# SUPPLEMENTARY INFORMATION

**FastQC reports**

(Separate FastQC reports for read_1 and read_2; first 16 bp of read_1 are the "droplet" barcodes in linked-read sequencing.)

# FastQC Report

## Summary

✅ [Basic Statistics](#)

✅ [Per base sequence quality](#)

✅ [Per tile sequence quality](#)

✅ [Per sequence quality scores](#)

⚠️ [Per base sequence content](#)

⚠️ [Per sequence GC content](#)

✅ [Per base N content](#)

✅ [Sequence Length Distribution](#)

⚠️ [Sequence Duplication Levels](#)

✅ [Overrepresented sequences](#)
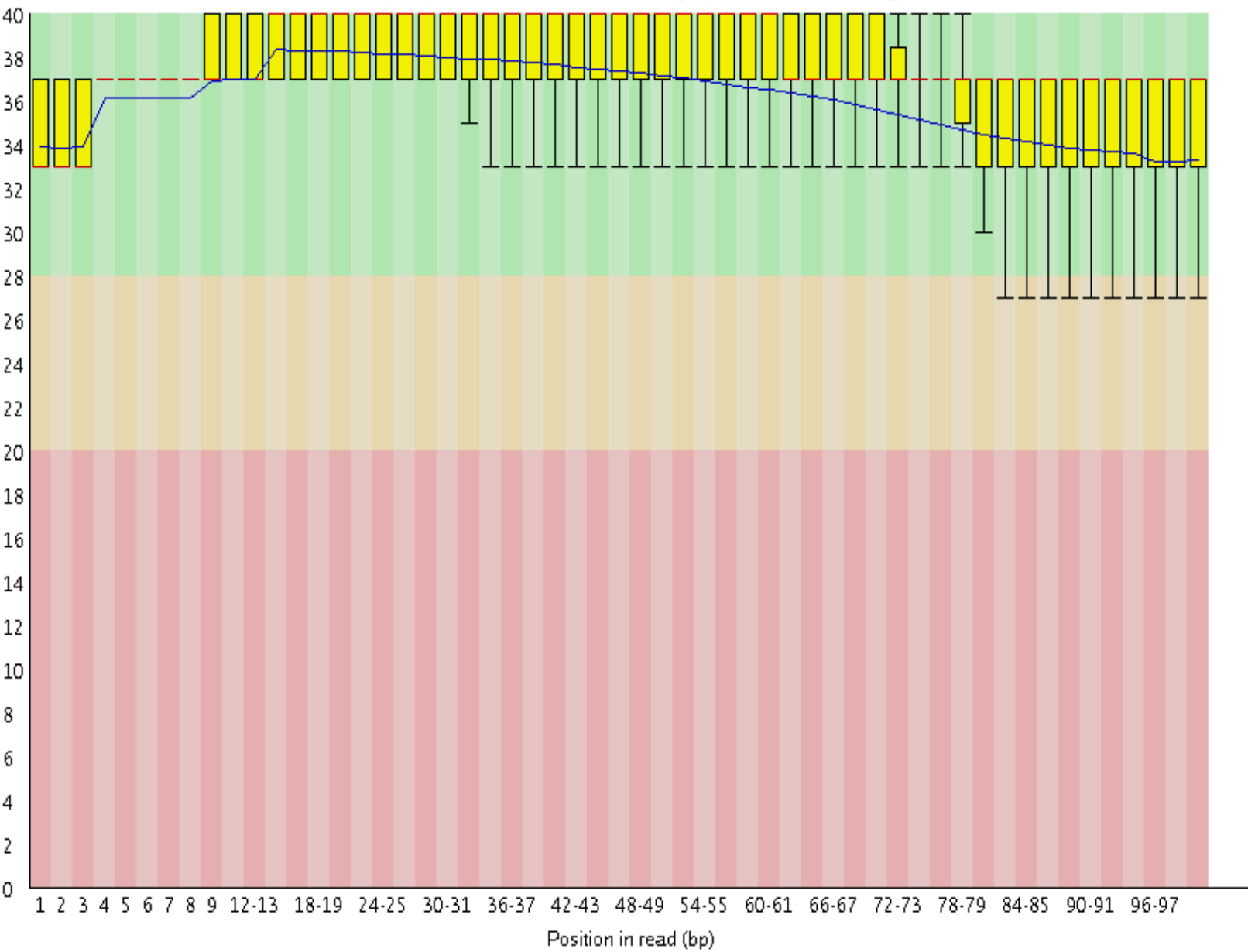
✅ [Adapter Content](#)

## ✅ Basic Statistics

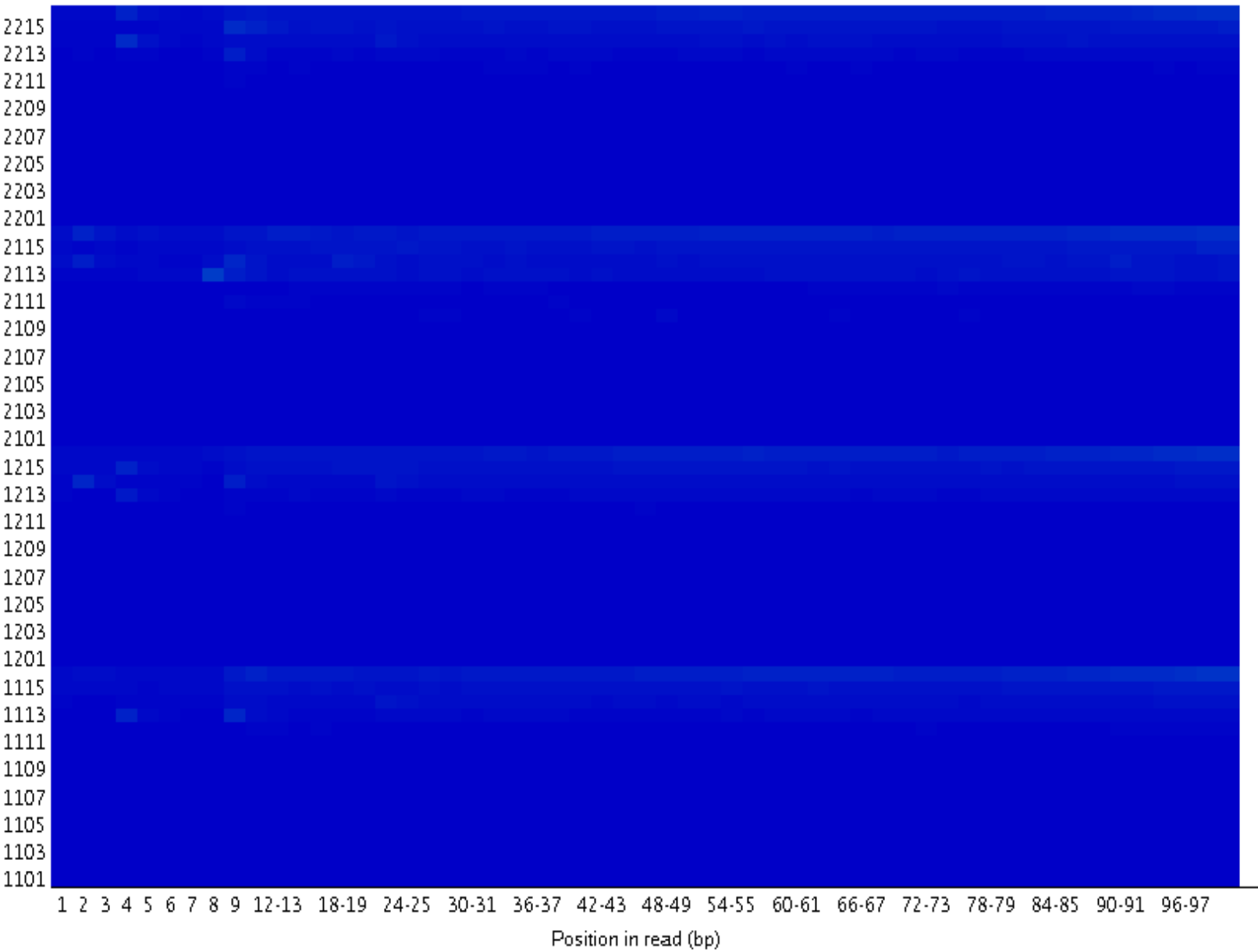| Measure | Value |
|---|---|
| Filename | stdin |
| File type | Conventional base calls |
| Encoding | Sanger / Illumina 1.9 |
| Total Sequences | 3214623356 |
| Sequences flagged as poor quality | 0 |
| Sequence length | 100 |
| %GC | 40 |

## ✅ Per base sequence quality

Quality scores across all bases (Sanger / Illumina 1.9 encoding)



Position in read (bp)

✅ **Per tile sequence quality**

Quality per tile

Position in read (bp)

# Per sequence quality scores

Quality score distribution over all sequences

Average Quality per read

Mean Sequence Quality (Phred Score)
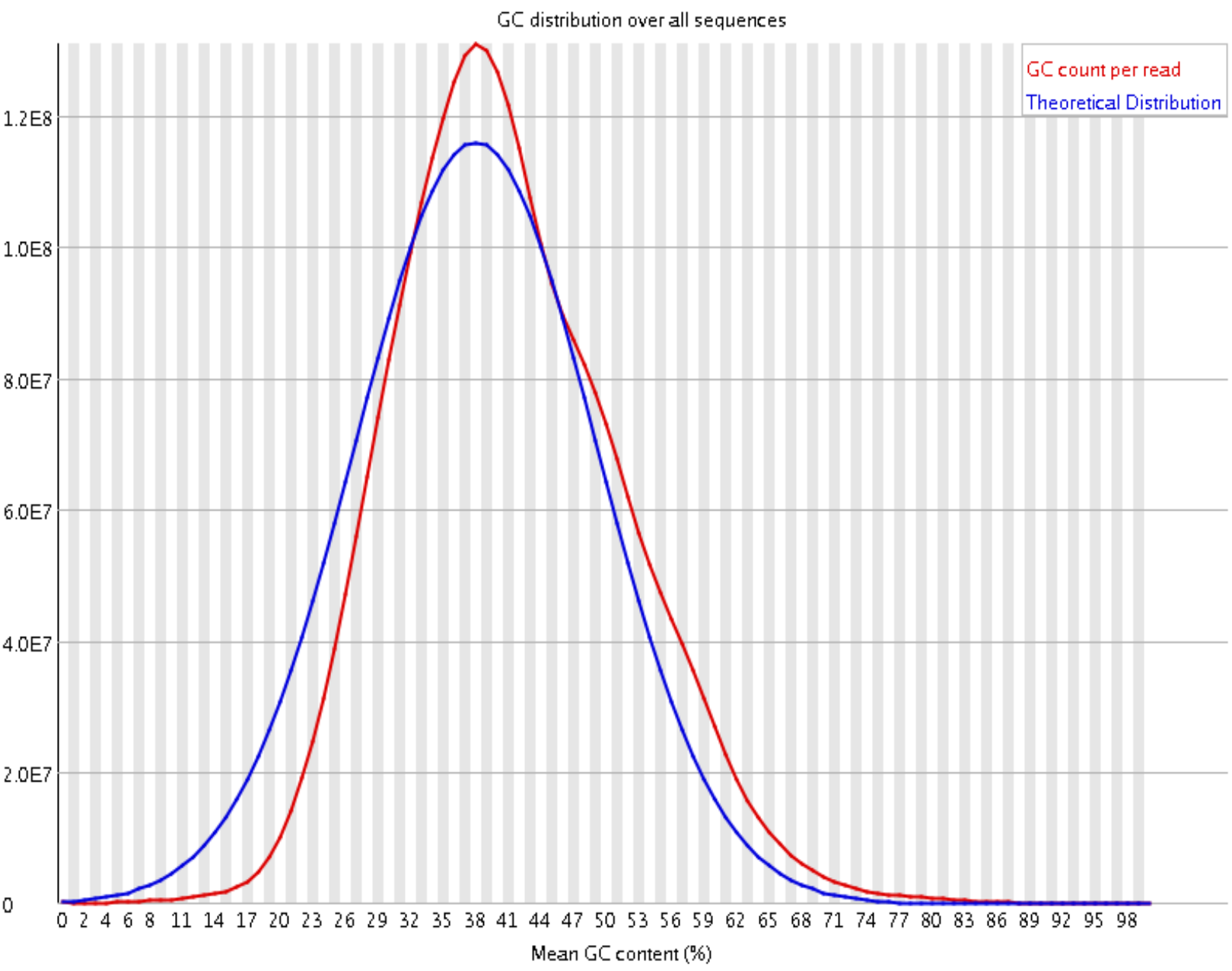
⚠️**Per base sequence content**

Sequence content across all bases

## Per sequence GC content

GC distribution over all sequences

**Per base N content**

N content across all bases

%N

Position in read (bp)

## ✅ Sequence Length Distribution

Distribution of sequence lengths over all sequences

Sequence Length (bp)

## ⚠ Sequence Duplication Levels

Percent of seqs remaining if deduplicated 67.74%

![Image: Line graph showing Sequence Duplication Level]

Legend:
% Deduplicated sequences
% Total sequences

X-axis: Sequence Duplication Level (1, 2, 3, 4, 5, 6, 7, 8, 9, >10, >50, >100, >500, >1k, >5k, >10k)

## ✅ Overrepresented sequences

No overrepresented sequences

## ✅ Adapter Content

% Adapter

Illumina Universal Adapter
Illumina Small RNA 3' Adapter
Illumina Small RNA 5' Adapter
Nextera Transposase Sequence
SOLID Small RNA Adapter

Position in read (bp)
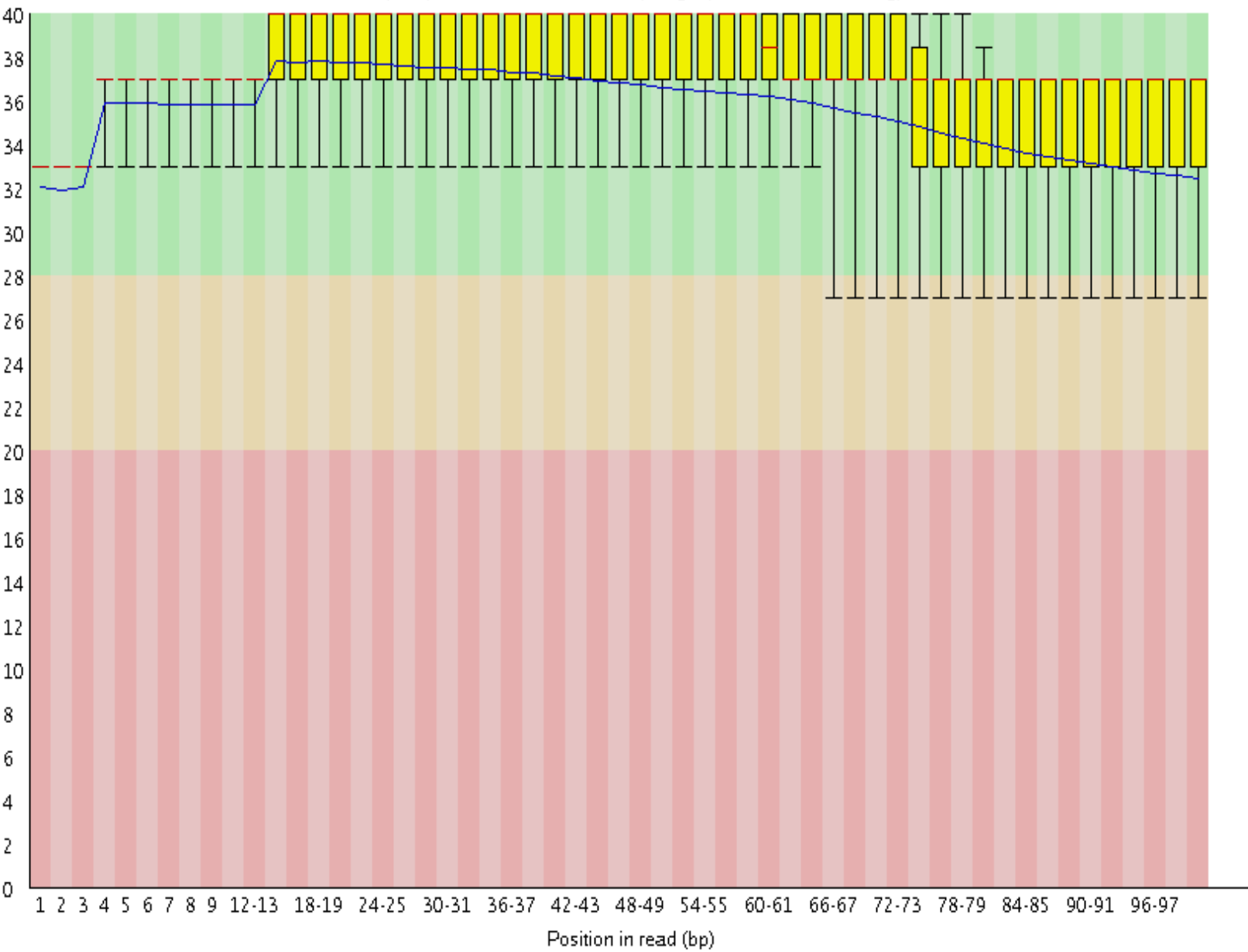
**Produced by [FastQC](FastQC) (version 0.11.7)**

# FastQC Report

## Summary

✅ Basic Statistics

✅ Per base sequence quality

✅ Per tile sequence quality

✅ Per sequence quality scores

❌ Per base sequence content

⚠️ Per sequence GC content

✅ Per base N content

✅ Sequence Length Distribution

✅ Sequence Duplication Levels

✅ Overrepresented sequences

✅ Adapter Content

## Basic Statistics

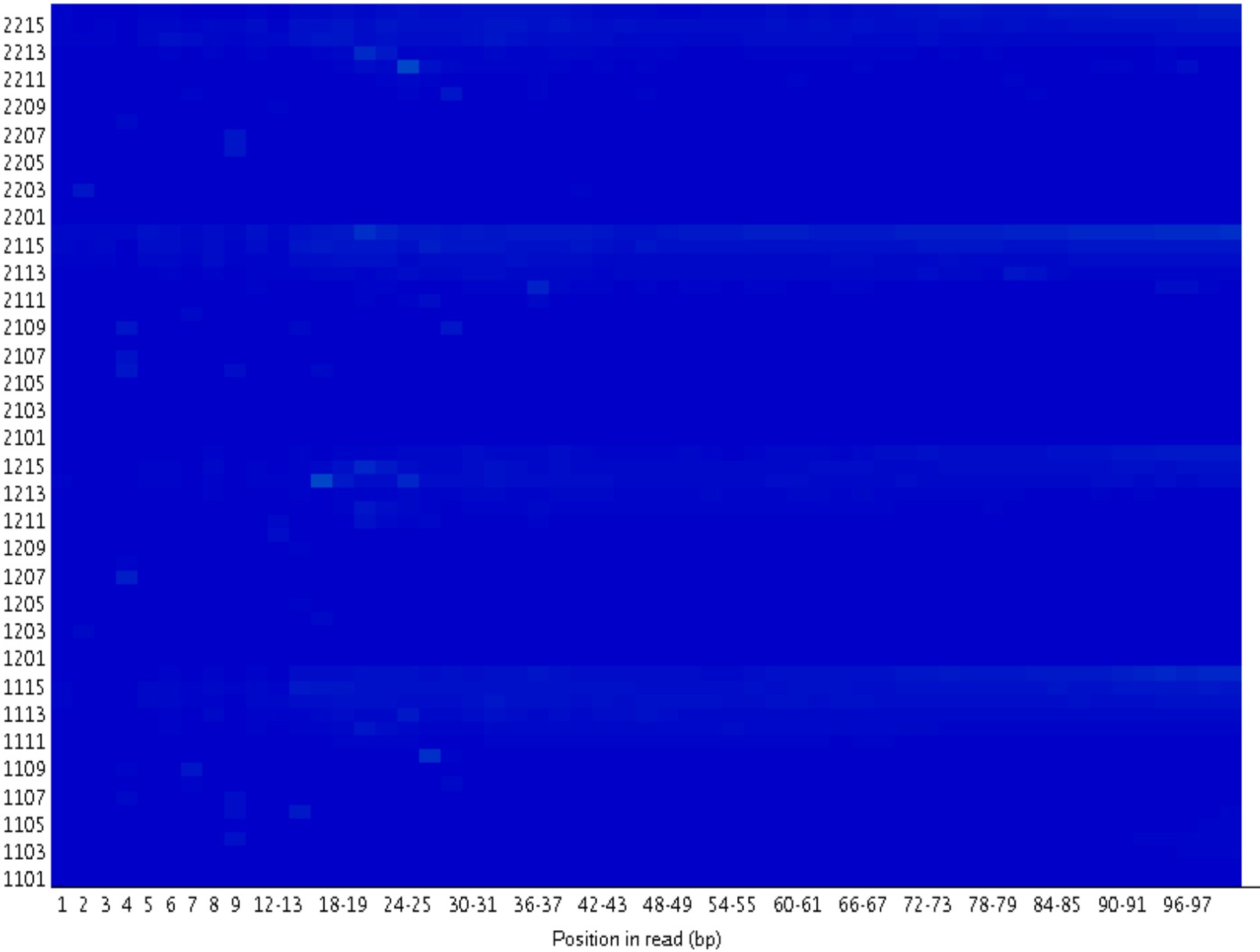| Measure | Value |
|---|---|
| Filename | stdin |
| File type | Conventional base calls |
| Encoding | Sanger / Illumina 1.9 |
| Total Sequences | 1280576580 |
| Sequences flagged as poor quality | 0 |
| Sequence length | 100 |
| %GC | 40 |

## Per base sequence quality
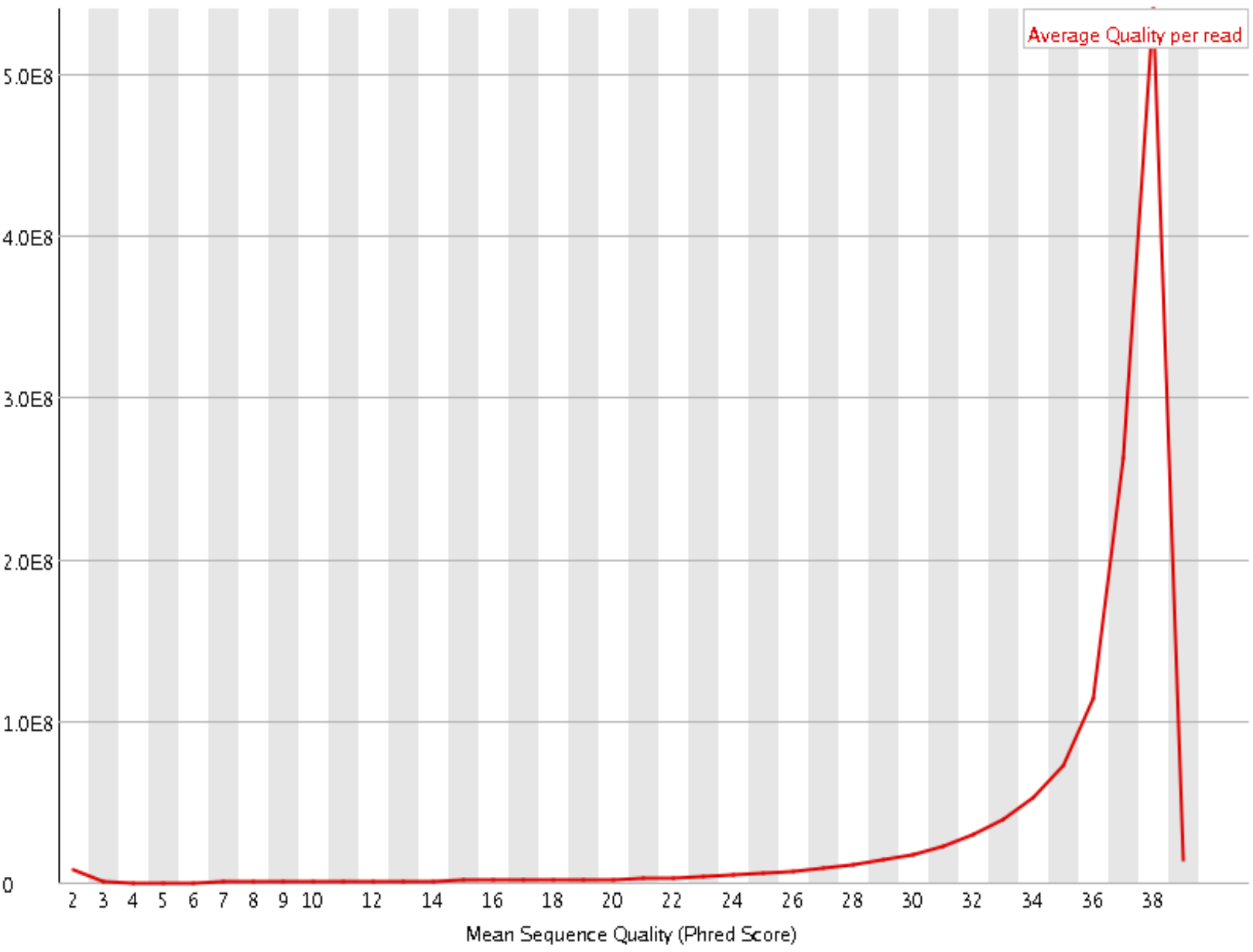
Quality scores across all bases (Sanger / Illumina 1.9 encoding)

Position in read (bp)

✅ **Per tile sequence quality**

Quality per tile



Position in read (bp)

