1 **A Complete Workflow for scmtDNAseq in CHO cells, from Cell Culture to Bioinformatic**
2 **Analysis**

3 **Alan Foley[1,3], Nga Lao[1], Colin Clarke[2,3], Niall Barron[1,3]**

4 [1]Cell Engineering Group, NIBRT, Dublin, Ireland

5 [2]Bioinformatics Group, NIBRT, Dublin, Ireland

6 [3]School of Chemical and Bioprocess Engineering, University College Dublin, Ireland.

7 **\* Correspondence:**
8 niall.barron@nibrt.ie

## 1    Abstract

12 Chinese Hamster Ovary (CHO) cells have a long history in the biopharmaceutical industry and
13 currently produce the vast majority of recombinant therapeutic proteins. The key step in controlling
14 process and product consistency is the development of a producer cell line derived from a single cell
15 clone. However, it is recognised that genetic and phenotypic heterogeneity between individual cells
16 in a clonal CHO population tends to arise over time. Previous bulk analysis of CHO cell populations
17 has revealed considerable variation within the mtDNA sequence (heteroplasmy) which could have
18 implications for the performance of the cell line. By analysing heteroplasmy of single cells within the
19 same population, this heterogeneity can be characterised with greater resolution. Such analysis may
20 identify heterogeneity in the mitochondrial genome which impacts the overall phenotypic
21 performance of a producer cell population, and potentially reveal routes for genetic engineering. A
22 critical first step is the development of robust experimental and computational methods to enable
23 single cell mtDNA sequencing (termed scmtDNAseq). Here, we present a protocol from cell culture
24 to bioinformatic analysis and provide preliminary evidence of significant mtDNA heteroplasmy
25 across a small panel of single CHO cells.

26

## 2    Introduction

28 Chinese Hamster Ovary (CHO) cells are the most commonly used mammalian host for the
29 production of recombinant proteins (Walsh et al., 2022). Optimisation of biopharmaceutical
30 production in CHO has led to titers routinely in the 3-8 g/L range (Kelley et al., 2018). Due to their
31 importance in energy production, understanding mitochondrial function in product-producing CHO
32 cell lines is of particular importance. While most mitochondrial proteins are encoded by nuclear
33 DNA, a small number are encoded by mitochondrial DNA (mtDNA). The CHO mitochondrial
34 genome contains 37 genes, all of which support oxidative phosphorylation (OXPHOS). Thirteen
35 protein-encoding subunits are accompanied by 2 rRNAs and 22 tRNAs in a 16,283bp plasmid-like
36 circular structure (NCBI, 2023). mtDNA is highly compact, with the only significant non-coding
37 region in the D-loop (**Fig. 1A**).

38  Assuming a CHO cell has typical numbers of mitochondria per cell (100-10,000), each with 2-10
39  copies of mtDNA, the total genome copy per cell is large (Dhiman et al., 2019). In 'homoplasmy' all
40  copies of mtDNA within a cell are identical; however, mitochondria can also exist in a state of
41  'heteroplasmy' where mutated versions of mtDNA co-exist with wild-type mtDNA within the same
42  cell; and possibly even within the same mitochondrion. When the proportion of mutant mtDNA is
43  above a particular threshold, mitochondrial dysfunction can occur (Dimauro and Davidzon, 2005). In
44  human disease, this could mean development of metabolic disease including neurodegenerative
45  disorders (Keogh and Chinnery, 2015); in CHO cell culture it could manifest as a change in
46  bioreactor performance.

47  Previous bulk analysis identified several levels of heteroplasmy between CHO cell lines (Kelly. P et
48  al., 2017); laying a theoretical explanation for the metabolic heterogeneity often observed in CHO
49  cell cultures (Gilbert et al., 2013). Single cell sequencing of mtDNA (scmtDNAseq) has previously
50  been employed in non-CHO cell lines using high PCR cycle numbers of 40 (Zambelli et al., 2017)
51  and 45 (Maeda et al., 2020). Higher PCR cycle numbers are associated with greater risk of
52  undesirable secondary products such as PCR artefacts (Lorenz, 2012). In single cell mtDNA analysis,
53  starting mtDNA copy number is low (<100,000) therefore even small contaminations can confound
54  true mutation nucleotide identification. In fact, to call heteroplasmic mutations at 0.015 allele
55  frequency (a conservative level) PCR amplification should ideally not exceed 30 cycles (Zambelli et
56  al., 2017). Also, Maeda et al. focused on specific mutations, not the whole mtDNA genome,
57  precluding identification of as yet unknown mutations. There is real value in novel whole mtDNA
58  single cell analysis with a low PCR cycle number.

59  Here, we sought to develop an optimised method to amplify mtDNA and sequence from single cells.
60  To demonstrate the method, we analysed four single CHO cells and a bulk (multiple cells) sample for
61  comparison. Single cells were isolated by FACS into lysis buffer with an emphasis on simple and
62  reproducible gating (**Fig. 1B**). After optimisation of the lysis buffer, PCR kit and purification system,
63  long-range PCR (LRPCR) cycle number (35x) was kept lower than previously reported methods.
64  Importantly, this provides more confidence in calling low-frequency heteroplasmy. To ensure
65  exclusion of contaminating nuclear mitochondrial DNA (Numts), primers were designed to
66  exclusively map to CHO mtDNA and amplicons size-selected via gel electrophoresis. By confirming
67  mtDNA amplification by agarose gel, we were able to improve the efficiency of our sequencing –
68  since only successful reactions were brought forward for library preparation. Illumina DNA libraries
69  were generated and iSeq100-derived sequencing output was processed and analysed using a bespoke
70  bioinformatics pipeline. Preprocessing was performed in Linux and data analysis in R.

71  **Fig. 1**

72

73  **3      Materials and Equipment**

74  **2.1 CHO cell culture**

75      1.  125mL bioreactor flasks (Nalgene 10266432).

76      2.  Appropriate CHO cell culture medium (e.g. Gibco CD FortiCHO 10887640).

77   3. CHO cell lines of interest (e.g. **Table 1**).

78 **2.2 Immunolabelling and Staining**

79   1. DPBS

80   2. Nuclease-free water

81   3. Trypan Blue 0.4% (Gibco 15250061)

82   4. Luna II (or other appropriate cell counter)

83   5. DAPI (Invitrogen D1306)

84   6. Goat F(ab')2-Fluorescein anti-Human IgG (Sigma Aldrich SAB3701254-2MG) to label IgG-

85     producing cells if desired. Other appropriate fluorescent stains could also be employed (e.g.

