

# Анализ данных MNase seq

---

Кожевникова Дарья

Авг - Сен 2020

# Цели и задачи

## Цели:

- Научиться обрабатывать данные MNase seq
- Сопоставить результаты анализа с имеющимся теоретическим знанием

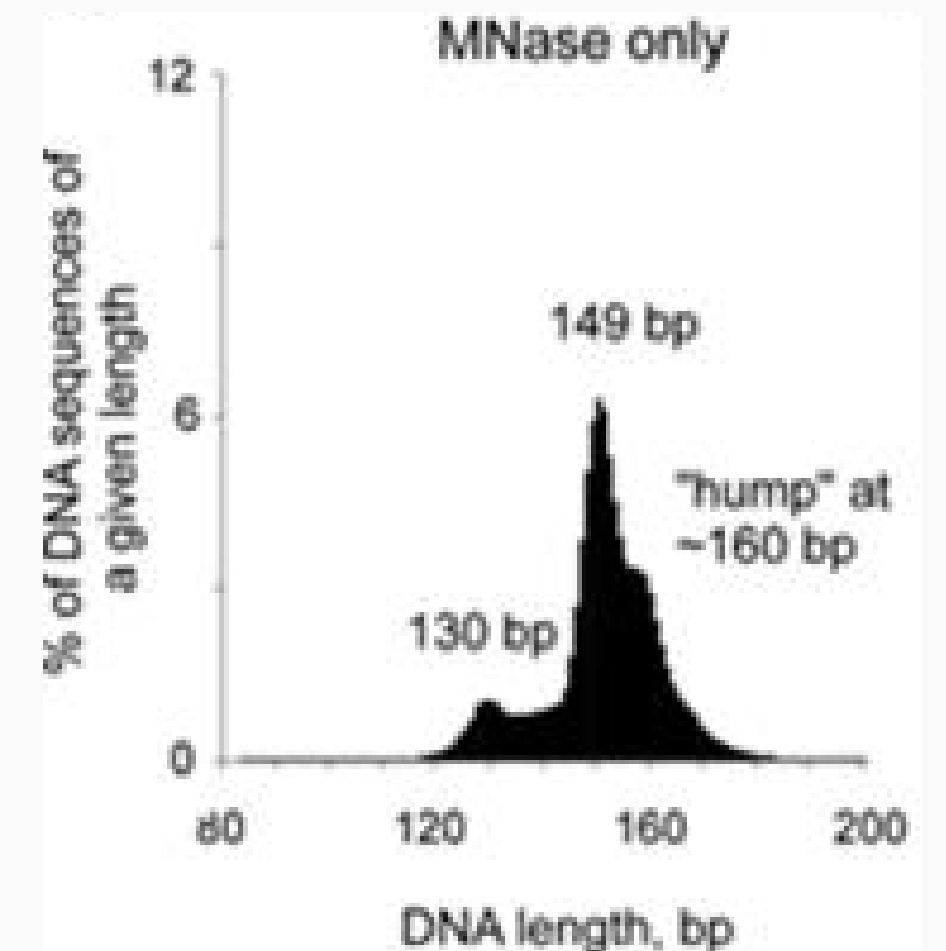
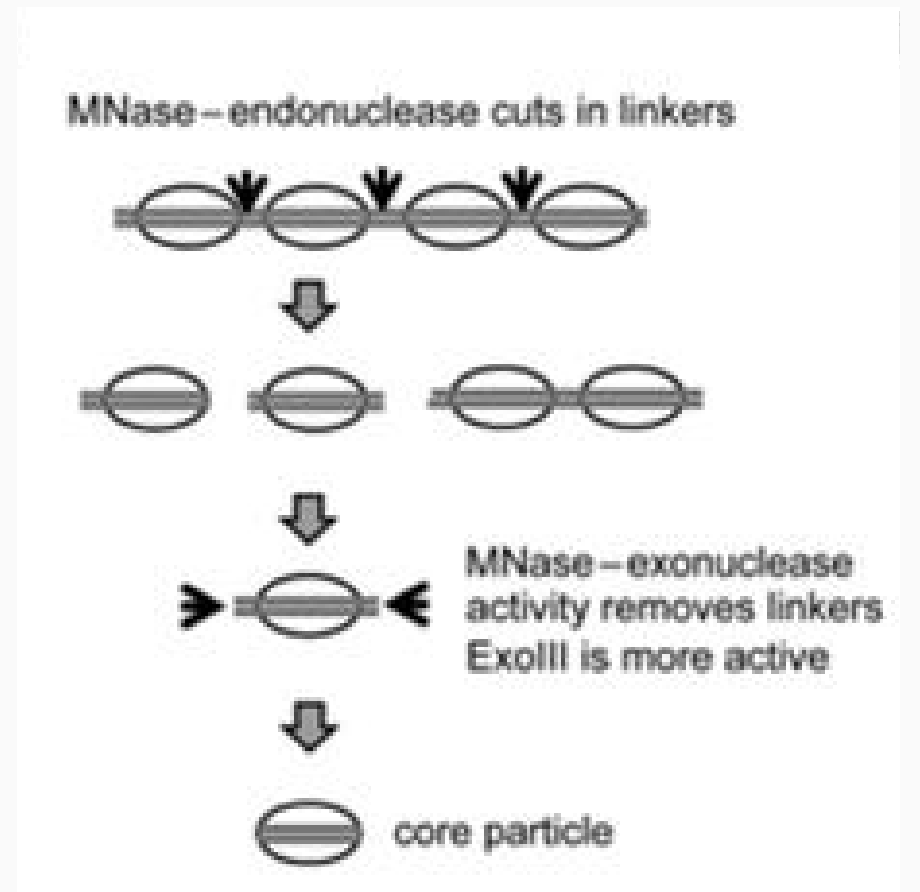
## Задачи:

- Сделать обзор литературы
- Скачать данные из статьи Feng Cui, Hope A. Cole, David J. Clark, Victor B. Zhurkin, Transcriptional activation of yeast genes disrupts intragenic nucleosome phasing, Nucleic Acids Research, Volume 40, Issue 21, 1 November 2012, Pages 10753–10764.
- Обработка данных с использованием инструментов биоинформатики
- Формулировка выводов

# MNase seq data

The most frequently used method of mapping nucleosome positions and occupancy involves:

- Digestion of chromatin with micrococcal nuclease (MNase), an endo- and exo-nuclease that preferentially digests the naked DNA between nucleosomes, releases the nucleosomes from chromatin, and enriches the nucleosome-protected DNA fragments.
- To determine nucleosome positions and occupancy, the resulting undigested DNA is subjected to high throughput sequencing (MNase-seq).
- Mapping reads to the reference genome.



Feng Cui, Hope A. Cole, David J. Clark, Victor B. Zhurkin,

# Transcriptional activation of yeast genes disrupts intragenic nucleosome phasing

*Nucleic Acids Research*, Volume 40, Issue 21, 1 November 2012, Pages 10753–10764,

<https://doi.org/10.1093/nar/gks870>

---

Nucleosomes often undergo extensive rearrangement when genes are activated for transcription.