✅ **Per sequence quality scores**

Quality score distribution over all sequences

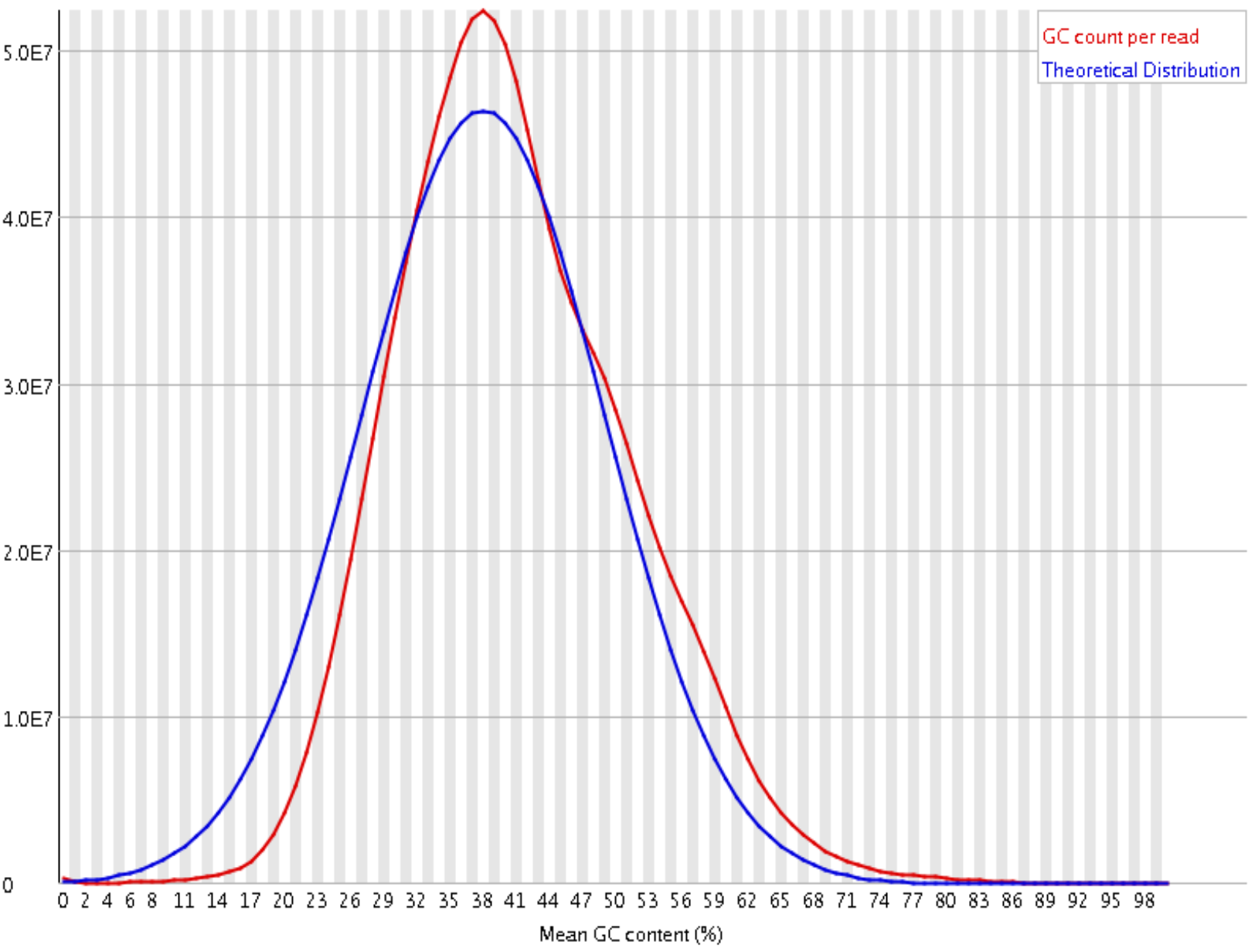Average Quality per read

Mean Sequence Quality (Phred Score)

❌ **Per base sequence content**

Sequence content across all bases

**Per sequence GC content**
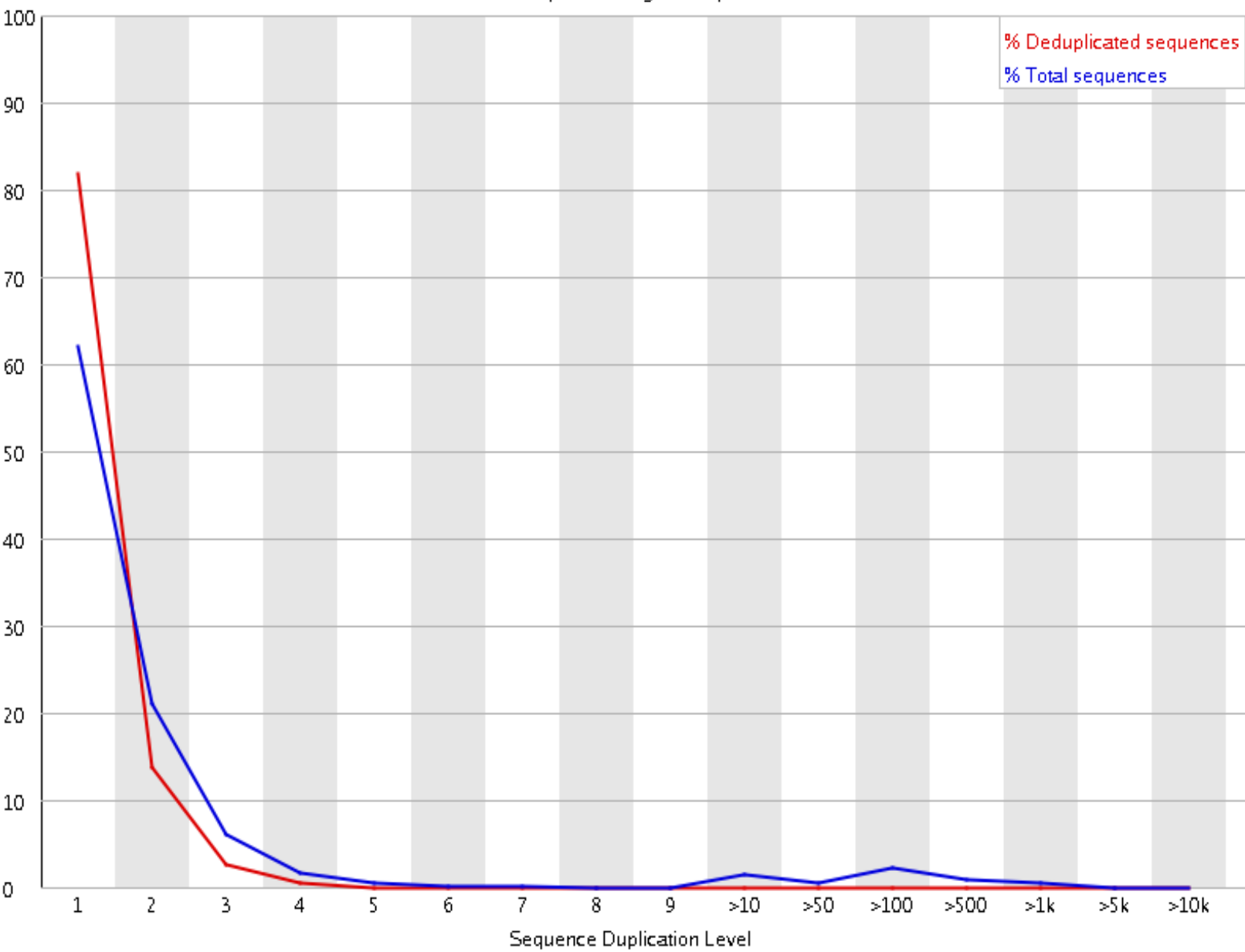
GC distribution over all sequences

**Per base N content**

N content across all bases

## Sequence Length Distribution

Distribution of sequence lengths over all sequences

## Sequence Duplication Levels

Percent of seqs remaining if deduplicated 75.85%

% Deduplicated sequences
% Total sequences

Sequence Duplication Level

✅ **Overrepresented sequences**

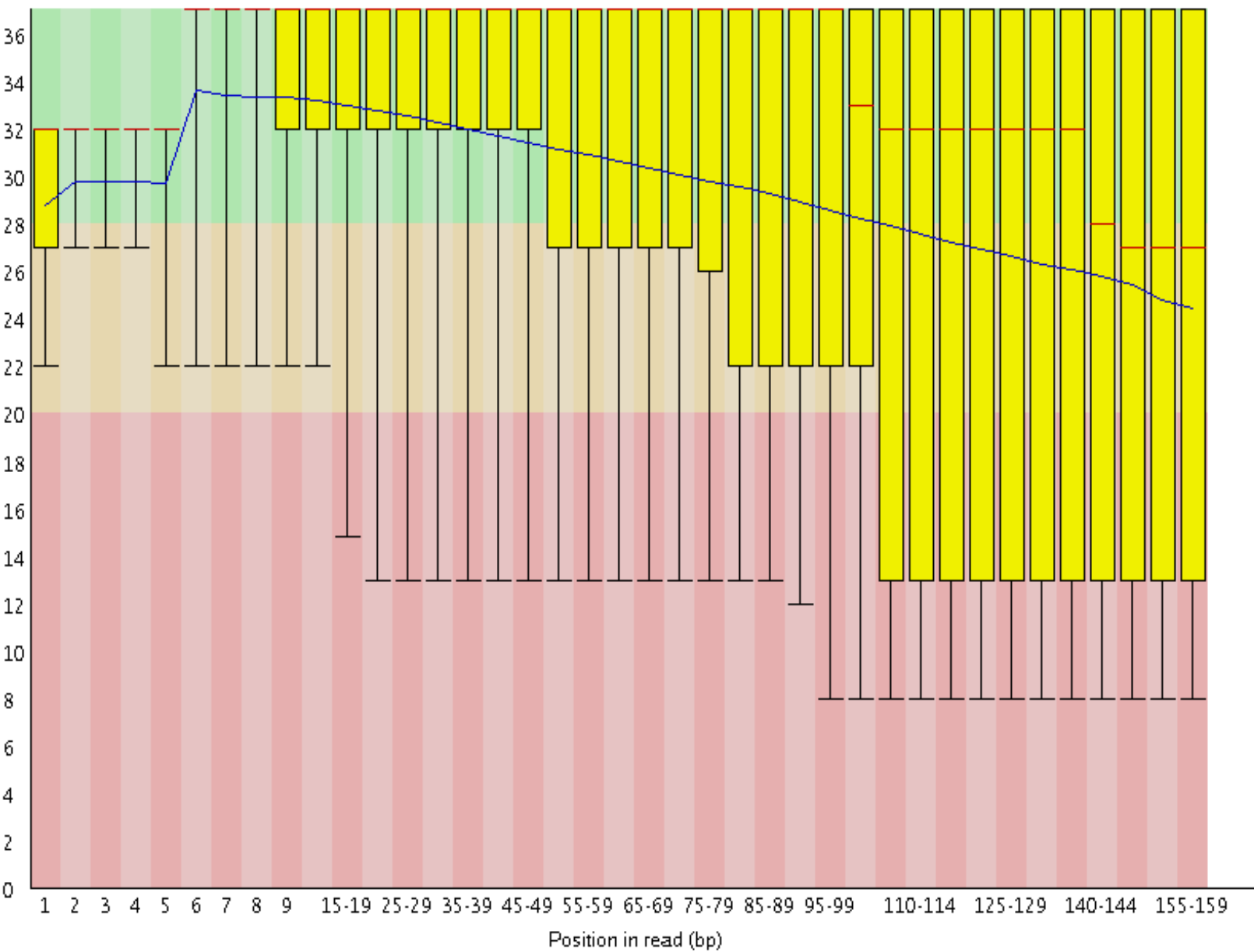No overrepresented sequences

✅ **Adapter Content**

# FastQC Report

## Summary

✅ Basic Statistics

✅ Per base sequence quality

❌ Per tile sequence quality

✅ Per sequence quality scores

⚠️ Per base sequence content

⚠️ Per sequence GC content

✅ Per base N content

⚠️ Sequence Length Distribution

⚠️ Sequence Duplication Levels

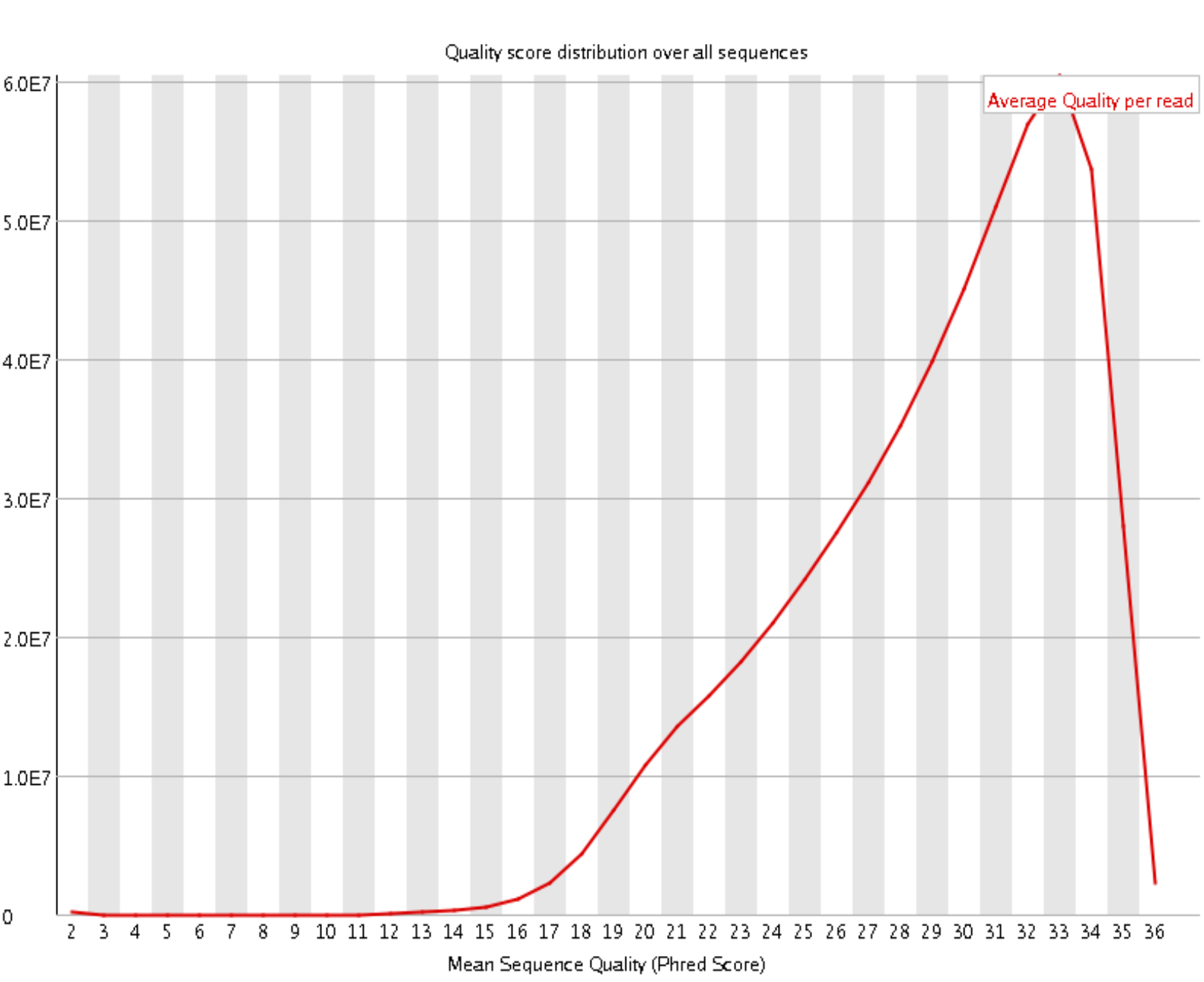⚠️ Overrepresented sequences

❌ Adapter Content

## ✅ Basic Statistics

| Measure | Value |
|---|---|
| Filename | stdin |
| File type | Conventional base calls |
| Encoding | Sanger / Illumina 1.9 |
| Total Sequences | 553277870 |
| Sequences flagged as poor quality | 0 |
| Sequence length | 35-159 |
| %GC | 41 |

## ✅ Per base sequence quality
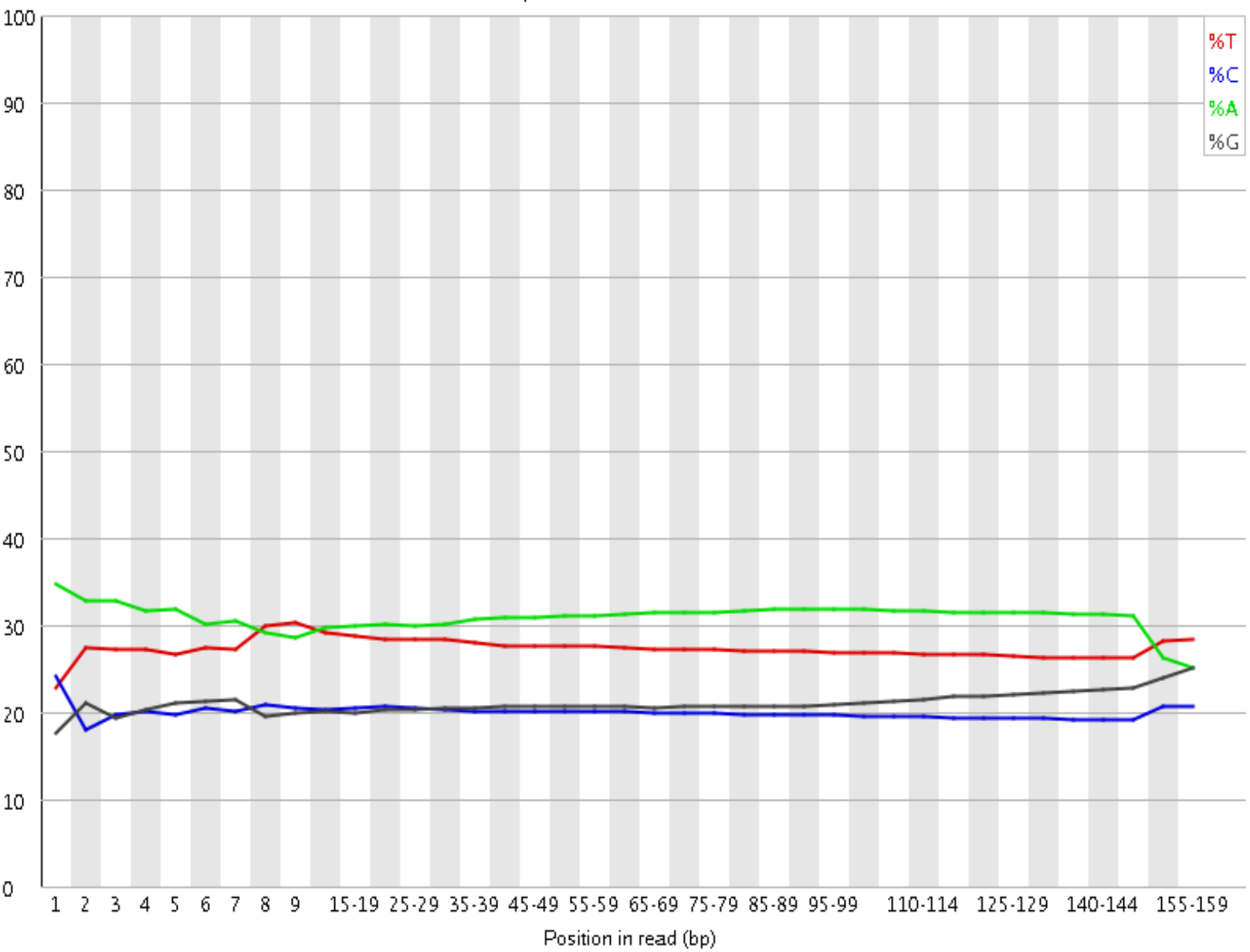
Quality scores across all bases (Sanger / Illumina 1.9 encoding)
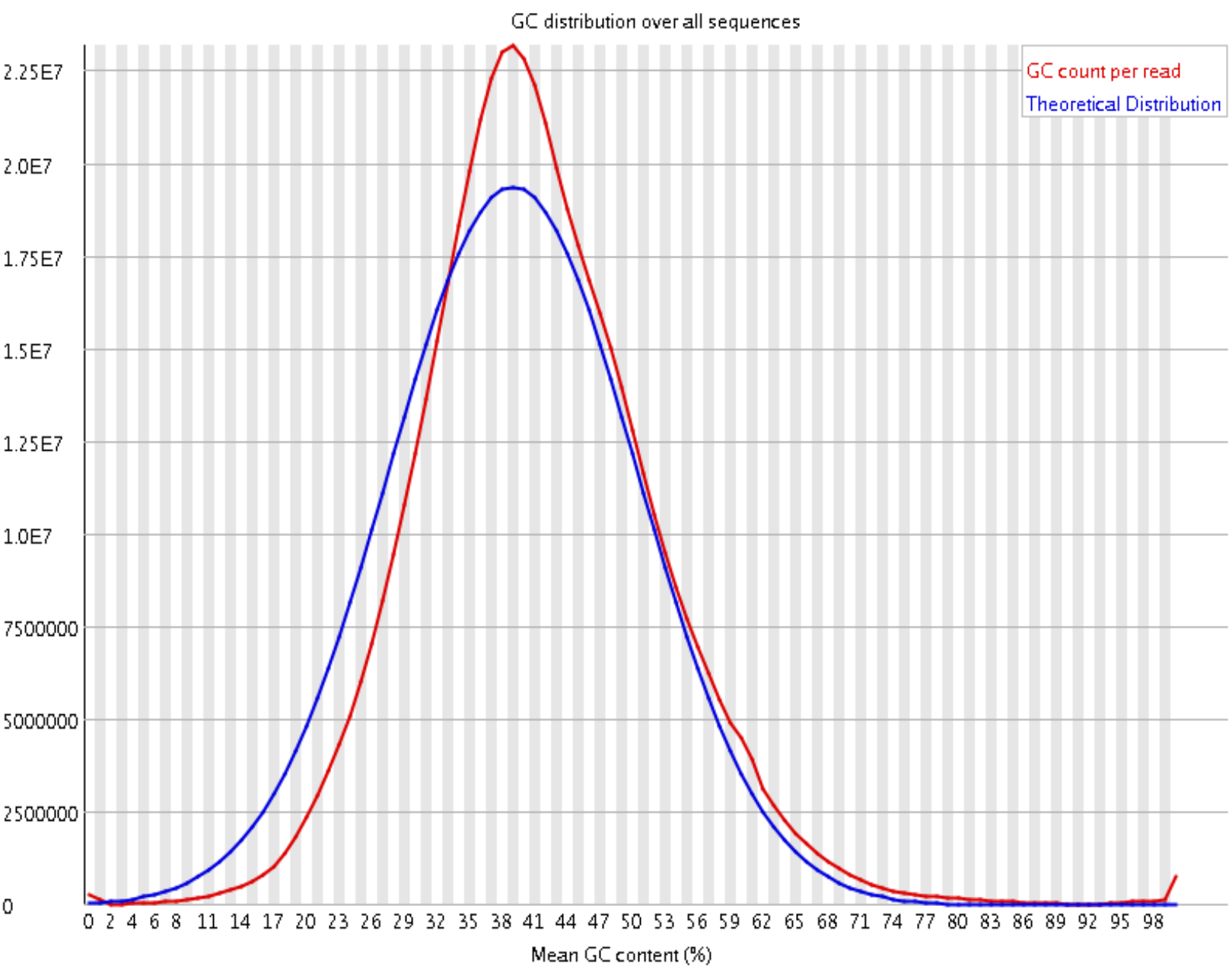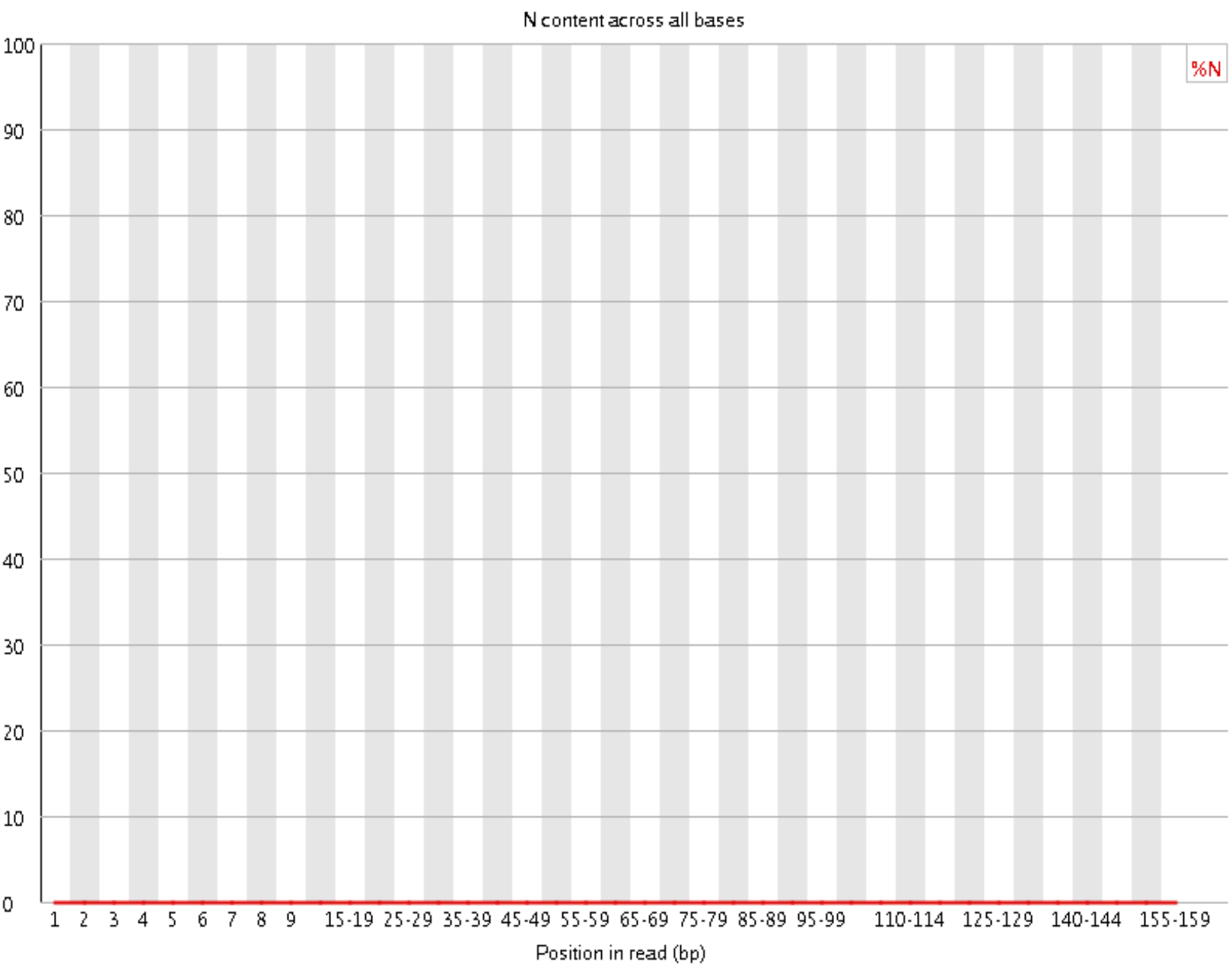
**❌ Per tile sequence quality**

Quality per tile

Position in read (bp)

![Green check icon] **Per sequence quality scores**

Quality score distribution over all sequences

Average Quality per read

Mean Sequence Quality (Phred Score)