86     CellTracker Green Invitrogen 11570166).

87 **2.3 FACS**

88   1. 70% IPA.

89   2. FACS with appropriate lasers for DAPI and FITC detection. Here, a BD FACS Melody was

90     used.

91   3. FACS polystyrene tubes (Falcon Corning 1018640)

92   4. U-bottom 96-well plates (Corning 3799)

93   5. Parafilm

94   6. TCL Buffer (QIAGEN 1070498)

95 **2.4 AMPure purification**

96   1. AMPure XP Beads (10136224)

97   2. 70% ethanol

98   3. Elution buffer (QIAGEN 19086)

99   4. Sterile PCR tubes (autoclaved)

100   5. 0.2mL tube magnetic stand (New England Biolabs S1515S)

101    6.  10uL multichannel pipette (optional)

102  **2.5 Long Range PCR**

103  Primers were designed using NCBI Primer-BLAST to specifically bind to mtDNA, and not to any

104  known CHO nuclear DNA sequences to minimise Numt contamination.

105    1.  SuperFi II Plat Taq (Invitrogen 12361010).

106    2.  PCR thermocycler

107    3.  10uM forward and reverse primers (**Table 2**) (IDT)

108    4.  10mM dNTP Mix (Thermo Scientific R0192)

109  **2.6 Agarose Gel**

110    1.  Agarose powder.

111    2.  TAE buffer.

112    3.  SafeView (NBS Biologicals). Ethidium Bromide is an alternative.

113    4.  GeneRuler 1kb Plus Ladder (Thermo Scientific SM1333).

114    5.  Gel Viewer/transilluminator.

115    6.  Disposable lab scalpel.

116    7.  Eppendorf tubes.

117  **2.7 Gel Purification**

118    1.  QIAquick Gel Extraction Kit (QIAGEN 28706X4). Other gel extraction kits could also be

119      utilised.

120  **2.8 Qubit**

121    1.  Qubit 4 Fluorometer (Invitrogen)

122    2.  Qubit 1x dsDNA HS Kit (Invitrogen Q33230)

123  **2.9 Sequencing**

124     1. iSeq100 (Illumina) PE150, 8 million reads

125     2. Illumina DNA Prep, (M) Tagmentation (24 Samples) (Illumina 20018704)

126     3. IDT for Illumina DNA/RNA UD Indexes Set A, Tagmentation (96 Indexes, 96 Samples)

127        (Illumina 20027213)

128     4. iSeq 100 i1 Reagent v2 (300-cycle) (Illumina 20031371)

129     5. PhiX v3 (Illumina FC-110-3001)

130

131  **4      Methods**

132  All steps up to the completion of the LRPCR for the 4 samples (**Table 1**) were performed in sterile
133  conditions (BSC).

134  **Table 1**

135  **4.1   PCR Component Storage**

136  Since the LRPCR amplifies from less than 5,000 copies of mtDNA, PCR components must have
137  optimal efficacy. This was ensured by making small (20uL) aliquots of dNTPs (Thermo Scientific
138  R0192) and primers (IDT) and storing at -80°C. New aliquots were used for each lot of PCR
139  performed and subsequently discarded.

140  **Table 2**

141  **4.2   CHO cell culture**

142  CHO-GS cells were cultured in FORTICHO (Gibco CD FortiCHO 10887640) at 37°C, 5% $CO_2$,
143  85% humidity, 125rpm with 25mm orbit in a shaking incubator in 125mL bioreactor flasks (Nalgene
144  10266432). Every 3-4 days, cells were passaged at $0.2*10^5$ cells/mL in 30mL media in 125mL
145  culture shaking flasks. A growth curve was established to ensure samples were taken at exponential
146  cell phase (**Table 1**).

147  **4.3   DAPI Stain**

148  A working concentration of 0.1ug/mL DAPI was determined as optimal for CHO cells. DAPI
149  (Invitrogen D1306) solutions were protected from light wherever possible. In a BSC, 10mg DAPI
150  powder was completely dissolved in 2mL sterile deionised water to make a 5 mg/mL DAPI stock
151  solution. This was aliquoted and stored at -20°C. Solutions were stable for at least six months. 1uL of
152  DAPI stock solution was added to 5mL DPBS for a 1ug/mL stock 2 DAPI working solution. 1mL of
153  1ug/mL stock 2 solution was added to 9mL of DPBS to prepare a 0.1ug/mL DAPI working solution.

154  **4.4   Staining Cells**

155 Here, an AB-FITC (Sigma Aldrich SAB3701254-2MG) conjugate was used which at 4°C can bind to
156 IgG on the cell membrane in the process of being excreted by the cell as previously demonstrated for
157 CHO cells (Gallagher and Kelly, 2017). This allowed the sorting of cells based on the productivity of
158 an IgG-based antibody. Cell samples were prepared as per **Table 1**. Cells were counted using trypan
159 blue (Gibco 15250061) and a hemacytometer as per the manufacturer's instructions. 1*10^6 of viable
160 cells were centrifuged at 200 x g for 5 minutes and supernatant discarded. Cells were washed in 1mL
161 DPBS, centrifuged at 200 x g for 5 minutes and supernatant discarded. This was repeated for a total
162 of 2 washes. Cells were resuspended in 1mL of DPBS using 2uL of anti-human IgG (Sigma Aldrich
163 SAB3701254-2MG). Cells were incubated at 4°C for 30 mins at 1000rpm, protected from light. Cells
164 were washed twice with DPBS as per steps 6 & 7 for a total of 2 washes. Cells were resuspended in
165 1mL of cold DPBS or cold DAPI working solution, incubated on ice for 5 mins and immediately
166 transferred on ice to the Fluoresence-activated cell sorting (FACS) lab for immediate analysis.

