

## Supplementary material

### **Brain transcriptome sequencing of a natural model of Alzheimer's disease**

Francisco Altimiras<sup>1,2,\*,#</sup>, Barbara Uszczynska-Ratajczak<sup>3,4,\*</sup>, Francisco Camara<sup>3,4</sup>, Anna Vlasova<sup>3,4</sup>, Emilio Palumbo<sup>3,4</sup>, Stephen Newhouse<sup>5</sup>, Robert M. J. Deacon<sup>6,7</sup>, Leandro A. E. Farias<sup>1</sup>, Michael J. Hurley<sup>8</sup>, David E. Loyola<sup>9</sup>, Rodrigo A. Vásquez<sup>10</sup>, Richard Dobson<sup>5</sup>, Roderic Guigó<sup>3,4,#</sup>, and Patricia Cogram<sup>6,7,#</sup>.

<sup>1</sup>Facultad de Ingeniería y Ciencias, Universidad Adolfo Ibañez, Santiago, Chile

<sup>2</sup>Telefonica Research and Development, Santiago, Chile

<sup>3</sup>Centre for Genomic Regulation, Barcelona Institute of Science and Technology, Barcelona, Spain

<sup>4</sup>Universitat Pompeu Fabra, Barcelona, Spain

<sup>5</sup> Institute of Psychiatry, Psychology & Neuroscience, King's College London, London, UK

<sup>6</sup> Laboratory of Molecular Neuropsychiatry, Institute of Cognitive and Translational Neuroscience (INCYT), INECO Foundation, Favaloro University, National Scientific and Technical Research Council (CONICET), Buenos Aires, Argentina

<sup>7</sup> GeN.DDI Ltd, London, UK

<sup>8</sup> Division of Brain Sciences, Centre for Neuroinflammation and Neurodegeneration, Imperial College, London, UK

<sup>9</sup> National Center for Genomics and Bioinformatics, Santiago, Chile

<sup>10</sup> Faculty of Sciences, Institute of Ecology and Biodiversity, Universidad de Chile, Santiago, Chile

\*Shared first authorship

#Corresponding authors

### Supplementary tables

**Table S1. Accuracy of gene prediction.** Accuracy of gene prediction on an *O. degus* “artificial scaffold” consisting of 238 concatenated *O. degus* test sequences (with 800 nucleotides of sequence between each of the gene models) using the ab initio programs geneid, AUGUSTUS and SNAP with their pre-existing Mammalian/*H. sapiens* parameter files (i.e. “mam/hs”). The accuracy of SGP2 (homology evidence-based prediction tool that used the genome of *H. sapiens* as reference) and that of Augustus (using RNASeq and transcript evidence i.e. “AUGUSTUS+hints”) were also tested for accuracy on the same set of sequences. Geneid (geneid+introns) and SGP2 (SGP2+introns) using introns as external evidence were also evaluated. (SN & SP: sensitivity & specificity at nucleotide level; SNe & SPe: sensitivity & specificity at exon level; SNg & SPg: sensitivity & specificity at gene level).

Program/param	SN	SP	SNe	SPe	SNg	SPg
Geneid <b>mam/hs</b>	0.83	0.75	0.65	0.69	0.09	0.06
Geneid+intron <b>mam/hs</b>	0.92	0.82	0.83	0.79	0.24	0.17
SGP2 <b>odegus / Hs</b>	0.90	0.82	0.77	0.73	0.12	0.08
SGP2+intron <b>odegus / mam/hs</b>	0.95	0.86	0.86	0.79	0.26	0.17
Augustus+hints <b>mam/hs</b>	0.87	0.94	0.81	0.90	0.33	0.35
Augustus <b>mam/hs</b>	0.81	0.84	0.68	0.75	0.06	0.07
SNAP <b>mam/hs</b>	0.83	0.45	0.57	0.32	0.03	0.01

**Table S2. Weights used by EVM to create a consensus CDS model *O. degus*.** The shortcuts interpretation: SPLAN2 uniprot90: SPLAN2 search against Uniprot90 proteins; SPALN2 uniprot-swissprot: SPALN2 against rodent uniprot/swissprot curated proteins; Exonerate uniprot-swissprot: exonerate against against rodent uniprot/swissprot curated proteins;