⚠️**Per base sequence content**

Sequence content across all bases

## Per sequence GC content

GC distribution over all sequences
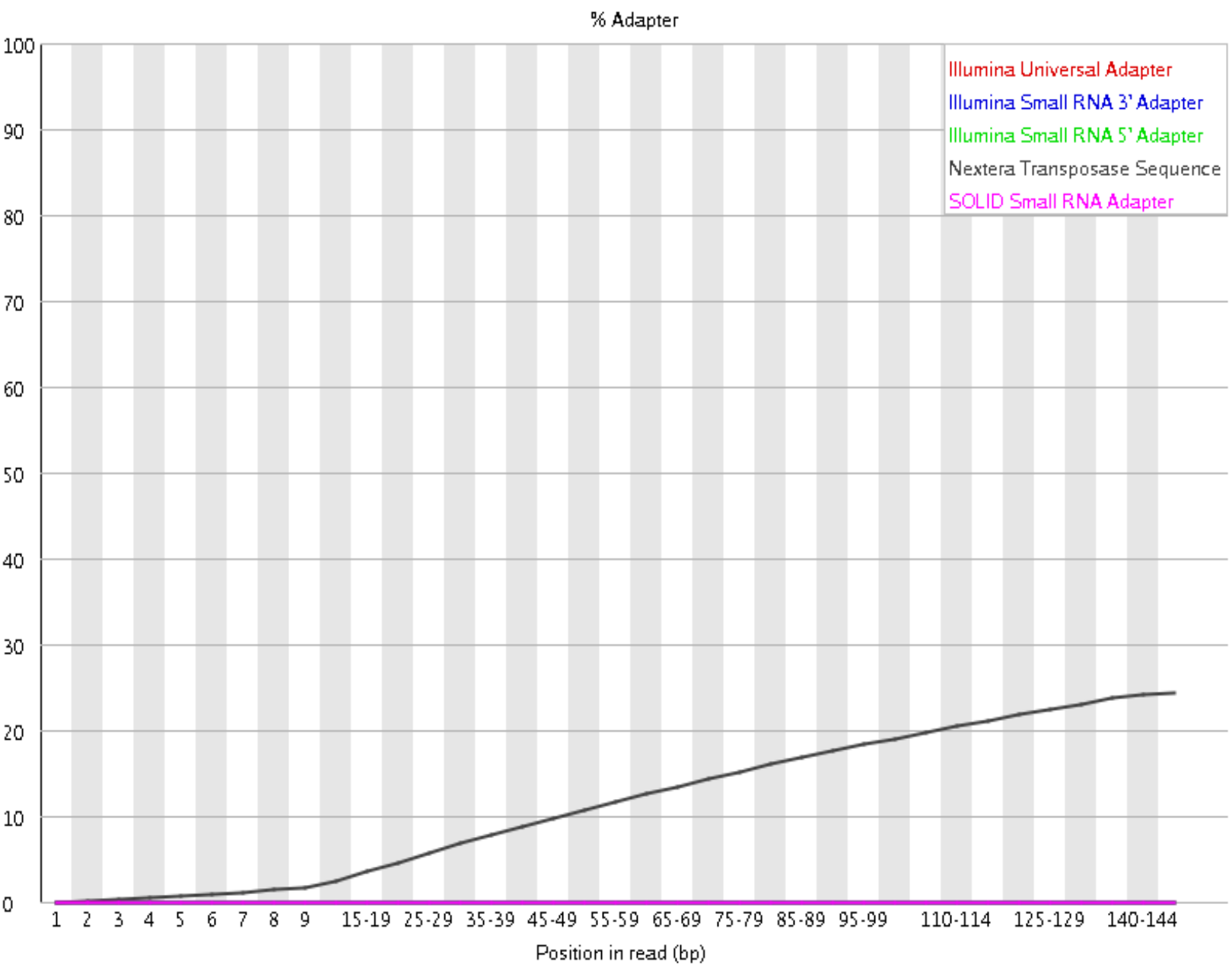
**Per base N content**

N content across all bases

Position in read (bp)

## Sequence Length Distribution

Distribution of sequence lengths over all sequences

## ⚠ Sequence Duplication Levels

Percent of seqs remaining if deduplicated 52.18%



% Deduplicated sequences
% Total sequences

Sequence Duplication Level

# ⚠️Overrepresented sequences

| Sequence | Count | Percentage | Possible Source |
|---|---|---|---|
| GGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGG | 1464007 | 0.2646061010898556 | No Hit |

# ❌Adapter Content

% Adapter

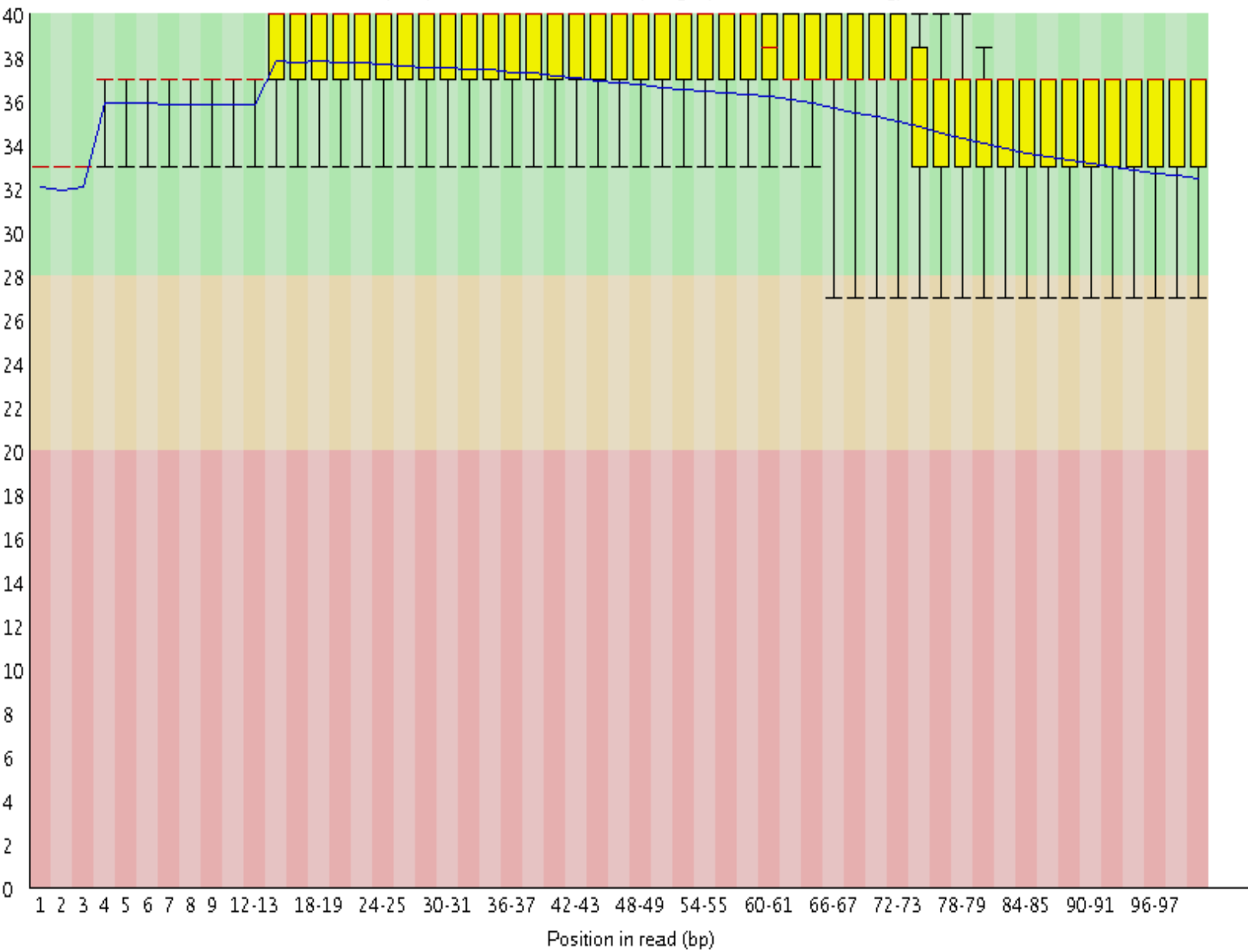# FastQC Report

## Summary

✅ [Basic Statistics](#)

✅ [Per base sequence quality](#)

✅ [Per tile sequence quality](#)

✅ [Per sequence quality scores](#)

❌ [Per base sequence content](#)

⚠️ [Per sequence GC content](#)

✅ [Per base N content](#)

✅ [Sequence Length Distribution](#)

✅ [Sequence Duplication Levels](#)

✅ [Overrepresented sequences](#)

✅ [Adapter Content](#)

## ✅ Basic Statistics

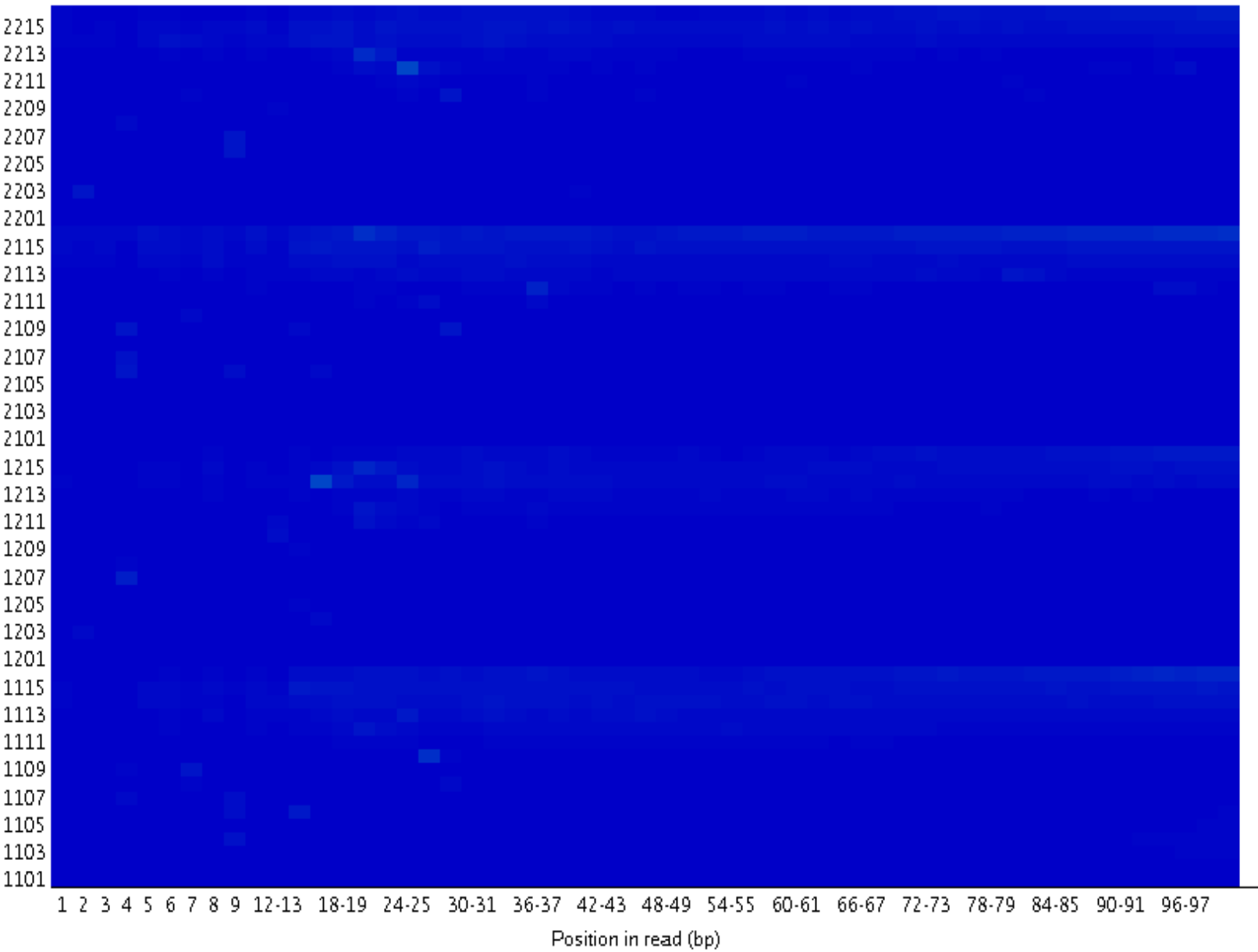| Measure | Value |
|---|---|
| Filename | stdin |
| File type | Conventional base calls |
| Encoding | Sanger / Illumina 1.9 |
| Total Sequences | 1280576580 |
| Sequences flagged as poor quality | 0 |
| Sequence length | 100 |
| %GC | 40 |

## ✅ Per base sequence quality

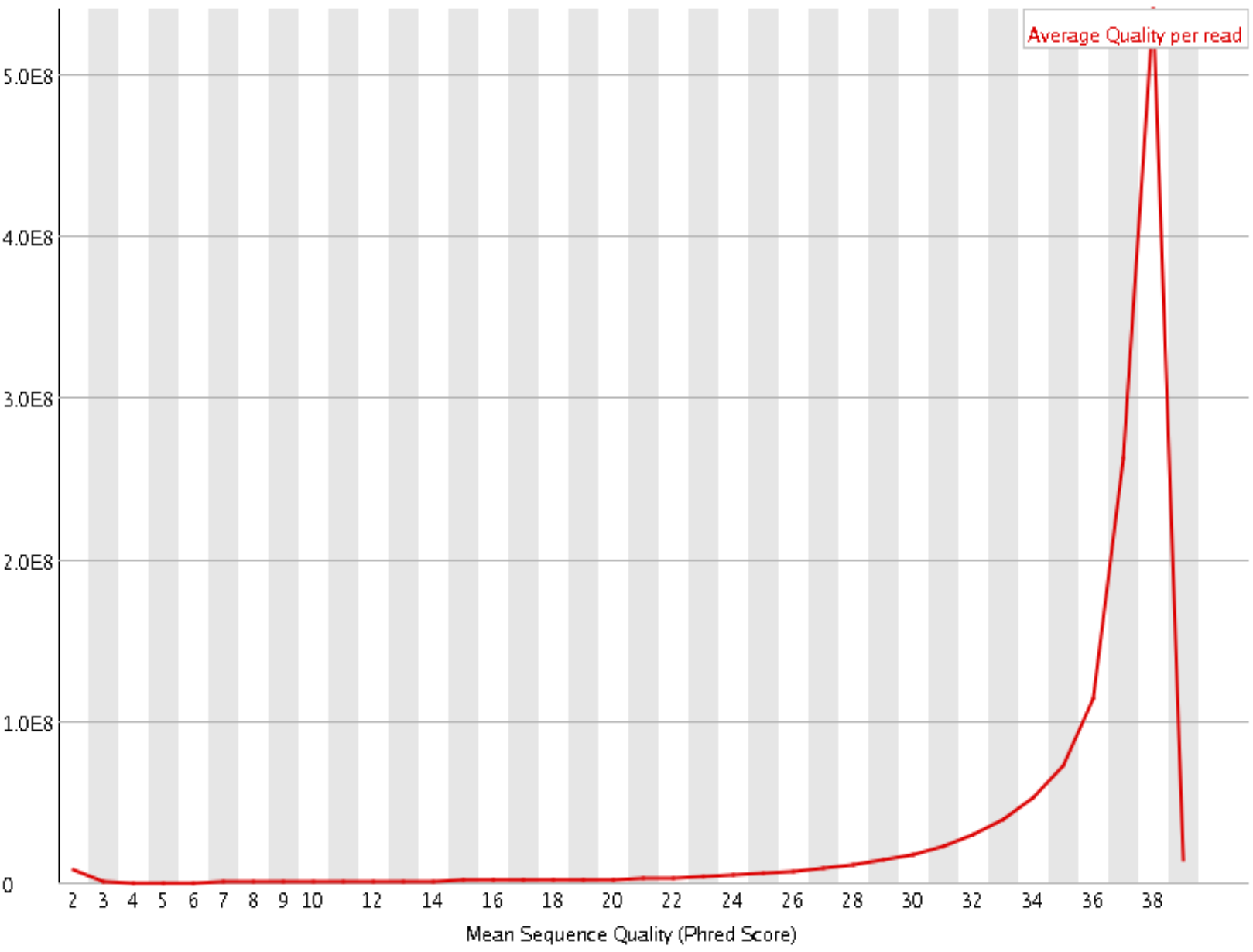Quality scores across all bases (Sanger / Illumina 1.9 encoding)

Position in read (bp)

![Green checkmark] **Per tile sequence quality**

Quality per tile

Position in read (bp)

## Per sequence quality scores

Quality score distribution over all sequences

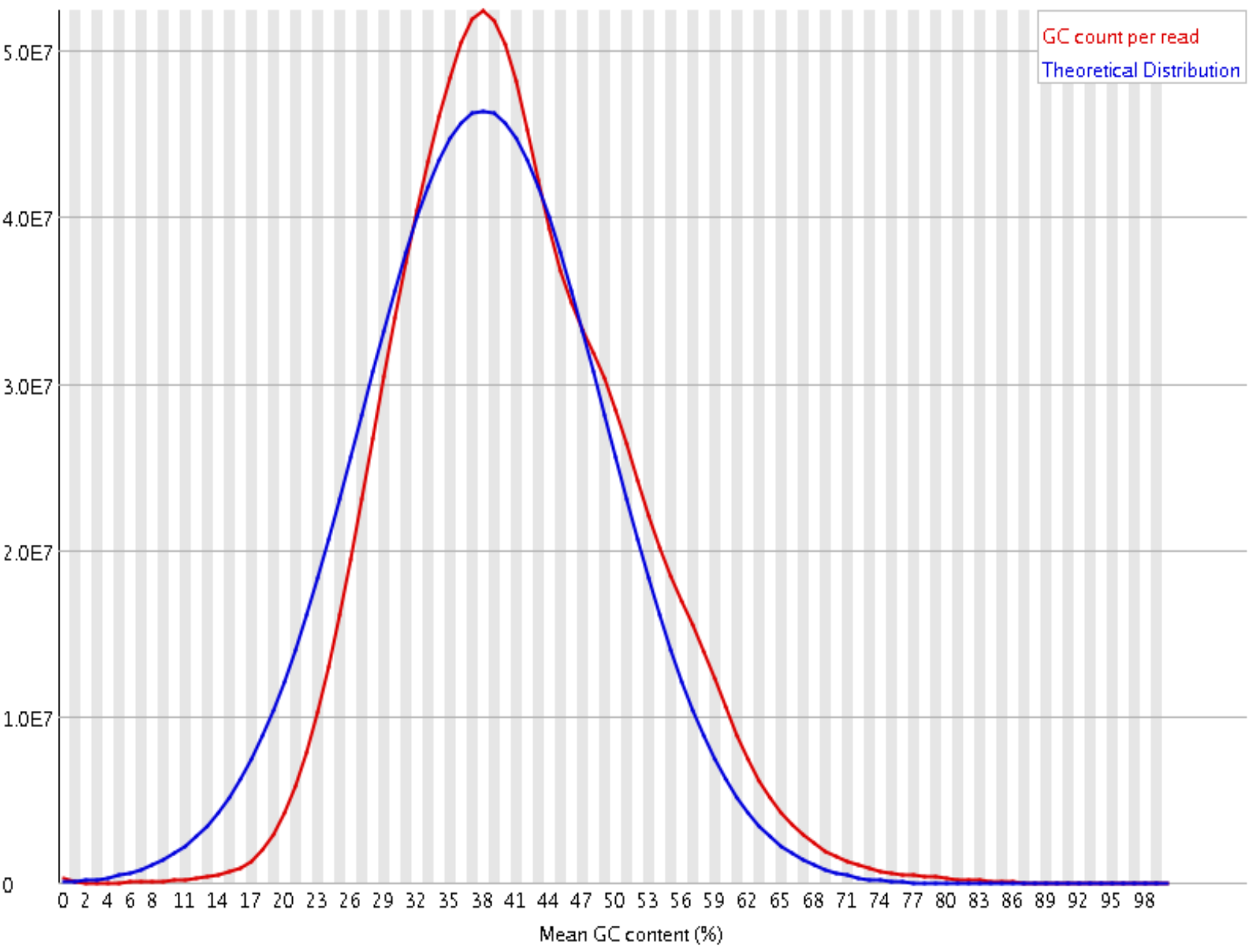Mean Sequence Quality (Phred Score)

Average Quality per read

## ❌ Per base sequence content

Sequence content across all bases

## Per sequence GC content

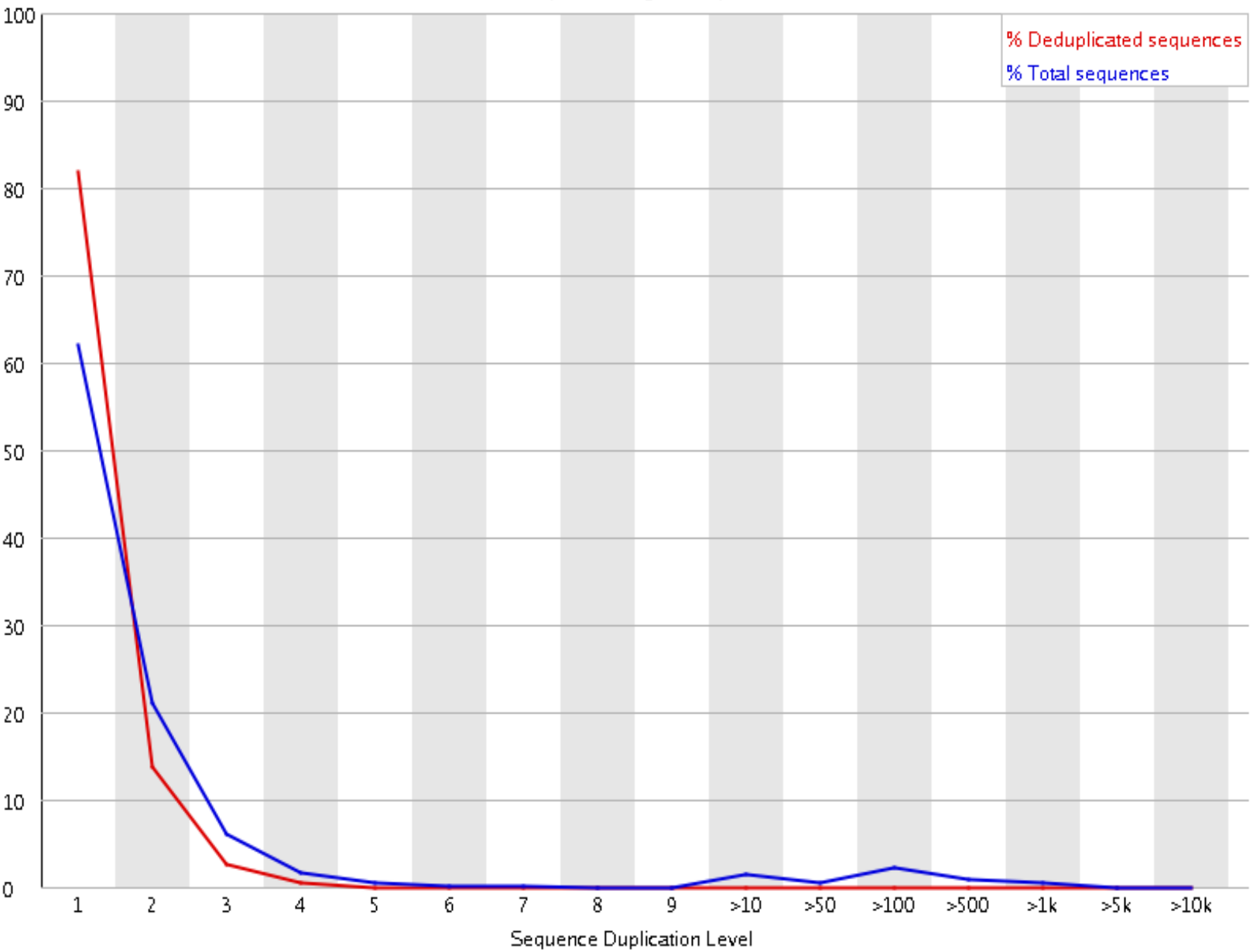GC distribution over all sequences

✅**Per base N content**

N content across all bases

## Sequence Length Distribution
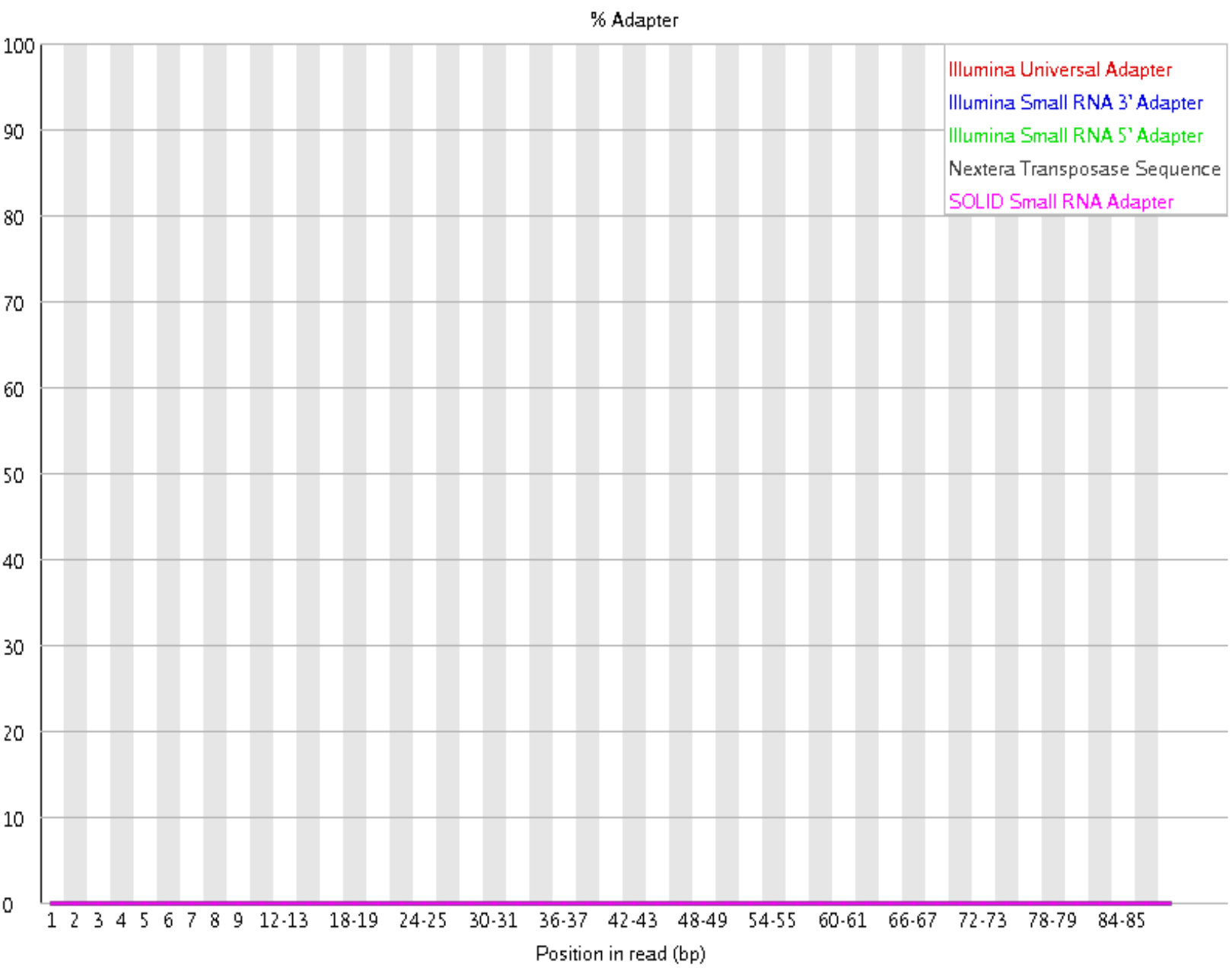
Distribution of sequence lengths over all sequences

