

RiboSeq

Лев Мазаев | МАДБМ18

Считал все на локальной машине.

Этап 1: скачивание и распаковка файлов образцов

Первая задача - скачать SRA-файлы. Я решил выбрать образцы с нокдауном метилтрансферазы METTL3. На GEO-странице эксперимента внизу можно перейти в SRA Run Selector:

Platforms (1) [GPL9185](#) Illumina Genome Analyzer (Mus musculus)

Samples (14) [GSM2717358](#) Ribo_seq_Scram_rep1
[GSM2717359](#) Ribo_seq_Scram_rep2
[GSM2717360](#) Ribo_seq_shMETTL3_rep1
[GSM2717361](#) Ribo_seq_shMETTL3_rep2
[GSM2717362](#) Ribo_seq_shABCF1_rep1
[GSM2717363](#) Ribo_seq_shABCF1_rep2
[GSM2717364](#) Ribo_seq_Torin
[GSM2717365](#) RNA_seq_Scram_rep1
[GSM2717366](#) RNA_seq_Scram_rep2
[GSM2717367](#) RNA_seq_shMETTL3_rep1
[GSM2717368](#) RNA_seq_shMETTL3_rep2
[GSM2717369](#) RNA_seq_shABCF1_rep1
[GSM2717370](#) RNA_seq_shABCF1_rep2
[GSM2717371](#) RNA_seq_Torin

Relations

BioProject [PRJNA395723](#)
SRA [SRP113540](#)

Download family

[SOFT formatted family file\(s\)](#)
[MINiML formatted family file\(s\)](#)
[Series Matrix File\(s\)](#)

Format

[SOFT](#) [?](#)
[MINiML](#) [?](#)
[TXT](#) [?](#)

Supplementary file	Size	Download	File type/resource
GSE101865_RAW.tar	2.1 Mb	(http) (custom)	TAR (of TXT)
SRA Run Selector ?			

Raw data are available in SRA

Processed data provided as supplementary file

Там мы видим все образцы (Samples) и их соответствующие SRR-записи:

Скачиваем SRR586579(4|5) (GSM271736(0|1) - Ribo_seq_shMETTL3_rep(1|2)) и SRR586580(1|2) (GSM271736(7|8) - RNA_seq_shMETTL3_rep(1|2)):

Далее, в процессе скачивания посмотрим, что делают команды `research` и `efetch`:

Из них при помощи ruby-скрипта можно получить табличку для выяснения номеров SRA, но не стал этого делать, так как всё есть в SRA Run Selector.

```
parallel-fastq-dump -s *.sra -t 8 --split-files
```

Сравнение размеров файлов:

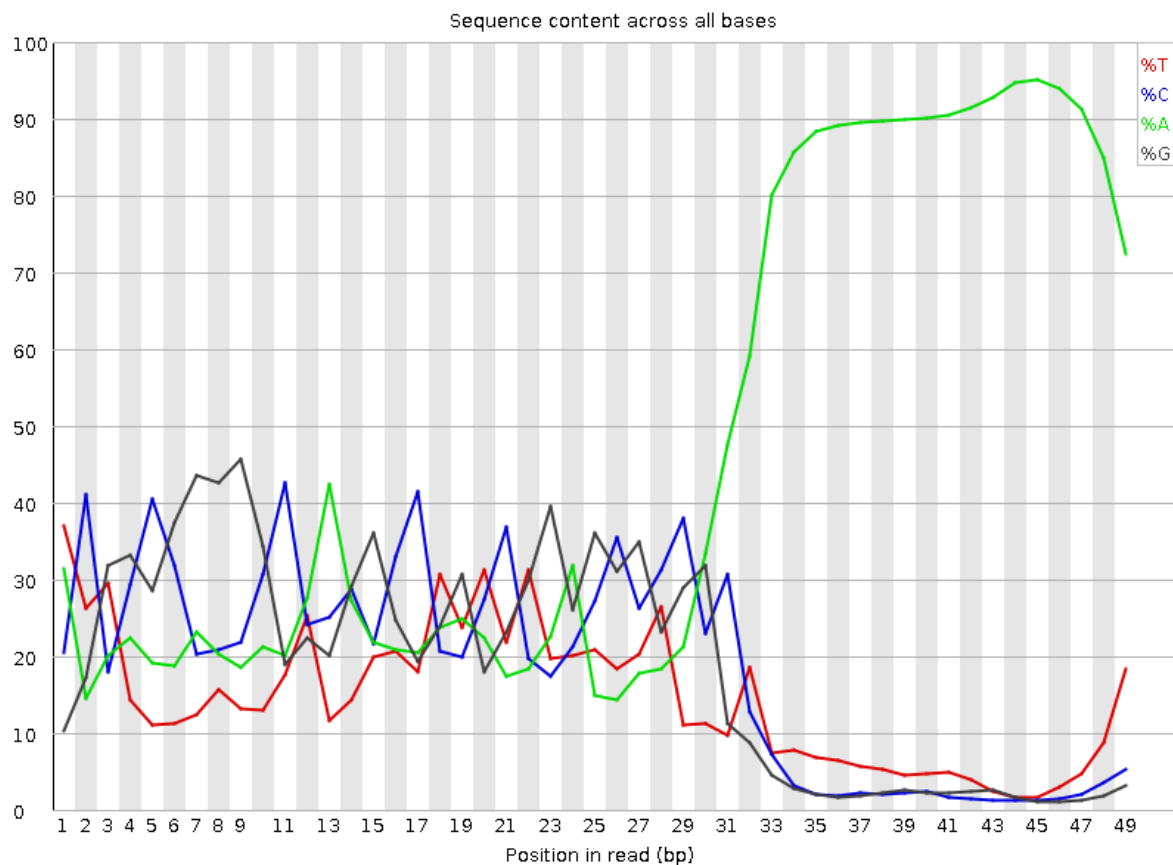
```
(riboseq) leo@MS7922:~/BioData/RiboSeq$ du -ah
12G      ./SRR5865794_1.fastq
586M     ./SRR5865802.sra
12G      ./SRR5865801_1.fastq
9,7G     ./SRR5865795_1.fastq
1,4G     ./SRR5865801.sra
1022M    ./SRR5865795.sra
1,4G     ./SRR5865794.sra
6,0G     ./SRR5865802_1.fastq
43G      .
```

Этап 2: отрезка адаптеров и QC

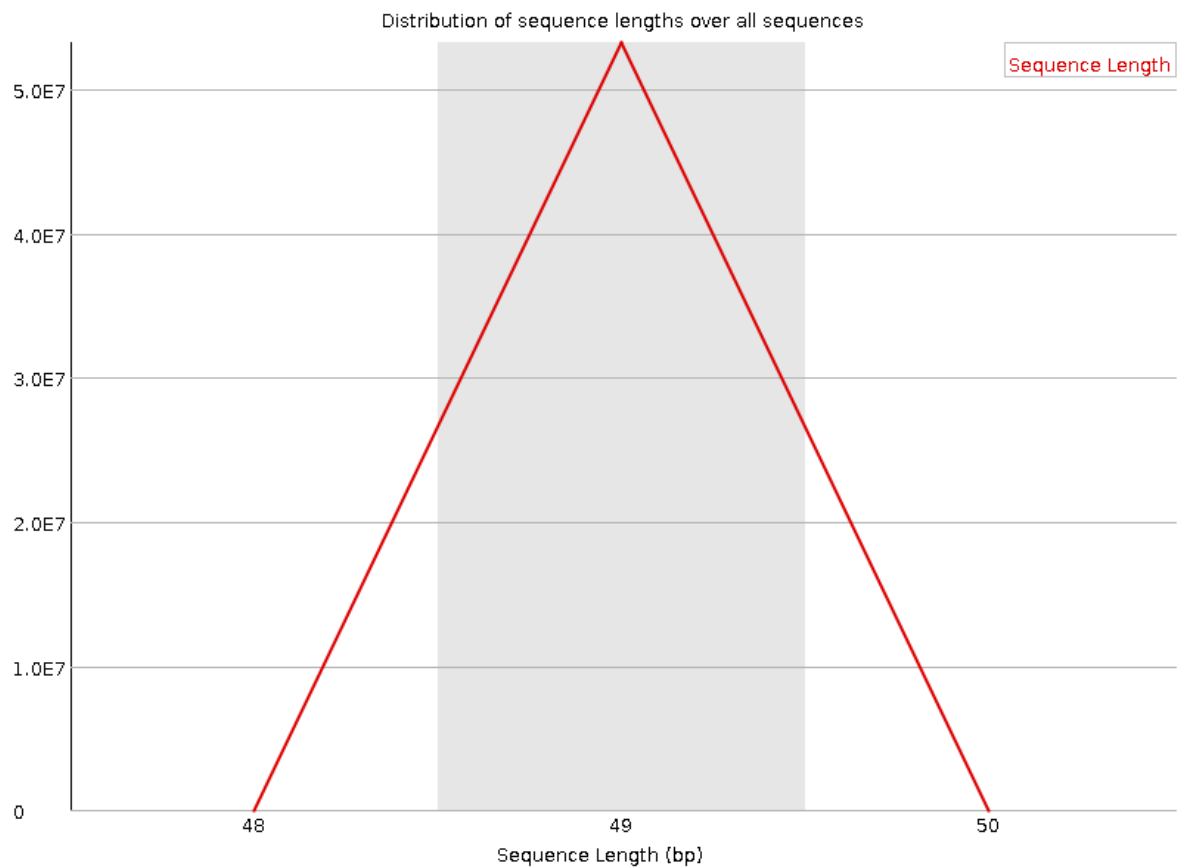
Первичный анализ последовательностей:

```
fastqc -t 8 ~/BioData/RiboSeq/*.fastq
```

