

# Eutherian comparative genomic analysis protocol

Marko Premzl (✉ [Marko.Premzl@alumni.anu.edu.au](mailto:Marko.Premzl@alumni.anu.edu.au))

Marko Premzl PhD, ANU Alumni

---

## Method Article

**Keywords:** comparative genomics, Eutheria, RRID:SCR\_014401

**Posted Date:** March 12th, 2018

**DOI:** <https://doi.org/10.1038/protex.2018.028>

**License:**  This work is licensed under a Creative Commons Attribution 4.0 International License.

[Read Full License](#)

---

**Version of Record:** A version of this preprint was published at Protocol Exchange on March 12th, 2018.  
See the published version at <https://doi.org/10.1038/protex.2018.028>.

# Abstract

The eutherian genomics momentum greatly advanced biology and medicine. Nevertheless, future revisions and updates of eutherian genomic sequence data sets were expected, due to potential genomic sequence errors and incompleteness of genomic sequences. The **eutherian comparative genomic analysis protocol** was established as guidance in protection against potential genomic sequence errors in public eutherian genomic sequence assemblies. The protocol revised, updated and published 10 major eutherian gene data sets, including 1293 complete coding sequences deposited in European Nucleotide Archive as curated third party data gene data sets under accession numbers: "FR734011-FR734074":<https://www.ebi.ac.uk/ena/data/view/FR734011-FR734074> , "HF564658-HF564785":<https://www.ebi.ac.uk/ena/data/view/HF564658-HF564785> , "HF564786-HF564815":<https://www.ebi.ac.uk/ena/data/view/HF564786-HF564815> , "HG328835-HG329089":<https://www.ebi.ac.uk/ena/data/view/HG328835-HG329089> , "HG426065-HG426183":<https://www.ebi.ac.uk/ena/data/view/HG426065-HG426183> , "HG931734-HG931849":<https://www.ebi.ac.uk/ena/data/view/HG931734-HG931849> , "LM644135-LM644234":<https://www.ebi.ac.uk/ena/data/view/LM644135-LM644234> , "LN874312-LN874522":<https://www.ebi.ac.uk/ena/data/view/LN874312-LN874522> , "LT548096-LT548244":<https://www.ebi.ac.uk/ena/data/view/LT548096-LT548244> and "LT631550-LT631670":<https://www.ebi.ac.uk/ena/data/view/LT631550-LT631670> .

## Introduction

**Background** The eutherian genomics momentum greatly advanced biological and medical sciences. For example, one major aim of initial sequencing and analysis of human genome was to update and revise human genes, as well as to uncover potential new drugs, drug targets, and molecular markers in medical diagnostics. Yet, future revisions and updates of eutherian genomic sequence data sets were expected, due to potential genomic sequence errors and incompleteness of reference genomic sequences. The potential genomic sequence errors included DNA sequencing errors (artefactual nucleotide deletions, insertions and substitutions) and bioinformatical errors (erroneous gene annotations, genomic sequence misassemblies). Among others, the potential genomic sequence errors could bias phylogenetic analyses. Thus, the **eutherian comparative genomic analysis protocol**

"RRID:SCR\_014401":[https://scicrunch.org/resources/Any/record/nlx\\_144509-1/SCR\\_014401/resolver?q=RRID:SCR\\_014401&l=RRID:SCR\\_014401](https://scicrunch.org/resources/Any/record/nlx_144509-1/SCR_014401/resolver?q=RRID:SCR_014401&l=RRID:SCR_014401) was established as guidance in protection against potential genomic sequence errors in public eutherian genomic sequence assemblies

"1":<https://dx.doi.org/10.1007/s10142-012-0291-2> , "2":<https://dx.doi.org/10.1007/s00438-013-0801-5> ,  
"3":<https://dx.doi.org/10.1016/j.gene.2014.02.049> , "4":<https://dx.doi.org/10.1016/j.gdata.2014.08.003> ,  
"5":<https://dx.doi.org/10.1016/j.mgene.2015.02.005> , "6":<https://dx.doi.org/10.1016/j.gdata.2015.09.007> ,  
"7":<https://dx.doi.org/10.1016/j.dib.2015.11.056> , "8":<https://dx.doi.org/10.1007/s00251-015-0887-5> ,  
"9":<https://dx.doi.org/10.1016/j.genrep.2016.10.009> ,  
"10":<https://dx.doi.org/10.1016/j.ymgmr.2017.01.009> .

