

1 **Construction of a new chromosome-scale, long-read reference**  
2 **genome assembly for the Syrian hamster, *Mesocricetus auratus***

- 3 R. Alan Harris<sup>1</sup> [rharris1@bcm.edu](mailto:rharris1@bcm.edu)  
4 Muthuswamy Raveendran<sup>1</sup> [raveendr@bcm.edu](mailto:raveendr@bcm.edu)  
5 Dustin T. Lyfoung<sup>2</sup> [lyfoung@wisc.edu](mailto:lyfoung@wisc.edu)  
6 Fritz J Sedlazeck<sup>1</sup> [fritz.sedlazeck@bcm.edu](mailto:fritz.sedlazeck@bcm.edu)  
7 Medhat Mahmoud<sup>1</sup> [medhat.mahmoud@bcm.edu](mailto:medhat.mahmoud@bcm.edu)  
8 Trent M. Prall<sup>3</sup> [prall@wisc.edu](mailto:prall@wisc.edu)  
9 Julie A. Karl<sup>3</sup> [jakarl@wisc.edu](mailto:jakarl@wisc.edu)  
10 Harshavardhan Doddapaneni<sup>1</sup> [doddapan@bcm.edu](mailto:doddapan@bcm.edu)  
11 Qingchang Meng<sup>1</sup> [qingchang.meng@bcm.edu](mailto:qingchang.meng@bcm.edu)  
12 Yi Han<sup>1</sup> [yhan@bcm.edu](mailto:yhan@bcm.edu)  
13 Donna Muzny<sup>1</sup> [donnam@bcm.edu](mailto:donnam@bcm.edu)  
14 Roger W. Wiseman<sup>2,3</sup> [rwwiseman@wisc.edu](mailto:rwwiseman@wisc.edu)  
15 David H. O'Connor<sup>2,3</sup> [dhoconno@wisc.edu](mailto:dhoconno@wisc.edu)  
16 Jeffrey Rogers<sup>1</sup> [jr13@bcm.edu](mailto:jr13@bcm.edu)  
17 (Corresponding author: [jr13@bcm.edu](mailto:jr13@bcm.edu); 713-798-7783)

18 <sup>1</sup>Human Genome Sequencing Center and Department of Molecular and Human Genetics,  
19 Baylor College of Medicine, Houston, TX 77030

20 <sup>2</sup>Wisconsin National Primate Research Center, University of Wisconsin, Madison, WI 53711

21 <sup>3</sup>Department of Pathology and Laboratory Medicine, University of Wisconsin, Madison, WI  
22 53711

23

## 24 **Abstract**

### 25 **Background**

26 The Syrian hamster (*Mesocricetus auratus*) has been suggested as a useful  
27 mammalian model for a variety of diseases and infections, including infection with  
28 respiratory viruses such as SARS-CoV-2. The MesAur1.0 genome assembly was  
29 generated in 2013 using whole-genome shotgun sequencing with short-read sequence  
30 data. Current more advanced sequencing technologies and assembly methods now  
31 permit the generation of near-complete genome assemblies with higher quality and  
32 greater continuity.

### 33 **Findings**

34 Here, we report an improved assembly of the *M. auratus* genome (BCM\_Maur\_2.0)  
35 using Oxford Nanopore Technologies long-read sequencing to produce a chromosome-  
36 scale assembly. The total length of the new assembly is 2.46 Gbp, similar to the 2.50  
37 Gbp length of a previous assembly of this genome, MesAur1.0. BCM\_Maur\_2.0 exhibits  
38 significantly improved continuity with a scaffold N50 that is 6.7 times greater than

39 MesAur1.0. Furthermore, 21,616 protein coding genes and 10,459 noncoding genes are  
40 annotated in BCM\_Maur\_2.0 compared to 20,495 protein coding genes and 4,168  
41 noncoding genes in MesAur1.0. This new assembly also improves the unresolved  
42 regions as measured by nucleotide ambiguities, where approximately 17.11% of bases  
43 in MesAur1.0 were unresolved compared to BCM\_Maur\_2.0 in which the number of  
44 unresolved bases is reduced to 3.00%.

## 45 **Conclusions**

46 Access to a more complete reference genome with improved accuracy and continuity  
47 will facilitate more detailed, comprehensive, and meaningful research results for a wide  
48 variety of future studies using Syrian hamsters as models.

## 49 **Keywords**

50 Syrian hamster, *Mesocricetus auratus*, genome, disease model, COVID-19

51

## 52 **Data Description**

### 53 **Introduction**

54 The Syrian hamster (*Mesocricetus auratus*, NCBI:txid10036) has been used in  
55 biomedical research for decades because it is a good model for studies of cancer [1],  
56 reproductive biology [2] and infectious diseases [3,4], including SARS-CoV-2, influenza  
57 virus, and Ebola virus [5–9]. The use of Syrian hamsters in research has declined [10],  
58 likely due to advances in the genetic and molecular tools available for other rodents,  
59 especially laboratory mice, and not to a reduction in the utility of hamsters in biomedical  
60 research [3].

61 Syrian hamsters are particularly important for COVID-19 research. They spontaneously  
62 develop more severe lung disease than other animal models, such as wild-type mice,  
63 macaques, marmosets, and ferrets [5,11–14]. After intranasal infection, Syrian hamsters  
64 consistently show signs of respiratory distress, including labored breathing, but typically  
65 recover after 2 weeks [15]. This is in stark contrast to wild-type laboratory mice that are  
66 minimally susceptible to most SARS-CoV-2 strains that were circulating in 2020, though  
67 laboratory mice may be more susceptible to certain variants of concern that began  
68 circulating in 2021 [8,16]. Furthermore, a recent analysis has suggested that Syrian  
69 hamsters fed a high-fat, high-sugar diet exhibit accelerated weight gain and pathological  
70 changes in lipid metabolism, as well as more severe disease outcomes when  
71 subsequently infected with SARS-CoV-2 [17]. This result has obvious parallels with  
72 observations of the effects of comorbidities in humans suffering from COVID-19.

