



<b>Title</b>	Ultra-deep next generation mitochondrial genome sequencing reveals widespread heteroplasmy in Chinese hamster ovary cells
<b>Authors(s)</b>	Kelly, Paul S., Clarke, Colin, Costello, Alan, Barron, Niall, et al.
<b>Publication date</b>	2017-05-01
<b>Publication information</b>	Kelly, Paul S., Colin Clarke, Alan Costello, Niall Barron, and et al. "Ultra-Deep next Generation Mitochondrial Genome Sequencing Reveals Widespread Heteroplasmy in Chinese Hamster Ovary Cells," May 1, 2017. <a href="https://doi.org/10.1016/j.ymben.2017.02.001">https://doi.org/10.1016/j.ymben.2017.02.001</a> .
<b>Item record/more information</b>	<a href="http://hdl.handle.net/10197/10335">http://hdl.handle.net/10197/10335</a>
<b>Publisher's statement</b>	This is the author's version of a work that was accepted for publication in Metabolic Engineering. Changes resulting from the publishing process, such as peer review, editing, corrections, structural formatting, and other quality control mechanisms may not be reflected in this document. Changes may have been made to this work since it was submitted for publication. A definitive version was subsequently published in Metabolic Engineering (41, (2017)) <a href="https://doi.org/10.1016/j.ymben.2017.02.001">https://doi.org/10.1016/j.ymben.2017.02.001</a>
<b>Publisher's version (DOI)</b>	<a href="https://doi.org/10.1016/j.ymben.2017.02.001">10.1016/j.ymben.2017.02.001</a>

Downloaded 2026-05-01 23:38:19

The UCD community has made this article openly available. Please share how this access benefits you. Your story matters! (@ucd\_oa)



© Some rights reserved. For more information

1 **Ultra-deep next generation mitochondrial genome sequencing reveals**  
2 **widespread heteroplasmy in Chinese hamster ovary cells.**

3 Paul S Kelly<sup>1†</sup>, Colin Clarke<sup>2†</sup>, Alan Costello<sup>1</sup>, Craig Monger<sup>1,2</sup>, Justine Meiller<sup>1</sup>, Heena Dhiman<sup>2</sup>, Nicole  
4 Borth<sup>3</sup>, Michael J Betenbaugh<sup>4</sup>, Martin Clynes<sup>1</sup> and Niall Barron<sup>1\*</sup>

5 <sup>1</sup> National Institute for Cellular Biotechnology, Dublin City University, Glasnevin, Dublin9, Ireland.

6 <sup>2</sup> National Institute for Bioprocessing Research and Training, Fosters Avenue, Blackrock, Co. Dublin,  
7 Ireland.

8 <sup>3</sup> University of Natural Resources and Life Sciences Vienna, Muthgasse 18, Vienna 1190, Austria.

9 <sup>4</sup> Department of Chemical and Biomolecular Engineering, Johns Hopkins University, Baltimore,  
10 Maryland, USA.

11 † Equal contribution

12 \*Corresponding author

13 **Email:** niall.barron@dcu.ie

14 **Phone:** +353 1 700 5700

15 **Fax:** +353 1 700 5484

16

17 **Keyword:** Chinese hamster ovary; Biopharmaceutical; Mitochondrial; Heteroplasmy; Next generation  
18 sequencing; Genomics;

19

20 **Abbreviations:** Chinese hamster ovary (CHO); Metabolic flux analysis (MFA); Oxidative  
21 phosphorylation (OXPHOS); nuclear DNA (nDNA); mitochondrial DNA (mtDNA); displacement loop (D-  
22 loop); reactive oxygen species (ROS); minor allele frequency (MAF); nuclear encoded mitochondrial  
23 sequences (NumtS); single nucleotide polymorphism (SNP); insertion/deletion (INDEL);

1 **Abstract**

2 Recent sequencing of the Chinese hamster ovary (CHO) cell and Chinese hamster genomes has  
3 dramatically advanced our ability to understand the biology of these mammalian cell factories. In this  
4 study, we focus on the powerhouse of the CHO cell, the mitochondrion. Utilizing a high-resolution next  
5 generation sequencing approach we sequenced the Chinese hamster mitochondrial genome for the  
6 first time and surveyed the mutational landscape of CHO cell mitochondrial DNA (mtDNA). Depths of  
7 coverage ranging from ~3,319X to 8,056X enabled accurate identification of low frequency mutations  
8 (>1%) revealing that mtDNA heteroplasmy is widespread in CHO cells. A total of 197 variants at 130  
9 individual nucleotide positions were identified across a panel of 22 cell lines with 81% of variants  
10 occurring at an allele frequency of between 1 and 99%. 89% of the heteroplasmic mutations identified  
11 were cell line specific with the majority of shared heteroplasmic SNPs and INDELs detected in clones  
12 from 2 cell line development projects originating from the same host cell line. The frequency of common  
13 predicted loss of function mutations varied significantly amongst the clones indicating that  
14 heteroplasmic mtDNA variation could lead to a continuous range of phenotypes and play a role in cell  
15 to cell, production run to production run and indeed clone to clone variation in CHO cell metabolism.  
16 Experiments that integrate mtDNA sequencing with metabolic flux analysis and metabolomics have the  
17 potential to improve cell line selection and enhance CHO cell metabolic phenotypes for  
18 biopharmaceutical manufacturing through rational mitochondrial genome engineering.

## 1 **1 Introduction:**

2 The continual improvement of bioprocesses over the last 20 years has enabled the production of g/L  
3 quantities of complex therapeutic proteins (e.g. monoclonal antibodies) from industrial scale Chinese  
4 hamster ovary (CHO) cell culture [1,2]. These dramatic improvements in performance have been  
5 achieved, in part, through understanding the nutrient requirements of CHO cells to optimise media  
6 formulations. Industrial scale cell culture processes, where possible, now utilise chemically defined  
7 media that maintain growth rate and increase titre as well as eliminate the batch variation associated  
8 with biological components such as serum [3]. The development of fed-batch cell culture strategies  
9 have also been central to achieving high product titres, counteracting the production of cellular waste  
10 products during cell culture and extending production runtimes. During exponential growth, CHO cells  
11 channel glucose and glutamine through the glycolytic pathway even in cases of high oxygen availability  
12 (aerobic glycolysis), a metabolic phenotype similar to the Warburg effect observed in cancer cells [4].  
13 The resulting secretion of lactate and ammonium inhibit cell growth and productivity as well as initiate  
14 apoptosis and decrease product quality [3]. In a fed-batch process, cells are initially grown to a high cell  
15 density before the bioreactor environment is altered through e.g. reducing the temperature or altering  
16 the pH of cell culture. A stationary phase of cell growth is induced to shift CHO cells from a lactate  
17 production to consumption phenotype extending, viability and maximising protein production [3].  
18 Metabolic flux analysis (MFA) and metabolomics approaches have proven to be powerful tools for  
19 understanding the molecular basis of CHO cell metabolic phenotypes to enable predictable, rapid and  
20 inexpensive optimisation of industrial bioprocesses [5]. These techniques have enabled the  
21 development of complex closed loop feeding strategies to limit glucose and glutamine concentrations  
22 as well as identify genetic engineering targets to drive CHO cells toward metabolically desirable  
23 phenotypes [6].

24 Recent studies of CHO cell metabolism have indicated that mitochondrial function is central to lactate  
25 production/consumption [7] and indeed the variability observed in CHO cell metabolic phenotypes [8].  
26 Mitochondria play a central role in eukaryotic cellular energy metabolism via oxidative phosphorylation  
27 (OXPHOS) and have important functions in biological processes such as intracellular calcium signalling  
28 [9] and apoptosis [10]. While the overwhelming majority of proteins required to carry out these functions  
29 are transcribed from nuclear DNA (nDNA) [11], mitochondria also contain a distinct, double stranded  
30 circular genome. Eukaryotic cells can contain more than 1,000 copies of mitochondrial DNA (mtDNA)  
31 packaged within DNA-protein structures known as mitochondrial nucleoids (each nucleoid contains 2-  
32 10 mtDNA molecules). mtDNA copy number varies according to cell type, for instance with myocardial  
33 muscle cells containing on average 6,000 copies while leukocytes may have as few as 350 per cell  
34 [12,13]. In humans, a significant degree of variation in mtDNA copy number has been observed between  
35 the same tissues of different individuals as well as across multiple tissues from the same individual [14].  
36 The mitochondrial genome is between 15 to 17 kb in length and contains 37 genes (28 on the guanine  
37 rich "heavy" or H-strand and 9 on the cytosine rich "light" or L-strand). mtDNA encodes 13 polypeptide  
38 subunits of OXPHOS complexes I, III, IV and V along with 2 ribosomal RNA subunits and 22 tRNAs  
39 required for intra-mitochondrial protein synthesis. The mtDNA genome is extremely compact, genes

1 lack introns, intergenic regions are limited to 1 or 2 nucleotides and in some cases genes can overlap  
2 (e.g. *ATP6* and *ATP8*). The only significant non-coding regulatory region is called the displacement loop  
3 (D-loop) and contains the origin of replication for the H-strand. Transcription is initiated from one of two  
4 H-strand promoters or a single promoter on the L-strand resulting in polycistronic RNA and  
5 subsequently processed to produce mRNAs, tRNAs and rRNAs. Mitochondria utilise a distinct genetic  
6 code for mRNA translation allowing translation of all codons using only 22 tRNAs [15,16]. mtDNA  
7 encodes for only two of the four nDNA stop codons (“AGA” and “AGG”), with the “UAA” stop codon  
8 added to transcribed mRNAs via polyadenylation. The nDNA stop codon “UGA” encodes tryptophan in  
9 mtDNA while the “AUA” codon that encodes for isoleucine in nDNA encodes methionine in mtDNA.

10 Relatively inefficient DNA repair mechanisms and close proximity to reactive oxygen species (ROS)  
11 contribute to a mitochondrial genome mutation rate at least 10 fold higher than that of the nuclear  
12 genome. This high mutation rate has seen widespread application of mtDNA sequencing for studies in  
13 evolutionary biology, population genetics and forensic science. The first pathogenic human  
14 mitochondrial mutations were identified nearly 30 years ago and since then more than 250  
15 polymorphisms, insertions and deletions have been implicated in metabolic disorders as well as cancer  
16 and diabetes. Mitochondrial genome polyploidy can give rise to two cellular states: 1) all mtDNA copies  
17 are identical known as homoplasmy or 2) a mixture of wild-type and mutated mtDNA copies are present,  
18 known as heteroplasmy. In healthy cells, wild-type and mutated mtDNA copies can co-exist;  
19 mitochondrial dysfunction occurs when the ratio reaches a particular level known as the mitochondrial  
20 threshold effect, and in some cases, the frequency of heteroplasmy correlates with the severity of a  
21 clinical phenotype [17]. In recent years, next generation sequencing (NGS) technologies have seen  
22 widespread application for the study of heteroplasmy due to increased specificity, sensitivity and  
23 throughput in comparison to traditional Sanger sequencing. Although initially thought to be a rare  
24 phenomenon, NGS has revealed the prevalence of mitochondrial heteroplasmy in the human  
25 population as well as the age related increase of heteroplasmic variants [18]. Studies utilizing ultra-  
26 deep sequencing to identify very low frequency variants have indicated that heteroplasmy is universal  
27 with each cell containing a complex mixture of mitochondrial genotypes [17].

28 Efforts to understand the biology of CHO cell factories and improve industrial scale biopharmaceutical  
29 manufacturing have been dramatically enhanced since the landmark publication of the CHO-K1  
30 genome in 2011 [19]. A wealth of sequence data is now freely available for several CHO cell lines as  
31 well as the Chinese hamster [19–22]. Direct analyses of these data have permitted the first studies of  
32 CHO cell genome instability [22], chromosomal rearrangement [23], and copy number variation [20].  
33 Methods for expression profiling have also seen marked improvement in the CHO cell post-genomic  
34 era and overcome the reliance on homology with model species that limited early studies in the field.  
35 CHO-cell-specific sequence databases have increased the number of identifications from mass  
36 spectrometry based proteomic analysis [24]. The combination of genome sequence and NGS  
37 technology to study RNA (termed RNA-Seq) has been employed to study mRNA and small RNA [25–  
38 27] expression patterns as well as to annotate transcripts [28] and identify promoter regions [29],  
39 providing novel insights in to the CHO cell transcriptome.

1 While the availability of nuclear genome sequences has undoubtedly advanced CHO cell biology, we  
2 know little about the mtDNA and the impact of mutations on cell metabolism and bioprocess  
3 performance. Here, we present the first comprehensive survey of the CHO cell mitochondrial genome,  
4 spanning a panel of cell lines originating from industry, the ATCC and our laboratory. The utilisation of  
5 next generation sequencing technology enabled high-resolution detection of mtDNA mutations  
6 including those occurring at low frequency. Our results indicate that heteroplasmy is widespread in CHO  
7 cell lines, tends to be cell line specific and that these mutations could play a role in metabolic phenotype  
8 variability.

## 9 **2 Materials and Methods**

### 10 **2.1 Extraction of DNA from Chinese hamster and mouse liver tissue**

11 Genomic DNA was extracted from 30 mg of liver tissue from an outbred Chinese hamster and a  
12 CB17/lcr-Prkdcscid/Crl mouse liver samples (Supplementary Table 1) using a DNeasy Blood and  
13 Tissue kit (QIAGEN, 69581). Tissue samples were sheared using a Dounce homogenizer in 180  $\mu$ l of  
14 ALT buffer. The purity and integrity of extracted tissue-derived genomic DNA was determined on a  
15 nano-drop and via a DNA-agarose gel stained with ethidium bromide (Supplementary Figure 1A).

### 16 **2.2 CHO cell culture and mtDNA extraction**

17 22 CHO cell lines (Supplementary Table 1) were grown in suspension unless indicated otherwise and  
18 harvested at 72 hours. All suspension cultured cell lines were seeded initially at  $2 \times 10^5$  cells/mL in 5  
19 mL of culture media. The 4 Biogen cell lines were cultured in suspension in proprietary chemically  
20 defined media supplied by the industry partner and cultured in-house. The 8 clones from Pfizer  
21 originating from 2 cell line development projects (CLD1 and CLD2), were grown in attached culture in  
22 DMEM supplemented with 5% serum. All remaining CHO cell lines were cultured in suspension in 5 mL  
23 of serum-free media at an initial density of  $2 \times 10^5$  cells/mL.  $17 \times 10^6$  cells were acquired at 72 hours  
24 for mitochondrial DNA extraction [30]. To reduce the contaminating nuclear DNA an additional step was  
25 carried out for CHO cell lines to enrich for double stranded mtDNA using a bacterial mini-prep kit  
26 (QIAGEN, 27104), as previously described [31]. Purity and integrity of isolated CHO mitochondrial DNA  
27 was determined using a NanoDrop and an agarose DNA gel (Supplementary Figure 1B).

