# DataQuality

*DI Cruz Dávalos*

*7/26/2019*

# Methods

## Samples

The samples included in this study were taken from 24 individuals whose remains belonged to the anthropological collection from the National Museum in Rio de Janeiro, Brazil. Throughout this document we keep the National Museum's identifiers to refer to such individuals. Twenty-two of the individuals were identified from the Museum's archival as belonging to a "Botocudo" group. Similarly, one mummified individual from a cave in Minas Gerais, Brazil, was not associated to any specific group, and one individual was excavated from a shell mound from Santa Catarina, Brazil. Regarding the twenty-two Botocudos, one tooth was sampled for twenty of them, and one piece of skull and petrous bone for the remaining two individuals (MN00019 and MN0008, respectively). One tooth was sampled from the Minas Gerais' mummy and one tooth from the shell mound's individual. Two of the Botocudo individuals presented here (MN00013 and MN00065) have been previously studied by [**?**], but no genomic data was reported at the time.

## DNA extraction, library preparation, and sequencing

```
Twenty-four DNA extracts were prepared: one from petrous bone, one from a piece of skull and twenty-two
```

## Genomic data quality assessment

### Mapping

Remnants of adapters, low-quality bases and nucleotides reported as "N" were trimmed from the reads with AdapterRemoval version 2.1.7 ([**?**]). Reads of 30 bp length and above were mapped to the human genome reference *built 19* with BWA aln version 0.7.15 ([**?**]), disabling the seed to avoid mapping bias due to damage at the 5' termini of the reads ([**?**]). Reads with a mapping quality score equal or greater than 30 were retained. Duplicate reads were identified and removed with picard tools MarkDuplicates version 2.9.0 (http://broadinstitute.github.io/picard/), and indel realignment was performed with GATK version 3.7 with default options ([**?**]). Molecular damage parameters were obtained with mapDamage2 ([**?**]).

### Contamination estimation

```
We estimated contamination using a Bayesian statistical approach on mitochondrial data (\cite{Fu2013}),
```

### Error rate estimation

### Molecular sex determination and uniparental markers

```
Molecular sex was determined by computing the ratio of reads mapping to the Y chromosome with respect t
To call Y-chromosome and mitochondrial haplogroups, we used ANGSD version 0.921 (\cite{Korneliussen2014}
```

# Results

## Processing of genomes and ancient DNA authentication

We shotgun-sequenced 24 samples to an average depth of coverage between $0.001\times$ and $9.2\times$ (Table **??**).

```
<!-- %Between 21,554,888 and 1,124,846,215 reads were sequenced per sample.  -->
```

After trimming adaptors and removing low-quality bases from the reads, between 69.1% and 95.6% of the reads were retained and used as input for mapping.% (that is, between 16,504,949 and 1,030,165,295 reads per sample). The clonality levels (percentage of mapped reads classified as PCR duplicates) within samples ranged from 0.50% to 26.0%. After removing duplicates, we obtained between 70,205 and 533,336,166 reads mapped per sample. Therefore, the percentage of retained reads that were uniquely mapped (i.e., endogenous content) per sample ranged from 0.03% to 51.8%.

The sequenced reads show the common signatures of molecular damage observed in ancient samples, such as short lengths and high rates of deamination (Figure **??**). Retained and mapped reads had similar lengths, with an average of 47.1 - 67.5 bp and 42.0 - 67.1 bp per sample, respectively. Non-USER-treated libraries showed average deamination rates between –% and –% at the termini of the reads.

We estimated contamination based on mitochondrial data. The estimates vary according to the number of reads used in the calculations, which also varies depending on whether we consider transitions as mismatches to the endogenous mitochondrial genome. We notice a larger dispersion in the posterior distribution when using ~2,000 reads or less.

When accounting for all polymorphism types and samples with more than 2,000 mitochondrial reads (n = 19, mitochondrial coverage: $9.9\times$ - $222.4\times$), the maximum a posteriori estimate for contamination is between 0.69% and 8.41%. If we remove transitions from the estimation, we have 14 samples with more than 2,000 (mitochondrial coverage: $35.6\times$ - $222.4\times$) for which we estimated between 0.03% and 3.10% of contaminant reads.

Regarding samples with mitochondrial coverage above $10\times$, we estimated less than 6% of contaminant reads (Table **??**).
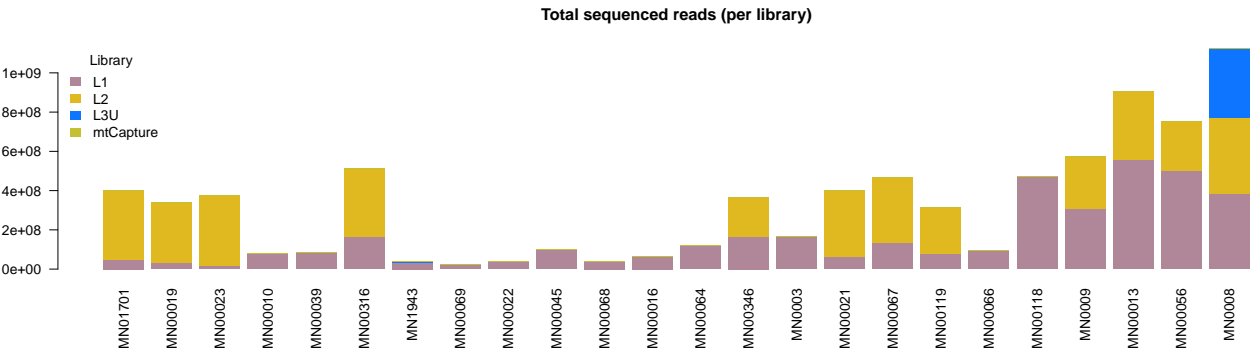
Data quality assessment