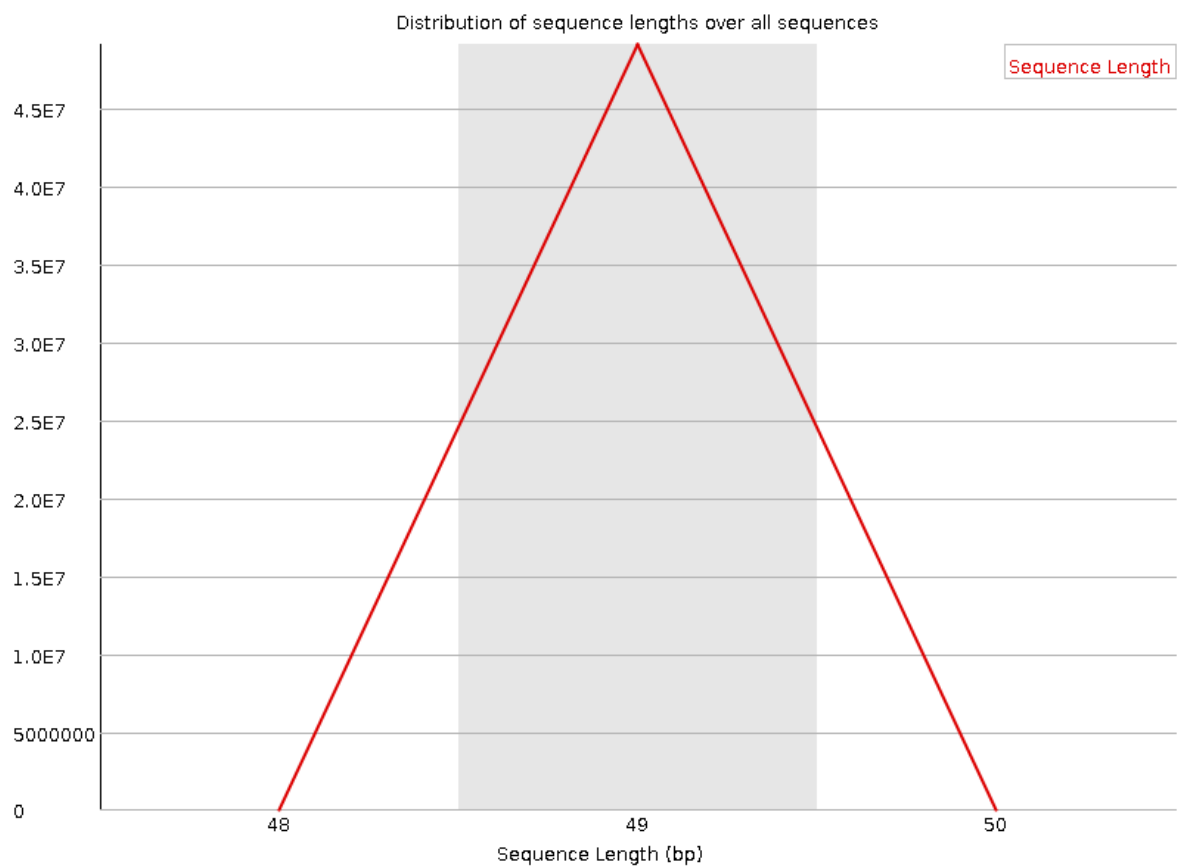
Длина ряда в первых репликах - 49, во вторых - 51. Заметен полиА-адаптер в первых репликах (рис. для RiboSeq rep1):



Распределение длин прочтений (RiboSeq rep1):



RNASeq rep1:



Для вторых реплик тоже самое с пиком на 51. Теперь произведём отрезку адаптеров.

Data processing For replicate 1, the 3' adaptor AAAAAAAAAA and low quality bases were trimmed by Cutadapt. The trimmed reads were mapped to Mouse transcriptome (GRCm38.83) by Bowtie, using parameters: -a --best -m1 --strata. Two mismatches were permitted.

For replicate 2, the 3' adaptor CTGTAGGCACCATCAAT were trimmed by Cutadapt. The first 6 nucleotide and the last 4 nucleotides were also clipped after adaptor removal. The trimmed reads were mapped to Mouse transcriptome (GRCm38.83) by Bowtie, using parameters: -a --best -m1 --strata. Two mismatches were permitted. The PCR duplicates were removed from the mapped reads using the methods in Lecanda et al., 2016 (PMID: 27450428)

Genome_build: GRCm38.83

Supplementary_files_format_and_content: RPKM files were generated. Each line represents one mRNAs. RPKM values were separated by Tabs.

Адаптер первой реплики - AAAAAAAAAA:

```
cutadapt -a AAAAAAAAAA -q 20 --minimum-length 20 -j 8 -o  
riboseq1.trimmed.fastq.gz --trimmed-only SRR5865794_1.fastq  
cutadapt -a AAAAAAAAAA -q 20 --minimum-length 20 -j 8 -o  
rnaseq1.trimmed.fastq.gz SRR5865801_1.fastq
```

Адаптер второй реплики - CTGTAGGCACCATCAAT (minimum length 30, так как потом будем ещё отрезать баркоды):

```
cutadapt -a CTGTAGGCACCATCAAT -q 20 --minimum-length 30 -j 8 -o  
riboseq2.trimmed.fastq.gz --trimmed-only SRR5865795_1.fastq  
cutadapt -a CTGTAGGCACCATCAAT -q 20 --minimum-length 30 -j 8 -o  
rnaseq2.trimmed.fastq.gz SRR5865802_1.fastq
```

Далее, нужно дедуплицировать вторые реплики:

```
seqkit rmdup -s -j 8 riboseq2.trimmed.fastq.gz | gzip -c >  
riboseq2.dedup.fastq.gz
```

[INFO] 18088740 duplicated records removed. То есть 46% записей было удалено из исходных 38801308 в riboseq2.trimmed.fastq.gz.

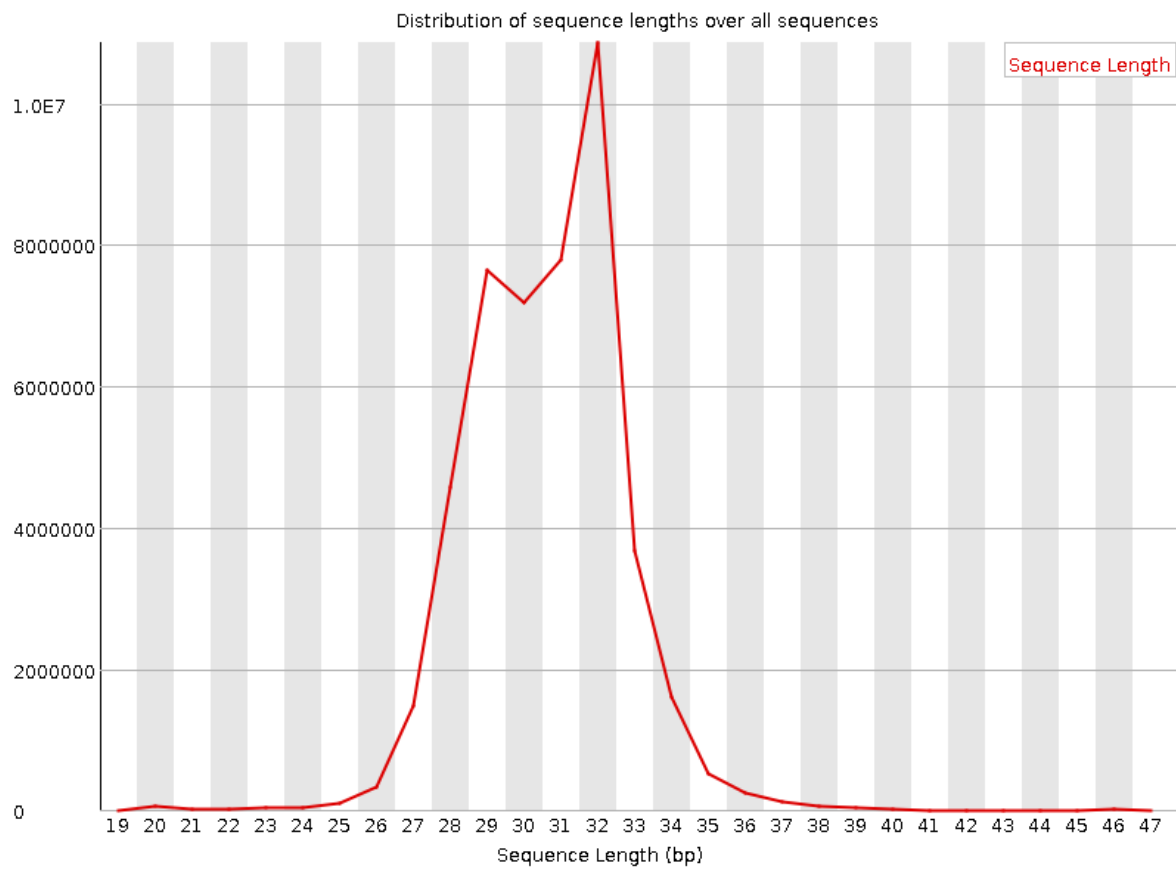
```
seqkit rmdup -s -j 8 rnaseq2.trimmed.fastq.gz | gzip -c >  
rnaseq2.dedup.fastq.gz
```

