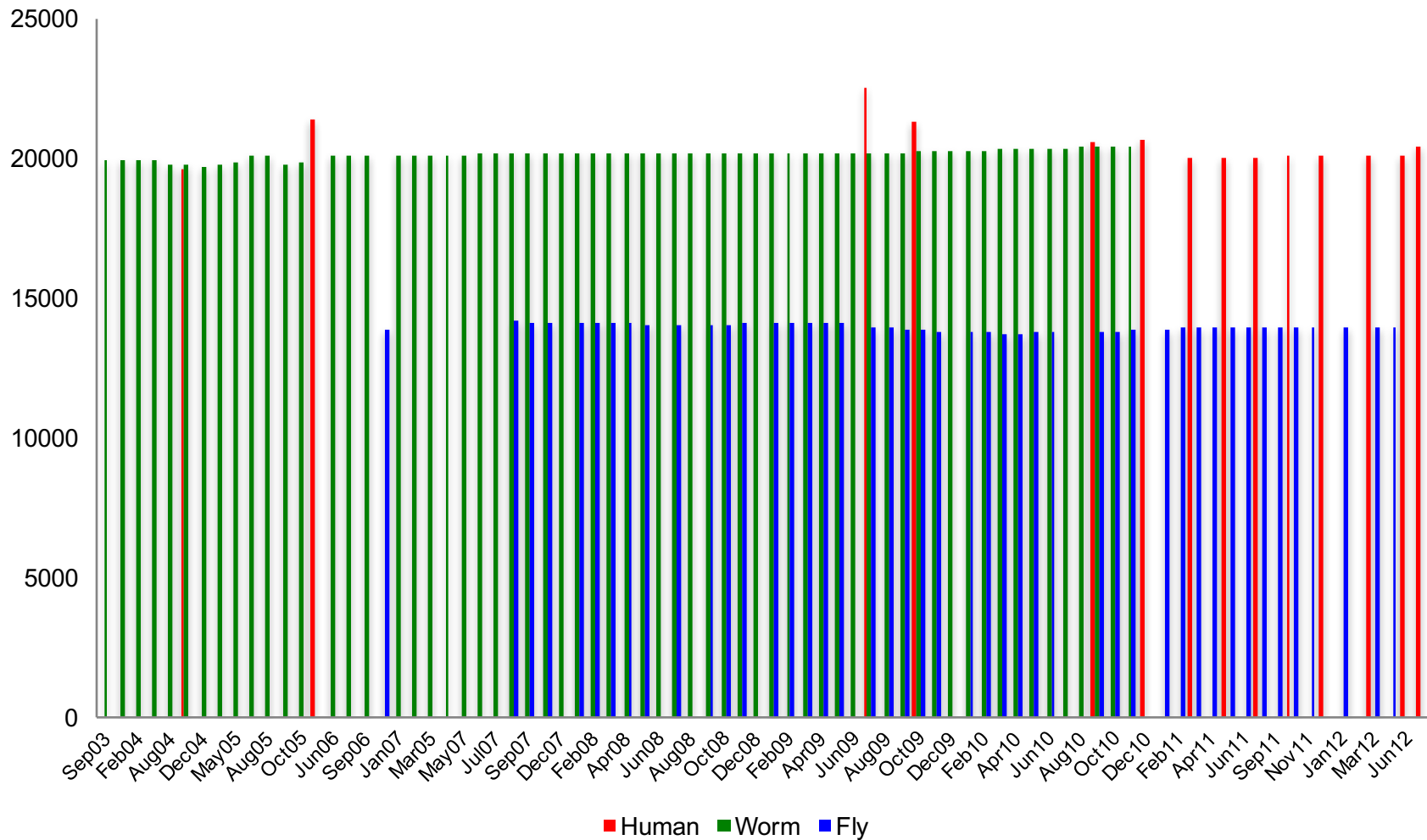
Yale

8 June 2016

Outline

- **The Discovery of Pervasive Transcription**
- **Pervasive Transcription, Take 2**
 - The advent of Next Gen Sequencing
- **Drilling into one type of pervasive transcription:
Transcribed Pseudogenes**

During the genome annotation era, protein-coding gene counts in worm, fly & human have remained fairly constant

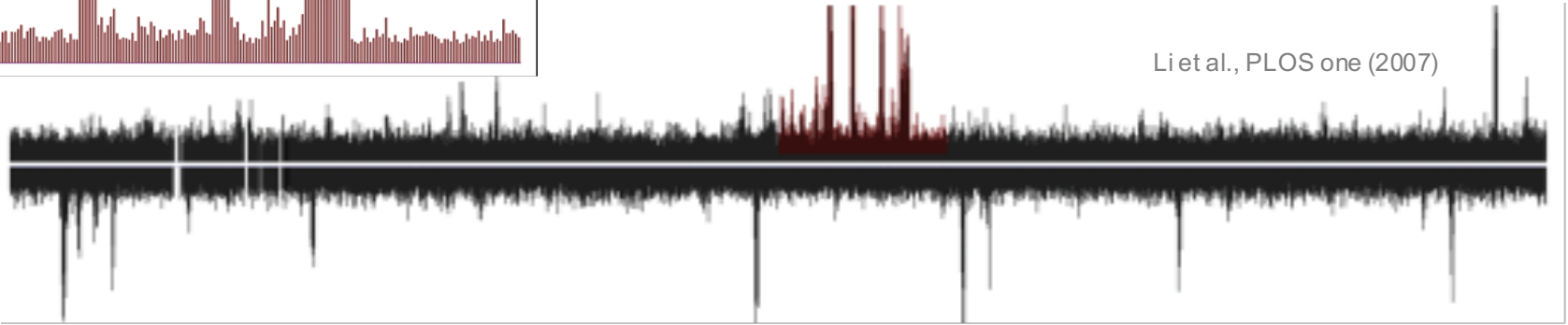
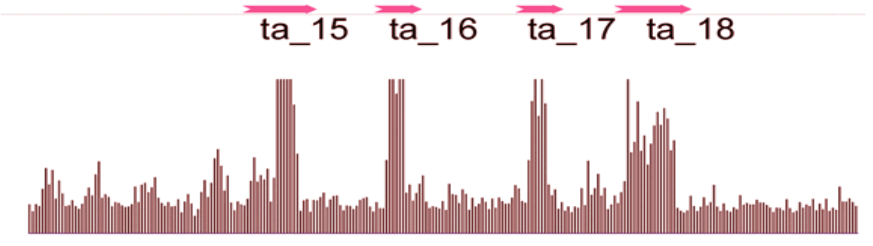
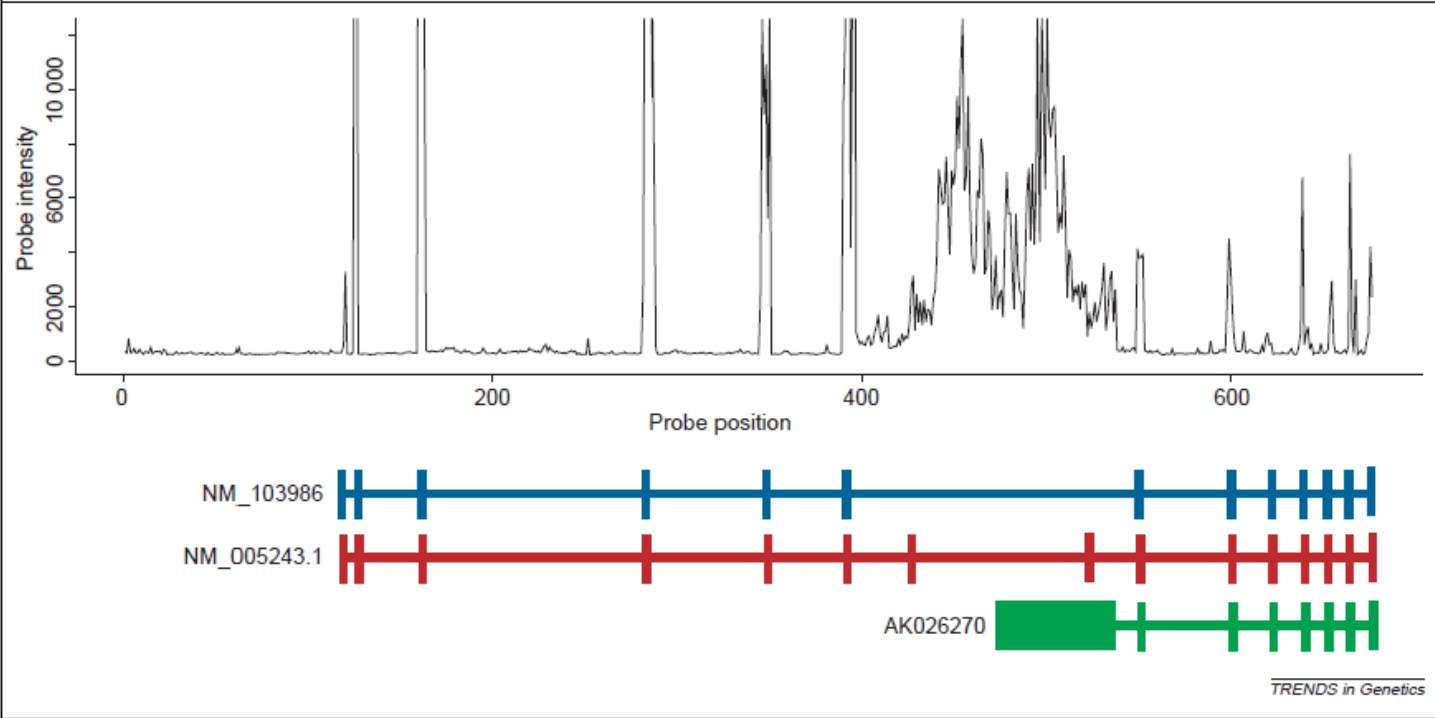


Discovery of Pervasive Transcription: Dark Matter of the Genome

- **With the advent of custom tiling arrays it was shown that a significant portion of the human genome (outside known protein coding genes) is transcribed**
 - **Chr 21/22: Kapranov et al....Gingeras (2002) Science**
 - “When compared with the sequence annotations available for these chromosomes, it is noted that as much as an **order of magnitude more of the genomic sequence is transcribed** than accounted for by the predicted and characterized exons.”
 - **Also (Chr 22): Rinn et al... Snyder (2003) Genes & Dev.**
 - **Whole Genome:**
 - Bertone et al. (2004) Science & Cheng et al. (2005) Science
 - **Pilot ENCODE (1% Genome) Nature (2007)**
 - “...our studies show that **14.7% of the bases** represented in the unbiased tiling arrays are transcribed in at least one tissue sample.”

