

1 INTRO

Pseudogenes have long been considered nonfunctional elements. However, recent studies indicate that pseudogenes can be transcribed, translated and can play key regulatory roles. In particular pseudogenes can regulate the expression of functional protein coding genes by serving as a source of siRNAs, antisense transcripts, microRNA binding sites, or competing mRNAs \cite{22726445,21080588,22990117}. The pseudogenization process is also closely linked to loss-of-function (LOF) events such as premature truncation of proteins, disruption of splicing and loss_of-functional or structural domains \cite{24026178,22344438,21205862}. Finally, the annotation of pseudogenes is important in the analysis of personal genomes, providing a means to avoid errors in genotyping assays and variant calling.

Pseudogenes are defined as disabled copies of functional genes. Depending on their formation mechanism, they can be referred to as unprocessed (originating through a gene duplication event) or processed (originating through a retrotransposition event). A functional gene may also become a pseudogene by acquiring a disabling mutation, if its function no longer confers a fitness advantage to the organism due to a change in the environment or genetic background. Such pseudogenes are called unitary pseudogenes. Pseudogenes provide valuable opportunities to study the dynamics and evolution of gene functions.

2 PRELIMINARY RESULTS & EXPERIENCE WITH PSEUDOGENE ANNOTATION

2.1 Pseudogene Annotation Pipelines

The Yale group has substantial experience in pseudogene annotation and analysis. In collaboration with the UCSC and Sanger group, we have developed a variety of methods to identify pseudogenes \cite{16574694,16925835,22951037}.

Pseudopipe, our in house automatic annotation pipeline, is fast and accurate \cite{22951037} (See Fig PG1). The pipeline takes as input all known protein sequences in the genome and using an homology search is able to identify disabled copies of functional paralogs (referred to as pseudogene parents). Based on their formation mechanism pseudogenes are classified into 3 different biotypes: processed, unprocessed and ambiguous. There is a good consensus overlap between the human pseudogene prediction set obtained with Pseudopipe and the set manually curated by the Gencode annotators \cite{22951037}. Even more, the Pseudopipe predictions fueled the manual curation of pseudogenes in GENCODE \cite{22951037}.

RCPedia, our newest pseudogene annotation pipeline focuses on the annotation of retrotransposed (processed) pseudogenes \cite{23457042} (see FIG PG1). This pipeline takes as input all known protein coding RNA transcripts and using sequence alignment is able to identify all possible retrocopies of functional genes. In the human genome there is an over 85% consensus between processed pseudogenes predicted by RCPedia and those annotated using Pseudopipe.

Formatted: Font: 18 pt

Deleted: o

Deleted: s

Deleted: .

... [1]

Deleted: comprehensive experiences

Deleted: a

Deleted: e

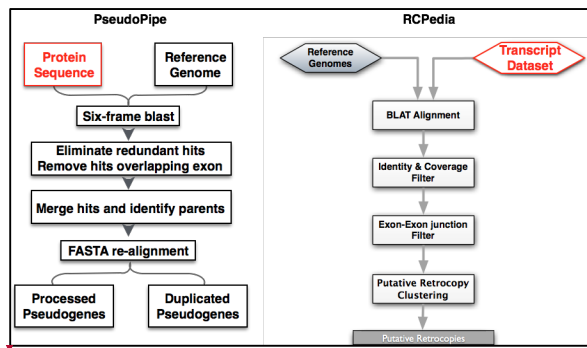
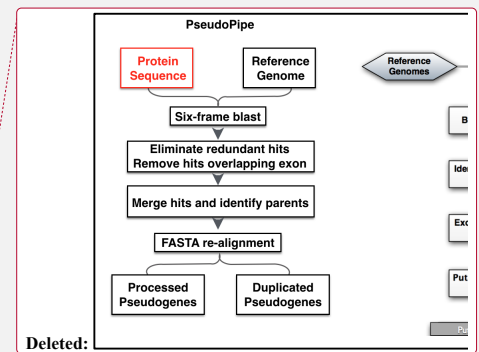


FIG PG1. Automatic pseudogene annotation pipelines.



Retrofinder is the UCSC retrogenes annotation pipeline. Retrogenes can be functional genes that have acquired a promoter, non-functional pseudogenes, or transcribed pseudogenes. Retrofinder finds retroposed messenger RNAs (mRNAs) in genomic DNA [\cite{18842134}](#). Candidate retrogenes overlapping by more than 50% with repeats identified by RepeatMasker [\cite{16093699,Smit}](#) and Tandem Repeat Finder [\cite{9862982}](#) are removed. Retrogenes are identified based on a score function using a weighted linear combination of features indicative of retrotransposition. These include: 1) Multiple contiguous exons with the parent gene introns removed; 2) Negatively scored introns as distinguished from repeat insertions (SVAs, LINEs, SINEs, Alus); 3) Lack of conserved splice sites; 4) Breaks in synteny with mouse and dog genomes (syntenic net alignments [\cite{14500911}](#)); and 5) Poly(A) tail insertion.

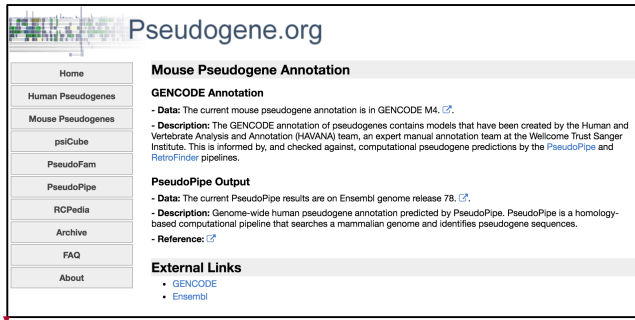
As a member of the GENCODE project, we used the pipelines to identify pseudogenes in human, mouse, worm, fly, and other model organisms [\cite{16925835,22951037,25157146}](#). The identified pseudogenes with related genomic and epigenomic data are available in our online databases [\cite{17099229,18957444,22951037,25157146}](#). Moreover, using data from the 1000 Genomes Project in addition to the pseudogene annotation resulting from our pipelines, we were able to study the impact of pseudogene in human population variation. To this end we evaluated 2,504 individuals across 26 human populations and we investigated the impact of coding and non-coding structural variants in the human genome [\cite{26432246}](#). We described processed pseudogenes as a novel class of gene copy number polymorphism that creates variability across populations. We were also able to associate their origin mechanism to cell division [\cite{24026178}](#).

2.2 Online Resources for Pseudogene Annotation and Analysis

Our experience in annotating and analyzing pseudogenes spans over a decade. Thus, we have built a number of tools to organize and analyze the available pseudogene data in a consistent and efficient manner.

We have built an online pseudogene repository, **pseudogene.org** \cite{17099229} (see Fig PG2), that provides information regarding annotation and functional characterization of pseudogenes. Currently pseudogene.org hosts the human (**psiDR** \cite{22951037}), and mouse pseudogene resources. It also provides a comparative pseudogene resource, **psiCUBE** \cite{25157146}, focused on cross species annotation and analysis of pseudogenes in a variety of model organisms. Both psiDR and psiCUBE also provide information regarding evolutionary and functional characterization of pseudogenes in the curated genomes.

Pseudogene.org also hosts **Pseudofam** \cite{18957444}, the pseudogene family database. Pseudofam resources focus on clustering pseudogenes into families based on their functional homolog protein family. Currently there are 10 eukaryotic genomes including human and mouse. Pseudofam also contains segmental duplication information associated with the human pseudogene dataset.



