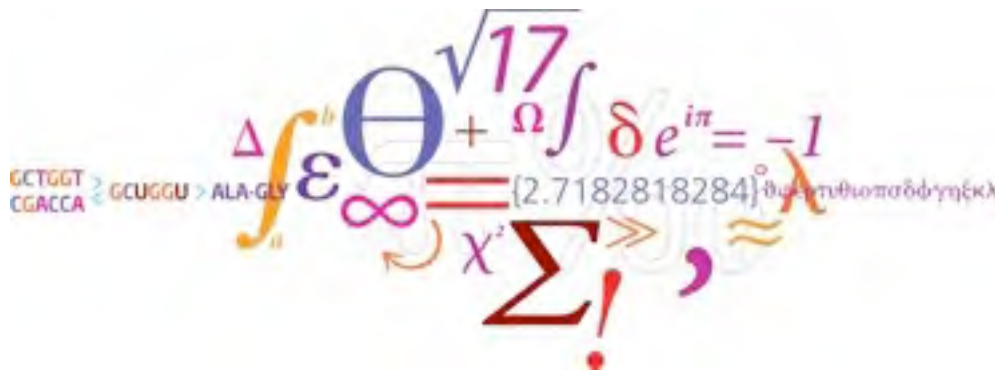




# Comparative Genomics

## Introduction to DNA atlases (DNA structures)



Dave Ussery

Genomics of Prokaryotes Workshop  
Universidad Miguel Hernández  
Alicante, Spain

Monday, 12 December, 2011





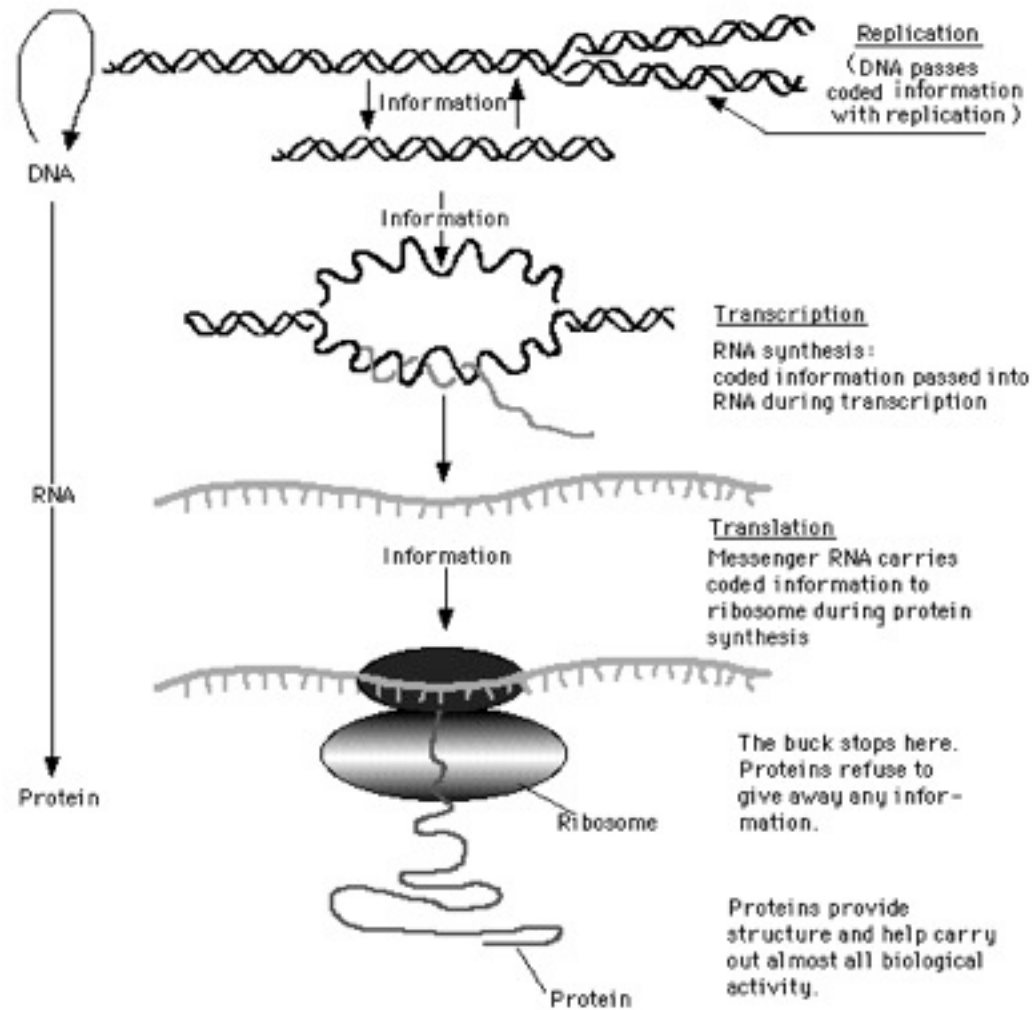
THE OXFORD  
 ENGLISH  
 DICTIONARY

**bioinformatics, n.**

The science of information and information flow in biological systems, esp. the use of computational methods in genetics and genomics.

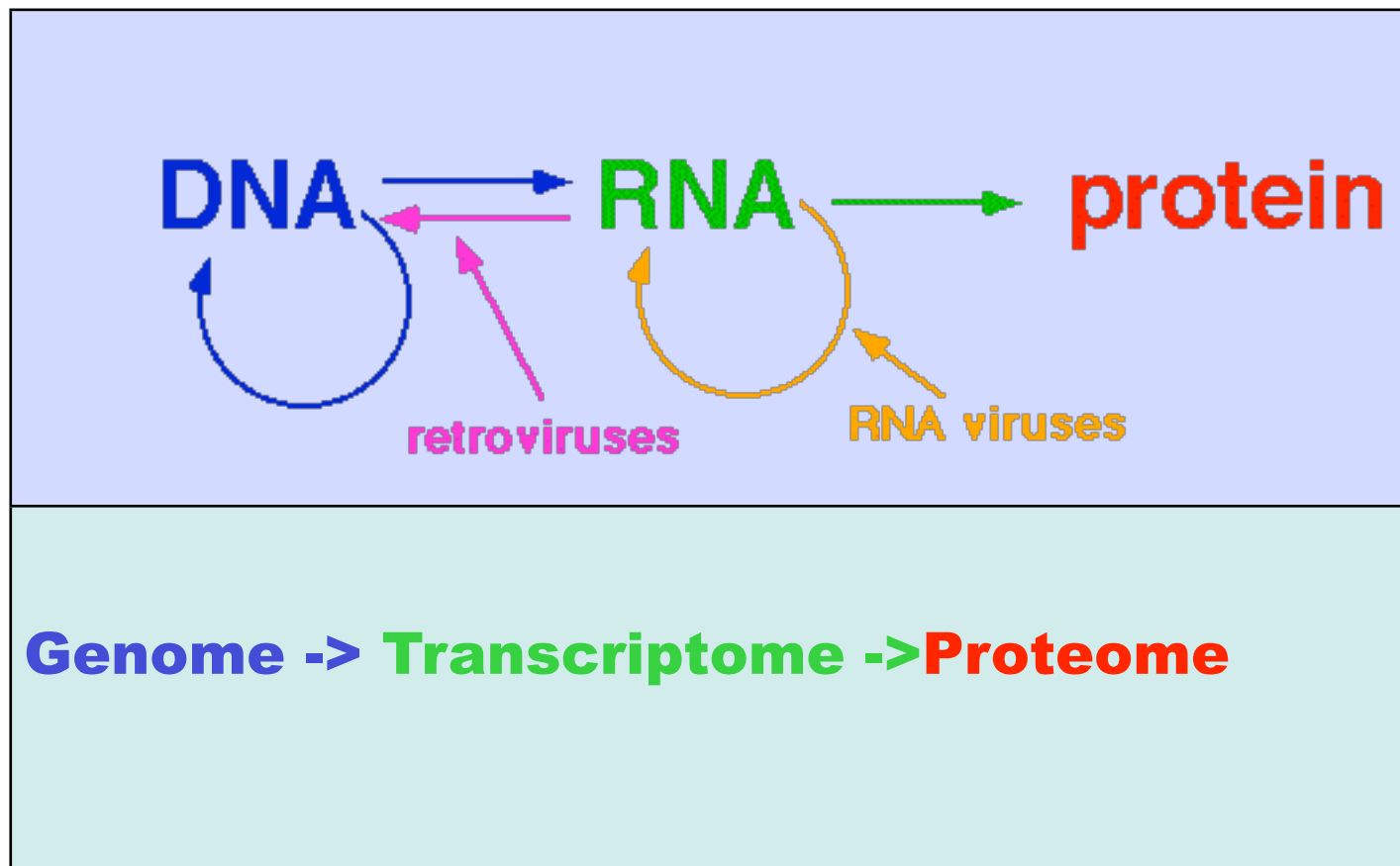
**1978** P. HOGEWEG in *Simulation* **31** 90/1 Since 1970 she has been a staff member at the Subfaculty of Biology of the University of Utrecht, with her main field of research in bioinformatics. **1985** *Jrnl. Theoret. Biol.* **113** 719 (*heading*) Tumor escape from immune elimination... R. J. De Beer, Bioinformatics Group, University of Utrecht. **1986** *Philos. Trans. Royal Soc. A.* **317** 324 The area of modelling mutants from a known structure has been revolutionized by the latest tools of molecular graphics... This is a key element in the whole technology and has attracted much interest (for example, the recent E.E.C. 'Bioinformatics' programme). **1987** *Science* 4 Sept. 1108/3 One of the latest developments [at the European Molecular Biology Laboratory] has been the creation of a new research program in bioinformatics. This is intended to bring together research in computing science, structural biology, and molecular genetics. **1996** *Fast Company* Aug.-Sept. 32/3 A lot of breakthroughs in medicine will come out of the efforts of bioinformatics. **2001** *N.Y. Times* 4 Jan. B6/2 The hope..is to make New York a leader in cutting-edge fields like bioinformatics, in which computers are used to decipher genes and proteins.