## Sequence Duplication Levels

Percent of seqs remaining if deduplicated 75.85%



## ✅ Overrepresented sequences

No overrepresented sequences

## ✅ Adapter Content

% Adapter

Illumina Universal Adapter
Illumina Small RNA 3' Adapter
Illumina Small RNA 5' Adapter
Nextera Transposase Sequence
SOLID Small RNA Adapter

Position in read (bp)

**Produced by [FastQC](FastQC) (version 0.11.7)**

# FastQC Report

## Summary

✅ Basic Statistics

✅ Per base sequence quality

❌ Per tile sequence quality

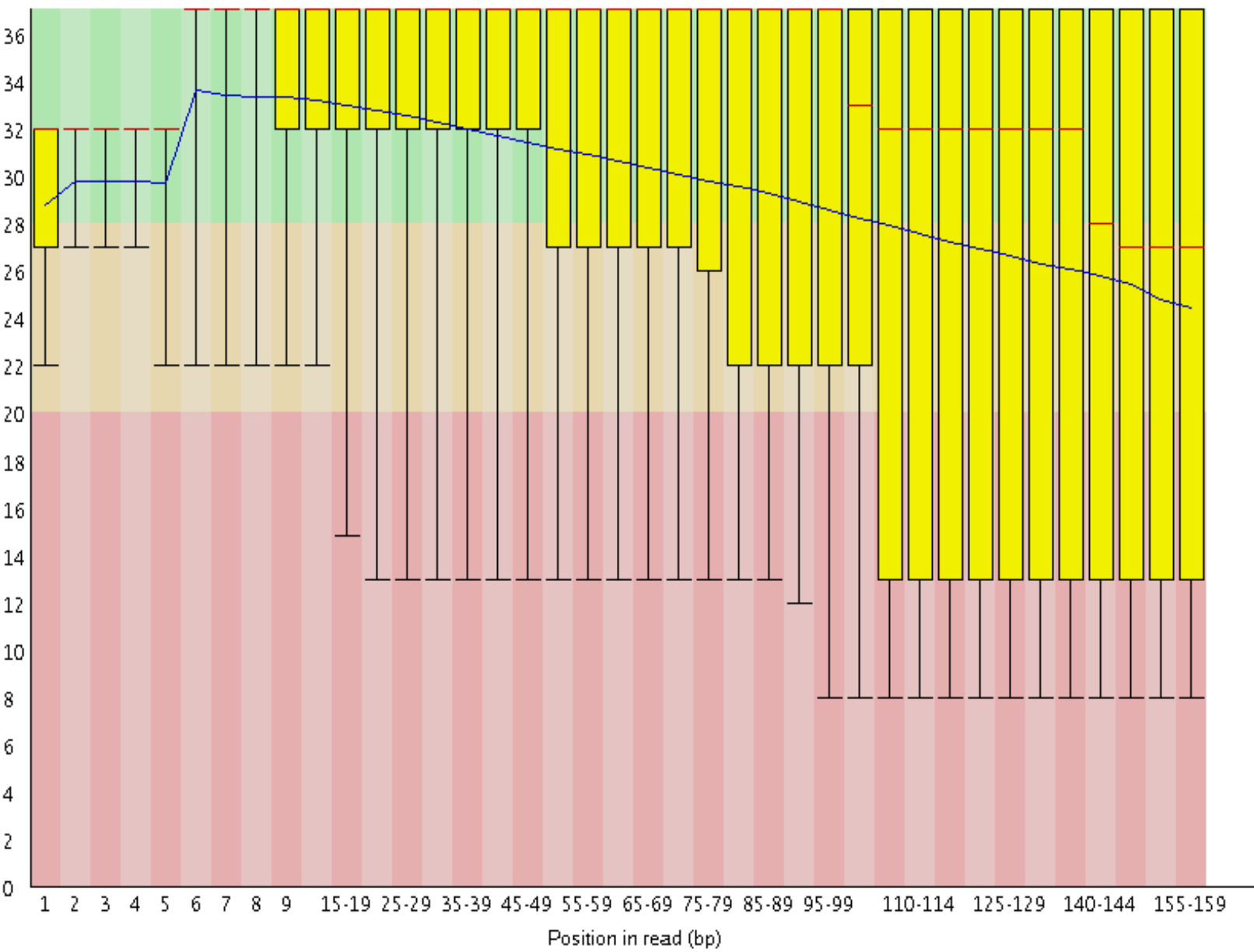✅ Per sequence quality scores

⚠️ Per base sequence content

⚠️ Per sequence GC content

✅ Per base N content

⚠️ Sequence Length Distribution

⚠️ Sequence Duplication Levels

⚠️ Overrepresented sequences

❌ Adapter Content

## ✅ Basic Statistics

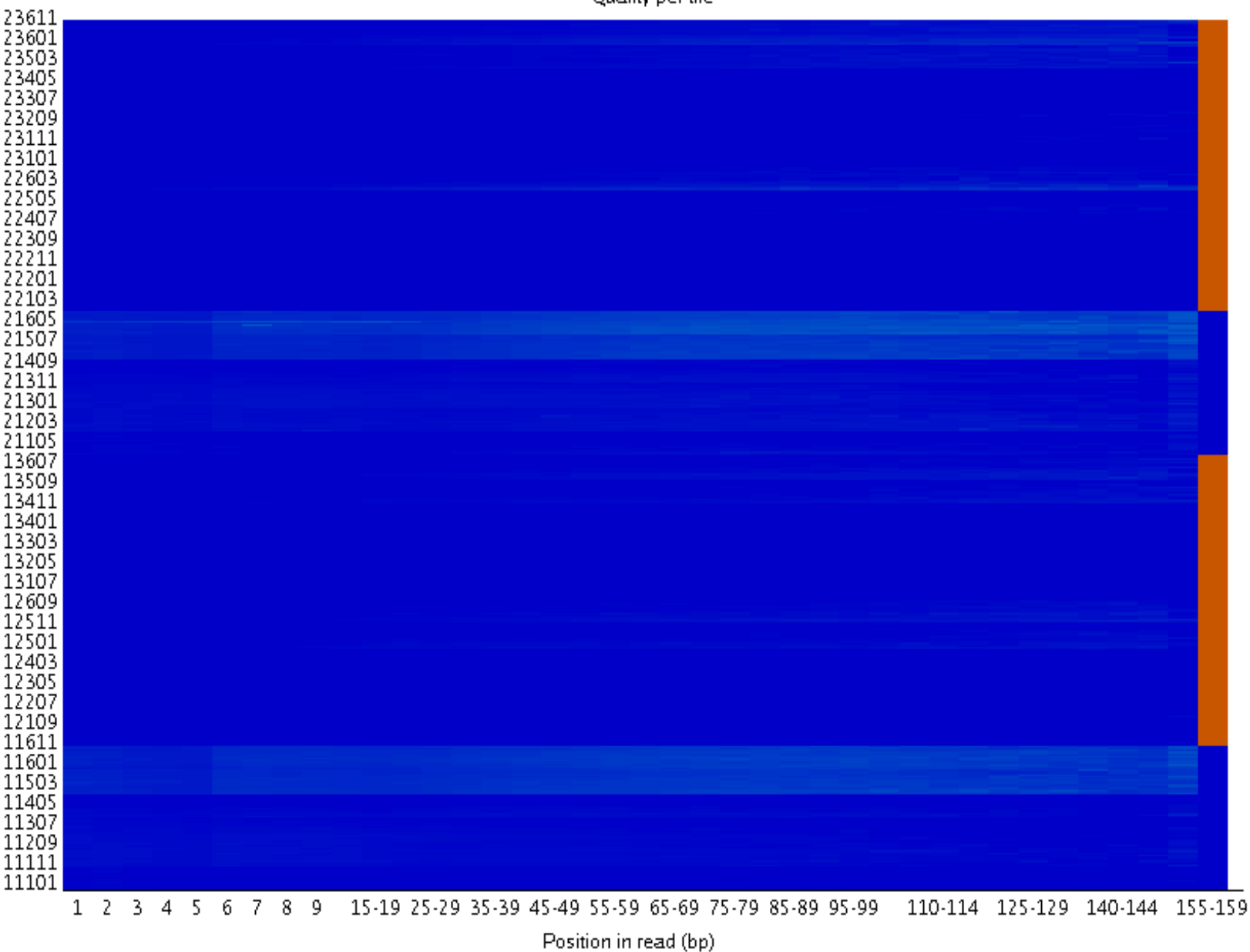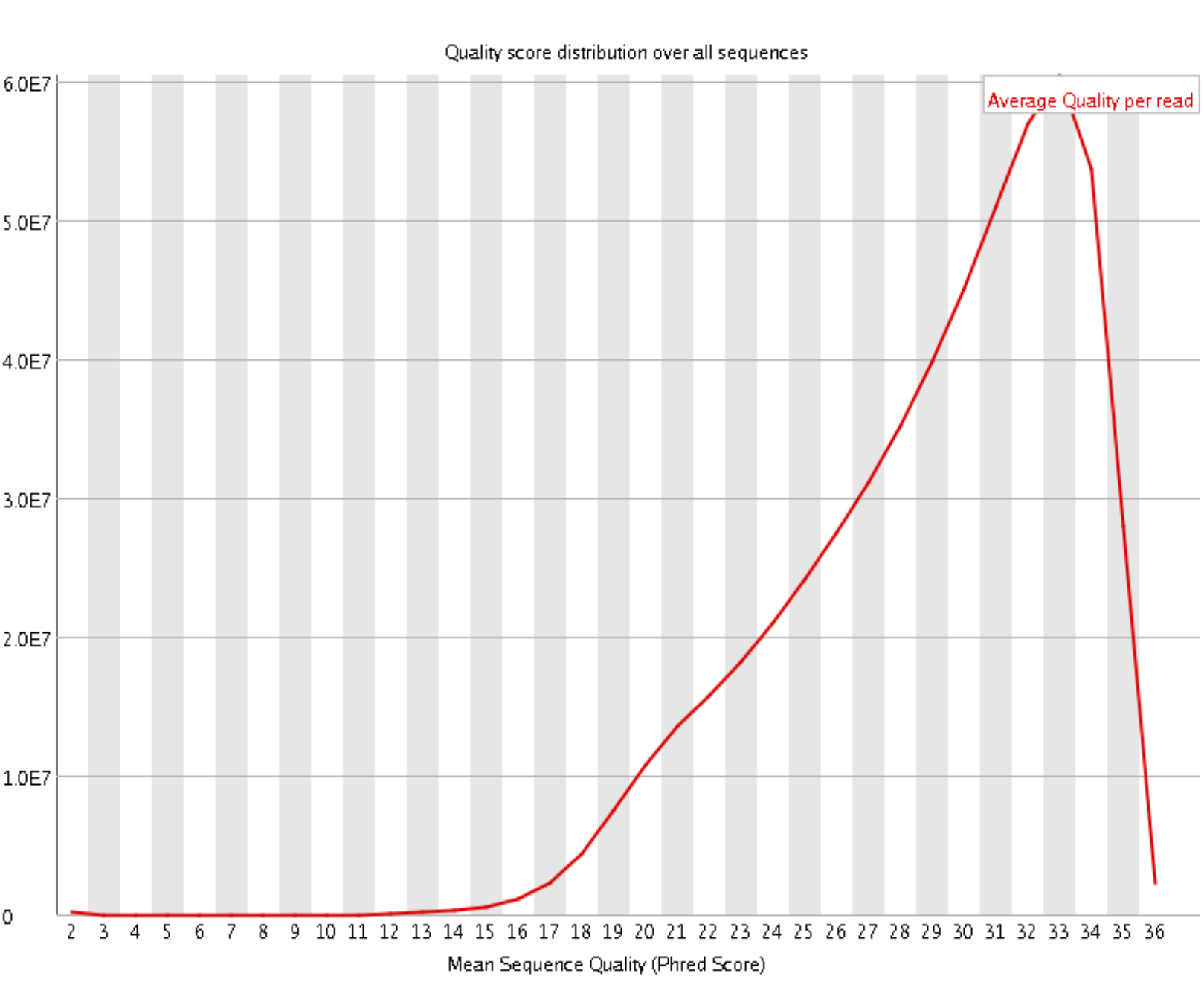| Measure | Value |
|---|---|
| Filename | stdin |
| File type | Conventional base calls |
| Encoding | Sanger / Illumina 1.9 |
| Total Sequences | 553277870 |
| Sequences flagged as poor quality | 0 |
| Sequence length | 35-159 |
| %GC | 41 |

## ✅ Per base sequence quality

Quality scores across all bases (Sanger / Illumina 1.9 encoding)

Position in read (bp)

**Per tile sequence quality**

Quality per tile

Position in read (bp)

✅ **Per sequence quality scores**

Quality score distribution over all sequences
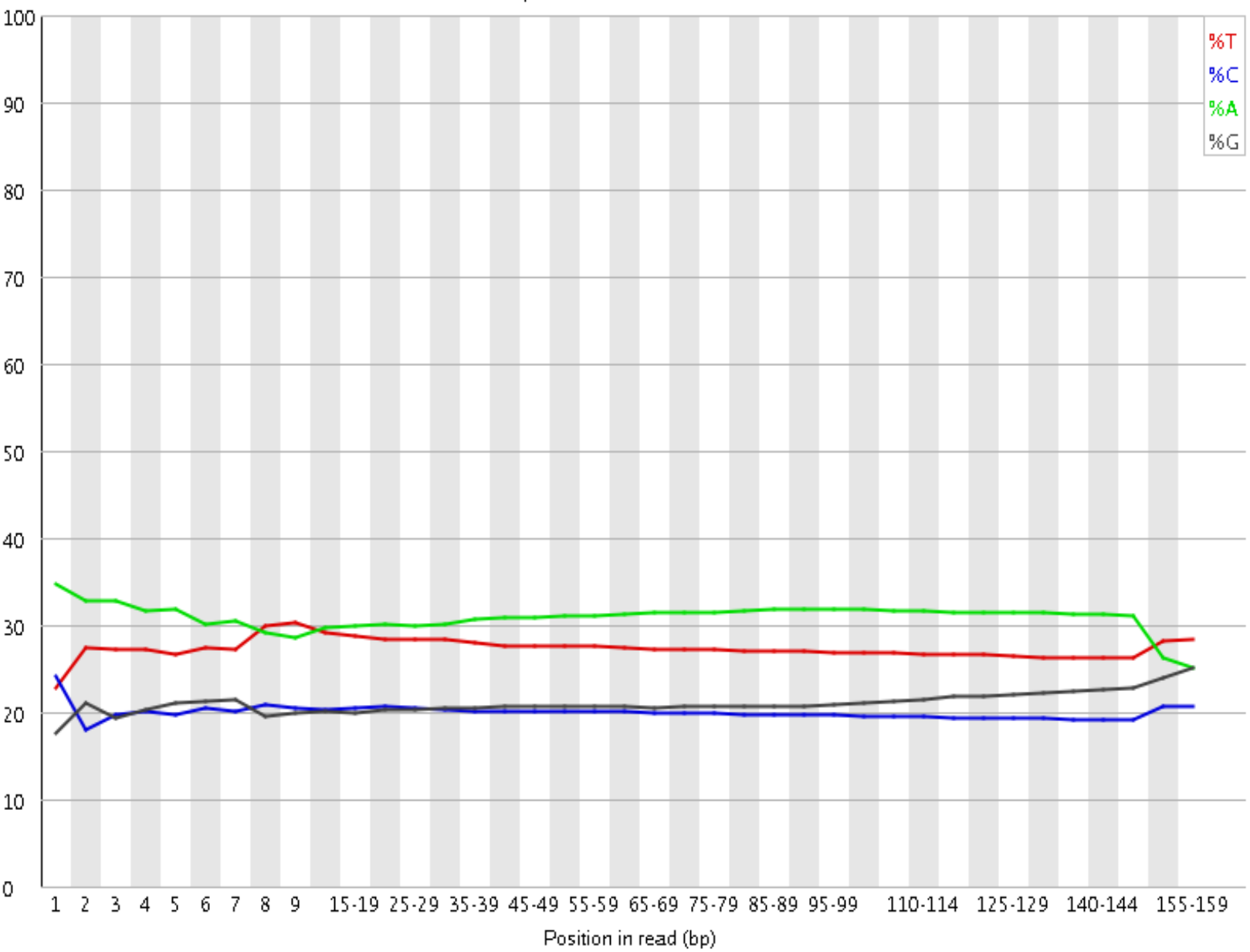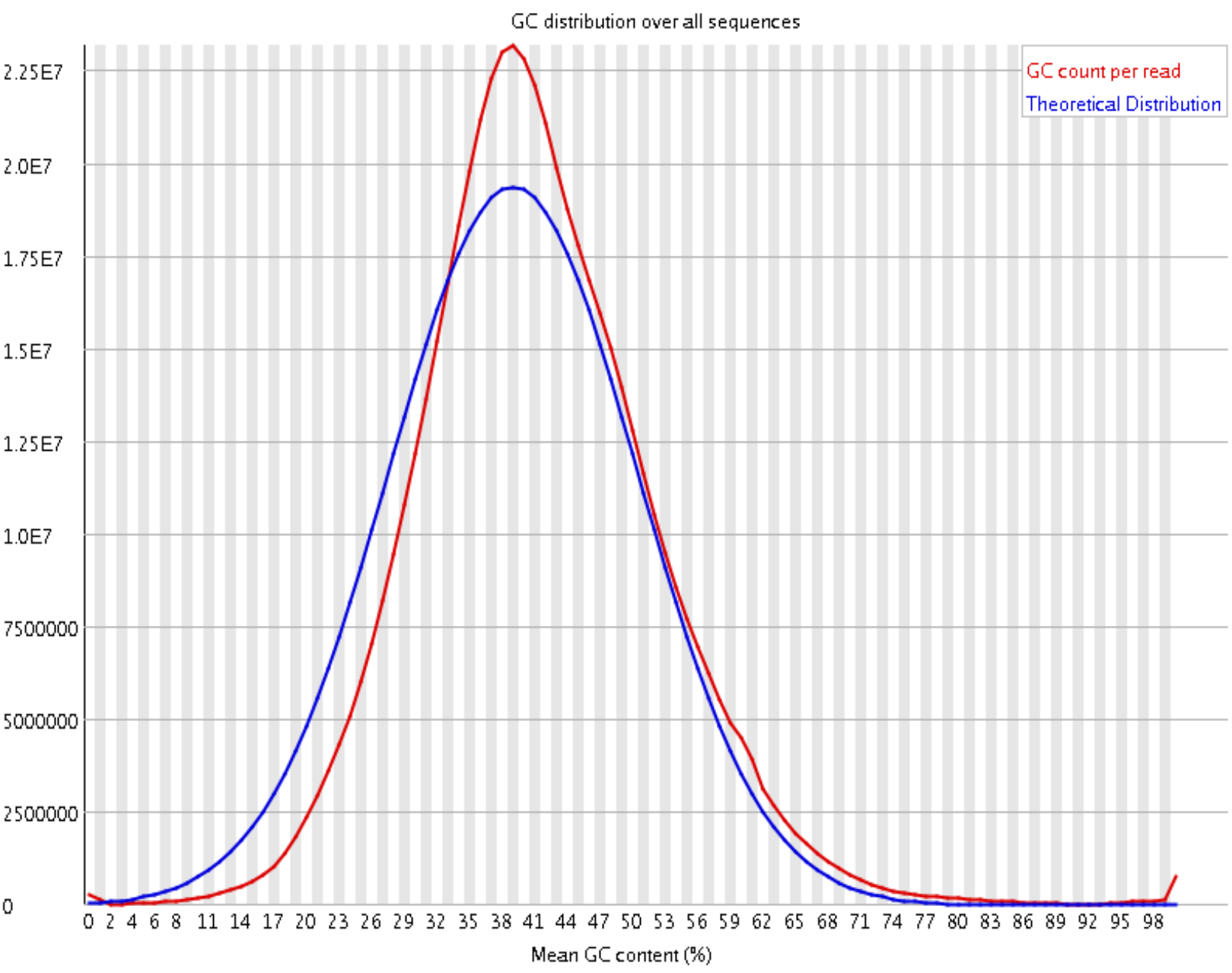
Average Quality per read

Mean Sequence Quality (Phred Score)

⚠ **Per base sequence content**

Sequence content across all bases

## Per sequence GC content

GC distribution over all sequences

GC count per read
Theoretical Distribution

Mean GC content (%)

✅ **Per base N content**

N content across all bases

Position in read (bp)
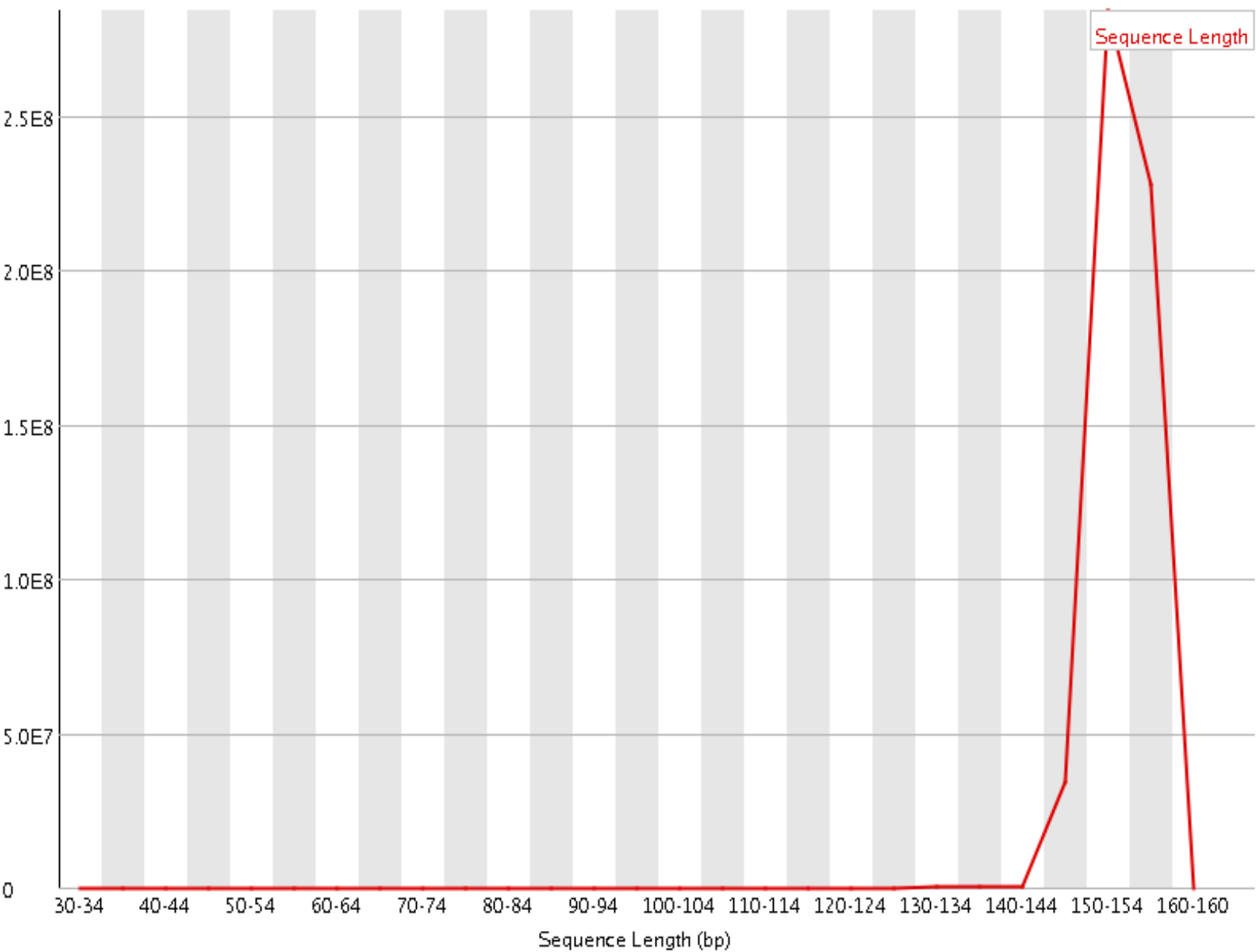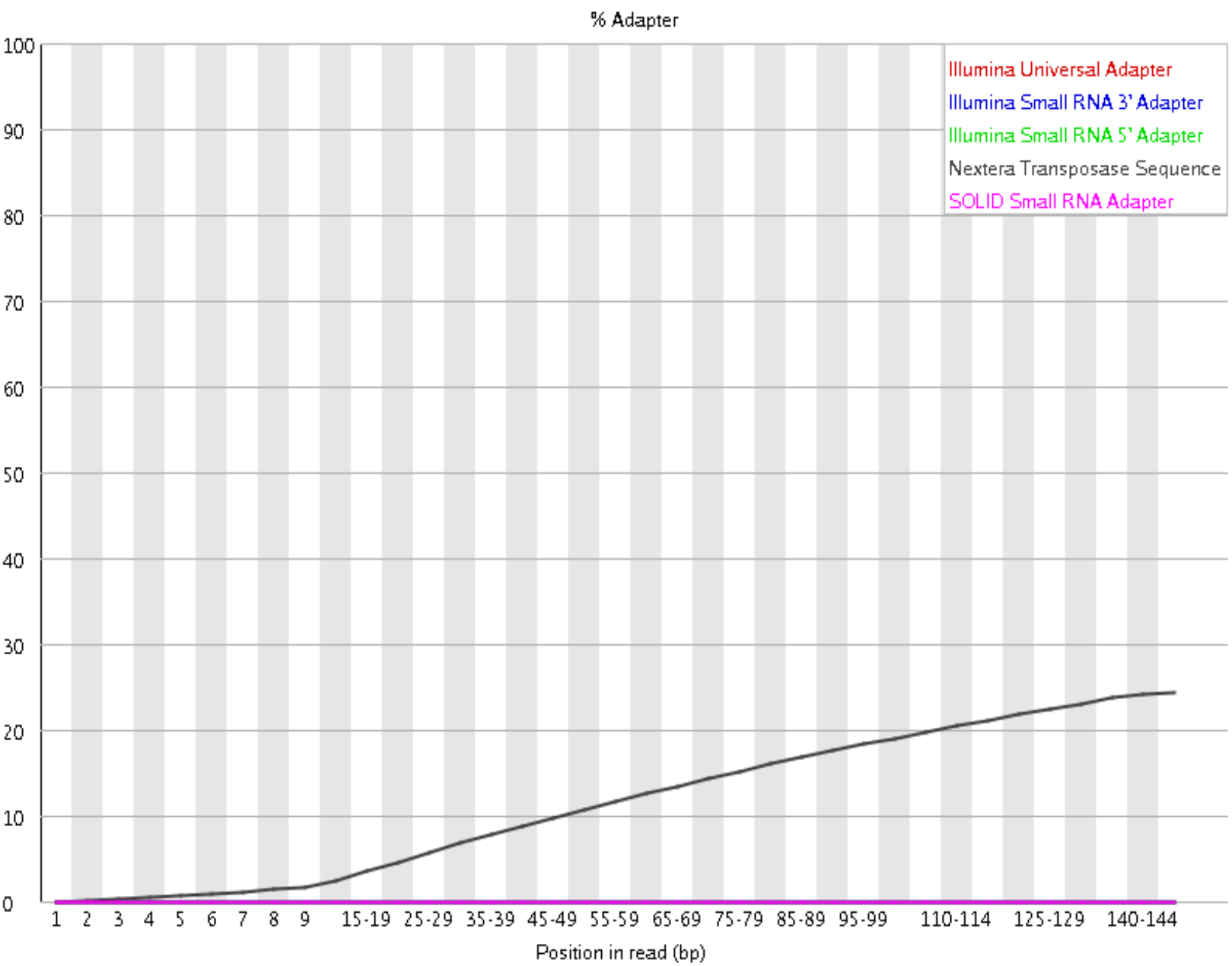
# Sequence Length Distribution