## Reagents

**\*\*On-line resources\*\*** "Ensembl genome browser":<https://www.ensembl.org/index.html> "European Nucleotide Archive":<https://www.ebi.ac.uk/ena> "National Center for Biotechnology Information":<https://www.ncbi.nlm.nih.gov>

## Equipment

Personal computer

## Procedure

**\*\*Eutherian comparative genomic analysis protocol\*\*** The **\*\*eutherian comparative genomic analysis protocol\*\*** integrated gene annotations, phylogenetic analysis and protein molecular evolution analysis into one framework of eutherian gene descriptions \ (Figure 1). **\*\*Gene annotations\*\*** The eutherian gene annotations included gene identifications in public reference genomic sequence assemblies, analyses of gene features, multiple pairwise genomic sequence alignments and tests of reliability of eutherian public genomic sequences. The protocol used sequence alignment editor "BioEdit":<http://www.mbio.ncsu.edu/BioEdit/bioedit.html> in all analyses of nucleotide and protein sequences. The identifications of potential coding sequence used eutherian reference genomic sequence assemblies downloaded from "National Center for Biotechnology Information \ (NCBI) GenBank":[ftp://ftp.ncbi.nlm.nih.gov/genomes/genbank/vertebrate\\_mammalian/](ftp://ftp.ncbi.nlm.nih.gov/genomes/genbank/vertebrate_mammalian/) or "Ensembl genome browser":<http://www.ensembl.org> , as well as "NCBI's BLAST programs":<ftp://ftp.ncbi.nlm.nih.gov/blast/> and "Ensembl genome browser's BLAST or BLAT web tools":<http://www.ensembl.org> . In analyses of gene features, the protocol used direct evidence of gene annotations available in "NCBI's nr, est\_human, est\_mouse and est\_others databases":<https://www.ncbi.nlm.nih.gov> . The protocol established **\*\*tests of reliability of eutherian public genomic sequences\*\*** that used potential coding sequences. The first test steps analysed nucleotide sequence coverages of each potential coding sequence, using BLASTN program and primary experimental genomic sequence information deposited in "NCBI's Trace Archive":<https://www.ncbi.nlm.nih.gov/Traces/trace.cgi> . The second test steps described potential coding sequences as complete coding sequences only if consensus trace coverages were available for every nucleotide in each potential coding sequence. Alternatively, the protocol described potential coding sequences as putative coding sequences \ (not used in analyses). The protocol used complete coding sequences in all analyses, and deposited them in "European Nucleotide Archive":<https://www.ebi.ac.uk/ena> as curated "third party data gene data sets":<https://www.ebi.ac.uk/ena/about/tpa-policy> . The guidelines of "human gene nomenclature":<http://www.genenames.org/about/guidelines> and "mouse gene nomenclature":<http://www.informatics.jax.org/mgihome/nomen/gene.shtml> were used in revisions and updates of gene classifications. In multiple pairwise genomic sequence alignments, the protocol used "mVISTA's AVID":<http://genome.lbl.gov/vista/index.shtml> . In base sequences used in multiple pairwise genomic sequence alignments, transposable elements were masked by