[INFO] 9371952 duplicated records removed. То есть 36% записей было удалено из исходных 25769506 в rnaseq2.trimmed.fastq.gz

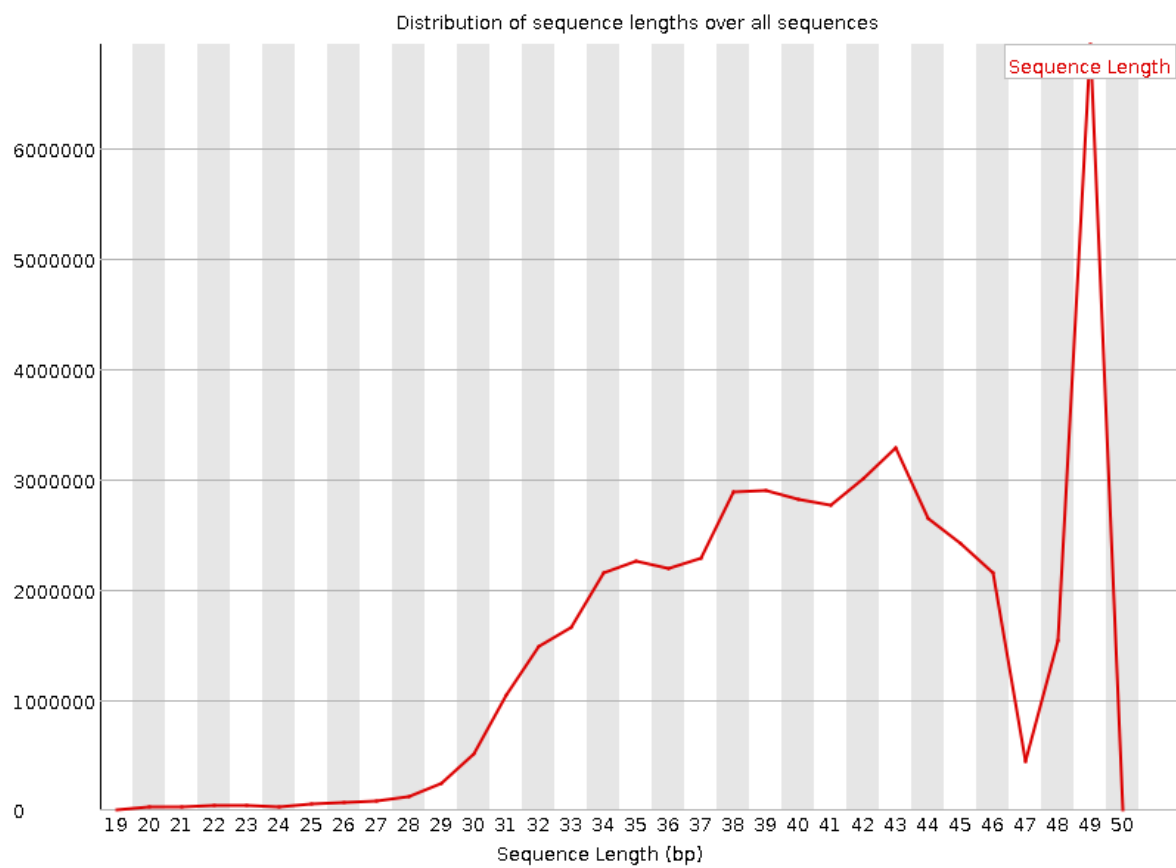
Теперь удалим баркоды:

```
cutadapt -u 6 -j 8 riboseq2.dedup.fastq.gz | cutadapt -u -4 -j 8 - -o  
riboseq2.final.fastq.gz  
cutadapt -u 6 -j 8 rnaseq2.dedup.fastq.gz | cutadapt -u -4 -j 8 - -o  
rnaseq2.final.fastq.gz
```

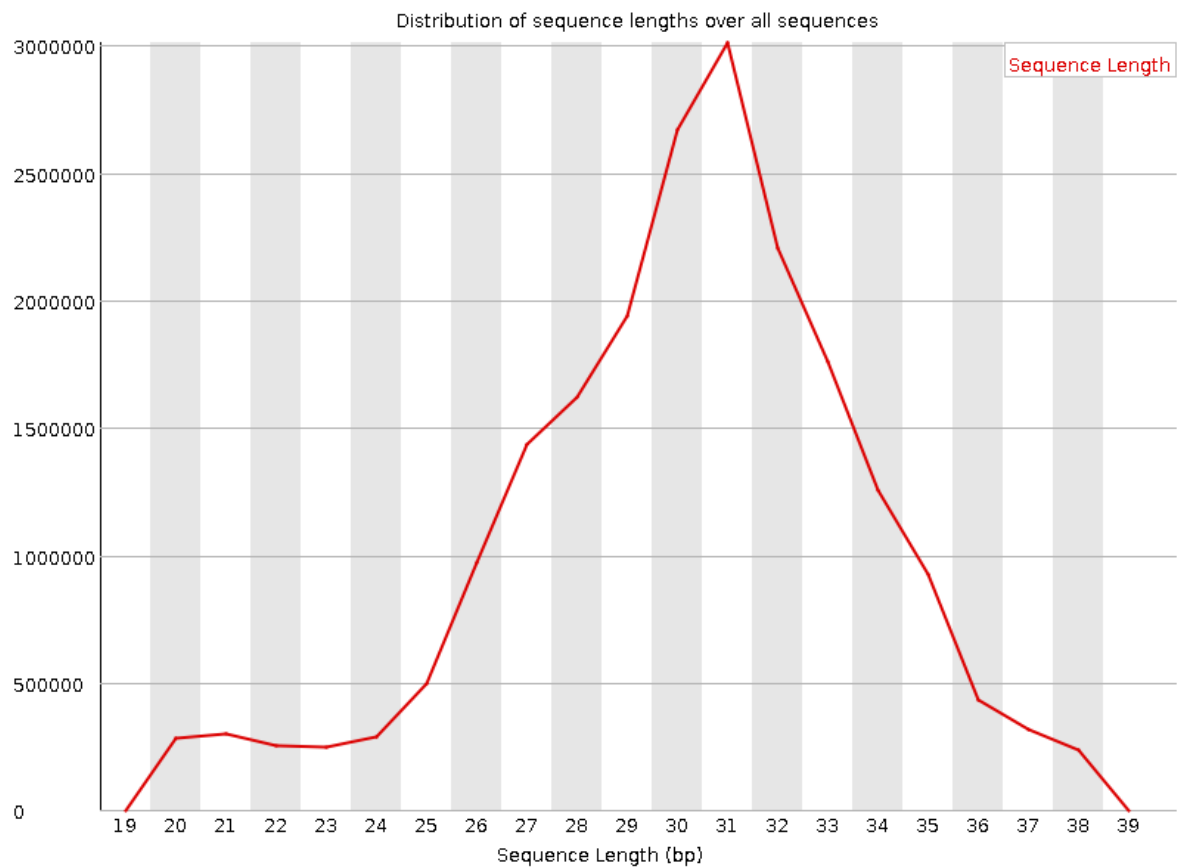
FastQC после всех операций. RiboSeq rep1:



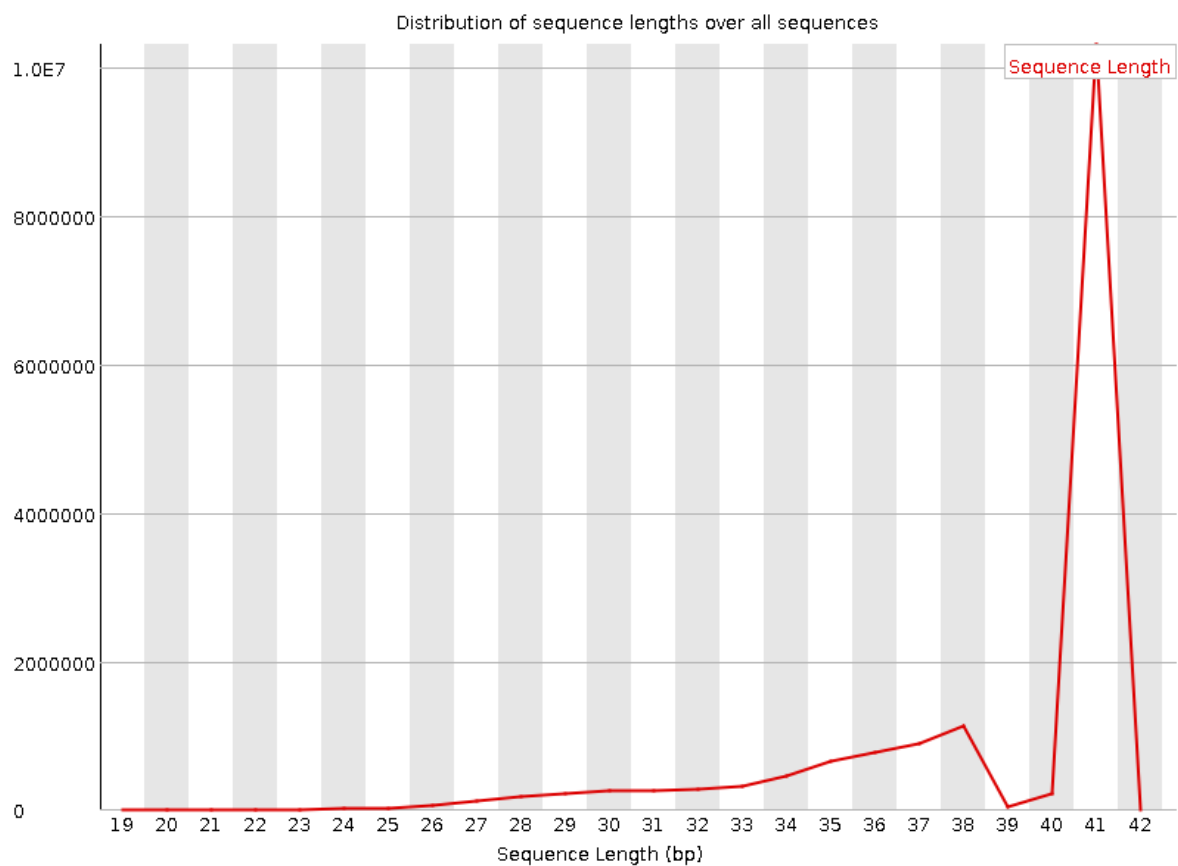
RNASeq rep1:



RiboSeq rep2:



RNASeq rep2:



Видим, что характерная длина рибосомных футпринтов - 25-32 пн. Из первой реплики RNASeq можно сделать вывод о том, что авторы старались сделать фрагменты РНК близкими по длине к рибосомным.

Этап 3: проверим долю рибосомной РНК

Построим индекс `bowtie` (скачал rRNA_euk.fasta с сервера):

```
bowtie-build --threads 8 rRNA_euk.fasta rRNA
```

Получим оценки доли рРНК для каждого из образцов, полученных выше:

```
(bio) leo@MS7922:~/BioData/RiboSeq$ bowtie --sam -p 8 ./bwt/rRNA riboseq1.trimmed.fastq.gz --chunkmbs 10000 > /dev/null
# reads processed: 46789965
# reads with at least one reported alignment: 23056356 (49.28%)
# reads that failed to align: 23733609 (50.72%)
Reported 23056356 alignments
(bio) leo@MS7922:~/BioData/RiboSeq$ bowtie --sam -p 8 ./bwt/rRNA rnaseq1.trimmed.fastq.gz --chunkmbs 10000 > /dev/null
# reads processed: 48273908
# reads with at least one reported alignment: 29740959 (61.61%)
# reads that failed to align: 18532949 (38.39%)
Reported 29740959 alignments
(bio) leo@MS7922:~/BioData/RiboSeq$ bowtie --sam -p 8 ./bwt/rRNA riboseq2.final.fastq.gz --chunkmbs 10000 > /dev/null
# reads processed: 20712568
# reads with at least one reported alignment: 12134426 (58.58%)
# reads that failed to align: 8578142 (41.42%)
Reported 12134426 alignments
(bio) leo@MS7922:~/BioData/RiboSeq$ bowtie --sam -p 8 ./bwt/rRNA rnaseq2.final.fastq.gz --chunkmbs 10000 > /dev/null
# reads processed: 16397554
# reads with at least one reported alignment: 10441188 (63.68%)
# reads that failed to align: 5956366 (36.32%)
Reported 10441188 alignments
```

Этап 4: картирование на реальный геном

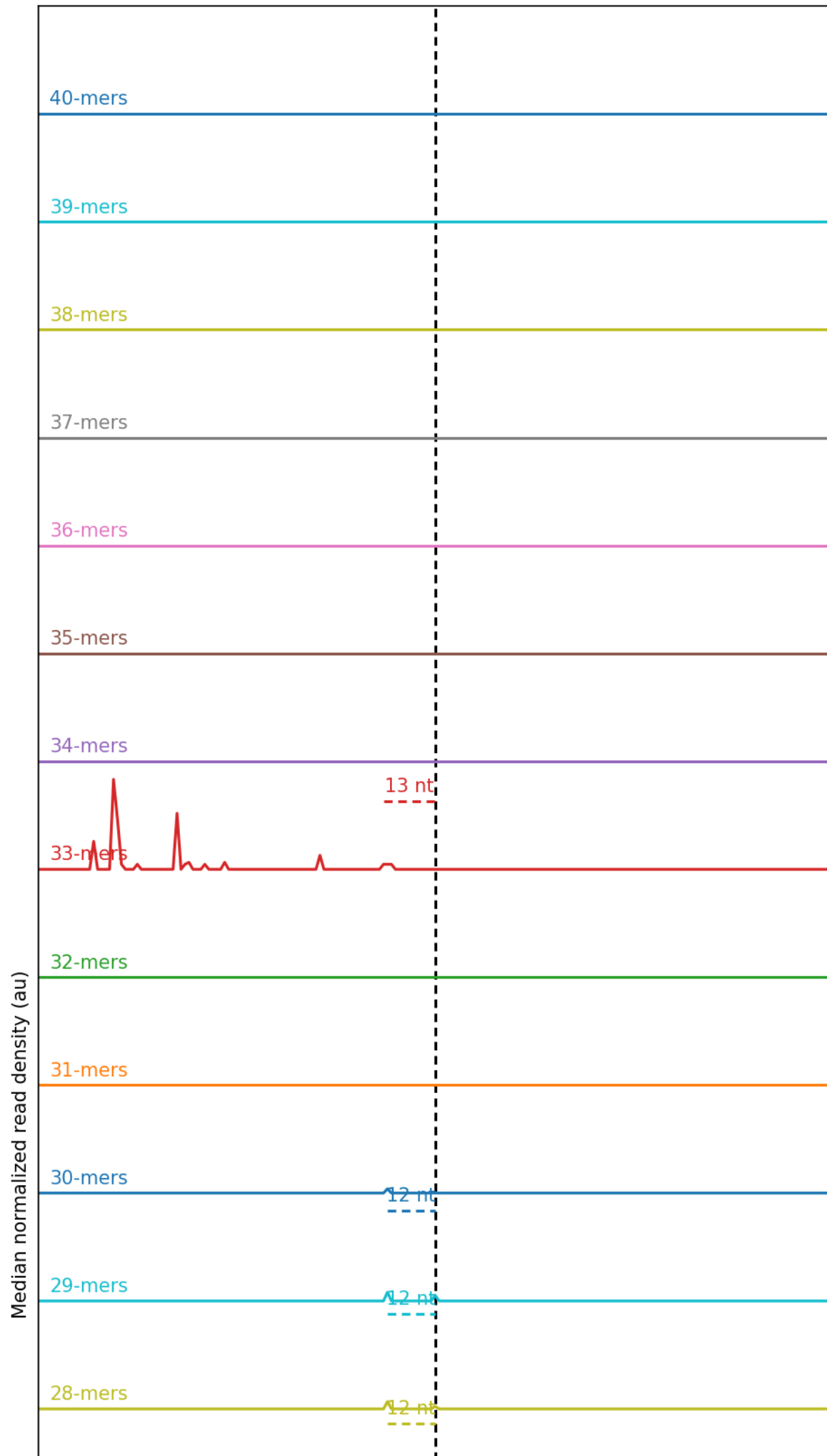
Использовал готовые файлы.