### 28 **2.3 Amplification of mitochondrial DNA**

29 To further eliminate nuclear DNA from the tissue and cell line samples, we amplified mtDNA fragments  
30 using a high fidelity PCR kit (Life Technologies, 11304-011) (Supplementary Figure 1C). CHO mtDNA  
31 primers were designed using the CHO cell mtDNA sequence available on GenBank (NC\_007936.1).  
32 Two overlapping ~8.5 kb mtDNA fragments were designed to span the ~16.5 kb mitochondrial genome  
33 sequence (Supplementary Table 2). Another set of overlapping primers was designed based on the  
34 *Mus musculus* mitochondrial genome sequence (NC\_005089.1). PCR amplification of each paired  
35 mitochondrial genome fragment was performed using the following thermo-cycler conditions: 94°C for  
36 2 min, 12 cycles at 94°C for 30 sec, 55°C for 30 sec and 68°C for 8.5 min. The resulting PCR products  
37 for each respective sample was cleaned using Ampure® XP DNA-binding magnetic beads (Agencourt,  
38 A63880). Two PCR primers were designed to amplify a short 454 bp amplicon of the gene CYTB

1 flanking two identified mutated sites (m14136 and m14378) which was used for Sanger sequencing  
2 (For: 5'-TTCAAAGATGTAGCCATACAACC-3' and Rev: 5- AACCGTAATAAACTCCTCGTCC-3').

### 3 **2.4 Nextera XT Mitochondrial DNA library preparation and sequencing**

4 For each mtDNA sample, a serial dilution was performed in nuclease-free water and quantified using  
5 the Qubit dsDNA HS assay kit (Invitrogen, Q32851) to obtain a 0.2 ng/μl stock. 1 ng of mitochondrial  
6 DNA library was prepared using the Nextera XT DNA Sample Preparation Kit (FC-131-1024) in  
7 accordance with the manufacturer's specifications. After each library was fragmented, adapters were  
8 added followed by the incorporation of sample indexes by PCR. Each of the 24 uniquely indexed  
9 samples was passed through a PCR clean-up using Ampure® XP DNA-binding magnetic beads  
10 (Agencourt, A63880). Libraries were quantified using the Qubit dsDNA HS assay kit and fragment size  
11 distribution was determined using the High Sensitivity DNA Bioanalyzer kit (Agilent, 5067-4626) to  
12 confirm the recommended 500-600bp range. Libraries were normalised to 4 nM in resuspension buffer  
13 prior to sequencing. Each 4 nM library was then pooled into a single sample and sequenced on an  
14 Illumina MiSeq (San Diego, CA) configured to produce 151bp paired end reads. Following sequencing,  
15 base calls were converted to 24 individual FASTQ format files for bioinformatics analysis (all raw data  
16 will be uploaded to NCBI's SRA database upon acceptance of publication).

### 17 **2.5 Reconstruction and annotation of the *Cricetulus griseus* mitochondrial genome sequence.**

18 To reconstruct the Chinese hamster mtDNA sequence the MITOBIM algorithm [32] was utilised in  
19 combination with the CHO cell mtDNA sequence available on GenBank (NC\_007936). Paired-end  
20 reads were merged using FLASH [33] prior to assembly. MITOBIM assembles a mitochondrial genome  
21 by mapping sequencing reads (from Chinese hamster liver tissue) to a closely related sequence, in this  
22 case CHO cell mtDNA. The newly assembled Chinese hamster mtDNA was initially annotated using  
23 the MITOS [34] and ARWEN [35] webservers. Annotations were verified and if necessary, refined via  
24 BLAST analysis and comparison to human, mouse and rat mtDNA.

### 25 **2.6 CHO cell line mapping and data pre-processing**

26 Reads corresponding to each of the 22 CHO cell lines sequenced were subjected to quality control  
27 assessment followed by the removal of adapter sequences and reads < 50bp using *trimmomatic* [36].  
28 The remaining sequence data was mapped to the Chinese hamster mitochondrial genome reference  
29 using the BWA-MEM algorithm [37]. Representation of a circular mitochondrial genome as a linear  
30 sequence (i.e. beginning at position 1 and ending at position 16,283) can give rise to incomplete read  
31 mapping due to the introduction of an artificial sequence break. Paired end reads that span the  
32 sequence break will be designated as unmapped and eliminated by algorithms such as BWA  
33 decreasing depth at the start and end of the mtDNA reference and affecting the ability to accurately  
34 detect variants in these regions. In this study the "double" alignment mapping strategy described by  
35 Ding *et. al.* [38] was utilised to ensure optimum alignment of reads to the Chinese hamster mitochondrial  
36 genome. Using this approach, CHO cell line reads were mapped to the original or "unshifted" reference  
37 sequence beginning at position 1 and ending at position 16,283. For the second mapping run a new  
38 reference mtDNA sequence was created by joining the start and ends of the original sequence and

- 1 introducing a new break point so that the sequence began at position 8,000 and ended at position 7999
- 2 on the original reference (Supplementary Figure 2).

## 1 **2.7 Variant pre-processing**

2 Following alignment against the “unshifted” and “shifted” reference sequence reads, with a MAPQ < 20  
3 were designated as “unmapped” and discarded. To pre-process the “unshifted” and “shifted” mapped  
4 data for variant calling, we followed the Genome Analysis ToolKit (GATK) [39] best practice guidelines.  
5 PCR duplicates can arise during the library preparation following the amplification of the multiple copies  
6 of identical DNA fragments. Duplicates propagate errors in sample and library preparation across the  
7 dataset, violate the assumption of independence during variant calling and potentially result in the  
8 identification of false positive variants. Reads corresponding to PCR duplicates were identified using  
9 *Picard* (<http://broadinstitute.github.io/picard/>) and eliminated from further analyses. The remaining  
10 reads were realigned around INDELS accounting for mapping artefacts that arise from the independent  
11 read by read alignment process. Base quality score recalibration reduces the effect of systematic  
12 sequencing biases by first determining covariation between factors including nucleotide context (e.g.  
13 AC dinucleotides are often lower quality than TG) and base position within the read (bases at the ends  
14 of the reads generally have more mismatches). The Phred scaled Q values are then adjusted  
15 accordingly to reduce false positives during variant calling.

## 16 **2.8 Variant discovery and annotation**

17 To identify CHO cell line mtDNA mutations using the dual mapping strategy, variant calls were made  
18 within specific regions of the “unshifted” and “shifted” Chinese hamster reference sequences  
19 (Supplementary Figure 2). For the original “unshifted” reference, variants are called between positions  
20 4,000 and 12,000, while sequence variants on the “shifted” reference sequence are called between  
21 positions 4,000 to 12,283, translating to the regions spanning 1 to 3,999 and 12,001 to 16,283 on the  
22 original reference sequence (encompassing the joined start and ends). The bioinformatics pipeline  
23 incorporated both VarScan [40] and LoFreq [41] for variant detection. Only those SNPs and INDELS  
24 identified by both algorithms with a minor allele frequency (MAF)  $\geq 1\%$ , minimum sequencing depth >  
25 1,500X at the variant position and an average Phred-scaled base quality ( $\geq Q25$ ) for the alternate allele  
26 were reported. Variants were eliminated if overrepresentation of reads supporting the variant was  
27 observed in either forward or reverse direction (i.e. strand bias) [42]. An additional threshold was  
28 employed for VarScan calls and only variants with  $p < 0.01$  were retained. Each identified INDEL was  
29 inspected manually to confirm potential false positives. Upon completion of the mutation detection  
30 pipeline, the coordinates of “shifted” sequence variants were transformed back to the original reference  
31 coordinates and combined with those variants identified following “unshifted” sequence analysis. To  
32 determine the putative effect of each mutation we first utilised snpEff [43] to annotate each mutation  
33 (i.e. frameshift, stop codon gained, start codon mutation, missense or non-synonymous). The  
34 PROVEAN algorithm [44] was utilised to predict the functional impact of missense variants on protein  
35 function. Those missense variants with a Provean score  $\leq -2.5$  were classified as deleterious.

## 36 **2.9 Estimation of contamination from nuclear mitochondrial sequences**

37 Heteroplasmy detection can be confounded, particularly from whole genome sequencing data, by the  
38 presence of nuclear encoded copies of mitochondrial sequences (NumtS) [45]. In this study, potential

1 contamination was reduced through the long range PCR to enrich for mitochondrial DNA and in the  
2 case of cell lines, the utilisation of a bacterial mini-prep kit to eliminate non-circular DNA prior to  
3 amplification. To confirm the effectiveness of the enrichment strategy for mitochondrial DNA, we utilised  
4 the mouse mitochondrial sequencing data to assess potential contamination from NumtS. The variant  
5 identification pipeline was first used to identify SNPs and INDELS against the mouse reference  
6 sequence (NC\_005089.1). To estimate the influence of NumtS, processed reads from mouse were  
7 separately aligned to the *Mus musculus* nuclear genome (mm9 assembly) using the BWA-MEM  
8 algorithm (the “-L” and “-T” parameters were set to “9,9” and “145” respectively). Those reads which  
9 mapped to known mm9 NumtS regions [46] with a MAPQ > 20 were extracted, reconverted to FASTQ  
10 files and remapped against the mouse mtDNA reference sequence using the variant discovery pipeline  
11 described above. The thresholds of the LoFreq and VarScan algorithms were modified to account for  
12 the lower depth of coverage in order to determine if NumtS were influencing variant class and  
13 heteroplasmy levels.

### 14 **3 Results**

#### 15 **3.1 Reconstruction of the *Cricetulus griseus* mitochondrial DNA sequence**

16 Examination of the NC\_007936 mtDNA sequence and the publication describing its acquisition [47]  
17 revealed that the sequence does not originate from *Cricetulus griseus* but the mitochondrial genome of  
18 either a CHO-K1 or CHO A<sub>L</sub> cell line. There are also a number of discrepancies between the annotation  
19 described in the publication and the GenBank entry. For example, the GenBank annotation states that  
20 tRNA<sup>Asn</sup> is encoded on mtDNA H-strand while the original publication places tRNA<sup>Asn</sup> on the L-strand.  
21 In this study, we sequenced the Chinese hamster mitochondrial genome to produce an accurate  
22 reference sequence for comparability of CHO cell line mtDNA as well as resolving any ambiguities in  
23 annotation (Figure 1). Chinese hamster mtDNA was isolated from liver tissue and sequenced on the  
24 Illumina MiSeq platform yielding 640,142 paired-end reads. We reconstructed the mitochondrial  
25 genome from these data using the MITOBIM [32] algorithm with the NC\_007936 mtDNA sequence used  
26 as the “backbone” sequence. The assembled *C.griseus* mitochondrial genome is 16,283bp in length  
27 (A=33.7%, C=22.8%, G=13.0%, T=30.5%) with an overall GC content of 35.7%.

28 Annotation was initially performed using MITOS [34] and ARWEN [35] and further refined through  
29 comparison of the Chinese hamster mtDNA with human, mouse and rat reference mtDNA sequences  
30 (Supplementary Table 3). The Chinese hamster mtDNA has conserved synteny with mammalian  
31 mitochondrial genomes with 13 protein-coding genes, 22 tRNAs and 2 ribosomal RNAs as well as a  
32 non-coding control region (D-loop). Nine genes are encoded on the mtDNA light strand (*ND6*, tRNA<sup>Gln</sup>,  
33 tRNA<sup>Ala</sup>, tRNA<sup>Asn</sup>, tRNA<sup>Cys</sup>, tRNA<sup>Tyr</sup>, tRNA<sup>Ser</sup>, tRNA<sup>Glu</sup> and tRNA<sup>Pro</sup>) with the remaining 28 genes  
34 encoded by the H-strand. 9 protein-coding genes start with ATG initiation codon (*COX1*, *COX2*, *ATP8*,  
35 *ATP6*, *COX3*, *ND4L*, *ND4*, *ND6*, *CYT6*), 2 with an ATT codon (*ND2* and *ND5*), 1 with a GTG (*ND1*) and  
36 1 with an ATA codon (*ND3*). 8 genes terminated with the TAA codon with the *ND1*, *ND2*, *COX3*, *ND6*,  
37 and *ND4* stop codons predicted to be completed via transcript polyadenylation. Comparison of the  
38 Chinese hamster mtDNA sequence to the currently available CHO cell line mitochondrial genome  
39 identified 7 variants, 5 in protein coding sequences with 2 mutations identified in the mtDNA D-loop

1 comprising 4 SNPs, 2 deletions and an insertion (Table 1). The *C.griseus* mitochondrial genome  
2 sequence and corresponding annotation have been submitted to GenBank (accession no: KX576660  
3 [release date 28.09.2016]).

### 4 **3.2 Mapping of CHO cell line next generation sequencing data to the *C.gresius* mitochondrial** 5 **genome**

6 A total of 22 CHO cell lines sourced from industry partners, the ATCC and from our laboratory were  
7 sequenced (Supplementary Table 1). Adapter sequences were trimmed from each of the CHO cell line  
8 datasets followed by removal of low quality reads and those less than 50bp. Upon completion of this  
9 initial pre-processing stage the number of reads remaining in each sample ranged from 1,547,006 (DCU  
10 CHO-K1 SEAP) to 465,194 (Biogen DG44 #1) cell line (Supplementary Table 4A). The BWA-MEM  
11 algorithm was used to align reads against the *C.gresius* mtDNA reference sequence as well as a  
12 modified version of the reference sequence where the original start and ends were joined and a new  
13 breakpoint introduced at 8000bp (Supplementary Figure 2). This dual mapping strategy was utilized to  
14 account for the circularity of the mitochondrial genome and remove bias arising from the use of a linear  
15 reference sequence (e.g. discarding reads that spanned the artificial breakpoint) and therefore improve  
16 our ability to detect mutations around the start and ends of the mtDNA sequence.

