

Anexo

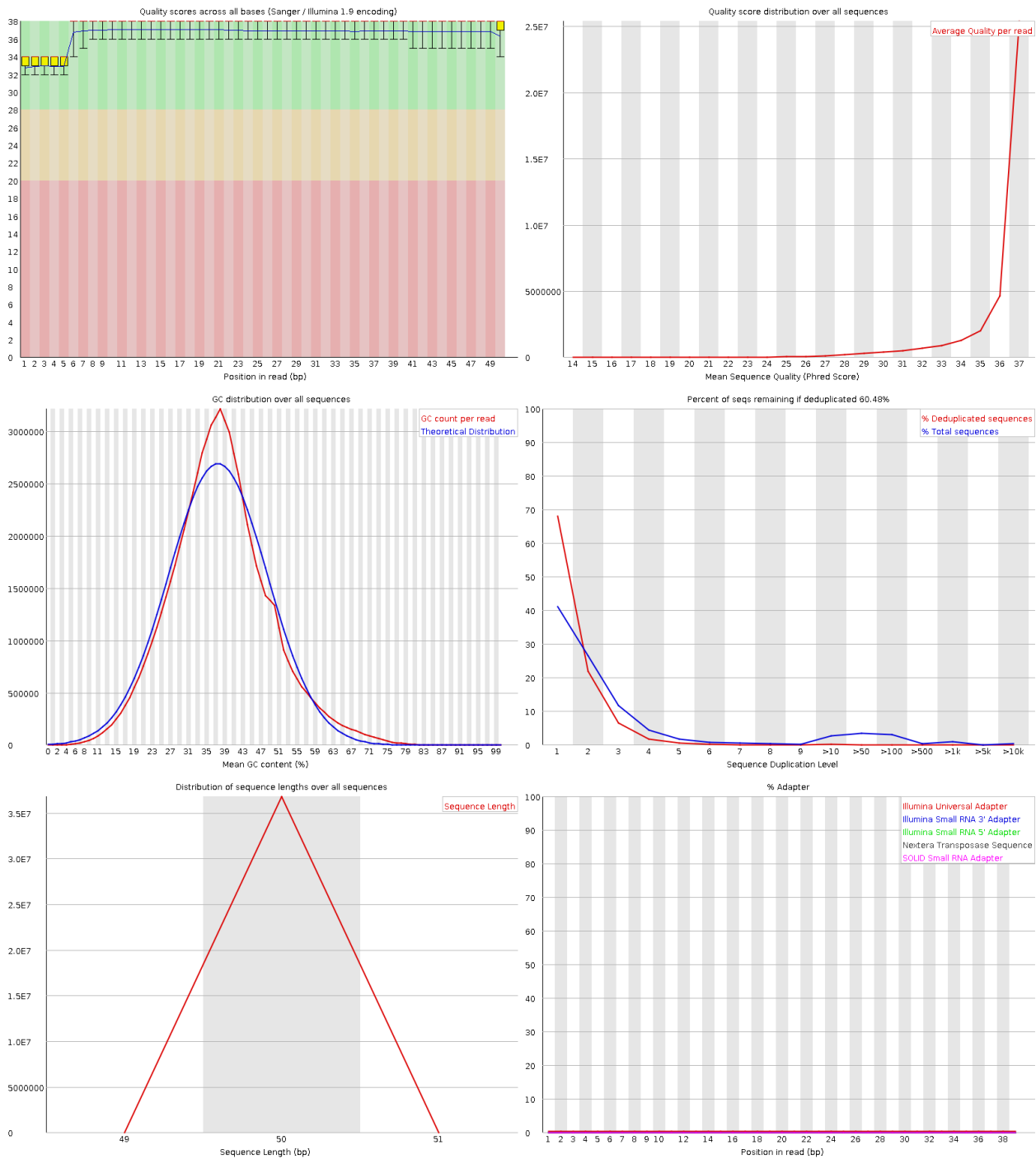


Figura 5.1. Análisis de calidad de las lecturas del control empleando fastQC. Muestra la calidad por base, calidad por secuencia, contenido en GC por secuencia, niveles de duplicación, distribución de la longitud de las lecturas y contenido de adaptadores.

