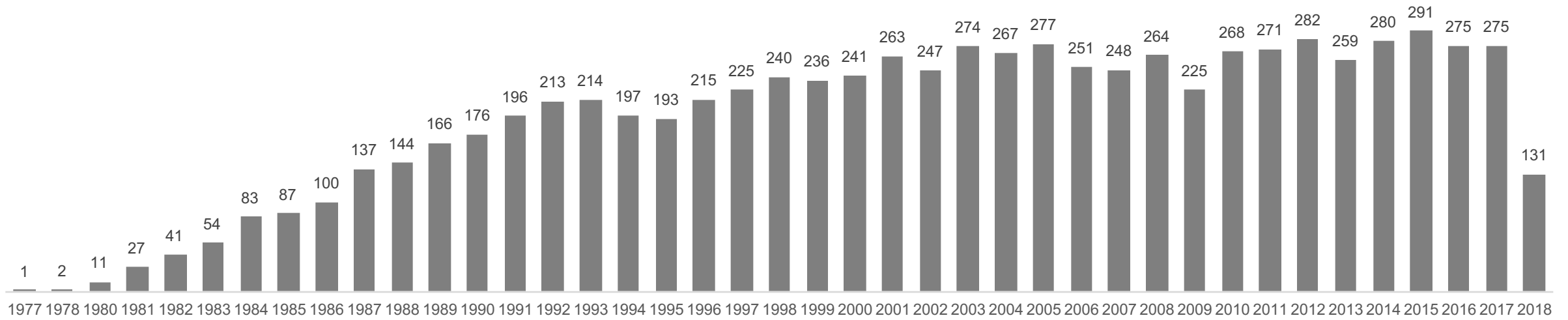# Updates on Pseudogene Analysis in Human and Mouse

**Mark Gerstein & Paul Muir**

**GENCODE meeting**

**21st June 2018**

Bar chart of publication counts by year:

| Year | Count |
|---|---|
| 1977 | 1 |
| 1978 | 2 |
| 1980 | 11 |
| 1981 | 27 |
| 1982 | 41 |
| 1983 | 54 |
| 1984 | 83 |
| 1985 | 87 |
| 1986 | 100 |
| 1987 | 137 |
| 1988 | 144 |
| 1989 | 166 |
| 1990 | 176 |
| 1991 | 196 |
| 1992 | 213 |
| 1993 | 214 |
| 1994 | 197 |
| 1995 | 193 |
| 1996 | 215 |
| 1997 | 225 |
| 1998 | 240 |
| 1999 | 236 |
| 2000 | 241 |
| 2001 | 263 |
| 2002 | 247 |
| 2003 | 274 |
| 2004 | 267 |
| 2005 | 277 |
| 2006 | 251 |
| 2007 | 248 |
| 2008 | 264 |
| 2009 | 225 |
| 2010 | 268 |
| 2011 | 271 |
| 2012 | 282 |
| 2013 | 259 |
| 2014 | 280 |
| 2015 | 291 |
| 2016 | 275 |
| 2017 | 275 |
| 2018 | 131 |

**1977**

**A Pseudogene Structure in 5S DNA of Xenopus laevis**
C. Jacq, J. R. Millier and G.G. Browniae

**Pseudogene** "has homologous structure, [is] nearly as long as, and almost an exact repeat of, the gene itself"

**1980s**
Pseudogenes' **structure** & **formation mechanisms**

**1990s**
Pseudogenes are **non functional, evolutionary fossils**

**2000-present**

**Systematic annotation** and analysis of **pseudogene** complements in genomes of **human** and **model organisms**

# The Gerstein lab has a long history in pseudogene annotation and analysis

## Pseudofam: the pseudogene families database

Hugo
Kei-

BIOINFORMATICS    ORIGINAL PAPER    *Vol. 22 no. 12 2006, p*
doi:10.1093/bio

*Genome analysis*

## PseudoPipe: an automated pseudogene identification pipeline

aul M. Harrison[5]

## Pseudogene.org: A comprehensive database and comparison platform for pseudogene annotation

John Karro[1,†], Yangpan Yan[2], Deyou Zheng[2], Zhaolei Zhang[3], Nicholas Carriero,[4] Paul Harrrison[5] and Mark Gerstein[2,‡]

Genome **Biology**

**RESEARCH**                                        Open Access

## The GENCODE pseudogene resource

Baikang Pei[1†], Cristina Sisu[1,2†], Adam Frankish[3], Cédric Howald[4], Lukas Habegger[1], Xinmeng Jasmine Mu[1], Rachel Harte[5], Suganthi Balasubramanian[1,2], Andrea Tanzer[6], Mark Diekhans[5], Alexandre Reymond[4], Tim J Hubbard[3], Jennifer Harrow[3] and Mark B Gerstein[1,2,7*]