### 4.5   Setting Single cell Gating

168 The FACS Melody was setup as per manufacturer's instructions. A U-bottom 96-well plate was
169 prepared (Corning 3799) with 5uL of 1x TCL buffer (QIAGEN 1070498) in the centre of each
170 functional well using a multichannel pipette. The plate was tapped firmly on a flat surface to
171 encourage the central location of the TCL buffer. The size threshold was set to >12um. Using Sample
172 4 (**Table 1**), voltages were set to allow the representation of cells in a SSC-A against FSC-A
173 logarithmic scale graph. Gate 1 (G1) excluded instrument noise and cell debris as per **Fig. 2A**. Using
174 Sample 4, data was brought forward and gate 2 (G2) set using FSC-H against FSC-A as per **Fig. 2B**
175 to exclude doublets. Using Samples 1 and 2, the G2 gate was brought forward and a range gate (G3)
176 was set to only include live cells as per **Fig. 2C** and **2.4D**. DAPI positive was considered dead cells.
177 Using Samples 1 and 3, the G3 gate was brought forward and a gate (G4) set for FITC-positive cells
178 as per **Fig. 2E** and **2F**. G4 was the sorting gate for live, singlet cells. After gates had been set, it was
179 important to record data for a large number of events (e.g. 10,000 cells) and to save FCS files.

### 4.6   Single Cell Sort

181 The flow rate was kept at a minimum to reduce the chance of doublets. Sample 1 was loaded FACS
182 set to "single cell" and "96-well plate" modes. Desired wells were selected for sorting with a splash
183 shield present. The lid was removed and immediately inserted into the FACS to proceed with sorting.
184 For the positive control, sort mode was changed to "purity". After the sort was complete, the well
185 plate was removed and immediately covered with the lid. An airtight seal was created around the
186 edges with parafilm and the plate immediately placed in a -80°C freezer. FCS files were saved for all
187 samples.

188 STOPPING POINT: Samples can be stored for up to 6 months at -80°C.

### 4.7   AMPure purification

190 AMPure beads (10136224) benefit from scalable purification – adapting to single cell samples,
191 volumes can simply be reduced. Sequences of mtDNA sometimes migrate and integrate into the
192 nuclear genome – known as Nuclear mitochondrial sequences (Numts). In a previous bulk analysis of
193 mtDNA, the miniprep step purified the plasmid-like mtDNA from contaminating linear nuclear DNA
194 (Kelly et al., 2017). Here, AMPure purification was used, leaving both mtDNA and nDNA in the
195 sample. Blast searching primer sequences against the CHO cell line reference genome (taxid: 10029)
196 and gel purifying 8.5kb bands provided additional protection against Numts. All steps were

197 performed in a BSC. The subsequent LRPCR is extremely sensitive and could potentially amplify
198 small contaminations. The 96-well plate was thawed at room temperature. Multiple samples were
199 taken through AMPure purification in batches (to a maximum of 12 samples). The 5uL lysed sample
200 was transferred to a labelled micro-centrifuge tube. AMPure beads were resuspended by vortexing
201 the bottle for 1 min. 9uL of AMPure beads were added per sample (if the lysed cell sample was
202 greater, 1.8x volume of AMPure beads was used) and pipette mixed 10 times. They were left at room
203 temperature for 5 mins. Tubes were placed on a magnetic stand (New England Biolabs S1515S) for 2
204 minutes. Keeping the tubes on the magnetic stand, the cleared solution was removed and discarded
205 leaving the beads. It was then washed with 40uL 70% ethanol. The supernatant was discarded,
206 leaving the beads. The ethanol wash was repeated. On the second wash, remaining ethanol was
207 removed by using a P10 pipette while avoiding removing any beads. Tubes were removed from the
208 magnetic stand and 18uL of elution buffer (QIAGEN 19086) was added to the bead aggregate and
209 pipette mix 10 times or until fully resuspended. Tubes were incubated for 5 minutes at room
210 temperature. Tubes were placed on the magnetic stand for 2 minutes. Eluate was split into two 8.5uL
211 aliquots – leaving the bead aggregate. Microcentrifuge tubes were labelled to identify which samples
212 came from the same single cell.

213 **4.8   SuperFi II Plat Taq LRPCR**

214 The bottleneck of single cell sequencing is the DNA amplification. Amplification techniques that
215 would work for bulk sequencing proved to be incompatible with single cell: mechanical purifications
216 took too much of the sample, bacterial lysis buffers did not release enough mtDNA, and components
217 lost effectiveness for the sensitive PCR. Once enough DNA is amplified, established protocols for
218 bulk sequencing can be followed (Kelly et al., 2017). The SuperFi II PCR kit (Invitrogen 12361010)
219 has 300x fidelity compared to Plat Taq. SuperFi II was better able to amplify from small samples
220 compared to Plat Taq. Higher fidelity also means greater confidence in lower-level heteroplasmy.

221 In addition to the below LRPCR protocol, single cell samples post AMPure purification were diluted
222 to 1/10, 1/100, 1/1,000 and 1/10,000; to demonstrate the limits of the high fidelity LRPCR kit. All
223 steps were performed in a BSC while maintaining samples at all steps on ice. Fresh aliquots of
224 primers and dNTPs were thawed at room temperature then stored on ice. SuperFi II 5x Buffer was
225 thawed and stored on ice. DNA Polymerase was maintained at -20°C and only removed briefly when
226 needed. Components were briefly vortexed and centrifuged before use – except for the DNA
227 Polymerase. mtDNA LRPCR was performed in 2 separate fragments (termed X and Y). The eluate
228 from a single cell had been split into 2 from AMPure Purification; 1 half was amplified using X
229 primers, the other half by Y primers (**Table 2**). A mastermix was generated with 10% overage, for
230 each X primer and Y primer, as per the example in **Table 3**. SuperFi II DNA Polymerase was added
231 last by briefly removing it from the -20°C freezer – to minimise the time spent at room temperature.
232 The mastermix was gently vortexed, centrifuged at 500xg and kept on ice. 16.5uL mastermix was
233 added to 8.5uL AMPure purified DNA. The sample was gently vortexed, centrifuged at 500xg and
234 kept on ice. Samples were placed in a PCR machine and set to a PCR cycle as per **Table 4**. Reaction
235 volume was set to 25uL with a lid temperature of 105°C. The cycle was run overnight. On
236 completion, samples were removed and stored at 4°C.

237 STOPPING POINT: Samples can be stored at -20°C for 2 weeks.

238 **Table 3**

239 **Table 4**

240 **4.9    Agarose Gel**

241   Limit of detection: an 8.5kb band was still observable when taking a 1/1000 dilution of a single cell.
242   You would expect around 100-10,000 mitochondria (Dhiman et al., 2019). Theoretically, this PCR
243   may be on the edge of viability for single-mitochondrial sequencing.