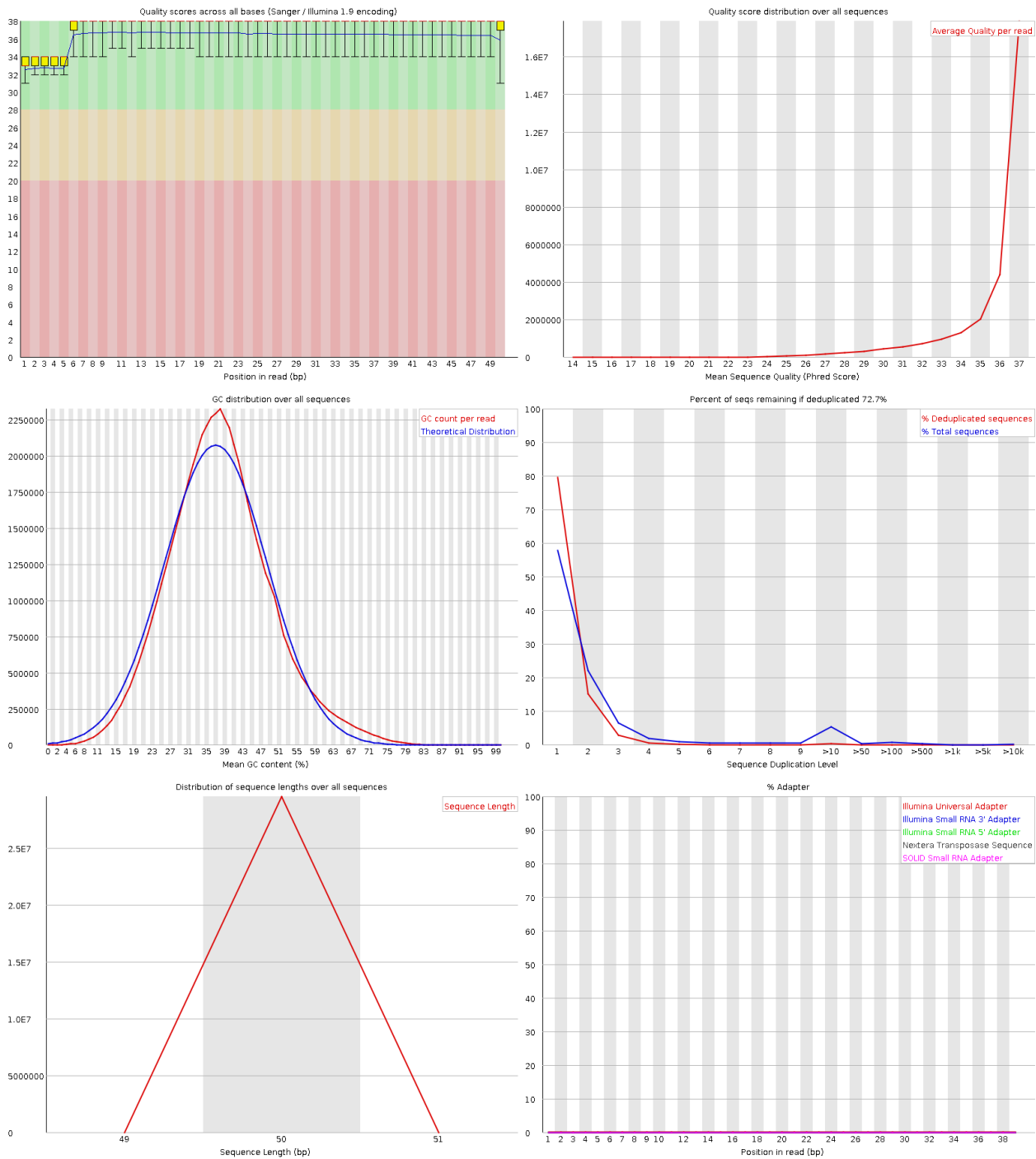


Figura 5.2. Análisis de calidad de las lecturas de la réplica 1 empleando fastQC. Muestra la calidad por base, calidad por secuencia, contenido en GC por secuencia, niveles de duplicación, distribución de la longitud de las lecturas y contenido de adaptadores.