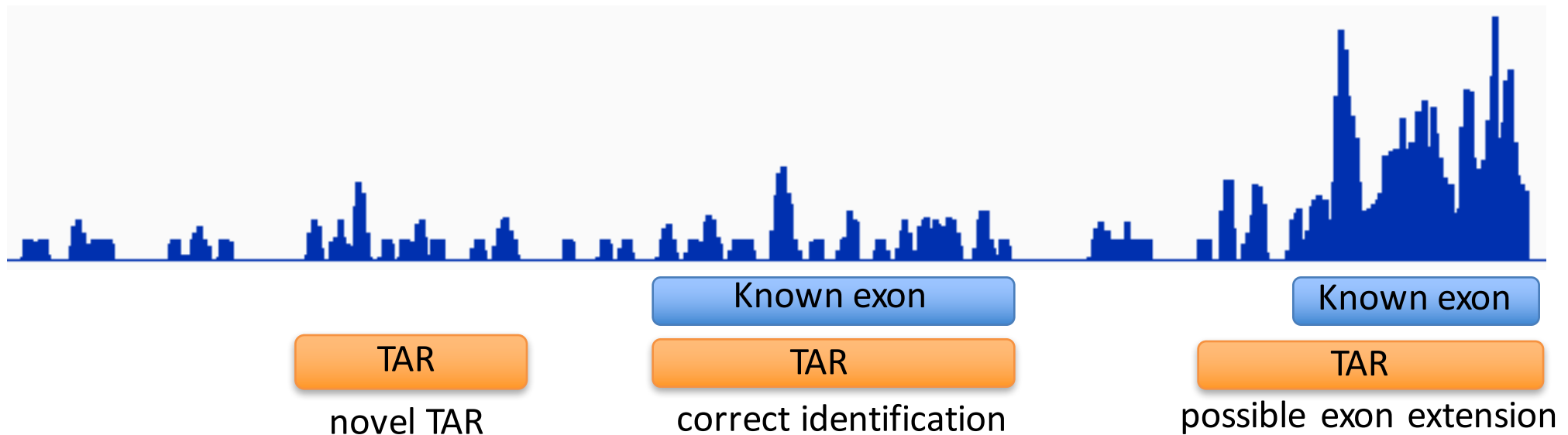
Noisy Raw Signal from Tiling Arrays (Transcription)

Johnson et al. (2005) TIG, 21, 93-102.



TARs (novel RNA contigs) from Segmenting Transcriptional Signal

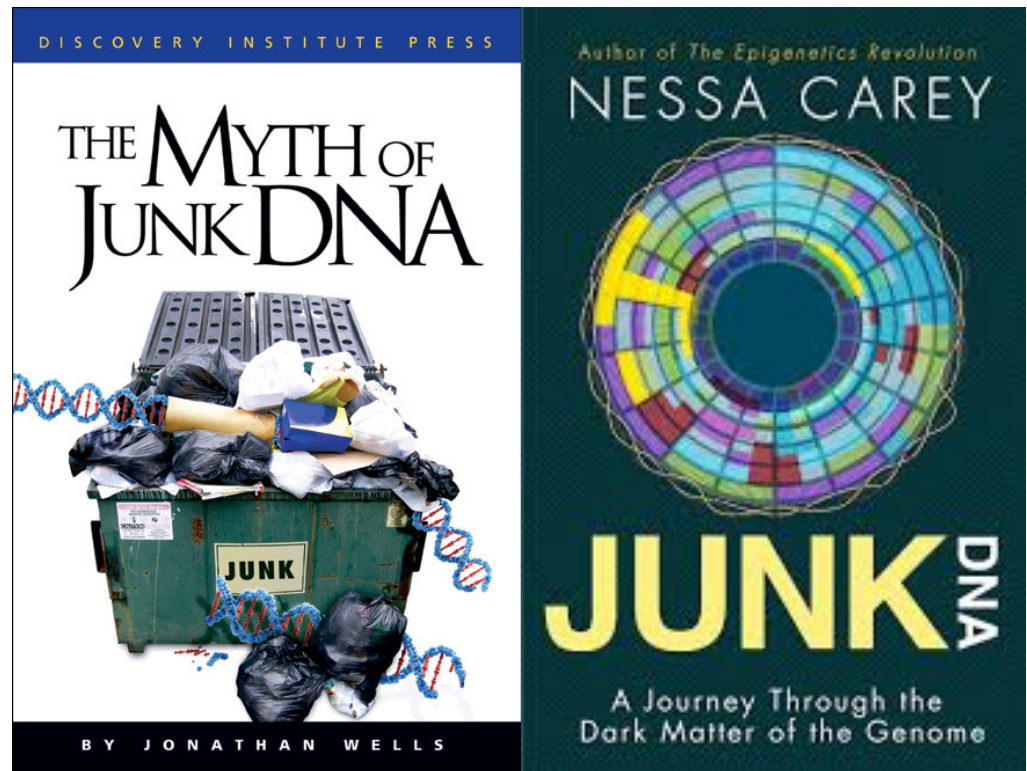
- Cluster reads setting minimum-run and maximum gap parameters for newly identified transcribed regions (TARs) [called TransFrag by Gingeras et al.]



Controversy of Pervasive Transcription

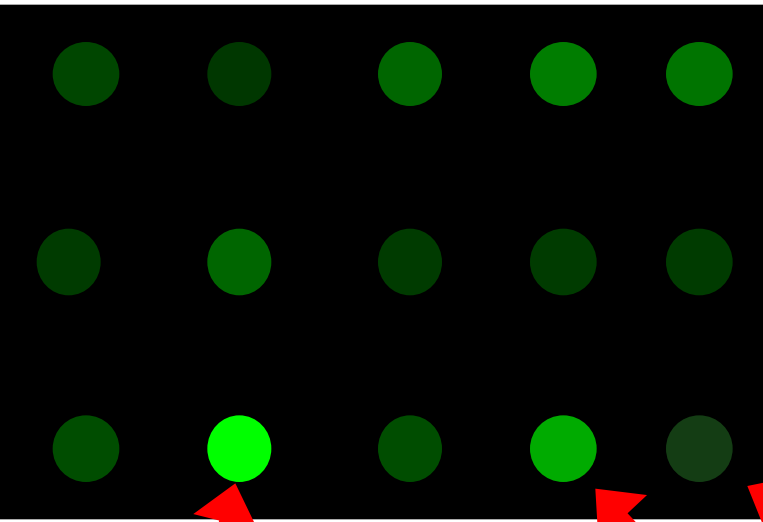
Over the last decade this result has been somewhat controversial
(Clark et al (2011) PLoS Bio)

“Current estimates indicate that only about 1.2% of the mammalian genome codes for amino acids in proteins. However, mounting evidence over the past decade has suggested that the vast majority of the genome is transcribed, well beyond the boundaries of known genes, a phenomenon known as pervasive transcription. Challenging this view, an article published in PLoS Biology by van Bakel et al. **concluded that ‘the genome is not as pervasively transcribed as previously reported’ and that the majority of the detected low-level transcription is due to technical artefacts** and/or background biological noise.”

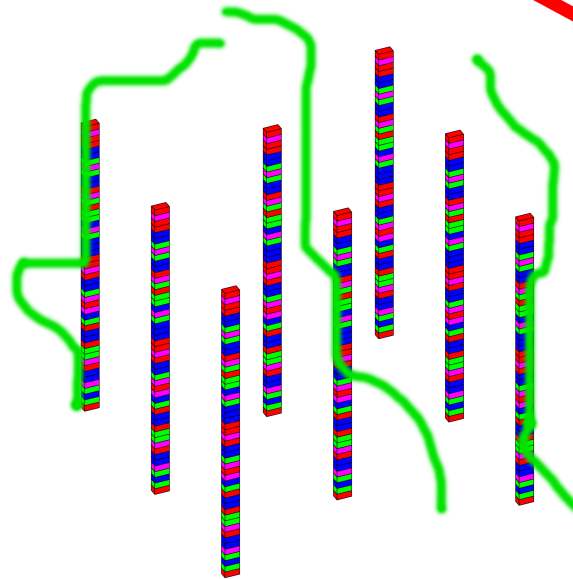


Cross-Hyb. – Specific & Non-specific

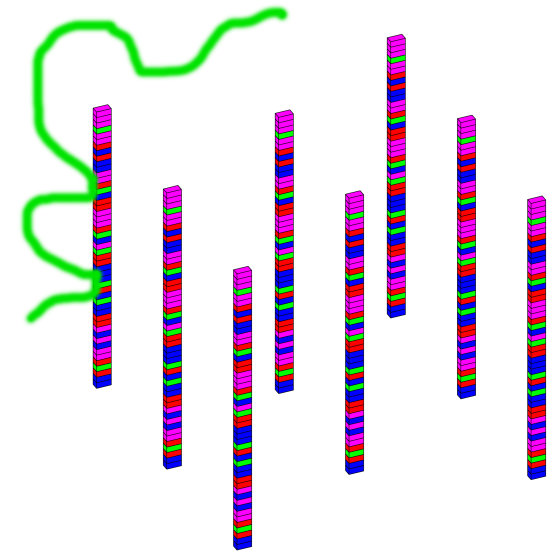
- Perfect match (PM): probe binding intended target
- Specific cross-hyb.: probes binding non-PM targets with a small number of mismatches
- Non-specific cross-hyb.: probes binding targets with many mismatches, due to general stickiness of oligos



Perfect Match



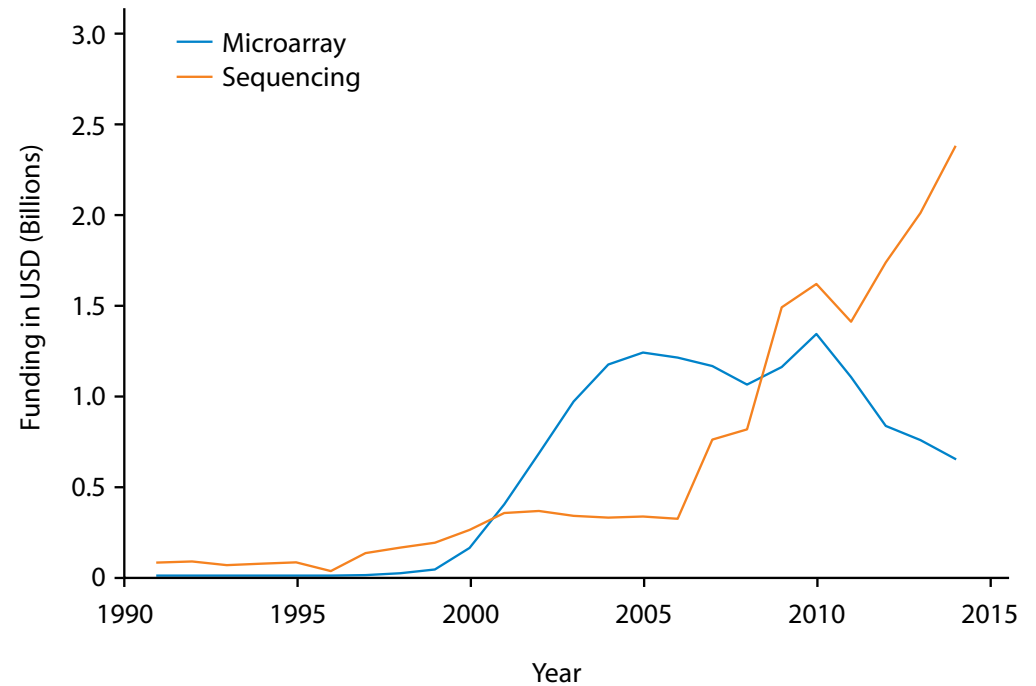
Specific Cross-hyb.



