

1 **Title:** *De novo* assembly of the Mongolian gerbil genome and transcriptome

2 **Authors:** Shifeng Cheng^{1,2*}, Yuan Fu^{1,3*}, Yaolei Zhang^{1,3}, Wenfei Xian^{1,2}, Hongli Wang^{1,3}, Benedikt
3 Grothe⁴, Xin Liu^{1,3}, Xun Xu^{1,3}, Achim Klug⁵, Elizabeth A McCullagh^{5#}

4

5 **Institutional Affiliations:** ¹BGI-Shenzhen, Beishan Industrial Zone, Yantian District, Shenzhen 518083,
6 China;

7 ²Agricultural Genome Institute at Shenzhen, Chinese Academy of Agricultural Sciences, Shenzhen 518124,
8 China;

9 ³China National GeneBank, BGI-Shenzhen, Shenzhen, 518083, China;

10 ⁴Division of Neurobiology, Ludwig-Maximilians-Universitaet Munich, Planegg-Martinsried 82152,
11 Germany

12 ⁵Department of Physiology and Biophysics, School of Medicine, University of Colorado Denver, Aurora,
13 CO, 80045, USA

14 *co first authors

15 #corresponding author

16

17

18 **Email Addresses for Authors:**

19 Shifeng Cheng: chengshifeng@caas.cn

20

21 Yuan Fu: fuyuan@genomics.cn

22

23 Yaolei Zhang: zhangyaolei@genomics.cn

24

25 Wenfei Xian: xianwenfei@caas.cn

26

27 Hongli Wang: wanghongli@genomics.cn

28

29 Benedikt Grothe: grothe@lmu.de

30

31 Xin Lu: liuxig@genomics.cn

32

33 Xun Xu: xuxun@genomics.cn

34

34 Achim Klug: achim.klug@ucdenver.edu

35

35 Elizabeth A McCullagh: elizabeth.mccullagh@ucdenver.edu

36

36

37

38 **Abstract:** (250 words)

39 **BACKGROUND:** The Mongolian gerbil (*Meriones unguiculatus*) has historically been used as a model
40 organism for the auditory and visual systems, stroke/ischemia, epilepsy and aging related research since
41 1935 when laboratory gerbils were separated from their wild counterparts. In this study we report genome
42 sequencing, assembly, and annotation further supported by transcriptome data from 27 different tissues
43 samples.

44 **FINDINGS:** The genome was assembled using Illumina HiSeq 2000 and resulted in a final genome size of
45 2.54 Gbp with contig and scaffold N50 values of 31.4 Kbp and 500.0 Kbp, respectively. Based on the k-mer
46 estimated genome size of 2.48 Gbp, the assembly appears to be complete. The genome annotation was
47 supported by transcriptome data that identified 36 019 predicted protein-coding genes across 27 tissue
48 samples. A BUSCO search of 3023 mammalian groups resulted in 86% of curated single copy orthologs
49 present among predicted genes, indicating a high level of completeness of the genome.

50 **CONCLUSIONS:** We report a *de novo* assembly of the Mongolian gerbil genome that was further
51 enhanced by annotation of transcriptome data from several tissues. Sequencing of this genome increases the
52 utility of the gerbil as a model organism, opening the availability of now widely used genetic tools.

53

54 **Keywords:** Gerbil genome, *Meriones unguiculatus*, transcriptome, model organism

55

56

57

58

59

60

61

62 DATA DESCRIPTION

63 Background information on *Meriones unguiculatus*

64 The Mongolian gerbil is a small rodent that is native to Mongolia, southern Russia, and northern China.
65 Laboratory gerbils used as model organisms originated from 20 founders captured in Mongolia in 1935 [1].
66 Gerbils have been used as model organisms for sensory systems (visual and auditory) and pathologies
67 (aging, epilepsy, irritable bowel syndrome and stroke/ischemia). The gerbil's hearing range covers the
68 human audiogram while also extending into ultrasonic frequencies, making gerbils a better model than rats
69 or mice to study lower frequency human-like hearing [2]. In addition to the auditory system, the gerbil has
70 also been used as a model for the visual system because gerbils are diurnal and therefore have more cone
71 receptors than mice or rats making them a closer model to the human visual system [3]. The gerbil has also
72 been used as a model for aging due to its ease of handling, prevalence of tumors, and experimental stroke
73 manipulability [1,4]. Interestingly, the gerbil has been used as a model for stroke and ischemia due to
74 variations in the blood supply to the brain due to an anatomical region known as the "Circle of Willis" [5].
75 In addition, the gerbil is a model for epileptic activity as a result of its natural minor and major seizure
76 propensity when exposed to novel stimuli [6,7]. Lastly, the gerbil has been used as model for inflammatory
77 bowel disease, colitis, and gastritis due to the similarity in the pathology of these diseases between humans
78 and gerbils [8,9]. Despite its usefulness as a model for all these systems and medical conditions, the utility
79 of the gerbil as a model organism has been limited due to a lack of a sequenced genome to manipulate. This
80 is especially the case with the increased use of genetic tools to manipulate model organisms.