244   1g of agarose was added to 100mL TAE buffer in a conical flask and microwaved for 2.5 mins, or
245   until fully dissolved. The flask was left to cool to about 50°C. 10uL of SafeView (NBS Biologicals)
246   was added and the mixture poured into the gel tray with a well comb. After a brief period, the gel
247   cooled and hardened at room temperature. The gel was placed in a gel box with TAE buffer just
248   covering the gel. Loading dye was added to all samples as per the manufacturer's instructions. Entire
249   samples were loaded into gel wells with an appropriate DNA ladder (Thermo Scientific SM1333).
250   Gels were run at 100V until bands were 70% down the gel. The power was turned off and the gel was
251   carefully placed in a gel viewer. Photos of the gel were taken.

252   **4.10   Gel excision**

253   Under a gel visualiser, 8.5kb bands were identified indicative of single cell reactions, as illustrated by
254   red rectangles in **Fig. 2G**. UV light exposure was minimised to limit degradation of DNA. Using a
255   new sterile disposable scalpel, the 8.5kb band was excised and placed in a 1.5mL Eppendorf tube.
256   The blade was thoroughly cleaned with 70% IPA and then reused. The Y single cell sample was
257   equally isolated and placed in a separate 1.5mL Eppendorf tube.

258   **4.11   Gel purification**

259   The QIAquick Gel Extraction protocol for "QIAquick Gel Extraction using a Microcentrifuge". 10uL
260   of elution buffer was used to encourage a higher final concentration.
261   STOPPING POINT: Samples can be stored at -20°C for 2 weeks.

262   **4.12   Equimolar combination**

263   The Qubit 1x dsDNA HS kit (Invitrogen Q33230) was used to quantify dsDNA. Kit components
264   were allowed to equilibrate to room temperature for 30 mins. 10uL of Standard 1 was added to a
265   Qubit tube, and 10uL Standard 2 to a separate Qubit tube. 190uL of 1x buffer was added to each. 1uL
266   of each X and Y fragments was added to the separate Qubit tubes and 199uL of 1x buffer was added
267   to each. Tubes were vortexed for 2-3 seconds and left at room temperature for 2 mins. The
268   concentration of standards 1 and 2 were measured using the Qubit (Invitrogen). The concentration of
269   samples was measured using the Qubit. The volume required to aliquot 1ng of the X fragment and Y
270   fragment from the same cell was calculated and these volumes combined in a new Lo-bind tube.
271   There should be a total of 2ng of mtDNA from each single cell. Low concentrations were expected
272   from our single cell samples.

273   **4.13   Library Prep**

274   The Illumina DNA Prep protocol was followed, using IDT for Illumina DNA/RNA UD Indexes Set
275   A, Tagmentation (96 Indexes, 96 Samples) (Illumina 20027213). Each single cell should have a
276   unique pair of indexes. The library quality of the cleaned-up library was checked by running 1uL on
277   a Tapestation D5000 microwell. The libraries were combined and diluted to a 2nM starting
278   concentration as per the manufacturer's instructions.
279   STOPPING POINT: Samples can be stored at -20°C for 30 days.

### 280  4.14  Sequencing

281  Libraries generated using Illumina DNA Prep were compatible with a wide range of Illumina
282  sequencers including HiSeq, iSeq100, MiniSeq, NextSeq and NovaSeq technologies.

283  The iSeq cartridge and flow cell were prepared as per the manufacturer's instructions (Illumina
284  20031371). 2% PhiX (Illumina FC-110-3001) spike-in was added. The sample sheet loaded onto
285  iSeq was checked to ensure correspondence to the sample sheet from Library Prep. The cartridge was
286  loaded, and the run performed as per manufacturer's instructions. After running, the data was
287  downloaded and backed-up on an external hard drive.

### 288  4.15  Data preprocessing

289  GitHub repository: https://github.com/alanfoleynibrt/SingleCellmtDNA

290  The bioinformatics pipeline is available in the above GitHub repository. Initial processing of data is
291  performed in Linux and figures are generated in R. All raw FASTQ data analysed is made available
292  in this pipeline. A step-by-step protocol and all materials are also available.

293  Briefly, trim_galore (0.4.3) used to trim adapter sequences in FASTQ files. Bowtie-2 (2.3.4.1) used
294  to map reads to the KX576660.1 CHO mtDNA reference genome. Picard (1.199) tools identified
295  duplicates (MarkDuplicates), added read groups (AddOrReplaceReadGroups) and built a BAM index
296  (BuildBamIndex). Gatk3.8-0 implemented to realign indels (IndelRealigner) and recalibrate bases
297  (BaseRecalibrator). Two separate mutation calling software programs were used: lofreq_star-2.1.2
298  and varscan.v2.3.9. When a mutation was called by both, it was brought forward for analysis. If a
299  mutation allele frequency was between 0.04 and 0.96, it was considered "heteroplasmic". The
300  potential impact of identified mutations was predicted using SnpEff. In tandem, analysis was
301  repeated using a shifted mtDNA reference genome to complete coverage over the D-loop region.
302  Unshifted mutation calls were concatenated with those from the shifted reference sequence to provide
303  full coverage. ggplot2 in R was used to generate figures.

304

## 305  5    Results

### 306  5.1  Single Cell Isolation

307  The overall aim of this project was to create a workflow for single cell mtDNA analysis in CHO
308  cells. We first sought to isolate single cells in a reproducible manner. Clonal populations of CHO
309  cells are often generated using the "serial dilution" method whereby a known number of cells is
310  progressively diluted to approximate a single cell per unit volume. However, the nature of cell
311  distribution in each dilution means the final dispensed sample could indeed have a single cell, but it
312  could also have 0 cells or multiple cells. We reasoned that a FACS-based method would be more
313  accurate and reproducible. Additionally, there is less manual work when scaling to generate large
314  numbers of clones. We loaded the sorter (FACS Melody, BD) with CHO-GS cells which had been
315  stained with an AB-FITC conjugate and DAPI.

### 316  Fig. 2

317  We followed basic guidelines for flow cytometry (Bio-Rad, 2022). We first used an SSC-A against
318  FSC-A dot plot with calibrated voltages (**Fig. 2A**). A gate of the main cell population was selected,

319   excluding instrument noise and cell debris. We next used an FSC-H against FSC-A graph (**Fig. 2B**).
320   Since forward scatter determines the "size" of particles, the "height" against the "area" determines
321   the ratio of the cell height against the cell area. A singlet will have 1x area, 1x width and 1x height.
322   Doublets would have 2x area, 2x width but 1x height. Thus to discriminate between singlets and
323   doublets, the ratio of area to height is considered (Bio-Rad, 2022). We used a gate to select only
324   singlets.