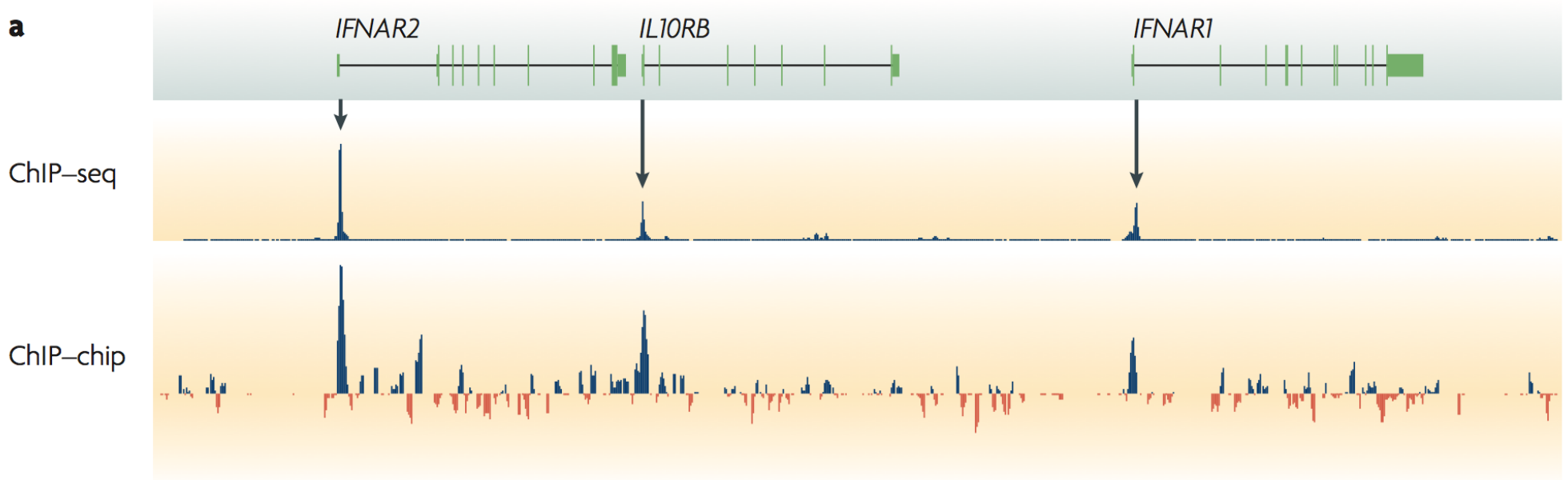
Non-specific Cross-hyb.

Advent of NGS: Much Cleaner Signals than Tiling Arrays, Supplanting this Technology

National Institutes of Health funding for
'microarray' and 'sequencing' projects



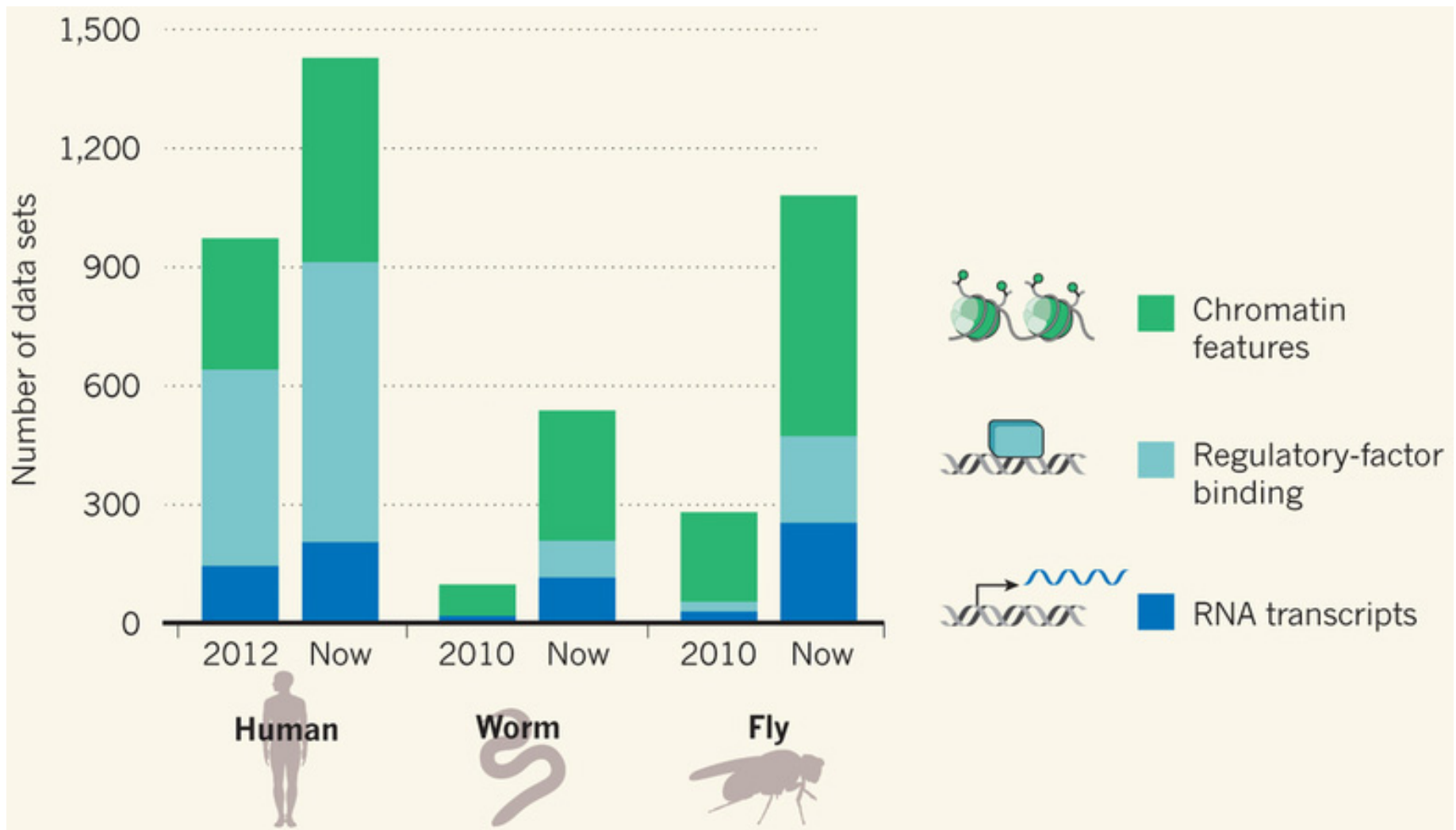
a



Comparative ENCODE Functional Genomics Resource

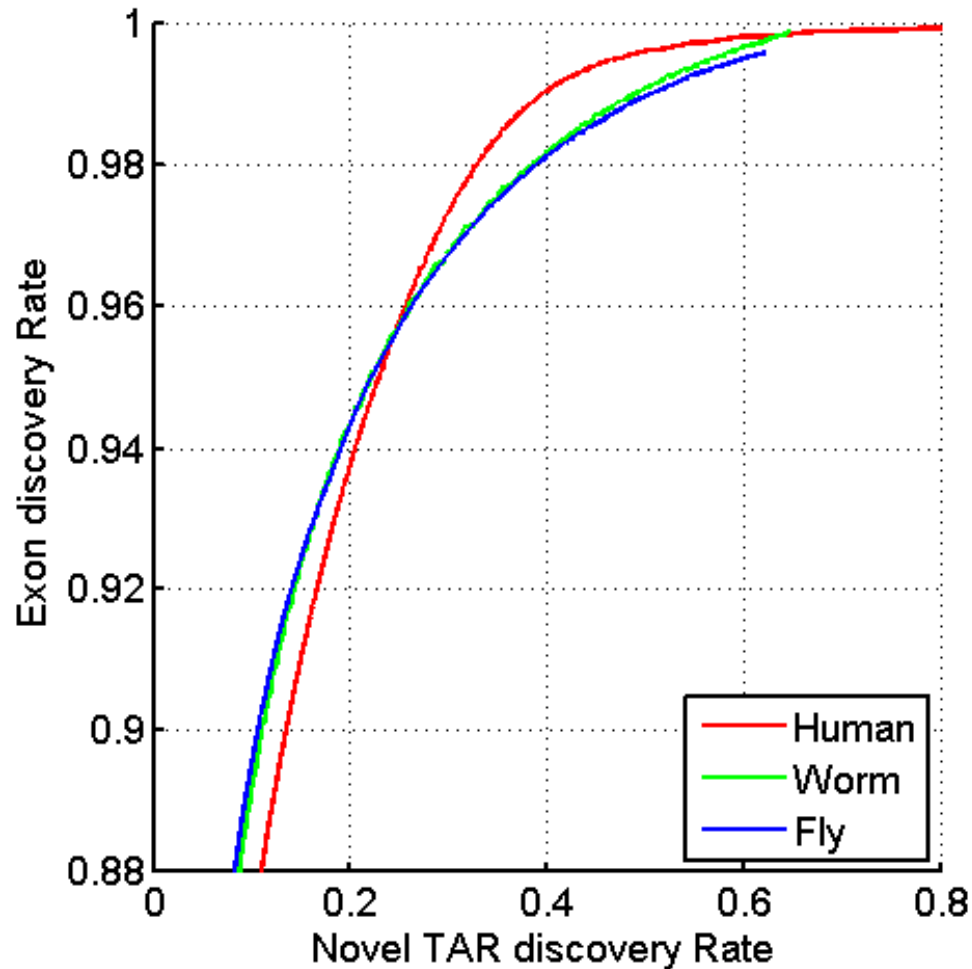
(EncodeProject.org/modENCODE.org)

- Broad sampling of conditions across transcriptomes & regulomes for human, worm & fly
 - embryo & ES cells
 - developmental time course (worm-fly)
- In total: ~3000 datasets (~130B reads)



