

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 Achim Klug: achim.klug@ucdenver.edu

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

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 The data sets supporting the results of this article are available in the China National GeneBank CNSA
57 repository, Accession id: CNP0000340.

58

59

60

61

62

63 DATA DESCRIPTION

64 Background information on *Meriones unguiculatus*

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

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

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

88 **Animals and Genome Sequencing**

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

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

99 **Assembly**

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

109 **Transcriptome Sequencing/Assembly/Annotation**

110 Gene expression data were produced to aid in the genome annotation process. Samples from 27 tissues were
111 collected from the seven gerbils described above (Supplementary Figure 1). The tissues were collected after
112 the animals were euthanized with isoflurane and stored on liquid nitrogen until homogenized with a pestle.
113 RNA was prepared using the RNeasy mini isolation kit (Qiagen, Venlo, Netherlands). RNA integrity was
114 analyzed using a Nanodrop Spectrophotometer (Thermo Fisher Waltham, MA, USA) followed by analysis
115 with an Agilent Technologies 2100 Bioanalyzer (Agilent Technologies, Santa Clara, CA, USA) and samples
116 with an RNA integrity number (RIN) value greater than 7.0 were used to prepare libraries which were
117 sequenced using an Illumina HiSeq2000 Genome Analyzer (Illumina, San Diego, CA, USA). The sequenced
118 libraries were assembled with Trinity (v2.0.6 parameters: "--min_contig_length 150 --min_kmer_cov 3 --
119 min_glue 3 --bfly_opts '-V 5 --edge-thr=0.1 --stderr'") generating 131,845 sequences with a total length of
120 130,734,893 bp. Quality of the RNA assembly was assessed by filtering RNA-seq reads using SOAPnuke
121 (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
122 assembled genome using HISAT2 (v2.0.4) and StringTie (v1.3.0). The initial assembled genes were then
123 filtered using CD-HIT (v4.6.1) with sequence identity threshold of 0.9 followed by a homology search
124 (human, rat, mouse proteins) and TransDecoder (v2.0.1) open reading frame (ORF) prediction. The RNA-
125 seq assembly resulted in 19,737 protein-coding genes with a total length of 29.4 Mbp, which is available in
126 the China National GeneBank CNSA repository, Accession id: CNP0000340.

127 **Genome Annotation**

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

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

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

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

159 **Quality Assessment**

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

178

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

187 **Availability of supporting data**

188 Genome annotation results are available at the China National GeneBank CNSA repository, Accession id:
189 CNP0000340, and supporting materials, which include transcripts and genome assembly, are available
190 under the same project.

191

192 **Additional files**

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

194

195 **Abbreviations**

196 bp: base pair

197 BUSCO: Benchmarking Universal Single-Copy Orthologs

198 CDS: coding sequence

199 LINES: long interspersed elements

200 LTRs: long terminal repeats

201 Myr: million years

202 NCBI: National Center for Biotechnology Information

203 RefSeq: Reference sequence

204 RNA-seq: high-throughput messenger RNA sequencing

205 RIN: RNA integrity number

206 SINEs: short interspersed elements

207

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

209

210 **Funding**

211 EAM is supported by NIH 3T32DC012280-05S1. AK is supported by NIH R01 DC 11582.

212

213 **Authors' contributions**

214 SC, EAM, and AK developed the ideas, methods, and, wrote and revised the manuscript. BG, YF, YZ, WX,

215 HW, XL, and XX advised and revised the manuscript. BG provided the old animal tissues from Munich,

216 Germany. SC, YF, YZ, WX, HW, XL, and XX performed the analysis and annotation of the genome and

217 transcriptome. EAM prepared the DNA and RNA samples for sequencing.

218

219 **Acknowledgements**

220 The authors would like to thank Hilde Wohlfrom for sending tissues from Germany.

221

222 **References**

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

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

226 3. Govardovskii VI, Röhlich P, Szél A, Khokhlova TV. Cones in the retina of the Mongolian gerbil,
227 *Meriones unguiculatus*: an immunocytochemical and electrophysiological study. *Vision Res.* 1992;32:19–
228 27.