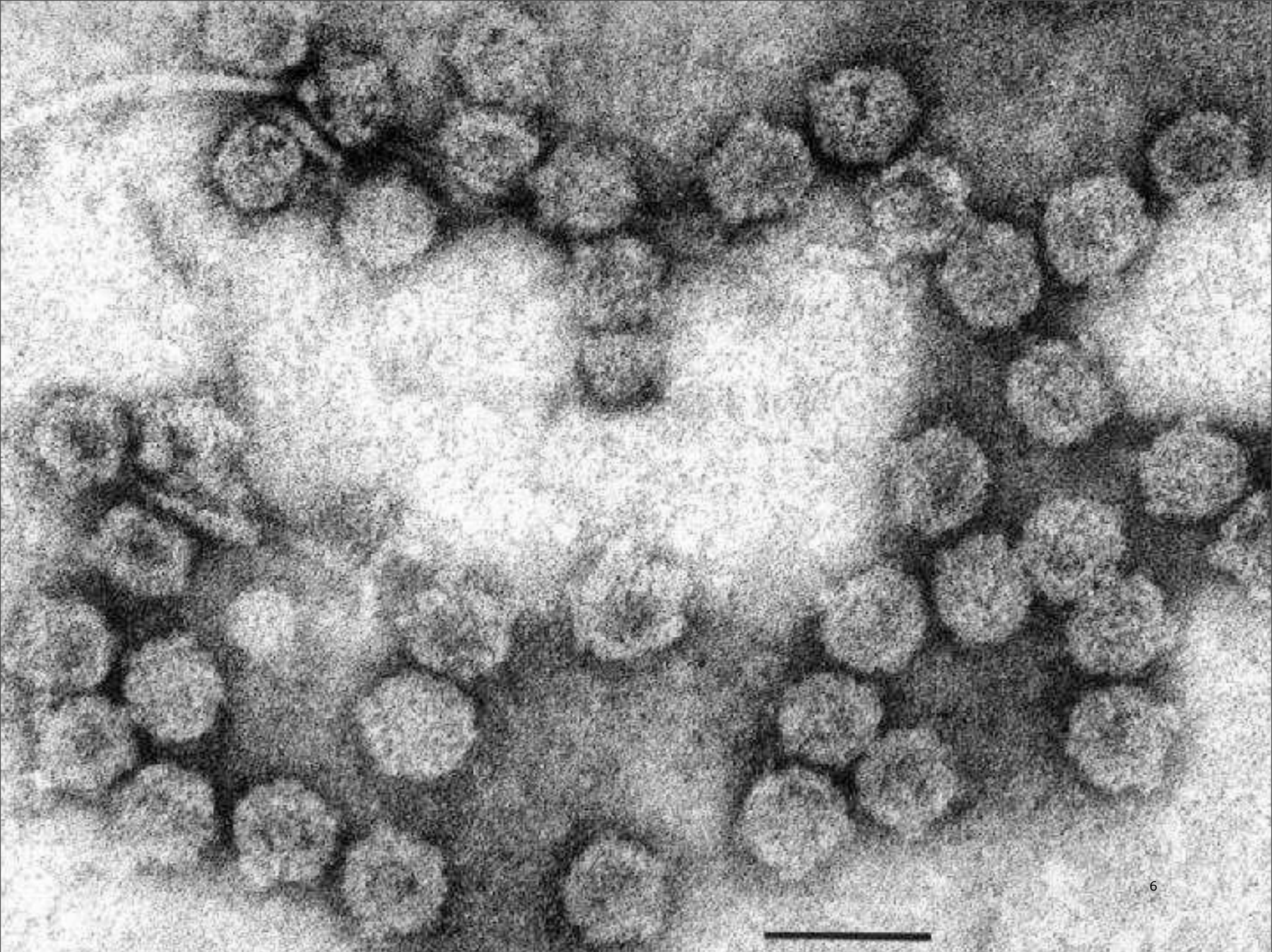
## The Central Dogma of Molecular Biology



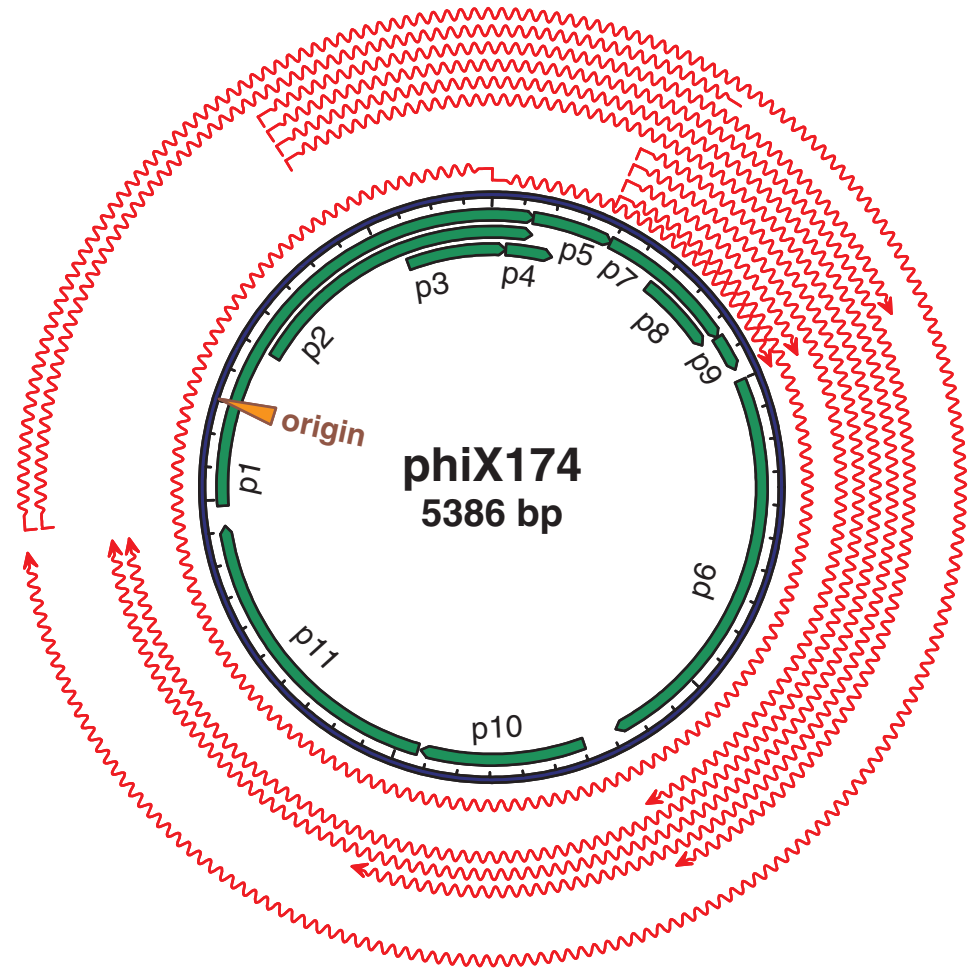
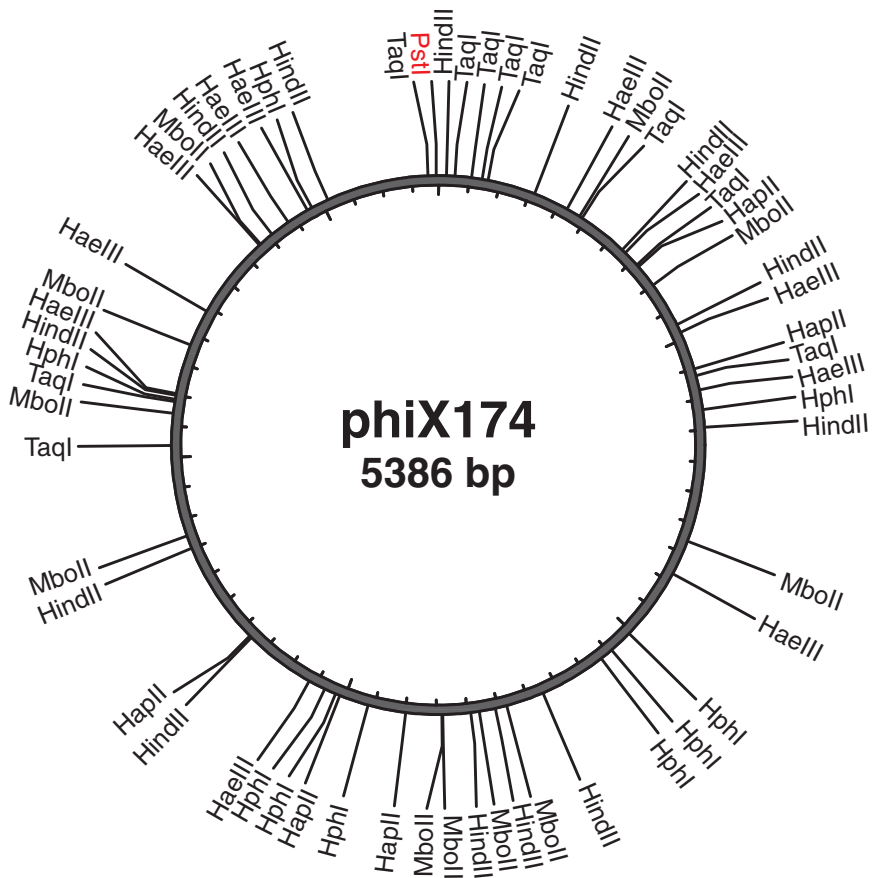