73 COVID-19 pathology in Syrian hamsters appears to be due to a dysregulated innate  
74 immune response involving signal transducer and activator of transcription factor 2  
75 (STAT2)-dependent type I (IFN-I) and type III interferon (IFN-III) signaling [18]. IFN-I  
76 signaling can limit virus replication and dissemination and it has been shown that  
77 intranasal administration of IFN-I in Syrian hamsters reduces viral load and tissue  
78 damage [19]. The human angiotensin-converting enzyme 2 (ACE2) was identified as  
79 the cell entry receptor of SARS-CoV-2 [20]. In addition, upon the engagement of ACE2  
80 with SARS-CoV2, cellular transmembrane protease 'serine 2' (TMPRSS2) mediates the  
81 priming of viral spike (S) protein by cleaving at the S1/S2 site and inducing the fusion of  
82 viral and host cellular membranes, thus facilitating viral entry into the cells [21]. Human  
83 ACE2 and hamster ACE2 receptors had previously been shown to share substantial  
84 sequence homology, which strongly points to interaction with SARS-CoV-2 receptor  
85 binding domain (RBD) structures and similar binding affinity [22]. *In silico* interaction  
86 prediction analysis suggests that human and hamster TMPRSS2 are structurally very  
87 similar. Even with slight differences in amino acid residue interactions, human and  
88 hamster TMPRSS2 activity are identical for residue interactions related to SARS-CoV-2  
89 infectivity [22]. As COVID-19 causes systemic disease in people, precision modeling of  
90 specific aspects of pathogenesis will require carefully evaluating similarities and  
91 differences across various biological processes in humans and Syrian hamsters which,  
92 in turn, will require extensive genomic comparisons between the two species.

93 The currently available reference genome sequence for the Syrian hamster was  
94 produced in 2013 using a whole-genome shotgun sequencing approach implementing  
95 short read sequencing technology. The resulting MesAur1.0 reference sequence

96 (Genbank accession number GCA\_000349665.1) is typical of those produced at that  
97 time, containing 237,699 separate contigs with contig N50 of 22,512 bp. The quality and  
98 research potential of the existing Syrian hamster genome is limited by the technology  
99 that was available at the time of its development; for example, the cluster of type I  
100 interferon genes was not resolvable with this technology. In this Data Note, we report  
101 the production of a new Syrian hamster reference genome that was sequenced using  
102 long-read methods on the Oxford Nanopore Technologies (ONT) PromethION platform  
103 and assembled into highly contiguous chromosomes using a combination of Flye [23]  
104 and Pilon [24] assembly software. The final assembly, BCM\_Maur\_2.0, improves upon  
105 quality and contiguity in comparison with MesAur1.0, with longer contigs and more  
106 contiguous sequence, allowing for a more complete reference genome with improved  
107 accuracy that will benefit a wide variety of future studies using the Syrian hamster  
108 reference genome.

## 109 **Methods**

### 110 **DNA isolation, library construction, and sequencing**

111 All genomic DNAs for this study were isolated from a single female LVG Golden Syrian  
112 hamster (SY011) that was purchased from Charles River, Inc. (Kingston, NY). All  
113 procedures were performed in accordance with the guidelines set by the Institutional  
114 Animal Care and Use Committee at the University of Wisconsin-Madison. The protocol  
115 was approved by the Institutional Animal Care and Use Committee at the University of  
116 Wisconsin-Madison (protocol number V00806). Data from this individual are available in  
117 NCBI BioProject [PRJNA705675](https://www.ncbi.nlm.nih.gov/bioproject/PRJNA705675), BioSamples [SAMN18096087](https://www.ncbi.nlm.nih.gov/biosamples/SAMN18096087) and [SAMN18096088](https://www.ncbi.nlm.nih.gov/biosamples/SAMN18096088).

118 Qiagen AllPrep DNA/RNA Mini kits were used to extract DNA from frozen liver while  
119 Qiagen Blood and Cell Culture DNA Midi Kits were used for extractions from frozen  
120 kidney. Ultra-high molecular weight DNA for optical mapping was purified from frozen  
121 liver using an Animal Tissue DNA Isolation Kit from Bionano Genomics, Inc. (San  
122 Diego, CA).

## 123 **Oxford Nanopore long-read sequencing**

124 We prepared three separate genomic DNA isolates from the same Syrian hamster  
125 (BioSample SAMN18096087). These aliquots were sheared to distinct target fragment  
126 lengths (10 kb, 20kb and 30kb) in order to assess the effect of fragment size on flowcell  
127 yield and improve efficiency. The two smaller length fragment libraries were sheared  
128 using Covaris gTube and the 30kb targeted size library was fragmented with Diagenode  
129 Megarupter 3, all following manufacturer's recommendations. The Oxford Nanopore  
130 sequencing libraries were prepared using the ONT 1D sequencing by ligation kit (SQK-  
131 LSK109). Briefly, 1-1.5ug of fragmented DNA was repaired with the NEB FFPE repair  
132 kit, followed by end repair and A-tailing with the NEB Ultra II end-prep kit. After a clean  
133 up step using AMPure beads, the prepared fragments were ligated to ONT specific  
134 adapters via the NEB blunt/TA master mix kit. Each library underwent a final clean up  
135 and was loaded onto a PromethION flow cell per manufacturer's instructions. One  
136 library was sequenced per flow cell with standard parameters for 72 hrs. Base-calling  
137 was done onboard the PromethION instrument using neuronal network based software  
138 (Oxford Nanopore Technologies, UK).

## 139 **Illumina sequencing**

140 500ng of input genomic DNA from a kidney sample (BioSample [SAMN18096088](#)) was  
141 used to generate standard PCR-free Illumina paired-end sequencing libraries. Libraries  
142 were prepared using KAPA Hyper PCR-free library reagents (KK8505, KAPA Bio-  
143 systems) in Beckman robotic workstations (Biomek FX and FXp models). Total genomic  
144 DNA was sheared into fragments of approximately 200-600 bp in a Covaris E220  
145 system (96-well format) followed by purification of the fragmented DNA using AMPure  
146 XP beads. A double size selection step was employed, with different ratios of AMPure  
147 XP beads, to select a narrow size band of sheared DNA molecules for library  
148 preparation. DNA end-repair and 3'-adenylation were then performed in the same  
149 reaction followed by ligation of the barcoded adaptors to create PCR-Free libraries. The  
150 resulting libraries were evaluated using the Fragment Analyzer (Advanced Analytical  
151 Technologies, Ames, Iowa) to assess library size and presence of remaining adaptor  
152 dimers. This was followed by qPCR assay using KAPA Library Quantification Kit and  
153 their SYBR FAST qPCR Master Mix to estimate the size and quantify fragment yield.