Uniform Annotation of non-coding Elements

- Uniformly processed the RNA-seq expression compendium

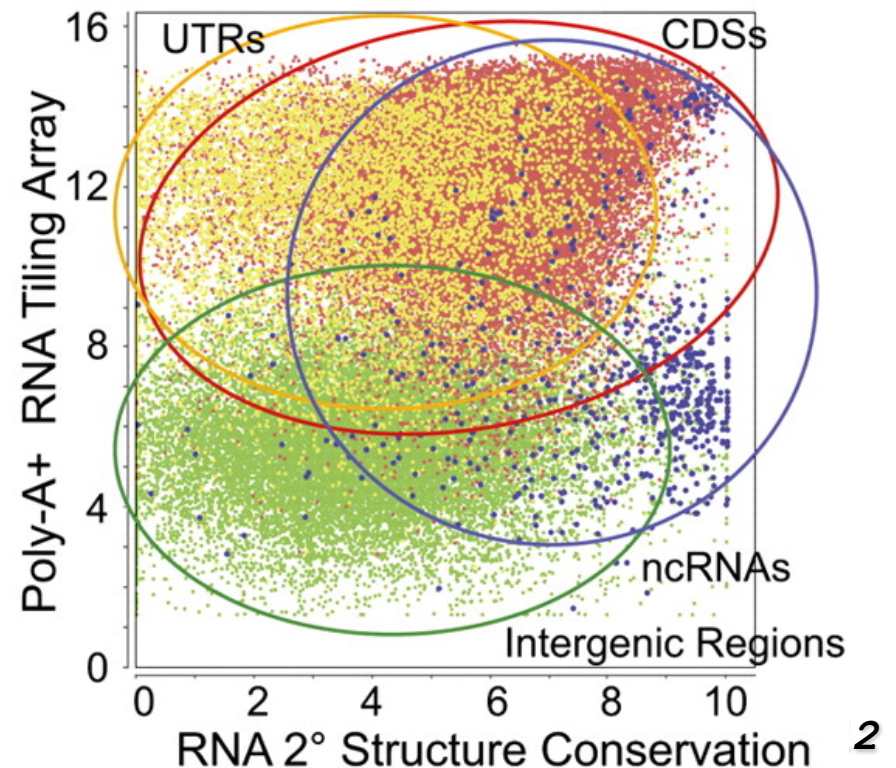
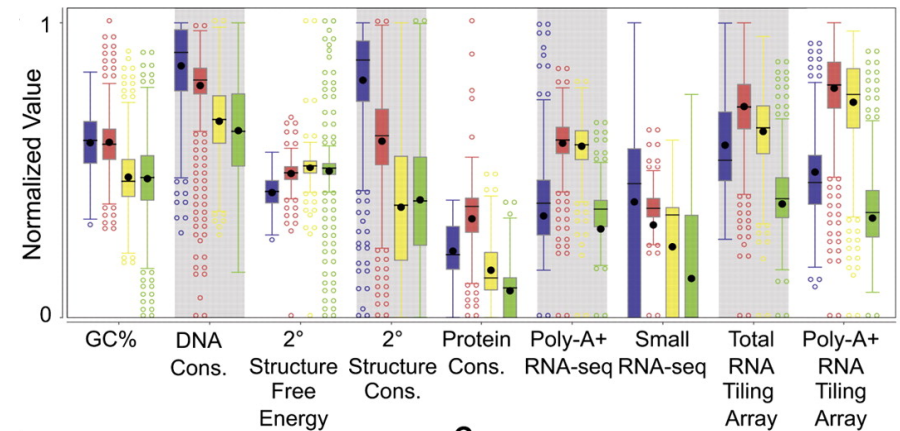
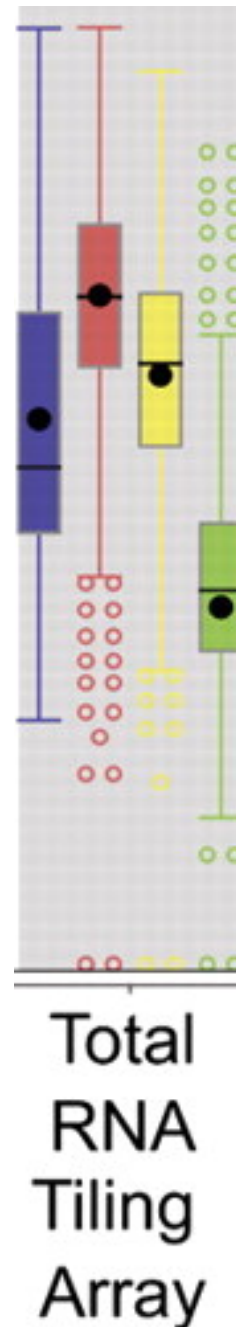


Gold-standard Set

■ Known ncRNAs ■ CDSs ■ UTRs ■ Intergenic Regions

lncRNA:
Machine-learning
Identification of
many candidate
ncRNAs through
evidence integration

- No single feature (e.g. expr. expts., conservation, or sec. struc.) finds all known ncRNAs => combine features in stat. model
- 90% PPV, 13 of 15 tested validate



Annotated ncRNAs

		Human			Worm			Fly			
		Elements	Genome Coverage		Elements	Genome Coverage		Elements	Genome Coverage		
			Kb	%		Kb	%		Kb	%	
mRNAs (exons)		20,007	86,560	3.0	21,192	34,437	34.3	13,940	35,970	28.0	
Pseudogenes		11,216	27,089	0.95	881	1,343	1.3	145	155	0.12	
Annotated ncRNAs	Comparable ncRNAs	pri-miRNA	58	1,158	0.04	44	16	0.02	43	300	0.23
		pre-miRNAs	1,756	162	0.006	221	20	0.02	236	22	0.02
		tRNAs	624	47	0.002	609	45	0.04	314	22	0.02
		snoRNAs	1,521	168	0.006	141	16	0.02	287	34	0.03
		snRNAs	1,944	210	0.007	114	14	0.01	47	7	0.006
		lncRNAs	10,840	10,581	0.37	233	184	0.18	852	868	0.68
	Other ncRNAs	5,411	3,268	0.11	40,104	2,329	2.3	376	2,103	1.6	
	nc-piRNA loci	88	1,272	0.04	35,329	449	0.45	27	1,473	1.1	
Total		22,154	17,770	0.62	41,466	2,611	2.6	2,155	3,279	2.6	

Identify non-canonical transcription in regions of the genome excluding mRNA exons, pseudogenes or annotated ncRNAs.

& Non-Canonical Transcription

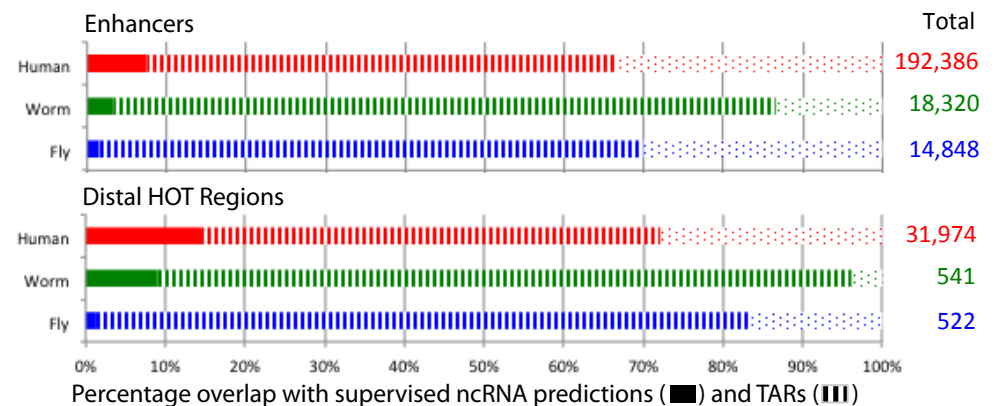
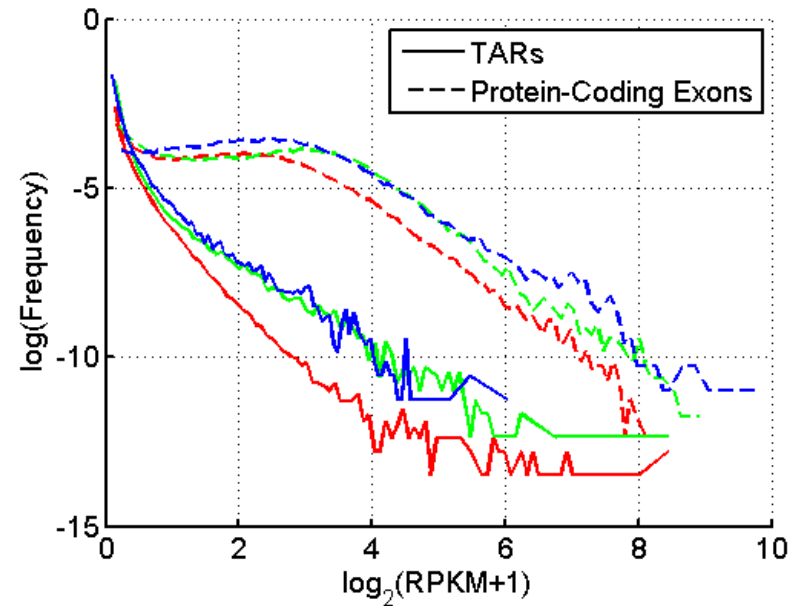
	Human			Worm			Fly		
	Elements	Genome Coverage		Elements	Genome Coverage		Elements	Genome Coverage	
		Kb	%		Kb	%		Kb	%
→ Total ncRNAs	22,154	17,770	0.62	41,466	2,611	2.6	2,155	3,279	2.6
Regions Excluding mRNAs, Pseudogenes or Annotated ncRNAs	283,816	2,731,811	95.5	143,372	63,520	63.3	60,108	89,445	69.6
Transcription Detected (TARs)	708,253	916,401	32.0	232,150	37,029	36.9	83,618	44,256	34.5
Supervised Predictions	104,016	13,835	0.48	2,525	392	0.39	599	164	0.13

- Similar fraction of non-canonical transcription of non-canonical transcription in human, worm and fly
 - 32-37% of each genome

TAR Characterization

Non-canonical transcription (TARs):

- Mostly transcribed at lower levels than protein-coding genes.
- Enrichment for overlap of TARs with ENCODE enhancers and distal HOT regions -> potential enhancer RNAs (eRNAs).