Molecular damage patterns
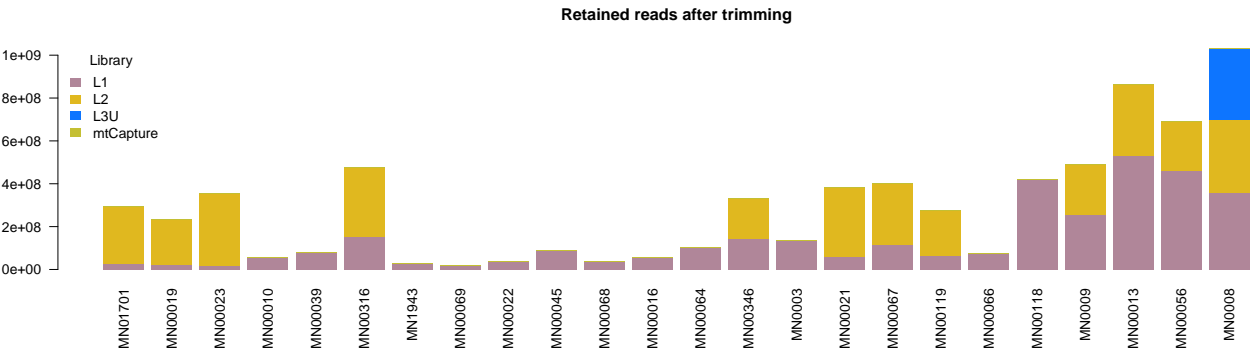
Contamination estimates
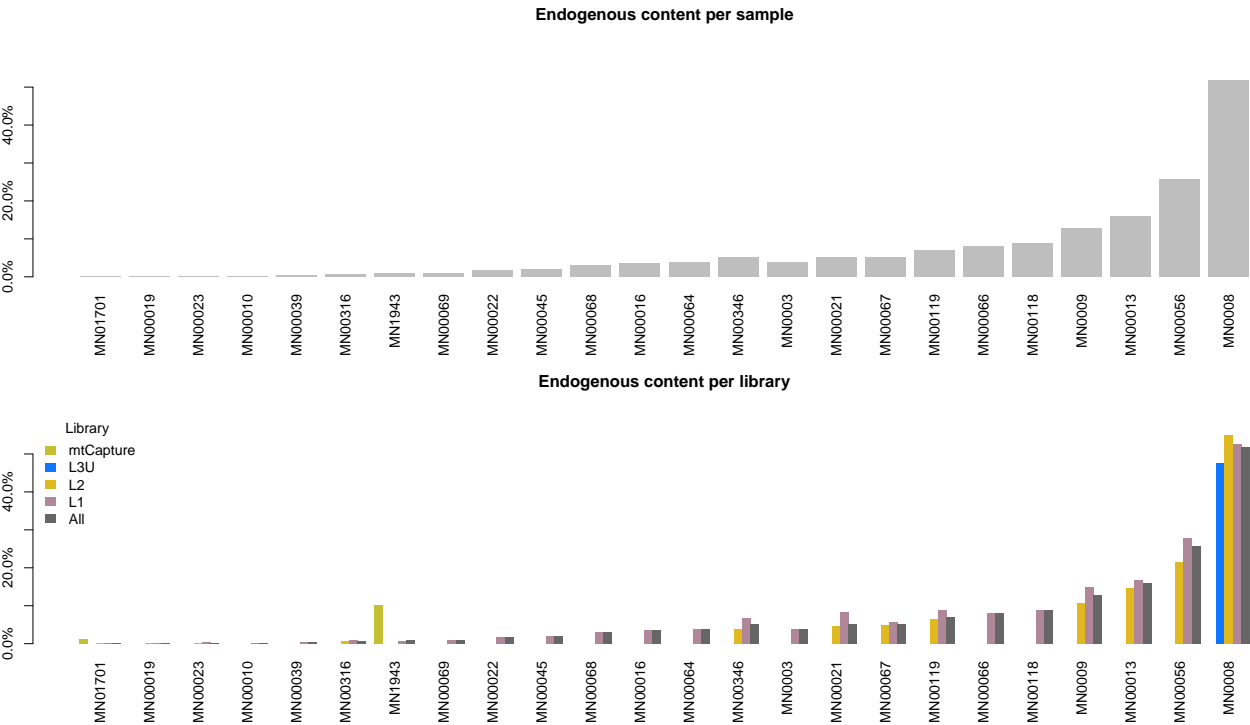
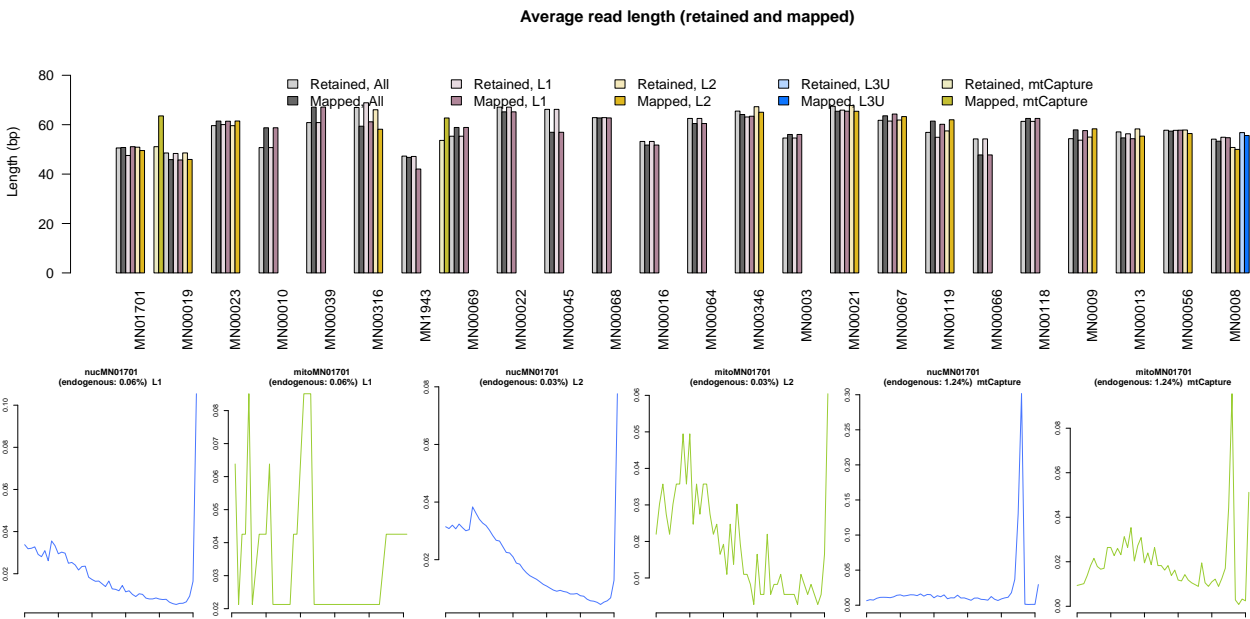Molecular sex determination