PG2. Pseudogene.org interface linking the available pseudogene tools and resources

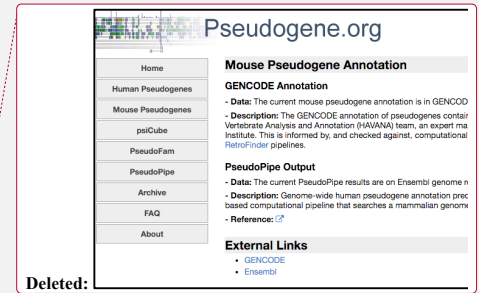
In order to record the structural and functional relationship between the pseudogenes within a family, we developed a **pseudogene ontology** \cite{20529940}. The pseudogene ontology is used in the generation of the GENCODE genomes annotation resource.

2.3 Current Results on Pseudogene as Part of GENCODE in Human and Model Organisms

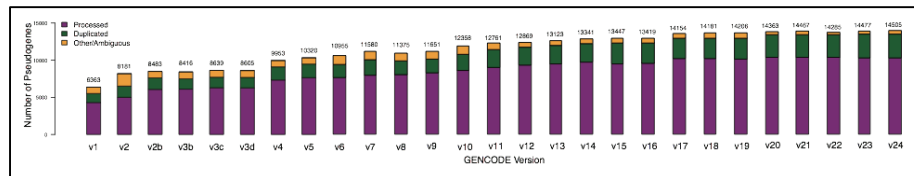
Our experience with annotating pseudogenes spans more than fifteen years. Over time we have annotated and reviewed pseudogenes in a variety of species ranging from prokaryotic organisms (archaea and bacteria) \cite{15345048,14583187}, to yeast \cite{11866506,12417195}, plants \cite{12083509}, worm \cite{11160906}, fly \cite{12034841,12560500}, and a wide range of vertebrates (e.g. zebrafish, mouse, rat, chimp, and human) \cite{19835609,12052146,12417195,12909341,18065488}. Our involvement in the GENCODE project started over a decade ago and ever since we have led and contributed to the identification and characterization of pseudogenes in human and model organisms (see Fig PG3).

Leveraging on the completed annotation of protein coding genes in human, worm and fly we were able to provide the complete and comprehensive set of pseudogenes in these organisms.

Deleted: es



In order to elucidate the role played by pseudogene in genome biology we integrated the annotation data with variation and functional genomics information.



PG3. GENCODE human pseudogene distribution in various releases.

In this respect we identified 14505 pseudogenes in human, 911 in worm, and 145 in fly [25157146,22951037]. A close comparison of the three genomes shows that pseudogenes complements do not follow the genome size or the number of protein coding genes in the genome, highlighting the species specific evolution of pseudogenes. This specificity is also reflected at pseudogene biotype level, where processed pseudogenes resulting from the burst of retrotransposition events that occurred at the dawn of primate lineage dominate the mammalian genomes, while the smaller fraction of duplicated pseudogenes hints at shared ancestry with more distant species [25157146].

We conducted a systematic analysis of human pseudogenes focusing on large groups of pseudogenes such as ribosomal pseudogenes [19123937,12417195,19835609], unitary [20210993] and polymorphic pseudogenes. The latter are a peculiar class of pseudogenes with a dual behavior – their sequence is disabled in the reference genome but in some individuals, it encodes a functional gene. We conducted a comprehensive review of polymorphic pseudogenes [21205862].

Despite the presence of disabling mutations such as premature stop codons, loss of promoters in the upstream sequence, numerous studies have shown that pseudogenes can be transcribed and even translated [15860774,16680195,15876366,17568002,16683022].

3 PRELIMINARY RESULTS AND EXPERIENCE WITH FUNCTIONAL CHARACTERIZATION OF PSEUDOGENES

We integrated ENCODE functional genomics data to obtain a comprehensive map of pseudogenes activity in human and other model organisms. We found that transcription signals have been observed for some pseudogenes and that the majority of pseudogenes (75% in human and 60% in worm and fly) have a large range in biochemical activity (e.g. presence of transcription factor or polymerase II binding sites in the upstream region, active chromatin, etc) (see Fig PG4). Moreover, we found 1,441, 143, and 23 transcribed pseudogenes in human, worm, and fly, respectively. We also identified 878 transcribed pseudogenes in mouse and 31 in zebrafish. These numbers represent a fairly uniform fraction (~15%) of the total pseudogene complement in each organism reflecting the similarity across phyla observed in their transcriptomes.

Deleted: in

Deleted: the presence of

Deleted: 3 Preliminary Results and Experience with Functional Characterization of Pseudogene Ac... [2]

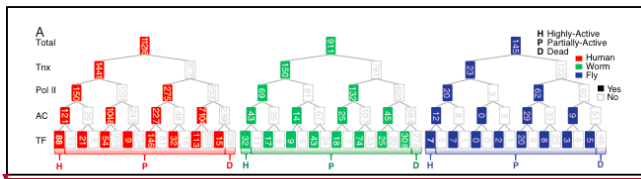


Fig PG4. Pseudogene activity. Distribution of pseudogenes as a function of various activity features: transcription (Tnx), active chromatin (AC), and presence of active Pol II and TF binding sites in the upstream region.

Among transcribed pseudogenes, ~13% in human and ~30% in worm and fly have a discordant transcription pattern with their parent genes over multiple samples. A large fraction of pseudogenes are associated with a few highly expressed gene families, e.g. the ribosomal proteins in human [25157146].

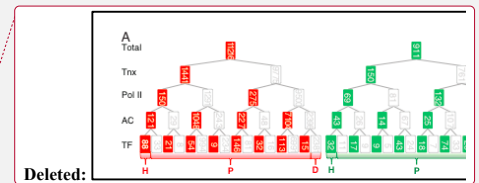
The parent genes of broadly expressed pseudogenes tend to be broadly expressed as well, but the reciprocal statement is not valid. Specifically, only 5.1%, 0.69%, and 4.6% of the total number of pseudogenes are broadly expressed in human, worm, and fly, respectively. However, in general, transcribed pseudogenes show higher tissue specificity than protein-coding genes [25157146].

We have also investigated pseudogene transcription by using the RNA-Seq data from the Illumina Human BodyMap data across 16 different tissues. Amongst all the transcribed pseudogenes identified, only a tiny proportion (~3%) are transcribed in all the 16 tissues, while the transcription of all the other pseudogenes show different degrees of tissue specificity. Furthermore, more than 50% of the transcribed pseudogenes are transcribed in one tissue only. While testis contained the largest number of transcribed pseudogenes, skeletal muscle contained the least [22951037].

4 PRELIMINARY RESULTS & EXPERIENCE WITH ANNOTATION AND ANALYSIS OF LOSS-OF-FUNCTION EVENTS

A loss-of-function (LOF) event is a genetic event that results in a severe disruption of the protein coding gene. Some LOFs can impact only one individual, resulting in the inactivation of an essential gene, which may lead to disease, while other LOFs can become fixed in the population as nonfunctional relics, through the pseudogenization process of the affected gene. The identification, analysis, and characterization of LOFs as either disease related or pseudogenization precursors is especially important in the era of personal genome annotation [21205862].