- 229 4. Vincent AL, Rodrick GE, Sodeman WA. The Mongolian gerbil in aging research. *Exp Aging Res.*
230 Routledge; 2007;6:249–60.
- 231 5. Small DL, Buchan AM. Animal models. *Br Med Bull.* Oxford University Press; 2000;56:307–17.
- 232 6. Bertorelli R, Adami M, Ongini E. The Mongolian gerbil in experimental epilepsy. *Ital J Neurol Sci.*
233 1995;16:101–6.
- 234 7. Löscher W. Genetic animal models of epilepsy as a unique resource for the evaluation of anticonvulsant
235 drugs. A review. *Methods Find Exp Clin Pharmacol.* 1984;6:531–47.
- 236 8. Bleich E-M, Martin M, Bleich A, Klos A. The Mongolian gerbil as a model for inflammatory bowel
237 disease. *Int J Exp Pathol.* Blackwell Publishing Ltd; 2010;91:281–7.
- 238 9. Hirayama F, Takagi S, Kusuhara H, Iwao E, Yokoyama Y, Ikeda Y. Induction of gastric ulcer and
239 intestinal metaplasia in mongolian gerbils infected with *Helicobacter pylori*. *J. Gastroenterol.* 1996;31:755–
240 7.
- 241 10. Zorio DAR, Monsma S, Sanes DH, Golding NL, Rubel EW, Wang Y. De novo sequencing and initial
242 annotation of the Mongolian gerbil (*Meriones unguiculatus*) genome. *Genomics.* Elsevier Inc; 2018.
- 243 11. Tarailo-Graovac M, Chen N. Using RepeatMasker to Identify Repetitive Elements in Genomic
244 Sequences. *Current Protocols in Bioinformatics.* Hoboken, NJ, USA: John Wiley & Sons, Inc;
245 2009;12:1269–4.10.14.
- 246 12. Jurka J, Kapitonov VV, Pavlicek A, Klonowski P, Kohany O, Walichiewicz J. Repbase Update, a
247 database of eukaryotic repetitive elements. *Cytogenet Genome Res.* 2005;110:462–7.
- 248 13. Benson G. Tandem repeats finder: a program to analyze DNA sequences. *Nucleic Acids Research.*
249 1999;27:573–80.
- 250 14. Smit AF. Interspersed repeats and other mementos of transposable elements in mammalian genomes.
251 *Current Opinion in Genetics & Development.* 1999;9:657–63.
- 252 15. Gibbs RA, Weinstock GM, Metzker ML, Muzny DM, Sodergren EJ, Scherer S, et al. Genome sequence
253 of the Brown Norway rat yields insights into mammalian evolution. *Nature.* 2004;428:493–521.
- 254
- 255 NCBI. NCBI Annotation Release *Meriones unguiculatus*, 2017;
256 https://www.ncbi.nlm.nih.gov/genome/annotation_euk/, (2 October 2017, date last accessed).
257
258
- 259
- 260
- 261
- 262

263

264

265

266

267

268

269

270

271

272

273

274 **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

275

276 **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

277
278
279
280
281
282
283
284
285
286
287

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

288
289

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

290

291

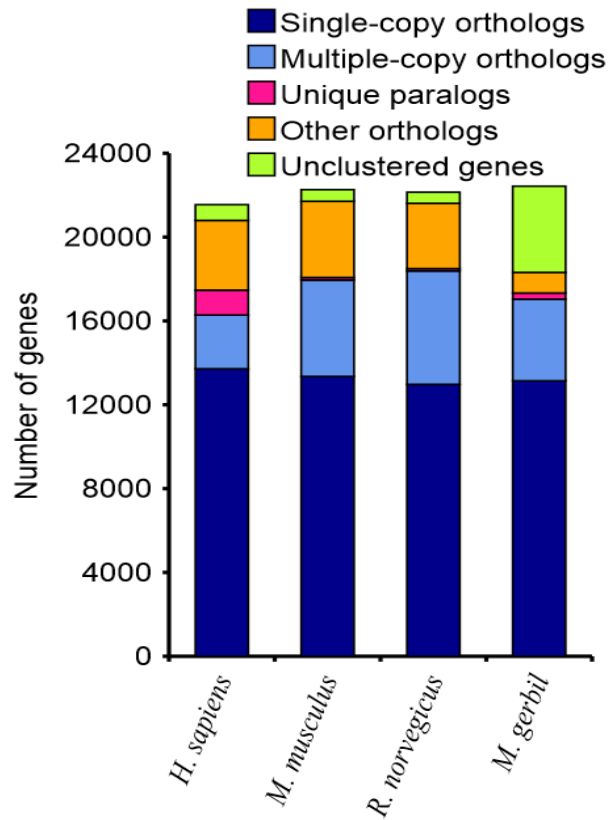
292

293

294

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



295

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

298

299

300 **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

301

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

Tissue	Run_accession	Sex	Age (postnatal)	Data size
--------	---------------	-----	-----------------	-----------

		day)	(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

303