81 Here we describe a *de novo* assembly and annotation of the Mongolian gerbil genome and transcriptome.
82 Recently, a separate group has sequenced the gerbil genome, however our work is further supported by
83 comparisons with an in-depth transcriptome analysis [10]. RNA-seq data were produced from 27 tissues that
84 were used in the genome annotation and deposited in the NCBI SRA database under the project _____. These

85 data provide a draft genome sequence to facilitate the continued use of the Mongolian gerbil as a model
86 organism and to help broaden the genetic rodent models available to researchers.

87 **Animals and Genome Sequencing**

88 All experiments complied with all applicable laws, NIH guidelines, and were approved by the University of
89 Colorado IACUC. Five young adult (postnatal day 65-71) gerbils (three males and two females) were used
90 for tissue RNA transcriptome analysis and DNA genome assembly. In addition, two old (postnatal day 1013
91 or 2.7 years) female gerbil's tissue was used for transcriptome analysis.

92 Genomic DNA was extracted from young adult animal tail and ear snips using a commercial kit (DNeasy
93 Blood and Tissue Kit, Qiagen, Venlo, Netherlands). We then used the extracted DNA to create different
94 pair-end insert libraries of 250 bp, 350 bp, 500 bp, 800 bp, 2 Kb, 4 Kb, 6 Kb, and 10 Kb. These libraries
95 were then sequenced using an Illumina HiSeq2000 Genome Analyzer (Illumina, San Diego, CA, USA)
96 generating a total of 322.13 Gb in raw data, from which a total of 287.4 Gb of 'clean' data was obtained
97 after removal of duplicates, contaminated reads, and low-quality reads.

98 **Assembly**

99 The gerbil genome was estimated to be approximately 2.48 Gbp using a k-mer-based approach. High-
100 quality reads were then used for genome assembly using the SOAPdenovo (version 2.04) package. The final
101 assembly had a total length of 2.54 Gb and was comprised of 31,769 scaffolds assembled from 114,522
102 contigs. The N50 sizes for contigs and scaffolds were 31.4 Kbp and 500.0 Kbp, respectively (Table 1).
103 Given the genome size estimate of 2.48 Gbp, genome coverage by the final assembly was likely complete
104 and is consistent with the previously published gerbil genome, which had a total length of 2.523 Gbp [10].
105 Completeness of the genome assembly was confirmed by successful mapping of the RNA-seq assembly
106 back to the genome showing that 98% of the RNA-seq sequences can be mapped to the genome with >50%
107 sequence in one scaffold.

108 **Transcriptome Sequencing/Assembly/Annotation**

109 Gene expression data were produced to aid in the genome annotation process. Samples from 27 tissues were
110 collected from the seven gerbils described above (Supplementary Figure 1). The tissues were collected after
111 the animals were euthanized with isoflurane and stored on liquid nitrogen until homogenized with a pestle.
112 RNA was prepared using the RNeasy mini isolation kit (Qiagen, Venlo, Netherlands). RNA integrity was
113 analyzed using a Nanodrop Spectrophotometer (Thermo Fisher Waltham, MA, USA) followed by analysis
114 with an Agilent Technologies 2100 Bioanalyzer (Agilent Technologies, Santa Clara, CA, USA) and samples
115 with an RNA integrity number (RIN) value greater than 7.0 were used to prepare libraries which were
116 sequenced using an Illumina HiSeq2000 Genome Analyzer (Illumina, San Diego, CA, USA). The sequenced
117 libraries were assembled with Trinity (v2.0.6 parameters: "--min_contig_length 150 --min_kmer_cov 3 --
118 min_glue 3 --bfly_opts '-V 5 --edge-thr=0.1 --stderr'") generating 131,845 sequences with a total length of
119 130,734,893 bp. Quality of the RNA assembly was assessed by filtering RNA-seq reads using SOAPnuke
120 (v1.5.2 parameters: "-l 10 -q 0.1 -p 50 -n 0.05 -t 5,5,5,5") followed by mapping of clean reads to the
121 assembled genome using HISAT2 (v2.0.4) and StringTie (v1.3.0). The initial assembled genes were then
122 filtered using CD-HIT (v4.6.1) with sequence identity threshold of 0.9 followed by a homology search
123 (human, rat, mouse proteins) and TransDecoder (v2.0.1) open reading frame (ORF) prediction. The RNA-
124 seq assembly resulted in 19,737 protein-coding genes with a total length of 29.4 Mbp, which is available in
125 the NCBI Nucleotide, Protein and Gene resources database.