154 Sequencing was performed on the NovaSeq 6000 Sequencing System using the S4  
155 reagent kit (300 cycles) to generate 2 x 150 bp paired-end reads. The final  
156 concentration of the libraries loaded on flowcells was 400-450 pM. Briefly, the libraries  
157 were diluted in an elution buffer and denatured in sodium hydroxide. The denatured  
158 libraries were loaded into each lane of the S4 flow cell using the NovaSeq Xp Flow Cell  
159 Dock. Each lane included ~1% of a PhiX control library for run quality control.



## 160 **Genome Assembly**

161 We generated 221 gigabases of sequence data using the ONT PromethION platform  
162 (NCBI BioProject PRJNA705675, SRA Experiment SRX11206953). This represents an  
163 anticipated 88X coverage of the expected 2.5 Gbp Syrian hamster genome. The raw  
164 sequencing reads exhibited an N50 length of 15,730 bp. We used the Flye assembler  
165 v2.8.1 [23] to generate an initial *de novo* genome assembly. Given the potential  
166 sequence error rate of PromethION reads, it is advisable to use higher quality Illumina  
167 short reads mapped to an assembly to correct sequence errors in initial contigs.  
168 Consequently, we used Pilon software v. 1.23 [24] with default settings and 30X  
169 genome coverage of Illumina data (SRX10928323) generated from a kidney sample  
170 (SAMN18096088) obtained from the same individual for this sequence polishing step.  
171 Pilon sequence polishing was performed one time prior to the optical mapping analyses.

## 172 **Optical mapping for scaffold improvement**

173 Ultra-high molecular weight (UHMW) DNA was extracted following manufacturer's  
174 guidelines ([Bionano Prep SP Tissue and Tumor DNA Isolation protocol](#)) from frozen  
175 liver tissues obtained from the same animal used for ONT PromethION sequencing  
176 (SAMN18096087). Briefly, a total of 15-20mg of liver tissue was homogenized in cell  
177 buffer and digested with Proteinase K. DNA was precipitated with isopropanol and  
178 bound with nanobind magnetic disk (Bionano Genomics, USA). Bound UHMW DNA  
179 was resuspended in the elution buffer and quantified with Qubit dsDNA assay kits  
180 (ThermoFisher Scientific). DNA labeling was performed following manufacturer's  
181 protocols ([Bionano Prep Direct Label and Stain protocol](#)). Direct Labeling Enzyme 1

182 (DLE-1) reactions were carried out using 750 ng of purified UHMW DNA. Labeled DNA  
183 was loaded on Saphyr chips for imaging. The fluorescently labeled DNA molecules  
184 were imaged sequentially across nanochannel arrays (Saphyr chip) on a Saphyr  
185 instrument (Bionano Genomics Inc, USA). Effective genome coverage of greater than  
186 100X was achieved for all samples. All samples also met the following QC metrics:  
187 labelling density of ~15/100 kbp; filtered (>15kbp) N50 > 230 kbp; map rate > 70%.

188

189 Genome analysis of the resulting data was performed using software solutions provided  
190 by Bionano Genomics Inc. Briefly, automated optical genome mapping specific  
191 pipelines consisting of Bionano Access v1.4.3 and Bionano Solve v. 3.6.1 were used for  
192 data processing ([BioNano Access Software User Guide](#)). Hybrid scaffolding was  
193 performed using Bionano's custom software program implementing the following steps:  
194 1) generate *in silico* maps for sequence assembly; 2) align *in silico* sequence maps  
195 against Bionano genome maps to identify and resolve potential conflicts in either data  
196 set; 3) merge the non-conflicting maps into hybrid scaffolds; 4) align sequence maps to  
197 the hybrid scaffolds; and 5) generate AGP and FASTA files for the scaffolds. Pairwise  
198 comparisons of all DNA molecules were made to generate the initial consensus genome  
199 maps (\*.cmap). Genome maps were further refined and extended with best matching  
200 molecules. Optical map statistics were generated using Bionano software producing the  
201 Bionano Molecule Quality Report (MQR).

202

203 The optical map N50 (including only maps  $\geq 150$  kbp and minSites  $\geq 9$ ) was 0.2341  
204 Mbp and the average label density (scaffolds  $\geq 150$  kbp) was 17.40/100 kbp. This

205 yielded an effective molecule coverage with optical mapping information of 125.38X.  
206 The optical mapping analysis identified 84 conflicts with the prior Flye/Pilon scaffolds  
207 and these initial scaffolds were broken at those 84 sites. The completed assembly was  
208 submitted to NCBI and is available under accession [GCA\\_017639785.1](#).

## 209 **Gene annotation**

210 NCBI performed gene annotation using RNA-Seq data from multiple tissues including  
211 lung, trachea, brain, olfactory bulb and small intestine that are targets for SARS-CoV-2  
212 infection (NCBI BioProject [PRJNA675865](#)) [19].

## 213 **Quality assessment**

214 To assess the quality of our assembly compared to the previous MesAur1.0 we used  
215 Quast v5.0.2 [25] together with MUMmer v3.23 [26]. These tools provided a detailed  
216 comparison between these assemblies. In addition, the Illumina reads from the original  
217 reference (NCBI SRA [SRR413408](#)) were mapped to our assembly and the MesAur1.0  
218 reference using BWA v0.7.17 [27]. Quast was used to obtain discordant pair statistics.

219 We next used the software Benchmarking Universal Single-Copy Orthologs (BUSCO)  
220 v5.2.2 [28] to assess the quality of the genome assembly. BUSCO is based on the  
221 concept that single-copy orthologs should be highly conserved among closely related  
222 species. BUSCO performs gene annotation on an assembly and reports the number of  
223 gene models generated. BUSCO was performed using the OrthoDB v10 (odb10)  
224 release consisting of 12,692 genes shared across the superorder Euarchontoglires  
225 [29], the appropriate test for the Syrian hamster.