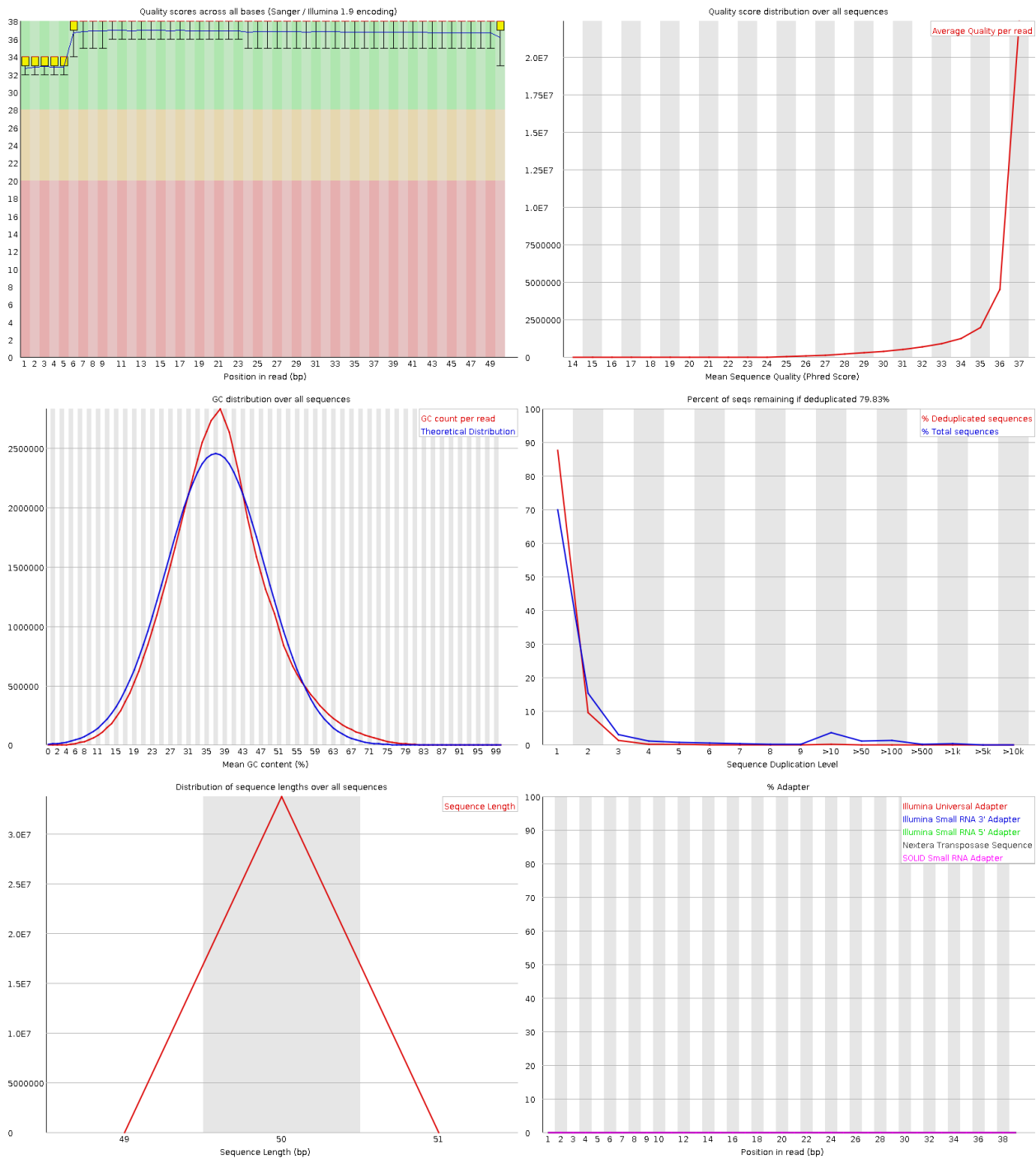


Figura 5.3. Análisis de calidad de las lecturas de la réplica 2 empleando fastQC. Muestra la calidad por base, calidad por secuencia, contenido en GC por secuencia, niveles de duplicación, distribución de la longitud de las lecturas y contenido de adaptadores.