Moreover, the identification of pseudogenization/LOF events in mouse provides a very useful resource for understanding LOF in humans, by using mouse LOF phenotypes as proxy for human LOF events. To this end, the identification of mouse-specific unitary pseudogenes (regions that are functional in human and non functional in mouse) is important in highlighting



Deleted:

Deleted: ,

Deleted: 3.1 Preliminary Results and Experience in Annotating and Analysing Loss of Function Events

human genes [that](#) can (have functional paralogs in mouse) or cannot (are paralogs to unitary pseudogenes in mouse) be studied in mouse models. [\cite{12909341.14746985}](#).

Deleted: suitable
Deleted: [\[add ref\]](#)

Taking advantage of the rich 1000 Genomes data, we have developed a tool, called Variant Annotation Tool (VAT) [\cite{22743228}](#) (see Fig PG5), to systematically annotate and catalogue LOF events in the human genome. This pipeline enables rapid and efficient annotation of genomic variations (SNPs, indels and SVs) with respect to a reference genome and a gene annotation model. VAT can be used to identify pseudogenization events such as premature STOPs as well as polymorphic pseudogenes where a pseudogene in the reference genome becomes functional in another genome due to genetic variability at the stop codon.

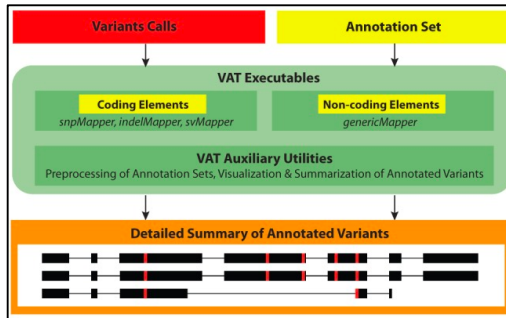


FIG PG5. Variant annotation tool (VAT) architecture.

We applied our tools to the 1000 Genomes [Phase 3](#) data and we were able to characterize putative LOF events from individuals [belonging to 26 different populations](#). [While earlier studies have suggested](#) that on average the human genome contains ~100 genuine LOF variants resulting in the total disablements of ~20 genes [\cite{22344438}](#), [we found this number to be higher](#). [On average the human genome contains 149–182 sites with protein truncating variants, ~11,000 sites with peptide-sequence-altering variants, and around 500,000 variant sites overlapping known regulatory regions \(untranslated regions, promoters, enhancers, etc.\) \cite{26432245}](#). Even more we were able to identify [24-30 sites per genome that are predicted severe disease-causing variants](#).

Deleted: 2951
Deleted: 185
Deleted: . Our results
Deleted: shown
Deleted: o
Deleted: }.
Deleted: 26 known and 21
Deleted: [\[Need to update to 1000G Phase3 but not incl. ALOFT\]](#)

In a similar manner we surveyed the impact of LOFs on personal genome annotation [\cite{21205862}](#). We found that LOFs variants that introduce premature STOPs resulting in a gene truncation in the reference genome can lead to an incorrect annotation of the gene. This highlights the importance of correct LOF identification for an accurate annotation. Finally we have studied the LOF events that results in a pseudogenization process. It is known that the loss_of-function in duplicated pseudogenes happens right after the gene duplication processes [\cite{20615899}](#). To this end we have developed a pipeline to identify unitary pseudogenes in human [\cite{20210993}](#) and we explored the functional constraints faced by different species and the timescale of functional gene loss [\cite{20210993}](#). All these results together with fully annotated sets of pseudogenes are deposited in our repository at [pseudogene.org](#).

5 PRELIMINARY RESULTS IN BUILDING AND USING PERSONAL GENOMES AND ANNOTATIONS

The human reference genome is a haploid sequence derived as a composite from multiple individuals. Current genome annotations are based on this reference and as such do not provide an accurate representation for the large genomic diversity of the human population. We have developed a computational tool, *vcf2diploid*^[21811232], which integrates an individual's variation data (SNVs, indels, and SVs) into the reference genome producing the maternal and paternal haplotypes of the individual's *personal genome* (see Fig PG6).

The tool's versatility to account for coordinate offsets between the reference and personal genome and to convert between them facilitates mapping of genomic annotated regions between the genomes. Thus, *personal annotation* can be generated by mapping Gencode annotations against the individual's personal genome. Using personal annotation in downstream analyses allows to account for differences due to impact of the personal variation on genes and other genomic elements between individuals as well as between haplotypes of the same individual.

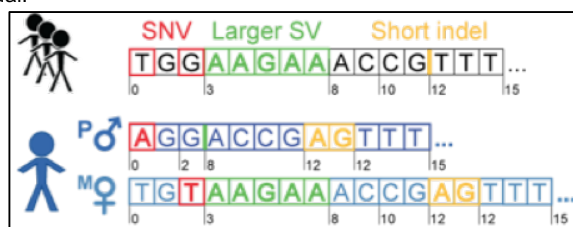


Figure PG6. Each haplotype in the diploid personal genome is derived by incorporating phased or unphased variants (SNVs, indels and SVs) into the human reference genome. The coordinates can be mapped back to the human reference coordinates to facilitate comparisons with other reference-based resources, such as gene annotations from Gencode.

We have a large experience with building personal genomes and annotations and using them in functional genomic analyses. We have previously constructed the personal diploid genome, splice-junction libraries and personalized gene annotations for NA12878. We have made this assembly available as a resource - alleleseq.gersteinlab.org - and have been updating it as new versions of the human reference genome, genomic annotations, and NA12878 genetic variation data are released.

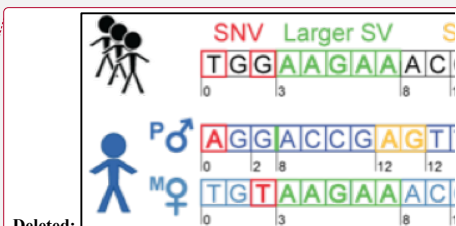
It has been demonstrated that using the diploid genome with individual's variants improves both mappability of the reads^[21811232] and the results of the downstream analyses^[26432246]. In particular, it was shown that using personal genome and annotation for NA12878 as opposed to the standard reference affected estimated expression of hundreds (525) of exons^[21811232].

Using personal genomes in analyses involving mapping of functional assay reads alleviates known biases associated with short read alignment to the reference genome: reduced

Deleted: 4 Preliminary Results In Building Personal Genomes and Annotations and Using Them in Analyses

Formatted: Font:18 pt

Formatted: Space After: 10 pt



Deleted:

mappability in regions with higher genomic variation and the preferential mapping of reads bearing the reference allele. Allele-specific analyses are particularly sensitive to these biases. For this purpose, the initial step of our *AlleleSeq* pipeline^{\cite{21811232}} involves construction of the personal diploid genome. We have spearheaded allele-specific analyses in several major consortia publications, including ENCODE and the 1000 Genomes Project^{\cite{22955620,22955619,24092746}}. The availability of an efficient computational tool enables construction of personal genomes and annotation in a high throughput fashion, as demonstrated in a recent publication^{\cite{27089393}} where we provided allele-specific annotation of variants in 382 individuals.

Formatted: Font color: Black

Deleted: 10.1038/NCOMMS11101

6 PLANS FOR PSEUDOGENE ANNOTATION

Deleted: 5

Formatted: Font:18 pt

6.1 Finishing the Mouse Annotation

Deleted: 5

6.1.1 Status on Mouse Reference Genome Pseudogene Annotation

Deleted: 5

Currently we are in the midst of completing the mouse reference genome pseudogene annotation, with plans to develop customised pseudogene annotations for the available 18 strains. Using Pseudopipe we are able to identify 18627 putative pseudogenes in the reference genome (MM8, Ensembl 83) that are classified into three groups based on their biotypes as follows: 9748 processed pseudogenes, 1940 duplicated and 6939 ambiguous. Using RCPedia we are able to identify 9755 processed pseudogenes in the reference genome. [Retrofinder predicts 18467 retrogenes in the mouse reference genome](#). The *tri-way* consensus between the [three pipelines with respect to the processed pseudogenes](#) is $\sim 80\%$. Integrating the automatic predictions with the manually curated pseudogenes we were able to annotate a comprehensive set of pseudogenes in the mouse reference genome. Table PG1 summarises the current state of pseudogene annotation.