"RepeatMasker":<http://www.repeatmasker.org/> . Finally, the pairwise nucleotide sequence identities of predicted promoter regions calculated using "BioEdit":<http://www.mbio.ncsu.edu/BioEdit/bioedit.html> were used in statistical analyses \ (Microsoft Office Excel). **\*\*Phylogenetic analysis\*\*** In phylogenetic analyses, the protocol included protein and nucleotide sequence alignments, calculations of phylogenetic trees and calculations of pairwise nucleotide sequence identity patterns. First, the complete coding sequences were translated using "BioEdit":<http://www.mbio.ncsu.edu/BioEdit/bioedit.html> , and then aligned at amino acid level using ClustalW implemented in "BioEdit":<http://www.mbio.ncsu.edu/BioEdit/bioedit.html> . After inspections and manual corrections of protein primary sequence alignments, the protocol prepared nucleotide sequence alignments. The "MEGA":<http://www.megasoftware.net> was used in phylogenetic tree calculations. The protocol used neighbour-joining, minimum evolution, maximum parsimony and unweighted pair group method with arithmetic mean methods. The pairwise nucleotide sequence identities of nucleotide sequence alignments calculated using "BioEdit":<http://www.mbio.ncsu.edu/BioEdit/bioedit.html> were used in statistical analyses \ (Microsoft Office Excel). For each nucleotide sequence alignment, the protocol calculated average pairwise identities and their average absolute deviations, as well as largest pairwise identities and smallest pairwise identities. Finally, the protocol discriminated between eutherian major gene clusters including and not including evidence of differential gene expansions. **\*\*Protein molecular evolution analysis\*\*** The protocol established **\*\*tests of protein molecular evolution\*\*** integrating patterns of nucleotide sequence similarities with protein primary structures. The protein and nucleotide sequence alignments were used in tests. First, for each nucleotide sequence alignment, the protocol calculated codon usage statistics using "MEGA":<http://www.megasoftware.net> . The ratios between observed and expected amino acid codon counts determined relative synonymous codon usage statistics \ ( $R_c$ ). Therefore, in nucleotide sequence alignments, the amino acid codons including  $R_c \leq 0.7$  were designated as not preferable amino acid codons. The protocol then described reference protein sequence amino acid sites as invariant amino acid sites \ (invariant alignment positions), forward amino acid sites \ (variant alignment positions that did not include amino acid codons with  $R_c \leq 0.7$ ) or compensatory amino acid sites \ (variant alignment positions that included amino acid codons with  $R_c \leq 0.7$ ), using protein and nucleotide sequence alignments.

## Anticipated Results

The **\*\*eutherian comparative genomic analysis protocol\*\*** published 1293 complete coding sequences that included: "64 interferon- $\gamma$ -inducible GTPase genes":<https://www.ebi.ac.uk/ena/data/view/FR734011-FR734074> <sup>1</sup>:<https://dx.doi.org/10.1007/s10142-012-0291-2> , "255 ribonuclease A genes":<https://www.ebi.ac.uk/ena/data/view/HG328835-HG329089> <sup>2</sup>:<https://dx.doi.org/10.1007/s00438-013-0801-5> , "119 Mas-related G protein-coupled receptor genes":<https://www.ebi.ac.uk/ena/data/view/HG426065-HG426183> <sup>3</sup>:<https://dx.doi.org/10.1016/j.gene.2014.02.049> , "116 lysozyme genes":<https://www.ebi.ac.uk/ena/data/view/HG931734-HG931849> <sup>4</sup>:<https://dx.doi.org/10.1016/j.gdata.2014.08.003> , "128 adenohipophysis cystine-knot

genes":<https://www.ebi.ac.uk/ena/data/view/HF564658-HF564785>

"5":<https://dx.doi.org/10.1016/j.mgene.2015.02.005> , "30 D-dopachrome tautomerase and macrophage migration inhibitory factor genes":<https://www.ebi.ac.uk/ena/data/view/HF564786-HF564815>

"5":<https://dx.doi.org/10.1016/j.mgene.2015.02.005> , "100 growth hormone genes":<https://www.ebi.ac.uk/ena/data/view/LM644135-LM644234>

"6":<https://dx.doi.org/10.1016/j.gdata.2015.09.007> , "211 tumor necrosis factor ligand genes":<https://www.ebi.ac.uk/ena/data/view/LN874312-LN874522>