# Sequenced reads

## Total number of reads

**Total sequenced reads (per library)**

## Reads retained after trimming

**Retained reads after trimming**

**Retained reads after trimming (fraction)**

# Endogenous content

**Endogenous content per sample**



**Endogenous content per library**



# Read length

**Average read length (retained and mapped)**

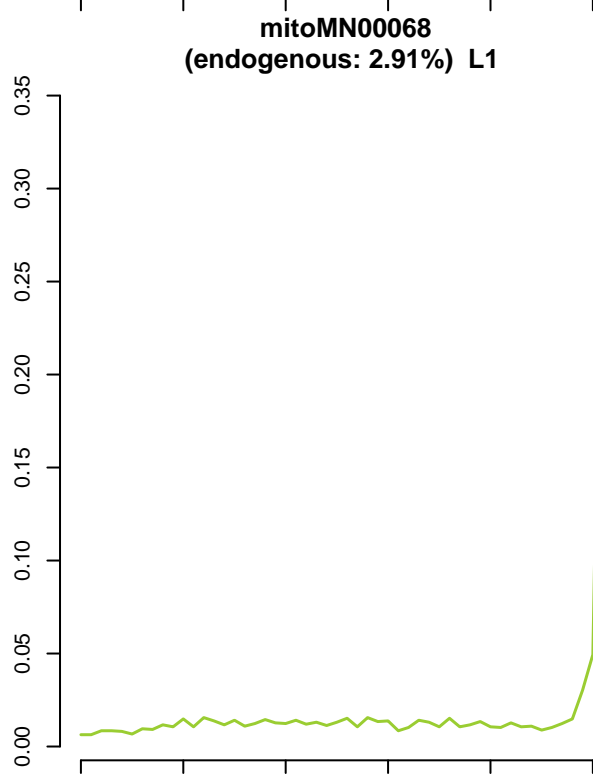nucMN00019
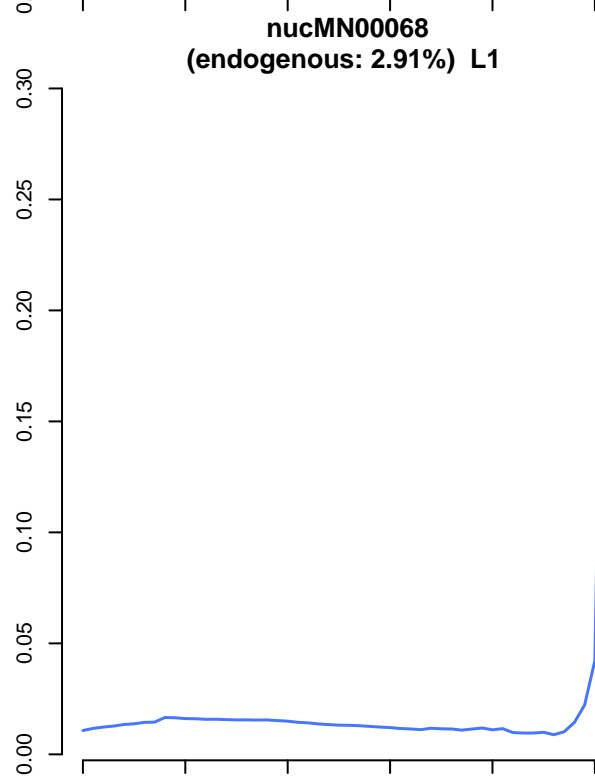(endogenous: 0.07%) L1