Deleted: two automatic annotation

Deleted: over 85

Formatted: Highlight

Table PG1. Mouse reference genome MM8 pseudogene annotation. PSSD = Processed pseudogenes, DUP = duplicated pseudogenes, FRAG = ambiguous pseudogenes

Pipeline	PSSD	PSSD Parents	DUP	DUP Parents	FRAG	FRAG Parents	Total
Pseudopipe	9748	2581	1940	1146	6939	2884	18627
RCPedia	9755	2731	-	-	-	-	9755
Retrofinder	18467						18467

Formatted Table

Deleted: 5

6.1.2 Status on Human-Mouse Pseudogene Comparison

Preliminary comparative analysis of human and mouse genomes have shown that while they are divergent enough to permit a reliable identification of species specific elements, they are also similar enough to allow a reliable comparative analysis^{\cite{14746985,25157146}}.

Comparing the two organisms we found that they exhibit a similar number of pseudogenes. The pseudogene complements of both human and mouse are dominated by processed pseudogenes. At family level we found that most of the pseudogenes are lineage specific and the majority of them arise from housekeeping genes (e.g. ribosomal proteins). Also, the age distribution of mouse processed pseudogenes closely resembles that of LINEs, in contrast to human, where the age distribution closely follows Alus (SINEs).

6.1.3 Plans to Improve the Mouse Reference and Mouse Strain Pseudogene Annotation

Using the GENCODE manually annotated pseudogenes as a gold standard we plan to determine the annotation accuracy of our automatic pipelines Pseudopipe and RCPedia. Next we are going to use the false positives to refine and improve the pseudogene identification process as well as the biotype assignment.

Leveraging on the availability of the improved 18 mouse strain assemblies, we are going to extend our automatic annotation pipelines to annotate the pseudogenes in these strains. The 18 available mouse strains are 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, PWK_PhJ, Pahari_EiJ, SPRET_EiJ, WSB_EiJ.

The evolutionary distance between these strains ranges from 400,000 to 2,000 years [25038446]. The relatively small divergence time frame will allow us to map the reference mouse annotation on each of the strains. To this end we are going to use the UCSC LiftOver tool to align the reference genome annotated pseudogenes to each of the strains. In addition, the UCSC group is currently in the midst of completing a first draft of the mouse strain specific annotation of protein coding genes. We are going to use these two draft annotation sets as a support structure in building each strain's pseudogene complement. Overall this will allow us to produce a better and more accurate pseudogene annotation and will facilitate our identification of conserved elements and loss-of-function events in mouse strains.

Next, we are going to develop extensions to the available in house automatic annotation pipelines in order to use them in predicting pseudogene models in the mouse strains. The details of the proposed pipeline updates are described below.

6.2 Extending the Automatic Annotation Pipelines

6.2.1 Extending Pseudopipe

Given the close evolutionary time scale between the strains, we expect that the expressed protein amino acid sequence to be preserved across all strains for the conserved protein coding genes. As such we are going to use the conserved protein coding genes between each strain and the reference genome as input for identifying pseudogenes. The extended Pseudopipe workflow is summarized in the following steps: 1) Identify the consensus protein coding genes between strain and reference; 2) extract the amino acid sequence of the

Deleted: 5

Deleted: and study

Deleted: conservation of pseudogenes as well as

Deleted: loss

Deleted: function events.

Deleted: the resulting

Formatted: Not Highlight

Formatted: Not Highlight

Deleted: 5

Deleted: 5

conserved proteins from the available ENSEMBL peptide database for the mouse reference genome; 3) mark and identify the coordinates of the consensus protein coding genes in the analyzed strain; 4) use a six frame blast homology search to match the consensus peptides to the strain sequence; 5) refine results and eliminate redundant hits (e.g. remove matches that overlap protein coding exons); 6) merge hits and identify parents; 7) align parents and pseudogenes and check for the presence of disablements (e.g. frameshifts, premature stop codons); 7) assign pseudogene biotype.

Deleted: frame shifts

6.2.2 Extending RCPedia

Deleted: 5

Similarly to the protein sequence approach, transcript sequence is expected to be conserved at short evolutionary time scales. RCPedia will be adapted to integrate gold standard transcript annotation, such as GENCODE mouse annotation, and annotations based on the strain genome. The extended RCPedia pipeline is summarized as follow: 1) Merge multiple annotations of parental transcripts using an hierarchical prioritization; 2) align transcripts sequences to the target genome and extract alignments blocks and their distances; 3) select alignments containing mostly intronless blocks; 4) refine results removing alignments with most of the parental introns and remove putative genomic duplications; 5) merge call sets and select the most likely parent transcript; 7) calculate properties of the putative pseudogene such as target site duplication sequence, identity and polyA length.

6.2.3 Using Cloud Environment to Update the Pseudogene Annotation Pipelines

Deleted: 5

Both Pseudopipe and RCPedia pipelines are broadly used by the pseudogene research community and both are available through our online resource pseudogene.org. Collectively they use many different standalone tools such as aligners, toolsets and well established annotation software such as repeatmasker. The invariably complex environment necessary to install and configure these pipelines can create difficulties to the end user. In order to mitigate dependency and compatibility issues, we plan to create docker images for both pipelines and make them publicly available after the mentioned extension. Docker images will contain all dependencies necessary to set up Pseudopipe and RCPedia as well as all in house scripts. Parameterization and fine tuning will be made by a single configuration file editable by user. We will also create amazon machine images (AMI) compatible with Amazon AWS and other major cloud services so users can easily annotate additional genomes.

Deleted: respectively

Deleted: at

Deleted: <http://>

Deleted: [/](#)

Deleted: . [\[\[use Yale url\]\]](#)

Formatted: Not Highlight

6.3 Developing a Cross Strain Pseudogene Analysis Database

Deleted: 5

The results from the pseudogene annotation in the mouse strains will be collected into an online database that is going to be made freely accessible. Our aim is to build a vcf-like format relational database that will facilitate the cross strain analysis (see Table PG2).

Table PG2: Relational database providing a cross strain annotation of pseudogenes.

Pseudogene ID	Reference	Strain Specific	Observations
ENSMUST00000050706	17:58593659-58594843	-	129S1:StrainLoc:PG;AJ:StrainLoc:LOF;AKR:na:na;...
ENSMUST00000078706	-	CAROLI:14:68538759-68594843	129S1:StrainLoc:PG;AJ:StrainLoc:PG;AKR:StrainLoc:PG;...

Formatted Table

For this we are going to create a superset containing unique pseudogene instances from all the strains. Each pseudogene will be described by both reference genome and strain specific coordinates. The location fields will link to the Genome Browser to facilitate a visual description of the pseudogene structure.

The observations field will contain information with respect to the genomic location of the pseudogene in each analysed strain as well as flags characterizing the pseudogene sequence in that particular strain. If the pseudogene is not present in the analysed strain the flag will be "NA". If the aligned sequence is a pseudogene in the analysed strain, this will be indicated by a "PG" flag. If the pseudogene aligns with a protein coding gene in the strain of interest, the flag will be "LOF" indicating that the pseudogenization process resulted in a loss-of-function (LOF) event with respect to the analysed strain.