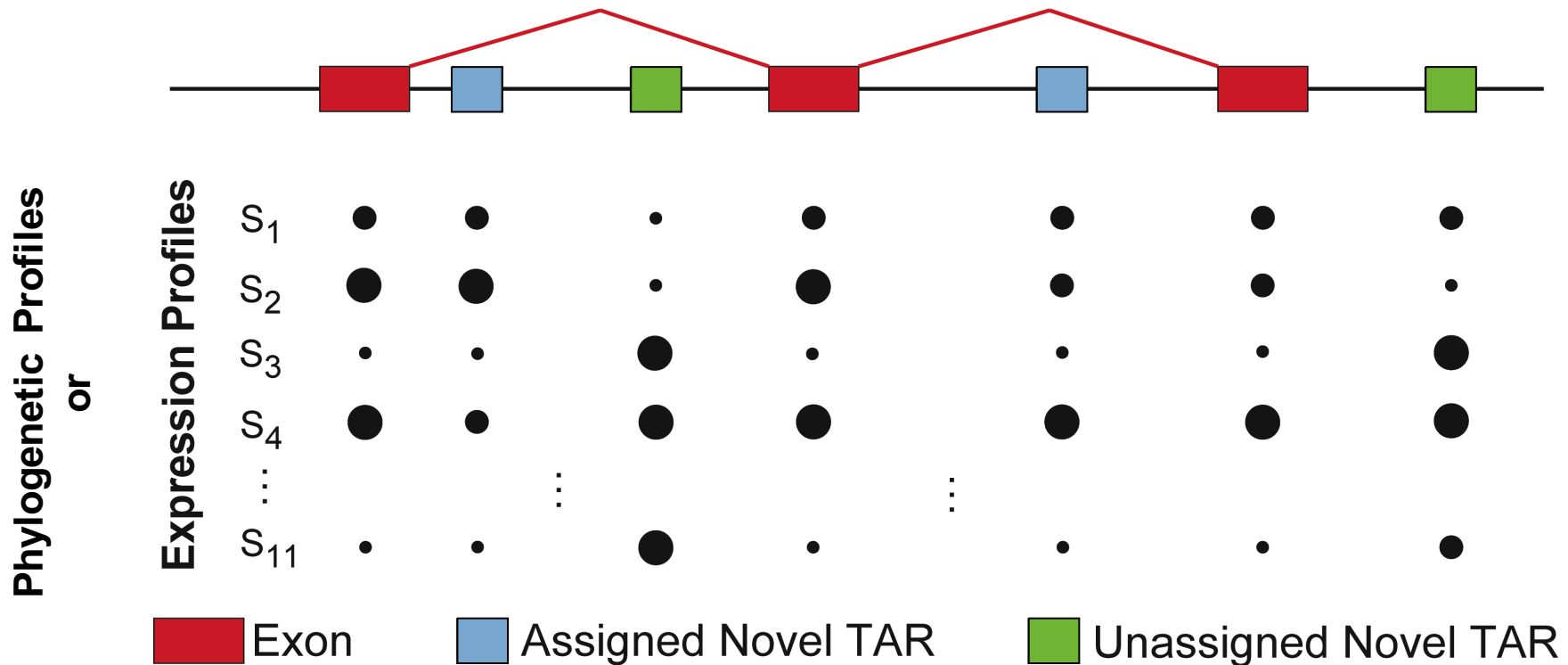


Human, Worm & Fly

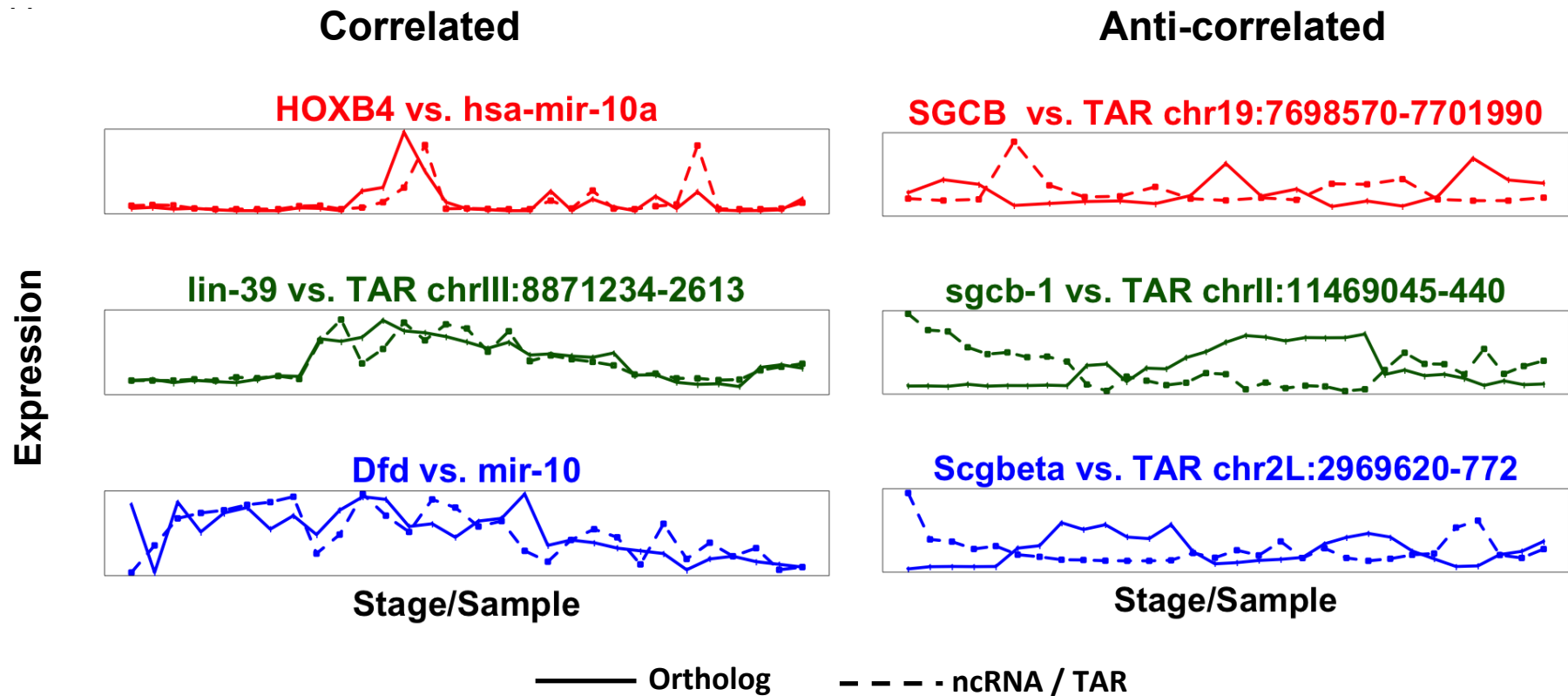
HOT Regions = High TF Co-occupancy

Clustering & Classifying Blocks of Un-annotated Transcription into larger units

Assignment of novel TARs to known gene loci



ncRNAs/TARS can be clustered with known genes



Pseudogenes are among the most interesting intergenic elements

- Formal Properties of Pseudogenes (Ψ G)
 - Inheritable
 - Homologous to a functioning element – ergo a repeat!
 - Non-functional
 - No selection pressure so free to accumulate mutations
 - Frameshifts & stops
 - Small Indels
 - Inserted repeats (LINE/Alu)
 - **What does this mean?** no transcription, no translation?...

Pseudogene Definition

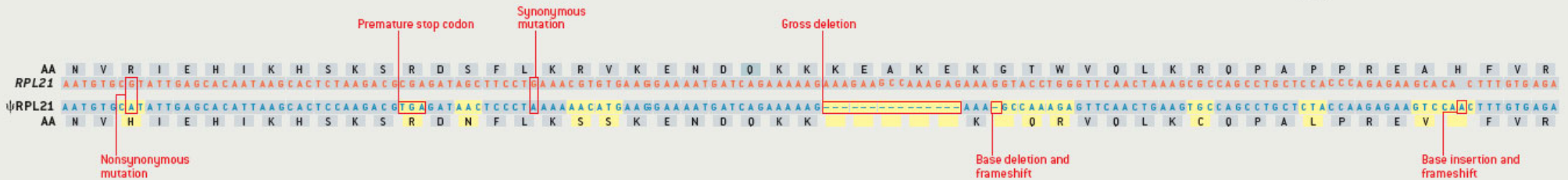
PSEUDOGENE BIRTH AND GENE DEATH

Two distinct processes can duplicate genes, and together they allow genomes to grow and diversify over evolutionary time. If errors in a copy destroy its ability to function as a gene, however, it becomes a pseudogene instead (*right*). The mutations that can kill a gene (*below*) range from gross deletions (such as the loss of the promoter region that signals the start of a gene sequence) to minute changes in the DNA sequence that skew the meaning of the gene's protein-encoding segments, called exons.