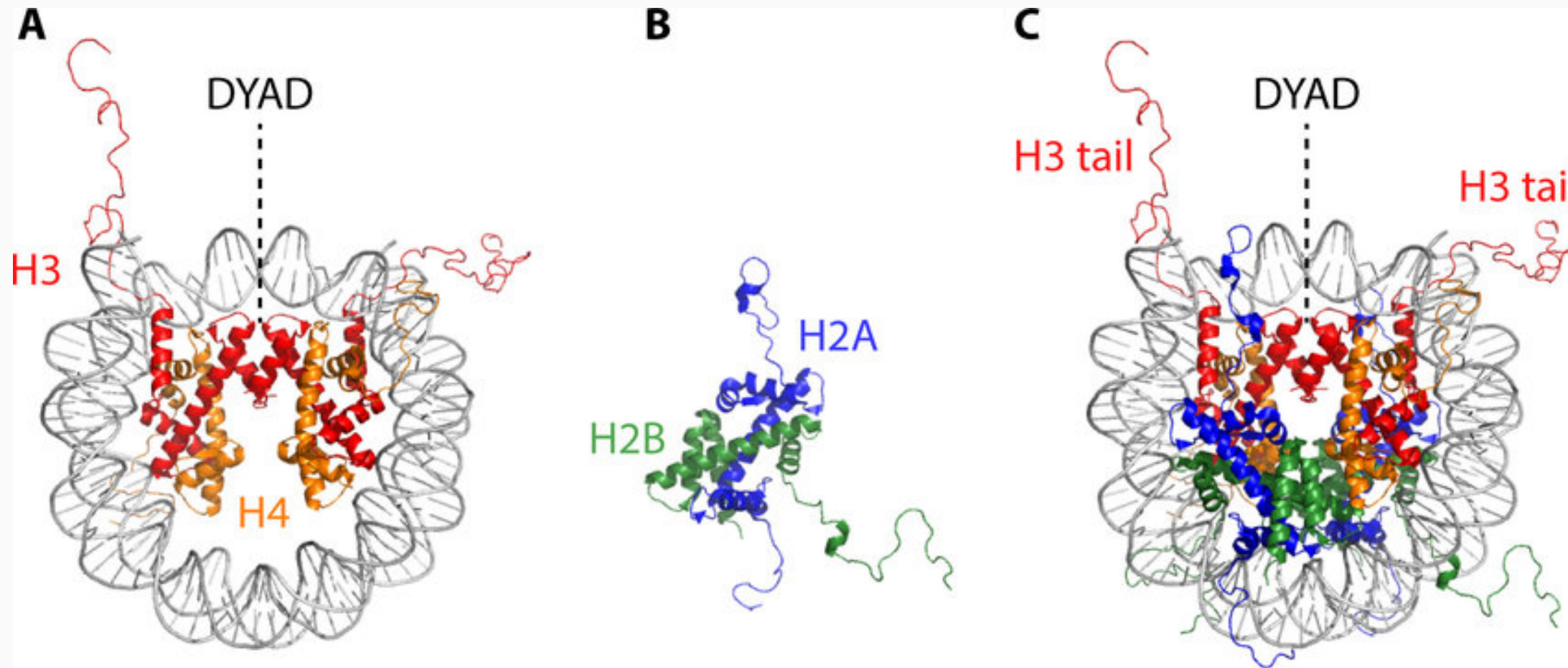
It was shown previously, using paired-end sequencing of yeast nucleosomes, that major changes in chromatin structure occur when genes are activated by 3-aminotriazole (3AT), an inducer of the transcriptional activator Gcn4.

At the genomic level, nucleosomes are regularly phased relative to the transcription start site.

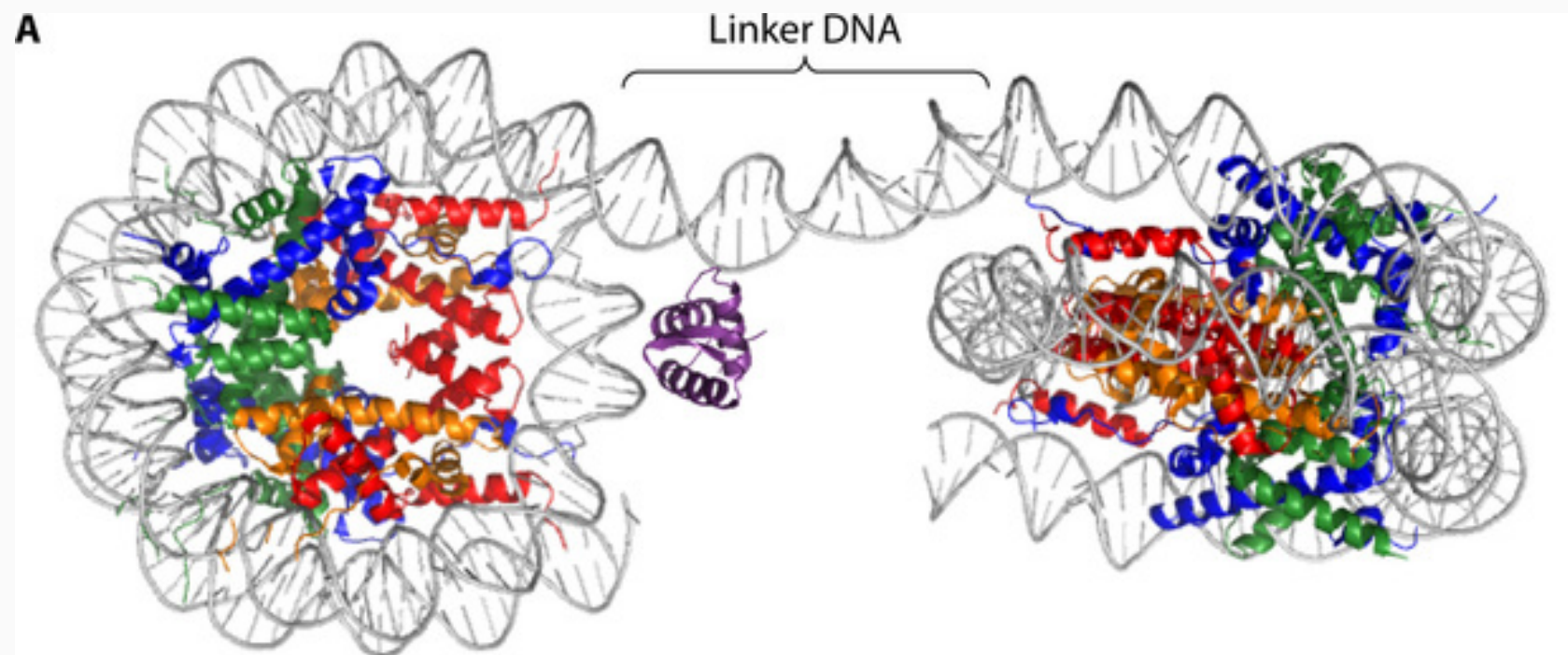
However, for a subset of 234 strongly induced genes, this phasing is much more irregular after induction, consistent with the loss of some nucleosomes and the re-positioning of the remaining nucleosomes.

DAC analysis of the 3AT-induced genes suggests that transcription activation coincides with rearrangement of nucleosomes into irregular arrays with longer spacing.

Sequence analysis of the +1 nucleosomes belonging to the 45 most strongly activated genes reveals a distinctive periodic oscillation in the A/T-dinucleotide occurrence that is present throughout the nucleosome and extends into the linker. This unusual pattern suggests that the +1 nucleosomes might be prone to sliding, thereby facilitating transcription.



Assembly of the nucleosome core particle. (A) Association of the (H3/H4)<sub>2</sub> tetramer with DNA nucleates nucleosome assembly and defines the dyad axis. (B) H2A/H2B dimer. (C) Two H2A/H2B dimers are deposited to generate the nucleosome core particle. H3 N-terminal tails emerge near the DNA entry/exit points (based on data reported under PDB accession number 1KX5).



Histone H1 associates with linker DNA. (A) The globular domain of histone H1 (purple) binds the nucleosome at the dyad.