17 Following alignment to both the shifted and unshifted reference sequence, reads with a MAPQ < 20  
18 were eliminated from further analysis. The DCU CHO-K1 SEAP sample was found to have the largest  
19 number of reads mapping with MAPQ < 20, yet this represented only ~1% of the total reads in that  
20 dataset. On average > 99% of reads mapped to the unshifted and shifted reference sequence with a  
21 MAPQ ≥ 20 (Figure 2A; Supplementary Table 4B & 4C) demonstrating the effectiveness of the mtDNA  
22 isolation and amplification method utilized in this study. To ensure PCR duplicates were ignored in  
23 downstream stages of the bioinformatics analysis duplicates were “marked” using the *Picard* tool  
24 (<http://broadinstitute.github.io/picard/>). The DCU CHO-K1 SEAP sample had the highest proportion of  
25 duplicates identified (~44%) in the sample set while the Biogen DG44 #1 cell sample had the lowest  
26 proportion of duplicates (~21%) (Figure 2A). The duplicate marked data was further pre-processed for  
27 variant calling by INDEL realignment and base recalibration in line with the *GATK* best practice  
28 guidelines [48,49].

### 29 **3.3 Identification of CHO cell mitochondrial genome variants**

30 The average depth of coverage and perbase coverage for each cell line sample was calculated using  
31 *samtools* [50]. The lowest average depth of coverage across the unshifted reference sequence was  
32 observed for the Biogen DG44 #1 cell line (3,319X) while the deepest coverage was observed for the  
33 miRNA-NC CHO-K1 #1 cell line (8,056X) (Figure 2B, Supplementary Table 4D & 4E). While a negligible  
34 difference in average coverage (~4X) across the entire shifted and unshifted reference sequences was  
35 observed, the effectiveness of the dual mapping strategy is illustrated by an average increase in  
36 coverage across the first and last 100bp of the mtDNA reference sequence of 1,128X and 1,079X  
37 respectively. The coverage at each individual nucleotide position of the reference sequence was found  
38 to be extremely deep (Figure 2C) permitting high resolution analysis of the CHO cell mitochondrial  
39 genome and confident identification of low frequency heteroplasmic variants across the 22 cell lines.

1 The Lofreq and VarScan algorithms were utilized in parallel to identify mutations in CHO cell lines when  
2 compared with Chinese hamster reference sequence. Only those SNPs and INDELS that (1) were  
3 identified by both algorithms, (2) had a minor allele frequency > 1%, (3) the sequencing depth at the  
4 mutant position was > 1,500X, (4) no strand bias was observed, and (5) the average Phred scaled base  
5 quality score for alternative allele was  $\geq$  Q25 were retained for further analysis. For INDELS, the Q  
6 scores of the 10 flanking bases surrounding the variant position were examined and inspected  
7 manually. This procedure was carried out for both the unshifted and shifted reference sequences, and  
8 mutations within the defined calling regions, combined upon completion to produce the final variant set  
9 (Supplementary Figure 2).

10 In total, 197 mutations (175 SNPs and 21 deletions and an insertion) were identified across the 22 CHO  
11 cell lines (Figure 3 and Supplementary Tables 5-7). The SNPs identified corresponded to 99 nucleotide  
12 transitions (A $\leftrightarrow$ G or C $\leftrightarrow$ T) and 22 nucleotide transversions (A $\leftrightarrow$ C, A $\leftrightarrow$ T, G $\leftrightarrow$ C or G $\leftrightarrow$ T), yielding a  
13 4.5:1 ts/tv ratio, similar to previous estimates for mutations in mammalian mtDNA [51]. 130 (121 SNPs  
14 and 9 deletions and an insertion) individual variant nucleotide positions were detected in one or more  
15 of the CHO cell lines. The largest number of variants was identified in the ATCC DG44 and Biogen  
16 DG44 #1 samples with 30 and 21 mutations identified respectively (Figure 4A & Supplementary Table  
17 7) while the Biogen DG44 #2 cell line contained 9 mutated mtDNA positions. The least mutated in the  
18 cell line panel was the ATCC CHO-S cell line with only a single SNP identified while 6 SNPs were  
19 identified in the mitochondrial genome of the Biogen CHO-S cell line. Cell lines from the CHO-K1  
20 lineages varied from 3 mutations to as many as 12 variants. The number of variant positions in the  
21 CHO-K1 cell lines originating from Pfizer cell lines derived from the same parental host ranged from 7  
22 to 12. The remaining CHO-K1 cell lines all had less than 10 variant positions.

23 Of the 37 genes in the mitochondrial genome, 23 were found to have at least one variant position in  
24 one of the CHO cell lines sequenced (Figure 4B). A SNP and an insertion were identified in the mtDNA  
25 D-loop in the ATCC DG44 and Pfizer CHO-K1 #F5 cell lines. The 16S rRNA gene had the largest  
26 number of mutated positions with variants identified at 17 separate nucleotides. All of the 13 protein  
27 coding genes had at least one mutation while 8 of 22 tRNAs harboured a mutation. The *CYTB* and  
28 *COX1* were found to have the largest number of mutations amongst protein coding genes while the  
29 *ND3* gene had only one variant. A homoplasmic (MAF > 99%) SNP in tRNA<sup>Val</sup> (m.1074C>T) and the  
30 16S rRNA gene (m.2235C>T) was identified in each of the 17 cell lines from the CHO-K1 lineage  
31 (Supplementary Table 7). The CHO cell lines from the S and DG44 lineages are identical to the Chinese  
32 hamster mtDNA reference sequence at these positions. The two Biogen DG44 cell lines had a shared  
33 homoplasmic mutation in tRNA<sup>Val</sup> (m1092A>G) yet this mutation was not detected in the ATCC DG44  
34 cell line. The Biogen CHO S cell line contains the only homoplasmic mutation in a protein-coding gene,  
35 a SNP identified in *CYTB* (m.14311C>T).

### 36 **3.4 Identification of CHO cell line heteroplasmy**

37 81% of all CHO cell line mutations identified in this study were heteroplasmic (i.e. the MAF>1% & <99%)  
38 with a minor allele frequency spanning from 1% to 96.2% (Figure 4C). While the majority of  
39 heteroplasmic variants were identified in a single cell line, 11 of these mutations were shared in two or

1 more cell lines (Table 2). The effectiveness of the dual mapping strategy is further demonstrated by the  
2 identification of two variants identified in the first 100bp of the mitochondrial genome with tRNA<sup>Phe</sup>  
3 (m.62C>T) and the 12S rRNA gene (m.74G>A). These low frequency heteroplasmies were not  
4 identified by the variant calling pipeline using the unshifted linear reference sequence. Next generation  
5 sequencing data from the Chinese hamster liver sample was also analysed using the variant detection  
6 pipeline. No heteroplasmic variants were identified in the Chinese hamster mitochondrial genome  
7 sequence. We did, however, identify a previously reported heteroplasmic insertion [52] and deletion  
8 [53] from the mouse liver mtDNA sequencing data using the mouse mtDNA reference sequence  
9 (GenBank accession: NC\_005089) and the variant identification pipeline (Supplementary Table 8).

10 The mouse mtDNA sequencing data was utilised to demonstrate the effectiveness of the bioinformatics  
11 pipeline and estimate the potential influence of NumtS contamination on variant calling and  
12 heteroplasmy levels. Reads originating from mouse mtDNA sequencing were first analysed using an  
13 identical bioinformatics pipeline to the Chinese hamster and CHO cell line analyses, incorporating the  
14 shifted and unshifted reference sequence mapping. From this analysis, 3 known [52,53] mouse mtDNA  
15 variants were identified - a homoplasmic SNP (m.9461T>C), a heteroplasmic deletion (m.5171delA)  
16 and an insertion (m.9820insAA) (Supplementary Table 8). To determine if these variants and their  
17 corresponding allele frequencies were influenced by NumtS, we stringently mapped reads against the  
18 mouse nuclear genome and extracted reads that aligned to known NumtS [46]. The NumtS aligned  
19 reads were remapped to the mouse mtDNA reference sequence using the same alignment and variant  
20 pre-processing approach to that of Chinese hamster and CHO cell for cell line data with the exception  
21 of variant calling thresholds that were modified to account for the lower depth of coverage. Reads were  
22 found to align predominantly to 5 regions in the mouse mtDNA reference sequence (Supplementary  
23 Figure 3), which did not overlap with the 3 variants called on the full dataset. Furthermore no SNPs or  
24 INDELS were identified in the 5 regions where NumtS reads aligned indicating that NumtS  
25 contamination was not a contributing factor in either variant detection or heteroplasmy measurement.

### 26 **3.5 Prediction of the effect of mitochondrial genome variations**

27 Of the variants identified in this study, 62% (81/130) lie within protein coding regions of the mitochondrial  
28 genome. To determine the putative effects of these variants we utilized the snpEff [43] tool to annotate  
29 each mutation. Missense variants were the most common, accounting for 60.5% (49/81) of mutations  
30 in protein coding regions in comparison to 25% (18/81) of mutations predicted to be synonymous.  
31 PROVEAN predicted that 55% (27/49) of the amino acid substitutions arising of missense mutations  
32 would affect protein function (Supplementary Table 7). The remaining variants were predicted to result  
33 in a frameshift mutation (7.5%) or premature stop codon (5%), with 1 start codon mutation identified. Of  
34 the 22 cell lines sequenced, 20 contained mutations that resulted in alteration of the protein coding  
35 sequence of at least 1 gene in the mitochondrial genome (Figure 5A). Each protein coding gene  
36 harboured at least 1 mutation that altered the amino acid sequence in at least one sample (Figure 5B).  
37 Frameshift mutations with a high probability of functional consequences were identified in *ND1*, *COX1*,  
38 *ND4*, *ND5* and *CYTB* while *COX1*, *ND4L*, *ND6* while mutations resulting in a premature stop codon  
39 were observed in *COX1*, *ND4L*, *ND6* and *CYTB*.

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31

3.6 Is there a relationship between heteroplasmic variants and cell phenotype?

In order to investigate whether the existence of these heteroplasmic variants could infer anything about the phenotypic behaviour of the cell lines, we picked 2 closely related lines for a more detailed analysis of their growth behaviour and oxidative phosphorylation potential. While several of the SNPs and indels were predicted to impact protein function as mentioned above we chose two of the CLD1 with heteroplasmic variants in the CytB gene, CLD1 #3 and CLD1 #5. Fig X shows that Cell line #3 only reached  $1.4 \times 10^6$  cells/ml in culture over 5 days whereas cell line #5 reached  $1.8 \times 10^6$  cells/ml. When stained using Mitotracker, which gives an indication of mitochondrial number, it was apparent that there was not a significant difference between the 2 lines. Interestingly the staining showed a strong anti-correlation with cell growth – with functional mitochondrial content increasing considerably once the cells entered stationary phase. This is in keeping with the theory that rapidly growing cells generate ATP by aerobic glycolysis but switch to oxidative phosphorylation during stationary phase – also referred to as the metabolic shift. This shift in energy generation is typically associated with greater productivity (Ref?). To further investigate whether there were other metabolic differences between the two lines we measured oxygen utilisation by respirometry (Oroboros). This was performed during the mid-exponential phase of growth and showed that the cells that went on to reach higher cell density had reduced oxygen consumption and lower oxphos potential than the cells that peaked at  $1.4 \times 10^6$  cells/ml. Fig Y shows that CLD1 #3 had greater leak, routine and maximum O<sub>2</sub> consumption, indicating a higher oxphos potential. This cell line (#5) had a 36% heteroplasmic frameshift mutation in CytB at nt14136 but did not have the stop codon-forming G>A SNP at nt14378. On the other hand cell line #3 had only 7% heteroplasmic variant at nt14236 but 35% heteroplasmy at nt14378. Sanger sequencing suggested that the 2 variants are unlikely to co-exist on the same individual mitochondrial genome (Suppl Fig – Pauls data) which by extension means that there would be a higher overall percentage of mutant CytB in the #3 cell line (43%). It is also interesting to note that this cell line showed an 18% heteroplasmic variant in Cox2 which is predicted to cause an initiator codon change from ATG to ATA and would be expected to reduce translation of that protein. While far from definitive, these observations could contribute to the reduced Oxphos capacity of these cells. Finally, measuring CytB and Cox2 in both lines showed that both were expressed at a higher level in the cells with enhanced oxphos.

**4 Discussion**

Analytical methods including metabolic flux analysis and metabolomics have shed light on CHO cell metabolic phenotypes, informing media design, feeding strategies and cell line engineering to increase the efficiency of industrial cell culture for biopharmaceutical production [5]. Despite the clear utility of these methods, our understanding of the origins of variation in cellular metabolism in different CHO cell lines, clones and in some cases, between and over the course of production runs have not yet been completely unravelled [6]. Comparison of nuclear genome sequences has shown that CHO cell lineages

1 harbour distinct mutations, and that millions of SNPs and INDELS arise during the development of a  
2 new cell line [20]. This plasticity of the CHO cell nuclear genome sequence undoubtedly plays an  
3 underlying role in the range of bioprocess phenotype variation observed amongst CHO cell clones  
4 during cell line development. Eukaryotic cells also contain a separate non-nuclear polyploid genome  
5 within each mitochondrion. mtDNA sequence variants in the mitochondrial genome have been widely  
6 studied in human biology, and the association of specific mutations with a number of metabolic  
7 disorders is well established. In recent years, the development of massively parallel sequencing  
8 technology has dramatically expanded our ability to study mitochondrial genomics and permitted  
9 analysis of low frequency mtDNA heteroplasmy. Here, we present the first survey of mtDNA  
10 heterogeneity across CHO cell lineages and cell lines as well as amongst clones produced during cell  
11 line development.

12 Considering the genomic instability of CHO cells, Chinese hamster mtDNA is an ideal common  
13 mitochondrial reference genome to compare sequence variants across cell lines. Assessment of the  
14 suitability of an existing Chinese hamster sequence available on GenBank (NC\_007936) for this  
15 purpose revealed that the sequence is from either a CHO<sub>A</sub>L or CHO-K1 CHO [47] cell line, not the  
16 Chinese hamster, and that a number of sequence features are incorrectly annotated. In order to ensure  
17 the accuracy of a reference sequence for comparison of cell lines, we sequenced mtDNA from the  
18 Chinese hamster. Comparative analysis of the Chinese hamster mtDNA sequence with the GenBank  
19 NC\_007936 identified 7 mutations comprising 4 SNPs and 3 INDELS within 4 protein-coding genes and  
20 the mitochondrial D-loop control region. The 16,283bp sequence acquired in this study represents the  
21 first accurate reference for the analysis of the mitochondrial genome, providing an essential resource  
22 for future studies of CHO cell mtDNA.