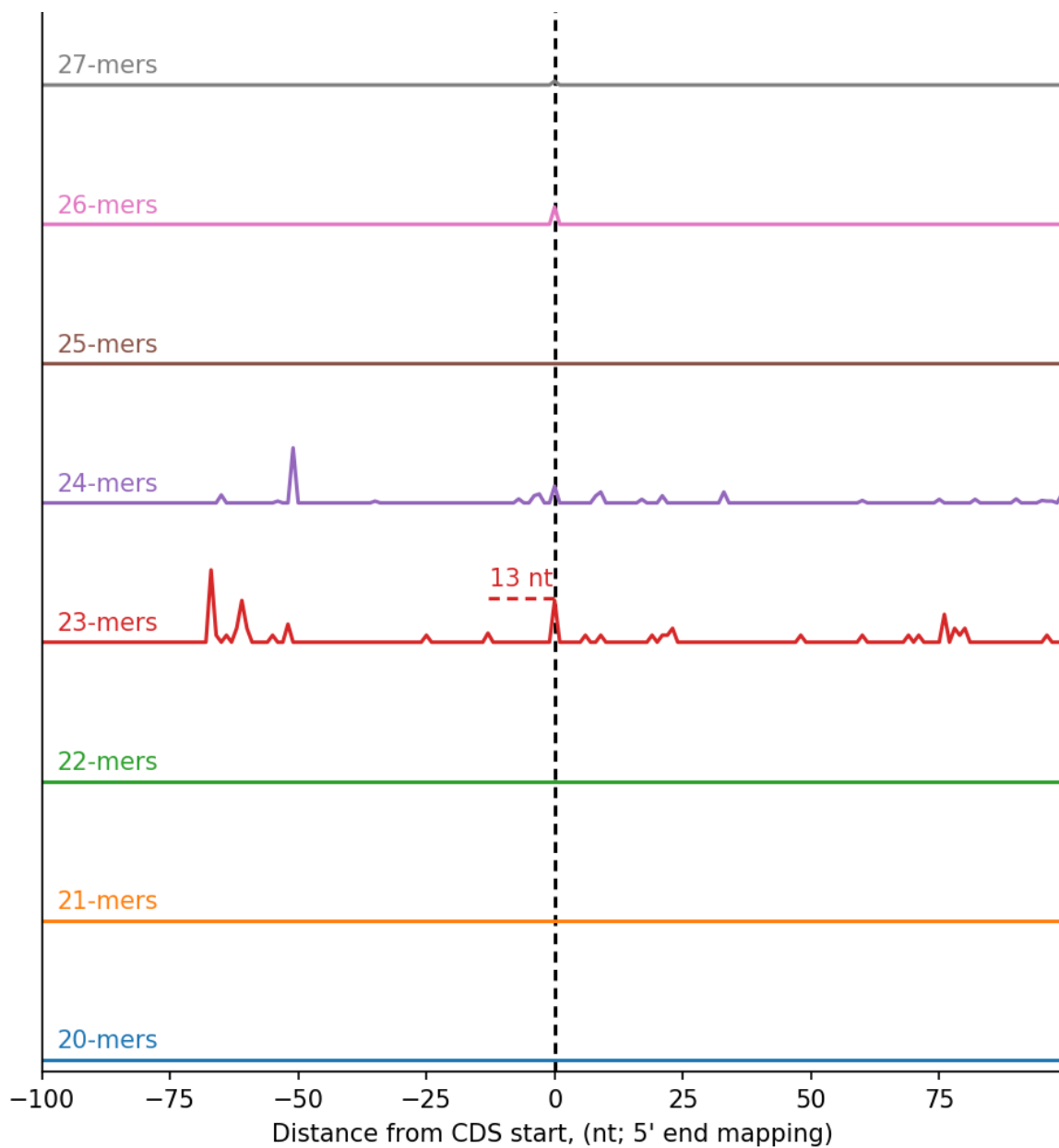
Этап 5: фазирование прочтений и метагенные профили

Фазирование по RiboSeq нокдауну METTL3 без aggregate:

```
psite ~/BioData/RiboSeq/plastidmetagene/mouse_start_rois.txt psite_test --
countfile_format BAM --count_files
~/BioData/RiboSeq/olbams/METTL3_ribo_Coots2017_m_r1.bam
~/BioData/RiboSeq/olbams/METTL3_ribo_Coots2017_m_r2.bam --min_length 20 --
max_length 40 --constrain 10 18 --min_count 10 --default 14
```