# Nucleosome positioning

## Translational positioning

- defined by the nucleosome midpoint (or dyad) with regard to the DNA sequence
- the DNA sequence patterns specifying translational nucleosome positioning are far from clear. The only well-established feature is the tendency of long A/T-rich fragments, and the A-tracts in particular, to be excluded from nucleosomes
- nucleosome positioning can be affected by DNA-binding transcription factors and chromatin remodeling enzymes

## Rotational positioning

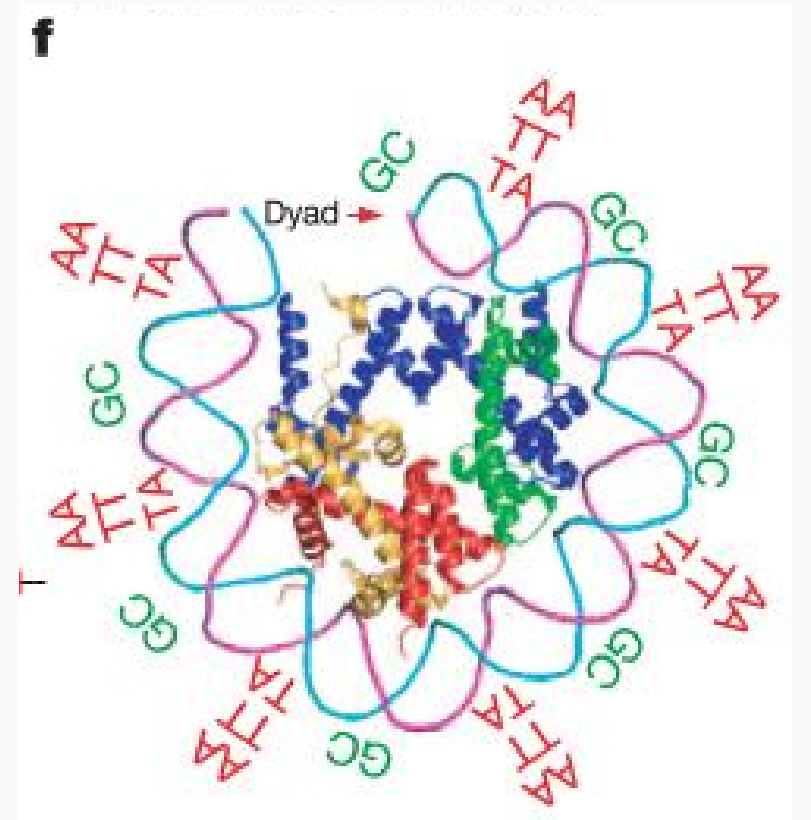
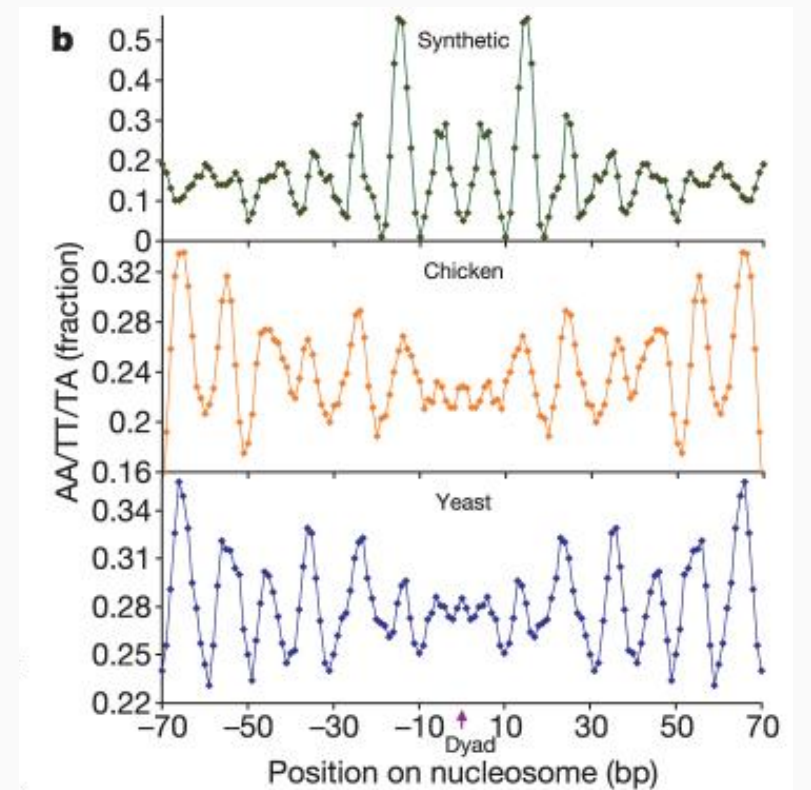
- defined by the side of the DNA helix that faces the histones
- related to the sequence-dependent preferences for DNA deformation, e.g. bending:  
In particular, the A/T-containing dimeric steps AA:TT, AT and TA preferentially occur where the DNA is bent into the **minor groove**, while G/C-containing dimers GG:CC, GC and CG are frequently situated at the sites where DNA is bent toward the **major groove**. The occurrences of AT and GC dimers in nucleosome core DNA both display **sinusoidal patterns** with **~10-bp periodicity**, but they are **~5 bp out of phase** with one another. These sequence patterns are observed in nucleosomal DNA from chicken, yeast, fruit fly, nematode and human, indicating that the sequence rules for rotational positioning are essentially the same across species.

# Nucleosome positioning code

Segal E, Fondufe-Mittendorf Y, Chen L, et al.: "Our findings demonstrate that eukaryotic genomes use a nucleosome positioning code, and link the resulting nucleosome positions to specific chromosome functions."

As expected for a nucleosome–DNA interaction model, the resulting model exhibits distinctive sequence motifs that recur periodically at the DNA helical repeat and are known to facilitate the sharp bending of DNA around the nucleosome<sup>3</sup>. These include 10-bp periodic AA/TT/TA dinucleotides that oscillate in phase with each other and out of phase with, 10-bp periodic GC dinucleotides.

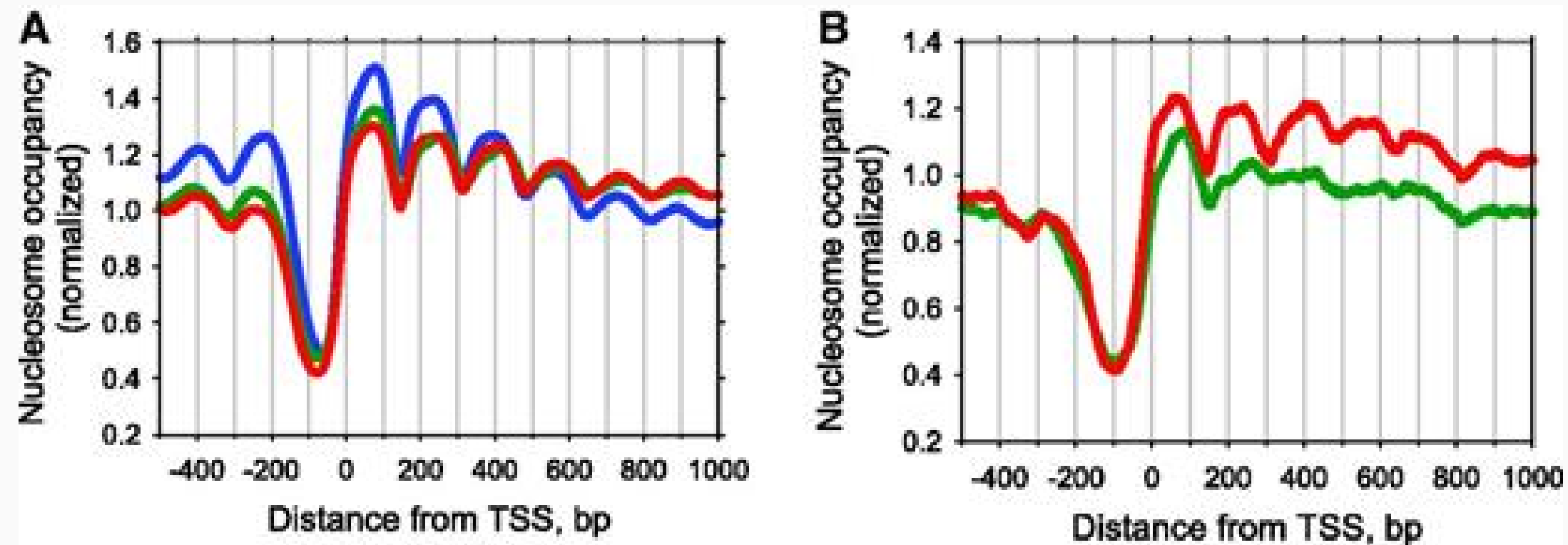
Improving the agreement of a sequence with these motifs increased its binding affinity to the nucleosome, whereas changing the periodicity or deleting the key motifs decreased that affinity.