23 To determine the prevalence of mutations in the CHO cell mitochondrial genome, 22 CHO cell lines  
24 derived from the CHO-K1, CHO-S and DG44 lineages, including industrial cell lines, engineered cell  
25 lines and clones generated from 2 cell line development projects were selected for mtDNA sequencing.  
26 The high depth of coverage achieved through massively parallel sequencing of the relatively small  
27 mitochondrial genome permits accurate identification of homoplasmic and heteroplasmic variants with  
28 minor allele frequencies as low as 1%. Established best practices were utilised to pre-process  
29 alignments before a conservative variant calling pipeline that required agreement between two  
30 algorithms for SNP and INDEL identification.

31 The heterogeneity of the CHO cell mitochondrial genome across the CHO cell lines analysed here is  
32 remarkable; cell lines were found to contain at least one to as many as 30 mutations in their mtDNA.  
33 We discovered both homoplasmic and heteroplasmic mutations at 130 nucleotide positions distributed  
34 across the entire mtDNA sequence with a total of 197 variants detected across the 22 CHO cell lines.  
35 A SNP or INDEL was identified in all 13 protein-coding genes, 8 tRNA genes, and both rRNA genes as  
36 well as with the D-loop region. Each of the protein coding genes were found to contain at least one  
37 heteroplasmic mutation in at least one cell line with *CYTB*, *COX1* and *ND5* each harbouring at least 10  
38 distinct heteroplasmic variants. Of the 130 variant positions identified, only 4 (37 mutations across the  
39 22 cell lines) were homoplasmic (MAF>99%). All 17 cell lines from the CHO-K1 lineage harboured a

1 homoplasmic SNP within both tRNA<sup>Val</sup> (m.1074C>T) and 16S rRNA (m.2235C>T) genes. We also  
2 identified a homoplasmic variant within tRNA<sup>Val</sup> (m.1092A>G) in the two CHO-DG44 cell lines provided  
3 by Biogen, yet this mutation was not identified in the CHO-DG44 cell line sourced from the ATCC. The  
4 fourth homoplasmic variant was identified in the *CYTB* gene (m.14311C>T) in the Biogen CHO-S cell  
5 line. The CHO-K1 and CHO-S were both derived from the original Chinese hamster ovary tissue isolate  
6 before being sent to two different laboratories [54]. We did not detect the tRNA<sup>Val</sup> (m.1074C>T) or 16S  
7 rRNA (m.2235C>T) common to CHO-K1 cell lines in the CHO-S cell lines or the cell line sequenced in  
8 the Partridge *et al.* study (CHO-K1 cell or derived from CHO<sub>AL</sub>). In addition, the homoplasmic *CYTB*  
9 mutation (m.14311C>T) observed in the Biogen CHO-S cell line was not identified in ATCC CHO-K1 or  
10 the ATCC CHO-S cell line. These findings indicate that CHO cell mtDNA mutations have arisen  
11 independently between cell lines following isolation from the original Chinese hamster tissue and over  
12 time, have become fixed in the mitochondrial genome. The tRNA<sup>Val</sup> (m.1902A>G) was identified in both  
13 CHO-DG44 cell lines provided by Biogen, and once again, this mutation was not present in the ATCC  
14 CHO-DG44. These mutations in CHO-DG44, which originated from a non-CHO-K1 or CHO-S lineage  
15 [54], imply that homoplasmic mutations in CHO cell mtDNA can occur spontaneously and differ between  
16 CHO cell lines of the same lineage.

17 The overwhelming majority of mutations identified in this study were heteroplasmic and were detected  
18 at 126 individual nucleotide positions. Heteroplasmic mutations were identified in tRNA, rRNA and  
19 protein coding regions as well as the D-loop control region. Considering the widespread heteroplasmy  
20 identified in the mitochondrial genome of CHO cells, it was somewhat surprising that no heteroplasmy  
21 was detected in the Chinese hamster mitochondrial genome. A recent study by Li *et al.* [18] reported  
22 tissue specific patterns of heteroplasmy in more than 150 humans across 12 tissue types and  
23 demonstrated that heteroplasmy is tissue-specific and, in some cases, tissues can be free of  
24 heteroplasmy while other tissues in the same individual can contain heteroplasmic mtDNA variants. It  
25 is also possible that very low frequency (MAF<1%) variations are present in the Chinese hamster liver  
26 tissue sample. While Li *et al.* utilised a > 0.5% MAF threshold, we felt MAF>1% threshold was an  
27 appropriate choice for this study, given the depth of mtDNA coverage obtained.

28 In comparison to homoplasmic mutations which tended to be lineage specific, 89% of heteroplasmic  
29 mutations were identified in a single cell line and ranged from 1-96%. For example, the 2 Biogen DG44  
30 CHO cell lines shared a common homoplasmic SNP in the tRNA<sup>Val</sup> gene (m.1092A>G) yet no shared  
31 heteroplasmic variants were common to both cell lines. There were also no shared heteroplasmic  
32 mutations identified between several of the CHO-K1 cell lines (i.e. Biogen CHO-K1, Pfizer CHO-K1  
33 2B6, Pfizer CHO-K1 114 cell lines, DCU CHO-K1 SEAP or the 4 DCU sponge transfected SEAP cell  
34 lines). Of those heteroplasmic mtDNA mutations that were identified in more than one cell line, the  
35 majority were shared amongst clones originating from 2 distinct cell line development projects (CLD1  
36 and CLD2).

37 PROVEAN predictions found that of the 49 amino acid substitutions arising from a missense  
38 mutation, 55% of these are highly likely to result in a functional effect. When comparing the  
39 three missense mutations in the ND5 coding sequence (m.11898G>A-18%, m.12078T>A-70%

1 and m.13065C>G20%), the wild-type Lysine (K) coded at position m.12078 is conserved  
2 between the Chinese hamster and Human whereas the other two amino acids are not. The 70%  
3 heteroplasmic shift from T>A thereby changing the amino acid coding sequence is likely to  
4 elicit some form of dysfunction given the level of conservation of the amino acid as well as the  
5 high frequency of heteroplasmy. This particular amino acid location has not been reported  
6 previously as a heteroplasmic variant in other cellular models, however a variety of mutations  
7 within ND5 have been reported such as m.13565C>T in human mitochondrial cybrids induced  
8 the decrease in Ca<sup>2+</sup> uptake to the mitochondria as well as a dependence on glycolysis for ATP  
9 production []. One conserved amino acid location between both human and CHO is Alanine  
10 (A11) which was identified to possess a 4.1% (m.9887G>A) missense heteroplasmic mutation  
11 in the publically available CHO-DG44 line. This very same mutation has been detected  
12 previously in the mitochondrial DNA of esophageal cancer [4]. Although a considerable  
13 number of amino acids found to be mutated in this study are conserved from CHO to Human  
14 suggesting a functional role, the specific amino acid in question has not yet been identified in  
15 other studies.

16 Seven heteroplasmic variants were identified in 2 or more of the CLD1 clones while 4 heteroplasmies  
17 were found in 2 or more of the 3 CLD2 clones. Four shared heteroplasmic variants were identified in at  
18 least one clone from CLD1 and CLD 2. A SNP in tRNA<sup>Lys</sup> (m.7721A>G) was identified in 1 clone from  
19 CLD1 (MAF=47%) and 2 clones from CLD2 (MAF=13% and MAF =9.5%). A SNP in *CYTB*  
20 (m.14849G>A) was found in a single CLD1 clone and in 3 clones from CLD2. Three of the clones had  
21 an average mutation frequency of ~45% yet the MAF of the fourth clone was 1.5%, a marked difference  
22 from its counterpart clones within CLD2. The only mutation present in all 8 clones from CLD1 and CLD2  
23 is a frameshift mutation in the *CYTB* gene (m.14136delA) with a MAF ranging from 7.7% to 52%. The  
24 m.14136delA mutation also appears in the ATCC CHO-K1 cell line (derived from an isolate of the  
25 original Chinese hamster ovary tissue) and it would seem that this mutation has been retained in the  
26 CLD project clones yet has been lost in the other CHO-K1 cell lines while the homoplasmic m.1074C>T  
27 and m.2235C>T have become fixed in all CHO-K1. The variation of mutation frequencies of  
28 heteroplasmic variants both within and across the two cell line development projects, following post  
29 single cell cloning, is in line with the model of random assortment of mtDNA upon cell division [55].

30 It is not possible at this point to determine the impact of rRNA or tRNA mutations without further  
31 experimentation. It would be expected, however, that perturbations within the mitochondria's  
32 translational machinery would have a considerable impact on the inner mechanics of an energy-  
33 producing mitochondrion. Lie *et al.* [\*] reported that a homoplasmic T10003C mutation in the tRNA<sup>Gly</sup>  
34 gene caused a 70% reduction in the steady state level of tRNA<sup>Gly</sup> with an associated 33% reduction in  
35 mitochondrial translation. In this instance, it was predicted that this mutation interfered with the  
36 formation of the tRNA secondary structure by forming a base pair at 13C-22G in the conserved stem.  
37 In our study, 8 tRNA genes were discovered to have heteroplasmic variations present to as far as 74%  
38 for tRNA<sup>Leu</sup> in the case of the Pfizer CHO-K1 2B6 line (m11699G>A). Given the short sequence length  
39 of tRNAs (76-90 nt) and their functional dependency on secondary structure, it would be expected the  
40 small changes observed in this study would have an effect. Our results do, however, indicate that the  
41 mtDNA mutations are not only widespread but also likely to influence mitochondrial function with each

1 of the protein coding genes in the mitochondrial genome found to contain at least one SNP or INDEL  
2 in one cell line and *CYTB*, *COX1* and *ND5* each harbouring at least 10 distinct mutations. A variety of  
3 effects were predicted for mutations in protein coding genes including synonymous, non-synonymous  
4 and mutations in the initiation codon as well as variants that are likely to have a more pronounced effect  
5 on the protein sequence, including premature stop and frameshift mutations. For instance, *CYTB* was  
6 found to harbour a SNP (m.14849G>A) predicted to result in a premature stop codon as well as a  
7 frameshift mutation (m.14136delA) shared between the CLD1 and CLD2 cell lines. The presence of  
8 mtDNA mutations detected here could also play a role in CHO cell metabolic phenotype variation. This  
9 gain of stop and frameshift mutation in the *CYTB* gene are predicted to result in a loss of function and  
10 could lead to a diminished efficiency in mitochondrial aerobic respiration in the electron transport system  
11 centred on complex III. Weakened oxidative phosphorylation (OXPHOS) could signal reprogramming  
12 of cellular metabolism to rely more heavily on glycolysis [56] and maintain cellular energy balance.  
13 Exclusive reliance on glycolytic metabolism, despite being energetically inefficient, has been shown to  
14 be associated with elevated cell growth in both cancer and CHO cells due to the intermediate  
15 metabolites of glycolysis that contribute to biomass accumulation [57]. All but one protein coding variant  
16 identified were heteroplasmic and in some cases, loss of function mutations spanned a wide range.  
17 Mutation frequencies were found to vary by as much as 1.5% to 50% e.g. a *ND1* frameshift  
18 (m.3205delCT), indicating that the potential effects of these mutations could lead to a continuous  
19 distribution of phenotypes. Extensive metabolic profiling of the panel of 22 CHO cell lines utilised in this  
20 study overlapped with this comprehensive mutational data would begin to answer these questions of  
21 variant translating to phenotype.

22 As mentioned previously, three protein-coding genes across the 22 cell lines harboured at least 10  
23 distinct mutations in the case of *CYTB*, *COX1* and *ND5*. However, in some instances, numerous  
24 mutations within the same cell line of the same gene was detected such as in the case of the *ND5* gene  
25 in the Biogen DG44 #1 cell line. In this case, *ND5* was observed to contain 4 mutations, 3 missense  
26 (m11,898G>A, m12,078T>C and m13,065C>G) occurring at a frequency of 18%, 70% and 20%,  
27 respectively, as well as a synonymous mutation at m13,157A>G (8.9%). These three missense  
28 mutations potentially changing the amino acid sequence could all impact negatively on the functional  
29 role of *ND5*. Taking into account the frequency of each mutation, if each mutant exists in isolation then  
30 a 100% dysfunctional *ND5* protein could prevail. If however, these mutations co-inhabit the same  
31 mtDNA genome, then wild-type protein could remain thereby not breaking the biochemical threshold. It  
32 has previously been shown that heteroplasmic mutations in *ND5* results in the disruption of NADH  
33 dehydrogenase assembly which is associated with a weakened OXPHOS and an increase in lactate  
34 production due to glycolytic dependency []. Cellular mitochondrial compartmentalisation would have to  
35 be determined in order to elucidate the partitioning of this mutational pool and predict whether low-  
36 frequency heteroplasmic variants remain segregated and functional []. To get a better understanding of  
37 this, we amplified the region of the *CYTB* gene that was identified to harbour two heteroplasmic  
38 mutations at m14136 and m.14378 in three clones from the cell line development panel (CLD1 #3,  
39 CLD1 #5 and CLD2 #2). TOPO cloning was performed and 10 clones from each cell line was Sanger  
40 sequenced as a means to determine the co-habitation of these mutations (Supplementary Figure S4).

1 When compared to the CH mitochondrial reference sequence, the early frameshift mutation, shared in  
2 all clones, was detected and reflected the heteroplasmic frequency identified through deep-sequencing.  
3 In the case of CLD1 #3 which contains both mutations, the second STOP mutation was not detected  
4 within the small panel of clones suggesting that these two mutations do not co-inhabit the same mtDNA  
5 copy further suggesting that the dysfunctional *CYTB* protein that results from these two mutants will  
6 have a synergistic effect. With such a small panel of clones, however, the later mutation was not  
7 detected. This highlights the potential impact that several low MAF variations can have on mitochondrial  
8 function if present on individual genomic copies and working in unison.

9 Phenotypic analysis of oxfhos potential, growth characteristics and mitochondrial content of two cell  
10 lines in particular provided some tantalizing evidence of how heteroplasmic variants in important  
11 mitochondrial encoded proteins could impact on the suitability of certain cell lines for recombinant  
12 protein production. However, more targeted mitochondrial genome engineering studies in future should  
13 ascertain whether these variants are directly responsible for the phenotypes observed and indeed  
14 whether they can be manipulated to improve these characteristics. Some the molecular tools required  
15 to achieve this have only recently been developed (cell paper – mito talens).