226 In addition, FRCbam [30] was used to compute Feature Response Curves (FRCurve)  
227 from the alignment of Illumina reads to the assembled contigs. FRC v1.3.0 was  
228 employed to evaluate both assemblies, using default parameters. BCM\_Maur\_2 was  
229 further evaluated using paired end mappings of the Illumina reads that had been used  
230 for Pilon polishing (SRX10928323). MesAur1.0 was then similarly evaluated using  
231 paired end mappings of Illumina reads used for the MesAur1.0 assembly (SRR413408).

## 232 Results

233 The initial Flye assembly consisted of 2.38 Gbp of sequence across 6,741 scaffolds with  
234 a scaffold N50 of 10.56 Mbp (**Table 1**). Pilon polishing of the Flye assembly had little  
235 effect on these metrics, but significant improvements were obtained when Bionano  
236 optical mapping results were used to improve scaffolding. As shown in **Table 1**, the  
237 optical mapping step reduced the total number of scaffolds in the final assembly by 395  
238 (5.9%) while increasing the N50 scaffold length by more than 8-fold to 85.18 Mbp. Our  
239 experience in comparing read lengths and total yield per flow-cell indicates that the  
240 optimal target size for fragmented DNA as input into Nanopore libraries and sequencing  
241 is 15-20 kb, which regularly yields 80-90 Gb of sequence data.

242 Of the 12,692 BUSCO gene models, 90.58% were annotated as complete genes in the  
243 initial Flye assembly (**Table 2**). Pilon polishing of this Flye-alone assembly added  
244 another 682 genes annotated completely and increased this proportion to 95.95% of the  
245 BUSCO gene model dataset. Improvements in assembly scaffolding resulting from the  
246 Bionano optical mapping step together with Pilon error correction decreased the  
247 proportions of fragmented and missing BUSCO gene models in the new assembly to

248 0.82% and 3.21% respectively, also improvements over the MesAur1.0 assembly. This  
249 advance translates to an additional 1189 complete BUSCO genes identified in the new  
250 assembly compared to MesAur1.0.

251

252 **Table 1.** Assembly statistics for BCM\_Maur\_2.0 versus the MesAur1.0 Syrian hamster  
253 assembly

Parameter	MesAur1.0	Flye	Flye + Pilon	Flye + Pilon + Bionano (BCM_Maur_2.0)
Assembly length (bp)	2,504,908,775	2,381,258,546	2,383,228,608	2,457,062,007
Ungapped length (bp)	2,076,159,990	2,381,254,546	2,383,226,373	2,383,228,883
Number of scaffolds	21,483	6,741	6,741	6,346
N50 scaffold length (bp)	12,753,307	10,564,357	10,573,641	85,184,847
Number of contigs	237,699	6,781	6,779	7,057
N50 contig length (bp)	22,512	10,022,145	10,097,207	9,471,653

254

255

256

257

258

259 **Table 2.** BUSCO statistics for BCM\_Maur\_2.0 versus the MesAur1.0 Syrian hamster

260 assembly

	MesAur1.0	Flye	Flye + Pilon	Flye + Pilon + Bionano (BCM_Maur_2.0)
Complete <sup>a</sup>	86.60%	90.58%	95.95%	95.97%
Complete and single-copy	85.75%	89.27%	94.43%	94.49%
Complete and duplicated	0.85%	1.31%	1.52%	1.47%
Fragmented	4.59%	3.23%	0.85%	0.82%
Missing	8.81%	6.19%	3.20%	3.21%

261 <sup>a</sup>12,692 gene models were included in this analysis

262

## 263 **Assembly Comparisons**

264 We also performed additional comparisons between the two assemblies. As

265 background, the karyotype of *M. auratus* is diploid  $2n = 44$ , including 14 pairs of

266 metacentric chromosomes, 3 pairs of telocentrics and 5 pairs of acrocentrics [31].

267 Illumina read k-mer analyses were performed to estimate the genome size using SGA

268 preqc [32] (2.57 Gbp) and Jellyfish [33] (2.90 Gbp). The total length of the  
269 BCM\_Maur\_2.0 assembly is 2.46 Gbp compared to the previous version's 2.50 Gbp.  
270 Despite having a similar total length, BCM\_Maur\_2.0 shows an improved continuity with  
271 a scaffold N50 that is 6.7 times greater than MesAur1.0 (**Table 1**); the L50 (i.e. the  
272 number of contigs longer than or equal to the N50 length) of BCM\_Maur\_2.0 is 22  
273 compared to MesAur1.0's 121. The longest scaffold of BCM\_Maur\_2.0 (187 Mb) is 2.35  
274 times larger than the longest scaffold from the previous assembly. N50 is calculated in  
275 the context of the assembly size rather than the genome size, so the NG50 statistic was  
276 used to directly compare the different assemblies. NG50 is the same as N50 except that  
277 it reports the length of the contig at which the size-ordered contigs (longest to shortest)  
278 collectively reaches 50% of the known or estimated genome size [34]. **Figure 1**  
279 illustrates the improved cumulative contig sequence length for any given NG50 value  
280 that is generated from the BCM\_Maur\_2.0 assembly as compared to MesAur1.0, based  
281 on the estimated genome size of 2.57 Gb calculated using SGA-preqc. The  
282 BCM\_Maur\_2.0 assembly further improves the unresolved regions as measured by  
283 nucleotide ambiguities (i.e. number of N's included in the final contigs). Approximately  
284 17.11% of bases in MesAur1.0 were unresolved. BCM\_Maur\_2.0 reduces the number  
285 of unresolved bases to 3.00%, with only very small gaps throughout the entire genome.  
286 **Figure 2** displays the overall increase in continuity of the BCM\_Maur\_2.0 assembly with  
287 longer contigs than the MesAur1.0 assembly and fewer short contigs. Finally, we  
288 compared feature response curves for BCM\_Maur\_2.0 and MesAur1.0 using  $FRC^{Bam}$   
289 [30].  $FRC^{Bam}$  shows that our new assembly is substantially more accurate based on the  
290 feature response approach (Supplementary Figure 1).