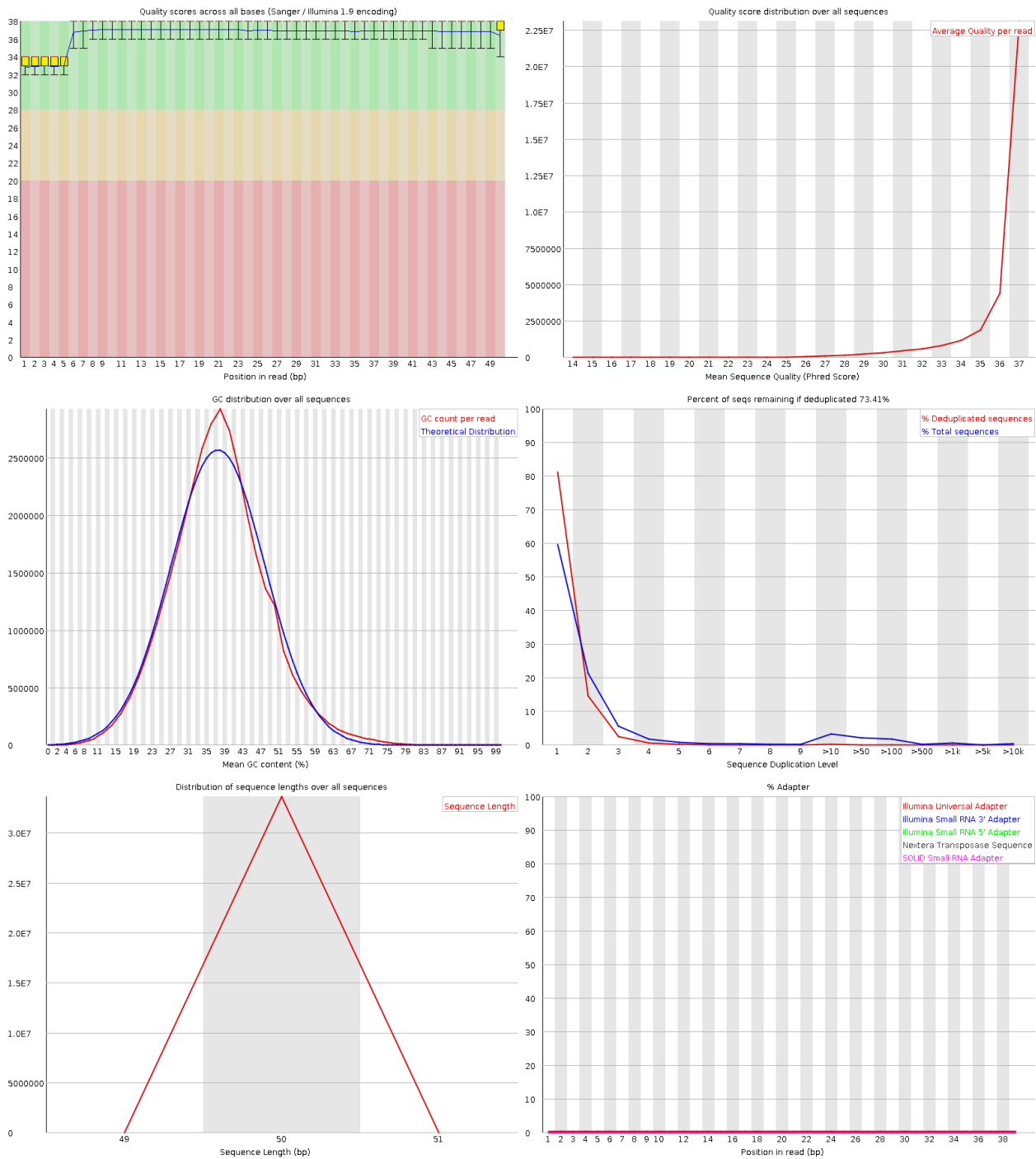


Figura 5.4. Análisis de calidad de las lecturas de la réplica 3 empleando fastQC. Muestra la calidad por base, calidad por secuencia, contenido en GC por secuencia, niveles de duplicación, distribución de la longitud de las lecturas y contenido de adaptadores.

```

1 #!/bin/bash
2
3 # Este script sirve para procesar lecturas de ChIP-seq
4 # desde el control de calidad hasta la llamada de picos
5
6 # Parametros
7 fastq_dir="fastq"
8 data_tsv="data.tsv"
9 control="ctrl_rep1"
10
11 genome_fasta="TAIR10.fa"
12 genome_size=119481543
13 threads=4
14
15 # 1. FastQC
16 mkdir -p results/fastqc logs/fastqc/
17
18 cat $data | tail -n +2 | while read name group rep fastq; do
19     fastqc \
20     -o results/fastqc/ \
21     $fastq_dir/$fastq \
22     &&> logs/fastqc/$name.log
23 done
24
25 # 2. Bowtie2
26 mkdir -p results/bamfiles logs/bowtie2_index logs/bowtie2
27
28 bowtie2-build \
29     $genome_fasta \
30     bowtie2_index &&> logs/bowtie2_index.log
31
32 cat $data | tail -n +2 | while read name group rep fastq; do
33     (bowtie2 -p 4 \
34     -x bowtie2_index \
35     -U $fastq_dir/$fastq \
36     -p $threads \
37     | samtools view -bh - \
38     | samtools sort -o results/bamfiles/$name.sorted.bam -) \
39     >> logs/bowtie2/$name.log
40 done
41
42 # 3. Picard y samtools
43 mkdir -p logs/picard
44
45 cat $data | tail -n +2 | cut -f1 | while read name; do
46     picard MarkDuplicates \
47     I=results/bamfiles/$name.sorted.bam \
48     O=results/bamfiles/$name.markdup.bam \
49     H=results/bamfiles/$name.picard.txt \
50     &&> logs/picard/$name.log