mitoMN00019
(endogenous: 0.07%) L1

nucMN00019
(endogenous: 0.06%) L2

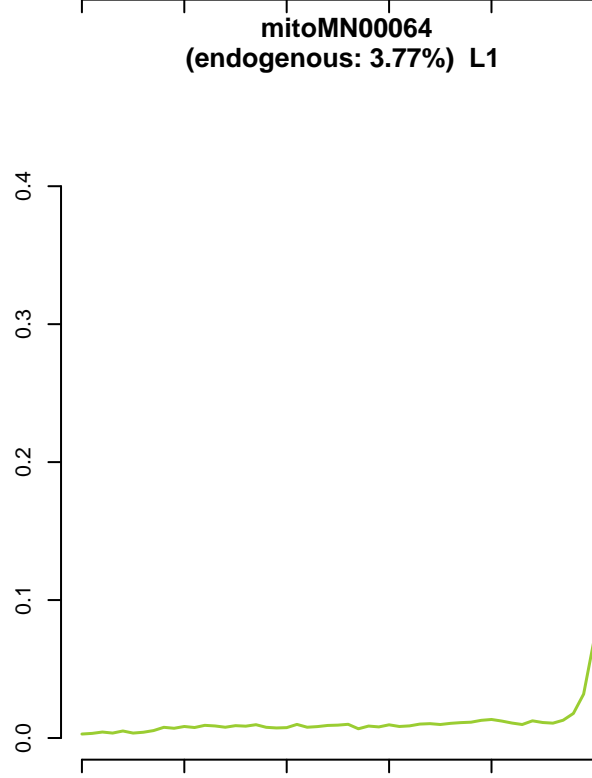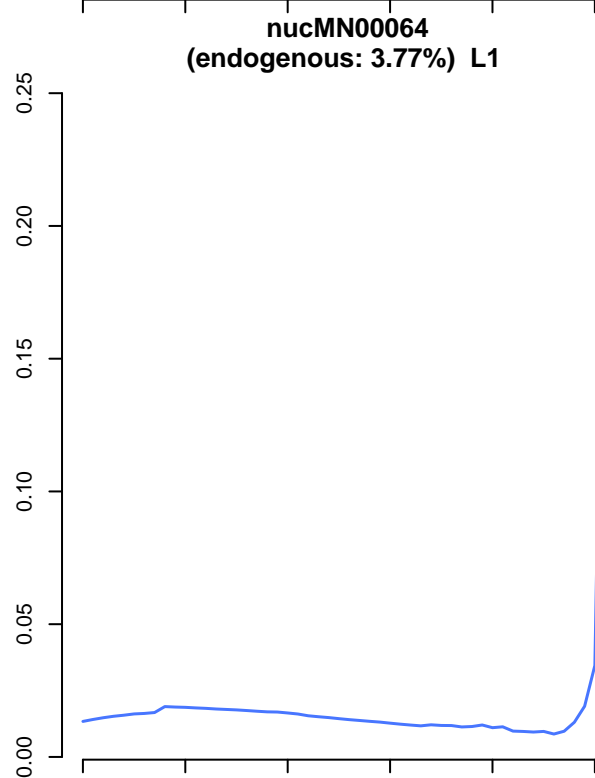mitoMN00019
(endogenous: 0.06%) L2
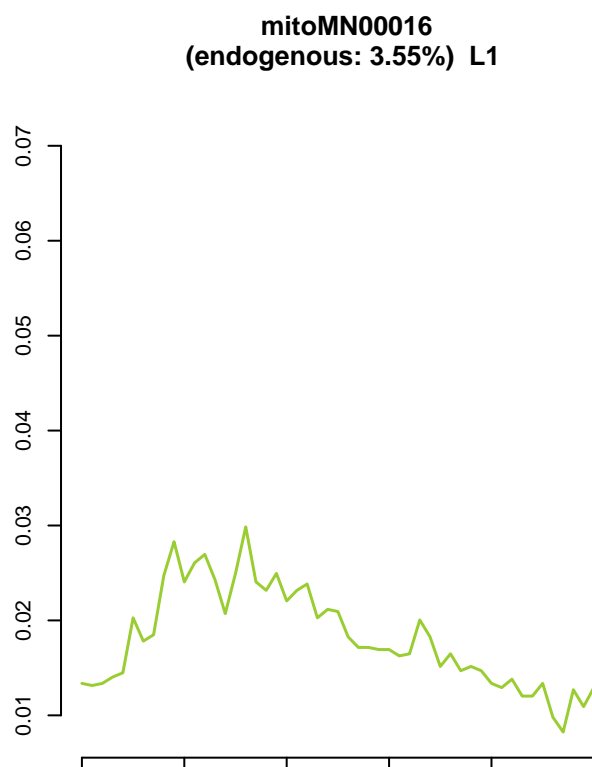
nucMN00023
(endogenous: 0.31%) L1