b, Fraction (3-bp moving average) of AA/TT/TA dinucleotides at each position of centre-aligned yeast, chicken or random chemically synthesized nucleosome-bound DNA sequences, showing, 10-bp periodicity of these dinucleotides.

f, Key dinucleotides inferred from the alignments are shown relative to the three-dimensional structure of one-half of the symmetric nucleosome.

# Local changes in nucleosome organization upon 3AT induction



Nucleosome organization around the 5'-end of yeast genes. ( A ) Overlays of nucleosome occupancy profiles of 4792 *S. cerevisiae* genes ( 30 ) relative to TSS (position 0). Nucleosome occupancy values are either taken directly from Kaplan et al. ( 29 ) (blue) or recalculated from Cole et al. (16 , 27 ): 3AT set (green) and CC set (red), respectively. In the latter two cases, the NCP fragments 147–152 bp in length were selected to calculate occupancy profiles. ( B ) Nucleosome occupancy map for 234 genes (out of 4792 genes) that are induced by 3AT by more than 2-fold ( 25 ). Note that the occupancy value at each nucleotide is normalized by summing all the nucleosome sequences covering this nucleotide and dividing that number by the average number of nucleosome sequences per base pair across the genome.



# SRR data from NCBI

SRA  Search

COVID-19 is an emerging, rapidly evolving situation. Get the latest public health information from CDC: <https://www.coronavirus.gov>. Get the latest research from NIH: <https://www.nih.gov/coronavirus>. Find NCBI SARS-CoV-2 literature, sequence, and clinical content: <https://www.ncbi.nlm.nih.gov/sars-cov-2/>.

Access: Public (4)  
Source: DNA (4)  
Library Layout: paired (4)  
Platform: Illumina (4)  
Strategy: EpiGenomics (4)  
Data in Cloud: GS (4), S3 (4)  
File Type: fastq (4)

Summary

Send results to Blast

Search results  
Items: 4

- [GSM651371: wild type +3AT 111A-3](#)  
1 ILLUMINA (Illumina Genome Analyzer II) run: 19.9M spots, 1.6G bases, 725.9Mb downloads  
Accession: SRX038810
- [GSM651370: wild type +3AT 111A-2](#)  
1 ILLUMINA (Illumina Genome Analyzer II) run: 19M spots, 1.5G bases, 495.3Mb downloads  
Accession: SRX038809
- [GSM651369: wild type control 111U-3](#)  
1 ILLUMINA (Illumina Genome Analyzer II) run: 18.3M spots, 1.5G bases, 506.3Mb downloads  
Accession: SRX038808
- [GSM651368: wild type control 111U-2](#)  
1 ILLUMINA (Illumina Genome Analyzer II) run: 24.4M spots, 2G bases, 672.6Mb downloads  
Accession: SRX038807

Filters: Manage Filters

Search in related databases

Database	Access		all
	public	controlled	
BioSample			
BioProject			
dbGaP			
GEO Datasets	<a href="#">1</a>		<a href="#">1</a>

Find related data  
Database: Select

Find items

Search details  
SRP005387[All Fields]

Search

[https://www.ncbi.nlm.nih.gov/sra/SRX038810\[accn\]](https://www.ncbi.nlm.nih.gov/sra/SRX038810[accn])

COVID-19 is an emerging, rapidly evolving situation. Get the latest public health information from CDC: <https://www.coronavirus.gov>. Get the latest research from NIH: <https://www.nih.gov/coronavirus>. Find NCBI SARS-CoV-2 literature, sequence, and clinical content: <https://www.ncbi.nlm.nih.gov/sars-cov-2/>.

Full

**SRX038810: GSM651371: wild type +3AT 111A-3**  
1 ILLUMINA (Illumina Genome Analyzer II) run: 19.9M spots, 1.6G bases, 725.9Mb downloads

Submitted by: Gene Expression Omnibus (GEO)

Study: Genome-wide nucleosome position maps in *Saccharomyces cerevisiae*  
[PRJNA136495](#) • [SRP005387](#) • [All experiments](#) • [All runs](#)  
[show Abstract](#)

Sample: wild type +3AT 111A-3  
[SAMN00191468](#) • [SRS153090](#) • [All experiments](#) • [All runs](#)  
Organism: *Saccharomyces cerevisiae*

Library:  
Name: GSM651371: wild type +3AT 111A-3  
Instrument: Illumina Genome Analyzer II  
Strategy: MNase-Seq  
Source: GENOMIC  
Selection: MNase  
Layout: PAIRED

Spot descriptor:  
1 forward 41 reverse

Experiment attributes:  
GEO Accession: GSM651371

Links:  
External link: [GEO Web Link](#)