# What is Biological Information?

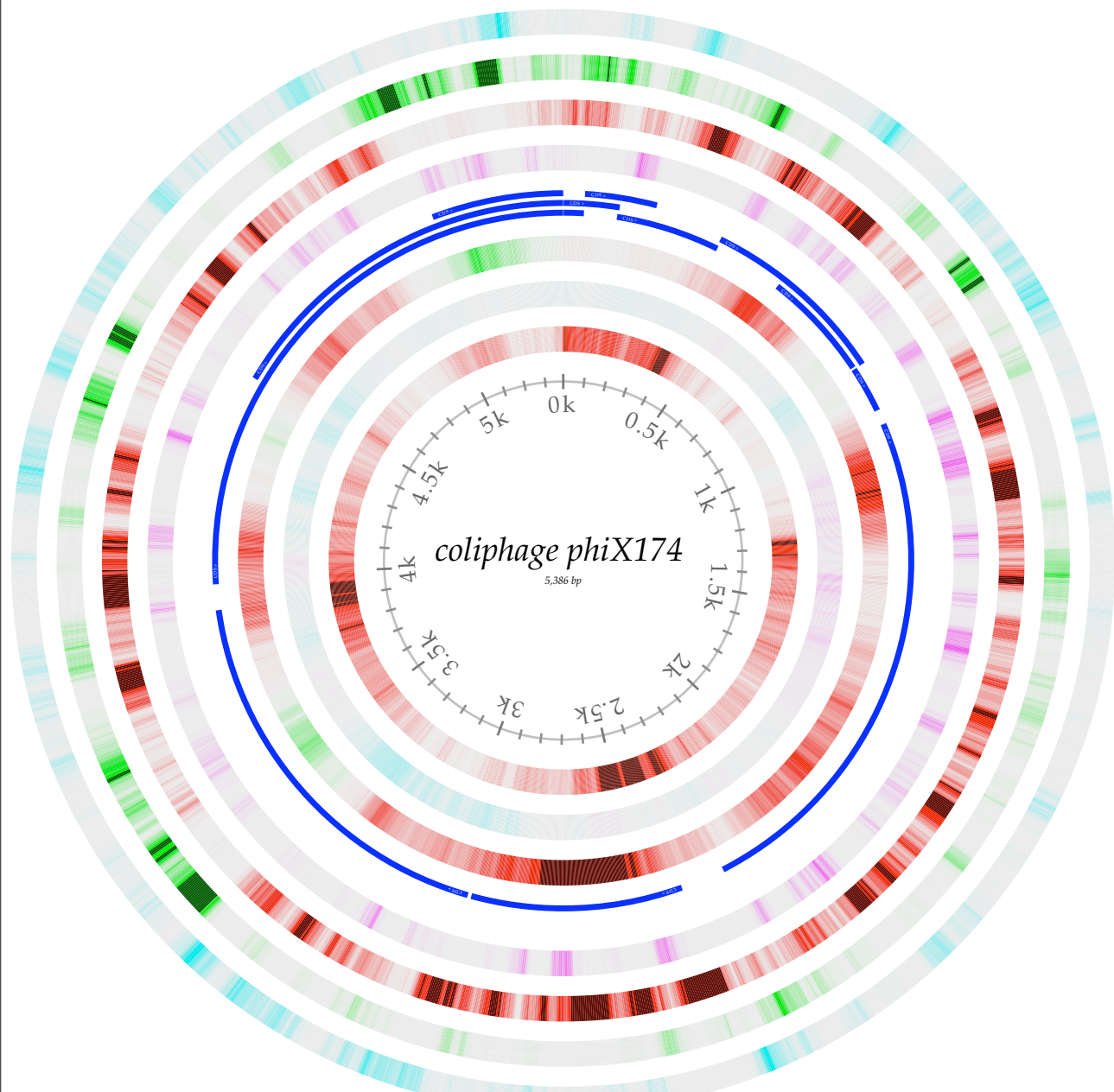




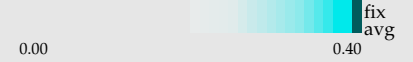
gagttttatc gcttccatga cgcagaagtt aacactttcg gatattttctg atgagtcgaa aaattatctt gataaagcag gaattactac tgcttgttta cgaattaaat cgaagtggac  
tgctggcggg aatgagaaa attcgaccta tccttgcgca gctcgagaag ctcttacttt gcgacctttc gccatcaact aacgattctg tcaaaaactg acgcttgga tgaggagaag  
tggttaata tgcttggcac gttcgtcaag gactggttta gatatgagtc acattttggt catggtagag attctcttgt tgacatttta aaagagcgtg gattactatc tgagtccgat  
gctgttcaac cactaatagg taagaaatca tgagtcaagt tactgaacaa tccgtacggt tccagaccgc tttggcctct attaagctca ttcaggcttc tgccgttttg gatttaaccg  
aagatgattt cgattttctg acgagtaaca aagtttggat tgctactgac cgctctcgtg ctctcgtcgt cgttgaggct tgcgtttatg gtacgctgga ctttgtggga taccctcgct  
ttcctgctcc tgttgagttt attgctgccg tcattgctta ttatgttcat cccgtcaaca ttcaaacggc ctgtctcatt atggaaggcg ctgaatttac ggaaaacatt ataatggcg  
tcgagcgtcc ggtaaagcc gctgaattgt tcgctgttac cttgctgtga cgcgcaggaa acactgacgt tcttactgac gcagaagaaa acgtgctca aaaattacgt gcggaaggag  
tgatgtaatg tctaaaggta aaaaacgttc tggcgtcgc cctggctgct cgcagccggt gcgaggctact aaaggcaagc gtaaaggcgc tcgtctttgg tatgtaggtg gtcaacaatt  
ttaattgcag gggcttcggc cccttacttg aggataaatt atgtctaata ttcaaacggc atgcccgtat acctttccca tcttggcttc cttgctggct agattggctg  
tcttattacc atttcaacta ctccggttat cgctggcgac tccttcgaga tggacgccgt tggcgtctc cgtctttctc cattgctgct tggccttgct attgactcta ctgtagacat  
ttttactttt tatgtccctc atcgtcacgt ttatggtgaa cagtggatta agttcatgaa ggatgggtt aatgccactc ctctcccgc tgtaaacact actggttata ttgacctgc  
cgttttctt ggcacgatta accctgatac caataaaatc cctaagcatt tgtttcaggg ttatttgaat atctataaca actattttaa agcgcctgg atgcctgacc gtaccgaggc  
taacccta atgacatta tgggtctgca agctgcttat gctaatttgc atactgacca agaacgtgat tacttcatgc agcgttacca tgatgttatt tcttcatttg gaggtaaac  
taccacatct attgacatta tgggtctgca agctgcttat gctaatttgc atactgacca agaacgtgat tacttcatgc agcgttacca tgatgttatt tcttcatttg gaggtaaac  
ctcttatgac gctgacaacc gtcctttact tgtcatgctc tctaactctt gggcatctgg ctatgatgtt gatggaactg accaaacgctc gttaggccag ttttctggct gtgttcaaca  
gacctataaa cattctgtgc cgcgtttctt tgttctgag catggcacta tgtttactct tgcgcttgtt cgttttccgc ctactgctgac taaagagatt cagtacctta acgctaaagg  
tgctttgact tataccgata ttgctggcga ccctgttttg tatggcaact tgccgcccgc tgaatttct atgaaggatg ttttccgctc tgggtgattcg tctaagaagt ttaagattgc  
tgagggtcag tggatcgtt atgcgccttc gttcagttg ttgcagtgga atagtcagggt taaattta atgaccgtt atcgcaatct gccgaccact cgcgattcaa tcatgacttc  
ccaccatgat tatgaccagt gttccagtc cgttcagttg ttgcagtgga atagtcagggt taaattta atgaccgtt atcgcaatct gccgaccact cgcgattcaa tcatgacttc  
gtgataaaag attgagtgtg aggttataac gccgaagcgg taaaaatctt aattttgc gctgaggggt tgaccaagcg aagcgcggta ggtttctgc ttaggagttt aatcatgttt  
cagactttta tttctgcca taattcaaac ttttttctg ataagctggt tctcactctt gttactccag ctcttctggc acctgtttta cagacaccta aagctacatc gtcaacgta  
tattttgata gtttgacggg taatgctggt aatgggtggt ttcttcattg cattcagatg gatacatctg tcaacgccgc taatcagggt gtttctgttg gtgctgatat tgctttgat  
gccgacccta aatttttgc ctggttggtt cgctttgagt ctcttctggg tccgactacc ctcccgactg cctatgatgt ttatcctttg aatggctgcc atgatgggtg ttattatacc  
gtcaaggact gtgtgactat tgacgtcctt ccccgtagc cgggcaataa cgtttatgtt ggtttcatgg tttggctaa ctttaccgct actaaatgcc gcggattggg ttcgctgaat  
aagagattat ttgtctccag ccacttaagt gaggtgattt atgtttgggt ctattgctgg cgttatgtct tctgctcttg ctggtgccgc catgtctaaa ttgtttggag gcggtcaaaa  
agccgcctcc ggtggcattc aaggtgatgt gcttgctacc gataacaata ctgtaggcat ggggtgatgct ggtattaaat ctgccattca aggcctaat gttcctaacc ctgatgaggc  
cgcccctagt tttgttctg gtgctatggc taaagctggt aaaggacttc ttgaaggtag gttgcaggct ggcacttctg ccgttctga taagttgctt gatttgggtg gacttgggtg  
caagtctgcc gctgataaag gaaaggatac tcgtgattat ctgctgctg catttctga gcttaatgct tgggagcgtg ctggtgctga tgcttctct gctggatgg ttgacgccgg  
atltgagaat caaaaagagc ttactaaaat gcaactggac aatcagaag agattgccga gatgcaaaat gagactcaaa aagagattgc tggcattcag tcggcgactt cacgccagaa  
tacgaaagac caggatgatg cacaaaatga gatgcttgc tatcaacaga aggagtctac tgctcgcgtt gcgtctatta tggaaaacac caatcttcc aagcaacagc aggttccga  
gattatgctc caaatgctta ctcaagctca aacggctggt cagtatttta ccaatgacca aatcaaaagaa atgactcgca aggttagtgc tgaggttgac ttagtctatc agcaaacgca  
gaatcagcgg tatggctctt ctcatattgg cgctactgca aaggatattt ctaatgtcgt cactgatgct gcttctgggt tggttgatatt ttttcatggg attgataaag ctggtgccga  
tacttggaac aatcttggga aagacggtaa agctgatggt attggctcta atttgtctag gaaataaccg tcaggattga cacctccca attgtatgtt tcatgcctc caaatcttg  
aggctttttt atggttcgtt ctatattacc ttctgaatgt cacgctgatt attttgactt tgagcgtatc gaggctctta aacctgctat tgaggcttgt ggcatttcta ctcttctca  
atccccatg ctggcttcc ataagcagat ggataaccgc atcaagctct tggaaagat tctgtctttt cgatgacagg gcgttgagtt cgataatggg gatattgatg ttgacggcca  
taaggctgct tctgacgttc gtgatgagtt tgtatctggt actgagaagt taatggatga atttggcaca tgctacaatg tgctccccca acttgatatt aataacacta tagaccaccg  
ccccgaagg gacgaaaaat ggtttttaga gaacgagaag acggttacgc agttttgccc caagctggct gctgaacgcc ctcttaagga tattcgcat gagtataatt acccaaaaa  
gaaaggattt aaggatgagt gttcaagatt gctggaggcc tccactatga aatcgcgtag aggccttctt attcagcgtt tgatgaaatgc aatgcgacag gctcatgctg atgggtggtt  
tatcgttttt gacactctca cgttggctga cgaccgatta gaggcgtttt atgataatcc caatgctttg cgtgactatt tctgtgatatt tggctgctg gttcttctg ccgagggtcg  
caaggcta at gattcacacg ccgactgcta tcagtatttt tgtgtgcctg agtatggtac agctaattggc cgtctcatt tccatgagggt gcactttatg cggacacttc ctacaggtag  
cgttgaccct aattttggc gtcgggtacg caatcgccgc cagttaaata gcttgcaaaa tacgtggcct tatggttaca gtatgccat gcagttcgc tacacgcagg acgcttttctc  
acgttctggt tggttgtggc ctggtgatgc taaaggtag ccgcttaaac ctaccagtta tatggctggt ggtttctatg tggctaaata cgtaacaaa aagtcagata tggacctgctc  
tgctaaagg ctaggagcta aagaatggaa caactcacta aaaaccaagc tgctcgtact tcccaagaag ctgttcagaa tcagaatgag ccgcaacttc gggatgaaaa tgctcacaat  
gacaaatctg tccacggagt gcttaatcca acttaccagg ctgggttacg acgcgacgcc gttcaaccag atattgaagc agaacgcaaa aagagagatg agattgaggc tgggaaaagt  
tactgtagcc gacgttttgg cggcgaacc tgtgacgaca aatctgctca aatttatgctg cgcttcgata aaaaatgattg gcgtatccaa cctgca