16 While the results of this study demonstrate the heterogeneity of CHO cell mtDNA, the polyploid nature  
17 of the mitochondrial genome presents considerable challenges to understanding the impact of  
18 mitochondrial mutations. A mtDNA variant might be spread across the entire population or confined to  
19 a subpopulation of cells. At the subcellular level, the mutation might be distributed across multiple  
20 mitochondria or indeed confined to a limited number of mitochondria. The emergence of new methods  
21 to sequence mtDNA at the single cell level [55] will play a valuable role in future studies and increase  
22 our understanding of the implications of particular mutation in a CHO cell population. It will also be  
23 important to integrate MFA and metabolomics analysis with mtDNA sequencing to understand the  
24 biochemical threshold at which individual mutations affect CHO cell behaviour. The knowledge gained  
25 in doing so has the potential to enable precise cell line selection and ultimately rational genetic  
26 engineering to improve CHO cell phenotypes. The recent development of mtDNA CRISPR-Cas9 based  
27 methods [58] and mitoTALENS [59] for site specific mitochondrial genome editing provide routes to  
28 rational engineering of CHO cell mitochondria to enhance the metabolic performance of CHO cells for  
29 biopharmaceutical manufacture.

## 30 **5 Conclusions**

31 Widespread heteroplasmy in the CHO cell mitochondrial genome raises intriguing questions about the  
32 genetics and selection of mitochondrial mutations in CHO cells during cell culture and cell line  
33 development for biopharmaceutical production. Closely related clones derived from the same parental  
34 host and even originating from the same cell line development project can harbour distinct  
35 heteroplasmic variations. These variations in the mitochondrial genome are likely to affect mitochondrial  
36 function and could play a role in cell to cell, production run to production run and indeed clone to clone  
37 variation observed in CHO cell culture and cell line development. The combination of mtDNA  
38 sequencing with established techniques in metabolic flux analysis and metabolomics will be necessary  
39 to associate these mutations with desirable or undesirable CHO cell metabolism. The understanding of

1 mtDNA variation could lead to new approaches to cell line screening and ultimately engineering of CHO  
2 cell mtDNA for more productive metabolic phenotypes.

### 3 **Acknowledgements**

4 The authors gratefully acknowledge funding from Science Foundation Ireland (grant refs:  
5 13/SIRG/2084, 13/IA/1841 and 13/IA/1963) and the eCHO systems Marie Curie ITN programme (grant  
6 ref: 642663). The authors would also to acknowledge Lin Zhang (Pfizer Inc.), Scott Estes, Chapman  
7 Wright and Brian St. Germaine (Biogen Inc.) for access to cell lines and comments on the manuscript.

## 1 References

- 2 [1] Wurm, F.M., Production of recombinant protein therapeutics in cultivated mammalian cells.  
3 *Nat. Biotechnol.* 2004, 22, 1393–1398.
- 4 [2] Li, F., Vijayasankaran, N., Shen, A. (Yijuan), Kiss, R., Amanullah, A., Cell culture processes for  
5 monoclonal antibody production. *mAbs* 2010, 2, 466–477.
- 6 [3] Butler, M., Animal cell cultures: recent achievements and perspectives in the production of  
7 biopharmaceuticals. *Appl. Microbiol. Biotechnol.* 2005, 68, 283–291.
- 8 [4] Vander Heiden, M.G., Cantley, L.C., Thompson, C.B., Understanding the Warburg effect: the  
9 metabolic requirements of cell proliferation. *Science* 2009, 324, 1029–1033.
- 10 [5] Quek, L.-E., Dietmair, S., Krömer, J.O., Nielsen, L.K., Metabolic flux analysis in mammalian cell  
11 culture. *Metab. Eng.* 2010, 12, 161–171.
- 12 [6] Young, J.D., Metabolic flux rewiring in mammalian cell cultures. *Curr. Opin. Biotechnol.* 2013,  
13 24.
- 14 [7] Zagari, F., Jordan, M., Stettler, M., Broly, H., Wurm, F.M., Lactate metabolism shift in CHO cell  
15 culture: the role of mitochondrial oxidative activity. *New Biotechnol.* 2013, 30, 238–245.
- 16 [8] Gilbert, A., McElearney, K., Kshirsagar, R., Sinacore, M.S., Ryll, T., Investigation of metabolic  
17 variability observed in extended fed batch cell culture. *Biotechnol. Prog.* 2013, 29, 1519–1527.
- 18 [9] Rizzuto, R., De Stefani, D., Raffaello, A., Mammucari, C., Mitochondria as sensors and  
19 regulators of calcium signalling. *Nat. Rev. Mol. Cell Biol.* 2012, 13, 566–578.
- 20 [10] Wang, C., Youle, R.J., The Role of Mitochondria in Apoptosis. *Annu. Rev. Genet.* 2009, 43, 95–  
21 118.
- 22 [11] Calvo, S., Jain, M., Xie, X., Sheth, S.A., et al., Systematic identification of human mitochondrial  
23 disease genes through integrative genomics. *Nat. Genet.* 2006, 38, 576–582.
- 24 [12] Miller, F.J., Rosenfeldt, F.L., Zhang, C., Linnane, A.W., Nagley, P., Precise determination of  
25 mitochondrial DNA copy number in human skeletal and cardiac muscle by a PCR-based assay:  
26 lack of change of copy number with age. *Nucleic Acids Res.* 2003, 31, e61.
- 27 [13] Liou, C.-W., Lin, T.-K., Chen, J.-B., Tiao, M.-M., et al., Association between a common  
28 mitochondrial DNA D-loop polycytosine variant and alteration of mitochondrial copy number in  
29 human peripheral blood cells. *J. Med. Genet.* 2010, 47, 723–728.
- 30 [14] Wachsmuth, M., Hübner, A., Li, M., Madea, B., Stoneking, M., Age-Related and Heteroplasmy-  
31 Related Variation in Human mtDNA Copy Number. *PLOS Genet* 2016, 12, e1005939.
- 32 [15] Anderson, S., Bankier, A.T., Barrell, B.G., de Bruijn, M.H.L., et al., Sequence and organization of  
33 the human mitochondrial genome. *Nature* 1981, 290, 457–465.
- 34 [16] Barrell, B.G., Anderson, S., Bankier, A.T., de Bruijn, M.H., et al., Different pattern of codon  
35 recognition by mammalian mitochondrial tRNAs. *Proc. Natl. Acad. Sci. U. S. A.* 1980, 77, 3164–  
36 3166.
- 37 [17] Payne, B.A.I., Wilson, I.J., Yu-Wai-Man, P., Coxhead, J., et al., Universal heteroplasmy of human  
38 mitochondrial DNA. *Hum. Mol. Genet.* 2013, 22, 384–390.
- 39 [18] Li, M., Schröder, R., Ni, S., Madea, B., Stoneking, M., Extensive tissue-related and allele-related  
40 mtDNA heteroplasmy suggests positive selection for somatic mutations. *Proc. Natl. Acad. Sci.*  
41 2015, 112, 2491–2496.
- 42 [19] Xu, X., Nagarajan, H., Lewis, N.E., Pan, S., et al., The genomic sequence of the Chinese hamster  
43 ovary (CHO)-K1 cell line. *Nat Biotech* 2011, 29, 735–741.
- 44 [20] Lewis, N.E., Liu, X., Li, Y., Nagarajan, H., et al., Genomic landscapes of Chinese hamster ovary  
45 cell lines as revealed by the *Cricetulus griseus* draft genome. *Nat Biotech* 2013, 31, 759–765.
- 46 [21] Brinkrolf, K., Rupp, O., Laux, H., Kollin, F., et al., Chinese hamster genome sequenced from  
47 sorted chromosomes. *Nat Biotech* 2013, 31, 694–695.
- 48 [22] Kaas, C.S., Kristensen, C., Betenbaugh, M.J., Andersen, M.R., Sequencing the CHO DCB11  
49 genome reveals regional variations in genomic stability and haploidy. *BMC Genomics* 2015, 16,  
50 160.

- 1 [23] Cao, Y., Kimura, S., Itoi, T., Honda, K., et al., Construction of BAC-based physical map and  
2 analysis of chromosome rearrangement in Chinese hamster ovary cell lines. *Biotechnol. Bioeng.*  
3 2012, 109, 1357–1367.
- 4 [24] Meleady, P., Hoffrogge, R., Henry, M., Rupp, O., et al., Utilization and evaluation of CHO-  
5 specific sequence databases for mass spectrometry based proteomics. *Biotechnol. Bioeng.*  
6 2012, 109, 1386–1394.
- 7 [25] Gerstl, M.P., Hackl, M., Graf, A.B., Borth, N., Grillari, J., Prediction of transcribed PIWI-  
8 interacting RNAs from CHO RNAseq data. *J. Biotechnol.* 2013, 166, 51–57.
- 9 [26] Diendorfer, A.B., Hackl, M., Klanert, G., Jadhav, V., et al., Annotation of additional evolutionary  
10 conserved microRNAs in CHO cells from updated genomic data. *Biotechnol. Bioeng.* 2015, n/a-  
11 n/a.
- 12 [27] Birzele, F., Schaub, J., Rust, W., Clemens, C., et al., Into the unknown: expression profiling  
13 without genome sequence information in CHO by next generation sequencing. *Nucleic Acids*  
14 *Res.* 2010, 38, 3999–4010.
- 15 [28] Rupp, O., Becker, J., Brinkrolf, K., Timmermann, C., et al., Construction of a Public CHO Cell Line  
16 Transcript Database Using Versatile Bioinformatics Analysis Pipelines. *PLoS ONE* 2014, 9,  
17 e85568.
- 18 [29] Wippermann, A., Rupp, O., Brinkrolf, K., Hoffrogge, R., Noll, T., The DNA methylation landscape  
19 of Chinese hamster ovary (CHO) DP-12 cells. *J. Biotechnol.* 2015.
- 20 [30] Clarke, C., Henry, M., Doolan, P., Kelly, S., et al., Integrated miRNA, mRNA and protein  
21 expression analysis reveals the role of post-transcriptional regulation in controlling CHO cell  
22 growth rate. *BMC Genomics* 2012, 13, 656.
- 23 [31] Quispe-Tintaya, W., White, R.R., Popov, V.N., Vijg, J., Maslov, A.Y., Fast mitochondrial DNA  
24 isolation from mammalian cells for next-generation sequencing. *BioTechniques* 2013, 55, 133–  
25 136.
- 26 [32] Hahn, C., Bachmann, L., Chevreux, B., Reconstructing mitochondrial genomes directly from  
27 genomic next-generation sequencing reads—a baiting and iterative mapping approach. *Nucleic*  
28 *Acids Res.* 2013, 41, e129.
- 29 [33] Magoč, T., Salzberg, S.L., FLASH: fast length adjustment of short reads to improve genome  
30 assemblies. *Bioinforma. Oxf. Engl.* 2011, 27, 2957–2963.
- 31 [34] Bernt, M., Donath, A., Jühling, F., Externbrink, F., et al., MITOS: improved de novo metazoan  
32 mitochondrial genome annotation. *Mol. Phylogenet. Evol.* 2013, 69, 313–319.
- 33 [35] Laslett, D., Canbäck, B., ARWEN: a program to detect tRNA genes in metazoan mitochondrial  
34 nucleotide sequences. *Bioinformatics* 2008, 24, 172–175.
- 35 [36] Bolger, A.M., Lohse, M., Usadel, B., Trimmomatic: a flexible trimmer for Illumina sequence  
36 data. *Bioinformatics* 2014, 30, 2114–2120.
- 37 [37] Li, H., Durbin, R., Fast and accurate short read alignment with Burrows-Wheeler transform.  
38 *Bioinforma. Oxf. Engl.* 2009, 25, 1754–1760.
- 39 [38] Ding, J., Sidore, C., Butler, T.J., Wing, M.K., et al., Assessing Mitochondrial DNA Variation and  
40 Copy Number in Lymphocytes of ~2,000 Sardinians Using Tailored Sequencing Analysis Tools.  
41 *PLOS Genet* 2015, 11, e1005306.
- 42 [39] Van der Auwera, G.A., Carneiro, M.O., Hartl, C., Poplin, R., et al., in: *Curr. Protoc. Bioinforma.*,  
43 John Wiley & Sons, Inc., 2002.
- 44 [40] Koboldt, D.C., Zhang, Q., Larson, D.E., Shen, D., et al., VarScan 2: somatic mutation and copy  
45 number alteration discovery in cancer by exome sequencing. *Genome Res.* 2012, 22, 568–576.
- 46 [41] Wilm, A., Aw, P.P.K., Bertrand, D., Yeo, G.H.T., et al., LoFreq: a sequence-quality aware, ultra-  
47 sensitive variant caller for uncovering cell-population heterogeneity from high-throughput  
48 sequencing datasets. *Nucleic Acids Res.* 2012, 40, 11189–11201.
- 49 [42] Guo, Y., Li, J., Li, C.-I., Long, J., et al., The effect of strand bias in Illumina short-read sequencing  
50 data. *BMC Genomics* 2012, 13, 666.