Type	Source	Weight
ABINITIO_PREDICTION	Augustus	1
ABINITIO_PREDICTION	AugustusHints	1.75
ABINITIO_PREDICTION	geneid	1
ABINITIO_PREDICTION	SGP2	1.25
ABINITIO_PREDICTION	geneid+introns	1.5
ABINITIO_PREDICTION	SGP2+introns	1.75
ABINITIO_PREDICTION	SNAP	0.3
PROTEIN	SPALN2 uniprot90	5
PROTEIN	SPALN2 uniprot-swissprot	4
PROTEIN	exonerate uniprot-swissprot	4
TRANSCRIPT	PASA	10

**Table S3. Comparison between EVM-based and GNOMON-based protein coding gene annotation.** Statistics for two protein-coding annotations for *O. degus*.

Annotation versions	<i>O. degus 2a (EVM-generated)</i>	<i>O. degus (ncbi GNOMON) -protein-coding only-</i>
Genome length (Mbases)	<b>2,995.89</b>	
number of scaffolds	<b>7,134</b>	
Number of protein-coding genes	<b>31,739</b>	<b>20,779</b>
Gene density (genes/Kbase)	<b>0.0106</b>	<b>0.007</b>
Number of protein-coding transcripts	<b>36,866</b>	<b>26,248</b>
Transcripts/gene (range) (% genes with more than 1 transcript)	<b>1.16 (SD 0.72) (1 – 32) (9.24%)</b>	<b>1.26 (SD 0.94) (1 – 31)(15%)</b>
Number of transcripts with UTRs	<b>10,648</b>	-
Number of proteins	<b>36,575</b>	<b>26,248</b>
Number of complete proteins (%)	<b>33,858 (92.57%)</b>	-
Number/(%) proteins with similarity to sequences in the NCBI NR database (E=10 <sup>-2</sup> ; min. identity=25%)	<b>35,475 (97%)</b>	-
Avg. length of proteins (range)	<b>461.96 aa. (SD 593.73) (25 – 34,458)</b>	<b>577.56 aa. (SD 641.03) (23 – 34,357)</b>
Avg. length of full-length proteins (range)	<b>478.57 aa. (SD 602.27) (25 – 34,458)</b>	-
Number of partial proteins (not starting with "M")	<b>1842 (5.04%)</b>	<b>259 (0.98%)</b>
Avg. length of partial proteins (not starting with "M")	<b>253.11 aa. (SD 431.8)</b>	-
Number of partial proteins (no terminal STOP codon)	<b>1589 (4.34%)</b>	<b>(can't determine as gnomon protein set has no clear STOP signal)</b>
Avg. length of partial proteins (no	<b>213.91 aa. (SD</b>	-

terminal STOP codon)	<b>350.85)</b>	
Number of partial proteins (not starting with an M -and- no terminal STOP codon)	<b>714 (1.95%)</b>	-
Avg. length of partial proteins (not starting with an M -and- no terminal STOP codon)	<b>158.83 aa. (SD 261.08)</b>	-
Number of partial proteins (not starting with an M -or- no terminal STOP codon)	<b>2,717 (7.43%)</b>	-
Avg. length of partial proteins (not starting with an M -or- no terminal STOP codon)	<b>254.96 aa. (SD 423.14)</b>	-
Number of protein-coding exons	<b>288,884</b>	<b>268,660</b>
Number of introns	<b>252,018</b>	<b>242,412</b>
Number of UTRs (spliced)	<b>19,003</b>	-
Number of single-exon genes	<b>10,114</b>	<b>3,156</b>
Number of multi-exonic transcripts (genes)	<b>26,752 (21,740)</b>	<b>23,092 (17,623)</b>
Exons/transcript (range) (excludes single-exon genes)	<b>10.42 (SD 10.50) (2 – 313)</b>	<b>11.49 (SD 10.35) (2 – 313)</b>
Introns/transcript (range)	<b>9.42 (SD 10.50) (1 – 312)</b>	<b>10.49 (SD 10.35) (1 – 312)</b>
“spliced” UTRs/transcript (range)	<b>1.785 (SD 0.74) (1 - 5)</b>	-
Avg. length of introns (range)	<b>5,998 (SD 19,994.1) (21 – 734,060)</b>	<b>5,613.03 (SD 19,909.6) (30 – 1,116,408)</b>
Avg. length of mono-exonic genes	<b>519.27 (SD 430.56)</b>	<b>872.88 (SD 618.70)</b>
Avg. length of exons (excludes mono-exonic genes)	<b>165.37 (SD 233.34)</b>	<b>161.25 (SD 230.37)</b>
Avg. length of first exons	<b>230.78 (SD 336.07)</b>	-
Avg. length of internal exons	<b>149.24 (SD 194.59)</b>	-
Avg. length of terminal exons	<b>235.72 (SD 352.41)</b>	-
Avg. length of CDS (range)	<b>1,392.9 (SD 1,782.42) (75 –</b>	<b>1,736.06 (SD 1,923.46)</b>