## Comparative analysis of pseudogenes across three phyla

Zhan
http://

Cristina Sisu[a,b,1], Baikang Pei[a,1], Jing Leng[a,1], Adam Frankish[c,1], Yan Zhang[a,1], Suganthi Balasubramanian[b], Rachel Harte[d], Daifeng Wang[a], Michael Rutenberg-Schoenberg[a], Wyatt Clark[a], Mark Diekhans[d], Joel Rozowsky[b], Tim Hubbard[c], Jennifer Harrow[c], and Mark B. Gerstein[a,b,e,2]
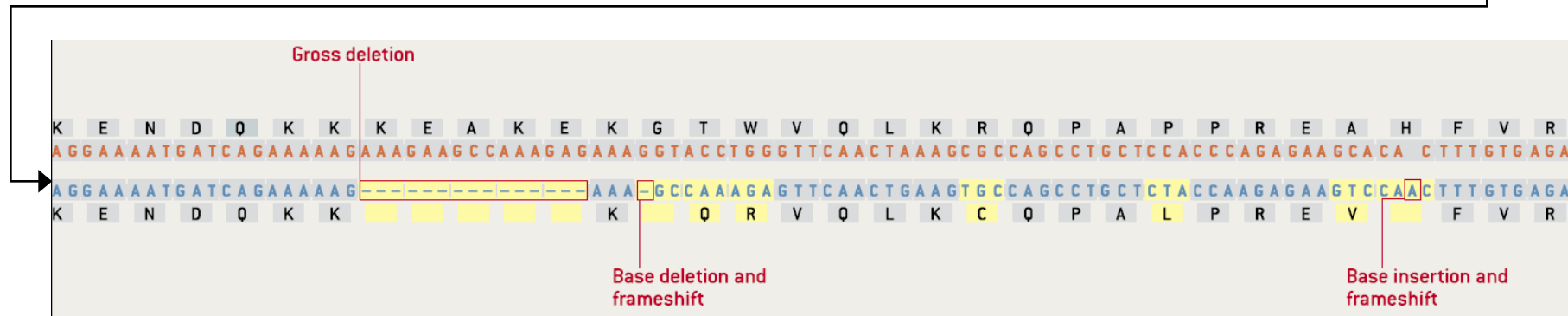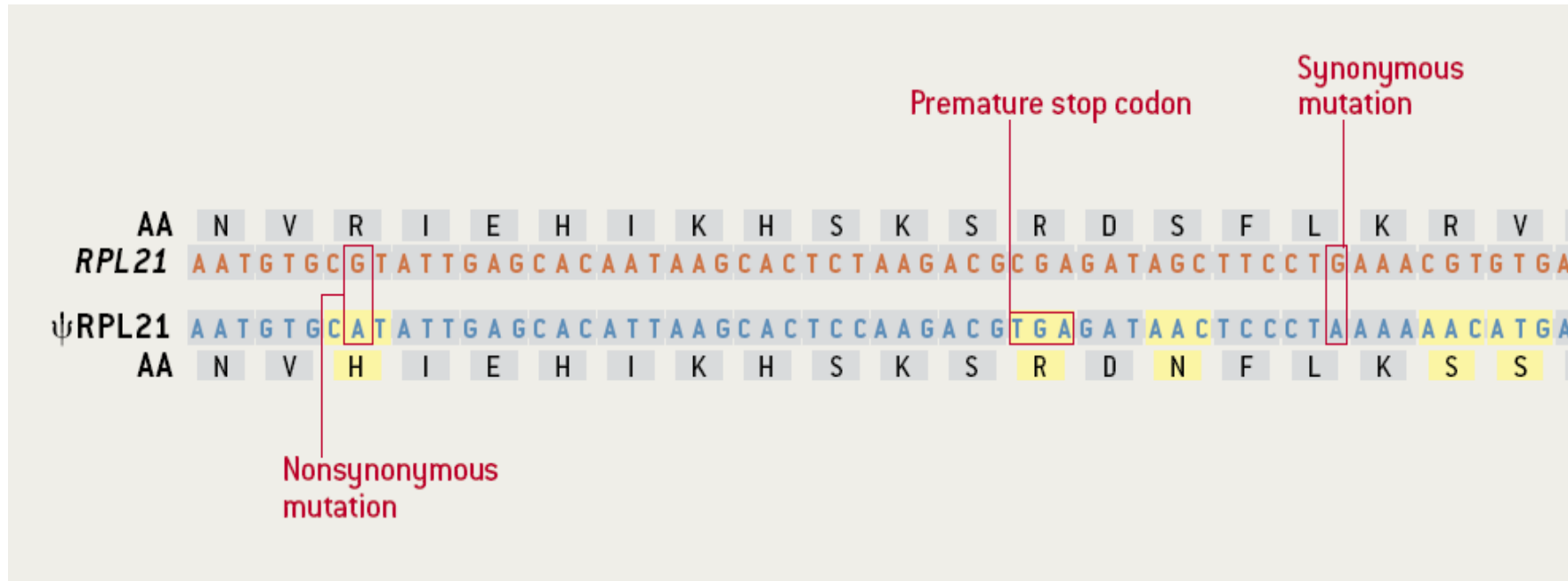
**RESEARCH**                                        Open Access
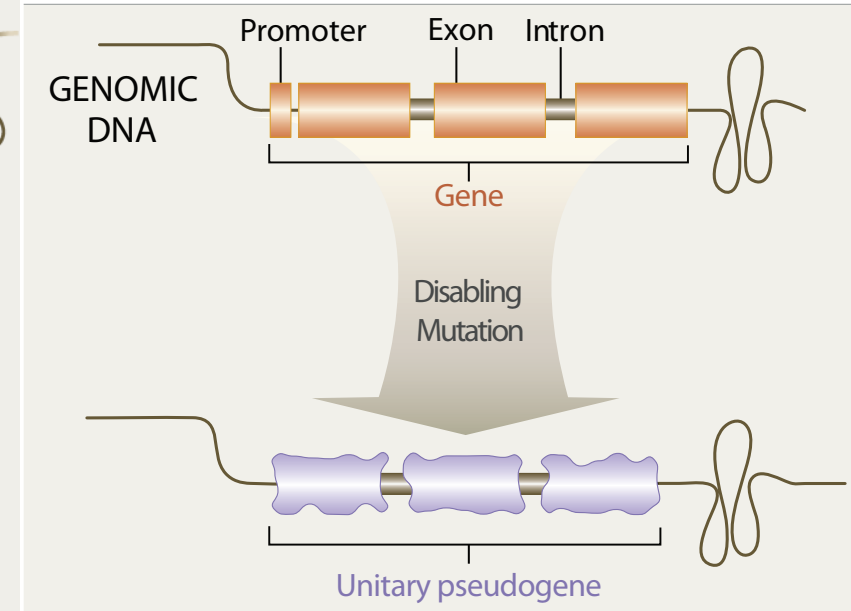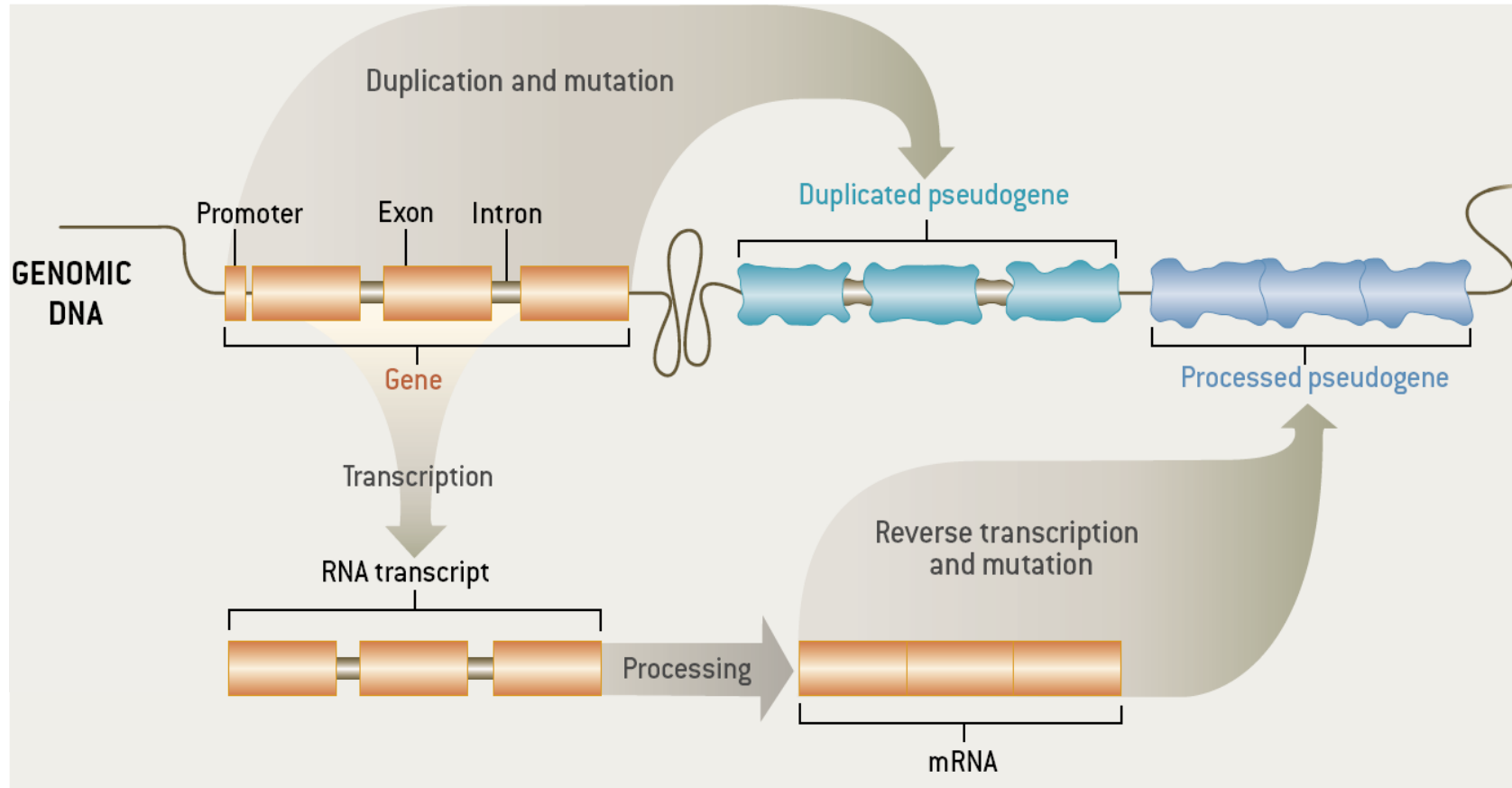
## Identification and analysis of unitary pseudogenes: historic and contemporary gene losses in humans and other primates