Deleted: L
Deleted: F

Further we will extend the Observation field with information regarding the functional characterization of pseudogenes (e.g. transcription flag, biochemical activity flags, etc.).

7 PLANS FOR ANNOTATING PSEUDOGENE ACTIVITY

Deleted: 6
Formatted: Font:18 pt

7.1 Plans for Analysing Pan-Tissue Pseudogene Transcription in Mouse Strains

Deleted: 6
Formatted: Heading 2
Deleted: ... [3]

We are going to leverage our experience in pseudogene transcription analysis in human, worm and fly to study the pseudogene transcription in mouse, significantly improving on previous efforts. We are going to focus on identifying tissue and strain specific transcriptionally active pseudogenes. In particular we will highlight pseudogenes that have a high coexpression correlation with their parents or are differentially transcribed with respect to their functional paralogs.

Deleted: the 18
Formatted: Not Highlight
Deleted: strain
Formatted: Not Highlight
Deleted: [need]
Formatted: Not Highlight
Deleted: diff. w/ previous

Our approach is summarized in the following steps: 1) calculate the genome mappability in each mouse strain; 2) remove low mappability regions from the pseudogene annotation in each mouse strain; 3) calculate the RPKM value of each pseudogene based on the RNA-Seq reads mapped to the remaining high mappability regions in that pseudogene locus; 4) quantile normalize the transcriptome signals across all the mouse strains and identify transcribed pseudogenes uniformly in each strain. We will combine the pseudogene transcription results with their annotation to study the strain-specificity of pseudogenes. For example, a pseudogene may exist and be transcribed in only one mouse strain, or a pseudogene may exist in all mouse strains but be transcribed in only one or a few closely related strains. Such

information and data resources will benefit the evolutionary studies on mouse and the comparative studies between mouse, human, and other model organisms.

Deleted: mouse

7.2 Integrating Functional Genomics Data to Characterize Pseudogene Activity in Mouse Strains

Deleted: 6.1

Formatted: Heading 2

We aim to integrate tissue specific transcription information and regulatory data with the pseudogene annotation in order to characterize pseudogene activity. In particular, we will focus on the transcriptomics (ENCODE, BrainSpan, TCGA), epigenomics (ENCODE, Roadmap Epigenomics) and cis-regulatory interactions data (GTEx, PsychENCODE). These datasets will allow us to provide annotation on tissue-specific pseudogene transcription and tissue-specific pseudogene regulation. Such information will be valuable for understanding the biological consequences from the pseudogene activities, such as the regulatory mechanisms of pseudogene transcription, and whether transcribed pseudogenes may perform regulatory roles through interaction with their functional paralogs.

Formatted: Space After: 10 pt

Deleted: ing

8 PLANS FOR ANNOTATING AND ANALYSING LOSS-OF-FUNCTION EVENTS IN MOUSE

Deleted: 6.2 Plans for Analysing Loss of Function Events in Mouse Strains

... [4]

8.1 Identifying Unitary and Polymorphic Pseudogenes Across Mouse Strains

Formatted: Heading 2

Building on our previously developed human unitary pseudogene annotation pipeline [20210993] we aim to develop a reliable framework for the identification of unitary pseudogenes across the 18 available strains. The unitary mouse pseudogene pipeline can be summarized as follows. First we will create a global inventory of orthologs between the mouse strains using the available multi sequence alignment data from UCSC. Next we are going to identify homologous regions between any two strains and annotate the syntenic ones. Finally we will conduct a survey of gene disablements in the syntenic regions. We are going to use the available mouse genome variation data to filter our false positives.

In order to create a comprehensive set of polymorphic pseudogenes in mouse we will focus on annotating variants that change the strain genome stop codon to a functional allele across another mouse strains. For this we will use our previously developed pipeline VAT to annotate the SNPs of interests. Next we will extend the VAT to annotate frame shifts that revert a stop codon across two mouse strains.

8.2 Building models to describe Loss-of-Function, Gene Death and Pseudogenization

Deleted: 6.2

Formatted: Heading 2

We will build on our experience in developing variant annotation tools to create a pipeline that will provide an extensive annotation of putative LOF variants. We will include variants causing

Deleted: o

premature stop codons, canonical splice-site disruption and frameshift-causing indels as putative LOF variants. The pipeline will feature (1) function-based annotations; (2) evolutionary conservation; and (3) biological network data. For comprehensive functional annotation we will integrate several annotation resources such as PFAM and SMART functional domains, signal peptide and transmembrane annotations, post-translational modification sites, NMD prediction, and structure-based features such as SCOP domains and disordered residues. For evolutionary conservation, our pipeline will output variant position-specific GERP scores, which is a measure of evolutionary conservation and dN/dS values. In addition, we will evaluate if the region removed due to the truncation of the coding sequence is evolutionarily conserved based on GERP and PhyloCSF constrained elements. Our model will also include network features to predict disease causing variants: we will use a proximity parameter that gives the number of disease genes connected to a gene in a protein-protein interaction network and the shortest path to the nearest disease gene. The pipeline will also include features to help identify erroneous LOF calls, potential mismapping, and annotation errors, because LOF variant calls have been shown to be enriched for annotation and sequencing artifacts.

Deleted: o

Deleted: o

Deleted: o

Deleted: o

To understand the impact of putative LOF variants on gene function we will develop a prediction model to classify premature stop causing variants into those that are benign, those that lead to recessive disease and those that lead to dominant disease using the annotation output as predictive features. To build the classifier, we will use benign homozygous premature stop variants, dominant heterozygous and recessive homozygous disease premature stop mutations. We will build the classifier to distinguish among the three classes and will provide class probability estimates for each mutation. To validate our classifier we plan to use LOF from Mendelian Diseases, Cancer samples and healthy control datasets such as 1000 Genomes and ExAC. The classifier will also be applied to variants in mouse strains and will be used to classify variations into Loss-of-Function, Gene Death and Pseudogenization.

9 PLANS FOR USE OF PERSONAL ANNOTATION

Annotation based on the current human reference genome does not account for variation in the number of functional genes between people and does not provide an accurate and complete set of an individual's genes and other genomic elements. We propose to use the diploid personal genome and personal annotation as a more accurate representation of an individual's genome. Individual variations can affect gene annotations and thus sequence variations need to be taken into account for annotation purposes.

Deleted: 7

Formatted: Font: 18 pt

For this we will develop a personal genome annotation resource containing a number of tools and utilities for constructing a personal genome and for creating a personal GENCODE annotation set. We will incorporate a personal genome construction step into our genome and variant annotation pipelines. Using well-characterized public genomes as well as matching variant calls and functional datasets, we will evaluate the impact of using personal genome and annotations for various types of genomic analyses. For these genomes we will generate a reference set of personal diploid (haplotype-resolved) annotations and make them publicly available.

Deleted: We will incorporate

Formatted: Highlight

Deleted: We will also produce tools for producing a personal genocode on new personal genomes as they are sequenced. [\[\[More on tools\]\]](#)

In particular, given an individual's variation data, the proposed annotation resource can be used to identify and further analyse GENCODE-annotated features characteristic to the individual, such as his/her distinctive set of functional genes or structures of variant-affected transcripts. Using our automated in house annotation pipelines we are going to create a comprehensive personal pseudogene complement. We will use the newly constructed personal annotations to identify LOF and pseudogenization events by comparison with the reference genome. We are going to assess the annotated personal SNPs for allele specific expression using the data from AlleleDB \cite{27089393}, an online repository that provides genomic annotation of cis-regulatory single nucleotide variants associated with allele-specific binding and expression.