mitoMN00023
(endogenous: 0.31%) L1

nucMN00023
(endogenous: 0.11%) L2

mitoMN00023
(endogenous: 0.11%) L2

## nucMN00010
## (endogenous: 0.13%)  L1

## mitoMN00010
## (endogenous: 0.13%)  L1

**nucMN00039**
**(endogenous: 0.37%) L1**

**mitoMN00039**
**(endogenous: 0.37%) L1**

**nucMN00316**
**(endogenous: 0.95%) L1**

**mitoMN00316**
**(endogenous: 0.95%) L1**

**nucMN00316**
**(endogenous: 0.66%) L2**

**mitoMN00316**
**(endogenous: 0.66%) L2**

**nucMN1943**
**(endogenous: 0.77%) L1**

**mitoMN1943**
**(endogenous: 0.77%) L1**

**nucMN1943**
**(endogenous: 10.19%) mtCapture**

**mitoMN1943**
**(endogenous: 10.19%) mtCapture**
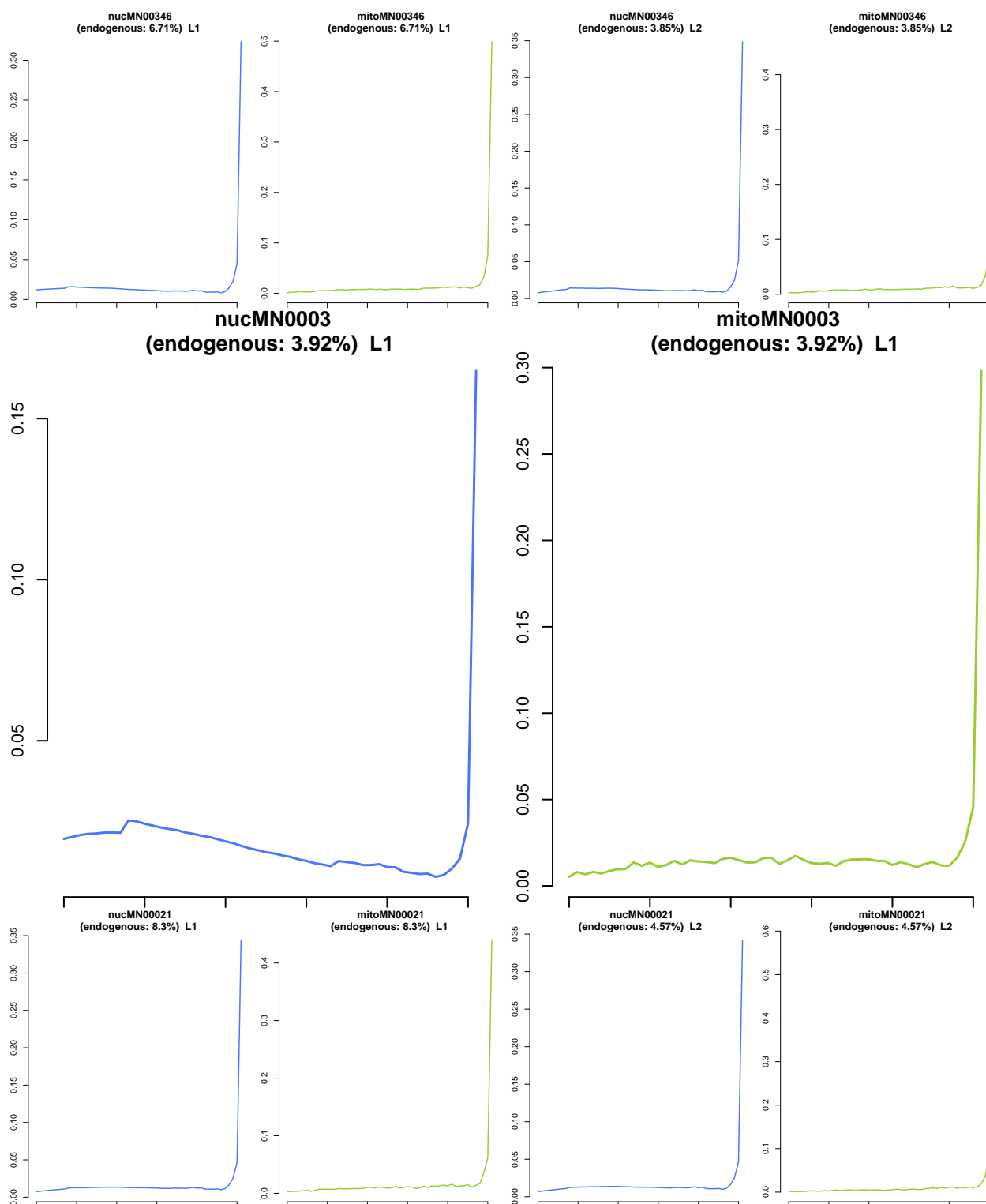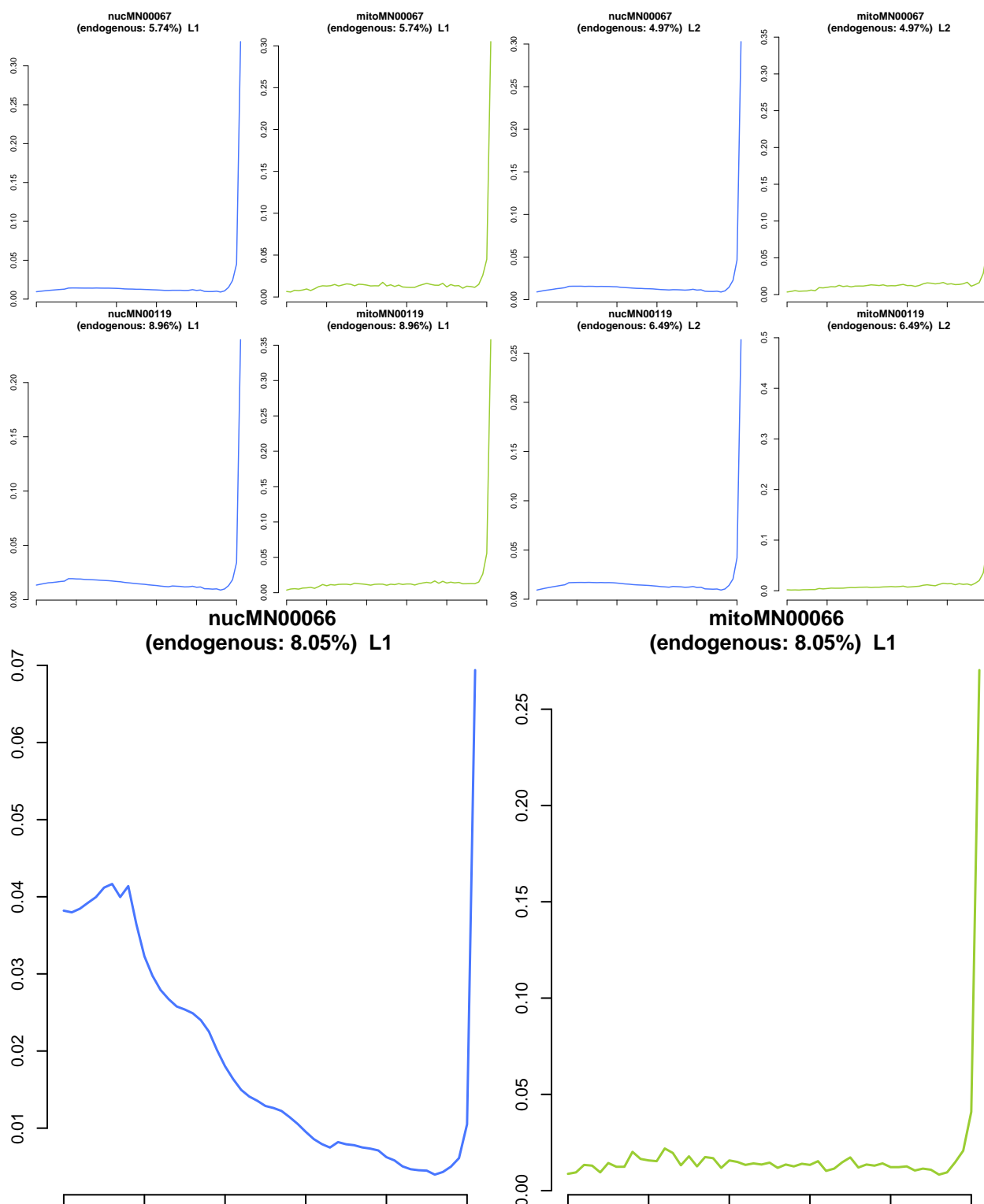
## nucMN00069
### (endogenous: 0.89%) L1