Zhengdong D Zhang[1], Adam Frankish[2], Toby Hunt[2], Jennifer Harrow[2], Mark Gerstein[1,3,4*]

The **Real Life** of **Pseudogenes**

By Mark Gerstein and Deyou Zheng

# A canonical pseudogene

Gerstein and Zheng SciAm 2006

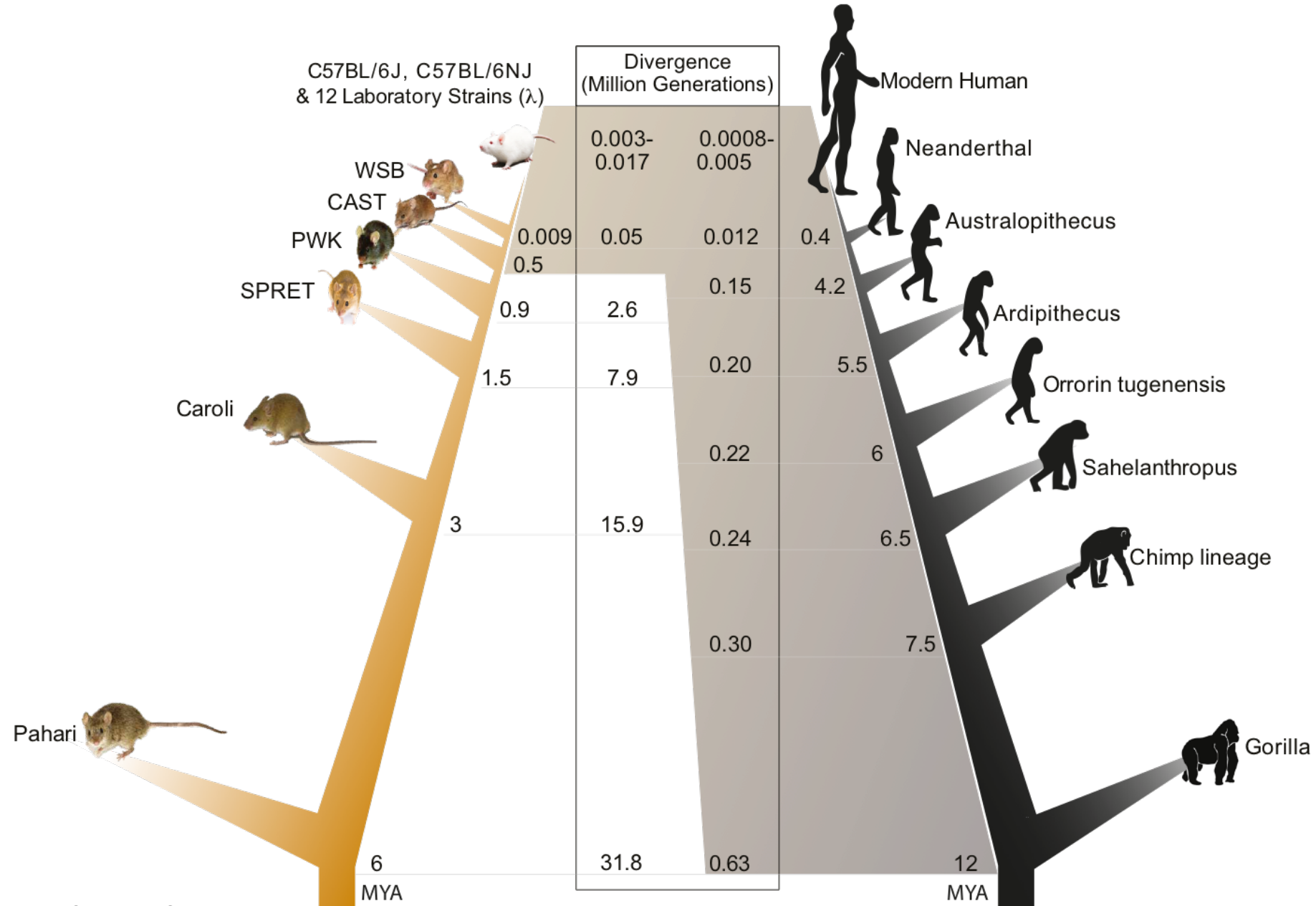# Review of pseudogene biogenesis



Gerstein and Zheng SciAm 2006

# Future work

- Finalize the annotation for mouse

- Improve the annotation in mouse strains

- Pseudogenes as personalized annotations

- Pseudogene annotation customized for human disease studies
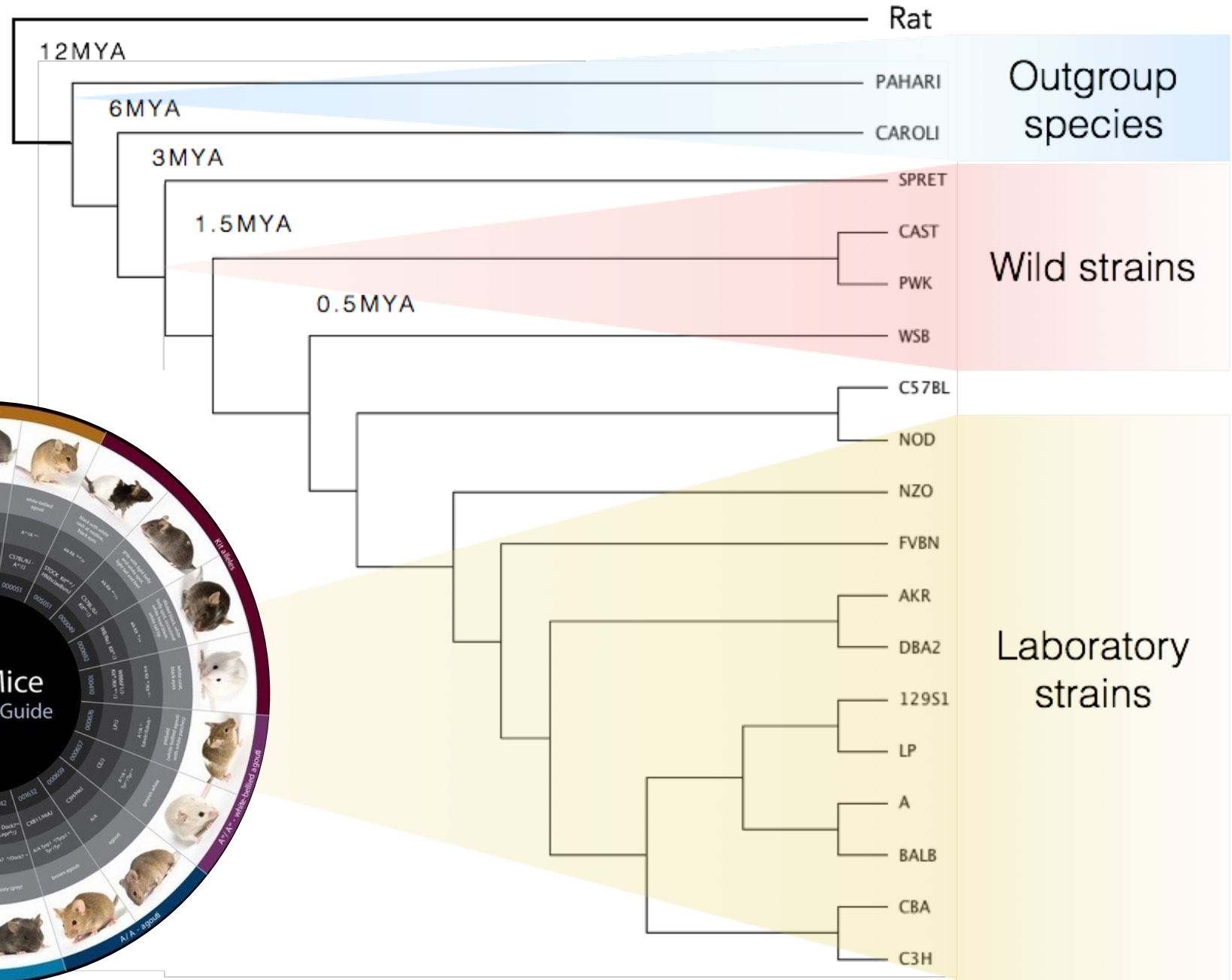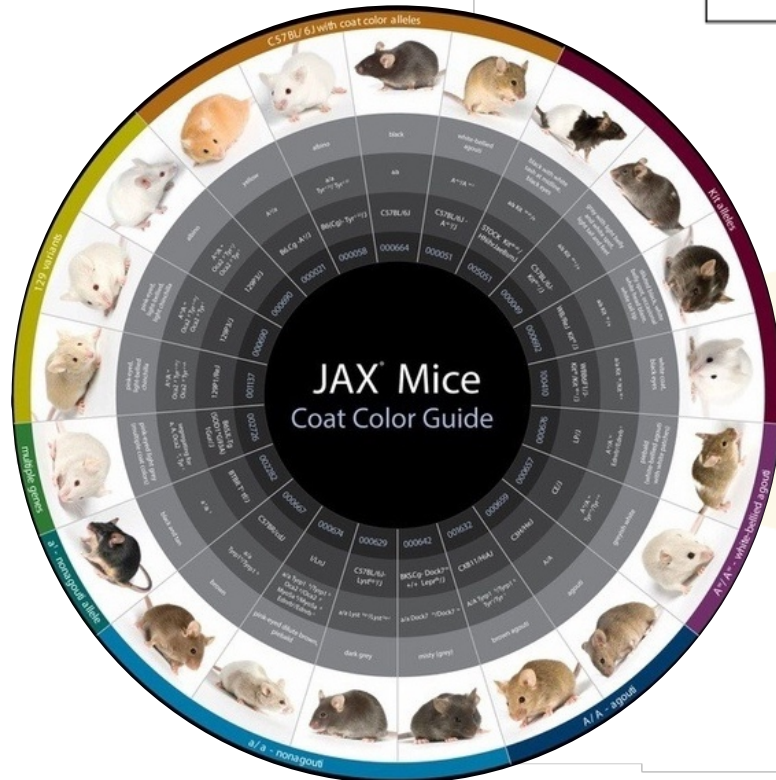
# Pseudogenes in the mouse lineage

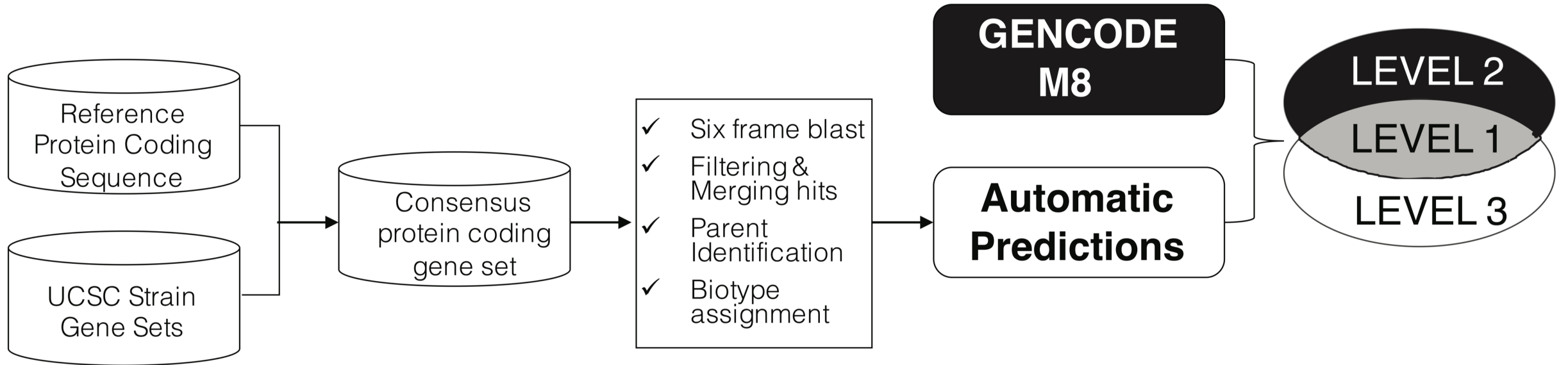# Comparisons across the mouse and primate lineages



C57BL/6J, C57BL/6NJ & 12 Laboratory Strains (λ)

| | Divergence (Million Generations) | |
|---|---|---|
| 0.003-0.017 | 0.0008-0.005 | |

Modern Human