291 To establish the correctness of the structure and completeness of BCM\_Maur\_2.0, we  
292 also leveraged the Illumina short-reads that were published as part of the MesAur1.0  
293 assembly project. When mapping the MesAur1.0 Illumina reads back to the MesAur1.0  
294 reference, only 92.19% reads mapped successfully. When the same Illumina reads  
295 were instead mapped to BCM\_Maur\_2.0, 97.32% mapped successfully. When  
296 considering only properly paired reads, 75.76% and 87.62% mapped to MesAur1.0 and  
297 BCM\_Maur\_2.0, respectively.

298 Alignments between the current and previous Syrian hamster assemblies performed by  
299 NCBI [35] show that BCM\_Maur\_2.0 covers 98.95% of MesAur1.0 while MesAur1.0  
300 only covers 86.67% of BCM\_Maur\_2.0. This together with the additional 307 Mbp of  
301 ungapped sequence in BCM\_Maur\_2.0 indicates that BCM\_Maur\_2.0 is a more  
302 complete representation of the Syrian hamster genome. The percent identity in the  
303 regions aligned between the two assemblies is 99.76%.

## 304 **Transcript and Protein Alignments and Annotation** 305 **Comparisons** 306

307 NCBI annotation of BCM\_Maur\_2.0 [35] with Syrian hamster transcript and protein data  
308 show this assembly to be of high quality. Transcript alignments of Syrian hamster  
309 RefSeq (n=273), Genbank (n=751), and EST (n=558) data to BCM\_Maur\_2.0 show  
310 99.44% or more average percent identity and 98.88% or more average percent  
311 coverage. Alignments of these same transcript datasets to MesAur1.0 show 99.13% or  
312 more average percent identity and 93.49% or more average percent coverage.



313 Alignments of RefSeq transcripts showed a similar average percent indels in the  
314 BCM\_Maur\_2.0 (0.10%) and MesAur1.0 (0.11%) assemblies. Protein alignments of  
315 Syrian hamster RefSeq (n=261) and Genbank (n=485) data to BCM\_Maur\_2.0 show  
316 80.95% or more average percent identity and 89.18% or more average percent  
317 coverage. Alignments of these same protein datasets to MesAur1.0 show 80.57% or  
318 more average percent identity and 84.87% or more average percent coverage.

319 NCBI annotated 21,616 protein coding genes and 10,459 noncoding genes in  
320 BCM\_Maur\_2.0 compared to 20,495 protein coding genes and 4,168 noncoding genes  
321 in MesAur1.0 [36]. Only 7% of gene annotations are identical between BCM\_Maur\_2.0  
322 and MesAur1.0, suggesting that a number of previous errors have been corrected,  
323 though some differences are likely to be real differences between the animals used for  
324 the different assemblies. Minor changes between BCM\_Maur\_2.0 and MesAur1.0 were  
325 made in 46% of gene annotations and major changes were made in 15% of gene  
326 annotations. We further note that, based on NCBI annotation feature counts,  
327 BCM\_Maur\_2.0 has only 33 RefSeq models that were filled using transcript sequence  
328 to compensate for an assembly gap [35]. This is compared to 5,050 RefSeq models  
329 similarly compensated in MesAur1.0.

### 330 **Interferon type 1 alpha gene cluster**

331 Given the importance of type I interferon responses during SARS-CoV-2 infection, we  
332 next compared the interferon type I alpha gene cluster in the BCM\_Maur\_2.0 assembly  
333 relative to this genomic region in the original MesAur1.0 assembly. The MesAur1.0  
334 scaffold NW\_004801649.1 includes annotations for four interferon type I alpha loci but

335 this genomic sequence is riddled with numerous gaps. Of these four candidate loci,  
336 only LOC101824534 appears to contain a complete interferon alpha-12-like coding  
337 sequence with the ability to encode a predicted protein (XP\_005074343.1). The  
338 LOC101824794 gene sequence can only encode a 162 amino acid protein due to a 5'  
339 truncation. The remaining pair of candidate genes (LOC101836618 and  
340 LOC101836898) appear to have aberrant transcript models that have fused putative  
341 exons from neighboring loci. In mice and humans, the interferon alpha gene cluster is  
342 flanked by single copy interferon beta 1 (*Ifnb1*) and interferon epsilon (*Ifne*) genes.  
343 Although neither of these genes are present on the MesAur1.0 scaffold  
344 NW\_004801649.1, this assembly does contain a *Ifne* gene on a short 2,408 bp contig  
345 that is predicted to code for a protein of 192 amino acids. These observations  
346 emphasize the need for an improved genomic assembly for Syrian hamsters given that  
347 the interferon alpha gene cluster includes more than a dozen tightly linked functional  
348 genes plus multiple pseudogenes in a wide variety of species including mice and  
349 humans.

350 In the BCM\_Maur\_2.0 assembly, the interferon type I alpha gene cluster is contained on  
351 the NW\_024429197.1 super scaffold that spans nearly 75 Mbp. **Figure 3** illustrates this  
352 genomic region in comparison with the interferon type 1 alpha regions of MesAur1.0  
353 (NW\_004801649.1) and the well-characterized C57BL/6J mouse assembly  
354 (NC\_000070.7). Fourteen predicted interferon type I alpha genes as well as five  
355 presumptive pseudogenes lie within a span of 196 Kbp of the new Syrian hamster  
356 assembly (**Figure 3** and [Supplemental Table 1](#)). This genomic organization is quite  
357 comparable to that observed in the mouse genome where there are also fourteen

358 functional interferon alpha genes and four pseudogenes. This hamster gene cluster is  
359 flanked by *Ifnb1* and *Ifne* genes consistent with expectations from the mouse and other  
360 species. The NCBI annotations characterize twelve of these genes as interferon alpha-  
361 12-like ([Supplemental Table 1](#)). The remaining pair of functional genes  
362 (LOC101824794 and LOC121144100) are listed as interferon alpha-9-like and they  
363 encode shorter predicted proteins. The increased length of this genomic region in the  
364 mouse assembly is largely due to the presence of the interferon zeta gene family (*Ifnz*,  
365 Gm13271, Gm13272, etc.). This *Ifnz* gene family appears to be absent in Syrian  
366 hamsters since the closest matches to predicted hamster protein sequences are only  
367 28% identical at the amino acid level. The interferon type I alpha gene cluster in the  
368 BCM\_Maur\_2.0 assembly lies within more than 12 Mbp of contiguous genomic  
369 sequence with the nearest flanking gaps located 2.66 Mbp proximal and 9.07 Mbp distal  
370 to the *Ifne* and *Ifnb1* genes, respectively. The availability of a contiguous hamster  
371 genomic sequence and associated transcriptional regulatory elements for this complex  
372 immune gene region may be helpful for investigators who are interested in unravelling  
373 mechanisms that control interferon expression during infections with SARS-CoV-2 as  
374 well as challenges with other viral pathogens.