126 **Genome Annotation**

127 Genomic repeat elements of the genome assembly were also identified and annotated using RepeatMasker
128 (v4.0.5 RRID:SCR_012954)[11] and RepBase library (v20.04)[12]. In addition, we constructed a *de novo*
129 repeat sequence database using LTR-FINDER (v1.0.6) [13] and RepeatModeler (v1.0.8) [13] to identify
130 any additional repeat elements using RepeatMasker. A combination of both repeat element identification

131 approaches resulted in a total length of 1016.7 Mbp of the total *M. unguiculatus* genome as repetitive,
132 accounting for 40.0% of the entire genome assembly. The repeat element landscape of *M. unguiculatus*
133 consists of long interspersed elements (LINEs)(27.5%), short interspersed elements (SINEs)(3.7%), long
134 terminal repeats (LTRs)(6.5%), and DNA transposons (0.81%) (Table 2). This is consistent with other
135 rodent species including mouse [14] and rat [15].

136 Protein-coding genes were predicted and annotated by a combination of homology searching, *ab initio*
137 prediction (using AUGUSTUS (v3.1), GENSCAN (1.0), and SNAP (v2.0)), and RNA-seq data (using
138 TopHat (v1.2 with parameters: “-p 4 --max-intron-length 50000 -m 1 -r 20 --mate-std-dev 20 --closure-
139 search --coverage-search --microexon-search”) and Cufflinks (v2.2.1 [http://cole-trapnell-](http://cole-trapnell-lab.github.io/cufflinks/)
140 [lab.github.io/cufflinks/](http://cole-trapnell-lab.github.io/cufflinks/))) after repetitive sequences in the genome were masked using known repeat
141 information detected by RepeatMasker and RepeatProteinMask. Homology searching was performed using
142 protein data from *Homo Sapiens* (human), *Mus musculus* (mouse), and *Rattus norvegicus* (rat) from
143 Ensembl (v80) aligned to the masked genome using BLAT. Genewise (v2.2.0) was then used to improve the
144 accuracy of alignments and to predict gene models. The *de novo* gene predictions and homology-based
145 search were then combined using GLEAN. The GLEAN results were then integrated with the transcriptome
146 dataset using an in-house program (Table 3). This resulted in an identification of a total of 22,998 protein-
147 coding genes with an average transcript length of 23,846.58 bp. There were an average of 7.76 exons per
148 gene with an average length of 197.9 bp and average intron length of 3300.83 bp. The 22,998 protein-
149 coding genes were aligned to several protein databases to begin to identify their possible function.
150 InterProScan (v5.11) was used to align the final gene models to databases (ProDom, ProSiteProfiles,
151 SMART, PANTHER, PRINTS, Pfam, PIRSF, ProSitePatterns, SignalP_EUK, Phobius, IGRFAM, and
152 TMHMM) to detect consensus motifs and domains within these genes. Using the InterProScan results, we
153 obtained the annotations of the gene products from the Gene Ontology database. We then mapped these
154 genes to proteins in SwissProt and TrEMBL (Uniprot release 2015.04) using blastp with an E-value <1E-5.

155 We also aligned the final gene models to proteins in KEGG (release 76) to determine the functional
156 pathways for each gene (Table 4). This resulted in 20,760 protein-coding genes that had a functional
157 annotation, or 90.3% of the total gene set.