- 1 [43] Cingolani, P., Platts, A., Wang, L.L., Coon, M., et al., A program for annotating and predicting  
2 the effects of single nucleotide polymorphisms, SnpEff: SNPs in the genome of *Drosophila*  
3 *melanogaster* strain w1118; iso-2; iso-3. *Fly (Austin)* 2012, 6, 80–92.
- 4 [44] Choi, Y., Sims, G.E., Murphy, S., Miller, J.R., Chan, A.P., Predicting the Functional Effect of  
5 Amino Acid Substitutions and Indels. *PLOS ONE* 2012, 7, e46688.
- 6 [45] Just, R.S., Irwin, J.A., Parson, W., Mitochondrial DNA heteroplasmy in the emerging field of  
7 massively parallel sequencing. *Forensic Sci. Int. Genet.* 2015, 18, 131–139.
- 8 [46] Calabrese, F.M., Simone, D., Attimonelli, M., Primates and mouse NumtS in the UCSC Genome  
9 Browser. *BMC Bioinformatics* 2012, 13 Suppl 4, S15.
- 10 [47] Partridge, M.A., Davidson, M.M., Hei, T.K., The complete nucleotide sequence of Chinese  
11 hamster (*Cricetulus griseus*) mitochondrial DNA. *DNA Seq. J. DNA Seq. Mapp.* 2007, 18, 341–  
12 346.
- 13 [48] Van der Auwera, G.A., Carneiro, M.O., Hartl, C., Poplin, R., et al., From FastQ data to high  
14 confidence variant calls: the Genome Analysis Toolkit best practices pipeline. *Curr. Protoc.*  
15 *Bioinforma. Ed. Board Andreas Baxevanis AI* 2013, 43, 11.10.1-33.
- 16 [49] McKenna, A., Hanna, M., Banks, E., Sivachenko, A., et al., The Genome Analysis Toolkit: A  
17 MapReduce framework for analyzing next-generation DNA sequencing data. *Genome Res.*  
18 2010, 20, 1297–1303.
- 19 [50] Li, H., Handsaker, B., Wysoker, A., Fennell, T., et al., The Sequence Alignment/Map format and  
20 SAMtools. *Bioinformatics* 2009, 25, 2078–2079.
- 21 [51] Belle, E.M.S., Piganeau, G., Gardner, M., Eyre-Walker, A., An investigation of the variation in  
22 the transition bias among various animal mitochondrial DNA. *Gene* 2005, 355, 58–66.
- 23 [52] Fan, W., Lin, C.S., Potluri, P., Procaccio, V., Wallace, D.C., mtDNA lineage analysis of mouse L-  
24 cell lines reveals the accumulation of multiple mtDNA mutants and intermolecular  
25 recombination. *Genes Dev.* 2012, 26, 384–394.
- 26 [53] Bayona-Bafaluy, M.P., Acín-Pérez, R., Mullikin, J.C., Park, J.S., et al., Revisiting the mouse  
27 mitochondrial DNA sequence. *Nucleic Acids Res.* 2003, 31, 5349–5355.
- 28 [54] Wurm, F.M., CHO Quasispecies—Implications for Manufacturing Processes. *Processes* 2013, 1,  
29 296–311.
- 30 [55] Jayaprakash, A.D., Benson, E.K., Gone, S., Liang, R., et al., Stable heteroplasmy at the single-cell  
31 level is facilitated by intercellular exchange of mtDNA. *Nucleic Acids Res.* 2015, gkv052.
- 32 [56] Zheng, J., Energy metabolism of cancer: Glycolysis versus oxidative phosphorylation (Review).  
33 *Oncol. Lett.* 2012, 4, 1151–1157.
- 34 [57] Templeton, N., Dean, J., Reddy, P., Young, J.D., Peak antibody production is associated with  
35 increased oxidative metabolism in an industrially relevant fed-batch CHO cell culture.  
36 *Biotechnol. Bioeng.* 2013, 110, 2013–2024.
- 37 [58] Jo, A., Ham, S., Lee, G.H., Lee, Y.-I., et al., Efficient Mitochondrial Genome Editing by  
38 CRISPR/Cas9, Efficient Mitochondrial Genome Editing by CRISPR/Cas9. *BioMed Res. Int. BioMed*  
39 *Res. Int.* 2015, 2015, 2015, e305716.
- 40 [59] Bacman, S.R., Williams, S.L., Pinto, M., Peralta, S., Moraes, C.T., Specific elimination of mutant  
41 mitochondrial genomes in patient-derived cells by mitoTALENs. *Nat. Med.* 2013, 19, 1111–  
42 1113.
- 43 [60] Krzywinski, M.I., Schein, J.E., Birol, I., Connors, J., et al., Circos: An information aesthetic for  
44 comparative genomics. *Genome Res.* 2009.
- 45

## 1 **Figure legends**

2 **Figure 1: The *Cricetulus griseus* mitochondrial genome.** mtDNA was extracted from Chinese  
3 hamster liver tissue, deep sequenced and reconstructed using the MITOBIM algorithm. The resulting  
4 16,283bp mitochondrial DNA sequence had an average depth of coverage of ~6,417X (the depth of  
5 coverage at each base is shown as a grey histogram within the inner circle of the plot). 13 proteins, 22  
6 tRNAs, 2 ribosomal RNAs and the non-coding D-loop control region are encoded by the Chinese  
7 hamster mitochondrial genome. 28 genes are present on the H-strand while 9 genes are on the L-strand  
8 of the mtDNA. Comparison of the Chinese hamster mtDNA sequence to the previously sequenced CHO  
9 cell line sequence on GenBank (accession no. NC\_007936) revealed 7 mutations (4 SNPs and 3  
10 INDELs) in protein coding genes (*COX3*, *ND4*, *ND5*, *ND6*) and the D-loop. The plot was generated  
11 using circos v.0.67 [60].

12 **Figure 2: High resolution sequencing of the CHO cell mitochondrial genome. (A)** Number of  
13 mapped reads for the 22 CHO cell lines illustrating the number of unmapped reads (MAPQ < 20),  
14 number of PCR duplicates and number of uniquely mapped reads. Note: a negligible difference in the  
15 total numbers of reads mapped against the shifted reference sequence were observed (Supplementary  
16 Table 4B & 4C). **(B)** Average sequencing depth for each CHO cell line sequenced. The average depth  
17 of coverage combines both variant calling regions within the unshifted and shifted reference sequences  
18 used for variant detection. **(C)** The mean (blue), minimum and maximum depth of coverage at each  
19 base position across the 22 CHO cell mitochondrial genomes sequenced. The utilisation of the dual  
20 mapping approach resulted in a 1,128X and 1,079X increase in coverage for the first and last 100bp of  
21 the mtDNA reference sequence respectively.

22 **Figure 3: The mutational landscape of CHO cell line mitochondrial genome.** A total of 197  
23 mutations (175 SNPs and 22 INDELs) were identified for 22 CHO cell line samples in comparison the  
24 Chinese hamster mitochondrial genome at 130 positions in the mtDNA sequence. The nucleotide  
25 alternations detected for each cell line are shown along with the mtDNA feature harbouring the mutation.  
26 All protein coding genes along with the D-loop and 8 tRNA genes had a least 1 variant in 1 CHO cell  
27 line.

28 **Figure 4: Mitochondrial heteroplasmy is widespread in CHO cell lines. (A)** Number of SNPs and  
29 INDELs identified for 22 CHO cell line mitochondrial genomes sequenced. The mtDNA of 2 of the 3  
30 DG44 cell lines sequenced were found to have the most variant positions while the ATCC CHO-S cell  
31 line was had only a single mutated position. **(B)** Number of variant positions for the 24 genes found to  
32 harbour at least one mutation across the 22 CHO cell lines sequenced. **(C)** The majority of variants  
33 identified were heteroplasmic - 160 mtDNA mutations out of 197 mutation variants had an allele  
34 frequency < 99%.

35 **Figure 5: Prediction of the effects of CHO cell mitochondrial genome mutations.** The majority of  
36 mutations identified in this study occurred in protein coding genes. **(A)** 22 of the CHO cell lines  
37 sequenced contained at least one mutation that altered the protein sequence. **(B)** Each protein coding  
38 gene harboured a mutation that altered the protein sequence in at least one of the cell lines.

1 **Table legends**

2 **Table 1: Comparison of the Chinese hamster mitochondrial genome to the previously**  
3 **sequenced CHO cell line.** 7 variants were identified within 4 protein coding genes as well as the D-  
4 loop control region. The position, mtDNA gene, type and nucleotide change observed are shown for  
5 each variant. Chinese hamster reads were mapped against the CHO cell line mtDNA sequence  
6 (GenBank accession no. NC\_007936) to determine the number of reads supporting each variant. All  
7 variants were homoplasmic (allele frequency > 99%).

8 **Table 2: Heteroplasmic mitochondrial variants identified in 2 or more CHO cell lines.** In total 11  
9 individual heteroplasmic mutations were identified in two or more cell lines. A number of heteroplasmic  
10 sites were shared from the cohort of cells acquired from an identical cell line development project. 4  
11 heteroplasmic sites were shared amongst the 4 fast or 4 slow growing CHO cell lines clones. A  
12 heteroplasmic frameshift mutation in *Cytb* was identified in all clones analysed.

13 **Supplementary Tables**

14 **Supplementary Table 1: Sample Information.** The origin of each sample sequenced in this study  
15 along with the respective file names corresponding to mtDNA sequencing data.

16 **Supplementary Table 2: Primer sequences for fragment amplification of the CHO cell**  
17 **mitochondrial genome as well as mtDNA isolated from Chinese hamster and mouse mtDNA.**

18 Primers were designed for PCR amplification of genomic CHO DNA using the CHOgenome.org  
19 database. A beta-actin control was designed against the whole CHO genome sequence with the reverse  
20 primers spanning an intron/exon junction. A small positive control (PC) for the CHO mitochondrial  
21 genome was designed using the available CHO-K1 mitochondrial genome sequence on  
22 CHOgenome.org. Primer pairs (#1 and #2) were also designed using the available CHO-K1  
23 mitochondrial genome sequence for high-fidelity PCR of the mitochondrial genomic DNA fragments.  
24 Finally, to account for possible sequence variation between the CHO mt-DNA sequence and the  
25 Chinese hamster mt-DNA sequence, primer sets (CH1-6/8 and CH2-86/7) were designed used the  
26 available CHO-K1 mt-DNA sequence available and matched for 100% sequence similarity with both  
27 the mouse (mmu) and rat (rno) mt-DNA sequences.

28 **Supplementary Table 3: Chinese hamster mtDNA annotation.** Coordinates for each of the 38  
29 annotated features in the *C.griseus* mitochondrial genome.

30 **Supplementary Table 4: Pre-processed read counts, mapping rates and average depth of**  
31 **coverage for CHO cell line sequencing.** The number of reads remaining following pre-processing  
32 along mapping and average coverage rate for each sample following shifted and unshifted sequence  
33 alignment.

34 **Supplementary Table 5: SNP calling outputs.** Variant calling outs for 22 individual CHO cell lines  
35 including the depth at coverage at each SNP, number of forward and reverse reads supporting the  
36 reference and alternative allele, average base quality (Q) for each nucleotide at the SNP position and  
37 snpEff annotation.

1 **Supplementary Table 6: INDEL detection data.** Variant calling outputs for 22 individual CHO cell lines  
2 including the depth at coverage at each INDEL, number of forward and reverse reads supporting the  
3 reference and alternative allele, average base quality (Q) at the INDEL position as well 5 upstream and  
4 downstream flanking regions nucleotide and the snpEff annotation.

5 **Supplementary Table 7: mtDNA SNP and INDEL summary.** Base changes along with allele  
6 frequency at each variant position for the 22 CHO cell lines. Mutation annotations outputted by snpEff  
7 tool are provided for each SNP and INDEL. For missense variants the Provean classification and score  
8 are included.

9 **Supplementary Table 8: Detection of mutations in mtDNA isolated from the mouse liver sample.**  
10 Using an identical dual mapping approach aligning reads to the unshifted and shifted mouse  
11 mitochondrial genome reference sequence we identified 3 mutations, two of which were heteroplasmic.

## 12 **Supplementary Figures**

13 **Supplementary Figure 1: mtDNA extraction, isolation and amplification.** Ethidium bromide stained  
14 agarose gels were ran to assess the quality of **(A)** Genomic DNA isolated from liver samples of Chinese  
15 Hamster using a DNeasy Kit with high integrity DNA running at ~100kb. **(B)** mitochondrial plasmid DNA  
16 isolated from CHO cell lines using a modified plasmid mini-prep kit with plasmid DNA running out as  
17 two bands (coiled and super-coiled) and **(C)** Mitochondrial genomic DNA fragments amplified by high-  
18 fidelity PCR and visualised individually and pooled

19 **Supplementary Figure 2: Dual mapping strategy and variant calling.** In this study we utilized a dual  
20 mapping strategy to account for the mitochondrial genome circularity to maximize depth of coverage at  
21 the beginning and end of the mtDNA sequence. We also required agreement between two different  
22 variant calling algorithms for a SNP or INDEL to be reported.

23 **Supplementary Figure 3: Alignment of potential NumtS to mtDNA.** The mouse mtDNA genome  
24 was utilised as a control to determine if NumtS were influencing heteroplasmy measurements. Reads  
25 were first stringently aligned to know NumtS region the mm9 genome. Those reads that aligned to  
26 NumtS were extracted and remapped against the mouse mitochondrial region. The IGV diagram  
27 illustrates that potential NumtS reads aligned to 5 regions. The grey bars illustrate the location of reads  
28 that align to both the mouse nuclear and mitochondrial genomes demonstrating that Numt  
29 contamination did not affect the 3 variants identified. In addition, no SNPs or INDELS were identified  
30 within the 5 regions.

31 **Supplementary Figure 4: Alignment of amplified *CYTB* fragment sequences.** Three clones from  
32 the two cell line development projects were selected for the amplification and Sanger sequencing of a  
33 short fragment of the *CYTB* gene which included an identified frameshift and STOP mutation.  
34 Sequences were aligned using the online software MultAlin. Sequence variations are highlighted in blue  
35 with full conservation in red. Each sequences starts with the initiating start ATG codon.