## mitoMN00069
### (endogenous: 0.89%) L1

## nucMN00022
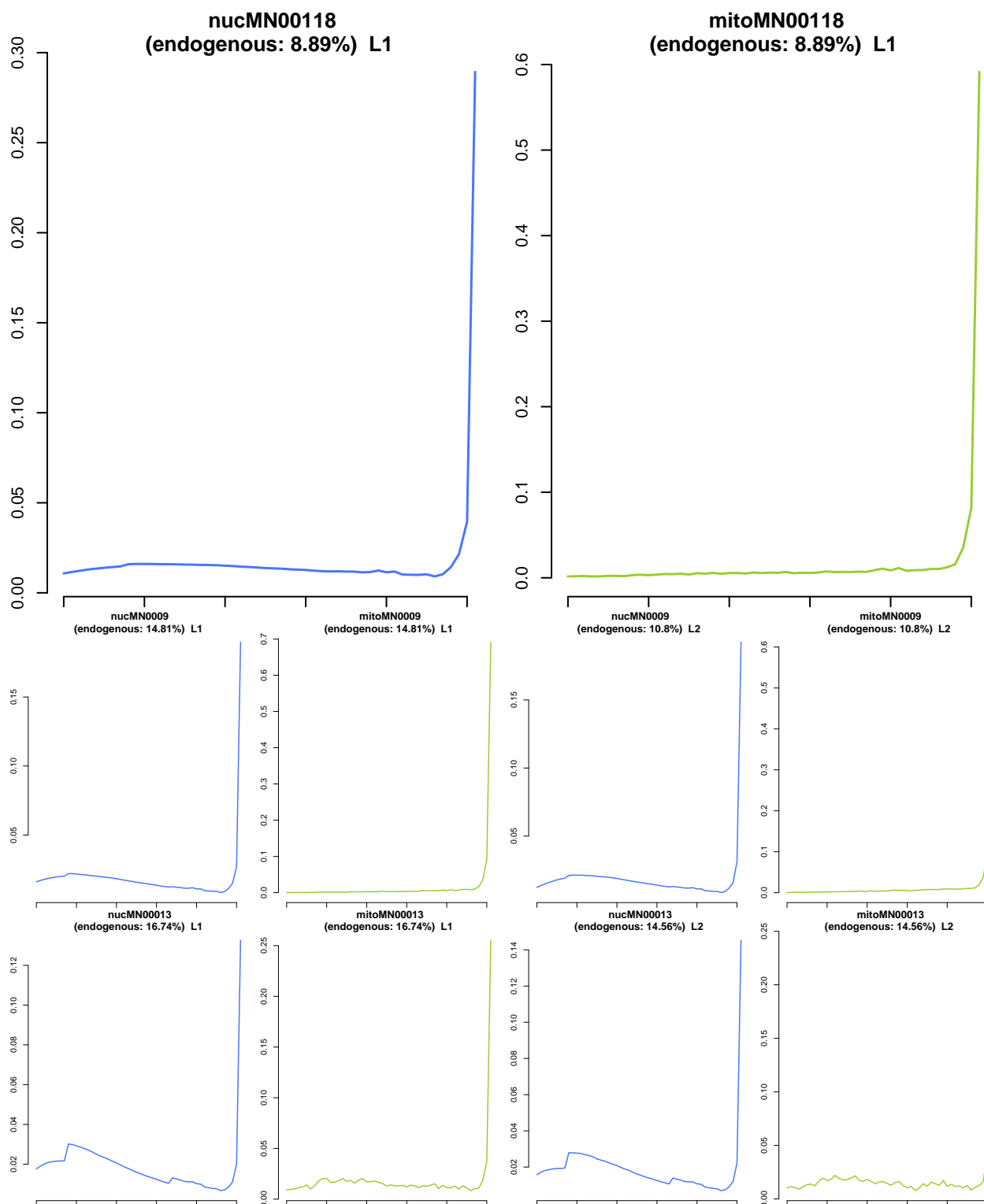### (endogenous: 1.74%) L1

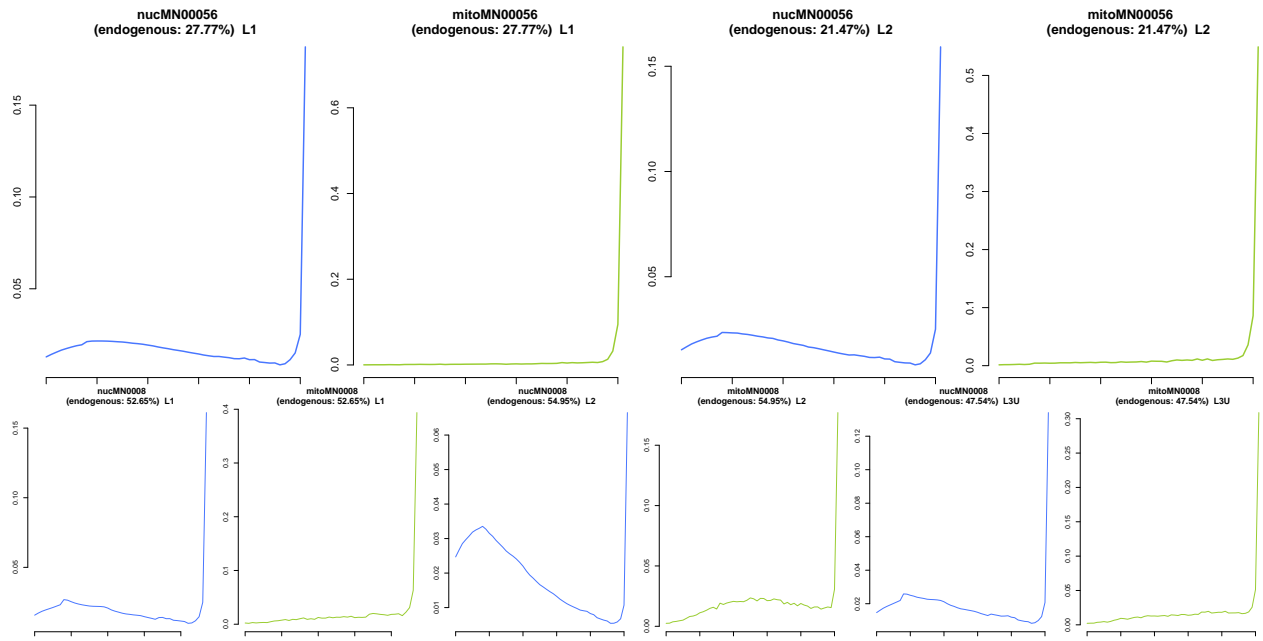## mitoMN00022
### (endogenous: 1.74%) L1