325   We next focused on FITC staining (due to our AB-FITC conjugate) and BV510-A (due to DAPI).
326   DAPI increases in fluorescence when binding to DNA; therefore it has applications for live/dead cell
327   gating in flow cytometry. Live cells would have lower BV510-A fluorescence, while dead cells
328   would have greater BV510-A fluorescence. We first loaded a population of dead/dying CHO cells
329   with 5% viability and viewed the population in a histogram of BV510-A fluorescence (**Fig. 2C**). We
330   calibrated the BV510-A voltage to allow space for "lower" BV510-A fluorescence to which we set a
331   5% gate. Thus, only 5% of the dead/dying population was within our lower BV510-A fluorescence
332   gate. We then ran live cells with 95% viability (**Fig. 2D**). As expected, the population migrated to our
333   "live" gate, thus allowing us to select for live cells.

334   To identify mAb-secreting CHO cells, we exploited the AB-FITC staining to select for cells with a
335   mAb in "stasis" in the CHO cell membrane. We first ran a sample with a non-producing CHO cell
336   line and ran the population on a histogram of FITC fluorescence (**Fig. 2E**). If the stain specifically
337   binds to mAb-producing CHO cells, there should be a shift of the population to greater FITC
338   fluorescence from non-producing to producing. We therefore set a gate for greater FITC fluorescence
339   with 0% of cells from the non-producing cell line. When we ran the same settings for our producer
340   CHO cell line (**Fig. 2F**), we observed an overall "shift" of the population towards greater FITC
341   fluorescence. This finalised our FACS-based method which selects for viable, singlet and mAb-
342   producing CHO cells. Cells that fulfilled our gating strategy were sorted into wells of a 96-well plate
343   with lysis buffer; choosing a U-bottom plate to encourage a central location for the 5uL of lysis
344   buffer within each well.

**5.2   DNA Purification**

346   Previous bulk analysis of CHO mtDNA used mini-prep kits (QIAGEN) to enrich the plasmid-like
347   mtDNA and reduce capture of the linear nuclear DNA (Kelly. P et al., 2017). This method proved to
348   be inefficient for single cell samples whose mtDNA mass was much smaller. The physical filtration
349   system required quenching meaning that low initial volumes, as with single cell samples, were lost.
350   AMPure bead purification emerged as a viable alternative since adaptation to lower volume samples
351   simply required volume reduction of all reagents. Here, purification was performed by adding
352   AMPure beads to single cell samples, applying a magnet, washing with ethanol and eluting with
353   elution buffer. However, caution is advised since nuclear DNA is also captured.

354   The miniprep kit had concomitantly provided some protection against Numts since it is designed to
355   purify circular mtDNA away from linear nuclear DNA (Kelly. P et al., 2017). Having eliminated the
356   miniprep step, we sought to incorporate additional protection against Numts. We performed a
357   BLAST search of our mtDNA amplification primer sequences against the nuclear CHO reference
358   genomes and found no matches, suggesting there are no nuclear sequences to which our primers
359   should bind. Also, Numts tend to be shorter sequences, with 78% shorter than 500bp in human
360   mtDNA (Wei et al. 2022). Therefore, we reasoned that specific gel purification of 8.5kb amplicons
361   would be unlikely to be contaminated with Numts.

### 5.3    DNA Amplification

The biggest bottleneck for single cell mtDNA sequencing was the DNA amplification step. Adaptation of the DNA amplification centred on compatibility with as much as a million-fold lower starting DNA material compared to extracting from a standard cell population sample. We first selected a high-fidelity, long-range LR-PCR kit (SuperFi II. Invitrogen 12361010)). We then designed primers to amplify the mtDNA in two separate overlapping fragments, as previously done for bulk analysis (Kelly et al., 2017). The mini-prep kit, as previously used by Kelly et al., utilised a bacterial-specific lysis method. We anticipated that a mammalian-specific lysis buffer might liberate more mtDNA than a bacterial lysis buffer, accounting for the presence of both the outer cell membrane and the double-membrane of the mitochondria (**Fig. 1A**).

We found that PCR component storage and utilisation was a critical element of the protocol being successful. Immediately on component arrival, dNTPs and primers were diluted to the desired concentrations, aliquoted into microtubes and stored at -80°C. For each LRPCR reaction, a new aliquot was thawed, used and the aliquot discarded. As a precaution, we sorted 0 cells into a 96-well plate (**Fig. 2G**) to test for sources of contamination during the protocol that could lead to non-specific amplification. We did not observe any amplicons from these wells. We then applied the method to a single cell sample, and also to a 1000 cell sample (**Fig. 2G**). We observed successful amplification of both amplicons of mtDNA from both 1000 cells and a single cell.

### 5.4    Sample Generation and Library Preparation

Once amplification of mtDNA from single cells was achieved, we reasoned the later steps of library preparation, sequencing and bioinformatic analysis should be largely unchanged. To verify this, four single cells and a bulk sample (4000 cells) for comparison were sorted into 5uL of TCL lysis buffer. Samples were purified, split into two equal portions and separately amplified by LRPCR; which were then visualised on agarose gel. mtDNA-specific bands were excised for both amplified fragments. Amplicons were recovered by gel purification. After quantifying this dsDNA using the Qubit 1x dsDNA HD kit, equimolar quantities of each fragment from the same cell were added into a single tube.

For library preparation, the Illumina DNA Prep protocol (20018705) with IDT for Illumina DNA/RNA UD Indexes Set A was implemented. Each single cell was separately assigned indices to be used to demultiplex single cells later. Separate libraries were combined, diluted to 2nM and loaded onto the Illumina iSeq for PE150 sequencing. A 2% PhiX library control spike-in was added. The iSeq had the option to include a "sample sheet" to which the index combinations were added. The iSeq was run to completion with an output of fastq.gz files ready for the bioinformatics pipeline.

### 5.5    Bioinformatic Analysis

Adapter sequences were removed from the reads which were then mapped against the CHO KX576660.1 mtDNA reference genome. During the PCR step of the Illumina DNA library prep, adapter-ligated fragments are PCR amplified. This can lead to multiple sequencing reads deriving from the same original fragment; possibly resulting in overrepresentation of certain alleles. As was performed previously in bulk analysis of CHO mtDNA, PCR duplicates were identified (**Fig. 3A**) and removed from the analysis (Kelly et al. 2017). The range of duplicate reads was 23.8% to 29.3% with the highest proportion in the mixed population. Further, indel realignment and base recalibration were used to cater for the effects of INDELS on read mapping.