Neanderthal

Australopithecus

Ardipithecus

Orrorin tugenensis

Sahelanthropus

Chimp lineage

Gorilla

WSB
CAST
PWK
SPRET
Caroli
Pahari

| Mouse (MYA) | λ | Divergence | Divergence | Primate |
|---|---|---|---|---|
| | 0.009 | 0.05 | 0.012 | 0.4 |
| 0.5 | | | 0.15 | 4.2 |
| 0.9 | | 2.6 | 0.20 | 5.5 |
| 1.5 | | 7.9 | 0.22 | 6 |
| 3 | | 15.9 | 0.24 | 6.5 |
| | | | 0.30 | 7.5 |
| 6 | | 31.8 | 0.63 | 12 |
| MYA | | | | MYA |

# Mouse strains



Rat

12MYA

6MYA — PAHARI

3MYA — CAROLI

Outgroup species

1.5MYA — SPRET

CAST

PWK

0.5MYA — WSB

Wild strains

C57BL

NOD

NZO

FVBN

AKR

DBA2

129S1

LP

A

BALB

CBA

C3H

Laboratory strains

JAX® Mice
Coat Color Guide

# Pseudogene annotation pipeline



Reference Protein Coding Sequence

UCSC Strain Gene Sets

Consensus protein coding gene set

- ✓ Six frame blast
- ✓ Filtering & Merging hits
- ✓ Parent Identification
- ✓ Biotype assignment