	<b>103,074)</b>	<b>(69 – 103,074)</b>
Avg. length of UTRs (range)	<b>653.40 (SD 942.07)</b> <b>(1 - 11,857)</b>	-
Avg. length of primary transcripts	<b>43,714.8 (SD 107,530)</b>	<b>56,055.1 (SD 117,349)</b>
G+C content exonic (mono-exonic genes)	<b>49.72% (SD 7.63%)</b>	<b>51.85% (SD 8.56%)</b>
G+C content exonic (excludes mono-exonic genes)	<b>52.53% (SD 7.47%)</b>	<b>53.52% (SD 7.42%)</b>
G+C content exonic (first exons)	<b>53.62% (SD 10.85%)</b>	-
G+C content exonic (internal exons)	<b>51.27% (SD 9.67%)</b>	-
G+C content exonic (terminal exons)	<b>53.59% (SD 10.84%)</b>	-
G+C content intronic	<b>45.05% (SD 11.54%)</b>	<b>45.45% (SD 11.61%)</b>
G+C content genomic	<b>40.16% (SD 5.63%)</b>	
G+C content UTRs	<b>53,76% (SD 5%)</b>	-

**Table S4. Non-default parameters for RNA-seq mappings.** Non-default parameters used during mapping step of pair-end reads of human brain AD subjects and control samples with STAR 2.4.0.1. First column refers to the name of the parameter, while the second to its value.

<b>Parameter</b>	<b>Value</b>
outSAMunmapped	Within
outFilterType	BySJout
outFilterMultimapNmax	20
outFilterMismatchNmax	999
outFilterMismatchNoverReadLmax	0.04
alignIntronMin	20
alignIntronMax	1000000
alignSJDBoverhangMin	1
readFilesCommand	zcat

**Supplementary Table S5. Functional annotation statistics.** Abbreviation KO – KEGG orthology groups.

Number of proteins/genes	36,575 / 31,739
<b>Annotated proteins/genes</b>	<b>34,571 (94.5%) / 30,336 (95.5%)</b>
Proteins with Interpro signatures	33,800 (92.4%)
Proteins with Blast2GO or KEGG definition	23,936 (65.4%)
Proteins with Blast2GO definition	16,737 (45.7%)
Proteins with KEGG definition	14,756 (40.3%)
Proteins assigned to KO groups	14,879 (40.6%)
Proteins with GO terms association	28,988 (79.2%)
Conserved domains signatures	31,874 (87.1%)
Conserved features signatures	15,017 (41%)



**Supplementary Table S6. GO term annotation.** **A.** Number of GO terms associated to each ontology **B.** Top 10 GO terms more frequently associated to proteins grouped by GO term type.

**A.**

<b>Term type</b>	<b>Number of proteins</b>
Biological process	22,200
Cellular component	20,110
Molecular function	27,220
<b>All</b>	<b>29,847</b>

**B**

<b>GO term id</b>	<b>GO term description</b>	<b># Proteins</b>
	<b>Biological process</b>	
GO:0006412	translation	2938
GO:0006355	regulation of transcription, DNA-templated	1865
GO:0007186	G-protein coupled receptor signaling pathway	1505
GO:0006414	translational elongation	1073
GO:0055114	oxidation-reduction process	1016
GO:0006413	translational initiation	904
GO:0007165	signal transduction	902
GO:0006468	protein phosphorylation	882
GO:0000184	nuclear-transcribed mRNA catabolic process, nonsense-mediated decay	830
GO:0006614	SRP-dependent cotranslational protein targeting to membrane	823
	<b>Molecular function</b>	
GO:0005515	protein binding	6905
GO:0003735	structural constituent of ribosome	3023
GO:0005524	ATP binding	2431
GO:0003676	nucleic acid binding	2166
GO:0008270	zinc ion binding	1982
GO:0000166	nucleotide binding	1785
GO:0003677	DNA binding	1755
GO:0046872	metal ion binding	1706
GO:0004930	G-protein coupled receptor activity	1311
GO:0003723	RNA binding	1280
	<b>Cellular component</b>	
GO:0005622	intracellular	3658
GO:0005634	nucleus	3517
GO:0016021	integral component of membrane	3502
GO:0005840	ribosome	2927