Deleted: Furthermore we

Next, by integrating Mendelian disease and cancer data we will be able to filter the LOF and pseudogenization variants and characterize them with respect to their disease driver potential. In particular we are going to use VAT and the newly proposed LOF analysis pipeline as described above (See Sections 4 and 8.2). Further, we are going to test the presence of mouse orthologs for all the curated genes affected by LOFs and pseudogenization variants, determining where or not mouse is a suitable model organism to study them.

Deleted: 3.2

Deleted: 6.2.2). [more] Finally

Deleted: will curate a comprehensive set

Deleted: disease related

Deleted: and identify their [????] orthologs in mouse genomes,

We aim to integrate all the personal annotation tools in an online framework that can easily be applied to newly sequenced individual genomes.

REFERENCES

Formatted: Font:18 pt

1. Kalyana-Sundaram, S., Kumar-Sinha, C., Shankar, S., Robinson, D. R., Wu, Y.-M., Cao, X., Asangani, I. A., Kothari, V., Prensner, J. R., Lonigro, R. J., Iyer, M. K., Barrette, T., Shanmugam, A., Dhanasekaran, S. M., Palanisamy, N., & Chinnaiyan, A. M. (2012) Expressed pseudogenes in the transcriptional landscape of human cancers. *Cell* 149, 1622-34, PMID:22726445.
2. Swami, M. (2010) Epigenetics: Demethylation links cell fate and cancer. *Nat Rev Cancer* 10, 740, PMID:21080588.
3. Poliseno, L. (2012) Pseudogenes: newly discovered players in human cancer. *Sci Signal* 5, re5, PMID:22990117.
4. Abyzov, A., Iskow, R., Gokcumen, O., Radke, D. W., Balasubramanian, S., Pei, B., Habegger, L., 1000 Genomes Project Consortium, Lee, C., & Gerstein, M. (2013) Analysis of variable retroduplications in human populations suggests coupling of retrotransposition to cell division. *Genome Res* 23, 2042-52, PMID:24026178.
5. MacArthur, D. G., Balasubramanian, S., Frankish, A., Huang, N., Morris, J., Walter, K., Jostins, L., Habegger, L., Pickrell, J. K., Montgomery, S. B., Albers, C. A., Zhang, Z. D., Conrad, D. F., Lunter, G., Zheng, H., Ayub, Q., DePristo, M. A., Banks, E., Hu, M., Handsaker, R. E., Rosenfeld, J. A., Fromer, M., Jin, M., Mu, X. J., Khurana, E., Ye, K., Kay, M., Saunders, G. I., Suner, M.-M., Hunt, T., Barnes, I. H. A., Amid, C., Carvalho-Silva, D. R., Bignell, A. H., Snow, C., Yngvadottir, B., Bumpstead, S., Cooper, D. N., Xue, Y., Romero, I. G., Wang, J., Li, Y., Gibbs, R. A., McCarroll, S. A., Dermitzakis, E. T., Pritchard, J. K., Barrett, J. C., Harrow, J., Hurles, M. E., Gerstein, M. B., & Tyler-Smith, C. (2012) A systematic survey of loss-of-function variants in human protein-coding genes. *Science* 335, 823-828, PMID:22344438.
6. Balasubramanian, S., Habegger, L., Frankish, A., MacArthur, D. G., Harte, R., Tyler-Smith, C., Harrow, J., & Gerstein, M. (2011) Gene inactivation and its implications for annotation in the era of personal genomics. *Genes Dev* 25, 1-10, PMID:21205862.
7. Zhang, Z., Carriero, N., Zheng, D., Karro, J., Harrison, P. M., & Gerstein, M. (2006) PseudoPipe: an automated pseudogene identification pipeline. *Bioinformatics* 22, 1437-9, PMID:16574694.
8. Zheng, D. & Gerstein, M. B. (2006) A computational approach for identifying pseudogenes in the ENCODE regions. *Genome Biol* 7 Suppl 1, S13.1-10, PMID:16925835.
9. Pei, B., Sisu, C., Frankish, A., Howald, C., Habegger, L., Mu, X. J., Harte, R., Balasubramanian, S., Tanzer, A., Diekhans, M., Reymond, A., Hubbard, T. J., Harrow, J., & Gerstein, M. B. (2012) The GENCODE pseudogene resource. *Genome Biol* 13, R51, PMID:22951037.
10. Navarro, F. C. P. & Galante, P. A. F. (2013) RCPedia: a database of retrocopied genes. *Bioinformatics* 29, 1235-7, PMID:23457042.

11. [Baertsch, R., Diekhans, M., Kent, W. J., Haussler, D., & Brosius, J. \(2008\) Retrocopy contributions to the evolution of the human genome. *BMC Genomics* 9, 466, PMID:18842134.](#)
12. [Jurka, J., Kapitonov, V. V., Pavlicek, A., Klonowski, P., Kohany, O., & Walichiewicz, J. \(2005\) Repbase Update, a database of eukaryotic repetitive elements. *Cytogenet Genome Res* 110, 462-7, PMID:16093699.](#)
13. [Smit, A. F. A., Hubley, R., & Green, P. \(1996-2010\) RepeatMasker Open-3.0. <http://www.repeatmasker.org>, , PMID:.](#)
14. [Benson, G. \(1999\) Tandem repeats finder: a program to analyze DNA sequences. *Nucleic Acids Res* 27, 573-80, PMID:9862982.](#)
15. [Kent, W. J., Baertsch, R., Hinrichs, A., Miller, W., & Haussler, D. \(2003\) Evolution's cauldron: duplication, deletion, and rearrangement in the mouse and human genomes. *Proc Natl Acad Sci U S A* 100, 11484-9, PMID:14500911.](#)
16. [Sisu, C., Pei, B., Leng, J., Frankish, A., Zhang, Y., Balasubramanian, S., Harte, R., Wang, D., Rutenberg-Schoenberg, M., Clark, W., Diekhans, M., Rozowsky, J., Hubbard, T., Harrow, J., & Gerstein, M. B. \(2014\) Comparative analysis of pseudogenes across three phyla. *Proc Natl Acad Sci U S A* 111, 13361-6, PMID:25157146.](#)
17. [Karro, J. E., Yan, Y., Zheng, D., Zhang, Z., Carriero, N., Cayting, P., Harrision, P., & Gerstein, M. \(2007\) Pseudogene.org: a comprehensive database and comparison platform for pseudogene annotation. *Nucleic Acids Res* 35, D55-60, PMID:17099229.](#)
18. [Lam, H. Y. K., Khurana, E., Fang, G., Cayting, P., Carriero, N., Cheung, K.-H., & Gerstein, M. B. \(2009\) Pseudofam: the pseudogene families database. *Nucleic Acids Res* 37, D738-43, PMID:18957444.](#)
19. [Sudmant, P. H., Rausch, T., Gardner, E. J., Handsaker, R. E., Abyzov, A., Huddleston, J., Zhang, Y., Ye, K., Jun, G., Hsi-Yang Fritz, M., Konkil, M. K., Malhotra, A., Stütz, A. M., Shi, X., Paolo Casale, F., Chen, J., Hormozdiari, F., Dayama, G., Chen, K., Malig, M., Chaisson, M. J. P., Walter, K., Meiers, S., Kashin, S., Garrison, E., Auton, A., Lam, H. Y. K., Jasmine Mu, X., Alkan, C., Antaki, D., Bae, T., Cerveira, E., Chines, P., Chong, Z., Clarke, L., Dal, E., Ding, L., Emery, S., Fan, X., Gujral, M., Kahveci, F., Kidd, J. M., Kong, Y., Lameijer, E.-W., McCarthy, S., Flicek, P., Gibbs, R. A., Marth, G., Mason, C. E., Menelaou, A., Muzny, D. M., Nelson, B. J., Noor, A., Parrish, N. F., Pendleton, M., Quitadamo, A., Raeder, B., Schadt, E. E., Romanovitch, M., Schlattl, A., Sebra, R., Shabalina, A. A., Untergasser, A., Walker, J. A., Wang, M., Yu, F., Zhang, C., Zhang, J., Zheng-Bradley, X., Zhou, W., Zichner, T., Sebati, J., Batzer, M. A., McCarroll, S. A., 1000 Genomes Project Consortium, Mills, R. E., Gerstein, M. B., Bashir, A., Stegle, O., Devine, S. E., Lee, C., Eichler, E. E., & Korbelt, J. O. \(2015\) An integrated map of structural variation in 2,504 human genomes. *Nature* 526, 75-81, PMID:26432246.](#)
20. [Holford, M. E., Khurana, E., Cheung, K.-H., & Gerstein, M. \(2010\) Using semantic web rules to reason on an ontology of pseudogenes. *Bioinformatics* 26, i71-8, PMID:20529940.](#)