**GENCODE M8**

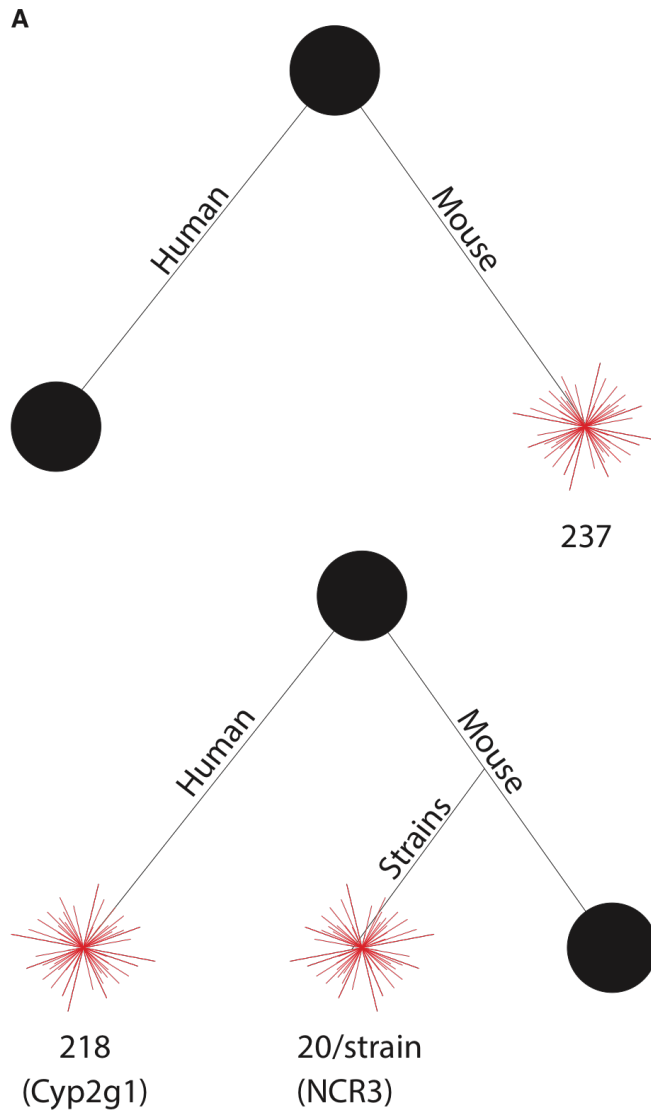**Automatic Predictions**

LEVEL 2

LEVEL 1

LEVEL 3

# Comparison of pseudogene annotations



Mouse strains have comparable pseudogene contents in both size and biotype distribution.

Fewer annotations in more divergent species due to use of reference mouse coding set.

Sisu C.*, Muir P.*  et al., *NComms.*, Submitted

# Unitary pseudogenes in human and mouse lineages



Foreign Strain/ Specie

Mouse Orthologs

Foreign Proteins without Mouse Orthologs

MAP to MOUSE
- ✓ Six frame blast
- ✓ Filtering & Merging hits
- ✓ Parent & Biotype Identification

**Automatic Predictions**

Filter known pseudogenes

**Unitary Pseudogenes**

Sisu C.*, Muir P.* et al., *NComms.*, Submitted

# Loss and gain of function in human and mouse lineages



A

B

**Cyp2g1 unitary pseudogene in human with respect to mouse**

Arg

**C G A**

↑

A T T G C C T G G A A A **T G A** A T G A A T A A G G C A G G

term

| | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|
| I | A | W | K | * | M | N | K | A | G | Human |
| I | A | W | K | R | M | N | K | A | G | Chimp |
| I | A | W | K | R | M | N | K | A | G | Orangutan |
| I | A | W | K | R | M | N | K | A | G | Rhesus |
| I | A | W | K | R | M | N | K | G | G | Marmoset |
| I | A | W | K | R | T | S | K | G | G | Mus Musculus |
| I | A | W | K | R | T | S | K | G | G | Lab Strains |
| I | A | W | K | R | T | S | K | G | G | Wild Strains (WSB & PWK) |
| I | A | W | K | W | T | S | K | G | G | Wild Strains (SPRET) |
| I | A | W | K | K | T | N | K | G | G | Guinea Pig |
| I | A | W | K | R | V | Q | K | P | G | Rabbit |

C

**NCR3 pseudogene unitary pseudogene with respect to Mus Caroli**

Trp

**T G G**

↑

T C C T G T G C T C T C **T G A** G T G T C T C A G C C C C C

term

| | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|
| S | C | L | A | * | V | S | Q | P | P | Mus Musculus |
| S | C | L | A | W | V | S | Q | P | P | Caroli |
| S | C | L | A | * | V | S | Q | P | P | Wild Strains |
| S | C | L | A | * | V | S | Q | P | P | Lab Strains |

237

218
(Cyp2g1)

20/strain
(NCR3)

Sisu C.*, Muir P.*  et al., *NComms.*, Submitted

# Pan-genome pseudogene annotation distribution



**A**

**B**

Out Group
6296
(254●)

2925★

Wild
8446
(10♦)

Lab+Ref
20420
(369✿)