**nucMN00045**
**(endogenous: 2.03%)  L1**

**mitoMN00045**
**(endogenous: 2.03%)  L1**

**nucMN00068**
**(endogenous: 2.91%)  L1**

**mitoMN00068**
**(endogenous: 2.91%)  L1**

**nucMN00346**
**(endogenous: 6.71%) L1**

**mitoMN00346**
**(endogenous: 6.71%) L1**

**nucMN00346**
**(endogenous: 3.85%) L2**

**mitoMN00346**
**(endogenous: 3.85%) L2**

# nucMN0003
# (endogenous: 3.92%) L1

# mitoMN0003
# (endogenous: 3.92%) L1

**nucMN00021**
**(endogenous: 8.3%) L1**

**mitoMN00021**
**(endogenous: 8.3%) L1**

**nucMN00021**
**(endogenous: 4.57%) L2**

**mitoMN00021**
**(endogenous: 4.57%) L2**

nucMN00067
(endogenous: 5.74%) L1

mitoMN00067
(endogenous: 5.74%) L1

nucMN00067
(endogenous: 4.97%) L2

mitoMN00067
(endogenous: 4.97%) L2

nucMN00119
(endogenous: 8.96%) L1

mitoMN00119
(endogenous: 8.96%) L1

nucMN00119
(endogenous: 6.49%) L2

mitoMN00119
(endogenous: 6.49%) L2

# nucMN00066
# (endogenous: 8.05%) L1

# mitoMN00066
# (endogenous: 8.05%) L1

## nucMN00118
## (endogenous: 8.89%) L1

## mitoMN00118
## (endogenous: 8.89%) L1

**nucMN0009**
**(endogenous: 14.81%) L1**

**mitoMN0009**
**(endogenous: 14.81%) L1**

**nucMN0009**
**(endogenous: 10.8%) L2**

**mitoMN0009**
**(endogenous: 10.8%) L2**

**nucMN00013**
**(endogenous: 16.74%) L1**

**mitoMN00013**
**(endogenous: 16.74%) L1**

**nucMN00013**
**(endogenous: 14.56%) L2**

**mitoMN00013**
**(endogenous: 14.56%) L2**

## Molecular damage

# MN00010, L1



# MN00039, L1



**MN00316, L1**

**MN00316, L2**



**MN1943, L1**

**MN1943, mtCapture**

## MN00069, L1



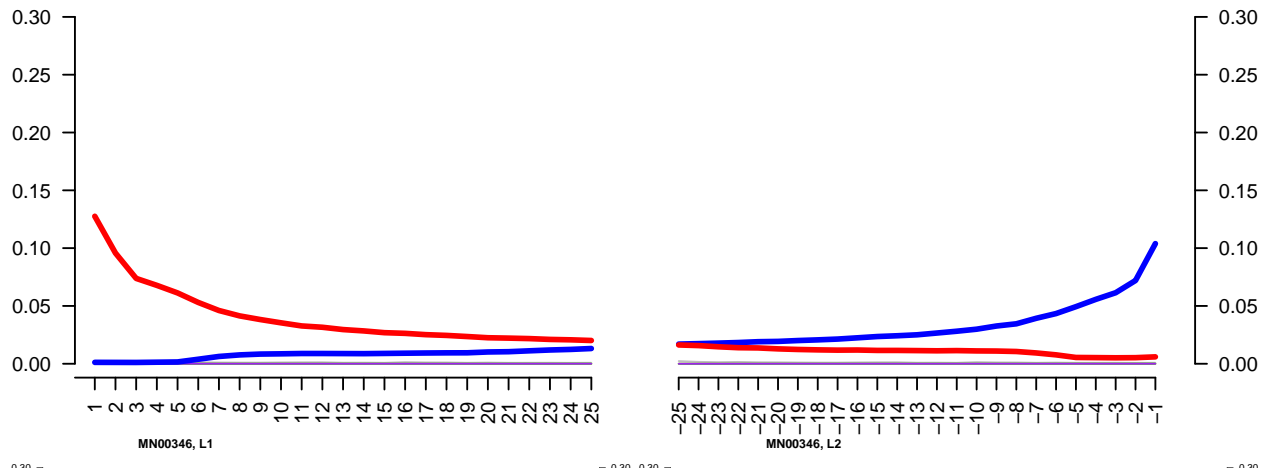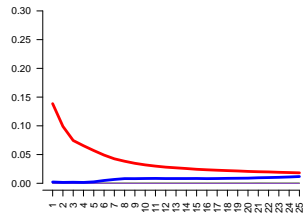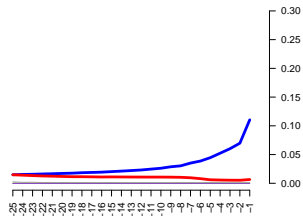## MN00022, L1



## MN00045, L1



15

## MN00068, L1



## MN00016, L1



## MN00064, L1



MN00346, L1

MN00346, L2

## MN0003, L1



### MN00021, L1



### MN00021, L2



### MN00067, L1



### MN00067, L2



### MN00119, L1



### MN00119, L2



## MN00066, L1

**MN00118, L1**

# Genome coverage



Genome coverage



Genome coverage (log-scale)

# Clonality



Clonality

# Sex determination

```
## [1] 120
```

Sex determination

# Contamination estimates

Work in progress

## Mitochondrial

Text on plots:

Mitochondrial coverage,

number of reads (all polymorphisms), number of reads (removing transitions)



Contamination estimates (mitochondrial data)