Runs: 1 run, 19.9M spots, 1.6G bases, 725.9Mb

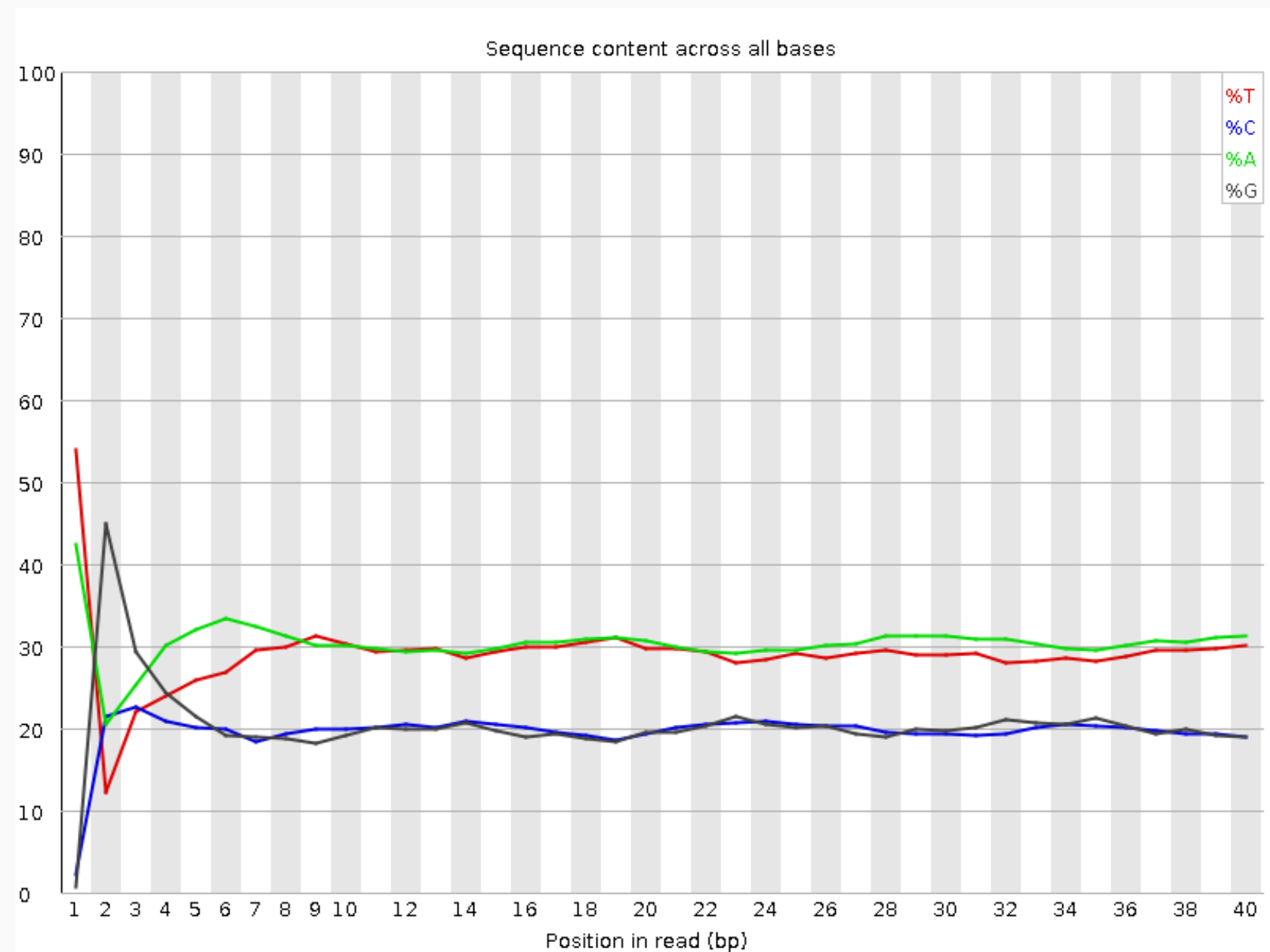
Run	# of Spots	# of Bases	Size	Published
<a href="#">SRR094652</a>	19,949,297	1.6G	725.9Mb	2011-07-20

CC: SRR094649, SRR094650  
3AT: SRR094651, SRR094652

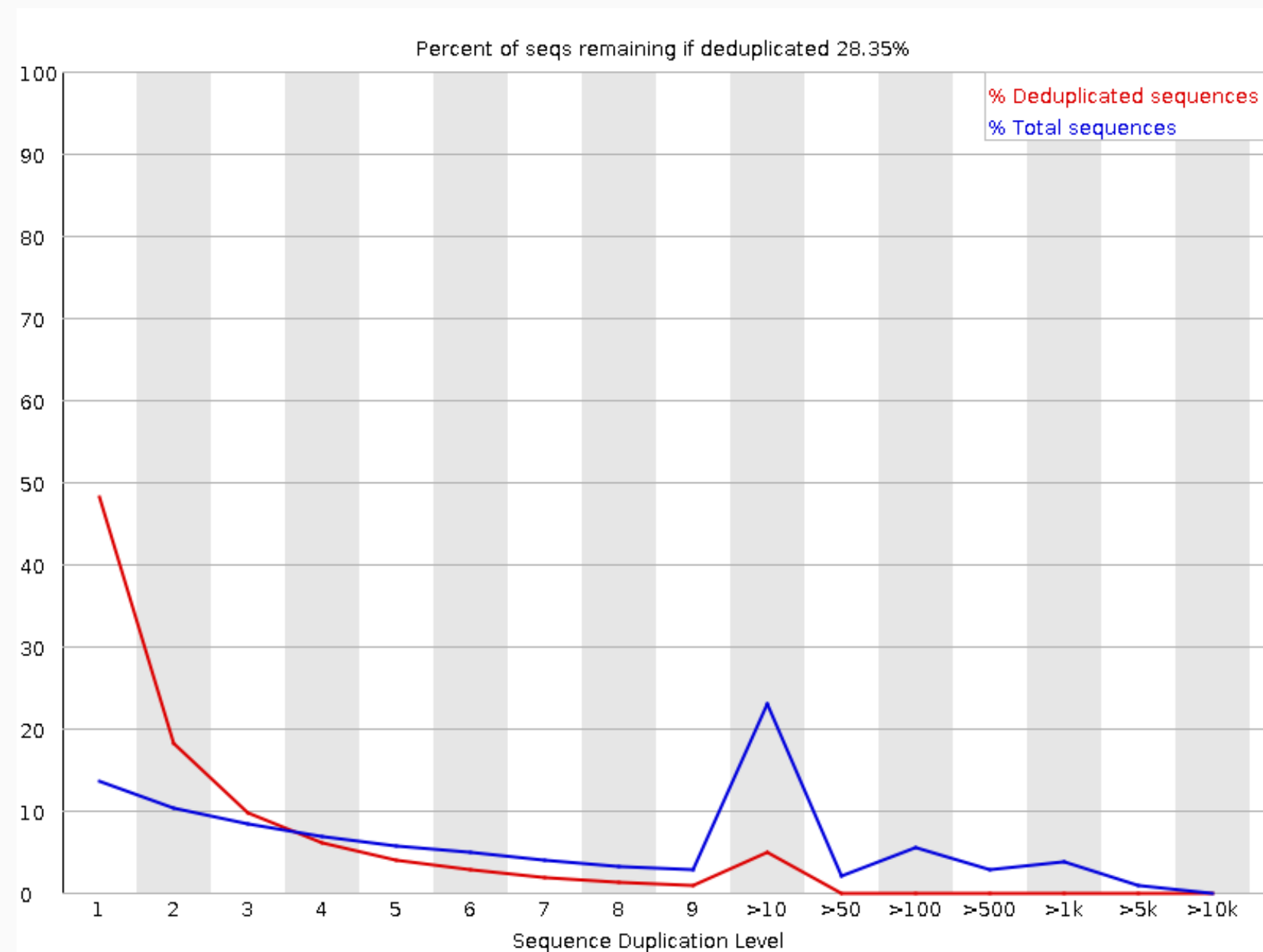


## ◆ Quality check with FastQC

Per base sequence content:



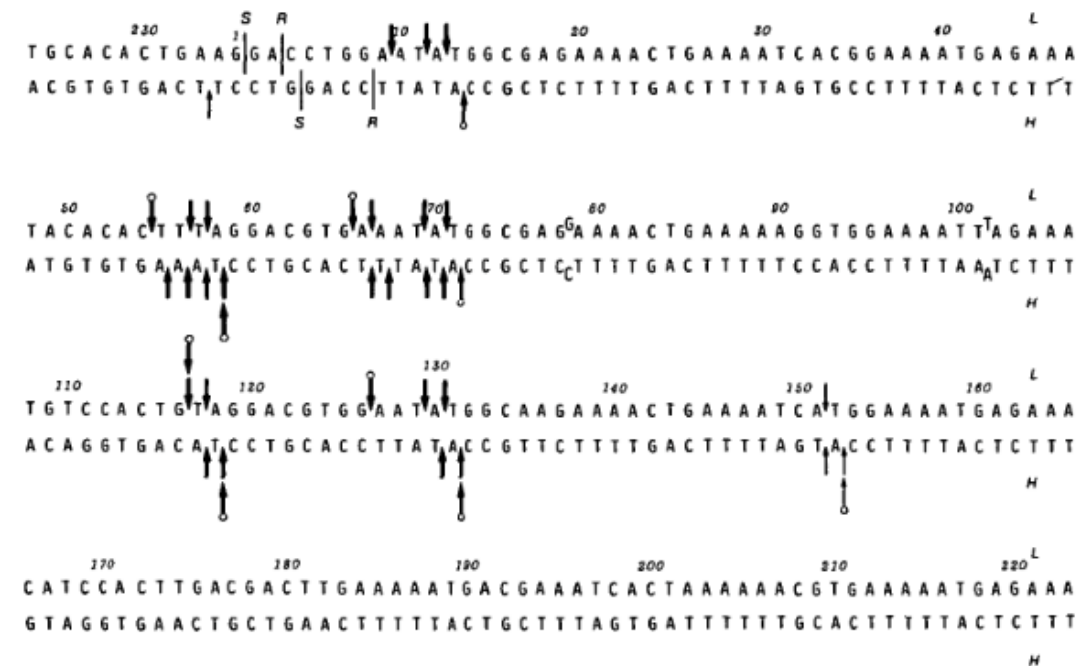
Sequence Duplication Levels:



# Sequence specificity in the digestion of double-stranded DNA by micrococcal nuclease

'...it was concluded that the nuclease cleaves initially in AT-rich regions and it was suggested that this specificity depended on the greater conformational motility of such regions in the DNA...