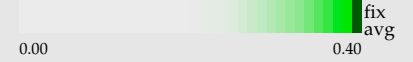




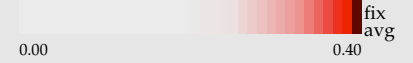
G Content



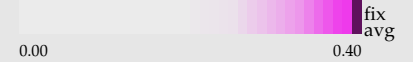
A Content



T Content



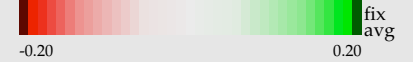
C Content



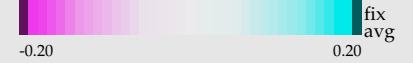
Annotations:



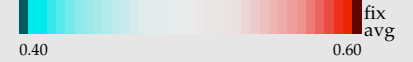
AT Skew



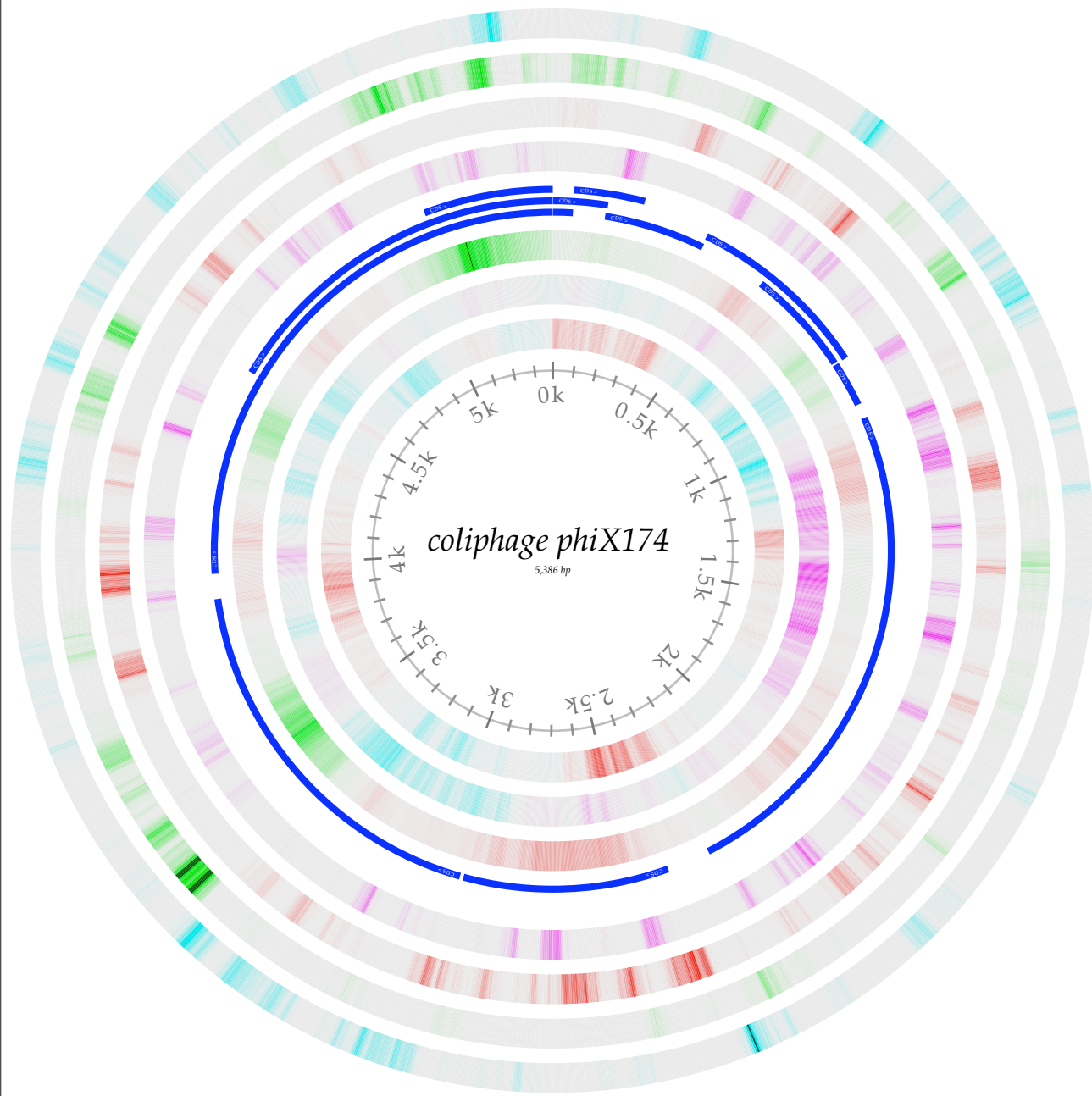
GC Skew



Percent AT



Resolution: 3



**G Content**  
0.07 dev avg 0.39