GO:0005737	cytoplasm	2338
GO:0016020	membrane	2165
GO:0005829	cytosol	1641
GO:0005886	plasma membrane	1487
GO:0005730	nucleolus	1429
GO:0022625	cytosolic large ribosomal subunit	1007

**Table S7. GO-terms enrichment for differentially expressed genes in *O. degus*.** Biological processes overrepresented by up- and down-regulated genes identified in *O. degus* brain samples. GO terms shown are those significantly overrepresented (pvalue < 0.05) by genes showing differential expression between AD-like subjects and controls. Categories are sorted by p-value.

**TABLE IN ADDITIONAL FILE 2**

**Table S8. GO terms enrichment for differentially expressed genes in human samples.** Biological processes overrepresented by up- and down-regulated genes identified in human brain samples. GO terms shown are those significantly overrepresented (pvalue < 0.05) by genes showing differential expression between AD subjects and controls. Categories are sorted by p-value.

**TABLE IN ADDITIONAL FILE 2**

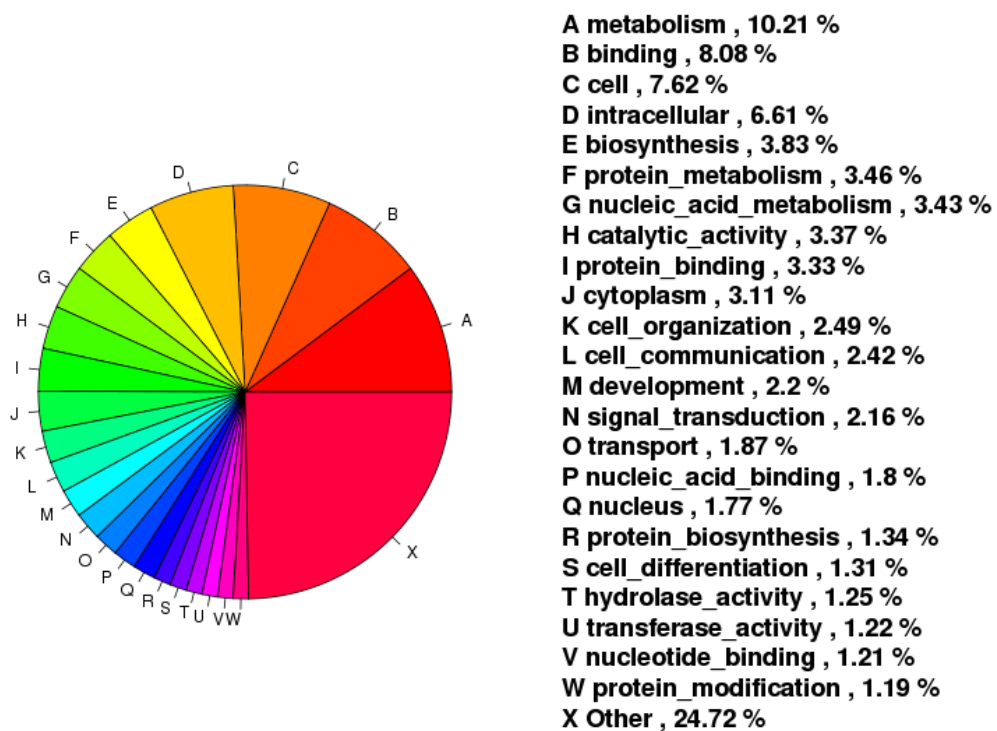
**Table S9. Genes differentially expressed between human samples.** The complete list of 2963 human genes displaying differential expression between AD subjects and controls. Genes are ranked by FDR (FDR < 0.05).