Nucleic Acids Research



**Figure 3. Regions of preferential attack by micrococcal nuclease on mouse satellite DNA.** The mouse satellite DNA repeat unit of 234 bp (8) is shown with the diester bonds cleaved by micrococcal nuclease (see Table 1) indicated by arrows. Thin arrows indicate low level of cleavage. Sites determined by cleavage of a micrococcal nuclease digests with Sau96I are indicated by  $\blacktriangleright$  while those derived from digestion of the Sau96I monomer with micrococcal nuclease are depicted as  $\circ\blacktriangleright$ . For details see text. The analysis of the region adjacent to the Sau96I site is incomplete since the very short fragments expected to occur in the micrococcal nuclease/Sau96I codigests were not analyzed. L and H designate the light and heavy strand of the satellite DNA. S and R denote the bonds cleaved by Sau96I and EcoRII, respectively. The L-strand (top) is 5' → 3', the complementary H-strand 3' → 5'.

Wolfram Hörz, Werner Altenburger (1981):

"Sequences of the type 5'CATA and 5'CTA are attacked preferentially, followed by exonucleotic degradation at the newly generated DNA termini. GC-rich flanking sequences further increase the probability of initial attack. Unexpectedly, long stretches containing only A and T are spared by the nuclease. These results, which were obtained with mouse satellite DNA and two fragments from the plasmid pBR322, do not support the previous contention that it is the regions of high AT-content which are initially cleaved by micrococcal nuclease. This specificity of micrococcal nuclease complicates its use in experiments intended to monitor the nucleoprotein structure of a DNA sequence in chromatin."

## ◆ HTSeq: Analysing high-throughput sequencing data with Python

The various classes of HTSeq:

–Sequences and FASTA/FASTQ files

In order to represent sequences and reads (i.e., sequences with base-call quality information), the classes *Sequence* and *SequenceWithQualities* are used. The classes *FastaReader* and *FastqReader* allow to parse FASTA and FASTQ files.

–Genomic intervals and genomic arrays

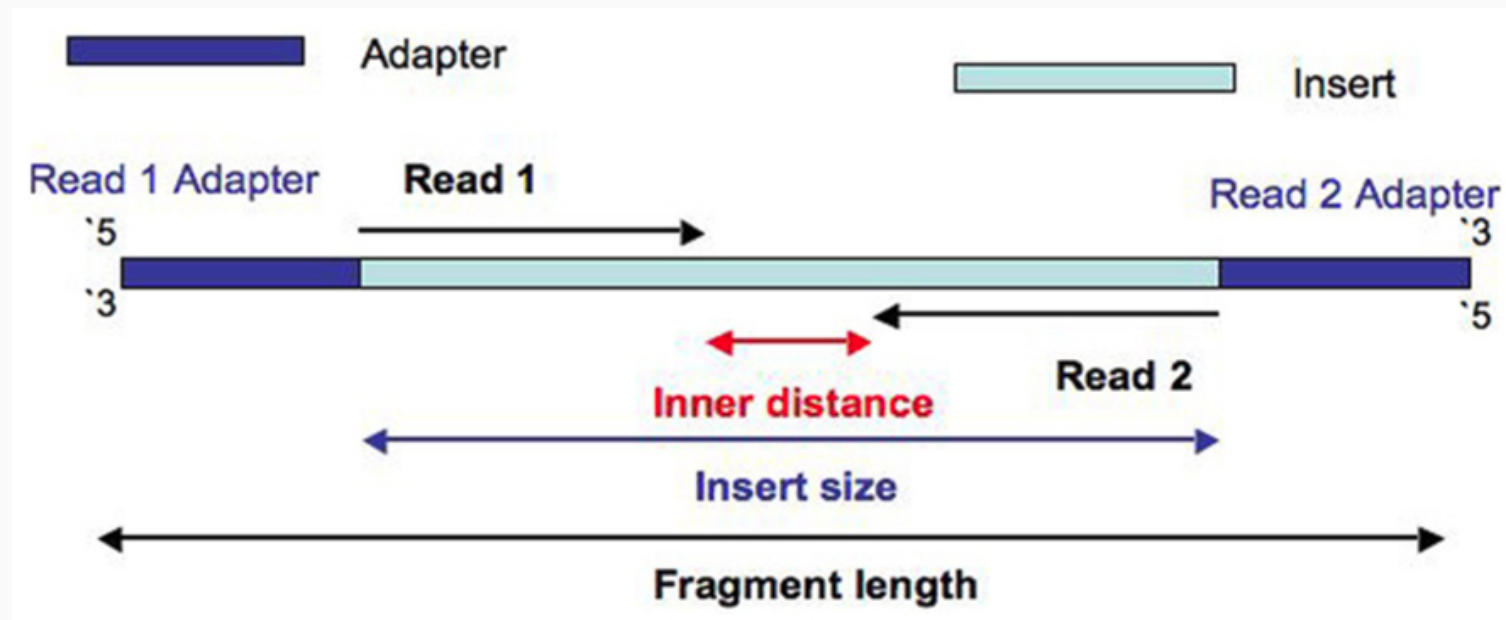
The classes *GenomicInterval* and *GenomicPosition* represent intervals and positions in a genome. The class *GenomicArray* is an all-purpose container with easy access via a genomic interval or position, and *GenomicArrayOfSets* is a special case useful to deal with genomic features (such as genes, exons, etc.)

–Read alignments

To process the output from short read aligners in various formats (e.g., SAM), the classes described here are used, to represent output files and alignments, i.e., reads with their alignment information.

–Features

The classes *GenomicFeature* and *GFF\_Reader* help to deal with genomic annotation data.