Distribution of sequence lengths over all sequences

## ⚠ Sequence Duplication Levels

## ⚠ Overrepresented sequences

| Sequence | Count | Percentage | Possible Source |
|---|---|---|---|
| GGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGG | 1464007 | 0.2646061010898556 | No Hit |

## ❌ Adapter Content

% Adapter

Produced by **FastQC** (version 0.11.7)

# FastQC Report

# Summary

✅ [Basic Statistics](#)
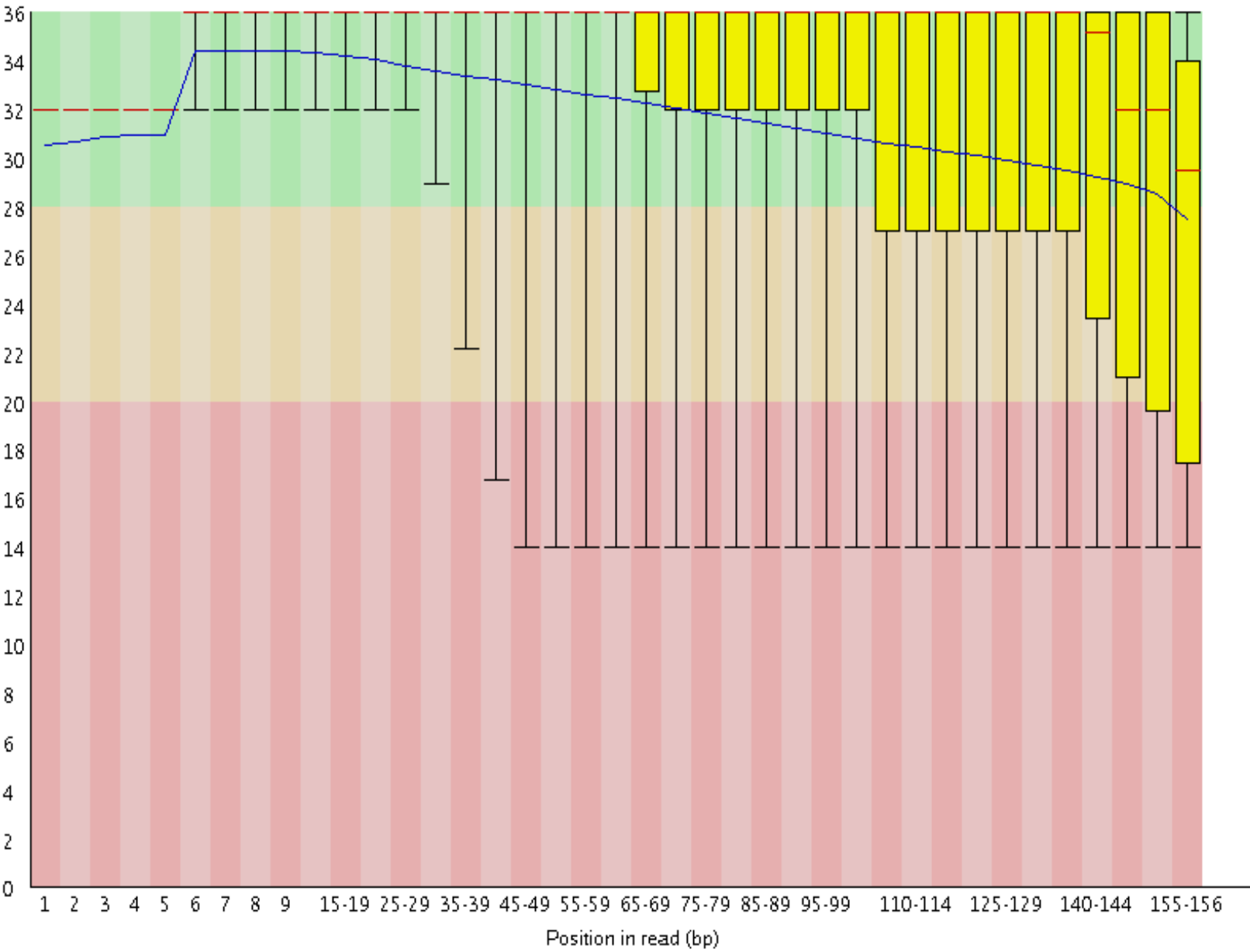
✅ [Per base sequence quality](#)

❌ [Per tile sequence quality](#)

✅ [Per sequence quality scores](#)

⚠️ [Per base sequence content](#)

⚠️ [Per sequence GC content](#)

✅ [Per base N content](#)

⚠️ [Sequence Length Distribution](#)

⚠️ [Sequence Duplication Levels](#)

⚠️ [Overrepresented sequences](#)

❌ [Adapter Content](#)

# ✅ Basic Statistics

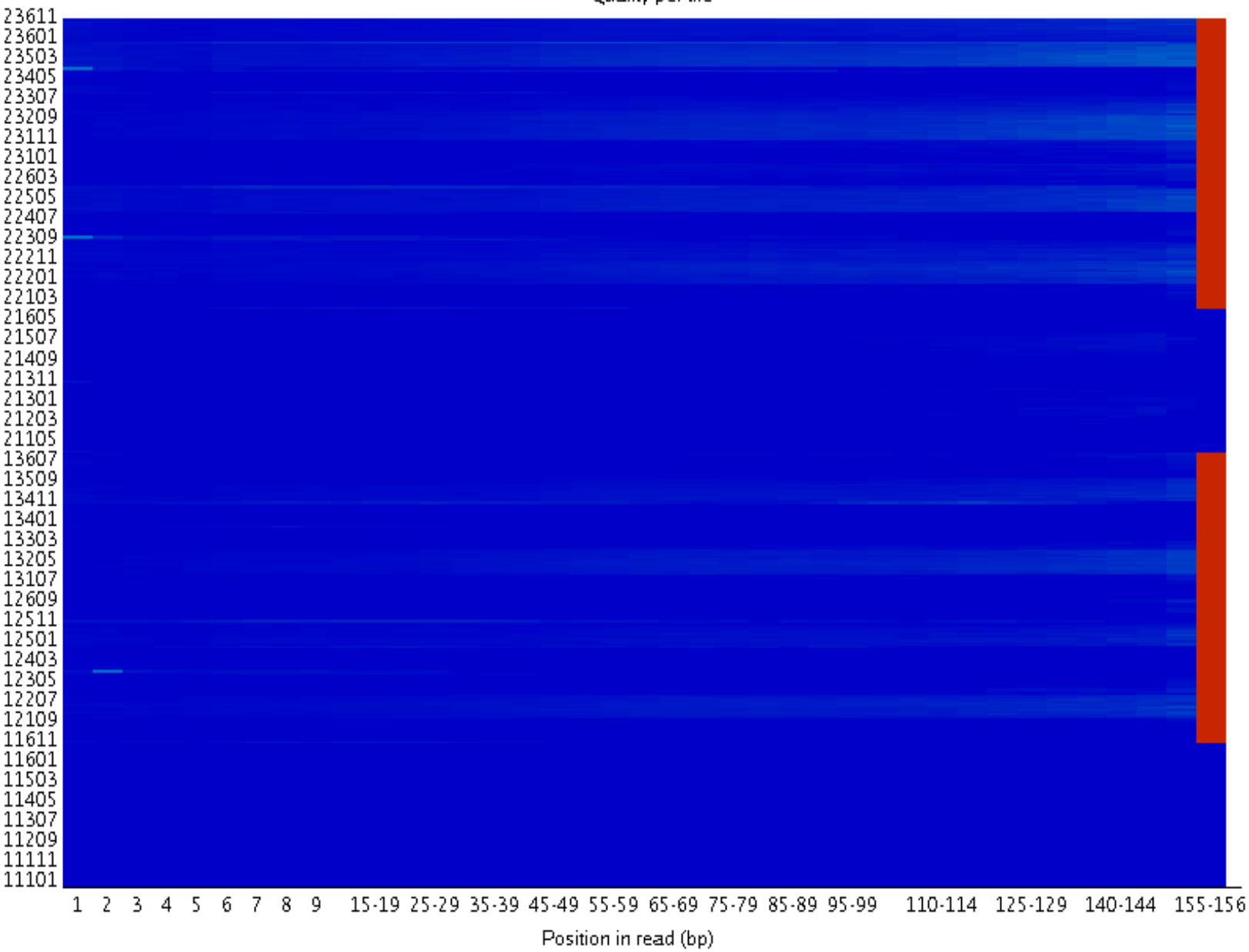| Measure | Value |
|---|---|
| Filename | stdin |
| File type | Conventional base calls |
| Encoding | Sanger / Illumina 1.9 |
| Total Sequences | 338626576 |
| Sequences flagged as poor quality | 0 |
| Sequence length | 35-156 |
| %GC | 43 |

# ✅ Per base sequence quality
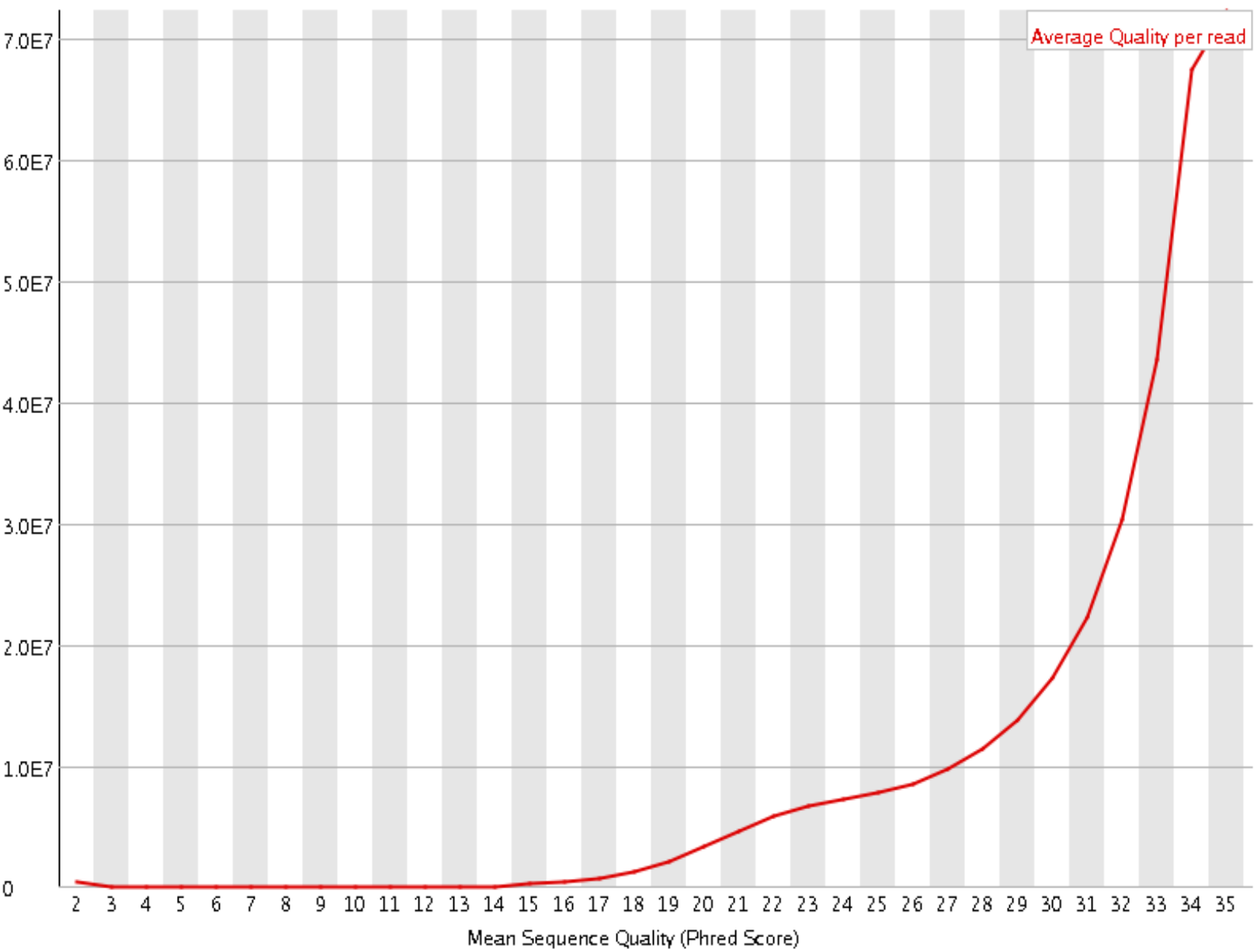
Quality scores across all bases (Sanger / Illumina 1.9 encoding)

Position in read (bp)

❌**Per tile sequence quality**
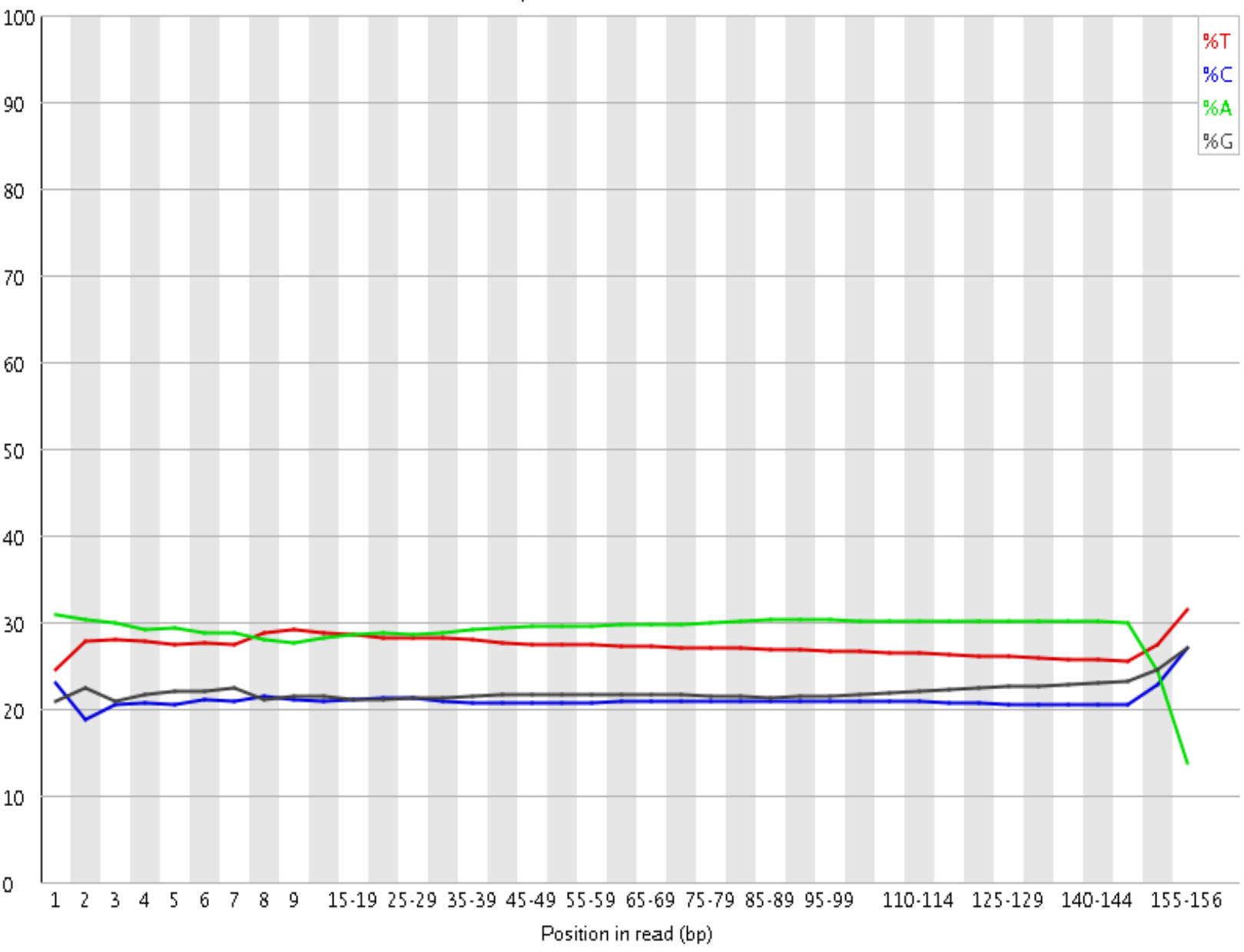
Quality per tile

Position in read (bp)

**Per sequence quality scores**
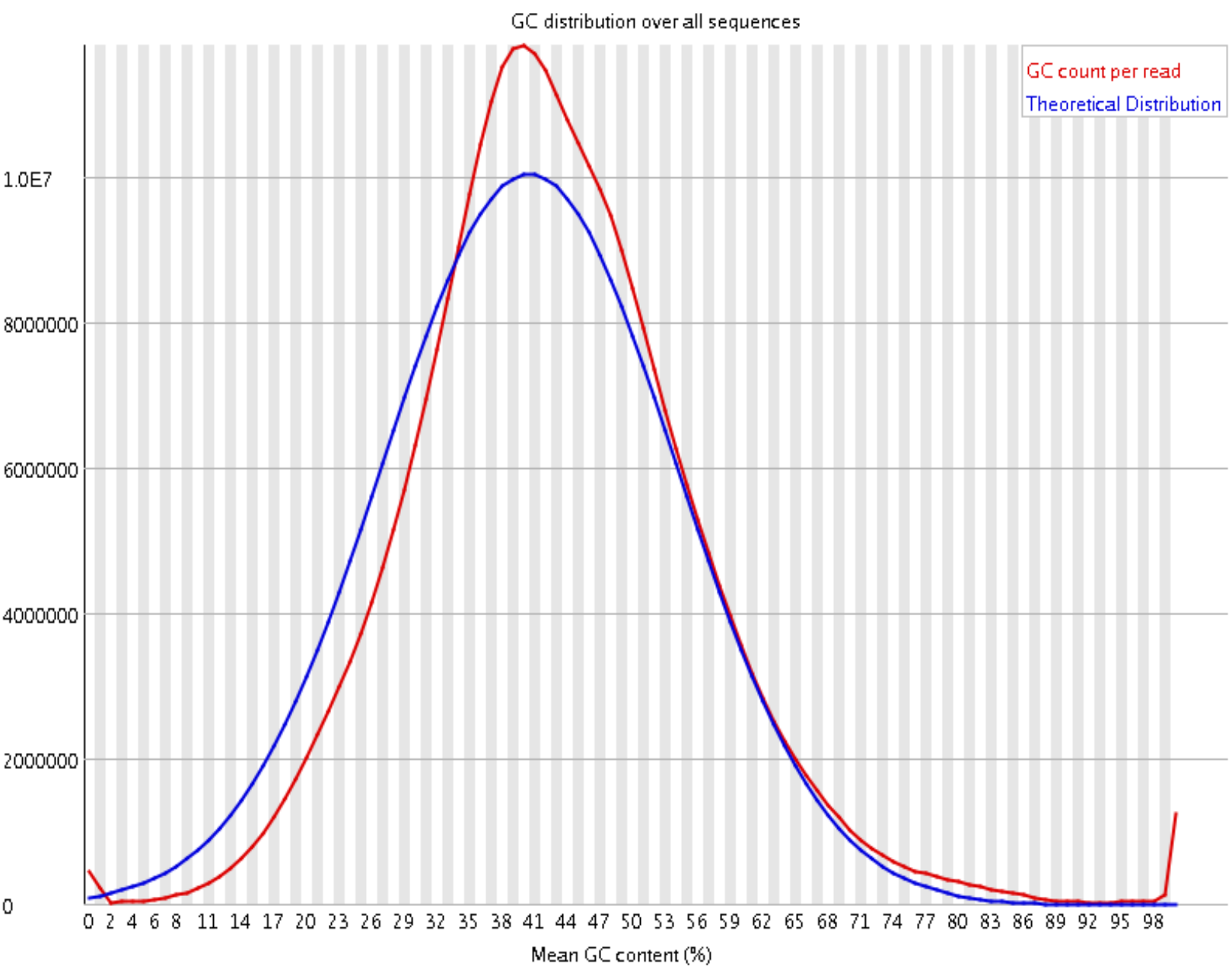
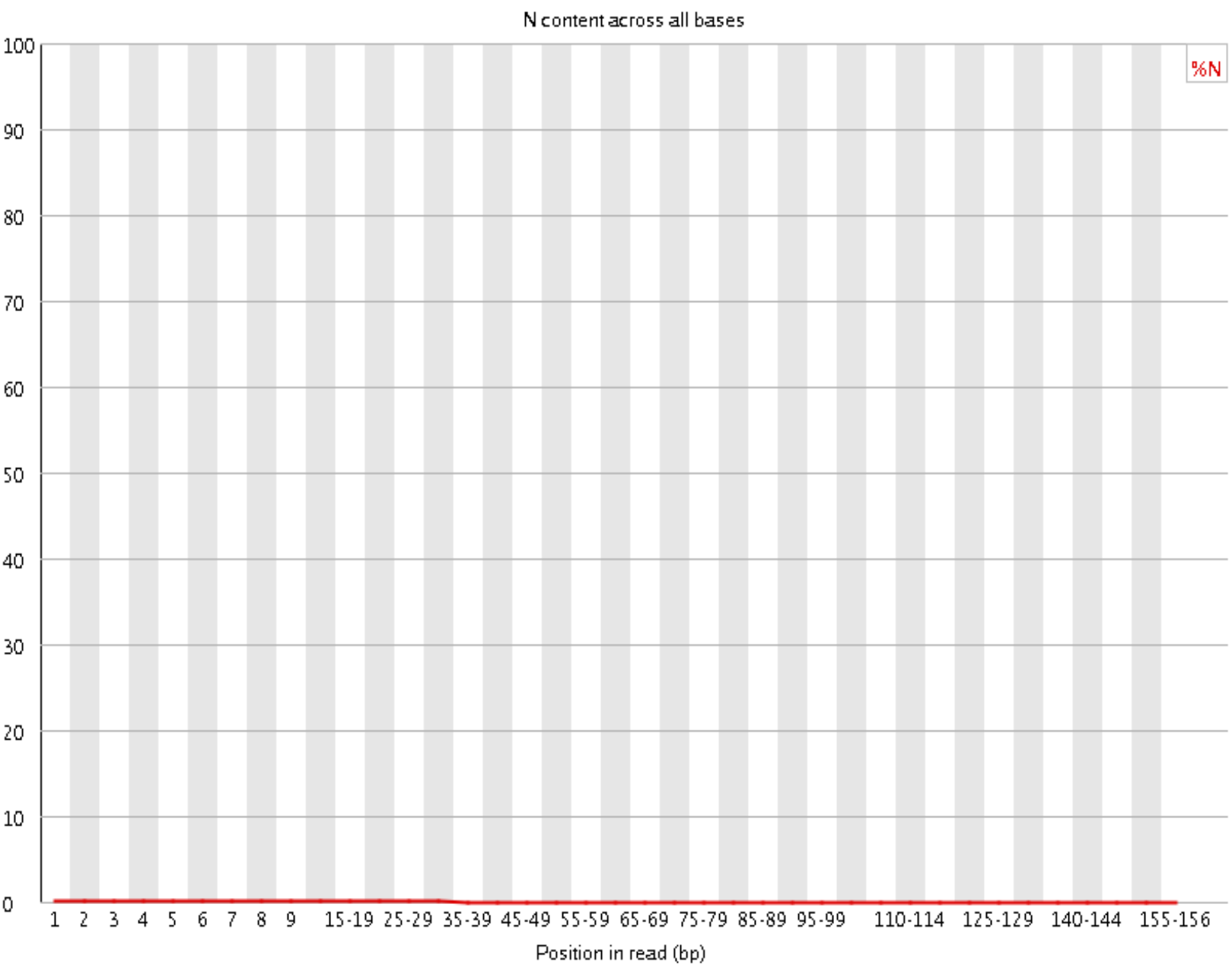Quality score distribution over all sequences

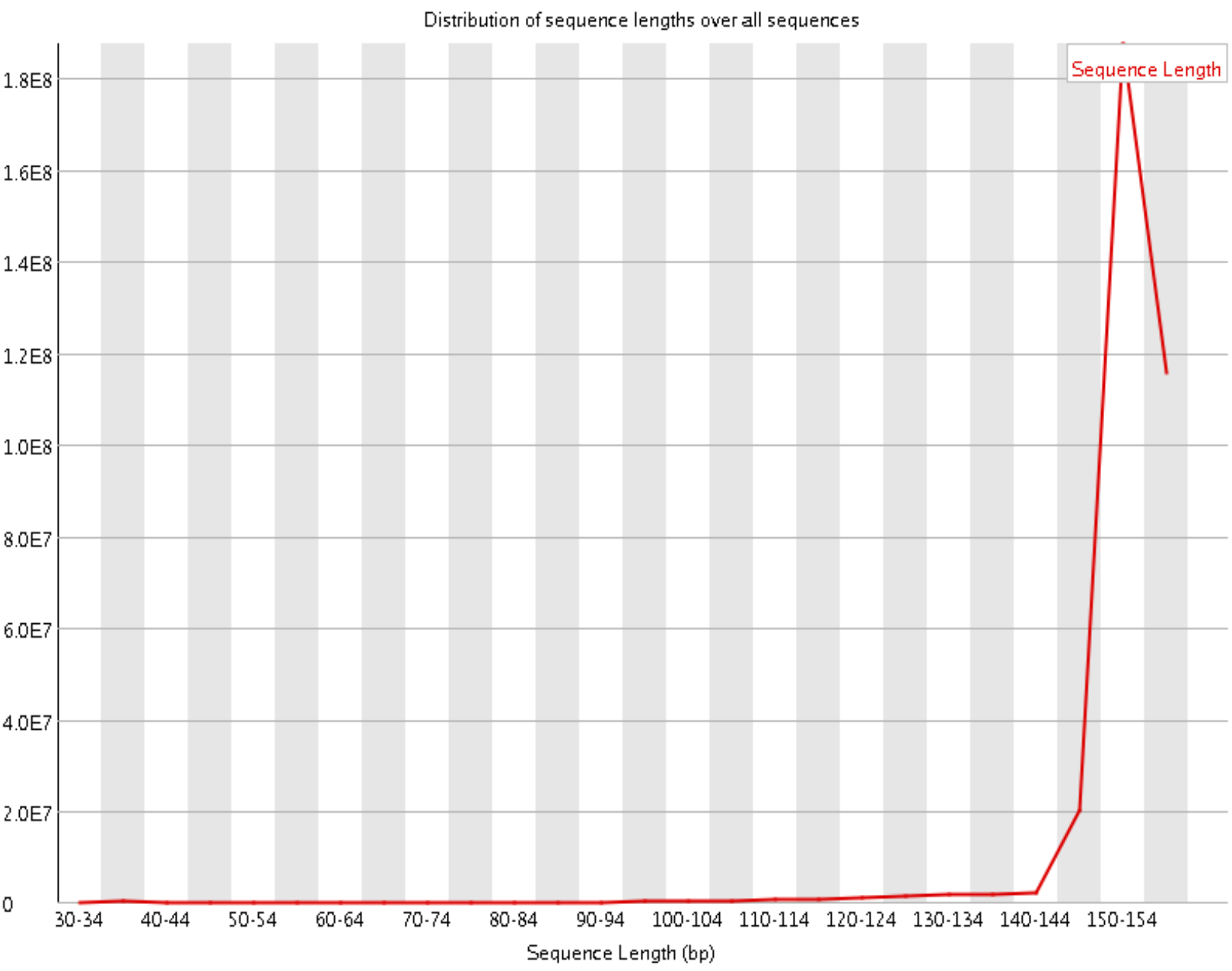**⚠️ Per base sequence content**

Sequence content across all bases

## Per sequence GC content

GC distribution over all sequences

## Per base N content

N content across all bases

## Sequence Length Distribution

Distribution of sequence lengths over all sequences

## Sequence Duplication Levels

Percent of seqs remaining if deduplicated 66.96%

## Overrepresented sequences

| Sequence | Count | Percentage | Possible Source |
|---|---|---|---|
| GGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGG | 2384371 | 0.7041299085751616 | No Hit |
| NNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNN | 414840 | 0.1225066280680817 | No Hit |