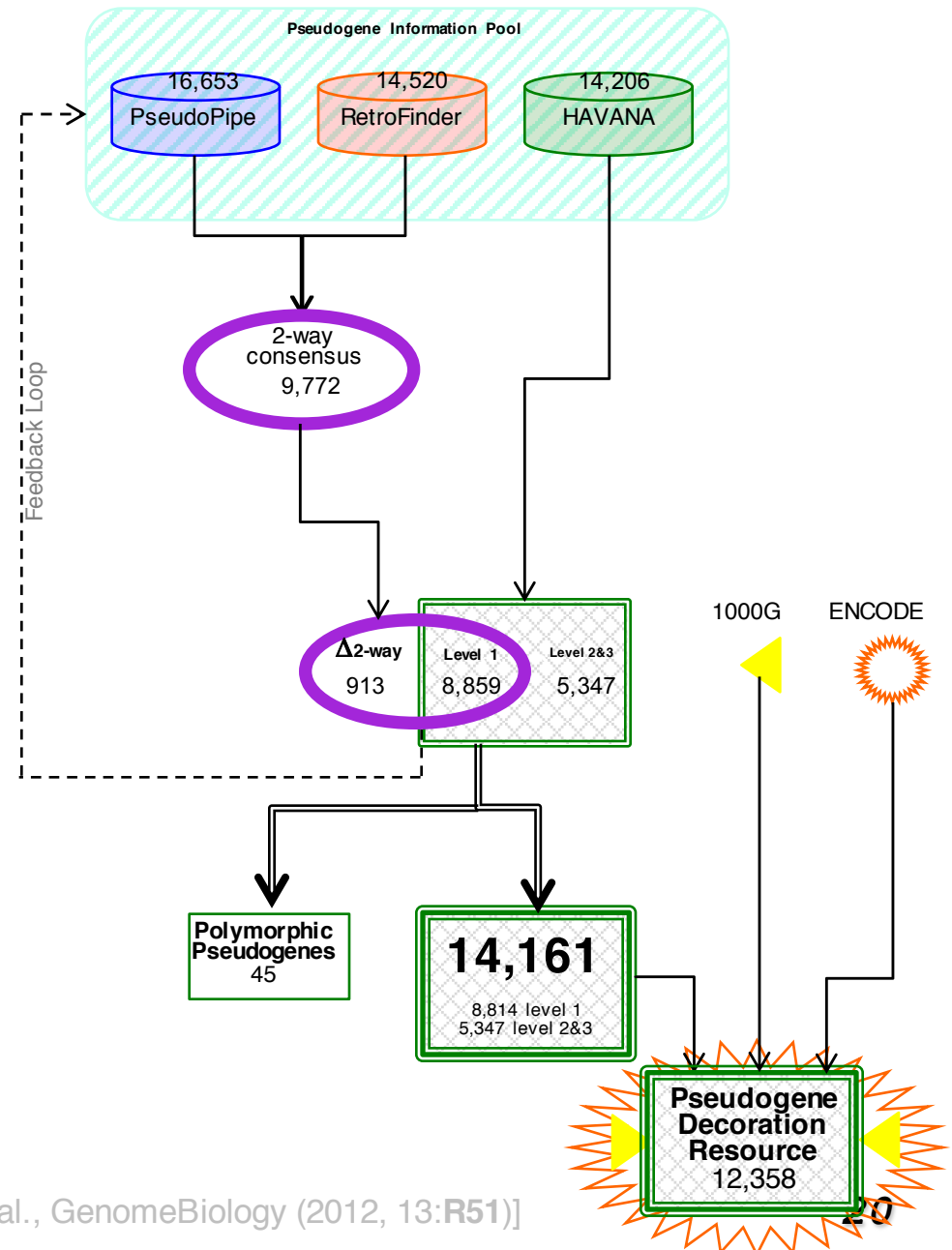
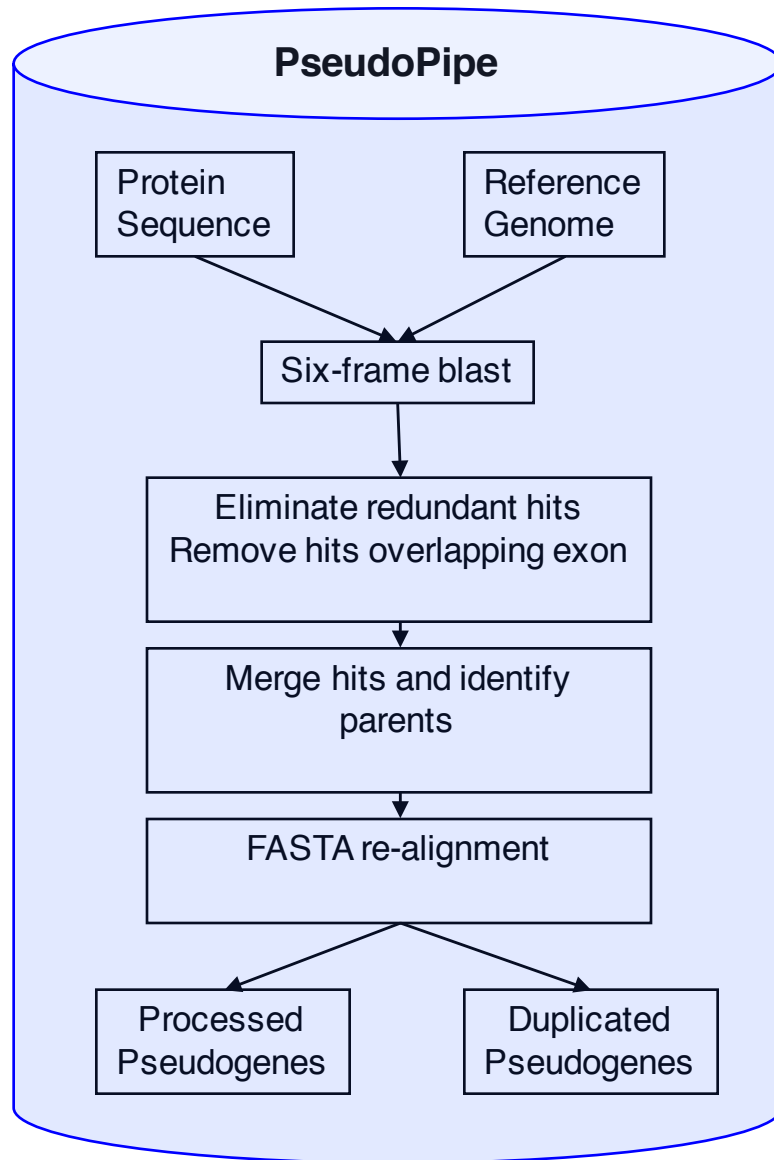
GENE DEATH

Genes die and become pseudogenes when mutations generated during the gene-copying process or accumulated over time render them incapable of giving rise to a protein. Cellular machinery reads the DNA alphabet of nucleotide bases (abbreviated A, C, G, T) in three-base increments called codons, which name an amino acid building block in a protein sequence or encode "stop" signals indicating the end of a gene. Even single-base mutations in codons

can alter their amino acid meaning, and base deletions or insertions can affect neighboring codons by shifting the cellular machinery's reading frame. The alignment shown here of a partial sequence for a human gene (*RPL21*) against one of its pseudogene copies (Ψ *RPL21*), along with each codon's corresponding amino acid (AA), illustrates some of the disabling mutations typically found in pseudogenes.

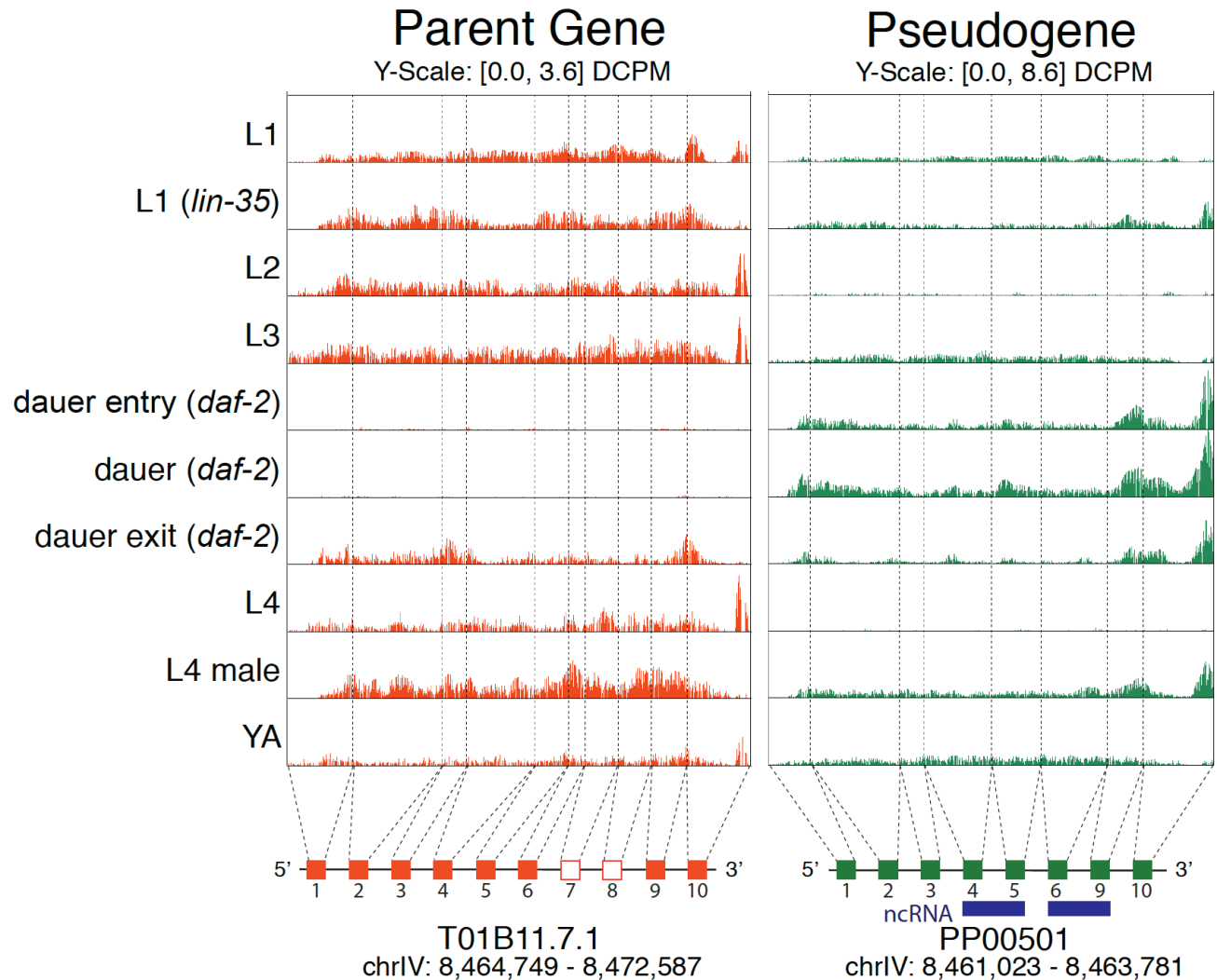


Genome-wide Annotation of Pseudogenes



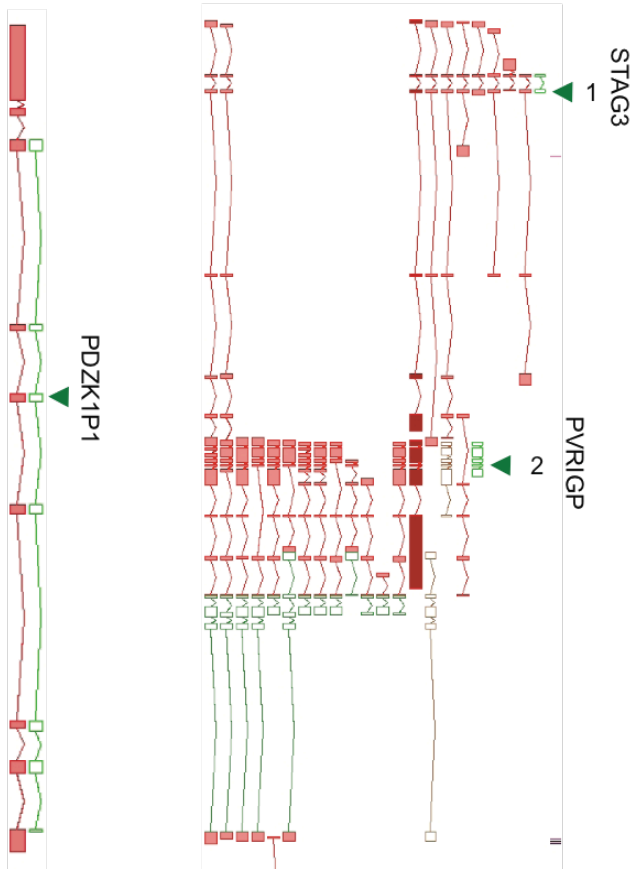
Calling transcribed pseudogenes (while guarding against mis- mapping)