C57BL/6NJ reference

Laboratory strains

Sisu C.*, Muir P.*  et al., *NComms.*, Submitted

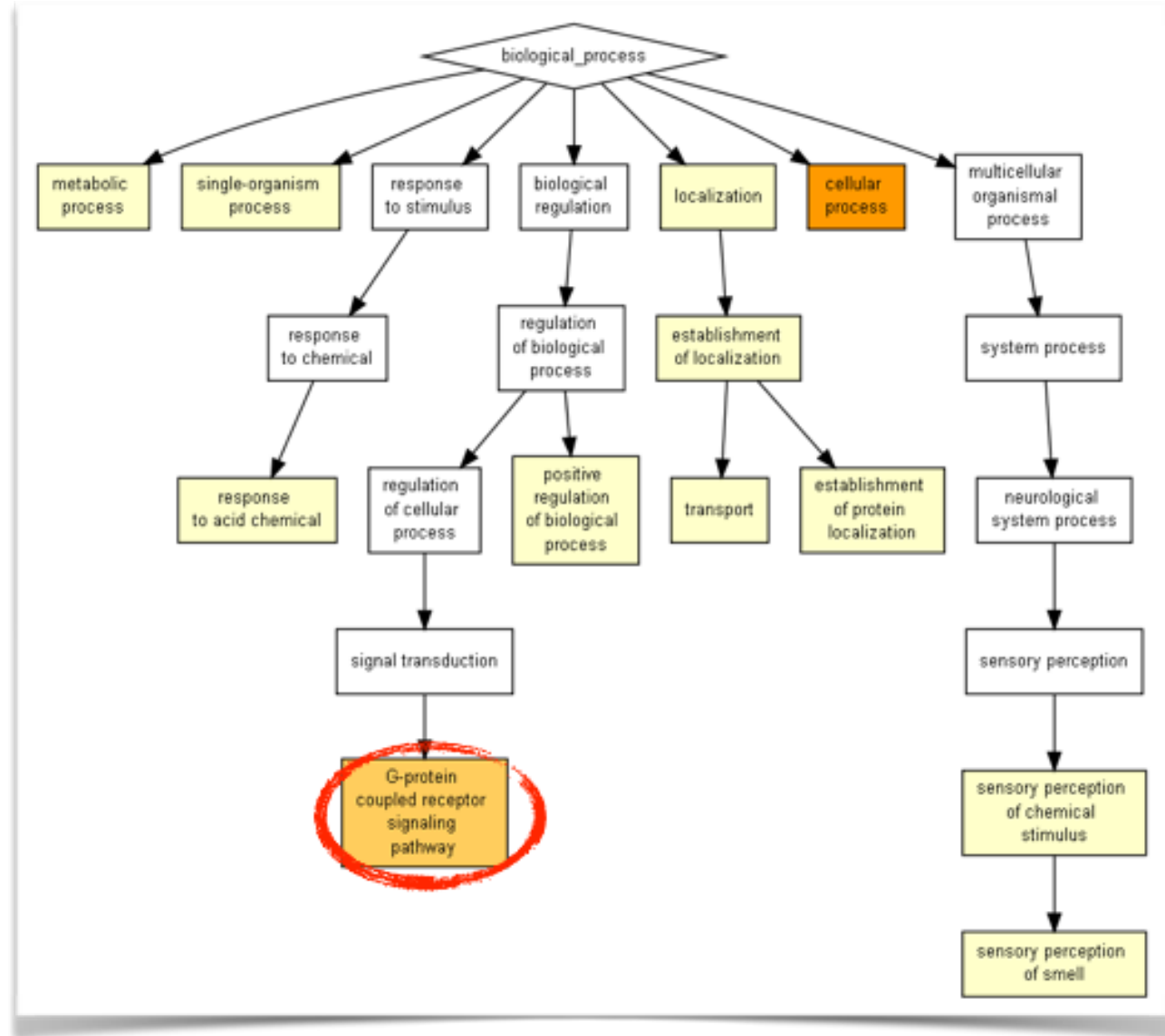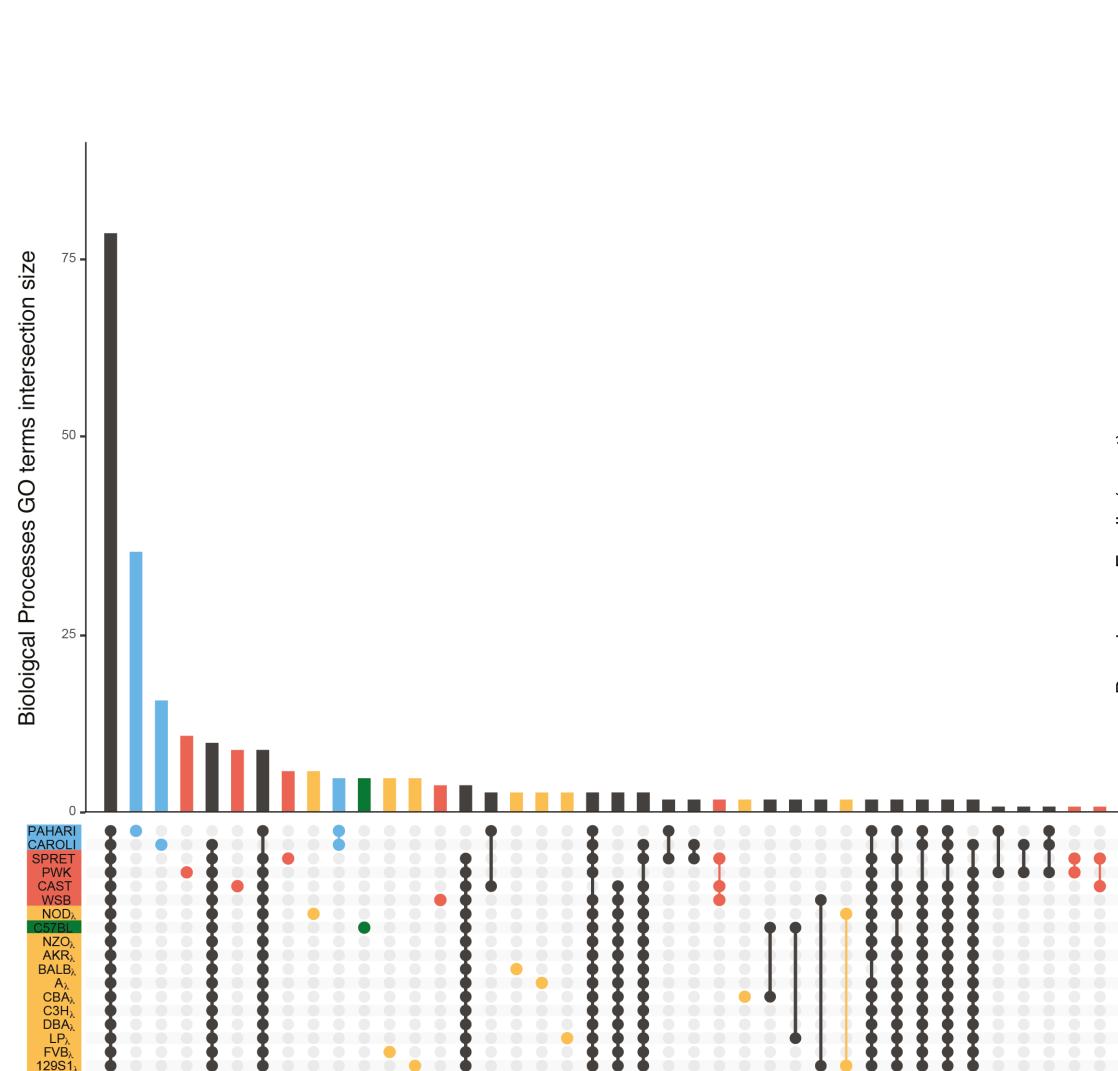# Historical patterns of transposon-mediated pseudogene genesis

# Gene Ontology enrichment analysis of parent genes

Highly abundant protein families show up in GO analysis of pseudogenes.