Fiveprime read offsets by length

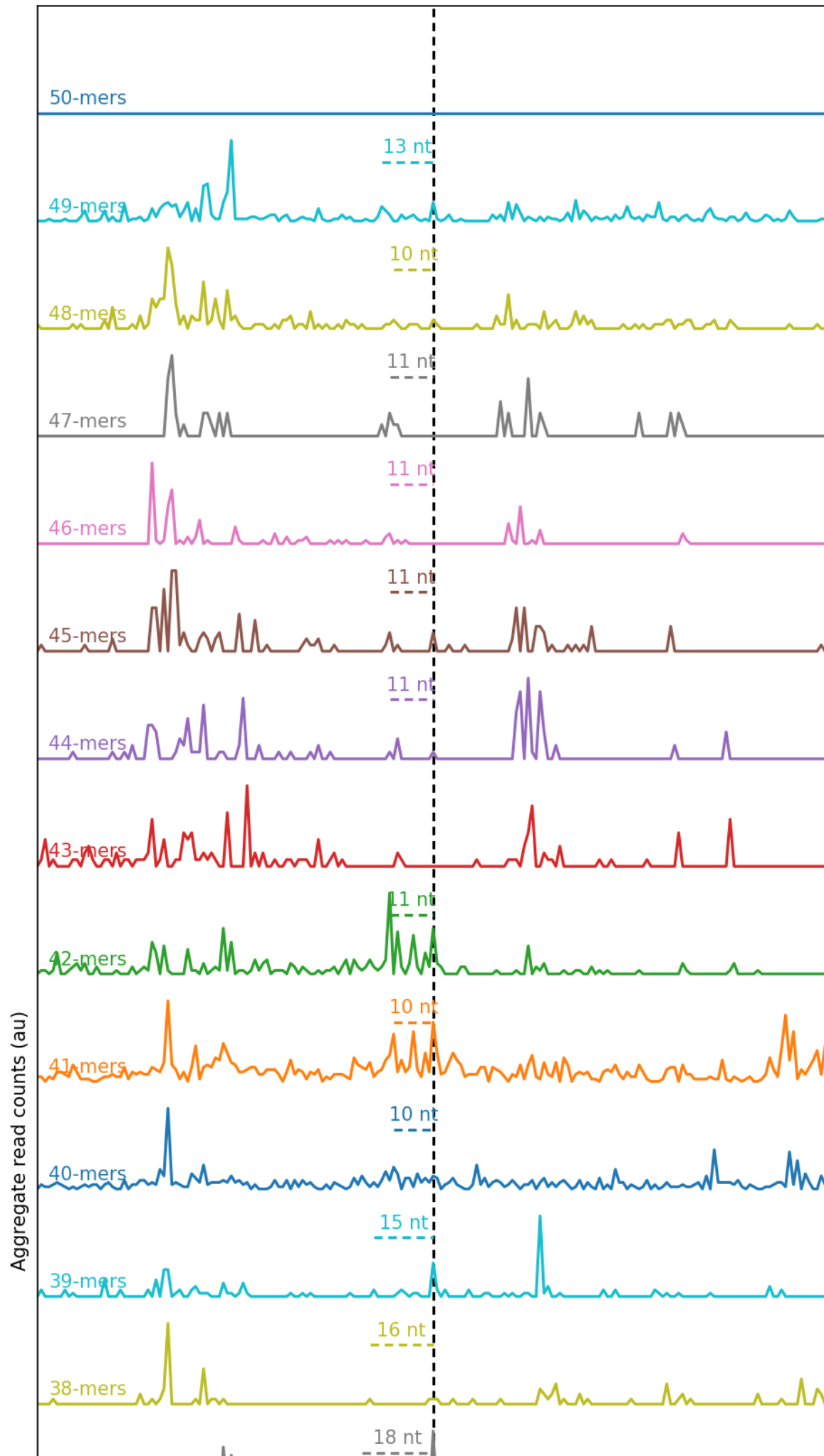


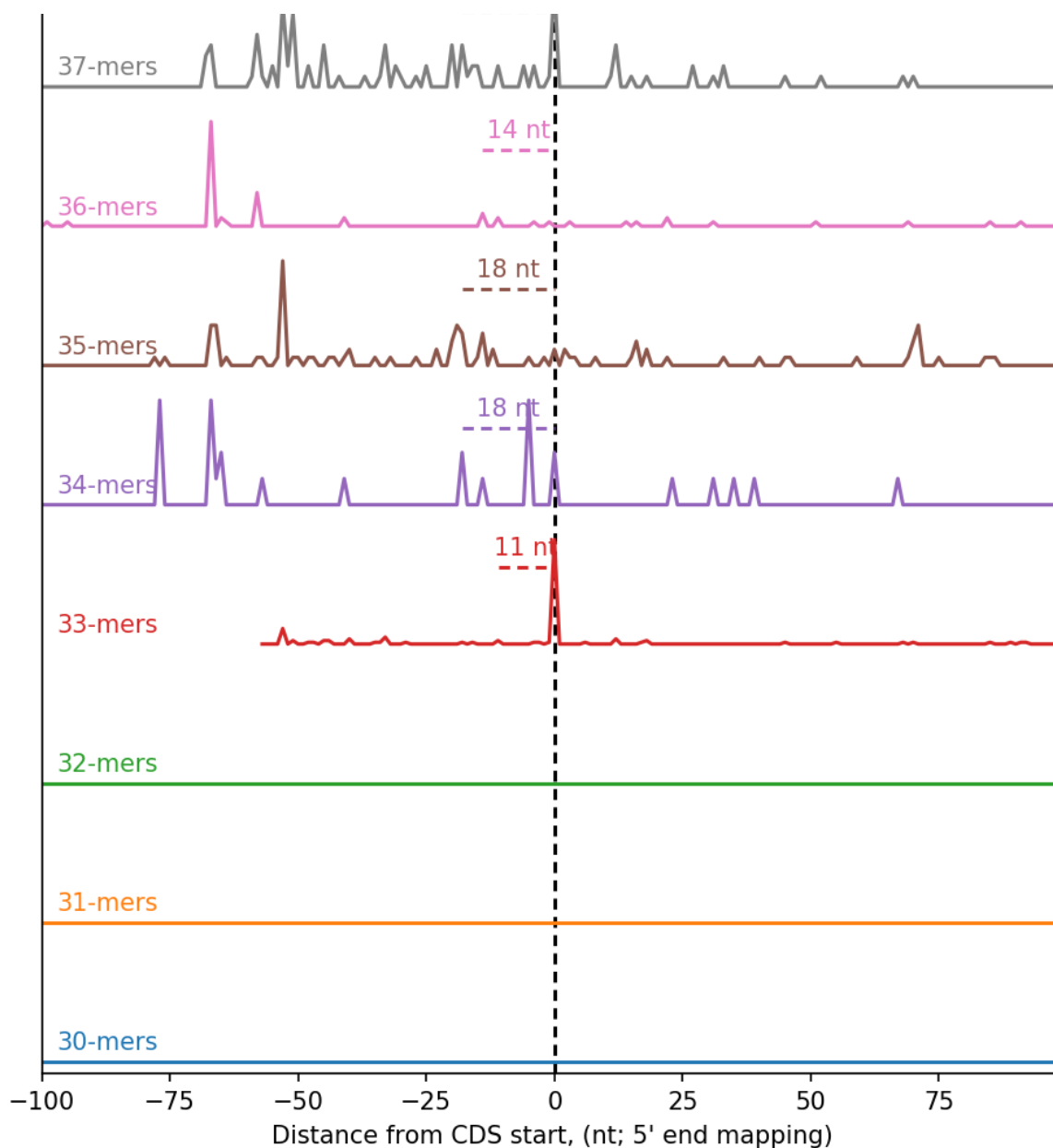


Фазирование по RNASeq (от 30 до 50):

```
psite ~/BioData/RiboSeq/plastidmetagene/mouse_start_rois.txt psite_test --
countfile_format BAM --count_files
~/BioData/RiboSeq/olbams/METTL3_rna_Coots2017_m_r1.bam
~/BioData/RiboSeq/olbams/METTL3_rna_Coots2017_m_r2.bam --min_length 30 --
max_length 50 --aggregate --constrain 10 18 --min_count 10 --default 14
```

Fiveprime read offsets by length

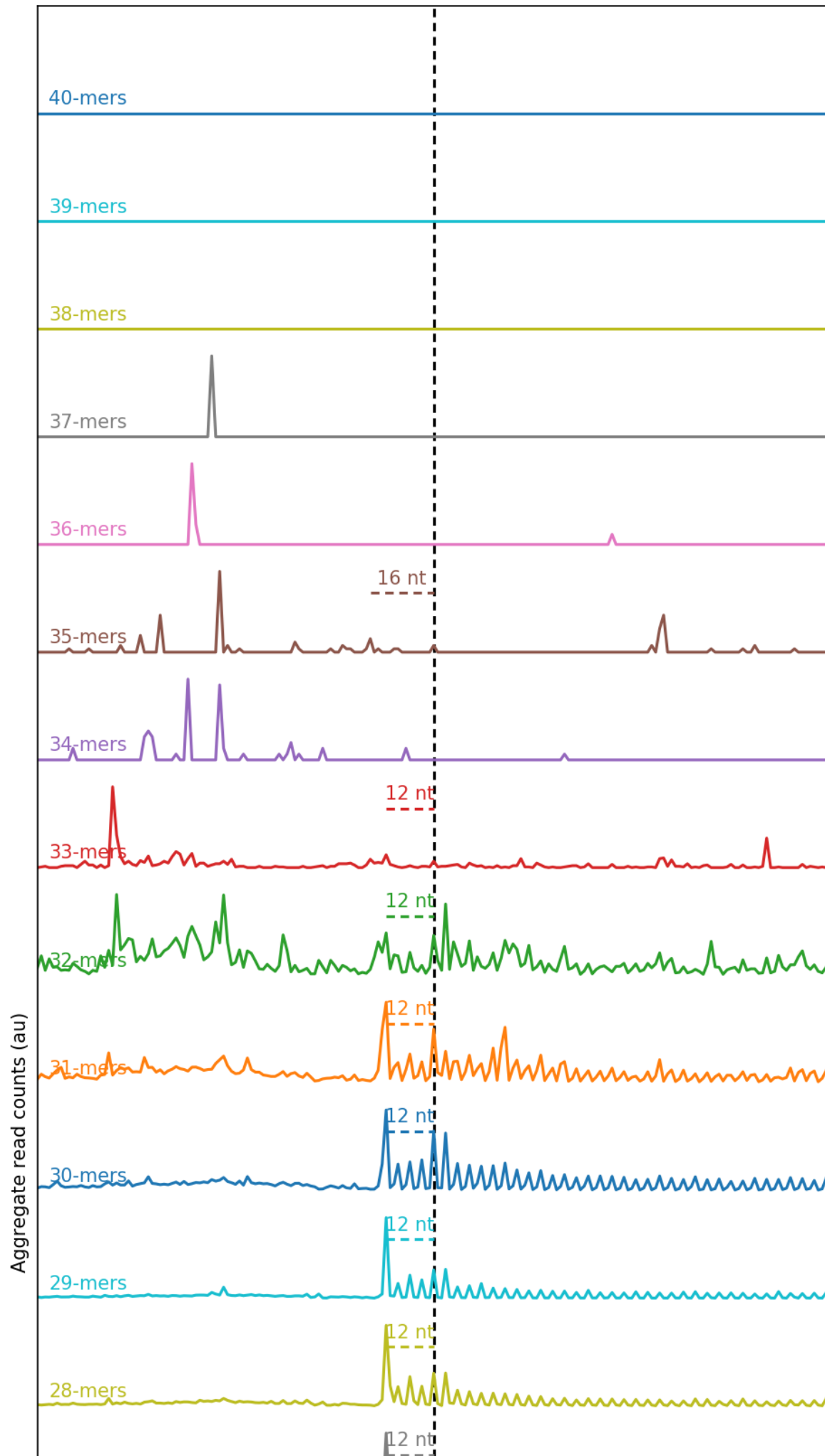


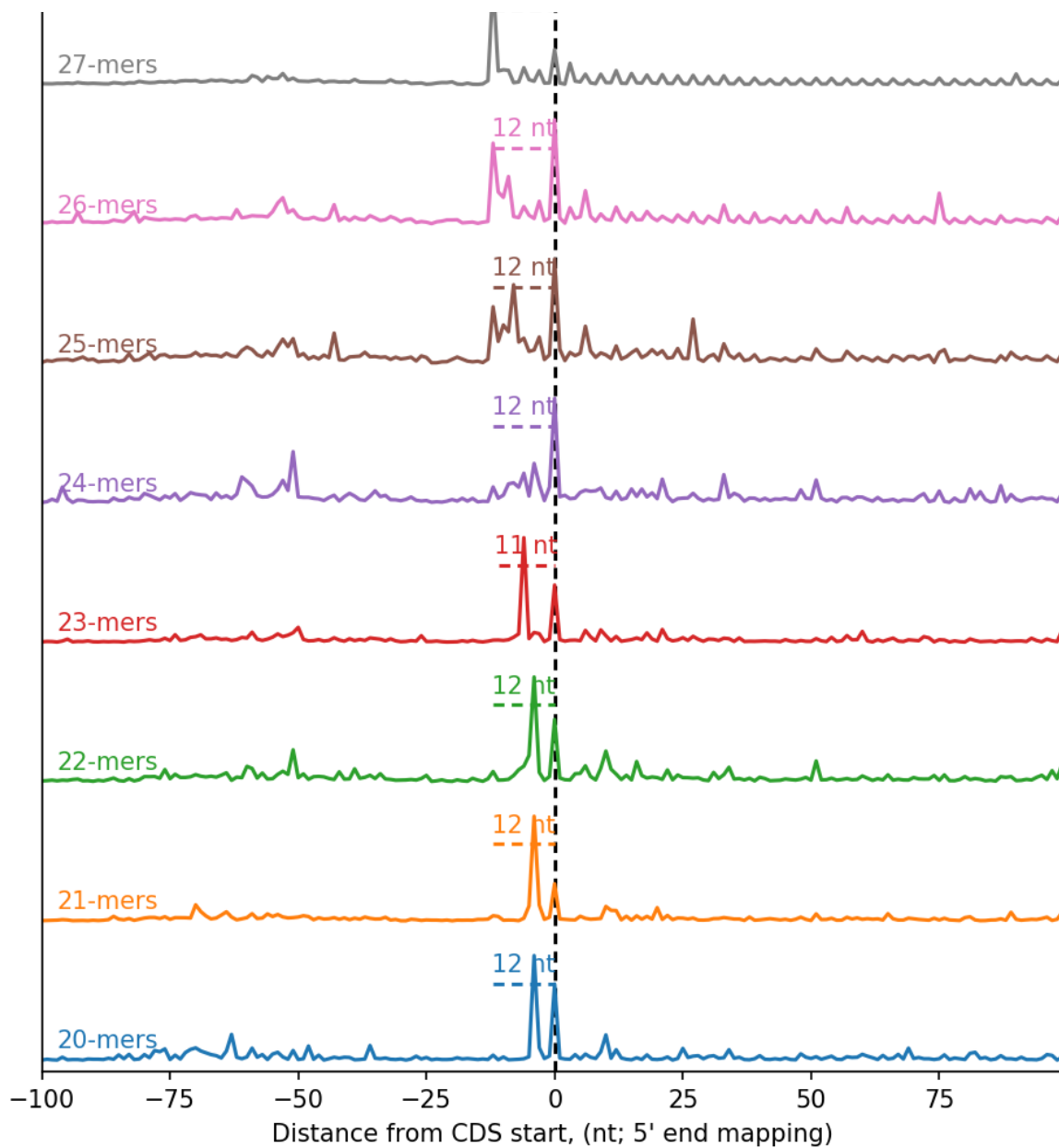


Фазирование по всем RiboSeq-образцам:

```
psite ~/BioData/RiboSeq/plastidmetagene/mouse_start_rois.txt psite_test --
countfile_format BAM --countfiles $(find ~/BioData/RiboSeq/bams/ -name
'*ribo*.bam' | tr '\n' ' ') --min_length 20 --max_length 40 --aggregate --
constrain 10 18 --min_count 10 --default 14
```