**TABLE IN ADDITIONAL FILE 2**

**Table S10. Pairwise comparisons to measure the gene expression level differences between AD-like degus and controls.**

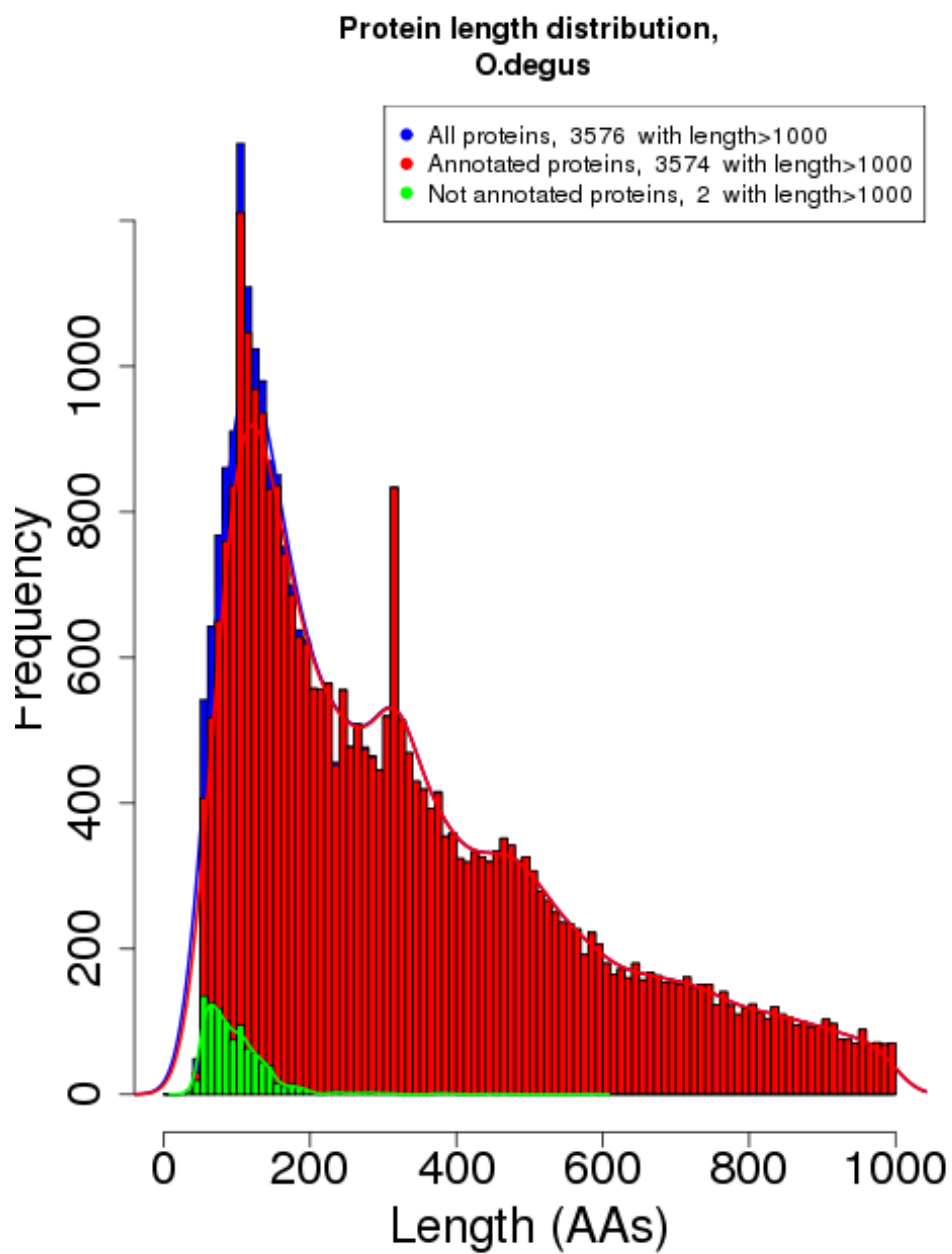
**TABLE IN ADDITIONAL FILE 2**

## Supplementary figures



**Figure S1. Gene ontology mapping of *O. degus* genes.** The GO terms were mapped into the general GO slim without top level categories – biological process, molecular function, cellular component.





**Figure S2. Distribution of functionally annotated and non-annotated proteins.** Number of annotated and non-annotated sequences in relation to their length. The blue color correspond to all proteins, red – annotated proteins and green to the non-annotated proteins.