1 **Tables**

2 **Table 1**

<b>Position</b>	<b>Gene</b>	<b>Variant type</b>	<b>CHO Cell allele</b>	<b><i>C.griseus</i> allele</b>	<b># reads supporting variant</b>
9306	<i>COX3</i>	SNP	A	G	5,201
11399	<i>ND4</i>	SNP	C	A	8,765
12683	<i>ND5</i>	SNP	G	A	11,464
13436	<i>ND6</i>	Deletion	AT	A	8,815
13456	<i>ND6</i>	Insertion	C	CA	8,461
15717	D-loop	Deletion	GC	C	7,522
15888	D-loop	SNP	G	T	8,056

3

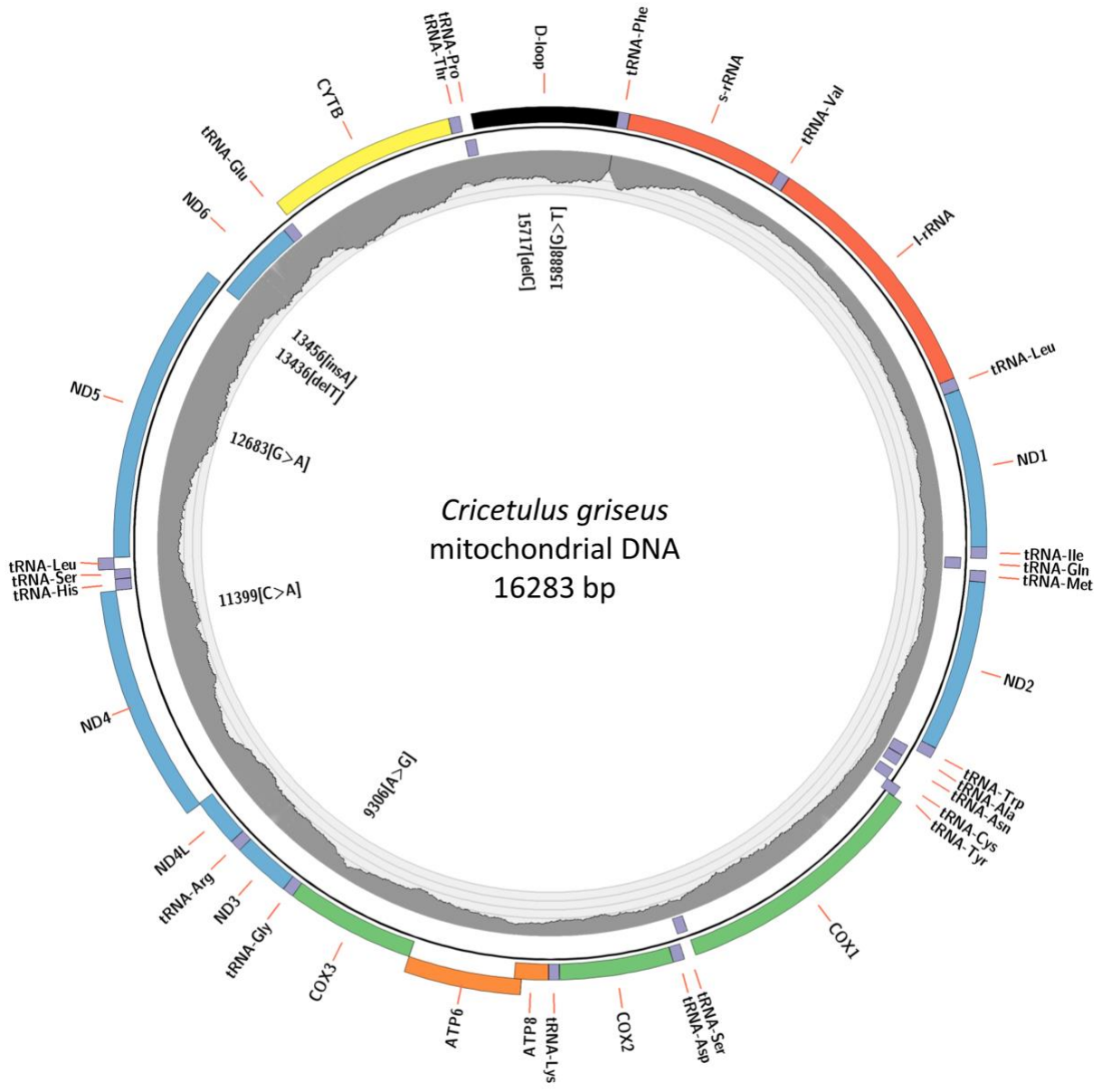
1 **Table 2.**

Gene	Position	Allele		Type	Minor Allele Freq. (%)		Effect	CHO Cell lines
		Ref	Alt		Min	Max		
<i>12S rRNA</i>	105	G	A	SNP	1.7	3	Non-coding	Pfizer CHO-K1 CLD1 #1 Pfizer CHO-K1 CLD1 #2 Pfizer CHO-K1 CLD1 #4
<i>16S rRNA</i>	1575	G	A	SNP	1.2	50	Non-coding	Pfizer CHO-K1 CLD2 #3 ATCC DGG44
<i>16S rRNA</i>	2151	C	T	SNP	1.2	9.9	Non-coding	Pfizer CHO-K1 CLD1 #1 Pfizer CHO-K1 CLD1 #2 Pfizer CHO-K1 CLD1 #3 Pfizer CHO-K1 CLD1 #4
<i>ND1</i>	3205	GCT	G	Deletion	1.5	50	Frameshift variant	Pfizer CHO-K1 CLD1 #5 Pfizer CHO-K1 CLD2 #2 Pfizer CHO-K1 CLD2 #3
<i>COX2</i>	6996	G	A	SNP	14	26	Initiation codon variation	Pfizer CHO-K1 CLD1 #1 Pfizer CHO-K1 CLD1 #2 Pfizer CHO-K1 CLD1 #3 Pfizer CHO-K1 CLD1 #4
<i>tRNA<sup>Lys</sup></i>	7721	A	G	SNP	9.5	47	Non-coding	Pfizer CHO-K1 CLD1 #5 Pfizer CHO-K1 CLD2 #1 Pfizer CHO-K1 CLD2 #3
<i>ATP6</i>	8067	C	T	SNP	1.8	6.8	Missense variation	Pfizer CHO-K1 CLD1 #1 Pfizer CHO-K1 CLD1 #2 Pfizer CHO-K1 CLD1 #3 Pfizer CHO-K1 CLD1 #4
<i>ND4</i>	11431	TCA	T	Deletion	1.7	2.0	Frameshift variant	Biogen CHO-K1 ATCC DGG44 Biogen DG44 #2
<i>tRNA<sup>Leu</sup></i>	11659	AG	A	Deletion	1.2	2.8	Non-coding	Pfizer CHO-K1 CLD1 #1 Pfizer CHO-K1 CLD1 #2
<i>CYTB</i>	14136	GA	G	Deletion	7.7	52	Frameshift variant	ATCC CHO-K1 Pfizer CHO-K1 CLD1 #1 Pfizer CHO-K1 CLD1 #2 Pfizer CHO-K1 CLD1 #3 Pfizer CHO-K1 CLD1 #4 Pfizer CHO-K1 CLD1 #5 Pfizer CHO-K1 CLD2 #1 Pfizer CHO-K1 CLD2 #2 Pfizer CHO-K1 CLD2 #3
	14378	G	A	SNP	11	36.0	Premature stop codon	Pfizer CHO-K1 CLD1 #1 Pfizer CHO-K1 CLD1 #2 Pfizer CHO-K1 CLD1 #3 Pfizer CHO-K1 CLD1 #4
	14849	G	A	SNP	1.5	53	Missense variation	Pfizer CHO-K1 CLD1 #5 Pfizer CHO-K1 CLD2 #1 Pfizer CHO-K1 CLD2 #2 Pfizer CHO-K1 CLD2 #3
	15136	C	A	SNP	1.7	19	Missense variation	Pfizer CHO-K1 CLD1 #3 Biogen DG44 #1