158 **Quality Assessment**

159 In addition to measuring standard assembly quality metrics, genome assembly and annotation quality were
160 further assessed by comparison with closely related species, gene family construction, evaluation of
161 housekeeping genes, and Benchmarking Universal Single-Copy Orthologs (BUSCO) search. The assembled
162 gerbil genome was compared with other closely related model organisms including mouse, rat, and hamster
163 (Table 5). The genomes from these species varied in size from 2.3 to 2.8 Gbp. The total number of
164 annotated proteins in gerbil (20,760) is most similar to mouse (22,598), followed by rat (23,347), and then
165 hamster (24,238). Gene family construction was performed using Treefam (<http://www.treefam.org/>)
166 (Figure 1). This analysis showed that single-copy orthologs in gerbil are similar to mouse and rat. To
167 examine housekeeping genes we downloaded 2169 human housekeeping genes from
168 (<http://www.tau.ac.il/~elieis/HKG/>) and extracted corresponding protein sequences to align to the gerbil
169 genome using blastp (v.2.2.26). We found there were 2141 genes consistent between human and gerbil
170 housekeeping genes (this is similar to rat (2153) and mouse (2146). Lastly, we employed BUSCO (v1.2) to
171 search 3023 mammalian groups. Of these groups, 86% complete BUSCO groups can be detected in the final
172 gene set. The presence of 86% complete mammalian BUSCO gene groups suggests a high level of
173 completeness of this gerbil genome assembly. A BUSCO search was also performed for the gerbil
174 transcriptome data resulting in detection of 82% complete BUSCO groups in the final transcriptome dataset
175 (Table 6). Based on the results from the quality metrics described above, we are confident of the quality of
176 the data for this assembly of the gerbil genome and transcriptome.

177

178 In summary, we report a fully annotated Mongolian gerbil genome sequence assembly enhanced by
179 transcriptome data from several different gerbils and tissues. The gerbil genome and transcriptome adds to
180 the availability of alternative rodent models that may be better models for diseases than rats or mice.
181 Additionally, the gerbil is an interesting comparative rodent model to mouse and rat since it has many traits
182 in common, but also differs in seizure susceptibility, low-frequency hearing, cone visual processing,
183 stroke/ischemia susceptibility, gut disorders and aging. Sequencing of the gerbil genome and transcriptome
184 opens these areas to molecular manipulation in the gerbil and therefore better models for specific disease
185 states.

186 **Availability of supporting data**

187 Genome annotation results are available at the NCBI, and supporting materials, which include transcripts
188 and genome assembly, are available at the *Gigascience* database, *GigaDB*.

189

190 **Additional files**

191 Additional file 1: Table S1 Tissues analyzed for RNA-seq data

192

193 **Abbreviations**

194 bp: base pair

195 BUSCO: Benchmarking Universal Single-Copy Orthologs

196 CDS: coding sequence

197 LINES: long interspersed elements

198 LTRs: long terminal repeats

199 Myr: million years

200 NCBI: National Center for Biotechnology Information

201 RefSeq: Reference sequence

202 RNA-seq: high-throughput messenger RNA sequencing

203 RIN: RNA integrity number

204 SINEs: short interspersed elements

205

206 **Competing Interests:** The authors declare that they have no competing interests.

207

208 **Funding**

209 **Authors' contributions**

210 **Acknowledgements**

211

212 **References**

213 1. Cheal ML. The gerbil: a unique model for research on aging. *Exp Aging Res.* 1986;12:3–21.

214 2. Ryan A. Hearing sensitivity of the mongolian gerbil, *Merionesunguiculatis*. *The Journal of the Acoustical*
215 *Society of America.* Acoustical Society of America; 1976;59:1222–6.

216 3. Govardovskii VI, Röhlich P, Szél A, Khokhlova TV. Cones in the retina of the Mongolian gerbil,
217 *Meriones unguiculatus*: an immunocytochemical and electrophysiological study. *Vision Res.* 1992;32:19–
218 27.

219 4. Vincent AL, Rodrick GE, Sodeman WA. *The Mongolian gerbil in aging research.* *Exp Aging Res.*
220 Routledge; 2007;6:249–60.

221 5. Small DL, Buchan AM. *Animal models.* Br Med Bull. Oxford University Press; 2000;56:307–17.

222 6. Bertorelli R, Adami M, Ongini E. The Mongolian gerbil in experimental epilepsy. *Ital J Neurol Sci.*
223 1995;16:101–6.