## Adapter Content

Produced by **FastQC** (version 0.11.7)

## Summary

✅ [Basic Statistics](#)

✅ [Per base sequence quality](#)

✅ [Per tile sequence quality](#)

✅ [Per sequence quality scores](#)

❌ [Per base sequence content](#)
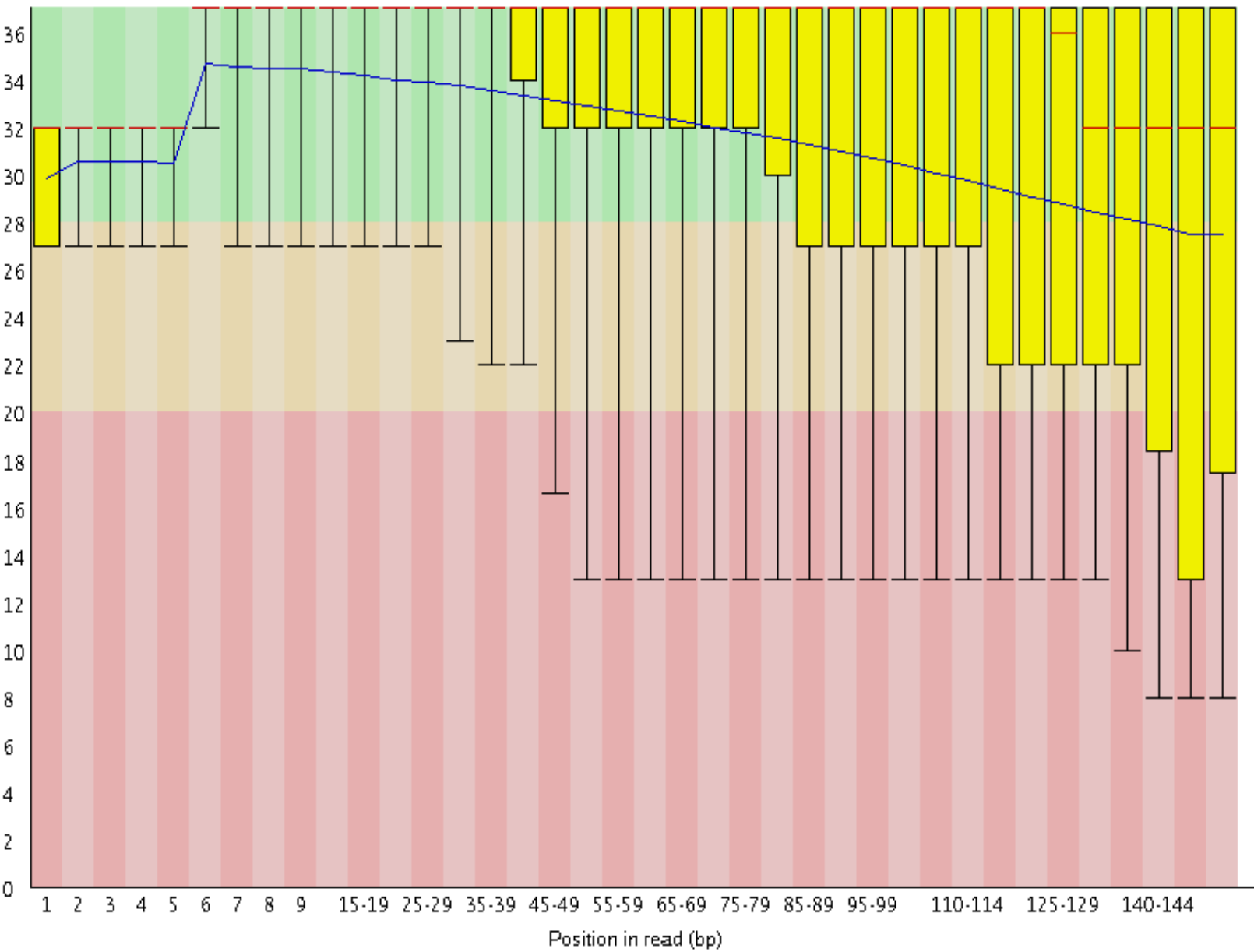
⚠️ [Per sequence GC content](#)

✅ [Per base N content](#)

⚠️ [Sequence Length Distribution](#)

❌ [Sequence Duplication Levels](#)

✅ [Overrepresented sequences](#)

❌ [Adapter Content](#)

## ✅ Basic Statistics

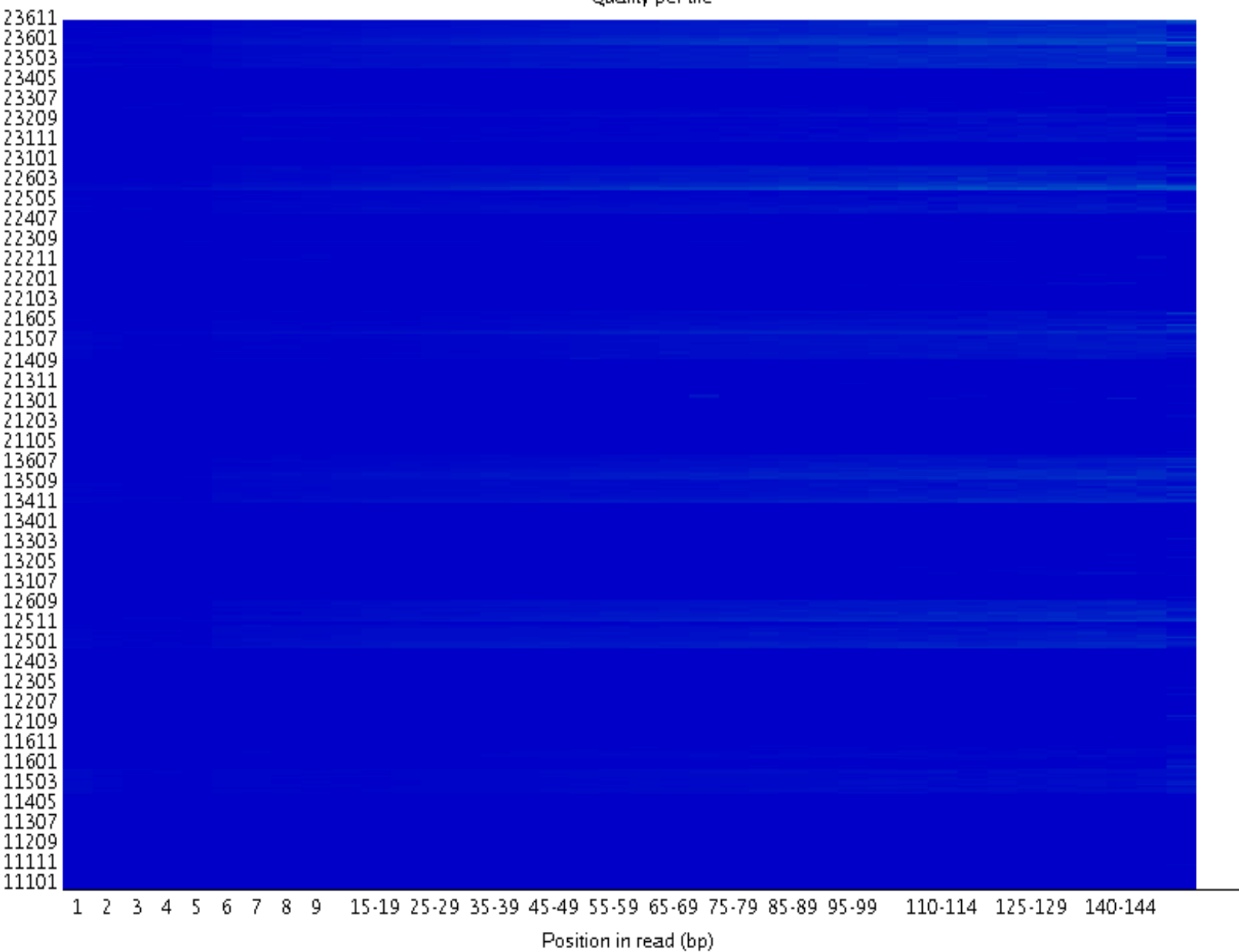| Measure | Value |
|---------|-------|
| Filename | stdin |
| File type | Conventional base calls |
| Encoding | Sanger / Illumina 1.9 |
| Total Sequences | 366199264 |
| Sequences flagged as poor quality | 0 |
| Sequence length | 35-151 |
| %GC | 40 |

## ✅ Per base sequence quality
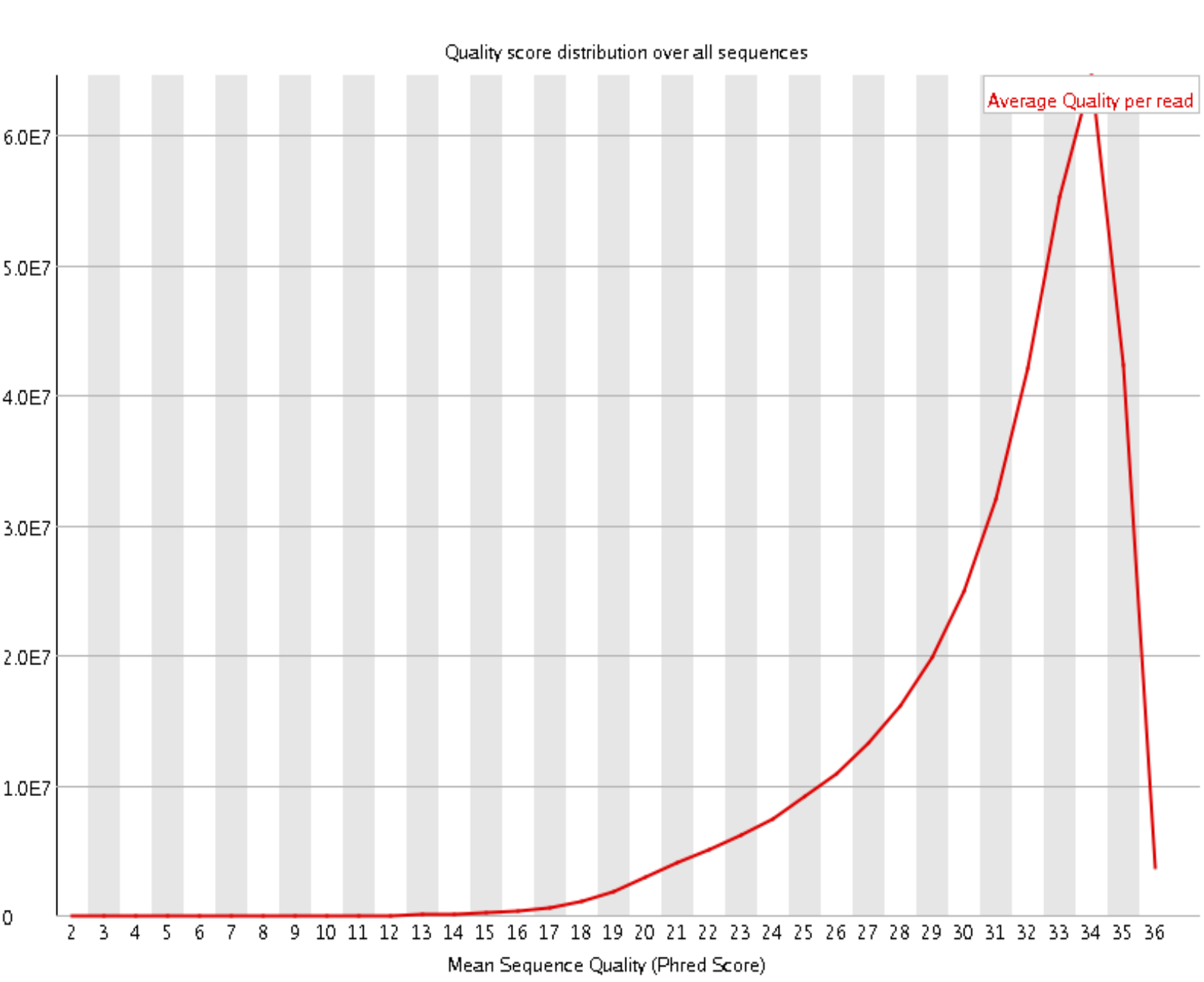
Quality scores across all bases (Sanger / Illumina 1.9 encoding)

Position in read (bp)

✅ **Per tile sequence quality**

Quality per tile

Position in read (bp)

✅ **Per sequence quality scores**

Quality score distribution over all sequences
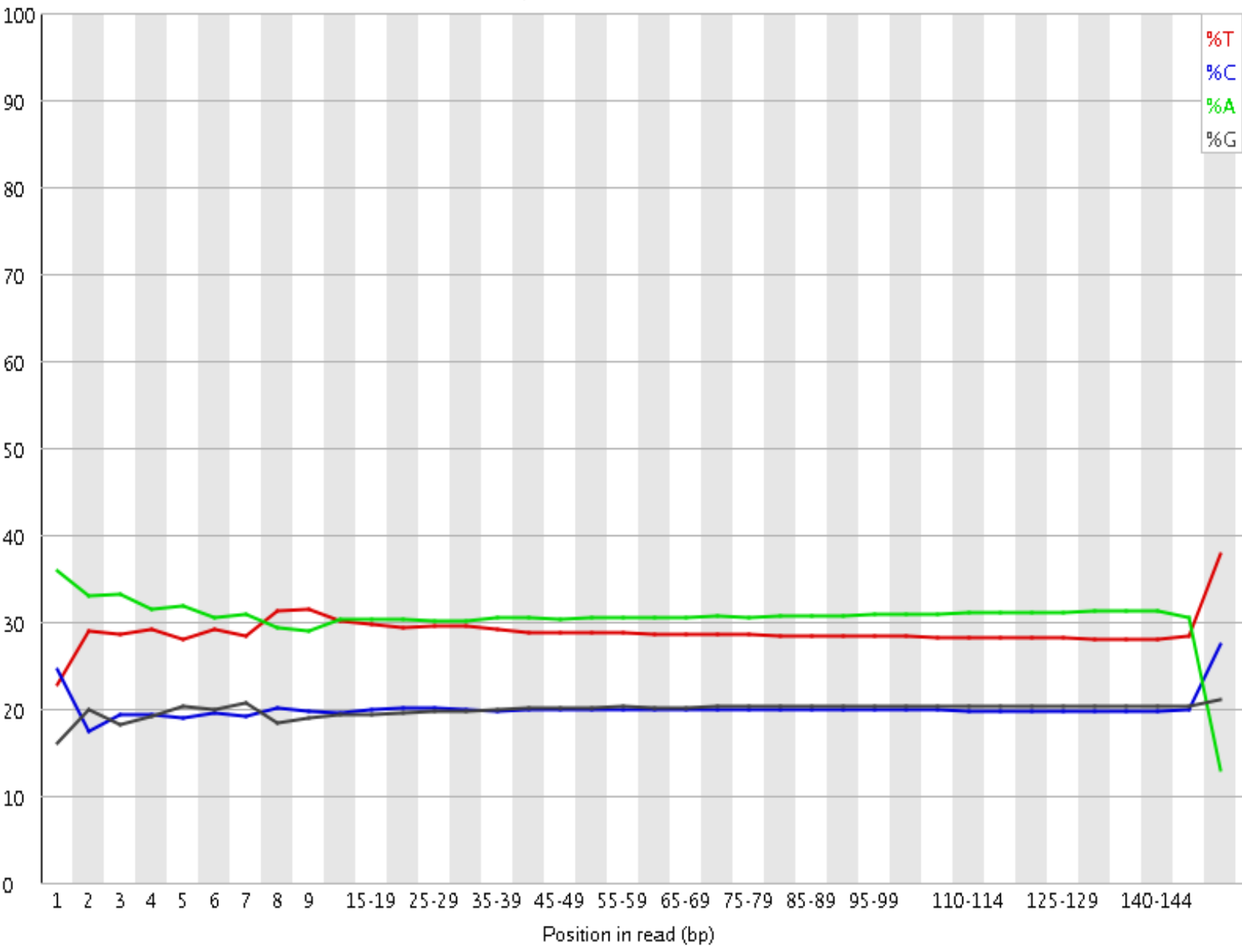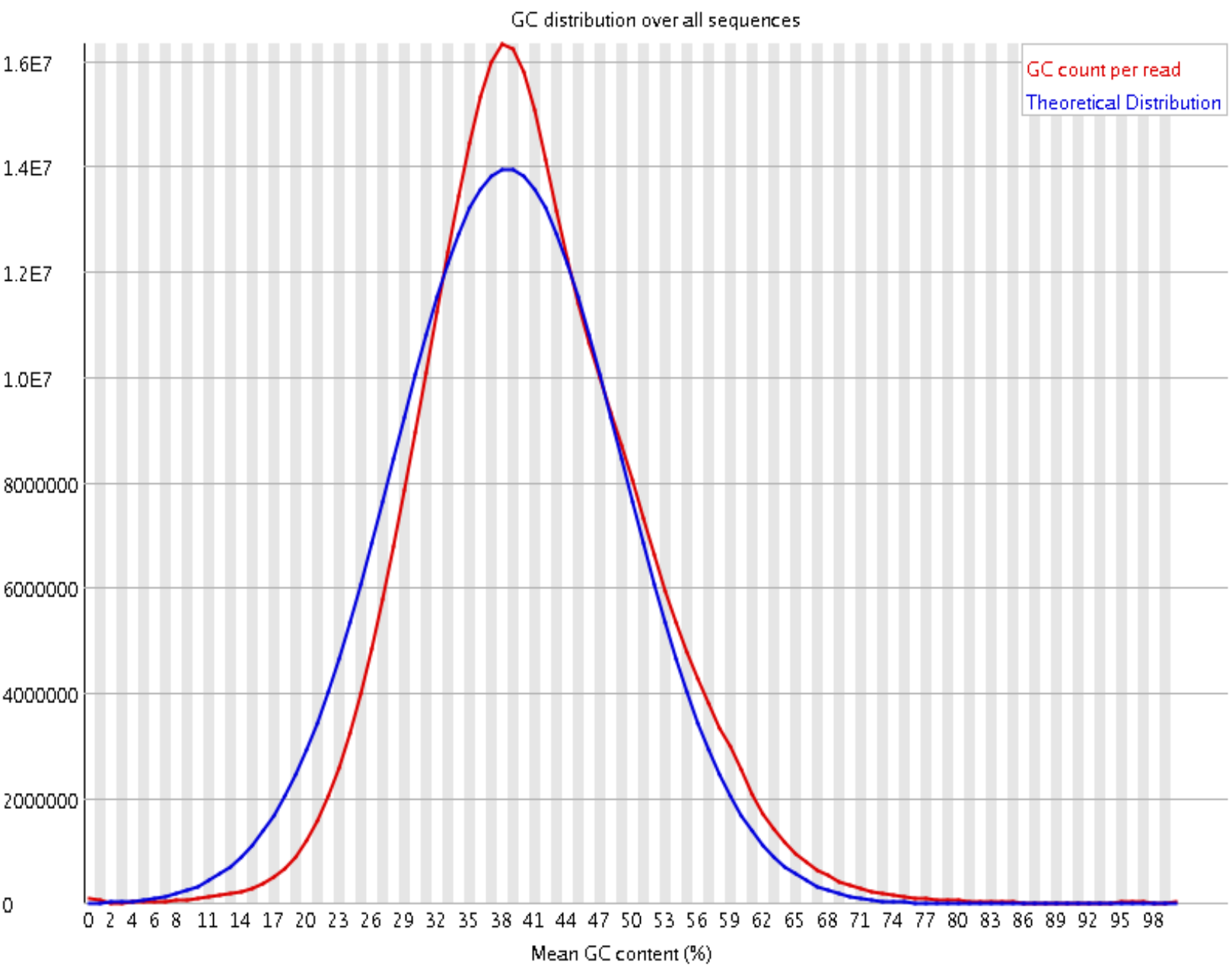
Average Quality per read

Mean Sequence Quality (Phred Score)

❌ **Per base sequence content**

Sequence content across all bases

## Per sequence GC content

GC distribution over all sequences

✅ **Per base N content**
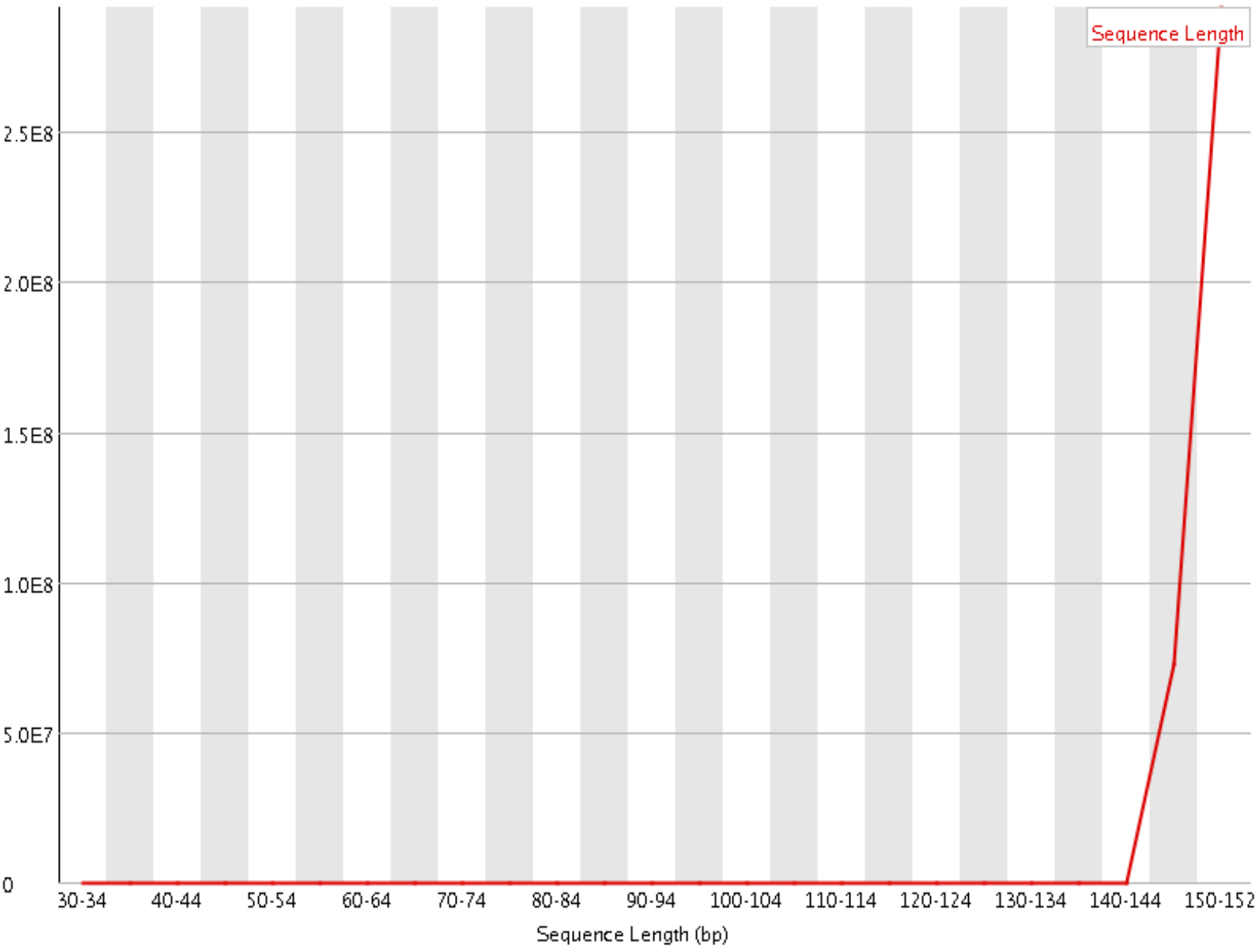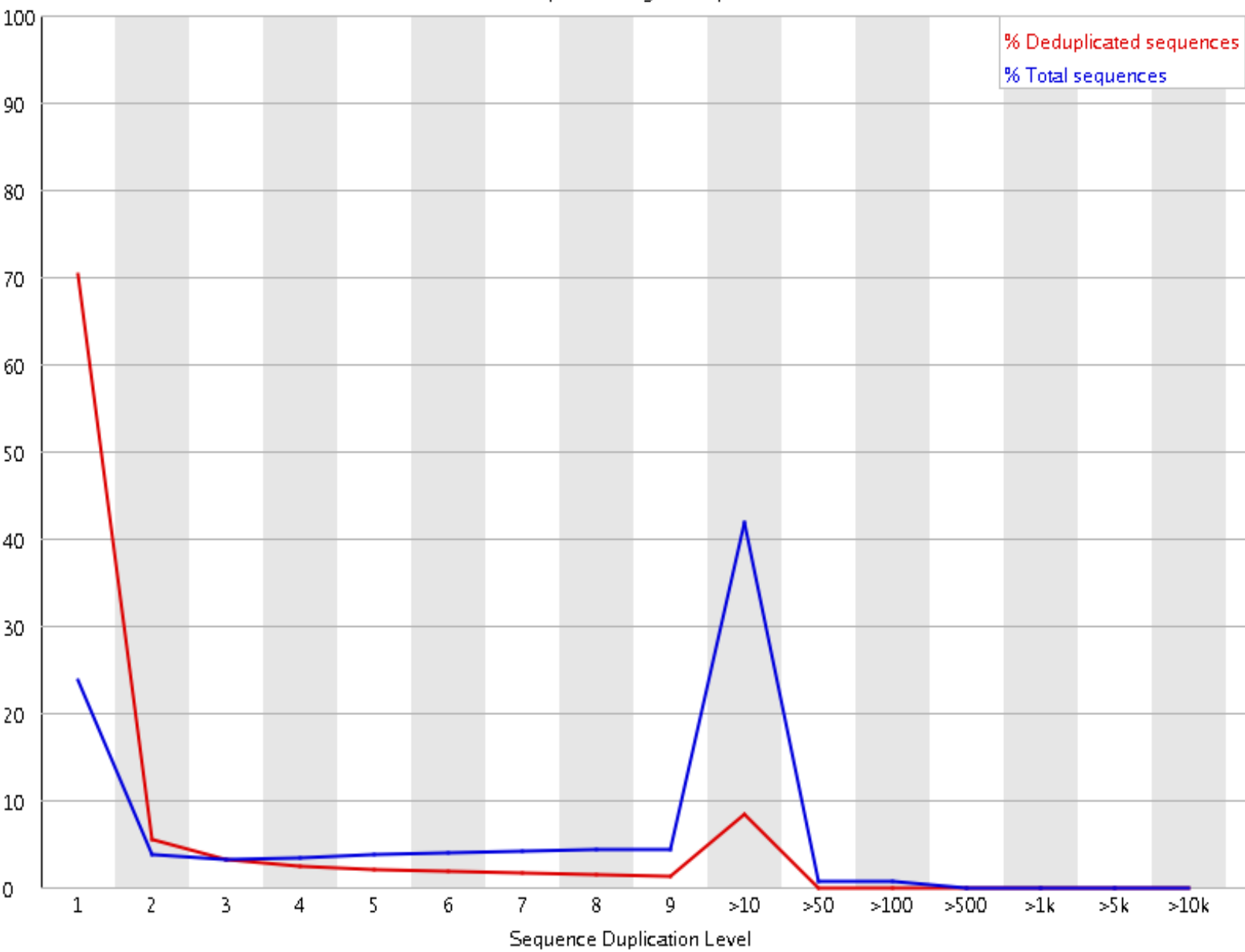
N content across all bases