2

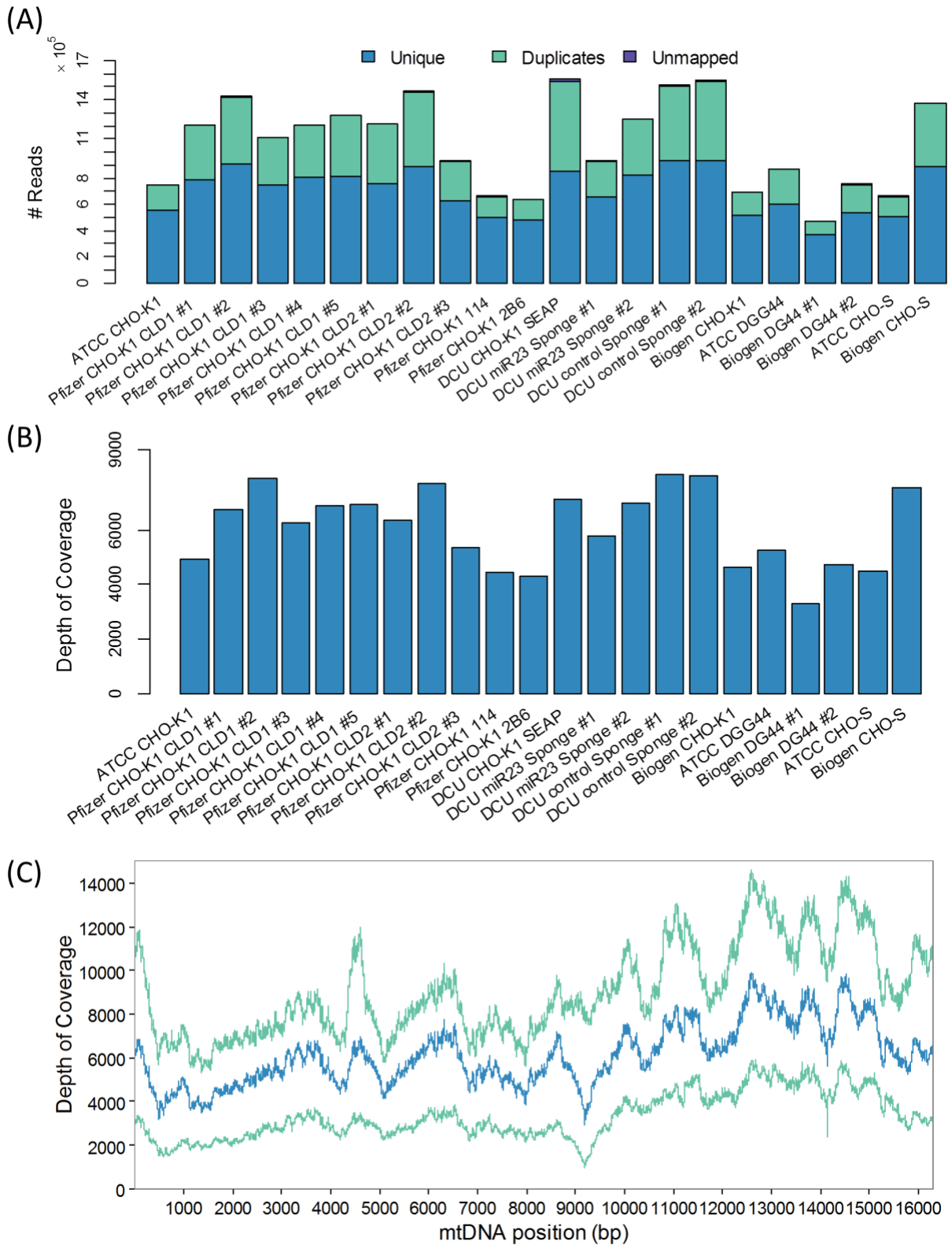
3

1 Figure 1



2  
3

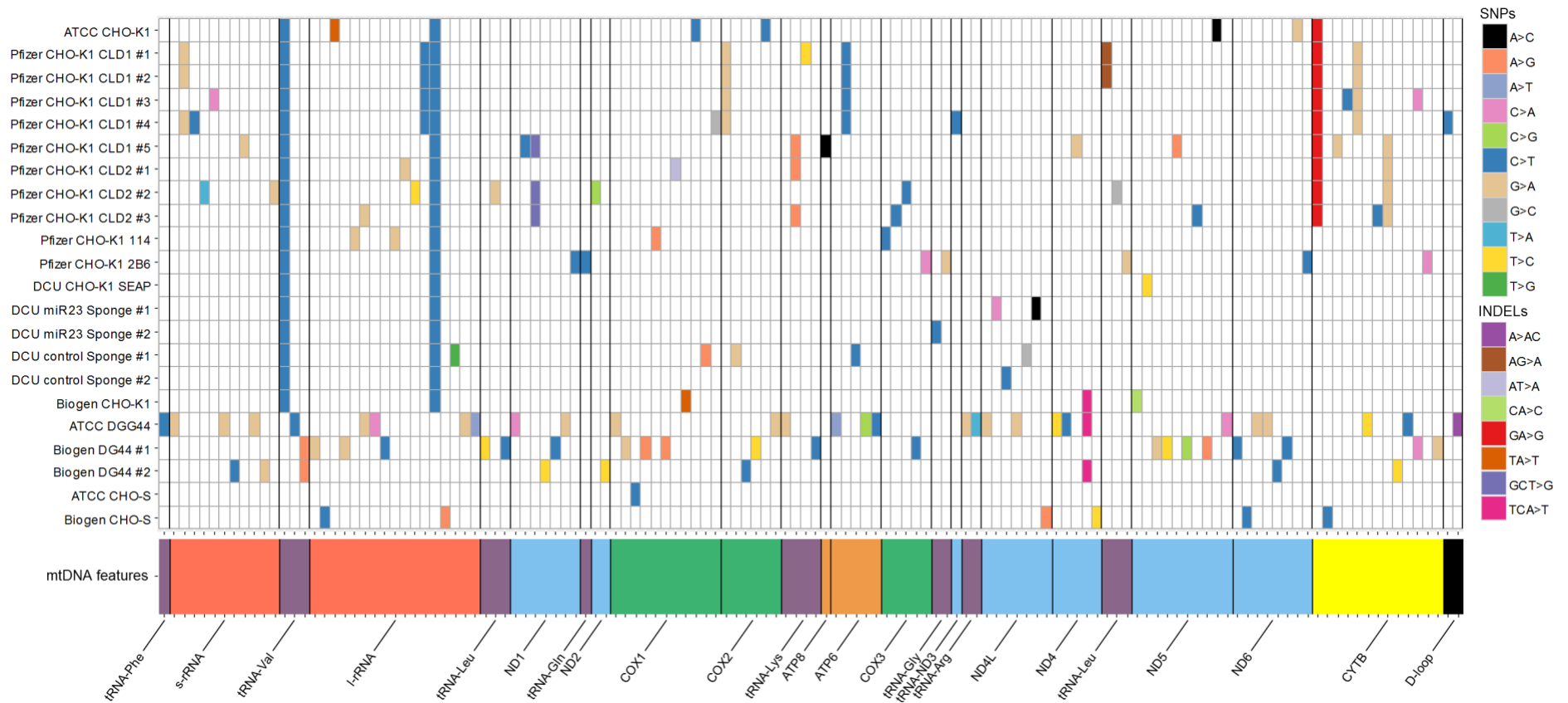
1 **Figure 2**



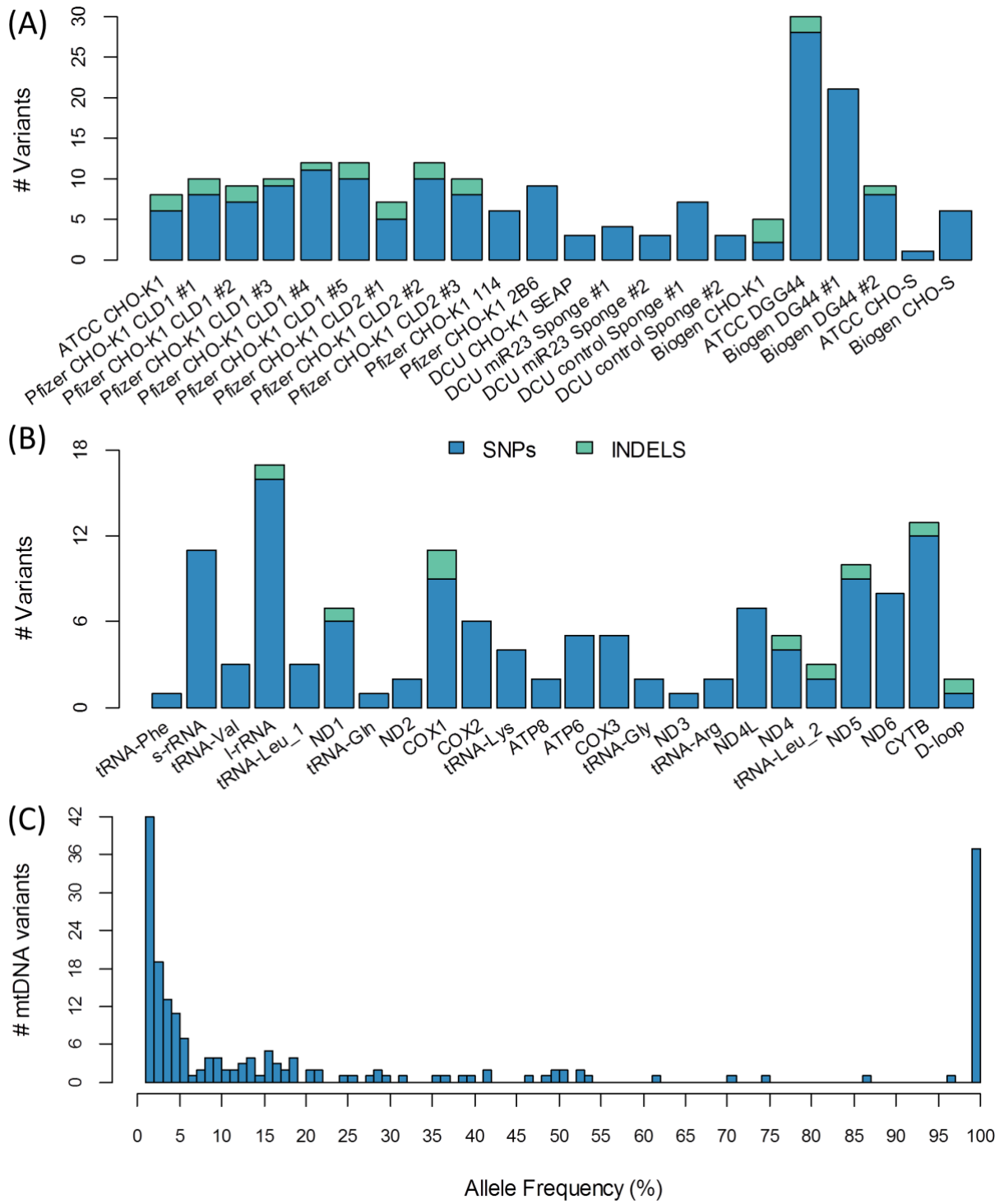
2

3

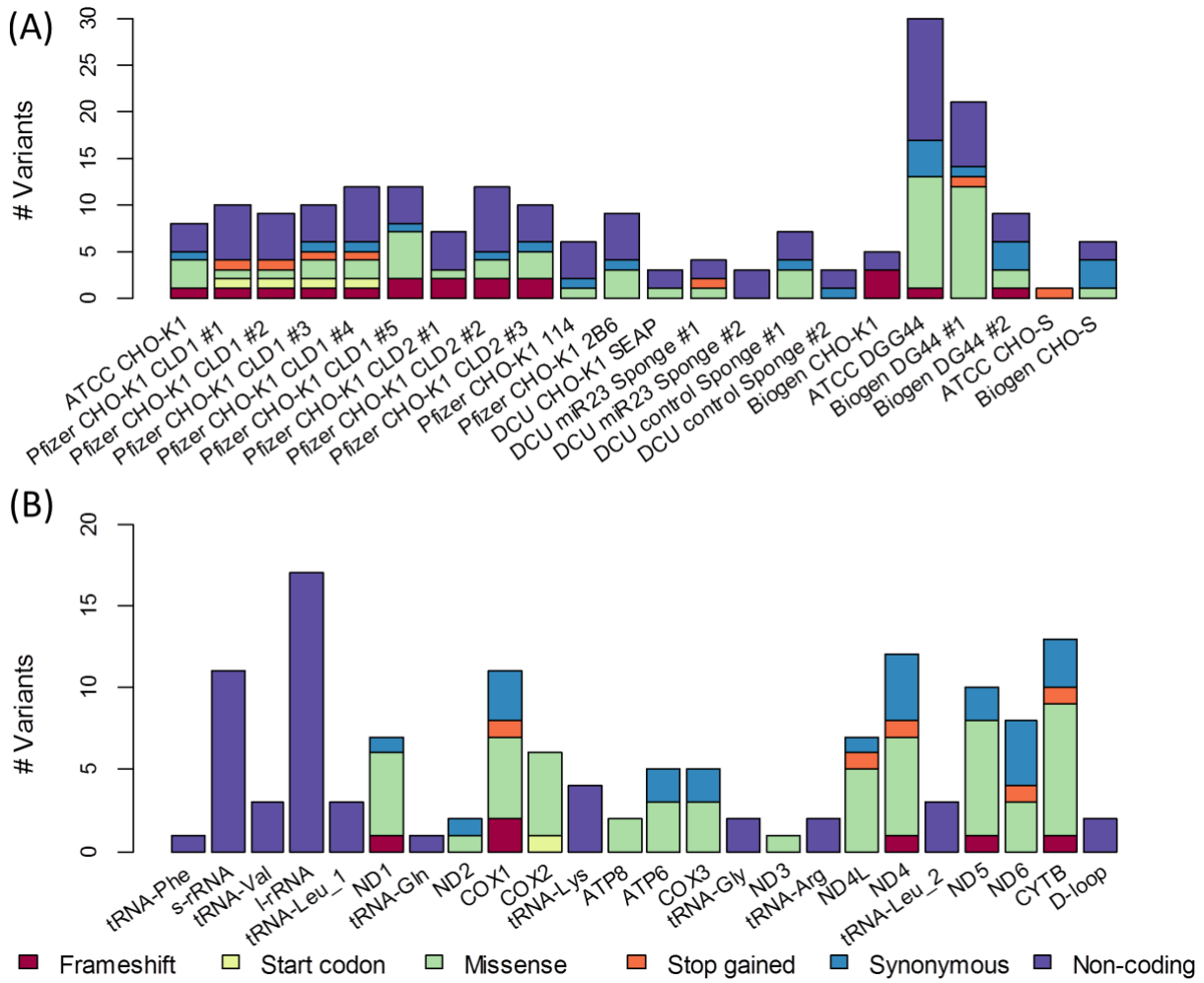
**Figure 3**



**Figure 4**

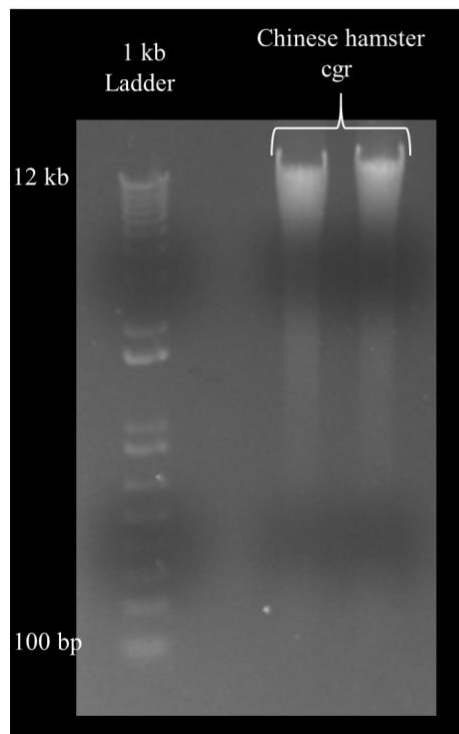


**Figure 5**

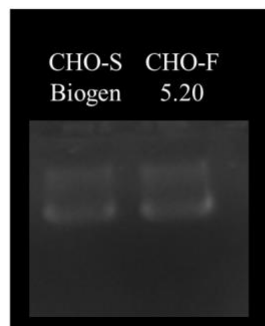


Supplementary Figure 1

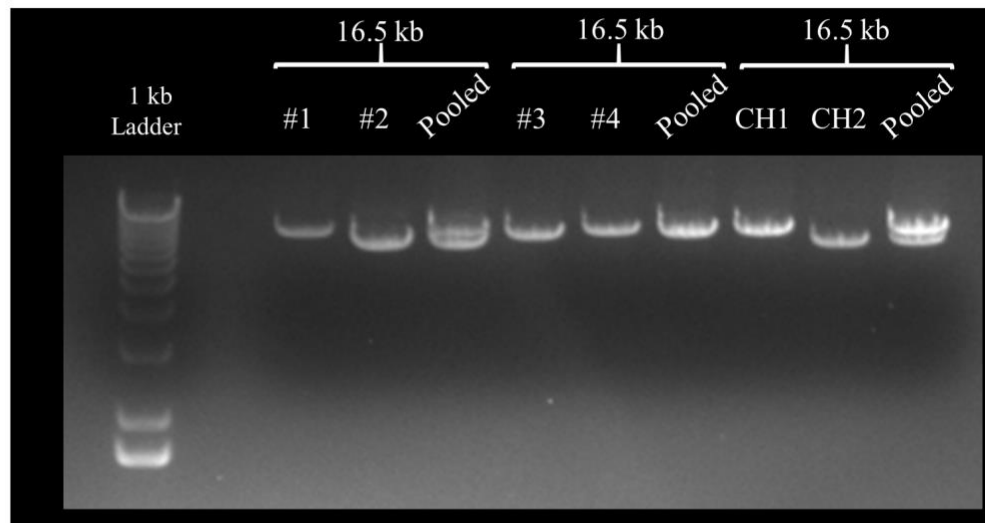
(A)



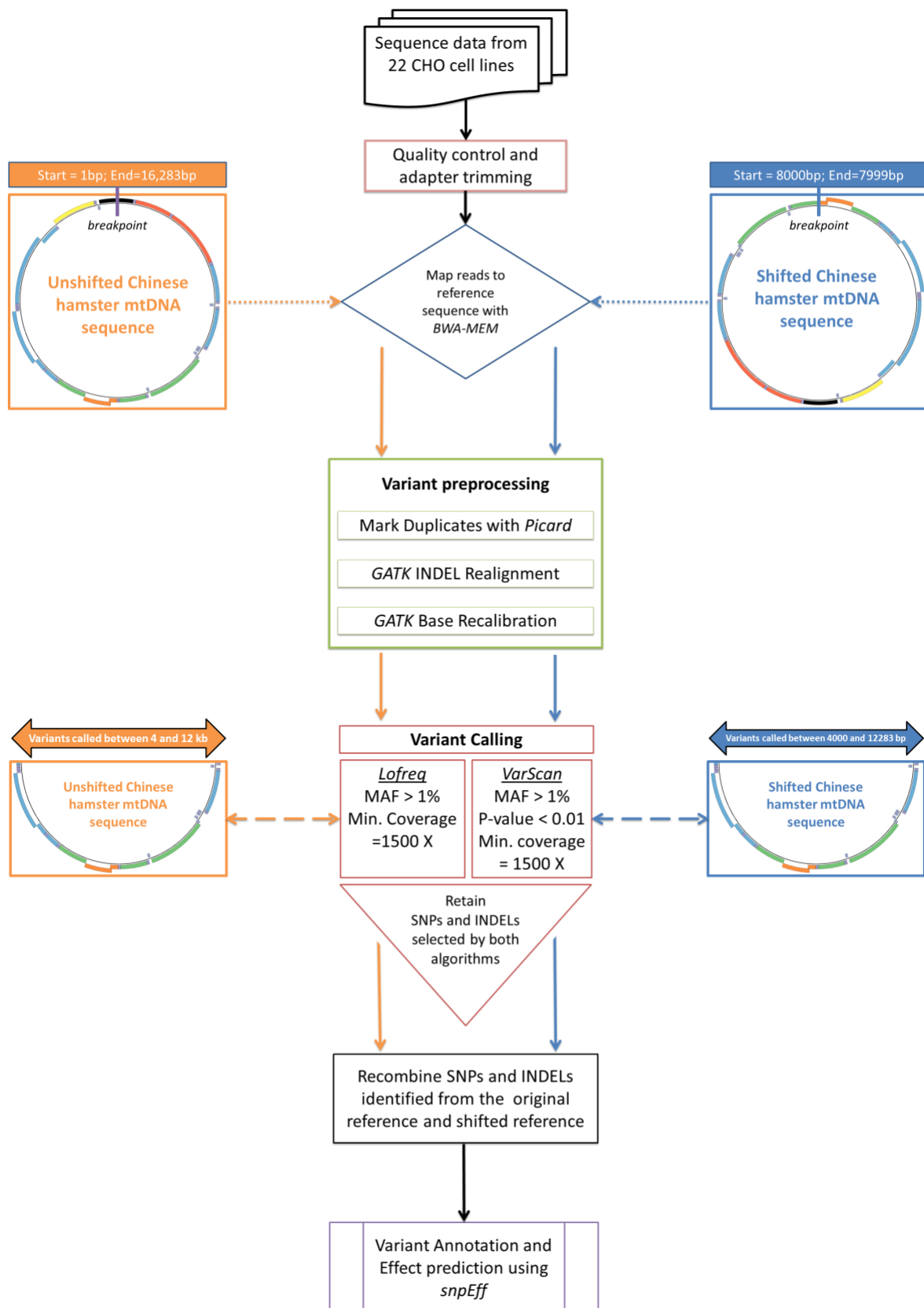
(B)



(C)



Supplementary Figure 2



### Supplementary Figure 3