"8":<https://dx.doi.org/10.1007/s00251-015-0887-5> , "149 globin genes":<https://www.ebi.ac.uk/ena/data/view/LT548096-LT548244>

"9":<https://dx.doi.org/10.1016/j.genrep.2016.10.009> and "121 kallikrein genes":<https://www.ebi.ac.uk/ena/data/view/LT631550-LT631670>

"10":<https://dx.doi.org/10.1016/j.ymgmr.2017.01.009> . These curated eutherian third party data gene data sets were applicable in gene annotations and genome analyses, phylogenetic analyses and analyses of protein molecular evolution.

## References

1. Premzl M \ (2012) Comparative genomic analysis of eutherian interferon- $\gamma$ -inducible GTPases. *Funct Integr Genomics* 12:599-607 "DOI: 10.1007/s10142-012-0291-2":<https://dx.doi.org/10.1007/s10142-012-0291-2>
2. Premzl M \ (2014) Comparative genomic analysis of eutherian ribonuclease A genes. *Mol Genet Genomics* 289:161-167 "DOI: 10.1007/s00438-013-0801-5":<https://dx.doi.org/10.1007/s00438-013-0801-5>
3. Premzl M \ (2014) Comparative genomic analysis of eutherian Mas-related G protein-coupled receptor genes. *Gene* 540:16-19 "DOI: 10.1016/j.gene.2014.02.049":<https://dx.doi.org/10.1016/j.gene.2014.02.049>
4. Premzl M \ (2014) Third party annotation gene data set of eutherian lysozyme genes. *Genom Data* 2:258-260 "DOI: 10.1016/j.gdata.2014.08.003":<https://dx.doi.org/10.1016/j.gdata.2014.08.003>
5. Premzl M \ (2015) Initial description of primate-specific cystine-knot Prometheus genes and differential gene expansions of D-dopachrome tautomerase genes. *Meta Gene* 4:118-128 "DOI: 10.1016/j.mgene.2015.02.005":<https://dx.doi.org/10.1016/j.mgene.2015.02.005>
6. Premzl M \ (2015) Third party data gene data set of eutherian growth hormone genes. *Genom Data* 6:166-169 "DOI: 10.1016/j.gdata.2015.09.007":<https://dx.doi.org/10.1016/j.gdata.2015.09.007>
7. Premzl M \ (2016) Curated eutherian third party data gene data sets. *Data Brief* 6:208-213 "DOI: 10.1016/j.dib.2015.11.056":<https://dx.doi.org/10.1016/j.dib.2015.11.056>
8. Premzl M \ (2016) Comparative genomic analysis of eutherian tumor necrosis factor ligand genes. *Immunogenetics* 68:125-132 "DOI: 10.1007/s00251-015-0887-5":<https://dx.doi.org/10.1007/s00251-015-0887-5>
9. Premzl M \ (2016) Comparative genomic analysis of eutherian globin genes. *Gene Rep* 5:163-166 "DOI: 10.1016/j.genrep.2016.10.009":<https://dx.doi.org/10.1016/j.genrep.2016.10.009>
10. Premzl M \ (2017) Comparative genomic analysis of eutherian kallikrein genes. *Mol Genet Metab Rep* 10:96-99 "DOI: 10.1016/j.ymgmr.2017.01.009":<https://dx.doi.org/10.1016/j.ymgmr.2017.01.009>

## Acknowledgements

The **eutherian comparative genomic analysis protocol** was established under research project "Comparative genomic analysis of eutherian genes":<https://www.ncbi.nlm.nih.gov/myncbi/mpremzl/cv/130205/>. **Project leader:** Marko Premzl PhD, ANU Alumni, 4 Kninski trg Sq., Zagreb, Croatia **E-mail contact:** Marko.Premzl@alumni.anu.edu.au