# X

Work in progress



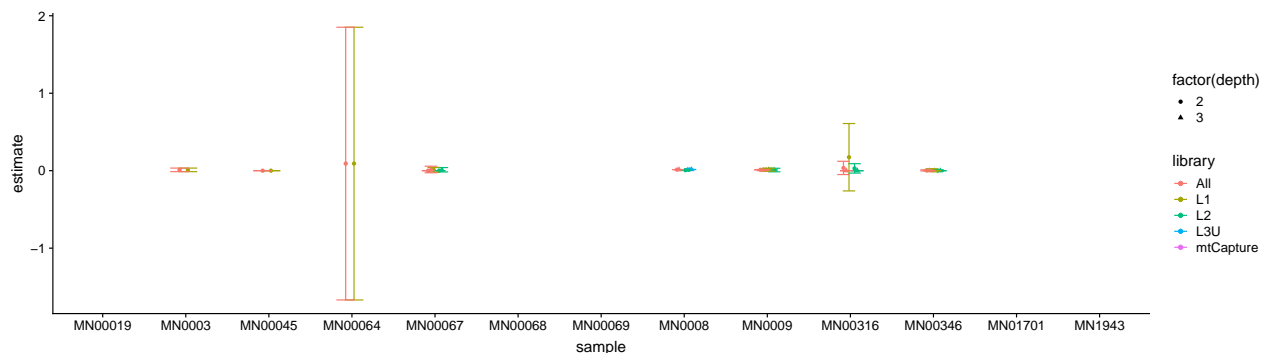| sample | library | mito_estimate | mito_lower | mito_upper | X_estimate | X_low | X_upper | mito_num | mito_coverage | mean_X_sites | min_depth | sex |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| MN00019 | A | 6% | 2% | 26% | NA% | NA% | NA% | 36 | 0.1240 | NA | 3 | XY |
| MN00019 | L1 | 4% | 0% | 66% | NA% | NA% | NA% | 7 | 0.0213 | NA | 3 | consistent with XY but not XX |
| MN00019 | L2 | 2% | 0% | 26% | NA% | NA% | NA% | 28 | 0.1030 | NA | 3 | XY |
| MN0003 | A | 2% | 0% | 4% | NA% | NA% | NA% | 7954 | 35.6000 | NA | 3 | XY |
| MN0003 | L1 | 2% | 0% | 4% | NA% | NA% | NA% | 7953 | 35.6000 | NA | 3 | XY |
| MN00045 | A | 2% | 0% | 2% | NA% | NA% | NA% | 10926 | 55.4000 | NA | 3 | XY |
| MN00045 | L1 | 2% | 0% | 2% | NA% | NA% | NA% | 10926 | 55.4000 | NA | 3 | XY |
| MN00064 | A | 2% | 2% | 4% | NA% | NA% | NA% | 11327 | 56.5000 | NA | 3 | XY |
| MN00064 | L1 | 2% | 2% | 4% | NA% | NA% | NA% | 11327 | 56.5000 | NA | 3 | XY |
| MN00067 | A | 4% | 2% | 6% | 2% | -2% | 6% | 21590 | 100.0000 | 82 | 3 | XY |
| MN00067 | L1 | 4% | 2% | 6% | NA% | NA% | NA% | 8362 | 37.8000 | NA | 3 | XY |
| MN00067 | L2 | 2% | 2% | 4% | 2% | -2% | 4% | 13228 | 62.5000 | 29 | 3 | XY |
| MN00068 | A | 4% | 2% | 10% | NA% | NA% | NA% | 2466 | 11.3000 | NA | 3 | XY |
| MN00068 | L1 | 4% | 2% | 8% | NA% | NA% | NA% | 2462 | 11.3000 | NA | 3 | XY |
| MN00069 | A | 12% | 8% | 20% | NA% | NA% | NA% | 1192 | 5.3400 | NA | 3 | XY |
| MN00069 | L1 | 12% | 8% | 18% | NA% | NA% | NA% | 1192 | 5.3400 | NA | 3 | XY |
| MN0008 | A | 0% | 0% | 2% | 2% | 2% | 2% | 47938 | 222.0000 | 47991 | 3 | XY |
| MN0008 | L1 | 0% | 0% | 2% | NA% | NA% | NA% | 16319 | 77.3000 | NA | 3 | XY |
| MN0008 | L2 | 0% | 0% | 0% | 0% | 0% | 0% | 15672 | 69.0000 | 12347 | 3 | XY |
| MN0008 | L3U | 2% | 0% | 2% | 2% | 2% | 2% | 15950 | 76.1000 | 10001 | 3 | XY |
| MN0009 | A | 4% | 4% | 6% | 0% | 0% | 2% | 37239 | 205.0000 | 1417 | 3 | XY |
| MN0009 | L1 | 2% | 2% | 4% | 2% | 0% | 4% | 18871 | 105.0000 | 314 | 3 | XY |
| MN0009 | L2 | 2% | 2% | 4% | 0% | -2% | 4% | 18368 | 100.0000 | 162 | 3 | XY |
| MN00316 | A | 2% | 0% | 2% | 0% | 0% | 0% | 26013 | 134.0000 | 1 | 3 | XY |
| MN00316 | L1 | 0% | 0% | 2% | NA% | NA% | NA% | 12198 | 63.9000 | NA | 3 | XY |
| MN00316 | L2 | 2% | 0% | 2% | 0% | 0% | 0% | 13820 | 70.5000 | 1 | 3 | XY |
| MN00346 | A | 2% | 0% | 2% | 0% | 0% | 0% | 29390 | 149.0000 | 97 | 3 | XY |
| MN00346 | L1 | 0% | 0% | 2% | 0% | 0% | 0% | 15392 | 78.6000 | 10 | 3 | XY |
| MN00346 | L2 | 2% | 0% | 2% | 0% | 0% | 0% | 14000 | 70.2000 | 20 | 3 | XY |
| MN01701 | A | 20% | 14% | 28% | NA% | NA% | NA% | 317 | 1.0700 | NA | 3 | consistent with XY but not XX |
| MN01701 | mtCapture | NA% | NA% | NA% | NA% | NA% | NA% | 1199 | 8.2400 | NA | 3 | consistent with XY but not XX |

| sample | library | mito_estimate | mito_low | mito_upper | Xer_estimate | low_Xest | upper_Xest | mito_num | num_Reads | average_X | sites | min_sex depth |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| MN19A3 | l | 20% | 14% | 30% | NA% | NA% | NA% | 11402 | 40.8000 | NA | 3 | consistent with XY but not XX |
| MN19L3 | l | 22% | 14% | 32% | NA% | NA% | NA% | 341 | 1.0100 | NA | 3 | consistent with XY but not XX |

**Supplementary table 1**