## Sequence Length Distribution

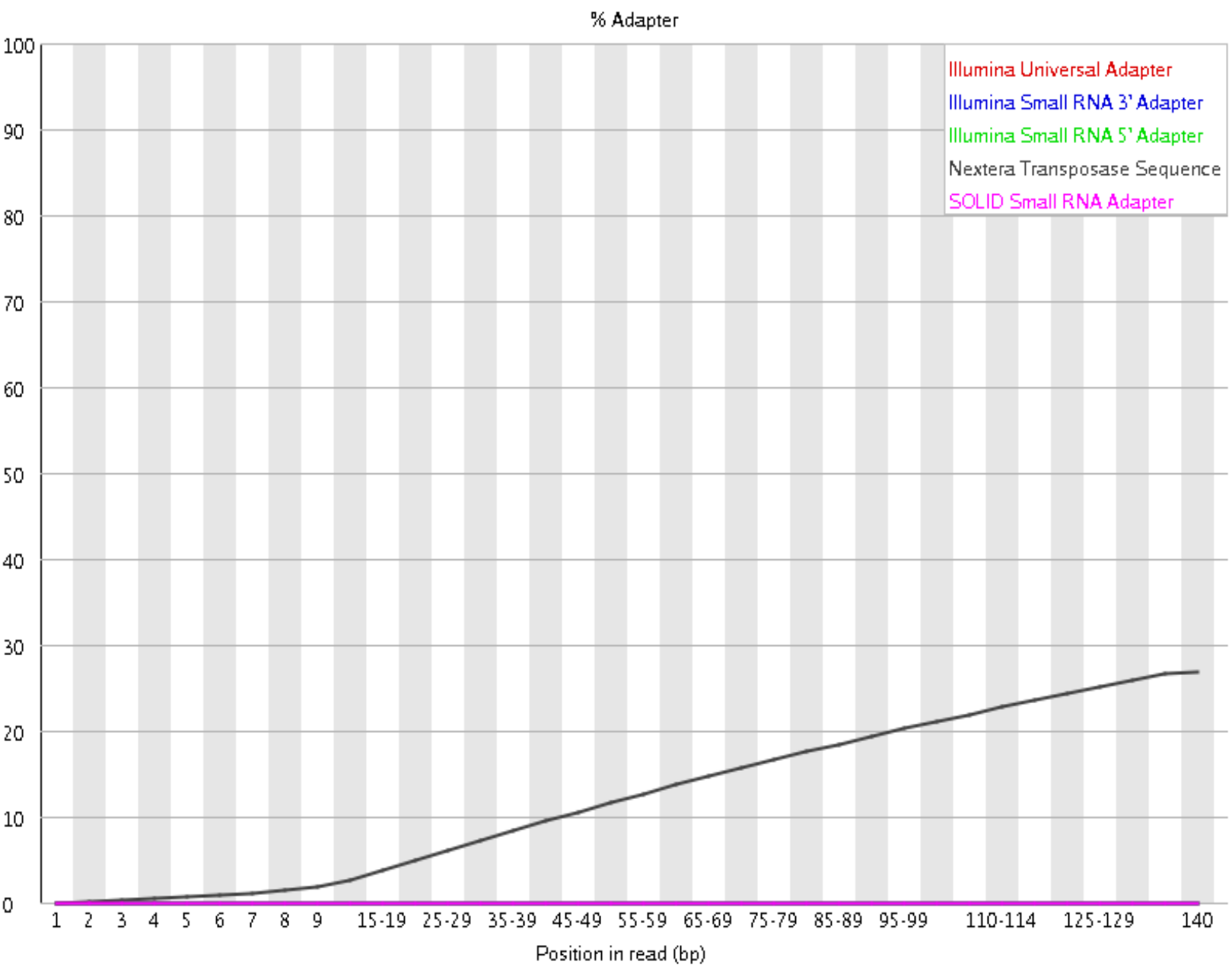Distribution of sequence lengths over all sequences

# Sequence Duplication Levels

Percent of seqs remaining if deduplicated 33.98%

% Deduplicated sequences
% Total sequences

Sequence Duplication Level

✅ **Overrepresented sequences**

No overrepresented sequences

❌ **Adapter Content**

% Adapter

# FastQC Report

# Summary

✅ Basic Statistics

✅ Per base sequence quality

❌ Per tile sequence quality

✅ Per sequence quality scores

✅ Per base sequence content

⚠️ Per sequence GC content

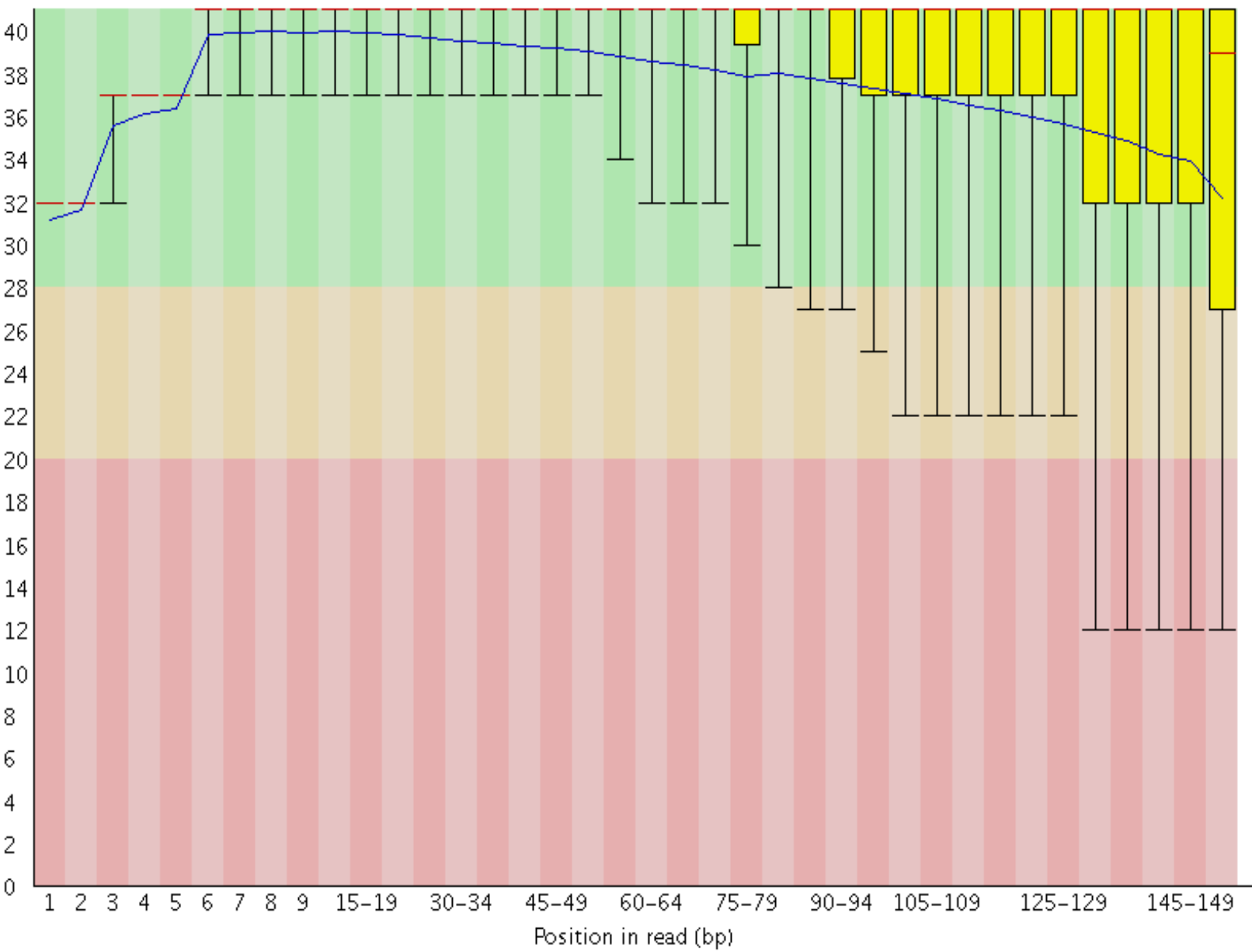✅ Per base N content

✅ Sequence Length Distribution

✅ Sequence Duplication Levels

✅ Overrepresented sequences

✅ Adapter Content

# ✅ Basic Statistics

| Measure | Value |
|---|---|
| Filename | Venter_S1_merge_R1.fastq.gz |
| File type | Conventional base calls |
| Encoding | Sanger / Illumina 1.9 |
| Total Sequences | 789239544 |
| Sequences flagged as poor quality | 0 |
| Sequence length | 151 |
| %GC | 44 |

# ✅ Per base sequence quality

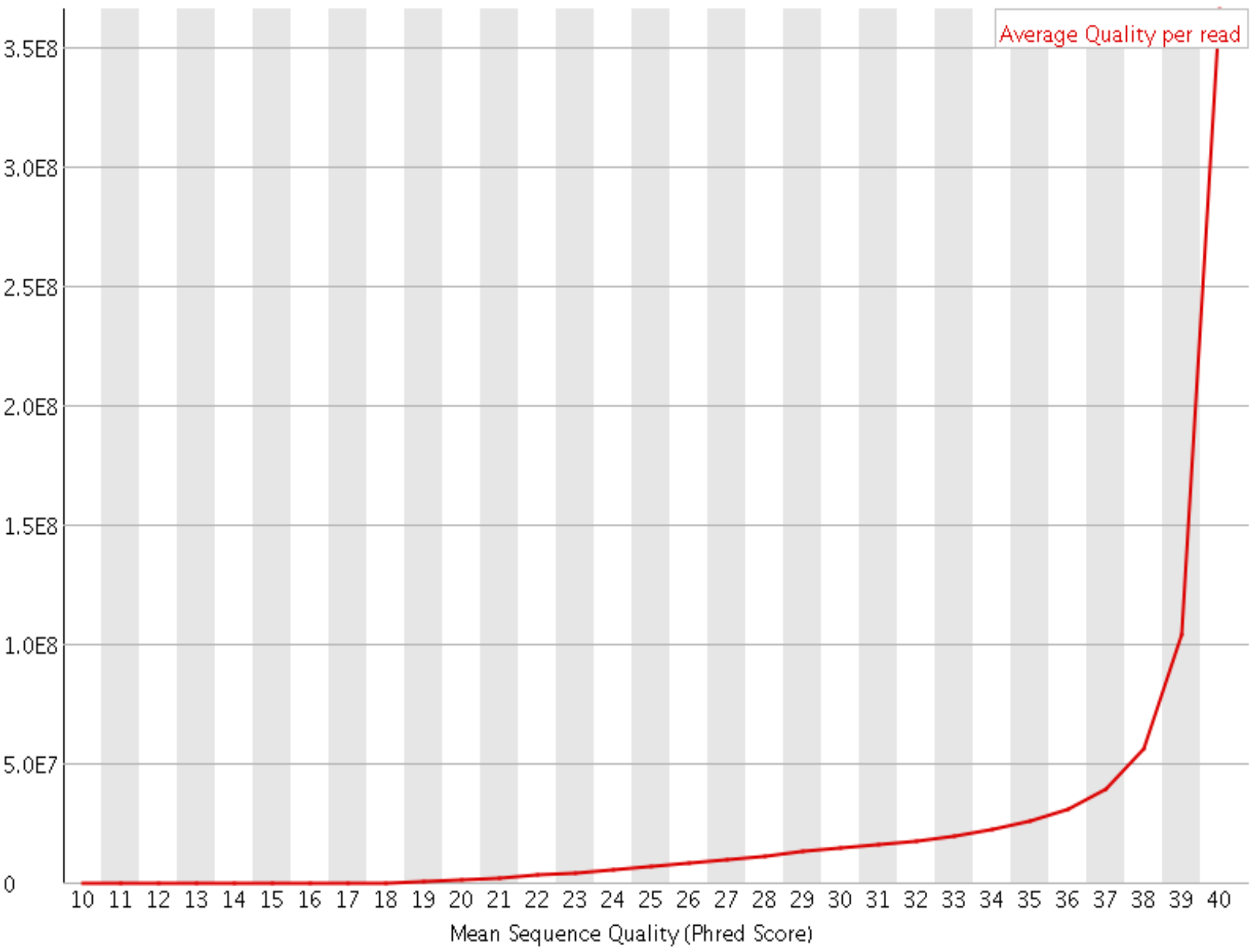Quality scores across all bases (Sanger / Illumina 1.9 encoding)

Position in read (bp)

❌ **Per tile sequence quality**
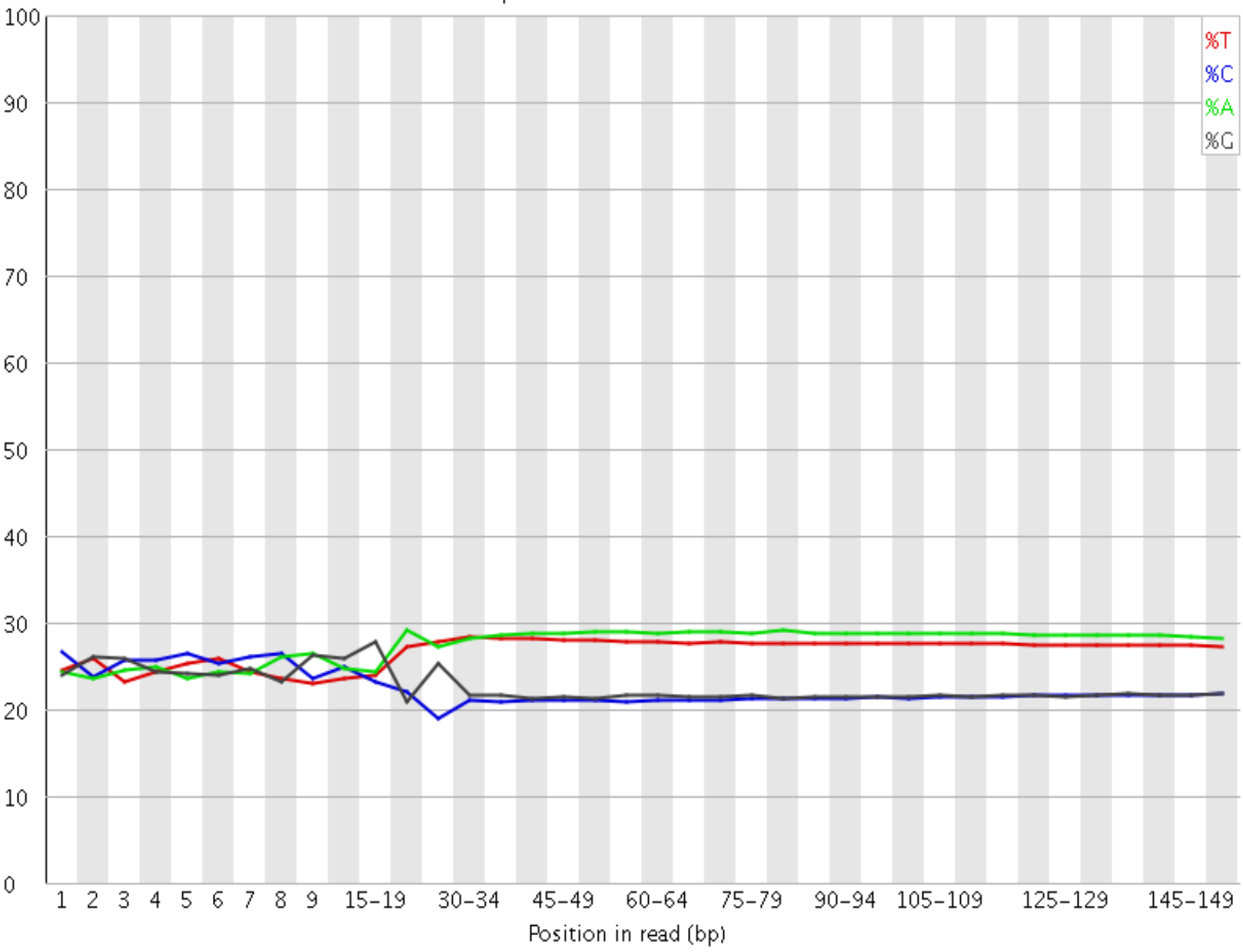
Quality per tile

Position in read (bp)

⊘ **Per sequence quality scores**

Quality score distribution over all sequences
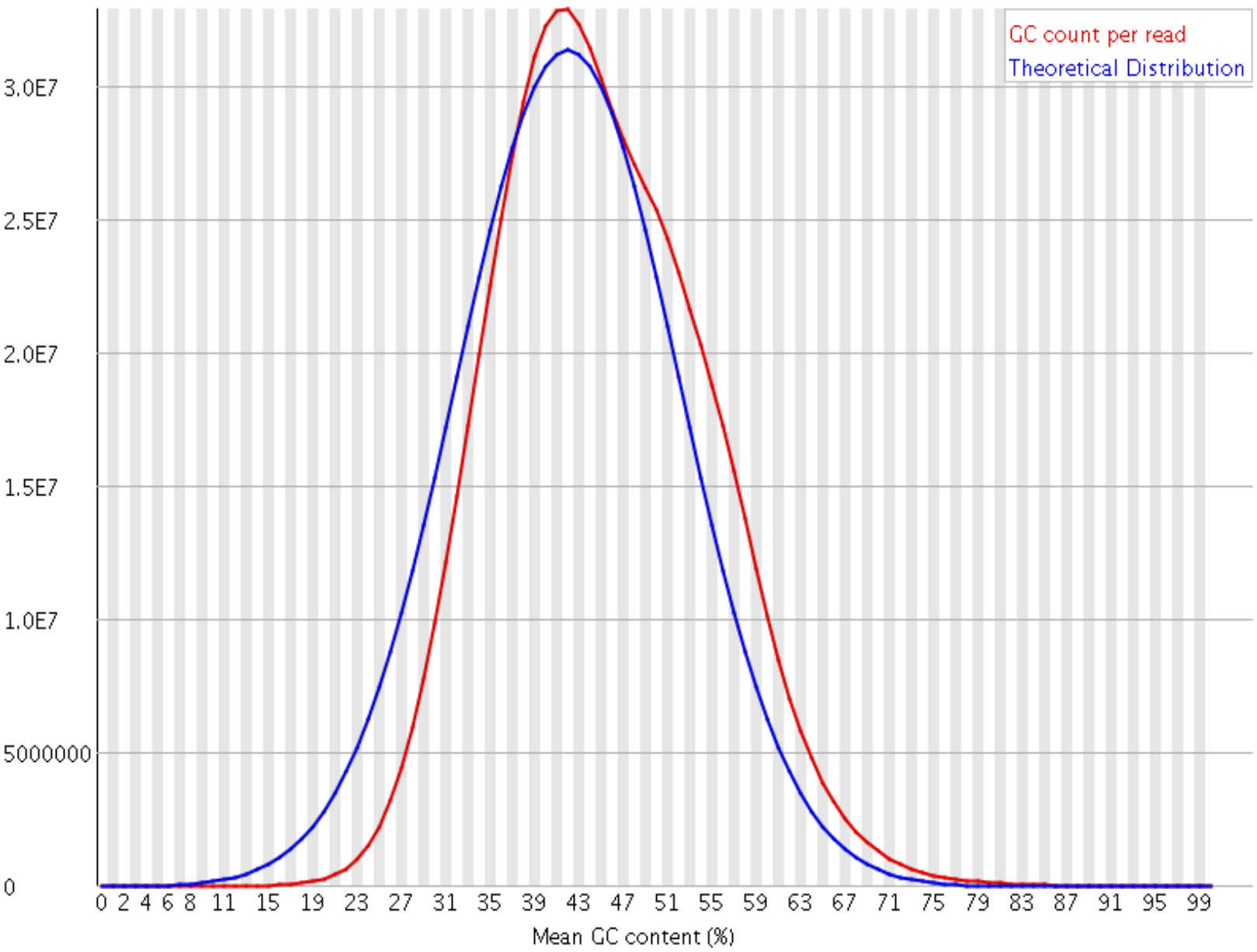
**Per base sequence content**

Sequence content across all bases
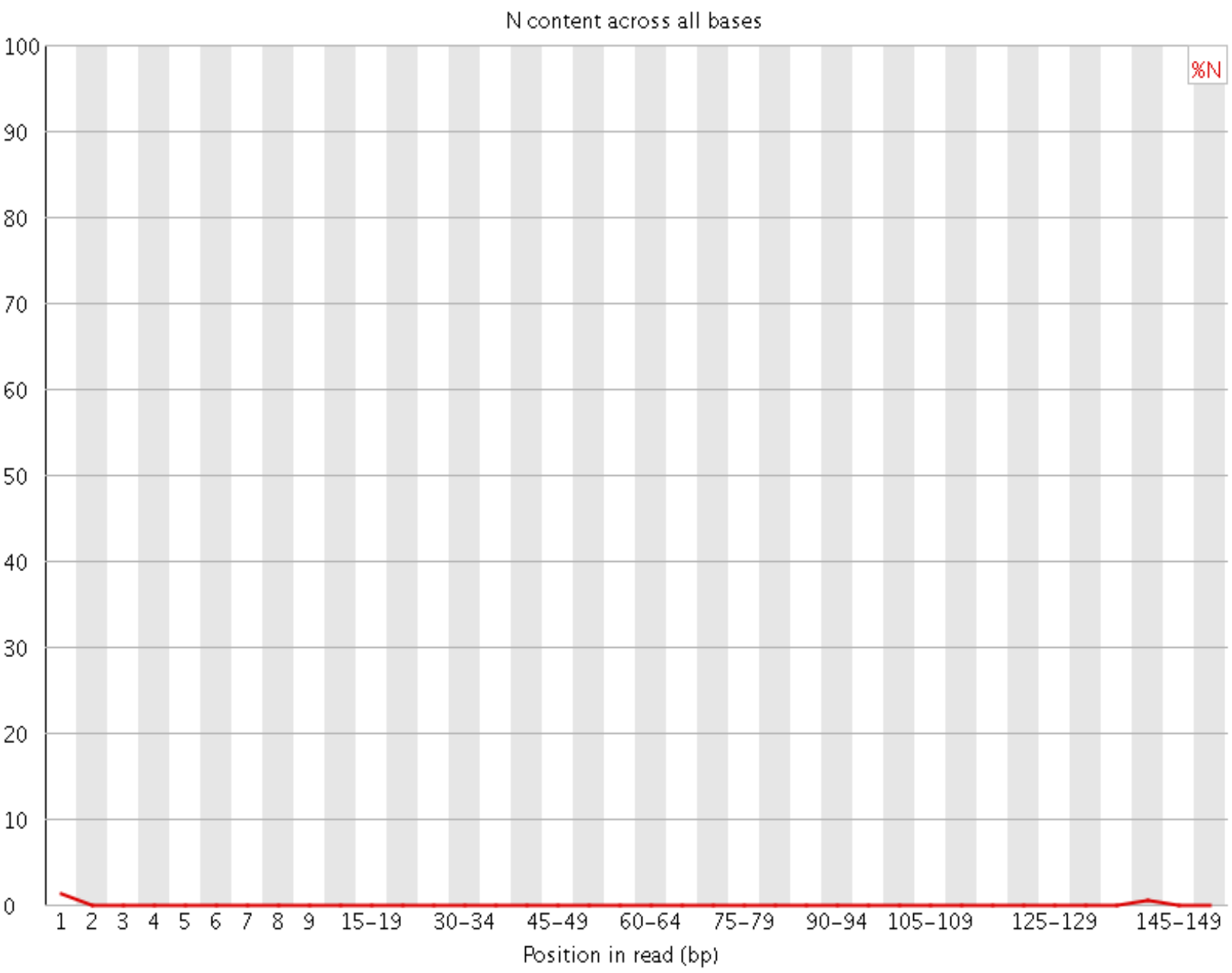
**Per sequence GC content**