**A Content**  
0.01 dev avg 0.47

**T Content**  
0.10 dev avg 0.53

**C Content**  
0.04 dev avg 0.39

**Annotations:**  
 CDS +

**AT Skew**  
-0.33 dev avg 0.18

**GC Skew**  
-0.10 dev avg 0.14

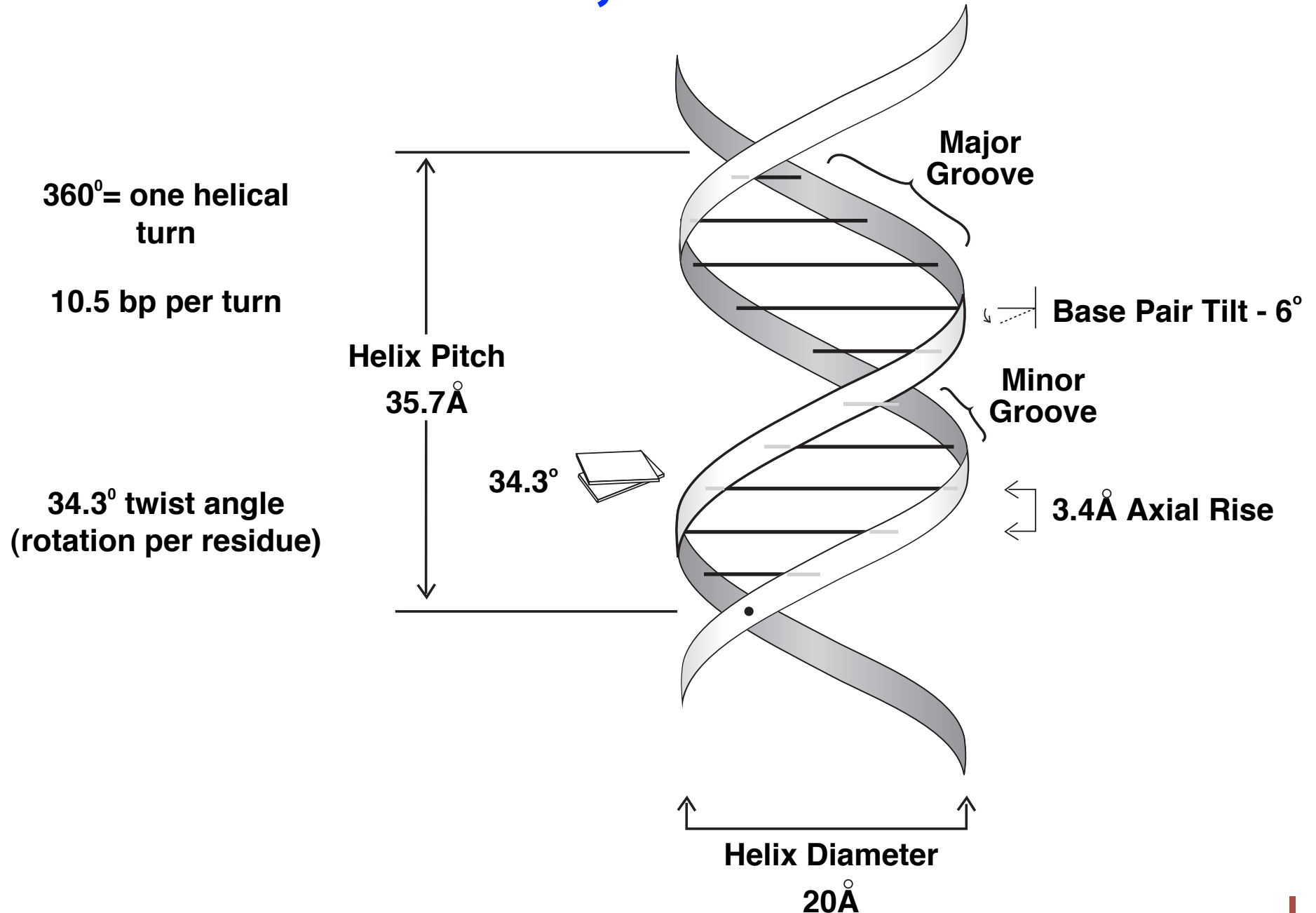
**Percent AT**  
0.47 dev avg 0.63

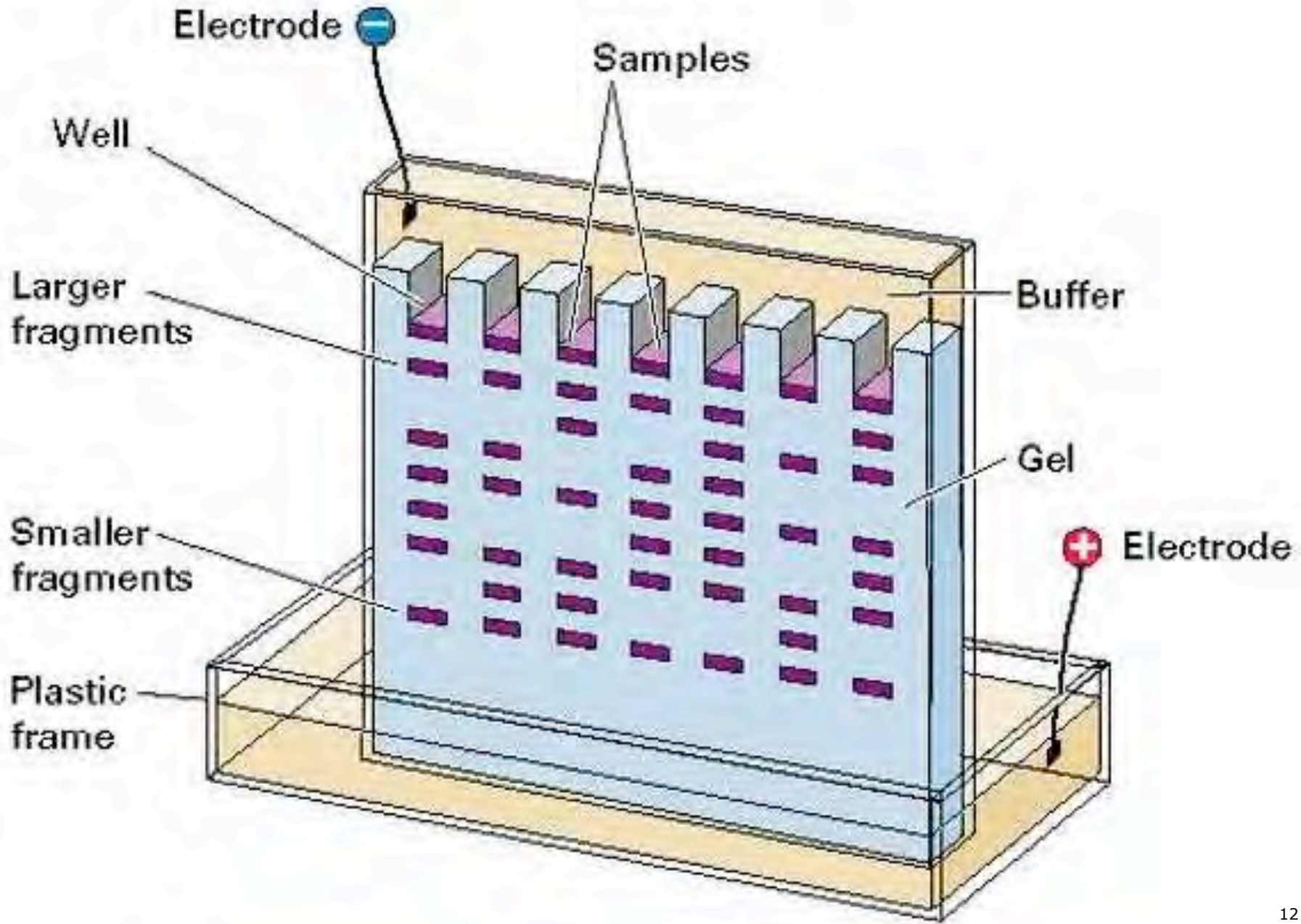
Resolution: 3

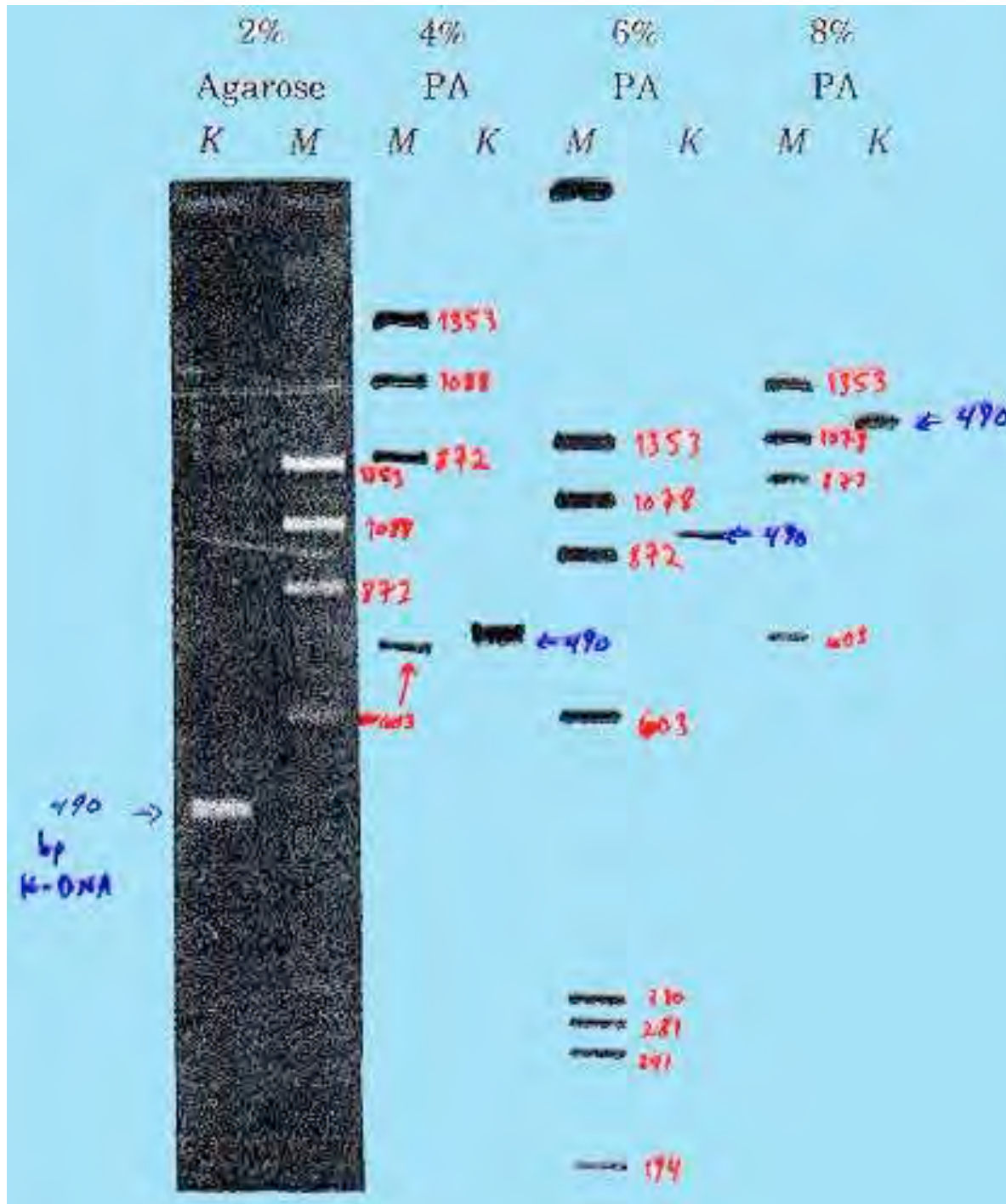
<http://www.cbs.dtu.dk/>  
Center for Biological Sequence Analysis