# Cross strain gene ontology and Pfam family analysis of pseudogenes



Sisu C.*, Muir P.* et al., *NComms.*, Submitted

# Gene Ontology term enrichment amongst pseudogenes (biological processes)

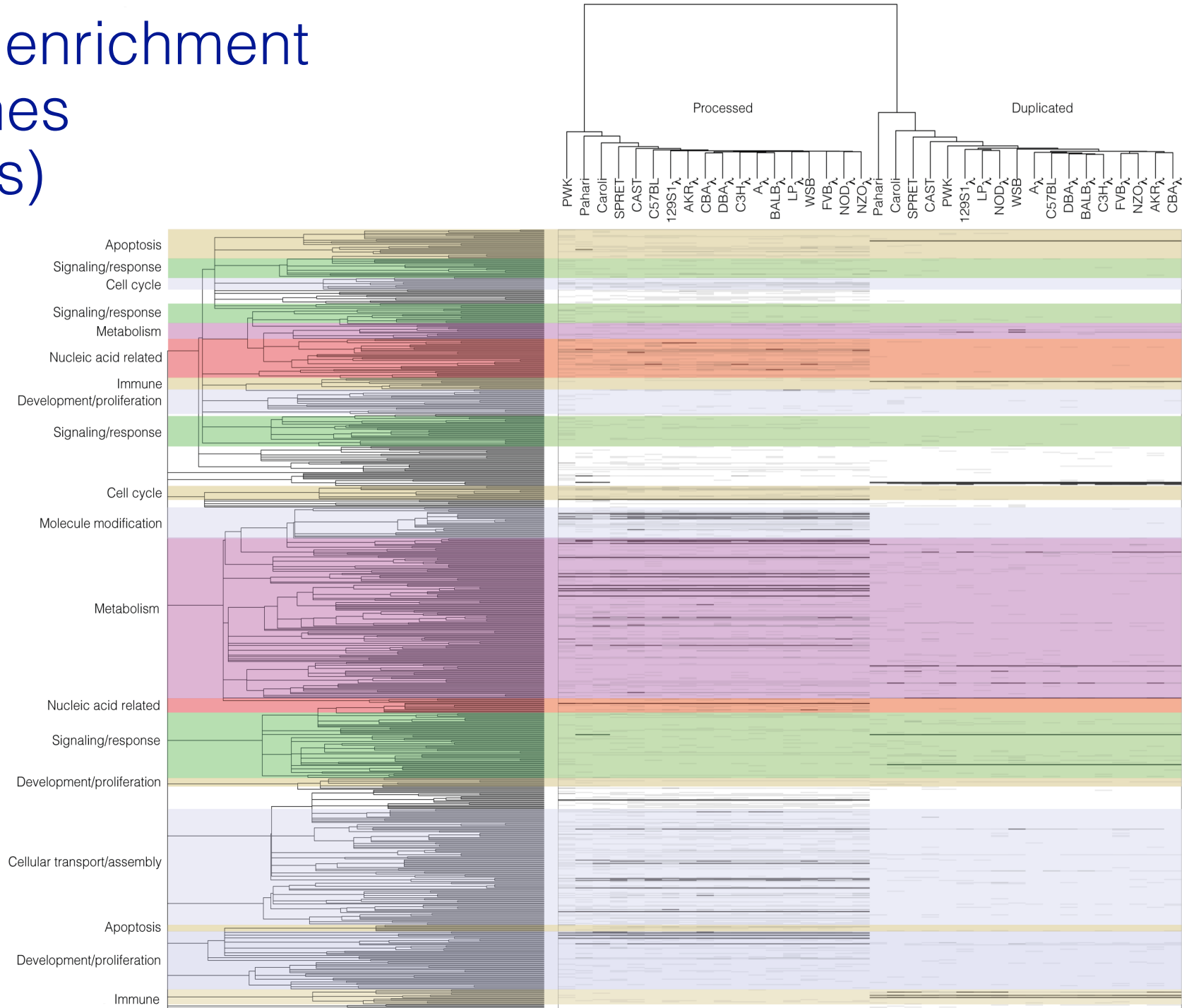Processed and duplicated pseudogenes enriched for different functions.

Processed pseudogenes enriched for:
- ribosomal functions
- cell cycle
- translation and RNA processing
- ubiquitination.

Duplicated pseudogenes enriched for:
- apoptosis
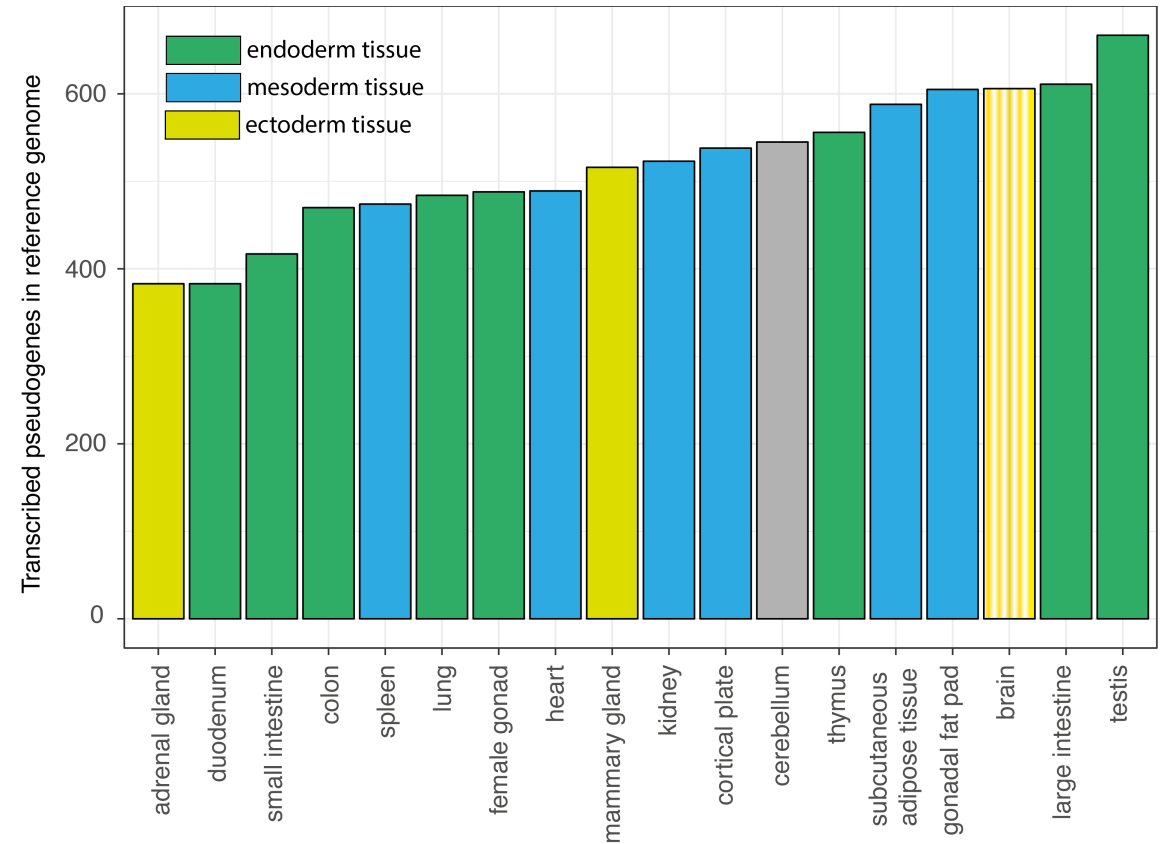- sensory and olfactory processes
- immune functions.

Sisu C.*, Muir P.* et al., *NComms.*, Submitted

# Transcriptional activity in reference genome



15% of mouse pseudogenes show evidence of residual transcription across multiple tissues

**Conserved pseudogenes with transcriptional activity** – this set of pseudogenes may need further review to ensure they are not misclassified functional elements.

**Strain-specific pseudogenes with transcriptional activity** – largely residual activity from pseudogenes with regulatory regions which have not decayed.



Sisu C.*, Muir P.*  et al., *NComms.*, Submitted

# Mouse Strains Pseudogenes

Welcome to the mouse strain pseudogene resource page!
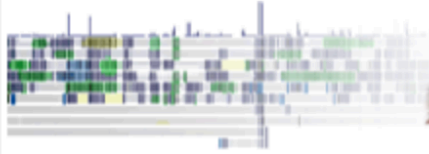
This database contains the latest annotation and characterization of pseudogenes in 18 related mouse strains. The pseudogene anotation was produced using a combination of automatic pipeline annotation using PseudoPipe and lift over of manually curated pseudogenes from the reference genome to each of the strains.

The resulting annotation set is characterised by 3 confidence levels. **Level 1** pseudogenes are identified by both PseudoPipe and manual lift over, **Level 2** pseudogenes are identified only by lifting over the manually curated set of the reference genome to the strain of interest; and **Level 3** pseudogenes are curated using just the automatic annotation pipeline.

- **Reference:** Sisu, Muir et al. **Pseudogenes in the mouse lineage: transcriptional activity and strain-specific history**. Submitted ⬀
- **Supplementary inforamtion:** All the supplementary information associated with the paper is available here.

## Annotation

### Reference Genome

The automatic pseudogene annotation for the mouse reference genome (Gencode vM12, Ensembl 87) is available here.