404    **Fig. 3**

405    After excluding duplicate and unmapped reads, all samples had an average sequencing depth of
406    >1500x – above the 1000x required for "ultra-deep sequencing" categorisation (**Fig. 3B**). Perbase
407    coverage of all samples confirmed complete and even mapping of sequencing reads (**Fig. 3C**). The
408    mapping indicated no strong bias for any particular region. Together, this confirmed our
409    scmtDNAseq protocol had been successful.

410    The great value of single cell sequencing mtDNA at such great depth is the potential to analyse at
411    high confidence the differences in the sequenced reads when compared to the reference genome; i.e.
412    a mutation. The multi-copy nature of mtDNA within each cell is represented by the proportion of
413    reads with a particular mutation. E.g. If 50/100 reads contain a mutation, the allele frequency is
414    determined as 0.5. This infers that out of all copies of mtDNA within that cell, 50% would contain
415    the mutation.

416    As was performed for the previous bulk analysis, LoFreq and VarScan were used to call mutations
417    (Kelly et al., 2017). When both tools called a mutation above 0.04, it was brought forward for
418    analysis. A total of 43 mutations were called among the 4 samples of which 17 were indels and 26
419    SNPs (**Fig. S1B**). Of mutations heteroplasmic in the bulk sample, the single cell average allele
420    frequency varied dramatically from bulk (**Fig. 4A**). For example, the 5462T>C mutation was 0.05 in
421    the bulk, but over 0.5 in the single cell average. This is likely a consequence of the small number of
422    single cell samples sequenced but also suggests that some mutations may exist sporadically at a high
423    frequency in a small number of cells within the population (scenario 2 in **Fig. 5**).

424    **Fig. 4**

425    To better assess the variability in allele frequency among single cells, we developed a list of "most
426    variable" mutations which must be heteroplasmic in at least 2 out of the 4 cells. There was a wide
427    range of allele frequencies among the single cells (**Fig. 4B**). Had only bulk analysis been performed,
428    this range of allele frequency would not have been captured; demonstrating the enhanced resolution
429    possible from single cell analysis.

430    We next considered whether the mutations among the samples were concordant. We generated a
431    heatmap of mutations from all samples, with each mutation type represented by a colour (**Fig. 4C**).
432    The mutations 5244TA>T and 14136T>A were present in all samples. Considering the nature of a
433    bulk sample, one might expect that all the mutations present in the single cells should be present in
434    the bulk sample. However, for the mutation to be called with confidence in the bulk sample, it must
435    be above 0.04; and it must therefore be present in individual cells in the population at sufficient
436    frequency to average above 0.04. All mutations from the bulk sample were found in at least 1 single
437    cell; apart from the intragenic 3733G>A. The limited number of single cells analysed likely resulted
438    in this observation. On the other hand, 73% (19/26) of mutation locations were exclusively found in
439    the single cells but not in the bulk sample, again evidencing the high degree of resolution from single
440    cell analysis.

441    The lowest number of mutations were in the bulk population with 13 (joint with Single cell 1) (**Fig.
442    4C**). This was not unexpected because single cells can contain additional rare mutations that average
443    out below 0.04 in the bulk sample. On the other hand, individual single cells do not necessarily
444    contain all the mutations found in the bulk sample. Ultimately, the greater the number of single cells

445 analysed, the greater the certainty about the nature of individual mutations across a population of
446 cells when compared to bulk cell sequencing.

447 **5.6    Predicted Phenotypic Impact of Mutations**

448 Although the presence of mutations is useful for demonstrating intercellular diversity, their
449 phenotypic impact may be limited by many factors including the mutation's effect on, for example,
450 amino acid sequence or tRNA structure. We therefore analysed mutations based on predicted impact
451 on phenotype using the snpEff software tool (**Fig. 4D**) (**Fig. S1D**). Of particular interest, frameshift
452 mutations were observed in protein-coding genes COX1, CYTB and ND4 (**Fig. S1E**). At least 1
453 frameshift mutation was called in all samples. Only 4 heteroplasmic mutations were above 0.5 allele
454 frequency, with most mutations present at low levels. Mapping of mutations against allele frequency
455 showed the vast majority to be in the 0.04-0.5 range (**Fig. S1F**).

456

457 **6    Discussion**

458 The phenotypic manifestations of disease typically occur when the responsible mtDNA mutation
459 allele frequency reaches a certain threshold within a cell. Previous bulk analysis of CHO mtDNA
460 identified heteroplasmy in clones derived from the same parental host, indicating at least three levels
461 of heterogeneity: (1) production run to production run, (2) cell line to cell line and (3) clone to clone
462 (Kelly. P et al., 2017). However, bulk analysis of heteroplasmy fails to identify the allele frequency
463 differences between individual cells in a population (**Fig. 1A**). Single cell analysis is therefore critical
464 to reveal the true phenotypic effect of heteroplasmy.

465 To illustrate this, consider if there was a hypothetical mutation with a threshold of 0.7, above which
466 phenotypic changes would manifest in an individual cell. If bulk analysis identified this critical
467 mutation at a frequency of 0.1 in 1 million cells (**Fig. 5**), three very different conclusions could be
468 arrived at: (1) All cells contain the mutation at a 0.1 allele frequency, and are therefore all unaffected
469 phenotypically; (2) 10% of cells contain the mutation at a 1.0 allele frequency, and therefore only
470 10% of cells are affected phenotypically or (3) cells contain the mutation at a variable rate (0-1.0),
471 and therefore the population is affected at a variable rate.

472 **Fig. 5**

473 Bearing in mind the strive for homogeneity in drug production, the implications of these scenarios
474 are significant. If a particular heteroplasmy profile affected product quality for example, perhaps only
475 a subset of the cells produce the product at a high quality; in which case the remainder could be
476 identified and potentially excluded to improve bioreactor performance. Equally, perhaps a cell line
477 could be engineered with a favourable heteroplasmy profile to improve bioreactor performance.
478 Further work is clearly needed to understand the link between mitochondrial heteroplasmy and
479 cellular behaviour in recombinant protein production, but single cell analysis should contribute
480 significantly in this regard.