12 December, 2011  
BASE ATLAS

# a digression - a few words about DNA structures, DNAatlases....

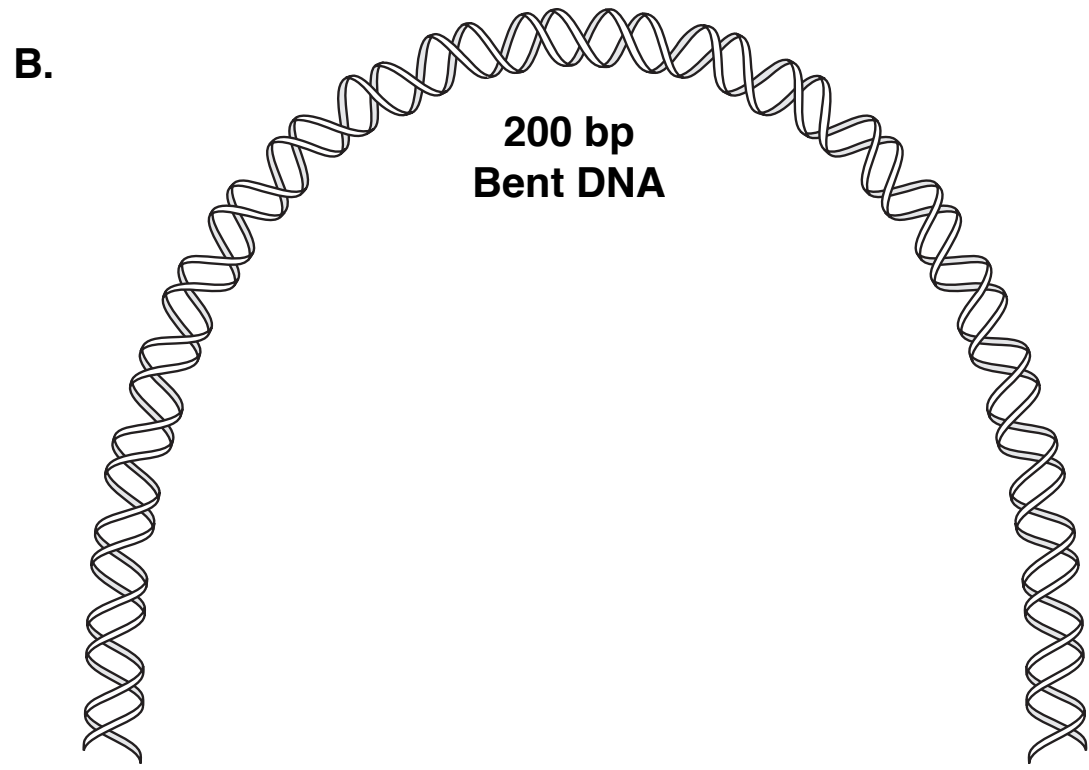








**200 bp  
Straight DNA**



**200 bp  
Bent DNA**

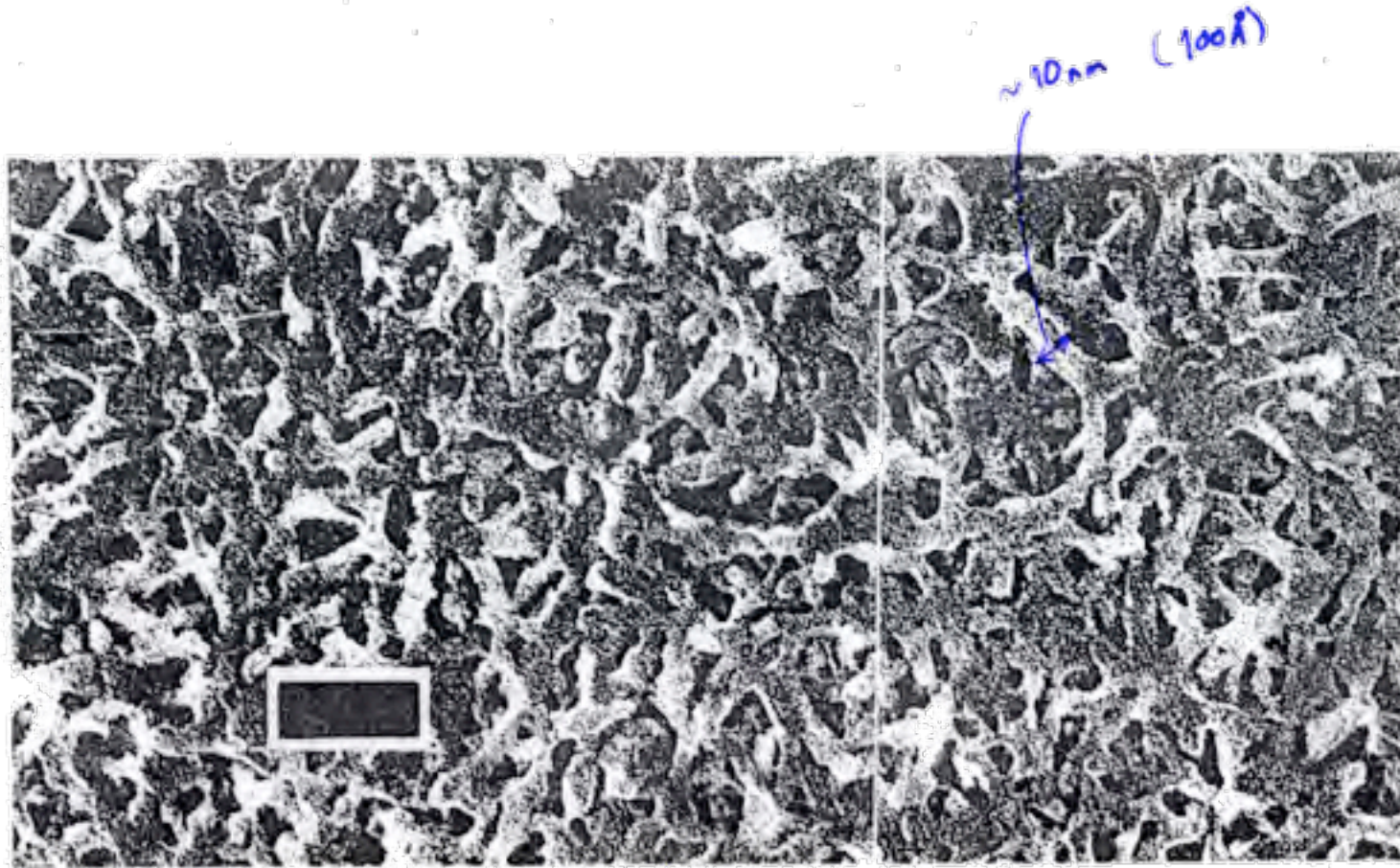
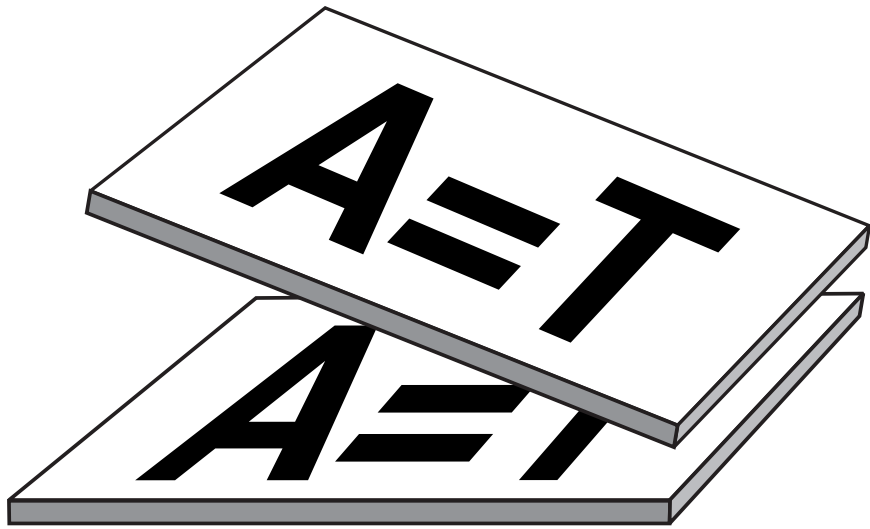


Figure 8.10 Electron micrograph of a portion of a 2% agarose gel,  $1\ \mu\text{m} \times 0.5\ \mu\text{m}$  overall: small black rectangle is  $1000\ \text{Å} \times 500\ \text{Å}$ . Individual gel fibers are about  $100\ \text{Å}$  wide. Courtesy of Sue Whytock and John Finch.



**Tilt**

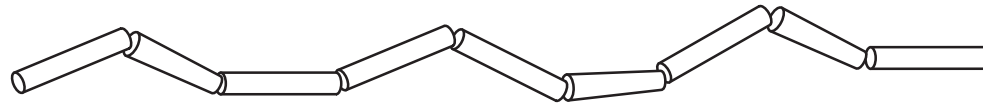


**Roll**

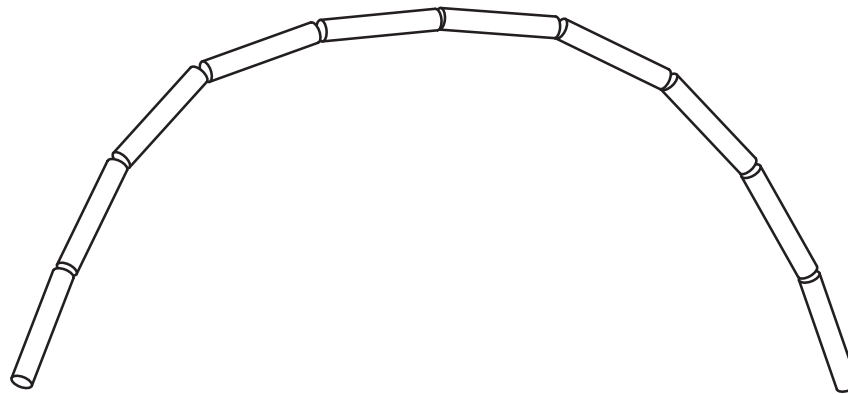




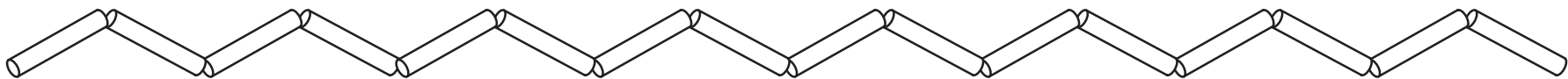
### Random Bends



### Bends Phased at 10.5 bp



### Bends Phased at 16 bp





1

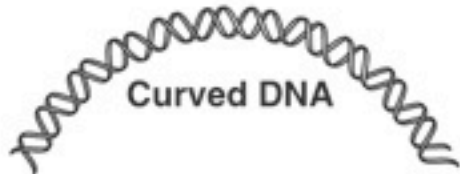
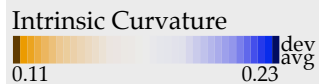
11

21

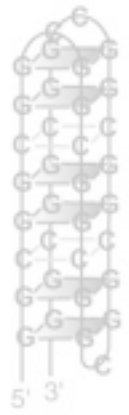
31

41

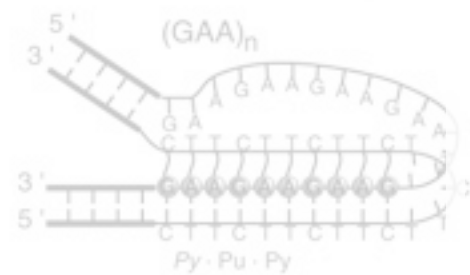
·                    ·                    ·                    ·                    ·  
CCCAAAATGTCAAAAAATAGGCAAAAAATGCCAAAAATCCCAAC