## 375 **Conclusions**

376 The improved Syrian hamster assembly and annotation described here will facilitate  
377 research into this important animal model for COVID-19. Specifically, reagents for  
378 studying immune responses in hamsters have lagged behind those available for  
379 laboratory mice. BCM\_Maur\_2.0 will facilitate the identification of cross-reactive  
380 reagents originally developed to study immunity in other species. Additionally, a more

381 accurate genome assembly will improve the analyses of host responses to infection by  
382 enabling more accurate interpretation of RNA-seq experiments.

383 Relative to other recent assemblies that use a combination of long-read sequencing and  
384 short-read polishing, this genome assembly and annotation compares very favorably.  
385 The scaffold N50 of >85 Mbp is quite consistent with other long read assemblies. The  
386 contig N50 and total number of scaffolds or contigs are likewise reasonable and  
387 consistent with other similar mammalian reference genomes. The number of protein  
388 coding genes identified is within the expected range, although additional attention will  
389 likely be needed to resolve duplicated, repetitive gene loci, potentially leveraging recent  
390 advances in ultralong read sequencing.

391 What additional genomic resources would be needed to make hamsters a better model  
392 for COVID-19? Deep long read transcriptome analysis of multiple tissues and ages  
393 would be the best next step, in order to define not just the genes expressed but the  
394 alternative splicing of genes across tissues and developmental stages. Also, long read  
395 RNA-seq of tissues following experimental challenge with SARS-CoV-2 and other  
396 viruses would facilitate improvements in the quality of antiviral gene models.

397 The availability of higher accuracy sequences should lead to the development of  
398 specific reagents for monitoring immune responses. For example, epitopes that are  
399 shared between hamsters and other rodents can be used to identify monoclonal  
400 antibody reagents for flow cytometry that are predicted to be cross-reactive. Additional  
401 reagent development will be enabled by creating synthetic versions of hamster proteins  
402 that can be used as immunogens to make hamster-specific antibodies.

403 One surprising motivation for this study is that Syrian hamsters, which were quickly  
404 identified as a high value model for COVID-19, did not have a higher quality reference  
405 genome at the start of the pandemic. While we worked quickly to generate this data and  
406 make it available to the scientific community, better preparedness will be critical for  
407 future unexpected epidemics. To this end, we would encourage investment in continued  
408 refinement and improvement of reference genomes for all of the rodent, bat and  
409 nonhuman primate models that are commonly used to study viruses in order to prevent  
410 this situation from recurring in the future. Such an investment would also yield improved  
411 genomic resources that would provide broad benefit to the entire scientific community.

412

413

## 414 **Availability of Supporting Data and** 415 **Materials**

416 The MesAur1.0 genome assembly is available in the NCBI database under BioProject  
417 [PRJNA77669](#) (GenBank accession [GCA\\_000349665.1](#)). The new BCM\_Maur\_2.0  
418 genome assembly is available in the NCBI data repository under BioProject  
419 [PRJNA705675](#) (GenBank accession [GCA\\_017639785.1](#)). Oxford Nanopore  
420 ([SRX11206953](#)) and Illumina ([SRX10928323](#)) sequencing data are available through  
421 the NCBI SRA. The Bionano data are available from the BioProject page as NCBI  
422 accession [SUPPF\\_0000004259](#). The Illumina RNA-Seq data from multiple tissues  
423 including lung, trachea, brain, olfactory bulb and small intestine are available under  
424 NCBI BioProject [PRJNA675865](#).

425

## 426 **Additional Files**

427 **Supplementary Table 1.** Predicted genes in the Interferon type 1 alpha cluster of the  
428 BCM\_Maur\_2.0 assembly.

## 429 **Abbreviations**

430 ACE2: angiotensin-converting enzyme 2; BCM: Baylor College of Medicine; bp: base  
431 pairs; BUSCO: Benchmarking Universal Single-Copy Orthologs; BWA: Burrows-  
432 Wheeler Aligner; COVID-19: coronavirus disease 2019; EST: expressed sequence tag;  
433 FFPE: formalin-fixed, paraffin-embedded; Gbp: gigabase pairs; GC: guanine-cytosine;  
434 IFN: interferon; kbp: kilobase pairs; Mbp: megabase pairs; MQR: Molecule Quality  
435 Report; NCBI: National Center for Biotechnology Information; NEB: New England  
436 BioLabs; ng: nanogram; ONT: Oxford Nanopore Technologies; PCR: polymerase chain  
437 reaction; RBD: receptor-binding domain; RNA-Seq: RNA-sequencing; SARS-CoV-2:  
438 severe acute respiratory syndrome coronavirus 2; STAT2: signal transducer and  
439 activator of transcription factor 2; TMPRSS2: transmembrane protease serine 2

## 440 **Competing interests**

441 The authors declare that they have no competing interests.

## 442 **Funding**

443 This research was supported by contract HHSN272201600007C awarded to DHO from the  
444 National Institute of Allergy and Infectious Diseases of the National Institutes of Health. The

445 content of this publication is solely the responsibility of the authors and does not  
446 necessarily represent the official views of the National Institutes of Health.

## 447 **Authors' Contributions**

448 R.A.H. performed genome assembly and quality assessment, data and metadata  
449 submission, and contributed to manuscript preparation. F.S. and M.M. performed  
450 assembly assessment and comparison analyses. T.M.P. and R.W.W. performed  
451 transcript and annotation comparisons. D.H.O. managed experimental design and  
452 oversight and coordinated manuscript preparation. H.D., Q.M. and Y.H. developed,  
453 optimized and implemented protocols for ONT PromethION sequencing. M.R., D.M.,  
454 J.A.K. and J.R. performed project and/or data management. R.A.H., D.H.O., D.T.L.,  
455 T.M.P., R.W.W., M.M., F.S. and J.R. wrote the manuscript. All authors approved the  
456 manuscript.