481 Though the 5 samples here are demonstrative, and not enough for strong statistical conclusions,
482 certain observations were made. The bulk population had the lowest number of reliably detectable
483 mutations (**Fig. S1B**). All mutations in the bulk population, bar one, were found in at least one single
484 cell (**Fig. 4C**). This demonstrated the improved resolution of mutation detection using a single cell
485 approach. A great range of allele frequencies in "most variable" mutations was observed (**Fig. 4B**);

486 further indicating an uneven spread of heteroplasmy among the 4 cells, reminiscent of scenario 3 in
487 **Fig. 5**.

488 High-impact mutations observed here in *CYTB* (**Fig. S1D**) would change the encoded amino acid
489 sequence. The phenotypic effects of *CYTB* mutations are well established in human disease where
490 patients experience highly variable severities of myopathy and muscle weakness (Blakely et al.,
491 2005). *CYTB* mutations in yeast models can cause severe decreases in respiratory function (Fisher et
492 al., 2004). In a bioreactor, *CYTB* mutated single cells (above a phenotypic threshold) may be one of
493 many contributing factors to the heterogeneity observed among clonally derived CHO populations.

494 In conclusion, a reliable method to amplify and analyse mtDNA from single CHO cells was
495 demonstrated (scmtDNAseq). This approach should help better understand the degree and likely
496 impact of heteroplasmy on recombinant protein production in CHO cells.

497

498 **7    Conflict of Interest**

499 *The authors declare that the research was conducted in the absence of any commercial or financial*
500 *relationships that could be construed as a potential conflict of interest.*

501

502 **8    Author Contributions**

503 Alan Foley: Developed novel method under supervision of Niall Barron. Adapted bioinformatics
504 code from Colin Clarke, and wrote manuscript.

505 Niall Barron: Supervisor of project. Guided method development. Manuscript review/editing.

506 Colin Clarke: Co-supervisor. Wrote original code for mtDNA analysis which was later adapted.
507 Manuscript review.

508

509 **9    Funding**

512

513 **10    Acknowledgments**

514 NIBRT, Ireland: a base for method development and bioinformatics.

515 Niall Barron: Supervised project.

516 Colin Clarke: Co-supervisor. Wrote original code for mtDNA analysis which was later adapted.

517    Nga Lao: Significant aid in wet lab work.

518

519    **11    References**

520     Bio-Rad (2022). Doublet Discrimination - Flow Cytometry Guide. Bio-Rad. Available at:
521    https://www.bio-rad-antibodies.com/flow-cytometry-doublet-discrimination.html [Accessed January
522    17, 2022].

523     Blakely, E. L., Mitchell, A. L., Fisher, N., Meunier, B., Nijtmans, L. G., Schaefer, A. M., et al.
524    (2005). A mitochondrial cytochrome b mutation causing severe respiratory chain enzyme deficiency
525    in humans and yeast. FEBS J 272, 3583–3592. doi: 10.1111/j.1742-4658.2005.04779.x.

526     Dhiman, H., Gerstl, M. P., Ruckerbauer, D., Hanscho, M., Himmelbauer, H., Clarke, C., et al.
527    (2019). Genetic and Epigenetic Variation across Genes Involved in Energy Metabolism and
528    Mitochondria of Chinese Hamster Ovary Cell Lines. Biotechnology Journal 14, 1800681. doi:
529    10.1002/biot.201800681.

530     Dimauro, S., and Davidzon, G. (2005). Mitochondrial DNA and disease. Ann Med 37, 222–232.
531    doi: 10.1080/07853890510007368.

532     Fisher, N., Castleden, C. K., Bourges, I., Brasseur, G., Dujardin, G., and Meunier, B. (2004).
533    Human disease-related mutations in cytochrome b studied in yeast. J Biol Chem 279, 12951–12958.
534    doi: 10.1074/jbc.M313866200.

535     Gallagher, C., and Kelly, P. S. (2017). Selection of High-Producing Clones Using FACS for CHO
536    Cell Line Development. Methods Mol Biol 1603, 143–152. doi: 10.1007/978-1-4939-6972-2_9.

537     Gilbert, A., McElearney, K., Kshirsagar, R., Sinacore, M. S., and Ryll, T. (2013). Investigation of
538    metabolic variability observed in extended fed batch cell culture. Biotechnol Prog 29, 1519–1527.
539    doi: 10.1002/btpr.1787.

540     Kelley, B., Kiss, R., and Laird, M. (2018). A Different Perspective: How Much Innovation Is Really
541    Needed for Monoclonal Antibody Production Using Mammalian Cell Technology? Adv Biochem
542    Eng Biotechnol 165, 443–462. doi: 10.1007/10_2018_59.

543     Kelly, P. S., Clarke, C., Costello, A., Monger, C., Meiller, J., Dhiman, H., et al. (2017). Ultra-deep
544    next generation mitochondrial genome sequencing reveals widespread heteroplasmy in Chinese
545    hamster ovary cells. Metabolic Engineering 41, 11–22. doi: 10.1016/j.ymben.2017.02.001.

546     Keogh, M. J., and Chinnery, P. F. (2015). Mitochondrial DNA mutations in neurodegeneration.
547    Biochimica et Biophysica Acta (BBA) - Bioenergetics 1847, 1401–1411. doi:
548    10.1016/j.bbabio.2015.05.015.

549     Lorenz, T. C. (2012). Polymerase chain reaction: basic protocol plus troubleshooting and
550    optimization strategies. J Vis Exp, e3998. doi: 10.3791/3998.

551 Ludwig, L. S., Lareau, C. A., Ulirsch, J. C., Christian, E., Muus, C., Li, L. H., et al. (2019). Lineage
552 Tracing in Humans Enabled by Mitochondrial Mutations and Single-Cell Genomics. Cell 176, 1325-
553 1339.e22. doi: 10.1016/j.cell.2019.01.022.

554 Maeda, R., Kami, D., Maeda, H., Shikuma, A., and Gojo, S. (2020). High throughput single cell
555 analysis of mitochondrial heteroplasmy in mitochondrial diseases. Sci Rep 10, 10821. doi:
556 10.1038/s41598-020-67686-z.

557 NCBI (2023). National Center for Biotechnology Information. Available at:
558 https://www.ncbi.nlm.nih.gov/ [Accessed June 5, 2023].