51 done
52
53 mkdir -p logs/samtools_view
54
55 cat $data | tail -n +2 | cut -f1 | while read name; do
56     samtools view -bh \
57     -F 1024 -F 4 -q 5 \
58     -@ $threads \
59     results/bamfiles/$name.markdup.bam \
60     >> results/bamfiles/$name.clean.bam \
61     >> logs/samtools_view/$name.log
62 done
63
64 cat $data | tail -n +2 | cut -f1 | while read name; do
65     samtools index results/bamfiles/$name.clean.bam
66 done
67
68 # 4. bamCoverage y bigwigCompare
69 mkdir -p results/bigwigs logs/bamCoverage
70
71 cat $data | tail -n +2 | cut -f1 | while read name; do
72     bamCoverage \
73     -b results/bamfiles/$name.clean.bam \
74     --normalizeUsing CPM \
75     --extendReads 200 \
76     -o results/bigwigs/$name.cpm.bw \
77     -p $threads \
78     &&> logs/bamCoverage/$name.cpm.log
79 done
80
81 mkdir -p logs/bigwigCompare
82
83 cat $data | tail -n +2 | cut -f1 | grep -v $control | while read name; do
84     bigwigCompare \
85     -t results/bigwigs/$name.cpm.bw \
86     -b2 results/bigwigs/$control.cpm.bw \
87     -o results/bigwigs/${name}_over_${control}.log2ratio.bw \
88     &&> logs/bigwigCompare/${name}_over_${control}.log2ratio.log
89 done
90
91 # 5. MACS2
92 mkdir -p results/peak_calling logs/mac2
93
94 cat $data | tail -n +2 | cut -f1 | grep -v $control | while read name; do
95     macs2 callpeak \
96     -t results/bamfiles/$name.clean.bam \
97     -c results/bamfiles/$control.clean.bam \
98     -n $name \
99     -g 119481543 \
100     --keep-dup all \
101     --outdir results/peak_calling \
102     --nomodel --extsize 200 \
103     &&> logs/mac2/$name.log
104 done