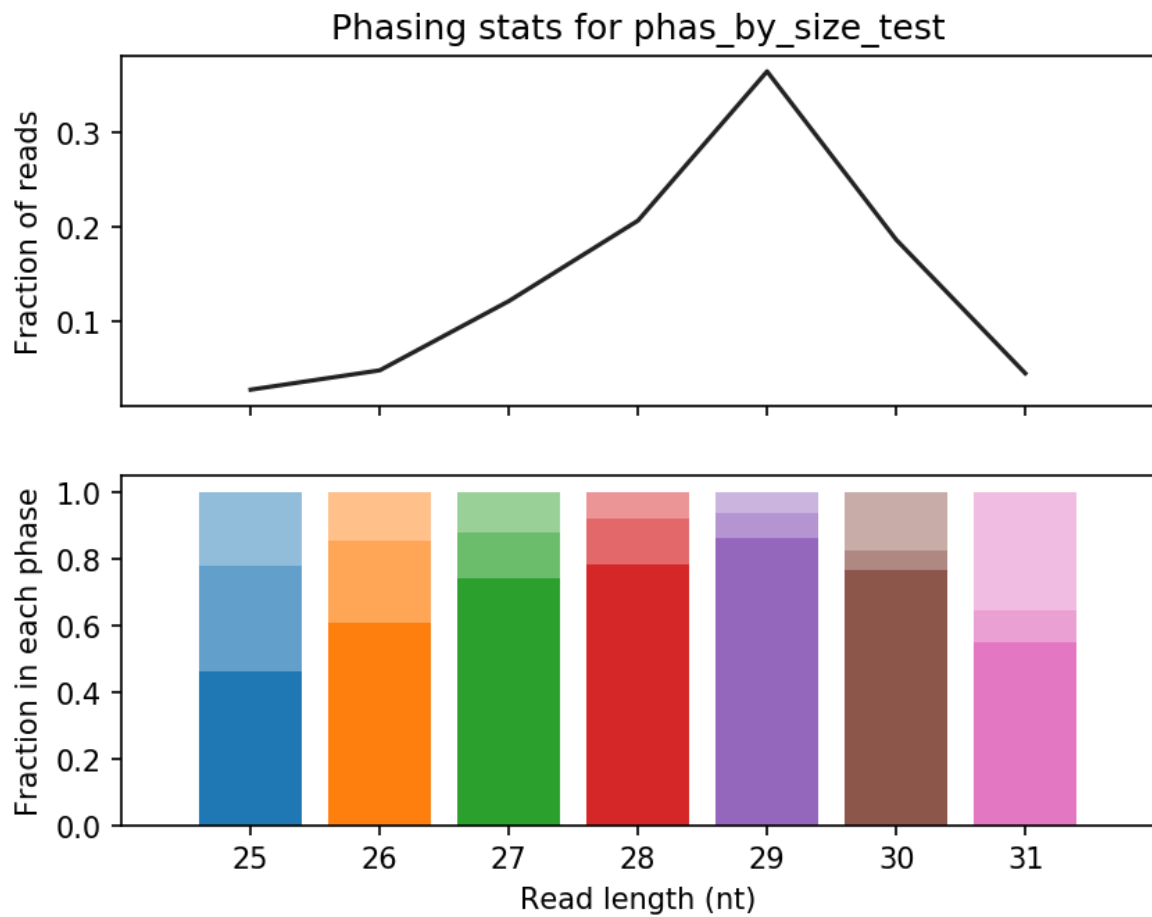
Fiveprime read offsets by length





Проверка фазирования:

```
phase_by_size --countfile_format BAM --count_files $(find
~/BioData/RiboSeq/bams/ -name '*ribo*.bam' | tr '\n' ' ') --fiveprime_variable
--offset ./psite/psite_test_p_offsets.txt --min_length 25 --max_length 31
~/BioData/RiboSeq/plastidmetagene/mouse_start_rois.txt phas_by_size_test
```

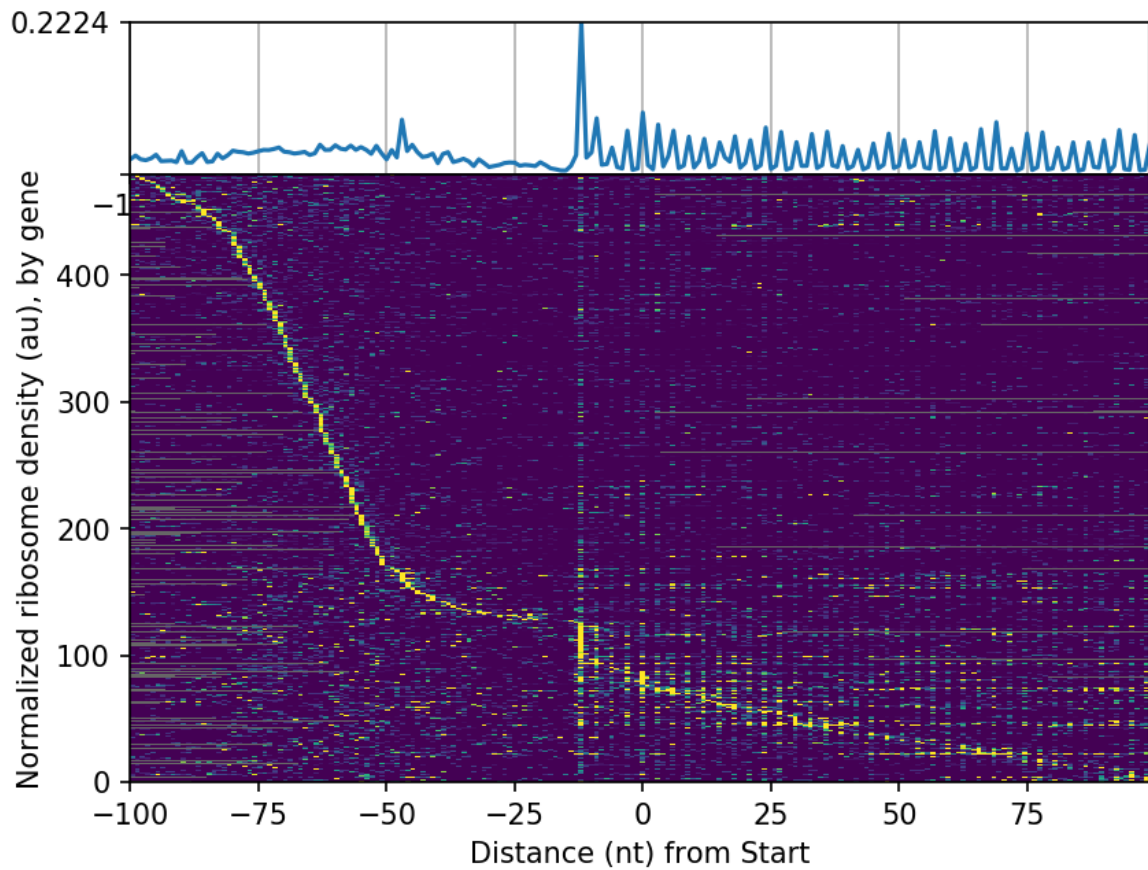


Получилось достаточно правильное, большая часть ридов ложится на frame0.