## 457 **Acknowledgements**

458 We are extremely grateful to Dr. Tadashi Maemura for collecting the Syrian hamster  
459 tissues that were used for the sequence analyses described here. We also thank Dr.  
460 Benjamin tenOever for sharing Syrian hamster RNA-Seq datasets generated by his  
461 group prior to publication. And we also wish to thank two reviewers for their helpful  
462 comments.

463

464

465 **Figure 1: Cumulative length and continuity comparison of MesAur1.0 and**  
466 **BCM\_Maur\_2.0.** This summarizes the length of contigs/scaffolds across the  
467 assemblies. Given the length of contigs, the NG50 (mid x-axis) summarizes the  
468 sequence length of the shortest contig/scaffold at 50% of the total genome length. For  
469 genome length, the SGA preqc estimate of 2.57 Gbp was used.

470 **Figure 2: Contig length and count comparison between BCM\_Maur\_2.0 and**  
471 **MesAur1.0.** Log length of contigs on the X axis and normalized count on the Y axis  
472 comparing BCM\_Maur\_2.0 assembly and the previous assembly. Contigs from  
473 BCM\_Maur\_2.0 are shown red and contigs for MesAur1.0 are shown in gray.

474 **Figure 3: Comparison of interferon type 1 alpha gene cluster between MesAur1.0,**  
475 **BCM\_Maur\_2.0 and GCRm39 mouse genome assembly.** The genomic intervals  
476 illustrated here are defined by the flanking interferon beta 1 and interferon epsilon  
477 genes except for MesAur1.0 which does not include an interferon epsilon or beta 1 gene  
478 in a continuous sequence with interferon type 1 alpha genes. White space within each  
479 scaffold represents gaps in the MesAur1.0 assembly. Accession numbers for each  
480 genomic sequence are indicated on the right with genomic coordinates for the extracted  
481 intervals shown below their respective accession numbers. Predicted interferon type 1  
482 alpha genes are highlighted in blue while putative pseudogenes are depicted with open  
483 symbols and labelled below each assembly.

484



## 485 **References**

- 486 1. LaRocca CJ, Han J, Gavrikova T, Armstrong L, Oliveira AR, Shanley R, et al.. Oncolytic  
487 adenovirus expressing interferon alpha in a syngeneic Syrian hamster model for the treatment  
488 of pancreatic cancer. *Surgery*. 2015; doi: 10.1016/j.surg.2015.01.006.
- 489 2. Pal S, Haldar C, Verma R. Photoperiodic modulation of ovarian metabolic, survival,  
490 proliferation and gap junction markers in adult golden hamster, *Mesocricetus auratus*. *Comp*  
491 *Biochem Physiol A Mol Integr Physiol*. 2021; doi: 10.1016/j.cbpa.2021.111083.
- 492 3. McCann KE, Sinkiewicz DM, Norvelle A, Huhman KL. De novo assembly, annotation, and  
493 characterization of the whole brain transcriptome of male and female Syrian hamsters. *Sci Rep*.  
494 2017; doi: 10.1038/srep40472.
- 495 4. Saini S, Rai AK. Hamster, a close model for visceral leishmaniasis: Opportunities and  
496 challenges. *Parasite Immunol*. 2020; doi: 10.1111/pim.12768.
- 497 5. Chan JF-W, Zhang AJ, Yuan S, Poon VK-M, Chan CC-S, Lee AC-Y, et al.. Simulation of the  
498 Clinical and Pathological Manifestations of Coronavirus Disease 2019 (COVID-19) in a Golden  
499 Syrian Hamster Model: Implications for Disease Pathogenesis and Transmissibility. *Clin Infect*  
500 *Dis*. 2020; doi: 10.1093/cid/ciaa325.
- 501 6. Prescott J, Falzarano D, Feldmann H. Natural Immunity to Ebola Virus in the Syrian Hamster  
502 Requires Antibody Responses. *J Infect Dis*. 2015; doi: 10.1093/infdis/jiv203.
- 503 7. Huo J, Mikolajek H, Le Bas A, Clark JJ, Sharma P, Kipar A, et al.. A potent SARS-CoV-2  
504 neutralising nanobody shows therapeutic efficacy in the Syrian golden hamster model of  
505 COVID-19. *Nat Commun*. 2021; doi: 10.1038/s41467-021-25480-z.
- 506 8. Mohandas S, Yadav PD, Shete A, Nyayanit D, Sapkal G, Lole K, et al.. SARS-CoV-2 Delta  
507 Variant Pathogenesis and Host Response in Syrian Hamsters. *Viruses*. 2021; doi:  
508 10.3390/v13091773.
- 509 9. Mifsud EJ, Tai CM, Hurt AC. Animal models used to assess influenza antivirals. *Expert Opin*  
510 *Drug Discov*. 2018; doi: 10.1080/17460441.2018.1540586.
- 511 10. Gao M, Zhang B, Liu J, Guo X, Li H, Wang T, et al.. Generation of transgenic golden Syrian  
512 hamsters. *Cell Res*. 2014; doi: 10.1038/cr.2014.2.
- 513 11. Imai M, Iwatsuki-Horimoto K, Hatta M, Loeber S, Halfmann PJ, Nakajima N, et al.. Syrian  
514 hamsters as a small animal model for SARS-CoV-2 infection and countermeasure development.  
515 *Proc Natl Acad Sci U S A*. 2020; doi: 10.1073/pnas.2009799117.
- 516 12. Rockx B, Kuiken T, Herfst S, Bestebroer T, Lamers MM, Oude Munnink BB, et al..  
517 Comparative pathogenesis of COVID-19, MERS, and SARS in a nonhuman primate model.  
518 *Science*. 2020; doi: 10.1126/science.abb7314.
- 519 13. Rogers TF, Zhao F, Huang D, Beutler N, Burns A, He W-T, et al.. Isolation of potent SARS-  
520 CoV-2 neutralizing antibodies and protection from disease in a small animal model. *Science*.  
521 2020; doi: 10.1126/science.abc7520.