Deleted: 2

Deleted: 3

Deleted: 4

Deleted: 15.

<p>21, Liu, Y., Harrison, P. M., Kunin, V., & Gerstein, M. (2004) Comprehensive analysis of pseudogenes in prokaryotes: widespread gene decay and failure of putative horizontally transferred genes.. <i>Genome Biol</i> 5, R64, PMID:15345048.</p>	<p>Deleted: 6</p>
<p>22, Harrison, P. M., Carriero, N., Liu, Y., & Gerstein, M. (2003) A "polyORFomic" analysis of prokaryote genomes using disabled-homology filtering reveals conserved but undiscovered short ORFs.. <i>J Mol Biol</i> 333, 885-892, PMID:14583187.</p>	<p>Deleted: 17.</p>
<p>23, Harrison, P., Kumar, A., Lan, N., Echols, N., Snyder, M., & Gerstein, M. (2002) A small reservoir of disabled ORFs in the yeast genome and its implications for the dynamics of proteome evolution.. <i>J Mol Biol</i> 316, 409-419, PMID:11866506.</p>	<p>Deleted: 18.</p>
<p>24, Zhang, Z. L., Harrison, P. M., & Gerstein, M. (2002) Digging deep for ancient relics: a survey of protein motifs in the intergenic sequences of four eukaryotic genomes.. <i>J Mol Biol</i> 323, 811-822, PMID:12417195.</p>	<p>Deleted: 19.</p>
<p>25, Harrison, P. M. & Gerstein, M. (2002) Studying genomes through the aeons: protein families, pseudogenes and proteome evolution. <i>J Mol Biol</i> 318, 1155-74, PMID:12083509.</p>	<p>Deleted: 0</p>
<p>26, Harrison, P. M., Echols, N., & Gerstein, M. B. (2001) Digging for dead genes: an analysis of the characteristics of the pseudogene population in the <i>Caenorhabditis elegans</i> genome.. <i>Nucleic Acids Res</i> 29, 818-830, PMID:11160906.</p>	<p>Deleted: 1</p>
<p>27, Echols, N., Harrison, P., Balasubramanian, S., Luscombe, N. M., Bertone, P., Zhang, Z., & Gerstein, M. (2002) Comprehensive analysis of amino acid and nucleotide composition in eukaryotic genomes, comparing genes and pseudogenes.. <i>Nucleic Acids Res</i> 30, 2515-2523, PMID:12034841.</p>	<p>Deleted: 2</p>
<p>28, Harrison, P. M., Milburn, D., Zhang, Z., Bertone, P., & Gerstein, M. (2003) Identification of pseudogenes in the <i>Drosophila melanogaster</i> genome.. <i>Nucleic Acids Res</i> 31, 1033-1037, PMID:12560500.</p>	<p>Deleted: 3</p>
<p>29, Liu, Y.-J., Zheng, D., Balasubramanian, S., Carriero, N., Khurana, E., Robilotto, R., & Gerstein, M. B. (2009) Comprehensive analysis of the pseudogenes of glycolytic enzymes in vertebrates: the anomalously high number of GAPDH pseudogenes highlights a recent burst of retrotrans-positional activity.. <i>BMC Genomics</i> 10, 480, PMID:19835609.</p>	<p>Deleted: 4</p>
<p>30, Balasubramanian, S., Harrison, P., Hegyi, H., Bertone, P., Luscombe, N., Echols, N., McGarvey, P., Zhang, Z., & Gerstein, M. (2002) SNPs on human chromosomes 21 and 22 - analysis in terms of protein features and pseudogenes.. <i>Pharmacogenomics</i> 3, 393-402, PMID:12052146.</p>	<p>Deleted: 25.</p>
<p>31, Zhang, Z. & Gerstein, M. (2003) The human genome has 49 cytochrome c pseudogenes, including a relic of a primordial gene that still functions in mouse.. <i>Gene</i> 312, 61-72, PMID:12909341.</p>	<p>Deleted: 26.</p>
<p>32, Zhang, Z. D., Cayting, P., Weinstock, G., & Gerstein, M. (2008) Analysis of nuclear receptor pseudogenes in vertebrates: how the silent tell their stories.. <i>Mol Biol Evol</i> 25, 131-143, PMID:18065488.</p>	<p>Deleted: 7</p>

33. Balasubramanian, S., Zheng, D., Liu, Y.-J., Fang, G., Frankish, A., Carriero, N., Robilotto, R., Cayting, P., & Gerstein, M. (2009) Comparative analysis of processed ribosomal protein pseudogenes in four mammalian genomes. *Genome Biol* 10, R2, PMID:19123937.

Deleted: 28.

34. Zhang, Z. D., Frankish, A., Hunt, T., Harrow, J., & Gerstein, M. (2010) Identification and analysis of unitary pseudogenes: historic and contemporary gene losses in humans and other primates.. *Genome Biol* 11, R26, PMID:20210993.

Deleted: 29.

35. Harrison, P. M., Zheng, D., Zhang, Z., Carriero, N., & Gerstein, M. (2005) Transcribed processed pseudogenes in the human genome: an intermediate form of expressed retrosequence lacking protein-coding ability. *Nucleic Acids Res* 33, 2374-83, PMID:15860774.

Deleted: 0

36. Svensson, O., Arvestad, L., & Lagergren, J. (2006) Genome-wide survey for biologically functional pseudogenes. *PLoS Comput Biol* 2, e46, PMID:16680195.

Deleted: 1

37. Zheng, D., Zhang, Z., Harrison, P. M., Karro, J., Carriero, N., & Gerstein, M. (2005) Integrated pseudogene annotation for human chromosome 22: evidence for transcription. *J Mol Biol* 349, 27-45, PMID:15876366.