224 7. Löscher W. Genetic animal models of epilepsy as a unique resource for the evaluation of anticonvulsant
225 drugs. A review. *Methods Find Exp Clin Pharmacol.* 1984;6:531–47.

226 8. Bleich E-M, Martin M, Bleich A, Klos A. The Mongolian gerbil as a model for inflammatory bowel
227 disease. *Int J Exp Pathol.* Blackwell Publishing Ltd; 2010;91:281–7.

228 9. Hirayama F, Takagi S, Kusuhara H, Iwao E, Yokoyama Y, Ikeda Y. Induction of gastric ulcer and
229 intestinal metaplasia in mongolian gerbils infected with *Helicobacter pylori*. *J. Gastroenterol.* 1996;31:755–
230 7.

231 10. Zorio DAR, Monsma S, Sanes DH, Golding NL, Rubel EW, Wang Y. De novo sequencing and initial
232 annotation of the Mongolian gerbil (*Meriones unguiculatus*) genome. *Genomics.* Elsevier Inc; 2018.

- 233 11. Tarailo-Graovac M, Chen N. Using RepeatMasker to Identify Repetitive Elements in Genomic
234 Sequences. *Current Protocols in Bioinformatics*. Hoboken, NJ, USA: John Wiley & Sons, Inc;
235 2009;12:1269–4.10.14.
- 236 12. Jurka J, Kapitonov VV, Pavlicek A, Klonowski P, Kohany O, Walichiewicz J. Repbase Update, a
237 database of eukaryotic repetitive elements. *Cytogenet Genome Res*. 2005;110:462–7.
- 238 13. Benson G. Tandem repeats finder: a program to analyze DNA sequences. *Nucleic Acids Research*.
239 1999;27:573–80.
- 240 14. Smit AF. Interspersed repeats and other mementos of transposable elements in mammalian genomes.
241 *Current Opinion in Genetics & Development*. 1999;9:657–63.
- 242 15. Gibbs RA, Weinstock GM, Metzker ML, Muzny DM, Sodergren EJ, Scherer S, et al. Genome sequence
243 of the Brown Norway rat yields insights into mammalian evolution. *Nature*. 2004;428:493–521.
- 244
- 245 NCBI. NCBI Annotation Release Meriones unguiculatis, 2017;
246 https://www.ncbi.nlm.nih.gov/genome/annotation_euk/, (2 October 2017, date last accessed).
247
248
- 249
- 250
- 251
- 252
- 253
- 254
- 255
- 256
- 257
- 258
- 259
- 260
- 261
- 262
- 263

264 **Table 1 Global statistics of the Mongolian gerbil genome**

Statistic	Value
Size (Gb)	2.54
Scaffold number (>2000bp)	31769
Scaffold N50 (Kb)	500.0
Contig number (>2000bp)	114522
Contig N50 (Kb)	31.4

265

266 **Table 2 Summary of mobile element types**

Type	Length (Kb)	Percentage of the genome (%)
DNA	20,498	0.81
LINE	697,185	27.5
SINE	94,229	3.7
LTR	164,504	6.5
Other	40,254	1.6
Total	1,016,671	40.0

267

268

269

270

271

272

273

274

275

276

277 **Table 3 General statistics of predicted protein-coding genes**

	Gene set	Number	Average transcript length (bp)	Average CDS length (bp)	Average exon per gene	Average exon length (bp)	Average intron length (bp)
De novo	SNAP	76858	42227.63	742.83	5.52	134.62	9182.18
	AUGUSTUS	24675	19838.68	1133.22	5.61	201.97	4056.79
	GENESCAN	49390	24183.55	1023.1	6.25	163.54	4406.54
Homolog	<i>Mus musculus</i>	22728	26977.32	1465.18	8.02	182.61	3632.46
	<i>Rattus norvegicus</i>	23686	23564.96	1336.56	7.43	179.83	3455.8
	<i>Homo sapiens</i>	17131	31217.18	1580.27	9.11	173.55	3656.27
	GLEAN	19893	18835.39	1418.26	7.72	183.69	2691.49
	Transcriptome	36019	33752.29	1758.58	10.74	163.77	3285.43
	Final set	22998	23846.58	1535.48	7.76	197.9	3300.83