## Figures

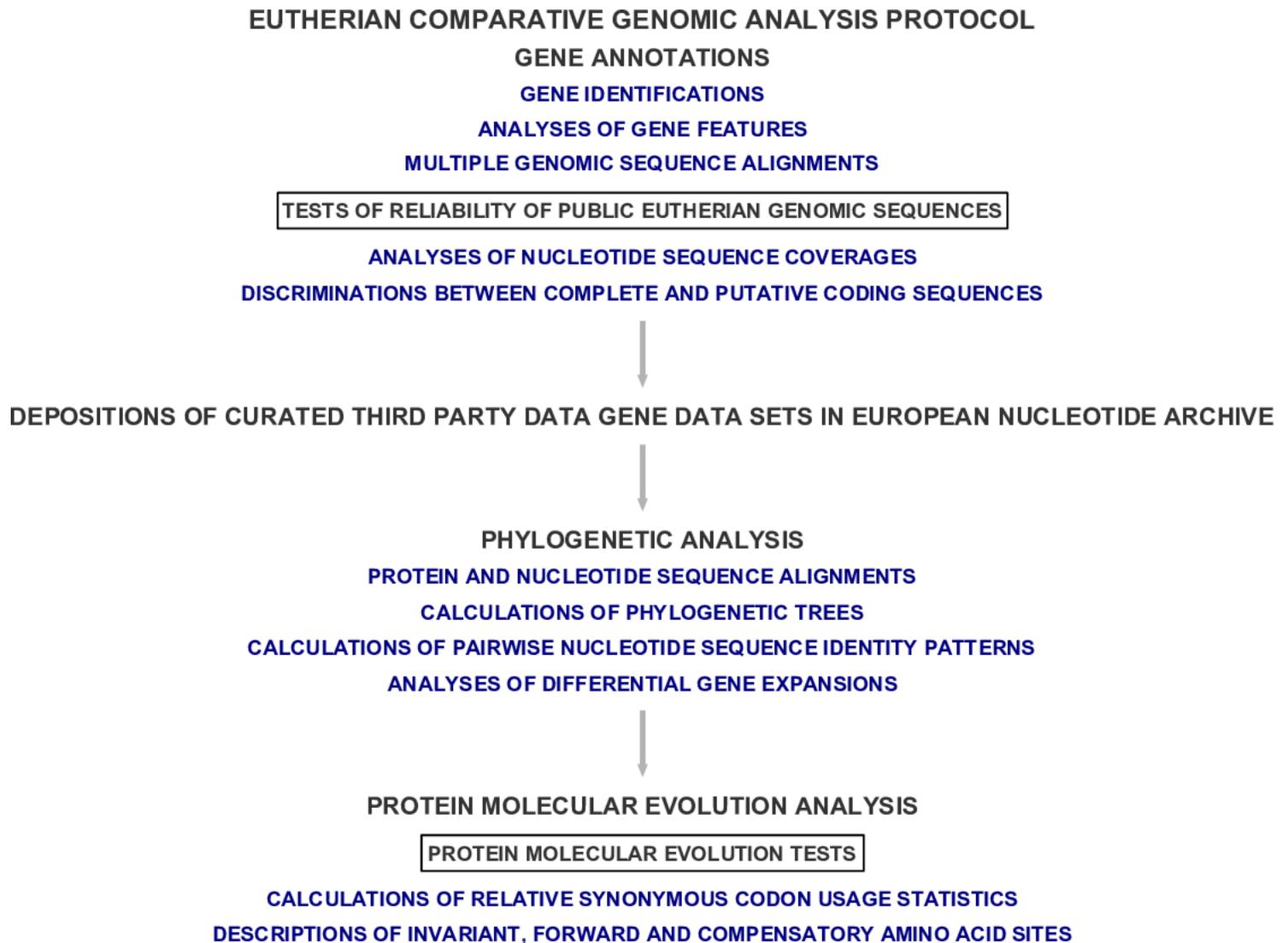


Figure 1

Figure 1 Eutherian comparative genomic analysis protocol flowchart.