### Individual Strains

| | | | | | |
|---|---|---|---|---|---|
| 129S1/SvImJ | AKR/J | A/J | BALB/cJ | C3H/HeJ | C57BL/6NJ |
| Caroli/EiJ | CAST/EiJ | CBA/J | DBA/2J | FVB/NJ | LP/J |
| NOD/ShiLtJ | NZO/HILtJ | Pahari/EiJ | PWK/PhJ | SPRET/EiJ | WSB/EiJ |

### Pangenome Set

The current pangenome pseudogene set comprising 18 mouse strains is available in data-frame and list file format.

### Unitary Pseudogenes

- **Mouse:** Annotated unitary pseudogenes in the mouse reference genome with respect to human ⬀.
- **Human:** Annotated unitary pseudogenes in the human reference genome with respect to mouse ⬀.
- **Strains:** Annotated unitary pseudogenes in the mouse strains with repsect to the reference laboratory strain C57BL/6NJ ⬀.

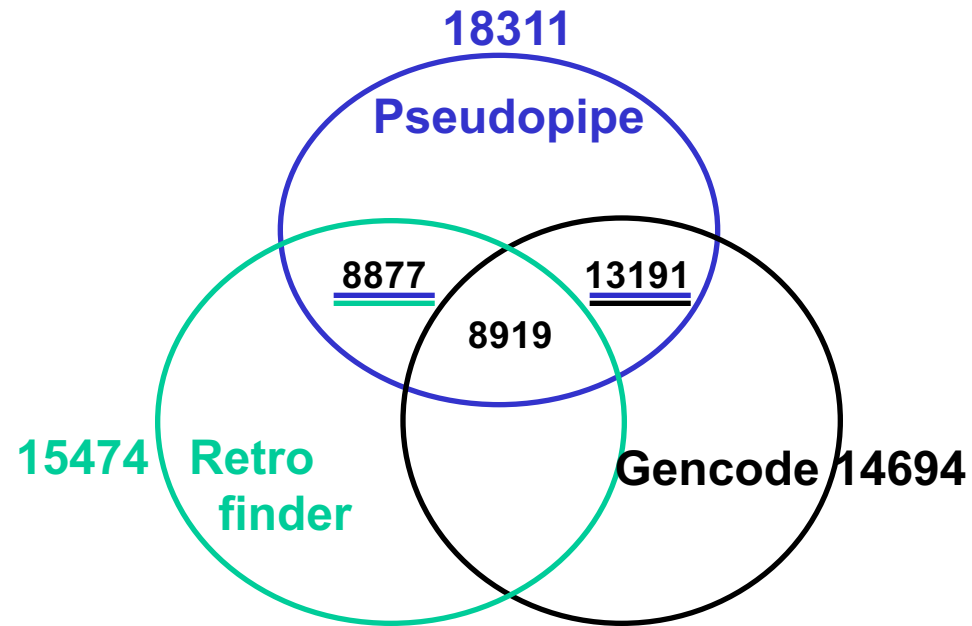Sisu C.*, Muir P.*  et al., *NComms.,* Submitted

# Summary

- The first draft of pseudogene annotation in 18 mouse strains and the reference genome

- On average 15-20% of  are strain specific and ~ 25% are ancestral, being conserved in all the strains.

- Top pseudogene families are matching closely the human counterparts.

- While human TE activity became silent after the retrotransposition burst, TE are still active in mouse strains.

- Similar to human, pseudogene prolific genes are not enriched in paralogs and vice versa.

- Pseudogene localization suggests multiple large scale genomic rearrangements between the out group - wild strains and the reference (lab strains) mouse genome.

- A significant proportion of  show signs of transcriptional activity.

## Acknowledgements

Cristina Sisu, Paul Muir, Adam Frankish, Ian Fiddes, Mark Diekhans, David Thybert, Duncan T. Odom, Paul Flicek, Thomas Keane, Tim Hubbard, Jennifer Harrow, Mark Gerstein
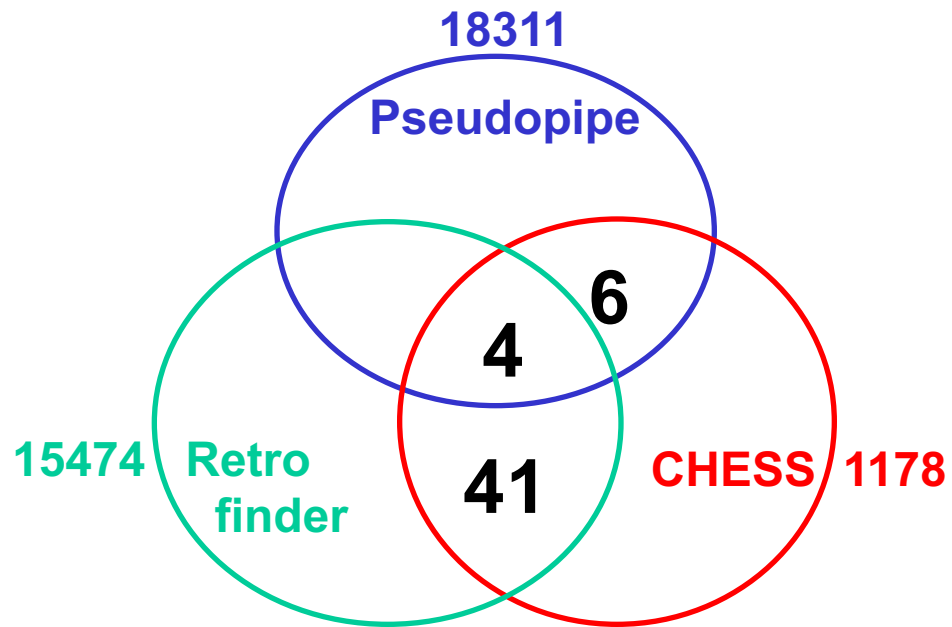
# CHESS – pseudogenes or coding genes?

18311

Pseudopipe

8877

13191

8919

15474 Retro finder

Gencode 14694

- CHESS genes exclude GENCODE pseudogenes

- Test the overlap between CHESS genes and PseudoPipe & Retrofinder genes

# 43 unique genes intersect pseudoexons with a 1bp minimum overlap



**18311**

Pseudopipe

15474  Retro finder

6

4

41

CHESS  1178

- 5/6 PseudoPipe pseudoexons have 100% sequence overlap
- 28/41 Retrofinder exons have 100% sequence overlap

- Others: 8-83% sequence overlap

# Is CHS.7402 a pseudogene?

• Similar to a protein from the crab eating macaque

>XP_005566708.1 PREDICTED: carbohydrate-binding protein AQN-1-like isoform X2 [Macaca fascicularis]

MRLSRAFAWSLLCSIATIVTAPFATAPSDCGGHYTDEYGRIFNYVGPKTECVWIIELNPGDIVVV
AIPELKGFVCGKEYVEVLDGPPGSESLGRICEAFSTFYHSSSNIITIKYSREPSHPPTFFEIYYF
VDAWSTH


**APSDCGGHYTDEYGRIFNYVGPKTECVWIIELNPGDIVVVAIPELK    KGFVCGKEYVEVLDGP**
**PGSESLGRICEAFSTFYHSSSNIITIKYSREPSHPPTFFEIYYFVDAWSTH (macaque)**
<span style="color:red">**APSDCGGHYTDEYGRIFNYAGPKTECVWIIELNPGEIVTVAIPDLK    RGFACGKEYVEVLDGP**</span>
<span style="color:red">**PGSESLDRICKAFSTFYYSSSNIITIKYSREPSHPPTFFEIYYFVDAWSTH  (human)**</span>

Misses the first 20 amino acids

Does not contain any indels or stop condon disablements

Potentially a duplicated pseudogene?