- 522 14. Shi J, Wen Z, Zhong G, Yang H, Wang C, Huang B, et al.. Susceptibility of ferrets, cats,  
523 dogs, and other domesticated animals to SARS-coronavirus 2. *Science*. 2020; doi:  
524 10.1126/science.abb7015.
- 525 15. Muñoz-Fontela C, Dowling WE, Funnell SGP, Gsell P-S, Riveros-Balta AX, Albrecht RA, et  
526 al.. Animal models for COVID-19. *Nature*. 2020; doi: 10.1038/s41586-020-2787-6.
- 527 16. Montagutelli X, Prot M, Levillayer L, Salazar EB, Jouvion G, Conquet L, et al.. The B1.351  
528 and P.1 variants extend SARS-CoV-2 host range to mice. bioRxiv.
- 529 17. Port JR, Adney DR, Schwarz B, Schulz JE, Sturdevant DE, Smith BJ, et al.. Western diet  
530 increases COVID-19 disease severity in the Syrian hamster. *bioRxiv*. 2021; doi:  
531 10.1101/2021.06.17.448814.
- 532 18. Boudewijns R, Thibaut HJ, Kaptein SJF, Li R, Vergote V, Seldeslachts L, et al.. STAT2  
533 signaling restricts viral dissemination but drives severe pneumonia in SARS-CoV-2 infected  
534 hamsters. *Nat Commun*. 2020; doi: 10.1038/s41467-020-19684-y.
- 535 19. Hoagland DA, Møller R, Uhl SA, Oishi K, Frere J, Golyunker I, et al.. Leveraging the antiviral  
536 type I interferon system as a first line of defense against SARS-CoV-2 pathogenicity. *Immunity*.  
537 2021; doi: 10.1016/j.immuni.2021.01.017.
- 538 20. Brooke GN, Prischi F. Structural and functional modelling of SARS-CoV-2 entry in animal  
539 models. *Sci Rep*. 2020; doi: 10.1038/s41598-020-72528-z.
- 540 21. Hoffmann M, Kleine-Weber H, Schroeder S, Krüger N, Herrler T, Erichsen S, et al.. SARS-  
541 CoV-2 Cell Entry Depends on ACE2 and TMPRSS2 and Is Blocked by a Clinically Proven  
542 Protease Inhibitor. *Cell*. 2020; doi: 10.1016/j.cell.2020.02.052.
- 543 22. Rizvi ZA, Dalal R, Sadhu S, Kumar Y, Shrivastava T, Gupta SK, et al.. Immunological and  
544 cardio-vascular pathologies associated with SARS-CoV-2 infection in golden syrian hamster.  
545 Cold Spring Harbor Laboratory.
- 546 23. Kolmogorov M, Yuan J, Lin Y, Pevzner PA. Assembly of long, error-prone reads using  
547 repeat graphs. *Nat Biotechnol*. 2019; doi: 10.1038/s41587-019-0072-8.
- 548 24. Walker BJ, Abeel T, Shea T, Priest M, Abouelliel A, Sakthikumar S, et al.. Pilon: an  
549 integrated tool for comprehensive microbial variant detection and genome assembly  
550 improvement. *PLoS One*. 2014; doi: 10.1371/journal.pone.0112963.
- 551 25. Gurevich A, Saveliev V, Vyahhi N, Tesler G. QUAST: quality assessment tool for genome  
552 assemblies. *Bioinformatics*. 2013; doi: 10.1093/bioinformatics/btt086.
- 553 26. Kurtz S, Phillippy A, Delcher AL, Smoot M, Shumway M, Antonescu C, et al.. Versatile and  
554 open software for comparing large genomes. *Genome Biol*. 2004; doi: 10.1186/gb-2004-5-2-r12.
- 555 27. Li H. Aligning sequence reads, clone sequences and assembly contigs with BWA-MEM.  
556 arXiv [q-bio.GN].
- 557 28. Seppey M, Manni M, Zdobnov EM. BUSCO: Assessing Genome Assembly and Annotation  
558 Completeness. *Methods Mol Biol*. 2019; doi: 10.1007/978-1-4939-9173-0\_14.
- 559 29. : BUSCO. [https://busco-data.ezlab.org/v4/data/lineages/euarchontoglires\\_odb10.2020-09-](https://busco-data.ezlab.org/v4/data/lineages/euarchontoglires_odb10.2020-09-)

- 560 10.tar.gz Accessed 2021 Jun 4.
- 561 30. Vezzi F, Narzisi G, Mishra B. Reevaluating assembly evaluations with feature response  
562 curves: GAGE and assemblathons. *PLoS One*. 2012; doi: 10.1371/journal.pone.0052210.
- 563 31. Lehman JM, Macpherson I, Moorhead PS. KARYOTYPE OF THE SYRIAN HAMSTER. *J*  
564 *Natl Cancer Inst*. 31:639–501963;
- 565 32. Simpson JT. Exploring genome characteristics and sequence quality without a reference.  
566 *Bioinformatics*. 2014; doi: 10.1093/bioinformatics/btu023.
- 567 33. Marçais G, Kingsford C. A fast, lock-free approach for efficient parallel counting of  
568 occurrences of k-mers. *Bioinformatics*. 2011; doi: 10.1093/bioinformatics/btr011.
- 569 34. Alhakami H, Mirebrahim H, Lonardi S. A comparative evaluation of genome assembly  
570 reconciliation tools. *Genome Biol*. 2017; doi: 10.1186/s13059-017-1213-3.
- 571 35. : Mesocricetus auratus Annotation Report.  
572 [https://www.ncbi.nlm.nih.gov/genome/annotation\\_euk/Mesocricetus\\_auratus/103/](https://www.ncbi.nlm.nih.gov/genome/annotation_euk/Mesocricetus_auratus/103/) Accessed  
573 2021 Jun 4.
- 574 36. : Mesocricetus auratus Annotation Report.  
575 [https://www.ncbi.nlm.nih.gov/genome/annotation\\_euk/Mesocricetus\\_auratus/102/](https://www.ncbi.nlm.nih.gov/genome/annotation_euk/Mesocricetus_auratus/102/) Accessed  
576 2021 Jun 4.