Deleted: 2

38. Zheng, D., Frankish, A., Baertsch, R., Kapranov, P., Reymond, A., Choo, S. W., Lu, Y., Denoeud, F., Antonarakis, S. E., Snyder, M., Ruan, Y., Wei, C.-L., Gingeras, T. R., Guigo, R., Harrow, J., & Gerstein, M. B. (2007) Pseudogenes in the ENCODE regions: consensus annotation, analysis of transcription, and evolution.. *Genome Res* 17, 839-851, PMID:17568002.

Deleted: 3

39. Frith, M. C., Wilming, L. G., Forrest, A., Kawaji, H., Tan, S. L., Wahlestedt, C., Bajic, V. B., Kai, C., Kawai, J., Carninci, P., Hayashizaki, Y., Bailey, T. L., & Huminiecki, L. (2006) Pseudomessenger RNA: phantoms of the transcriptome. *PLoS Genet* 2, e23, PMID:16683022.

Deleted: 4

40. Zhang, Z., Carriero, N., & Gerstein, M. (2004) Comparative analysis of processed pseudogenes in the mouse and human genomes.. *Trends Genet* 20, 62-67, PMID:14746985.

Deleted: 35.

Moved (insertion) [1]

41. Habegger, L., Balasubramanian, S., Chen, D. Z., Khurana, E., Sboner, A., Harmanci, A., Rozowsky, J., Clarke, D., Snyder, M., & Gerstein, M. (2012) VAT: a computational framework to functionally annotate variants in personal genomes within a cloud-computing environment.. *Bioinformatics* 28, 2267-2269, PMID:22743228.

42. Khurana, E., Lam, H. Y. K., Cheng, C., Carriero, N., Cayting, P., & Gerstein, M. B. (2010) Segmental duplications in the human genome reveal details of pseudogene formation. *Nucleic Acids Res* 38, 6997-7007, PMID:20615899.

Deleted: 36.

43. Rozowsky, J., Abyzov, A., Wang, J., Alves, P., Raha, D., Harmanci, A., Leng, J., Bjornson, R., Kong, Y., Kitabayashi, N., Bhardwaj, N., Rubin, M., Snyder, M., & Gerstein, M. (2011) AlleleSeq: analysis of allele-specific expression and binding in a network framework. *Mol Syst Biol* 7, 522, PMID:21811232.

Deleted: 7

Deleted: .

44. Djebali, S., Davis, C. A., Merkel, A., Dobin, A., Lassmann, T., Mortazavi, A., Tanzer, A., Lagarde, J., Lin, W., Schlesinger, F., Xue, C., Marinov, G. K., Khatun, J., Williams, B. A., Zaleski, C., Rozowsky, J., Röder, M., Kokocinski, F., Abdelhamid, R. F., Alioto, T., Antoshechkin, I., Baer, M. T., Bar, N. S., Batut, P., Bell, K., Bell, I., Chakraborty, S., Chen, X., Chrast, J., Curado, J., Derrien, T., Drenkow, J., Dumais, E., Dumais, J., Duttagupta, R.,

Deleted: 38.

Falconnet, E., Fastuca, M., Fejes-Toth, K., Ferreira, P., Foissac, S., Fullwood, M. J., Gao, H., Gonzalez, D., Gordon, A., Gunawardena, H., Howald, C., Jha, S., Johnson, R., Kapranov, P., King, B., Kingswood, C., Luo, O. J., Park, E., Persaud, K., Preall, J. B., Ribeca, P., Risk, B., Robyr, D., Sammeth, M., Schaffer, L., See, L.-H., Shahab, A., Skancke, J., Suzuki, A. M., Takahashi, H., Tilgner, H., Trout, D., Walters, N., Wang, H., Wrobel, J., Yu, Y., Ruan, X., Hayashizaki, Y., Harrow, J., Gerstein, M., Hubbard, T., Reymond, A., Antonarakis, S. E., Hannon, G., Giddings, M. C., Ruan, Y., Wold, B., Carninci, P., Guigó, R., & Gingeras, T. R. (2012) Landscape of transcription in human cells. *Nature* 489, 101-8, PMID:22955620.

45. Gerstein, M. B., Kundaje, A., Hariharan, M., Landt, S. G., Yan, K.-K., Cheng, C., Mu, X. J., Khurana, E., Rozowsky, J., Alexander, R., Min, R., Alves, P., Abyzov, A., Addleman, N., Bhardwaj, N., Boyle, A. P., Cayting, P., Charos, A., Chen, D. Z., Cheng, Y., Clarke, D., Eastman, C., Euskirchen, G., Frietze, S., Fu, Y., Gertz, J., Grubert, F., Harmanci, A., Jain, P., Kasowski, M., Lacroute, P., Leng, J., Lian, J., Monahan, H., O'Geen, H., Ouyang, Z., Partridge, E. C., Patacsil, D., Pauli, F., Raha, D., Ramirez, L., Reddy, T. E., Reed, B., Shi, M., Slifer, T., Wang, J., Wu, L., Yang, X., Yip, K. Y., Zilberman-Schapira, G., Batzoglou, S., Sidow, A., Farnham, P. J., Myers, R. M., Weissman, S. M., & Snyder, M. (2012) Architecture of the human regulatory network derived from ENCODE data. *Nature* 489, 91-100, PMID:22955619.

46. Khurana, E., Fu, Y., Colonna, V., Mu, X. J., Kang, H. M., Lappalainen, T., Sboner, A., Lochovsky, L., Chen, J., Harmanci, A., Das, J., Abyzov, A., Balasubramanian, S., Beal, K., Chakravarty, D., Challis, D., Chen, Y., Clarke, D., Clarke, L., Cunningham, F., Evani, U. S., Flicek, P., Fragoza, R., Garrison, E., Gibbs, R., Gümüs, Z. H., Herrero, J., Kitabayashi, N., Kong, Y., Lage, K., Liliushvili, V., Lipkin, S. M., MacArthur, D. G., Marth, G., Muzny, D., Pers, T. H., Ritchie, G. R. S., Rosenfeld, J. A., Sisu, C., Wei, X., Wilson, M., Xue, Y., Yu, F., 1000 Genomes Project Consortium, Dermitzakis, E. T., Yu, H., Rubin, M. A., Tyler-Smith, C., & Gerstein, M. (2013) Integrative annotation of variants from 1092 humans: application to cancer genomics. *Science* 342, 1235587, PMID:24092746.

47. Chen, J., Rozowsky, J., Galeev, T. R., Harmanci, A., Kitchen, R., Bedford, J., Abyzov, A., Kong, Y., Regand, L., & Gerstein, M. (2016) A uniform survey of allele-specific binding and expression over 1000-Genomes-Project individuals. *Nat Commun* 10.1038/NCOMMS11101, PMID:27089393.

48. Zheng, J., Chen, Y., Deng, F., Huang, R., Petersen, F., Ibrahim, S., & Yu, X. (2014) mtDNA sequence, phylogeny and evolution of laboratory mice. *Mitochondrion* 17, 126-31, PMID:25038446.

Deleted: 39.

Deleted: 0

Deleted: 1

Deleted: in press.

Deleted: 10.1038/NCOMMS11101

Deleted: 2

Moved up [1]: Zhang, Z., Carriero, N., & Gerstein, M. (2004) Comparative analysis of processed pseudogenes in the mouse and human genomes.. *Trends Genet* 20, 62-67, PMID:14746985. .

Deleted: 43.

2 Preliminary Results & Experience with Pseudogenes Annotation

3 Preliminary Results and Experience with Functional Characterization of Pseudogene Activity

6.1.1 Pseudogene Transcription

6.2 Plans for Analysing Loss of Function Events in Mouse Strains

6.2