GC distribution over all sequences
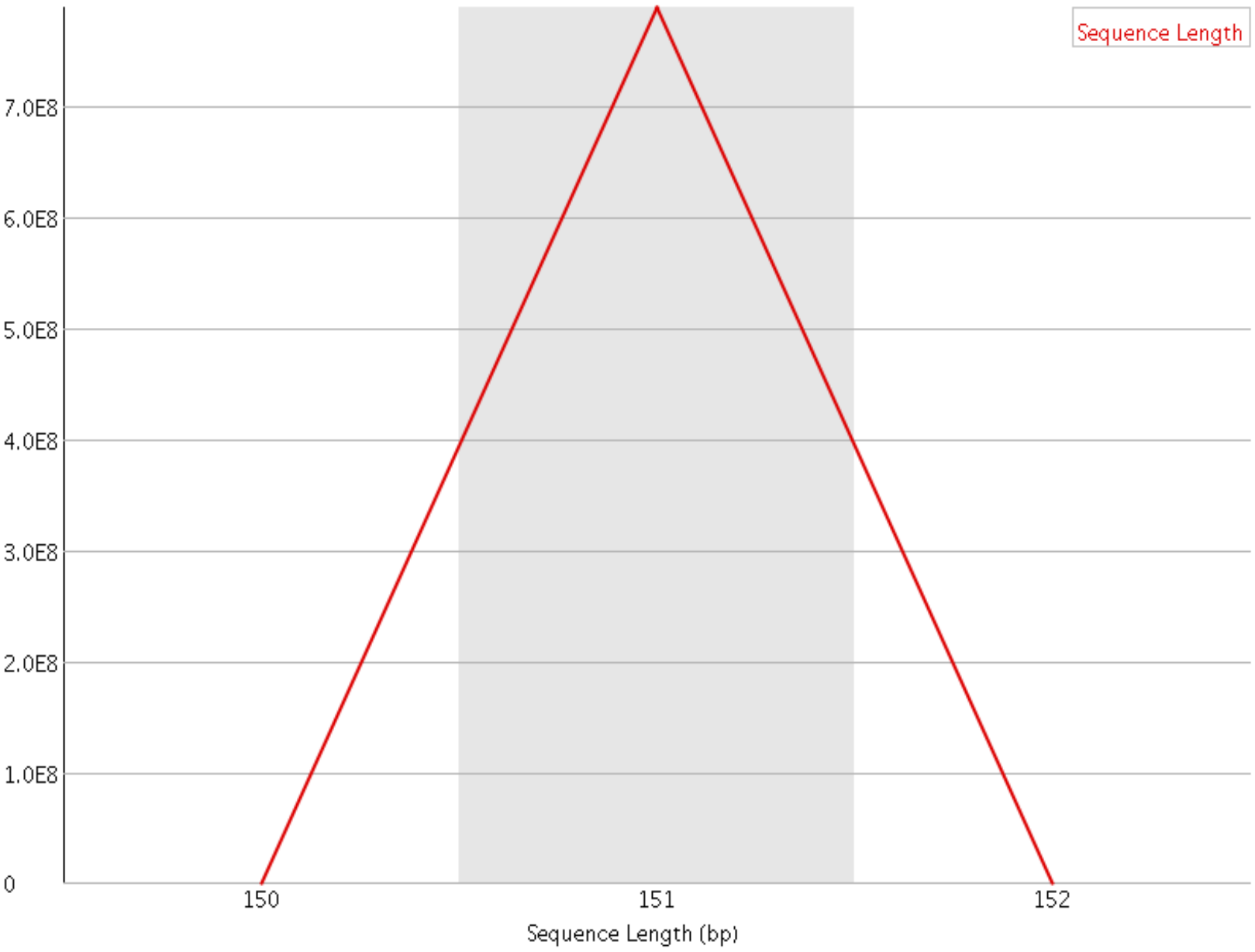
GC count per read
Theoretical Distribution

Mean GC content (%)

✅ **Per base N content**
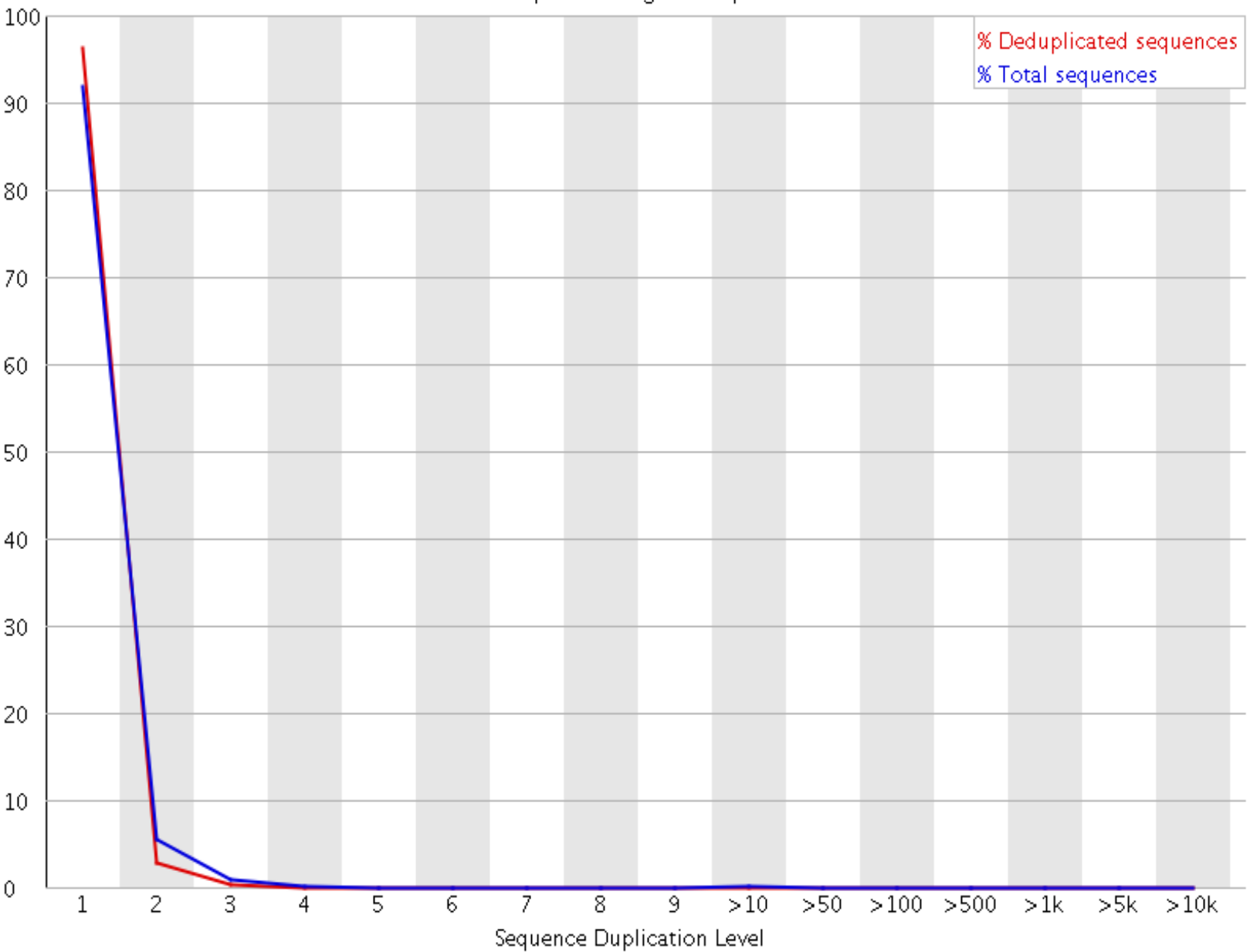
N content across all bases

## ✅ Sequence Length Distribution

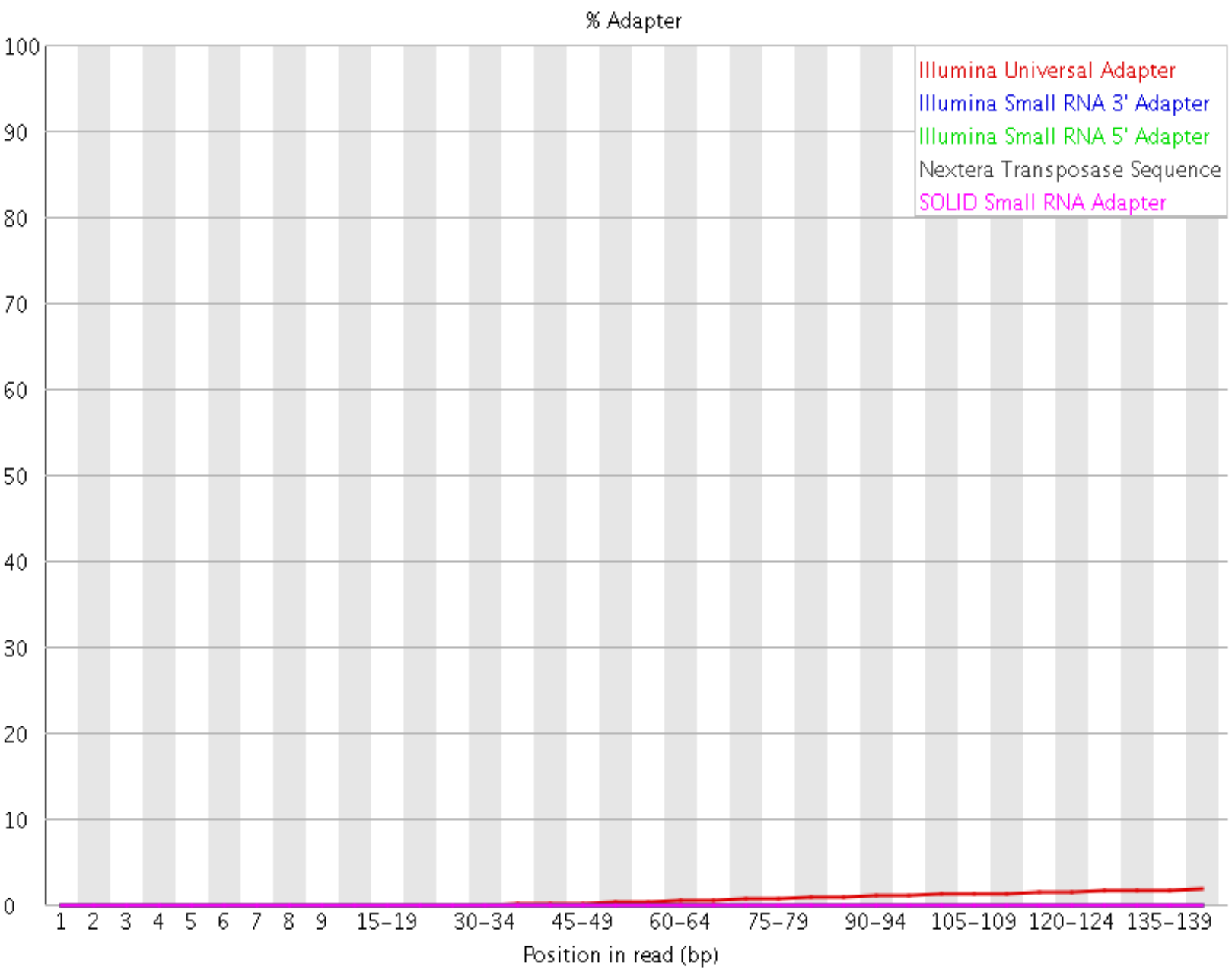Distribution of sequence lengths over all sequences

## Sequence Duplication Levels

Percent of seqs remaining if deduplicated 95.44%

% Deduplicated sequences
% Total sequences

Sequence Duplication Level

✅ **Overrepresented sequences**

No overrepresented sequences

✅ **Adapter Content**

% Adapter

Illumina Universal Adapter
Illumina Small RNA 3' Adapter
Illumina Small RNA 5' Adapter
Nextera Transposase Sequence
SOLID Small RNA Adapter

Position in read (bp)

**Produced by [FastQC](#) (version 0.11.7)**

# FastQC Report

## Summary

✅ [Basic Statistics](#)

⚠️ [Per base sequence quality](#)

✅ [Per tile sequence quality](#)

✅ [Per sequence quality scores](#)

✅ [Per base sequence content](#)

⚠️ [Per sequence GC content](#)

✅ [Per base N content](#)

✅ [Sequence Length Distribution](#)

✅ [Sequence Duplication Levels](#)
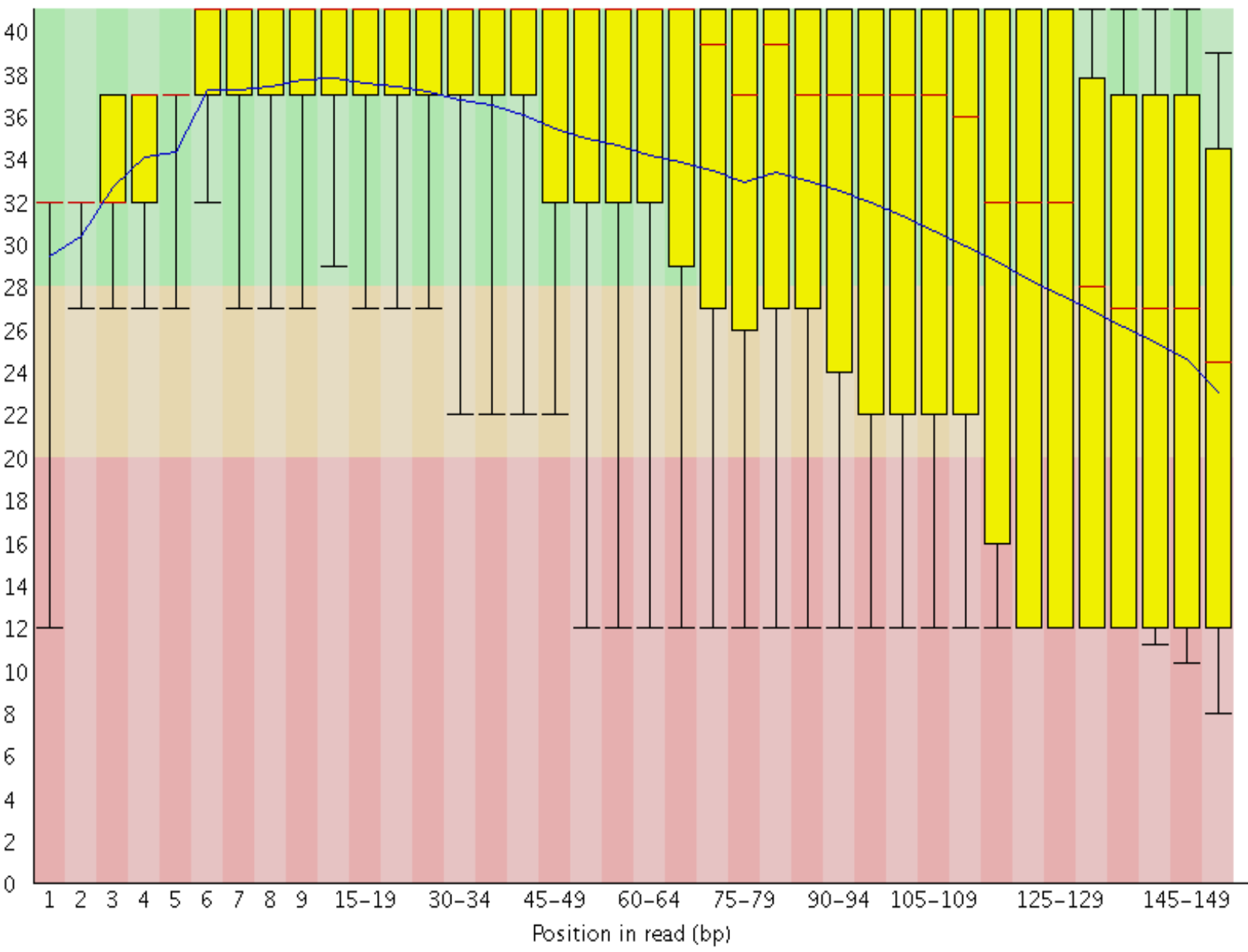
✅ [Overrepresented sequences](#)

✅ [Adapter Content](#)

## ✅ Basic Statistics

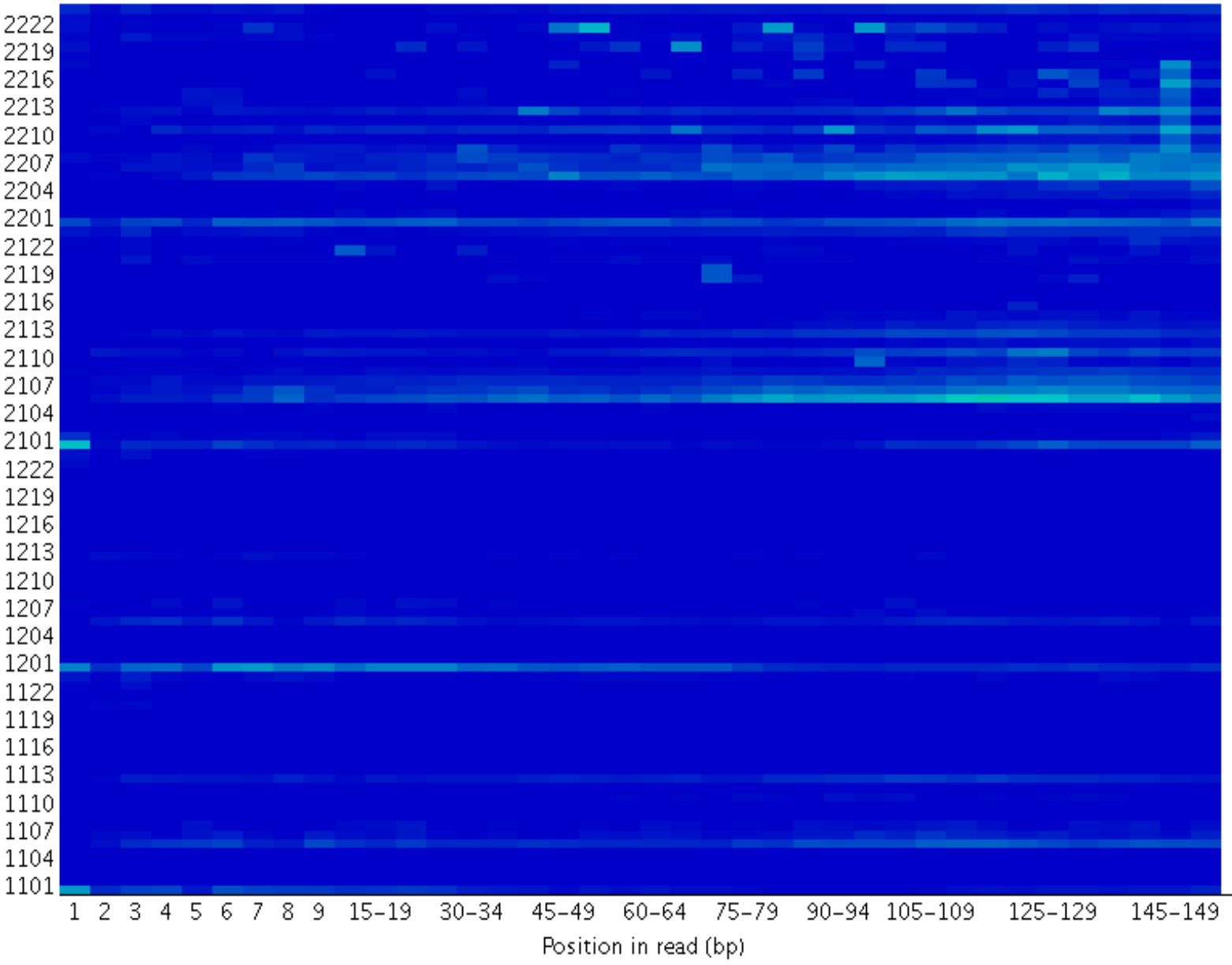| Measure | Value |
|---|---|
| Filename | Venter_S1_merge_R2.fastq.gz |
| File type | Conventional base calls |
| Encoding | Sanger / Illumina 1.9 |
| Total Sequences | 789239544 |
| Sequences flagged as poor quality | 0 |
| Sequence length | 151 |
| %GC | 43 |

## ⚠️ Per base sequence quality

Quality scores across all bases (Sanger / Illumina 1.9 encoding)

## Per tile sequence quality

Quality per tile

Position in read (bp)

✅ **Per sequence quality scores**

Quality score distribution over all sequences

Average Quality per read

**Per base sequence content**

Sequence content across all bases

Position in read (bp)

⚠ **Per sequence GC content**

GC distribution over all sequences

**Per base N content**
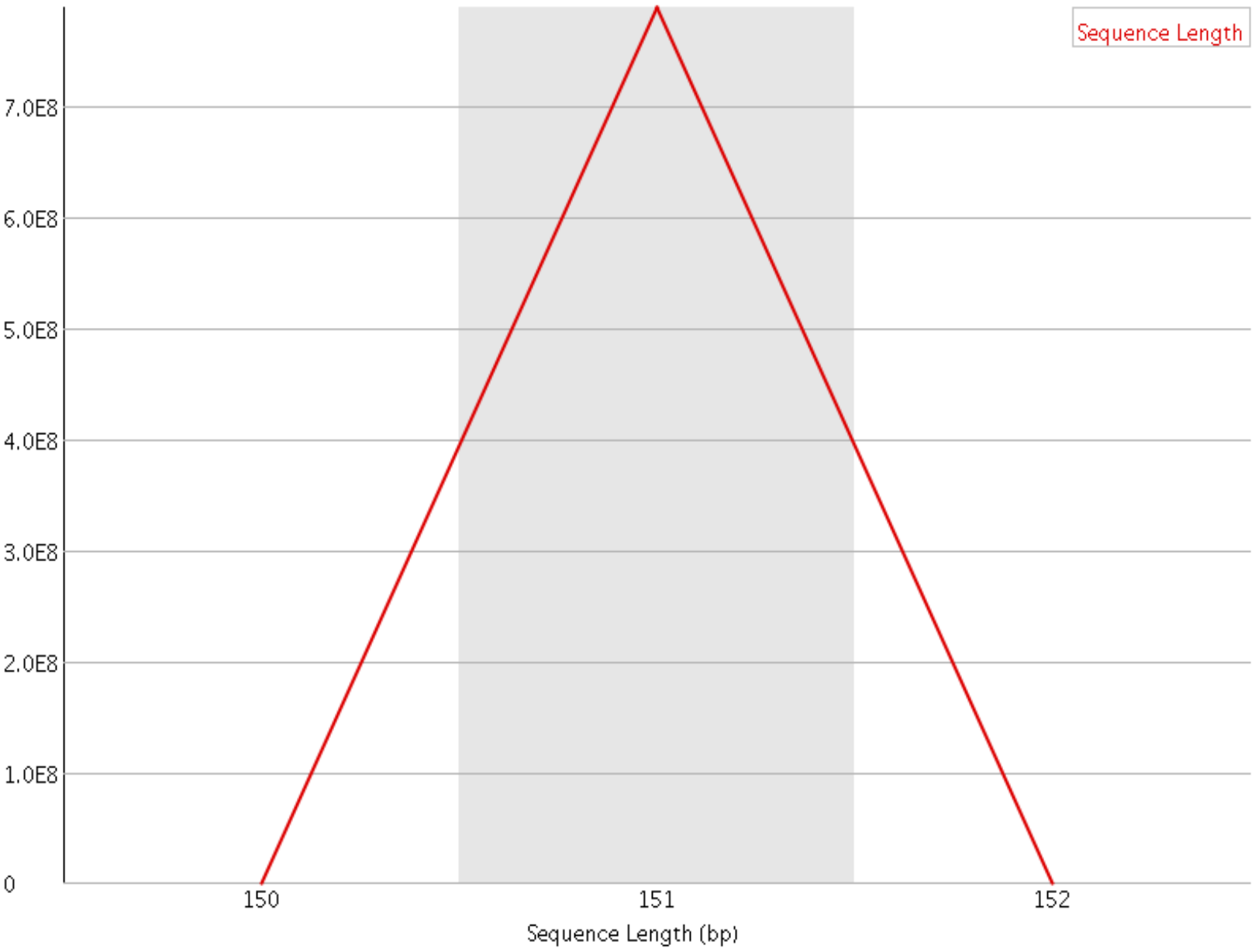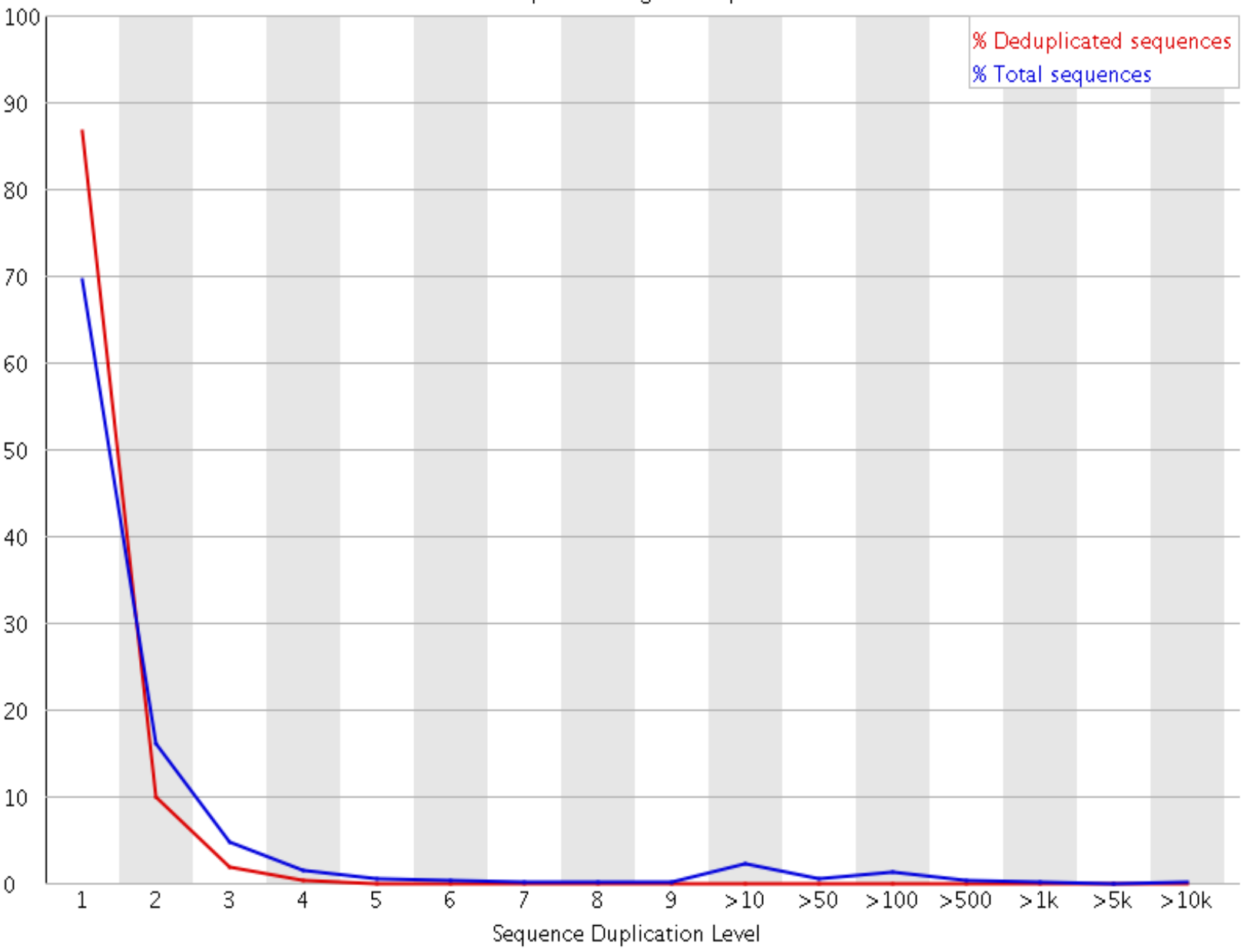
N content across all bases

## Sequence Length Distribution

Distribution of sequence lengths over all sequences

## Sequence Duplication Levels

Percent of seqs remaining if deduplicated 80.29%

![Overrepresented sequences]
## Overrepresented sequences
No overrepresented sequences

![Adapter Content]
## Adapter Content

% Adapter

Illumina Universal Adapter
Illumina Small RNA 3' Adapter
Illumina Small RNA 5' Adapter
Nextera Transposase Sequence
SOLID Small RNA Adapter

Position in read (bp)

**Produced by [FastQC](FastQC) (version 0.11.7)**