559 Walsh, G., and Walsh, E. (2022). Biopharmaceutical benchmarks 2022. Nat Biotechnol 40, 1722–
560 1760. doi: 10.1038/s41587-022-01582-x.

561 Wei, W., Schon, K. R., Elgar, G., Orioli, A., Tanguy, M., Giess, A., et al. (2022). Nuclear-
562 embedded mitochondrial DNA sequences in 66,083 human genomes. Nature 611, 105–114. doi:
563 10.1038/s41586-022-05288-7.

564 Zambelli, F., Vancampenhout, K., Daneels, D., Brown, D., Mertens, J., Van Dooren, S., et al.
565 (2017). Accurate and comprehensive analysis of single nucleotide variants and large deletions of the
566 human mitochondrial genome in DNA and single cells. Eur J Hum Genet 25, 1229–1236. doi:
567 10.1038/ejhg.2017.129.

568

## 12    Data Availability Statement

570 The datasets for this study can be found in the GitHub repository:
571 https://github.com/alanfoleynibrt/SingleCellmtDNA

572

## 13    Tables

574 **Table 1:** Samples 1-4 required for single cell sort. In this data, Late Exponential had a viability of
575 95%, Dead of 5%.

|   | Cells | Growth Phase | Stain | Function |
|---|---|---|---|---|
| 1 | Protein Producing CHO | Late Exponential | DAPI + FITC-AB | Sorting Sample |
| 2 | Protein Producing CHO | Dead | DAPI + FITC-AB | Gate Live/Dead Cells |
| 3 | Non-producing CHO | Late Exponential | DAPI + FITC-AB | Gate FITC negative |
| 4 | CHO | Late Exponential | None | Gate FSC, SSC. Gate FITC positive |

576

**Table 2**: Primer sequences for LRPCR. Other cell lines may need adaptation of these sequences.

| Primer | Sequence | |
|---|---|---|
| mt-490 F (X) | 5' - GGA TTA GAT ACC CCA CTA TGC TT – 3' | 578 |
| mt-9304 R (X) | 5' – ATG CTG CGG CTT CAA ATC CG – 3' | 579 |
| | | 580 |
| mt-9180 F (Y) | 5' – ATA GCA ACA GGT TTT CAC GG – 3' | 581 |
| | | 582 |
| mt-598 R (Y) | 5'- CGC CAA GTC CTT TGA GTT TTA – 3' | 583 |

584

**Table 3:** PCR components

| Reagent | Volume per Rx (uL) | 10x Mastermix (uL) X | 10x Mastermix (uL) Y |
|---|---|---|---|
| 5x Buffer | 5 | 50 | 50 |
| 10mM DNTP mix | 0.5 | 5 | 5 |
| 10uM Primer F | 1 | 10 (X primer) | 10 (Y primer) |
| 10uM Primer R | 1 | 10 (X primer) | 10 (Y primer) |
| Nuclease-free H2O | 8.5 | 85 | 85 |
| SuperFi II DNA Polymerase | 0.5 | 5 | 5 |
| **TOTAL** | **16.5** | **165** | **165** |

586

**Table 4:** PCR settings

| Step | Temperature (°C) | Time |
|---|---|---|
| | | |

| 1. Initial Denature | 94 | 2 mins |
|---|---|---|
| 2. (x35) Denature | 94 | 30 s |
| Annealing | 55 | 30 s |
| Extension | 68 | 9 mins |
| 3. Final extension | 68 | 10 mins |
| 4. Hold | 4 | Infinite hold |

588

589

590 **14    Figure Legends**

591 **Fig. 1**: (**A**) An explanation of the multi-copy nature (heteroplasmy) of mitochondrial DNA. Numbers
592 are true for CHO cells, though vary by eukaryotic cell type. (**B**) Method overview for scmtDNAseq.
593 Made with BioRender.

594 **Fig. 2**: (**A-F**) Gating strategy to sort alive, singlet, antibody-producing CHO cells. (**G**) Agarose gel
595 illustrating amplification of CHO cell mtDNA from a single cell. Also included is positive control of
596 1000 cells and negative control of 0 cells. "X" and "Y" refers to the 2 separate halves of the mtDNA
597 molecule. Together, 1X and 1Y represent amplification of the whole mtDNA molecule from a single
598 cell in two separate reactions. Red rectangles illustrate gel extraction boundaries to exclude bands
599 other than the desired 8.5kb amplicon. Made with BioRender.

600 **Fig. 3:** (**A**) Read mapping of samples to the KX576660 CHO mtDNA reference genome. (**B**)
601 Average per base sequencing depth of each sample. (**C**) Perbase coverage of 4 single cells and a bulk
602 sample with correction around 0 coordinate. Made with BioRender.

603 **Fig. 4:** (**A**) Comparison of bulk sample allele frequencies to the average of 4 single cells' allele
604 frequencies. The mutation must be heteroplasmic (between 0.04 and 0.96 allele frequency) in the
605 bulk sample. (**B**) Violin plot of "most variable" mutations which must be heteroplasmic in at least 2
606 of 4 single cells. (**C**) Base change heatmap of heteroplasmic mutations (between 0.04 and 0.96 allele
607 frequency) in 4 single cells and a bulk sample. (**D**) snpEff predicted impact of heteroplasmic
608 (between 0.04 and 0.96 allele frequency) mtDNA mutations of 4 single cells and a bulk sample.
609 Made with BioRender.

610 **Fig. 5**: How bulk analysis of heteroplasmy can obfuscate single cell orientations. Made with
611 BioRender.

612 **Fig. S1**: (**A**) Tapestation (4200) image of mtDNA LRPCR for dilutions of a single cell. Negative
613 control is 0 cells, positive control is 10 cells. Dilutions were made from 1/10 to 1/100,000. (**B**)

614    Number of heteroplasmic mutations (between 0.04 and 0.96 allele frequency) in 4 single cells and a

615    bulk sample. **(C)** Allele Frequency heatmap of heteroplasmic (between 0.04 and 0.96 allele

616    frequency) mutations in 4 single cells and a bulk sample. **(D)** snpEff predicted impact of

617    heteroplasmic (between 0.04 and 0.96 allele frequency) mutations in 4 single cells and a bulk sample.

618    **(E)** Number of mutations per gene in 4 single cells and a bulk sample. **(F)** Allele frequency

619    distribution of heteroplasmic (between 0.04 and 0.96 allele frequency) mutations in 4 single cells and

620    a bulk sample. Made with BioRender.