Нефазируемый метагеновый профиль METTL3 rep2 (длина 25-31):

```
metagene count --countfile_format BAM --count_files
~/BioData/RiboSeq/bams/METTL3_ribo_Coots2017_m_r2.bam --fiveprime --min_length
25 --max_length 31 --min_count 10 --use_mean --landmark Start
~/BioData/RiboSeq/plastidmetagene/mouse_start_rois.txt metagene_counts_test
```

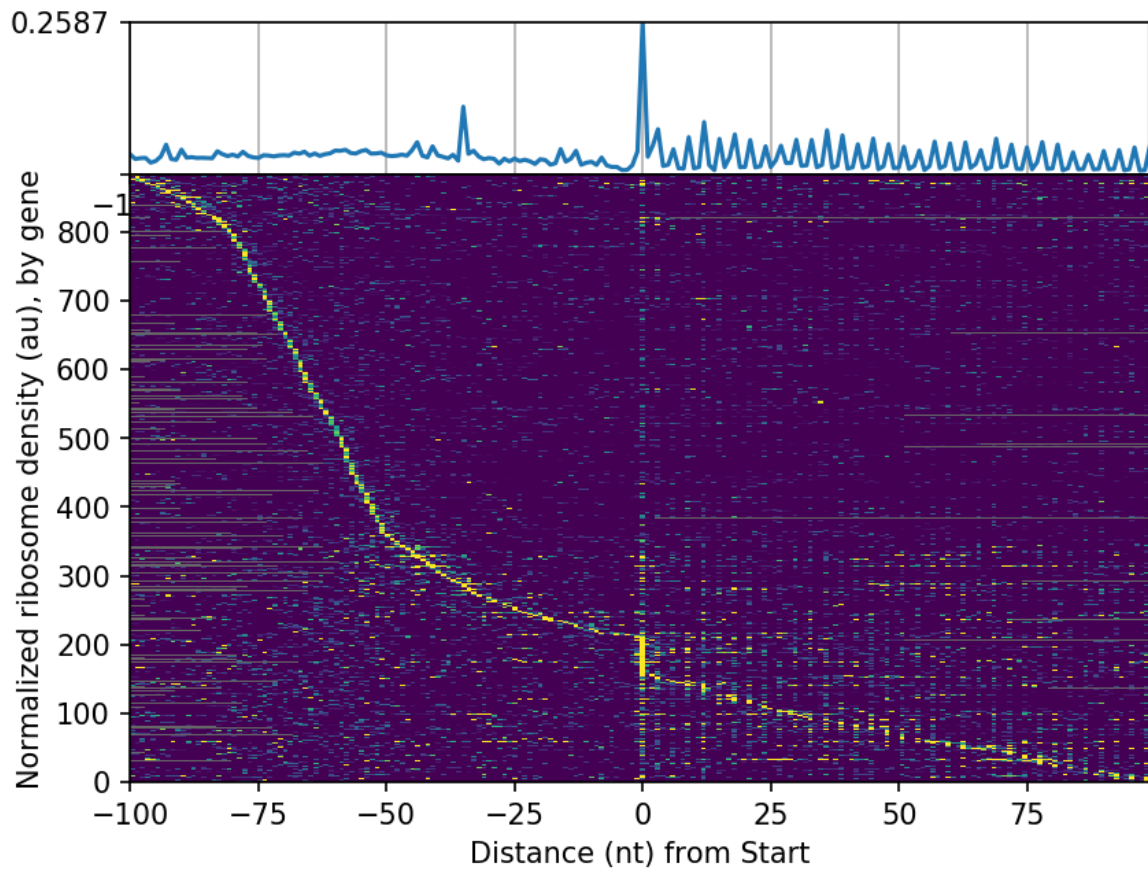
Metagene overview for metagene_counts_test



Фазированный метагеновый профиль того же образца:

```
metagene count --countfile_format BAM --count_files
~/BioData/RiboSeq/bams/METTL3_ribo_Coots2017_m_r2.bam --fiveprime_variable --
offset ./psite/psite_test_p_offsets.txt --min_length 25 --max_length 31 --
min_count 5 --use_mean --landmark Start
~/BioData/RiboSeq/plastidmetagene/mouse_start_rois.txt metagene_counts_test
```

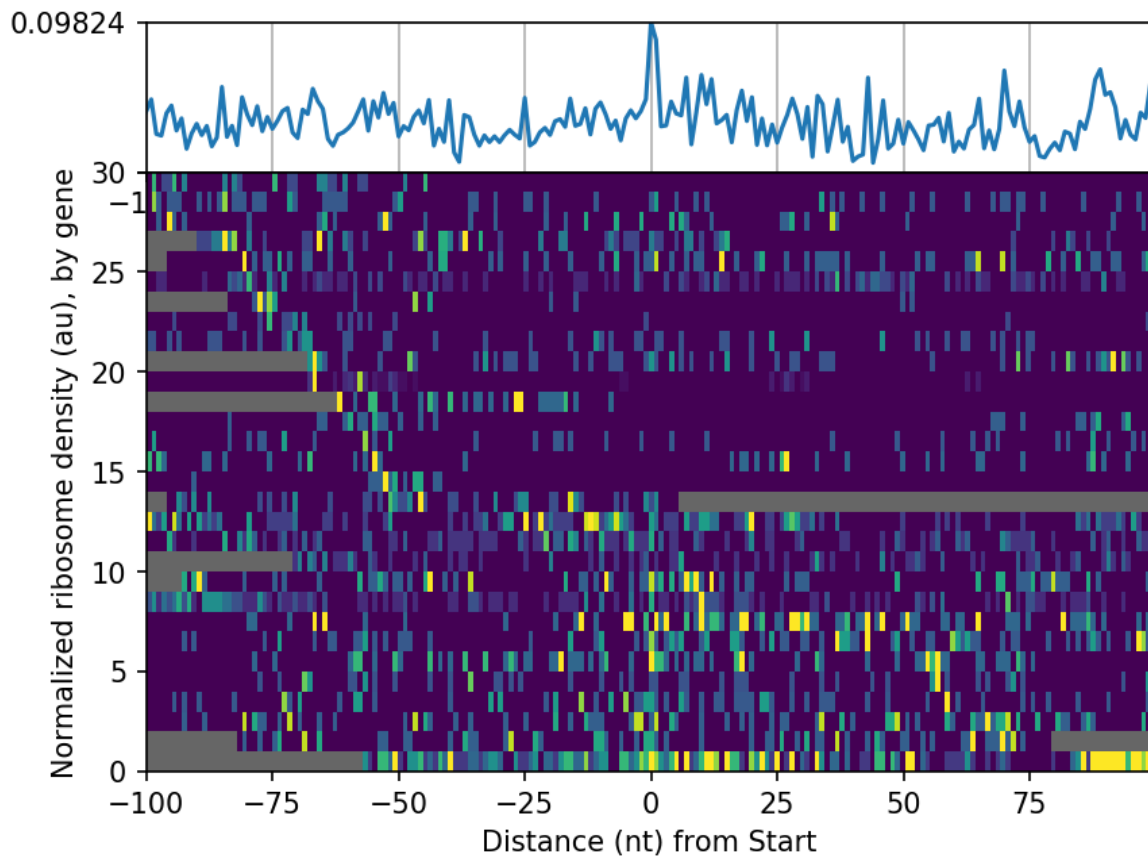

Metagene overview for metagene_counts_test



Метагенный профиль RNASeq METTL3 rep2:

```
metagene count --countfile_format BAM --count_files
~/BioData/RiboSeq/bams/METTL3_rna_Coots2017_m_r2.bam --fiveprime --min_length
35 --max_length 45 --min_count 10 --use_mean --landmark Start
~/BioData/RiboSeq/plastidmetagene/mouse_start_rois.txt metagene_counts_test
```

Metagene overview for metagene_counts_test



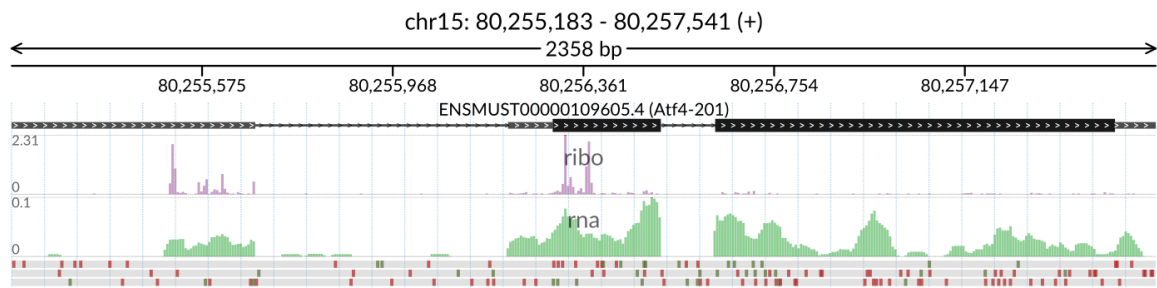
Этап 6: получение bedGraph и визуализация

Сделаем профили RiboSeq и RNASeq для METTL3 rep2:

```
make_wiggle -o METTL3_ribo2 --count_files
~/BioData/RiboSeq/bams/METTL3_ribo_Coots2017_m_r2.bam --normalize --min_length
25 --max_length 31 --fiveprime_variable --offset
../psite/psite_test_p_offsets.txt
make_wiggle -o METTL3_rna2 --count_files
~/BioData/RiboSeq/bams/METTL3_rna_Coots2017_m_r2.bam --normalize --center
bedtools unionbedg -i METTL3_ribo2_fw.wig METTL3_ribo2_rc.wig | csvtk mutate2
-H -t -L 5 -e '$4+$5' | cut -f 4-5 --complement > ribo.bedGraph
bedtools unionbedg -i METTL3_rna2_fw.wig METTL3_rna2_rc.wig | csvtk mutate2 -H
-t -L 5 -e '$4+$5' | cut -f 4-5 --complement > rna.bedGraph
```

Визуализируем ген *ATF4* с помощью `svist4get`:

```
svist4get -gtf ~/BioData/Annotations/gencode.vM23.basic.annotation.gtf -fa
~/BioData/GRCm38.primary_assembly.genome.fa -bg ribo.bedGraph rna.bedGraph -g
ENSMUSG00000042406.8
```



Как можно заметить, данный ген имеет преждевременную рамку считывания в 5'-UTR.