

Gencode Mar '10 Meeting: Pseudogene **Project Update**

Mark Gerstein

60

Overall Flow:

<u>Pipeline Runs, Coherent Sets,</u> <u>Annotation, Transfer to Sanger</u>

- Overall Approach
 - Overall Pipeline runs at Yale and UCSC, yielding raw pseudogenes
 - 2. Extraction of coherent subsets for further analysis and annotation
 - 3. Passing to Sanger for detailed manual analysis and curation
 - 4. Incorporation into final GENCODE annotation
 - 5. Pipeline modification

- Chronology of Sets
 - 1. Encode Pilot 1%
 - 2. Ribosomal Protein pseudogenes
 - 3. Unitary pseudogenes (Hard)
 - 4. Glycolytic Pseudogenes
 - 5. Polymorphic Pseudogenes
 - 6. Pseudogenes Associated with SDs

Specific Pseudogene Assignments: Glycolytic Pseudogenes (completed)



<u>Number of</u> pseudogenes for each glycolytic enzyme

[Liu et al. BMC Genomics ('09)]

Large numbers of processed GAPDH pseudogenes in mammals comprise one of the biggest families but numbers not obviously correlated with mRNA abundance.

Processed/Duplicated





<u>Number of</u> pseudogenes for each glycolytic enzyme

[Liu et al. BMC Genomics ('09)]

Large numbers of processed GAPDH pseudogenes in mammals comprise one of the biggest families but numbers not obviously correlated with mRNA abundance.

Processed/Duplicated





Distribution of human GAPDH pseudogenes



6



7 - Lectures.GersteinLab.org (o) 00



Specific Pseudogene Assignments: Unitary Pseudogenes (completed)



Pseudogenes

Pseudogenes: nongenic DNA segments with high sequence similarity to functional genes



 Unitary pseudogenes: unprocessed pseudogenes with no functional counterparts





zdz © mmix

Identification pipeline



Relativity of unitary pseudogenes

{ Unitary pseudogene







zdz © mmix

[Zhang et al. GenomeBiology (in press, '10)] 12

Unitary Pseudogene Families





Dating the pseudogenization events

Specific Pseudogene Assignments: Polymophic Pseudogenes (in process)



11 Polymorphic Pseudogenes

Gene	CDS-disruptive mutation		dbSNP ID ³	HanMan SNP ID		
	Change ¹	Location ²				
Nonsense mutation						
FBXL21	taT (Y) \rightarrow taA	chr5+:135,300,350	rs17169429	rs17169429 (+27)		
			(+27)			
FCGR2C	$Cag\ (Q) \to Tag$	chr1+:159,826,011	rs3933769 (–60)	rs3933769 (–60)		
GPR33	Cga (R) \rightarrow Tga	chr14-:31,022,505	rs17097921	rs17097921		
SEC22B	Caa (Q) \rightarrow Taa	chr1+:143,815,304	rs2794062	rs16826061 (+95)		
SERPINB11	Gaa (E) \rightarrow Taa	chr18+:59,530,818	rs4940595	rs4940595		
TAAR9	Aaa (K) → Taa	chr6+:132,901,302	rs2842899	rs2842899		
Frame-shift mutation						
CASP12	ΔCA	chr11-:104,268,39	rs497116 (–67)	rs497116 (-67)		
		4-5				
KRTAP7-1	ΔT	chr21-:31123841	rs35359062	rs9982775 (–20)		
PSAPL1	∇A	chr4–:7,487,457	rs58463471	rs4484302 (+441)		
TMEM158	∇A	chr3-:45,242,396	rs11402022	rs33751 (+725)		
TPSB2	ΔC	chr16-:1,219,240	rs2234647	rs2745145		
				(–1771)		

Table 2. Human polymorphic pseudogenes

Polymorphic pseudogenes (3 with allele frequency data)

CDS-disrupted gene	GPR33	SERPINBII	TAAR9
Disruptive mutation ¹	$Cga(R) \rightarrow Tga$	Gaa (E) \rightarrow Taa	Aaa (K) \rightarrow Taa
dbSNP ID	r\$17097921	rs4940595	rs2842899
Genomic location	chr14—:31,022,505	chr18+:59,530,818	chr6+:132,901,302
Disrupted codon position ²	140 (332)	89 (388)	61 (344)
Reference allele in human	Т	Т	Т
Reference allele in other primates ³	С	Т	Т
Allele frequency ⁴	CHJY	CHTJLKADGMY	CHTJLKADGMY
Test statistic for HWE in the meta-population ⁵	0.285 (P = 0.867)	8.659 (<i>P</i> = 0.013)	0.071 (P = 0.965)



3 SNPs not found to be under recent positive selection....

[Zhang et al. GenomeBiology (in press, '10)] 17

 $F_{\rm st}$ hierarchical clustering for rs4940595 in SERPINB11



••••but population structure at rs4940595—the difference in the allelic frequencies in different populations—could be result of different selective regimes that the same allele at rs4940595 is subjected to in different population subdivisions.

18

Specific Pseudogene Assignments: SD-associated Pseudogenes (in process)



Segmental duplications (SDs)

- Regions of the genome with ≥ 90% sequence identity and ≥ 1kb in length
- Based on neutral divergence correspond to last ~40 million years of human evolution
- Comprise ~5-6% of the human genome
- Enriched with genes (~18%) and pseudogenes (duplicated ~45%, processed ~22%)
 - Can the study of ψgenes in SDs provide information not obvious from individual dataset ?



Bailey et al, Science, 2002

Nucleotide substitutions in ψ genes and SDs containing them



K_{2m} : Nucleotide substitutions per site computed using Kimura's two parameter model

Most ψ genes show the same number of substitutions as larger SD region containing them

- Duplication accompanied by disablement
- Followed by neutral rate of evolution

Z Zhang E Khurana Y J Liu YK Lam S Balasubramanian G Fang N Carriero **R** Robilotto P Cayting **M** Wilson A Frankish **M** Diekhans **R** Harte T Hubbard **J** Harrow

Acknowledgements

Pseudogene.org