278

279 **Table 4 Functional annotation of the final gene set**

	Number	Percent (%)
Total	22,998	100
InterPro	18,570	80.7
GO	14,591	63.4
KEGG	17,572	76.4
Swissprot	20,113	87.5
TrEMBL	20,666	89.9
Annotated	20,760	90.3
Unannotated	2238	9.7

280

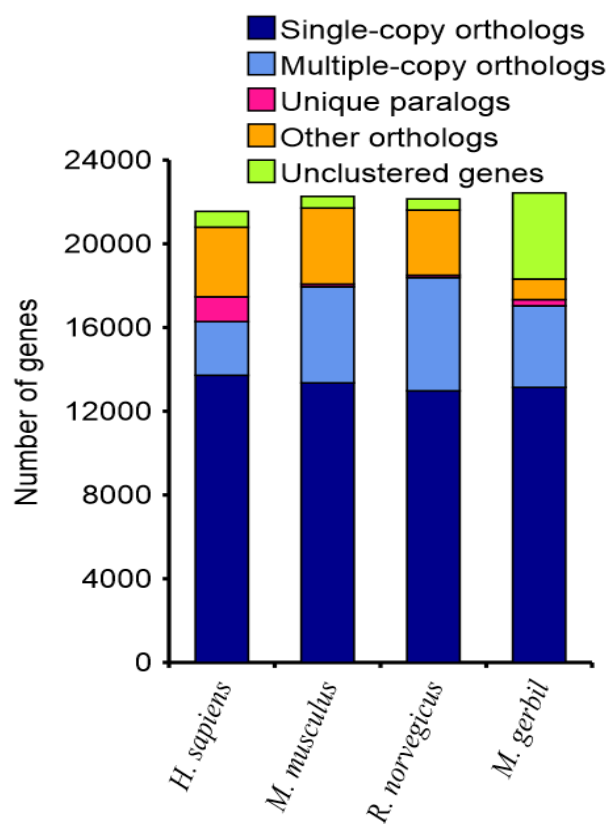
281

282

283

284 **Table 5 Genome annotation comparisons with other model organisms**

Species	Common name	Protein coding genes	Assembly Size	Divergence time to gerbils, Myr	RefSeq assembly accession	Annotation release ID
<i>Meriones unguiculatus</i>	Mongolian gerbil	20,760	2,537,533,819	--		
<i>Mus musculus</i>	mouse	22,598	2,818,974,548	22.5	GCF_000001635.26	106
<i>Rattus norvegicus</i>	rat	23,347	2,870,184,193	22.5	GCF_000001895.5	106
<i>Cricetulus griseus</i>	Chinese hamster	24,238	2,358,151,106	25	GCA_900186095.1	102



285

286 **Figure 1 Gene Family Construction.** The number of genes is similar between species compared (human,
287 mouse, rat, and gerbil.

288

289

290 **Table 6 Completeness of gerbil genome and transcriptome assembly as assessed by BUSCO**

	Genome	Transcriptome
Complete BUSCOs	2601	2508
Duplicated BUSCOs	55	46
Fragmented BUSCOs	170	293
Missing BUSCOs	252	222
Total BUSCO groups searched	3023	3023

291

292 **Supplementary Figure 1: Tissues sampled for RNA transcriptome**

Tissue	Run_accession	Sex	Age (postnatal day)	Data size (Mbp)
Lung		M	71	6733.54
Lung		F	1013	6347.26
Occipital lobe		F	1013	6231.73
Occipital lobe		F	70	5820.49
Kidney		F	1013	6412.73
Kidney		M	70	5609.90
Olfactory bulb		M	71	7467.99
Olfactory bulb		F	70	5576.19
Striatum		M	71	4596.98
Striatum		F	1013	5456.08
Striatum		M	71	6010.27
Striatum		F	71	8508.27
Cerebellum		F	1013	6021.12
Cerebellum		M	65	6724.73
Inferior colliculus		F	1013	5637.18
Inferior colliculus		M	71	6296.64

Liver	F	1013	5077.32
Liver	F	1013	6280.63
Spleen	M	71	9051.52
Spleen	F	1013	7943.03
Spleen	F	1013	6702.24
Frontal cortex	M	65	5895.65
Frontal cortex	F	1013	7202.13
Hippocampus	M	70	5189.69
Auditory brainstem	F	66	7332.74
Brainstem	M	65	5820.49
Parietal cortex	M	65	6786.95