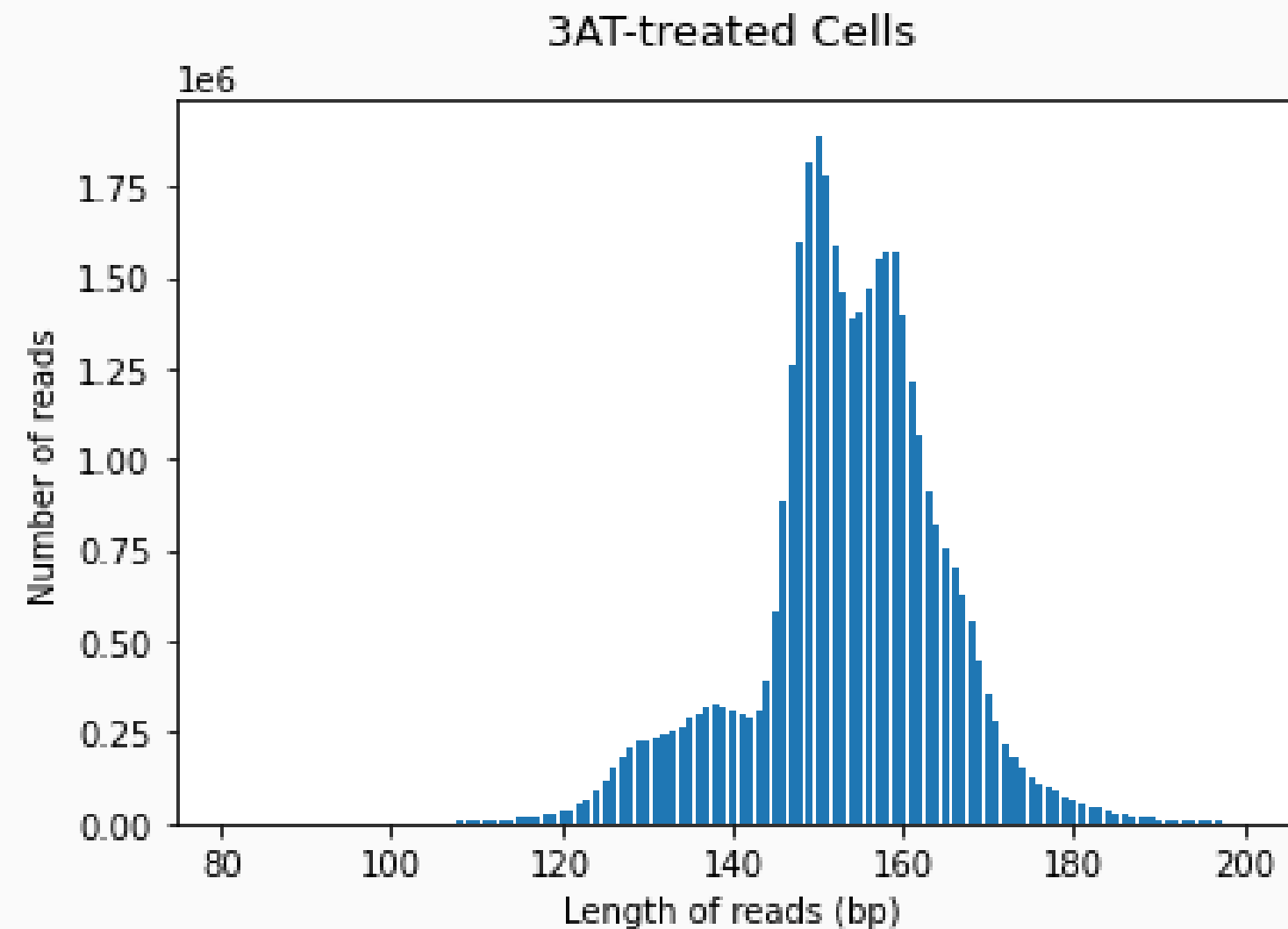
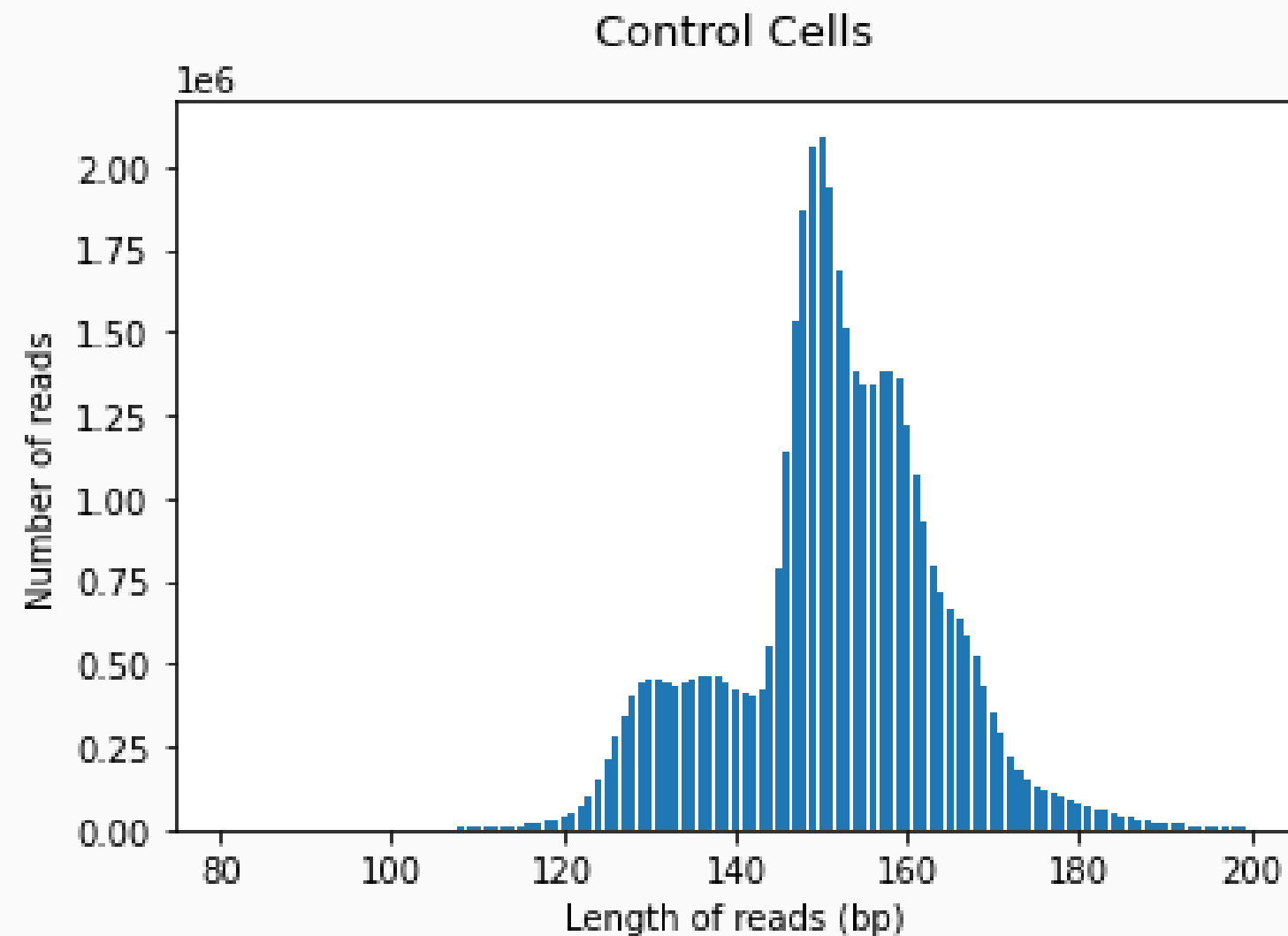