- Counting uniquely mapped reads in RNA-seq:
 - RPKM > 2
 - **1441** human transcribed pseudogenes

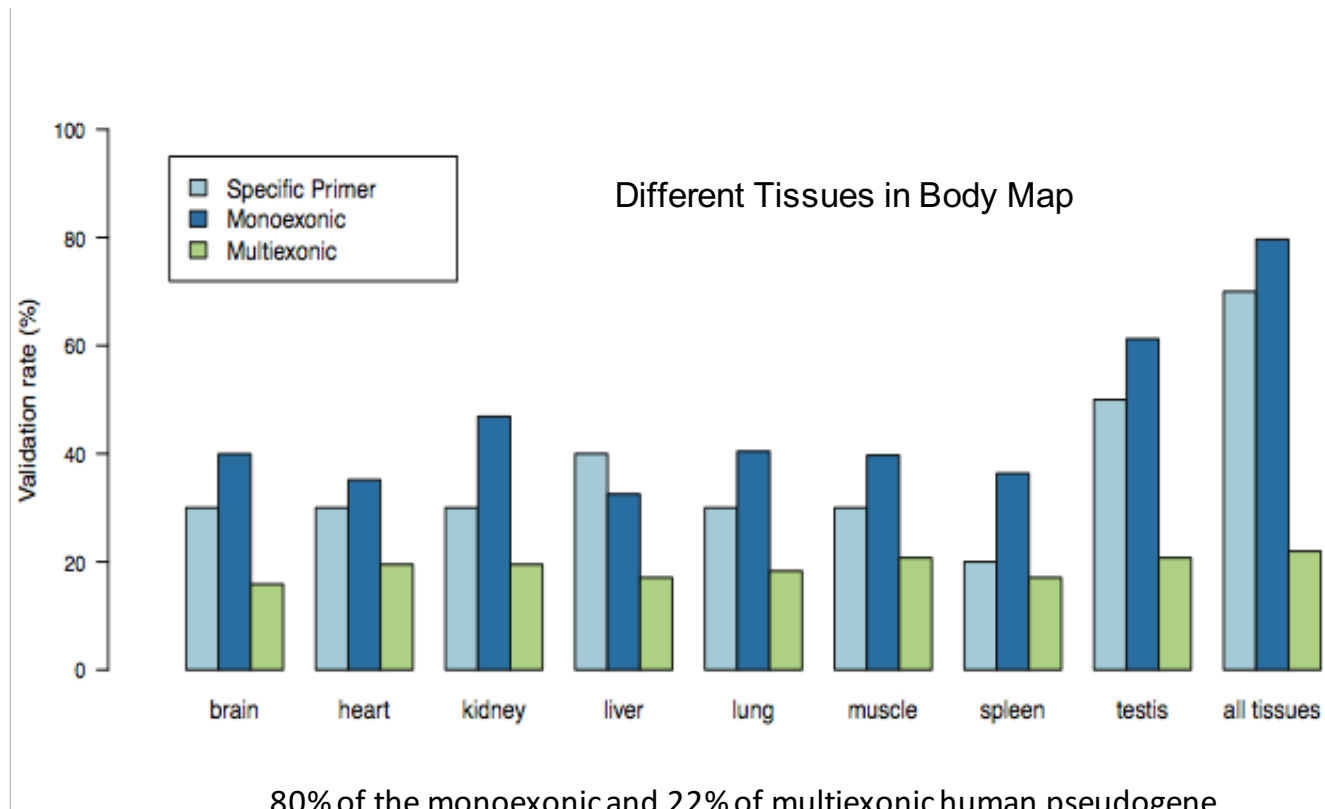


- Using **ESTs** Looking for a discordant expression
(misses pseudogenes co-expressed w. parent)

Validation of Pseudogene Transcription



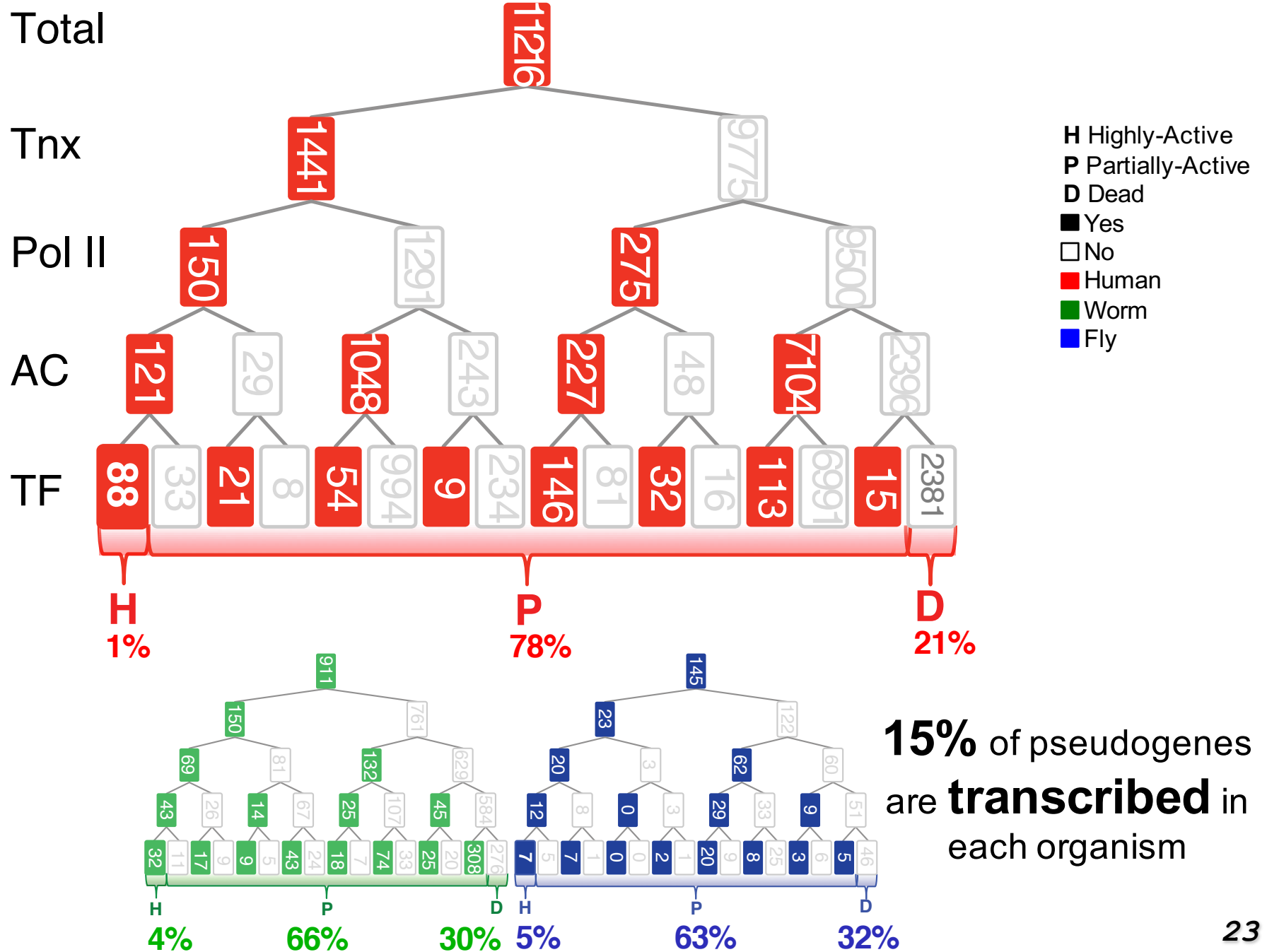
- Pseudogene model
- Protein-coding model
- Processed transcript model
- Indicates pseudogene locus



80% of the monoexonic and 22% of multiexonic human pseudogene models validated using RT-PCR-Seq; 57 of 76 in total