```

Figura 5.5. Script empleado para el procesamiento de las lecturas del ChIP-seq desde el control de calidad hasta la llamada de picos.

```

1 #!/bin/bash
2
3 # Este script lleva a cabo la interseccion entre picos de varias muestras,
4 # su anotacion, filtrado y la busqueda de motivos enriquecidos
5
6 mkdir -p results/peak_analysis
7
8 # Interseccion entre picos
9 cat results/peak_calling/*_peaks.narrowPeak | grep -v pooled \
10 | sortBed \
11 | mergeBed -i - -c 4 -o collapse \
12 | awk -v OFS="\t" '{ print $1, $2, $3, "rg_peak_"NR, $4 }' \
13 > results/peak_analysis/rg_peaks.bed
14
15 # Summits de la interseccion entre picos
16 bedtools intersect \
17 -a results/peak_analysis/rg_peaks.bed \
18 -b results/peak_calling/rg_pooled_summits.bed \
19 > results/peak_analysis/rg_summits.bed
20
21 # Filtrar los comunes a 2/3 replicas
22 python3 process_peaks.py
23 python3 filter_summits.py
24
25 # Anotar los picos con HOMER
26 mkdir -p logs/homer
27
28 annotatePeaks.pl rg_peaks_23.bed tair10 \
29 > results/peak_analysis/rg_peaks_23.annotation.tsv \
30 &&> logs/homer/annotatePeaks.log
31
32 # Buscar motivos de ADN en los summits
33 findMotifsGenome.pl rg_summits_23.bed tair10 \
34 results/peak_analysis/motifs \
35 &&> logs/homer/findMotifsGenome.log
36

```

Figura 5.6. Script empleado para el análisis de los picos. Lleva a cabo la intersección entre picos de varias muestras, su anotación, filtrado y la búsqueda de motivos enriquecidos.

```

1  #!/usr/bin/env python3
2
3  """
4  Este script sirve para procesar la interseccion entre picos.
5  La salida de bedtools merge se transforma en un fichero TSV
6  con el numero de picos, muestras y una columna para cada muestra
7  indicando si contiene el pico (True) o no (False)
8  """
9
10 def main():
11     in_file = open("results/peak_analysis/rg_peaks.bed", "r")
12     out_file = open("results/peak_analysis/rg_peaks.tsv", "w")
13     sample_names = ["rg_rep1", "rg_rep2", "rg_rep3"]
14
15     header = "\t".join(
16         ["peak_id", "peaks", "num_peaks", "num_samples"] + sample_names
17     )
18     out_file.write(header + "\n")
19
20     for line in in_file:
21         fields = line.strip().split("\t")
22
23         peaks = fields[4]
24         samples = [sample in peaks for sample in sample_names]
25         num_peaks = len(peaks.split(","))
26         num_samples = sum(samples)
27
28         new_line = "\t".join(
29             fields[3:5]
30             + [str(num_peaks), str(num_samples)]
31             + [str(sample) for sample in samples]
32         )
33         out_file.write(new_line + "\n")
34
35     in_file.close()
36     out_file.close()
37
38
39 if __name__ == "__main__":
40     main()
41

```

Figura 5.7. Script empleado para procesar la intersección entre picos. La salida de bedtools merge se transforma en un fichero TSV con el numero de picos, muestras y una columna para cada muestra indicando si contiene el pico (True) o no (False).

```

1  #!/usr/bin/env python3
2
3  """Este programa convierte ficheros GTF a BED con las coordenadas de cada gen"""
4
5  import sys
6
7  def main():
8      if len(sys.argv) < 3:
9          print("gtf2bed.py <gtf in> <bed out>")
10         exit()
11
12         input_gtf = sys.argv[1]
13         output_bed = sys.argv[2]
14
15         gtf = open(input_gtf, "r")
16         bed = open(output_bed, "w")
17
18         for line in gtf:
19             fields = line.strip().split("\t")
20
21             if fields[2] == "gene":
22                 info = fields[8].split(";")
23                 gene = info[1].strip().split(" ")[1].strip("")
24
25                 line_bed = "\t".join(
26                     [fields[0], fields[3], fields[4], gene, fields[5], fields[6]]
27                 )
28                 bed.write(line_bed + "\n")
29
30         gtf.close()
31         bed.close()
32
33 if __name__ == "__main__":
34     main()
35
36