## ◆ Length distribution

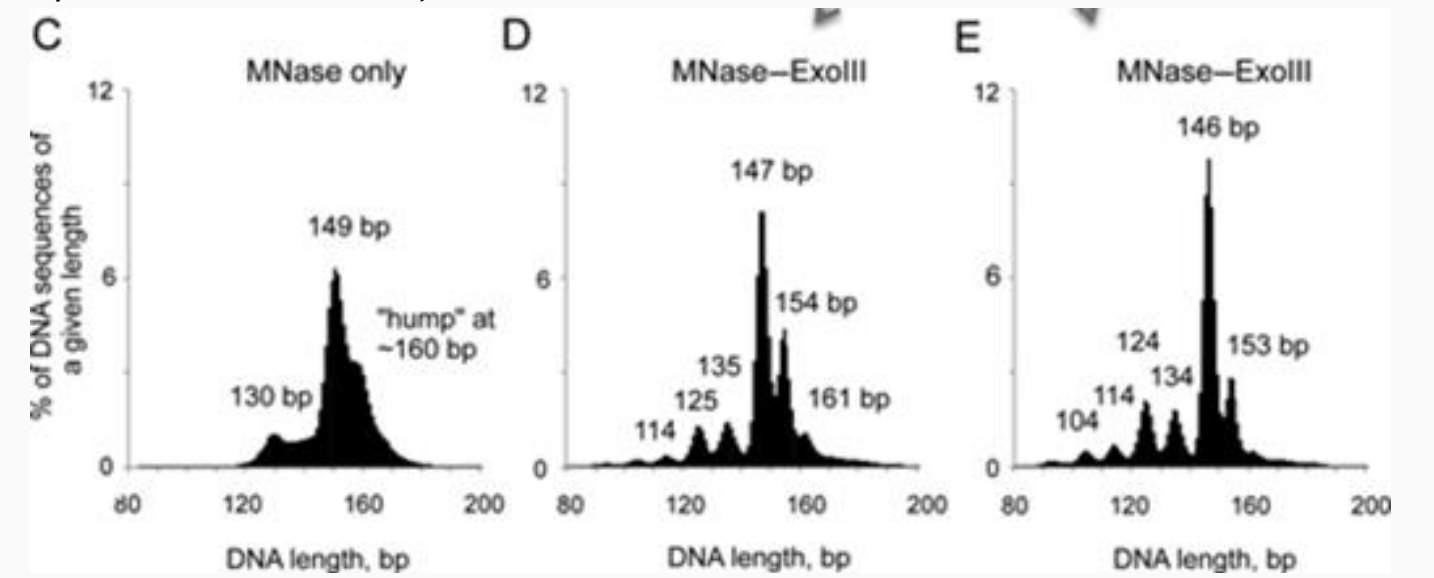
Tool: Jupiter Notebook

Python packages: HTSeq, pysam, matplotlib

Peak on 150 length both in CC and 3AT sets:



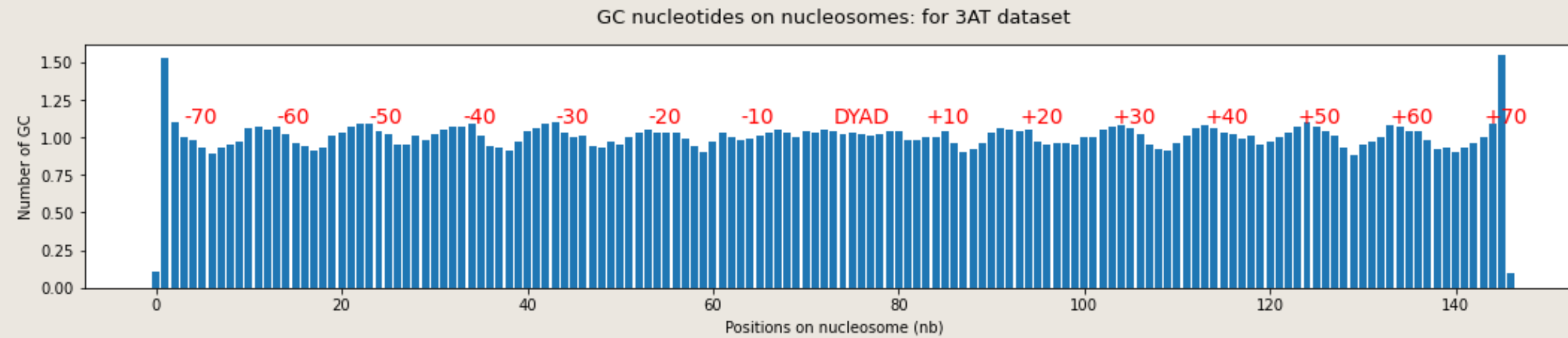
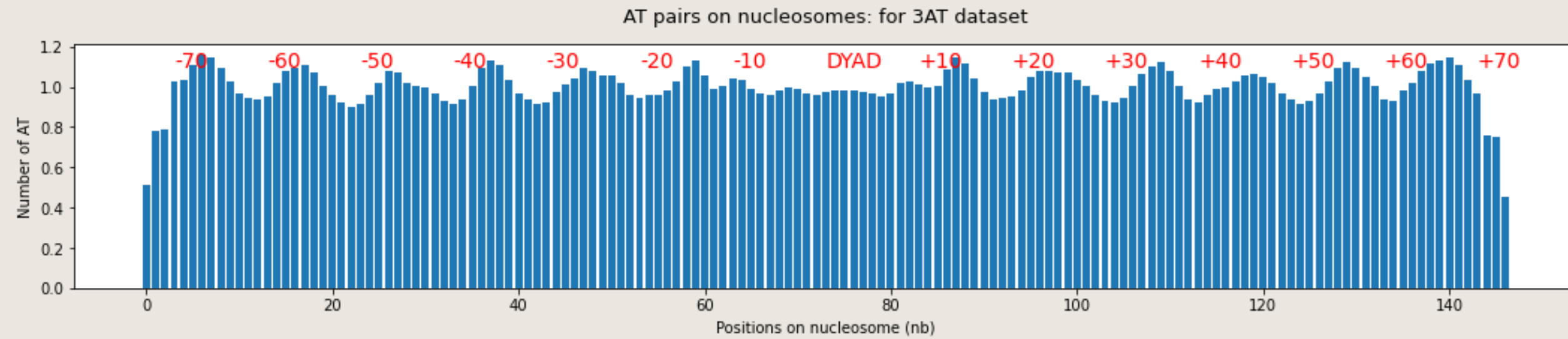
by Cole HA, Cui F, Ocampo J, et al



(C) Length distribution histograms for nucleosome sequences obtained by paired-end sequencing of the mono-nucleosome band from wild type strain JRY4012: MNase-only. (D, E) Length histograms for MNase-ExoIII at two different MNase concentrations

◆ A/T and G/C patterns

Done on dataset of reads with length 147



# Conclusions

- Weak exonuclease activity of MNase
- AT and GC sinusoidal patterns



# References

1. Cole HA, Cui F, Ocampo J, et al. Novel nucleosomal particles containing core histones and linker DNA but no histone H1. *Nucleic Acids Res.* 2016;44(2):573-581.
2. Feng Cui, Hope A. Cole, David J. Clark, Victor B. Zhurkin, Transcriptional activation of yeast genes disrupts intragenic nucleosome phasing, *Nucleic Acids Research*, Volume 40, Issue 21, 1 November 2012, Pages 10753–10764.
3. Segal E, Fondufe-Mittendorf Y, Chen L, et al. A genomic code for nucleosome positioning. *Nature.* 2006;442(7104):772-778. doi:10.1038/nature04979
4. Vladimir B. Teif, Nucleosome positioning: resources and tools online, *Briefings in Bioinformatics*, Volume 17, Issue 5, September 2016, Pages 745–757
5. Hope A. Cole, Bruce H. Howard, David J. Clark. The centromeric nucleosome of budding yeast is perfectly positioned and covers the entire centromere. *Proceedings of the National Academy of Sciences* Aug 2011, 108 (31) 12687-12692
6. Colin Dingwall, George P. Lomonosoff, Ronald A. Laskey, High sequence specificity of micrococcal nuclease, *Nucleic Acids Research*, Volume 9, Issue 12, 25 June 1981, Pages 2659–2674
7. Wolfram Hörz, Werner Altenburger, Sequence specific cleavage of DNA by micrococcal nuclease, *Nucleic Acids Research*, Volume 9, Issue 12, 25 June 1981, Pages 2643–2658