Simple & Complex Ex of Pseudo-gene Transcription

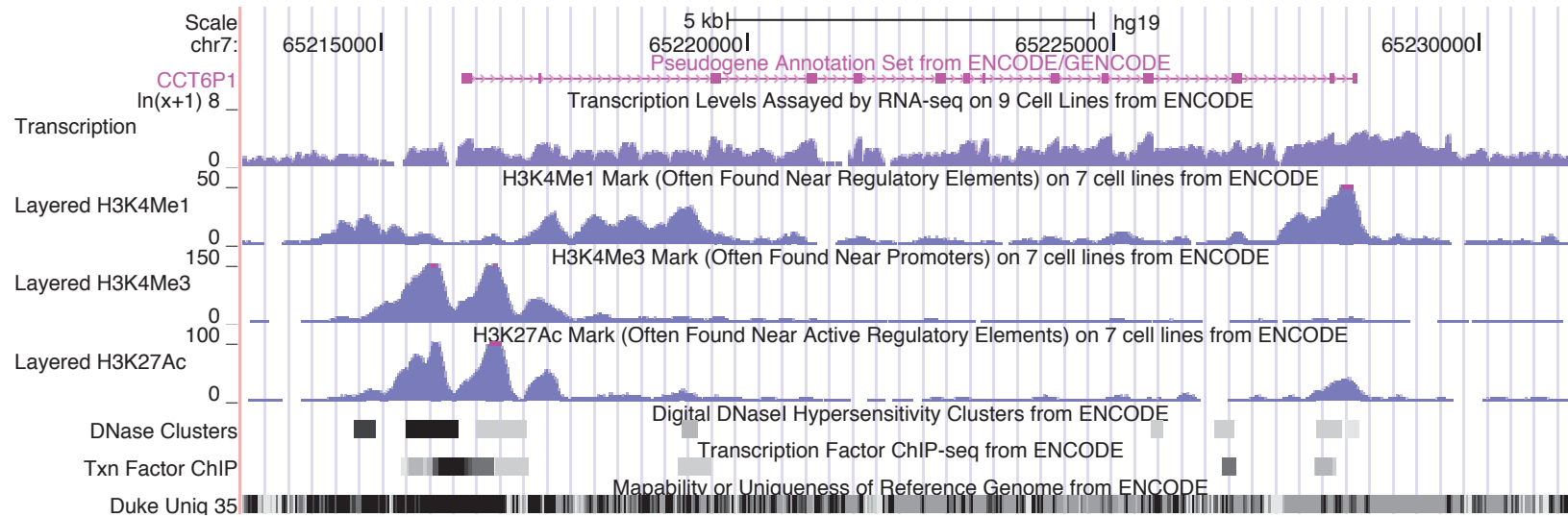
Pseudogene Activity



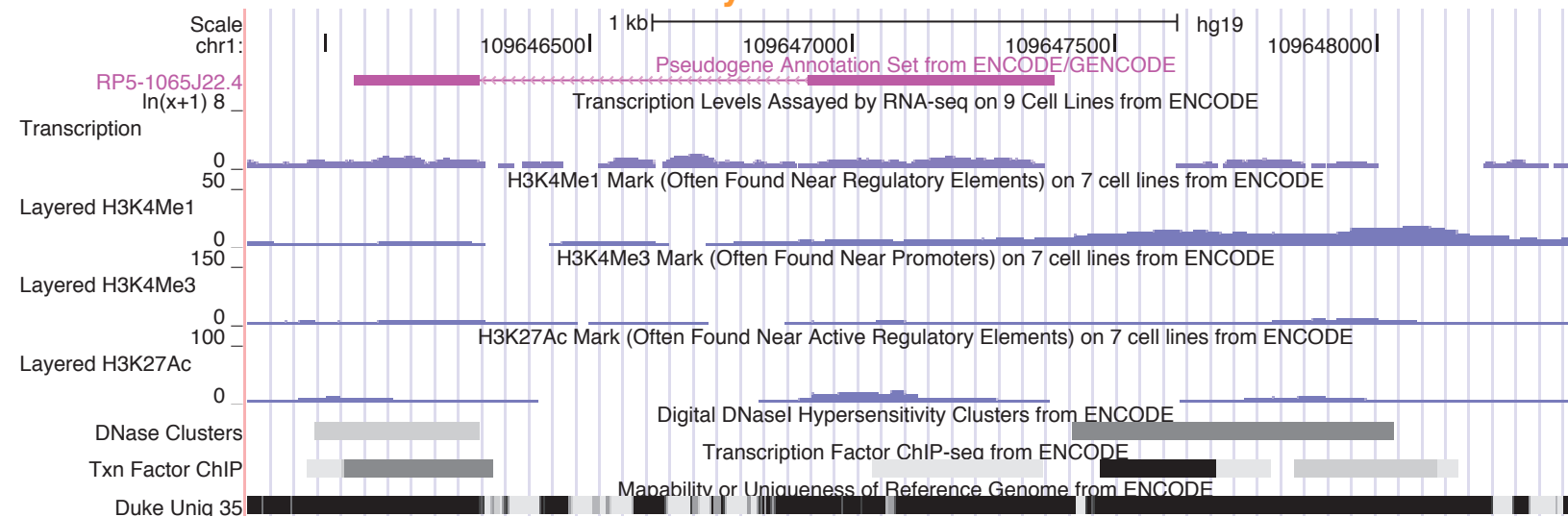
15% of pseudogenes are **transcribed** in each organism

Partial Pseudogene Activity

Transcribed with Additional Activity



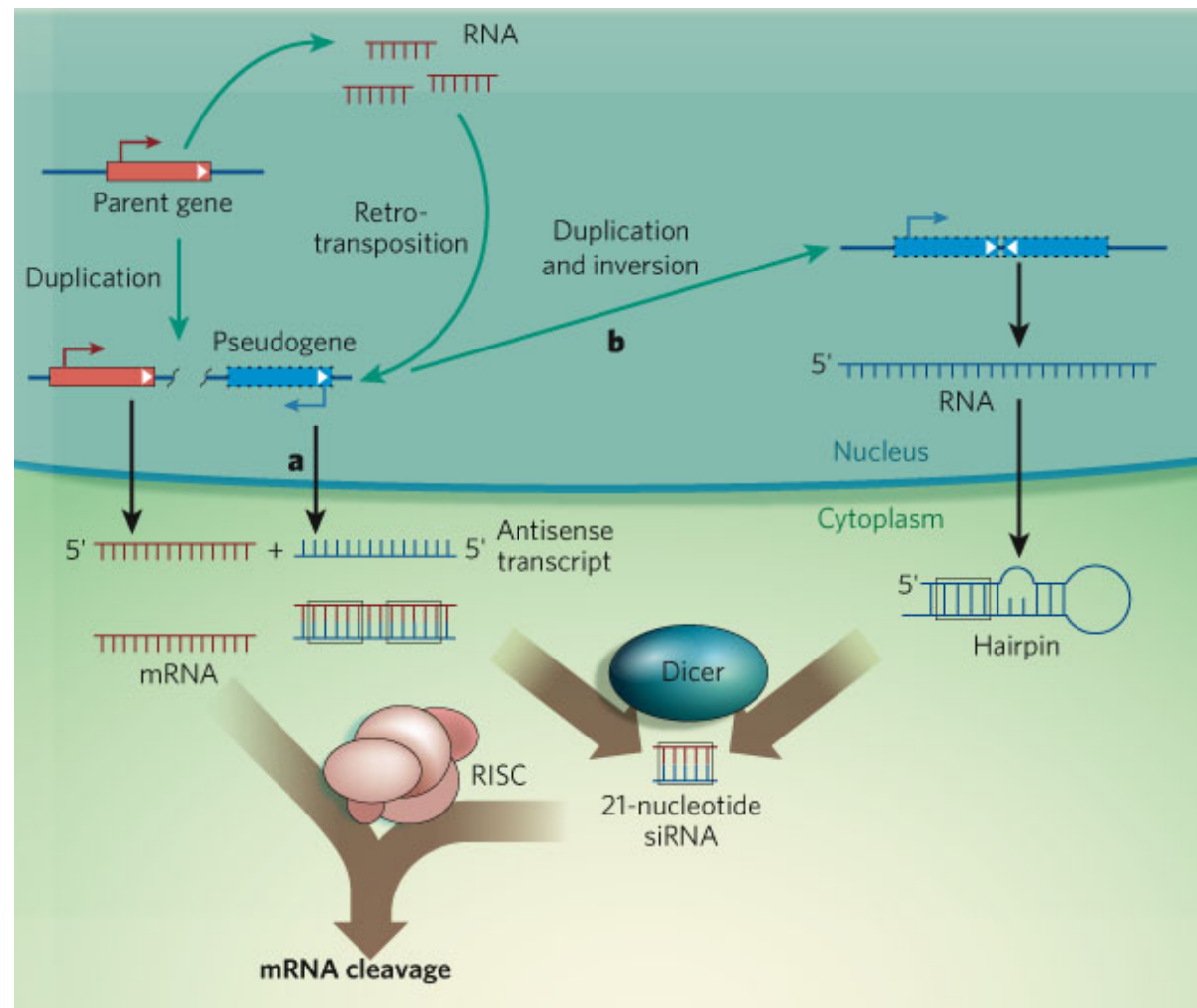
Partially Active



Examples & speculation on the function of pseudogene ncRNAs:

Regulating their parents

- via acting as **endo-siRNAs** [ex. in fly & mouse, '08 refs.]
- via acting as **miRNA decoys** [PTEN, BRAF]
- via **inhibiting degradation** of parent's mRNA [makorin]



[Sasidharan & Gerstein, Nature ('08)]

Alternatively,
just last gasps
of a dying gene

Czech *et al.* Nature 453: 798 ('08).
 Ghildiyal *et al.* Science 320: 1077 ('08).
 Kawamura *et al.* Nature 453: 793 ('08).
 Okamura *et al.* Nature 453: 803 ('08).
 Tam *et al.* Nature 453: 534 ('08).
 Watanabe *et al.* Nature 453: 539 ('08).

Poliseno *et al.* Nature 465:1033 ('10)
 Karreth *et al.* Cell ('15).

Prevalence of noncoding transcription in the human genome, in relation to lncRNAs: From noisy TARs to regulatory pseudogenes

- **The Discovery of Pervasive Transcription**

- Segmenting a noisy signal from Tiling Arrays into TARs
- The problem: were the results accurate?
- Cross-hyb.

- **Pervasive Transcription, Take 2**

- The advent of Nextgen seq.
- Evidence integration: lncRNA
- Now consistent cross-organism results: ~1/3 of the un-annotated genome is transcribed

- **Drilling into one type of pervasive transcription: Transcribed Pseudogenes**

- Tricky definition & great prevalence (~14K)
- Many transcribed: ~15%
 - Validation (RT-pcr)
 - Lots of other supporting evidence (other func. genomics expt. + selection)
- Ideas on a regulatory biological role (endo-siRNAs, sponges, &c)



Acknowledgements

EncodeProject.org/comparative

Mark Gerstein, Koon-Kiu Yan, Daifeng Wang,

Chao Cheng, James B. Brown, Carrie A. Davis, LaDeana Hillier, Cristina Sisu,

Jingyi Jessica Li, Baikang Pei, Arif O. Harmanci, Michael O. Duff, Sarah Djebali, Roger P. Alexander, Burak H. Alver, Raymond K. Auerbach, Kimberly Bell, Peter J. Bickel, Max E. Boeck, Nathan P. Boley, Benjamin W. Booth, Lucy Cherbas, Peter Cherbas, Chao Di, Alex Dobin, Jorg Drenkow, Brent Ewing, Gang Fang, Megan Fastuca, Elise A. Feingold, Adam Frankish, Guanjun Gao, Peter J. Good, Phil Green, Roderic Guigó, Ann Hammonds, Jen Harrow, Roger A. Hoskins, Cédric Howald, Long Hu, Haiyan Huang, Tim J. P. Hubbard, Chau Huynh, Sonali Jha, Dionna Kasper, Masaomi Kato, Thomas C. Kaufman, Rob Kitchen, Erik Ladewig, Julien Lagarde, Eric Lai, Jing Leng,

Zhi Lu, Michael MacCoss, Gemma May, Rebecca McWhirter, Gennifer Merrihew, David M. Miller, Ali Mortazavi, Rabi Murad, Brian Oliver, Sara Olson, Peter Park, Michael J. Pazin, Norbert Perrimon, Dmitri Pervouchine, Valerie Reinke, Alexandre Reymond, Garrett Robinson, Anastasia Samsonova, Gary I. Saunders, Felix Schlesinger, Anurag Sethi, Frank J. Slack, William C. Spencer, Marcus H. Stoiber, Pnina Strasbourger, Andrea Tanzer, Owen A. Thompson, Kenneth H. Wan, Guilin Wang, Huaien Wang, Kathie L. Watkins, Jiayu Wen, Kejia Wen, Chenghai Xue, Li Yang, Kevin Yip, Chris Zaleski, Yan Zhang, Henry Zheng, **Steven E. Brenner, Brenton R. Graveley,**

Susan E. Celniker, Thomas R Gingeras,

Robert Waterston

Pseudogene.org

C Sisu, B Pei, J Leng, A Frankish, Y Zhang, S Balasubramanian, R Harte, D Wang, M Rutenberg-Schoenberg, W Clark, M Diekhans, T Hubbard, **J Harrow, Mark Gerstein**