```

Figura 5.8. Script empleado para convertir ficheros GTF en ficheros BED con las coordenadas de cada gen.

```

1  #!/usr/bin/env python3
2
3  """
4  Este script selecciona los picos presentes en al menos 2 replicas
5  y escribe nuevos ficheros con las regiones y los summits de estos
6  """
7
8  def filter_and_write(input, output, peaks):
9
10     input_file = open(input, "r")
11     output_file = open(output, "w")
12
13     for line in input_file:
14         fields = line.strip().split("\t")
15
16         if fields[3] in peaks:
17             output_file.write(line)
18
19     input_file.close()
20     output_file.close()
21
22
23 def main():
24     num_reps = 2
25
26     selected_peaks = set()
27
28     with open("results/peak_analysis/rg_peaks.tsv") as file:
29         for line in file:
30             fields = line.strip().split("\t")
31
32             if fields[3] >= num_reps:
33                 selected_peaks.add(fields[0])
34
35     filter_and_write(
36         input="results/peak_analysis/rg_peaks.bed",
37         output="results/peak_analysis/rg_peaks_23.bed",
38         peaks=selected_peaks
39     )
40
41     filter_and_write(
42         input="results/peak_analysis/rg_summits.bed",
43         output="results/peak_analysis/rg_summits_23.bed",
44         peaks=selected_peaks
45     )
46
47
48
49 if __name__ == "__main__":
50     main()
51

```

Figura 5.9. Script empleado para seleccionar los picos presentes en al menos 2 de las 3 réplicas del ChIP-seq y escribir nuevos ficheros con las regiones y los summits de estos.


```

1  import sys      #rg_peaks.tsv
2
3  import matplotlib.pyplot as plt
4  from matplotlib_venn import venn3
5
6  input_file = sys.argv[1]
7
8  peaks_rep1 = set()
9  peaks_rep2 = set()
10 peaks_rep3 = set()
11
12 peaks_rep1_file = open("peaks_rep1.txt", "w")
13 peaks_rep2_file = open("peaks_rep2.txt", "w")
14 peaks_rep3_file = open("peaks_rep3.txt", "w")
15
16 with open(input_file) as f:
17     f.readline()
18     for line in f:
19
20         fields = line.strip().split("\t")
21         name = fields[0]
22         rep1, rep2, rep3 = fields[2:5]
23
24
25         if rep1 == "True":
26             peaks_rep1.add(name)
27             peaks_rep1_file.write(name + "\n")
28
29         if rep2 == 'True':
30             peaks_rep2.add(name)
31             peaks_rep2_file.write(name + "\n")
32
33         if rep3 == 'True':
34             peaks_rep3.add(name)
35             peaks_rep3_file.write(name + "\n")
36
37
38 venn3(
39     [peaks_rep1, peaks_rep2, peaks_rep3],
40     set_labels=["RGA rep1", "RGA rep2", "RGA rep3"],
41
42     plt.show()
43

```

Figura 5.10. Script empleado para la representación de un diagrama de Venn que muestra la intersección entre los picos de las diferentes réplicas del ChIP-seq.

```

1  tracks=$(ls *log2cpm.bw)
2
3
4
5  computeMatrix scale-regions \
6  -R genes.bed \
7  -S $tracks \
8  -a 2000 -b 2000 -m 2000 \
9  --missingDataAsZero --skipZeros \
10 -o matrix_2000.gz \
11 -p 3
12

```

Figura 5.11. Script empleado para la representación de un heatmap que muestra la distribución genómica de RGA en los genes diana.

```

1  # -*- coding: utf-8 -*-
2  """
3  Spyder Editor
4
5  This is a temporary script file.
6  """
7  #import numpy as np
8
9  import matplotlib.pyplot as plt
10
11  ann = {
12      "promoter-TSS":0,
13      "TTS":0,
14      "Intergenic":0,
15      "exon":0,
16      "intron":0
17  }
18
19
20  with open('rg_peaks_23.annotation.tsv') as f:
21      f.readline()
22      for line in f:
23          fields = line.strip().split('\t')
24          info = fields[7].split(" ")[0]
25
26          ann[info] += 1
27
28
29  sns.palplot (sns.color_palette ("hls", 9))
30
31  plt.pie(
32      x=ann.values(),
33      labels=ann.keys(),
34      autopct="%.2f%%",
35
36
37  plt.show()
38

```

Figura 5.12. Script empleado para anotación de los picos en las diferentes categorías genómicas y su representación.