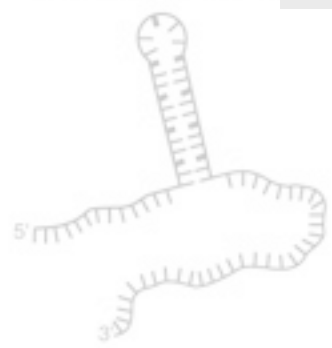
**CGG quadruplex**



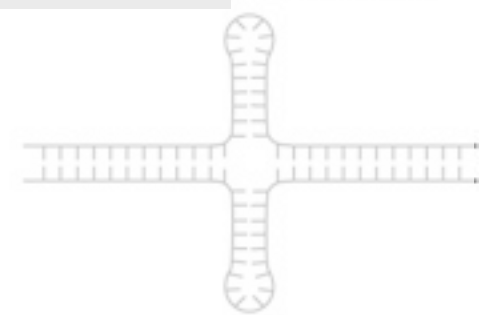
**Hy3-type intramolecular triplex**



**Hairpin structure**



**Y structure**



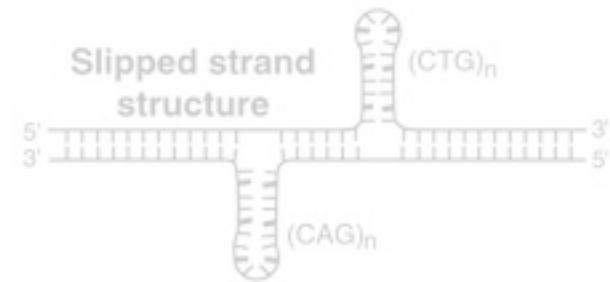
**Anti Parallel DNA**



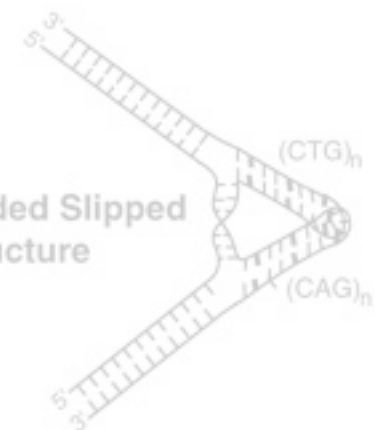
**Parallel DNA**



**Slipped strand structure**

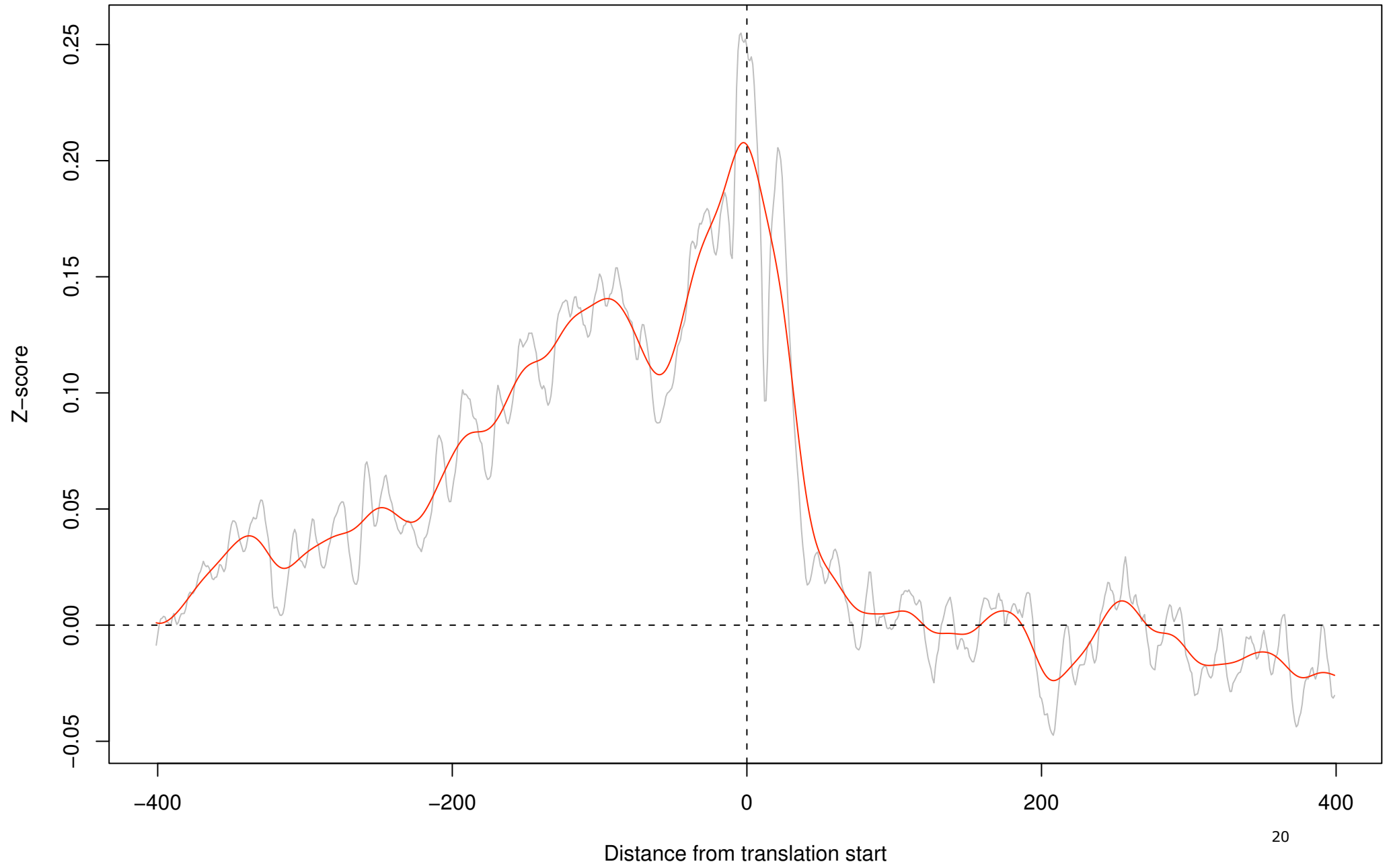


**Folded Slipped Structure**



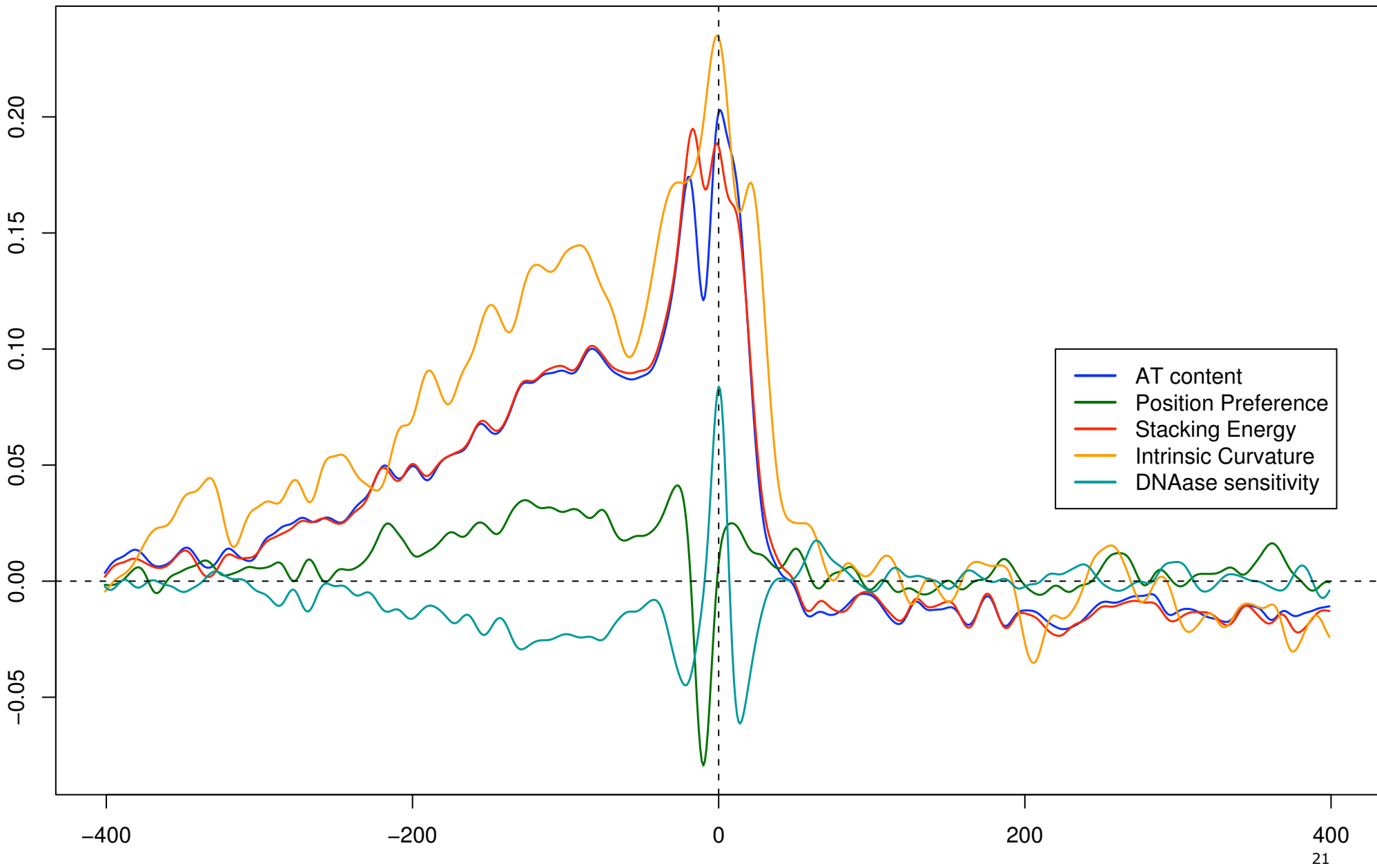


Ecoli\_K-12\_W3110\_Main: **CURVATURE**





### Ecoli\_K-12\_W3110\_Main: Structural Profile

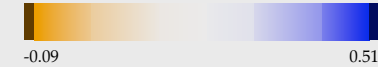


- AT content
- Position Preference
- Stacking Energy
- Intrinsic Curvature
- DNAase sensitivity

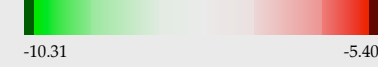
- Ori -



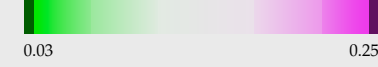
Intrinsic Curvature



Stacking Energy



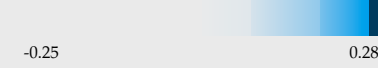
Position Preference



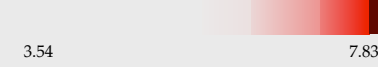
Annotations:



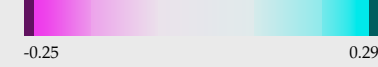
Perfect-palindromes



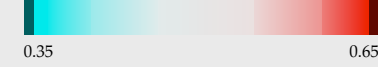
Inverted Repeats



GC Skew



Percent AT

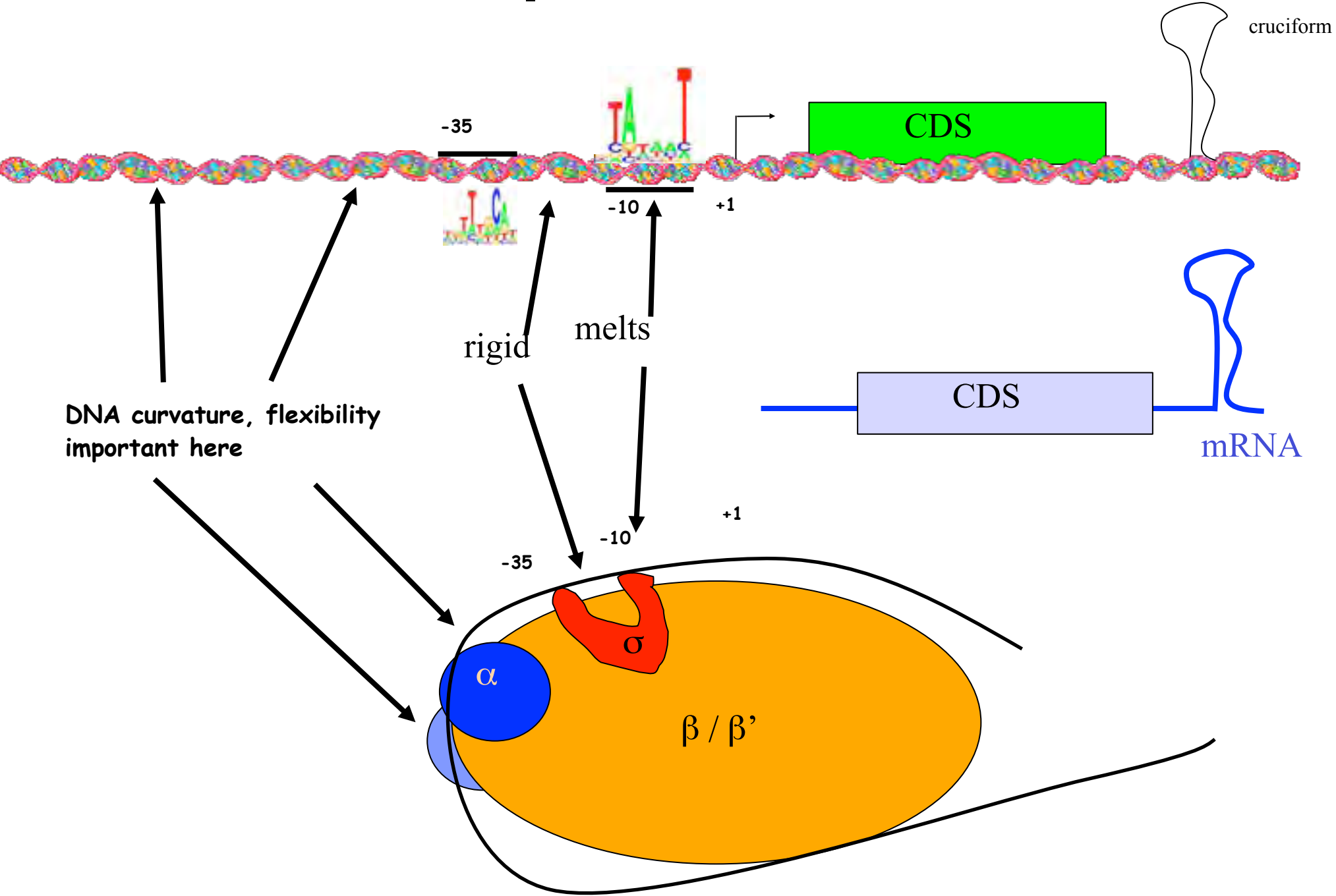


Resolution: 3

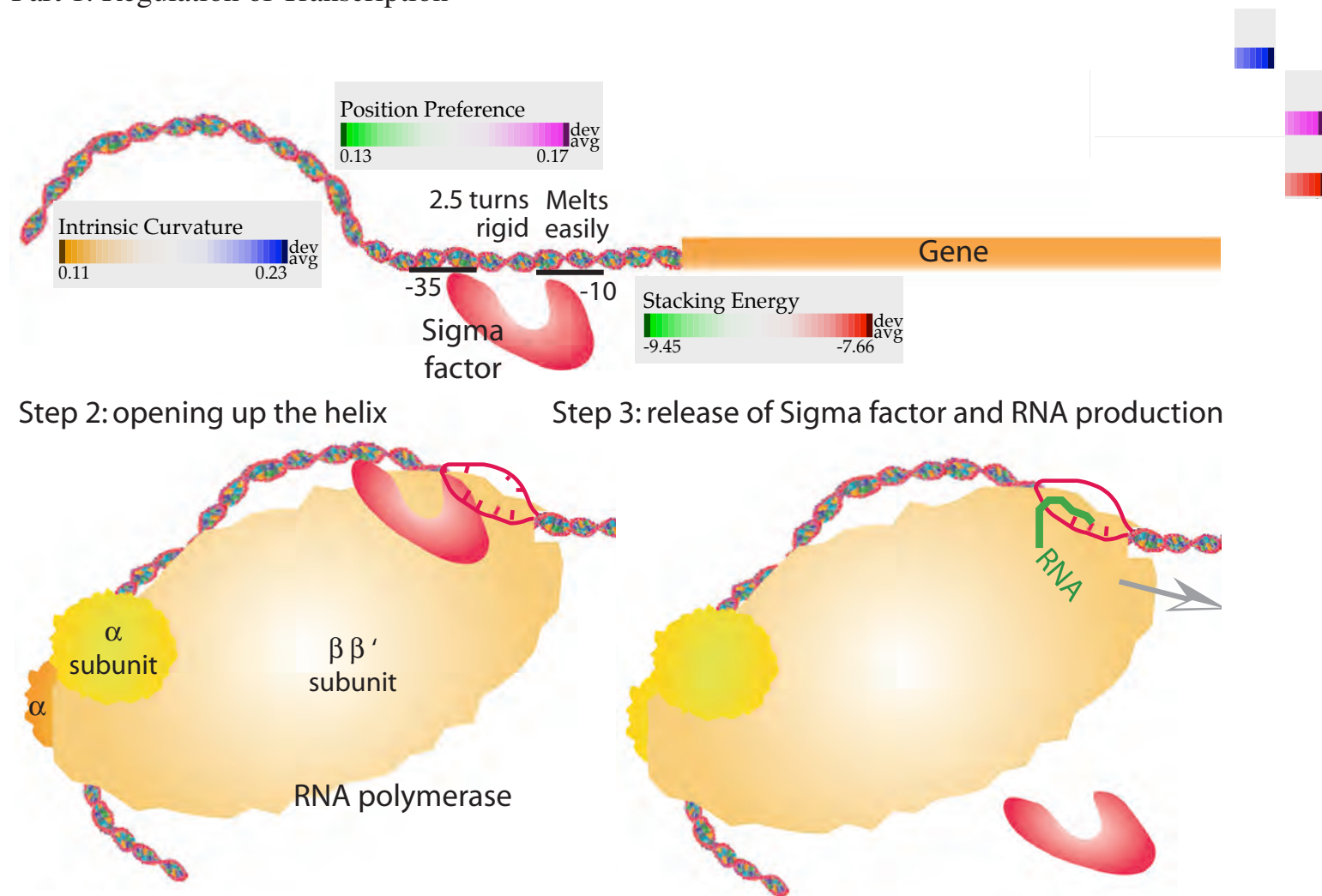
<http://www.cbs.dtu.dk/>  
Center for Biological Sequence Analysis

GENOME ATLAS

# Promoter Structural profile

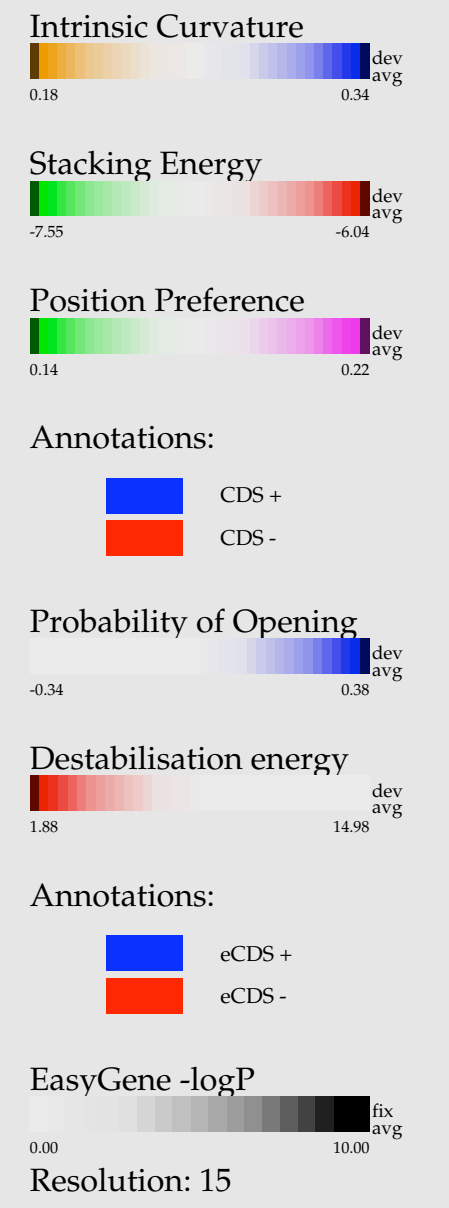
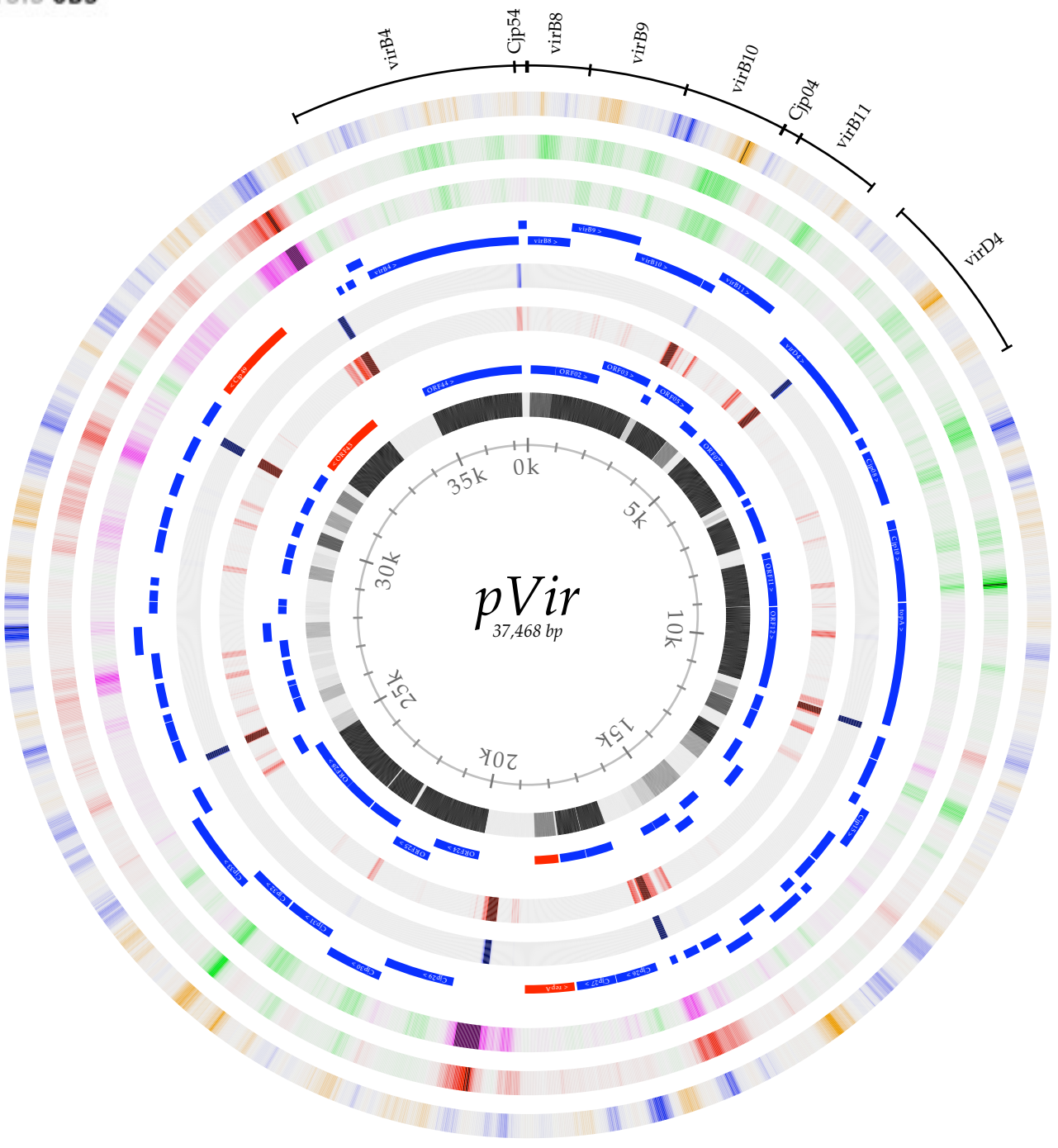


## Part 1: Regulation of Transcription



**Fig. 10.5** Initiation of transcription in bacteria. In the first step, Sigma factor binds to the DNA on two locations (in the case of Sigma 70 the -35 and -10 sites). RNA polymerase (a complex of two  $\alpha$ , one  $\beta$  and one  $\beta'$  subunit) binds next, after which the DNA wraps around the protein. Sigma induces local strand separation so that RNA polymerase starts producing RNA (*in green*). The Sigma factor is then released and RNA polymerase proceeds along the DNA, moving with a local bubble of melted DNA (*indicated by the arrows*)

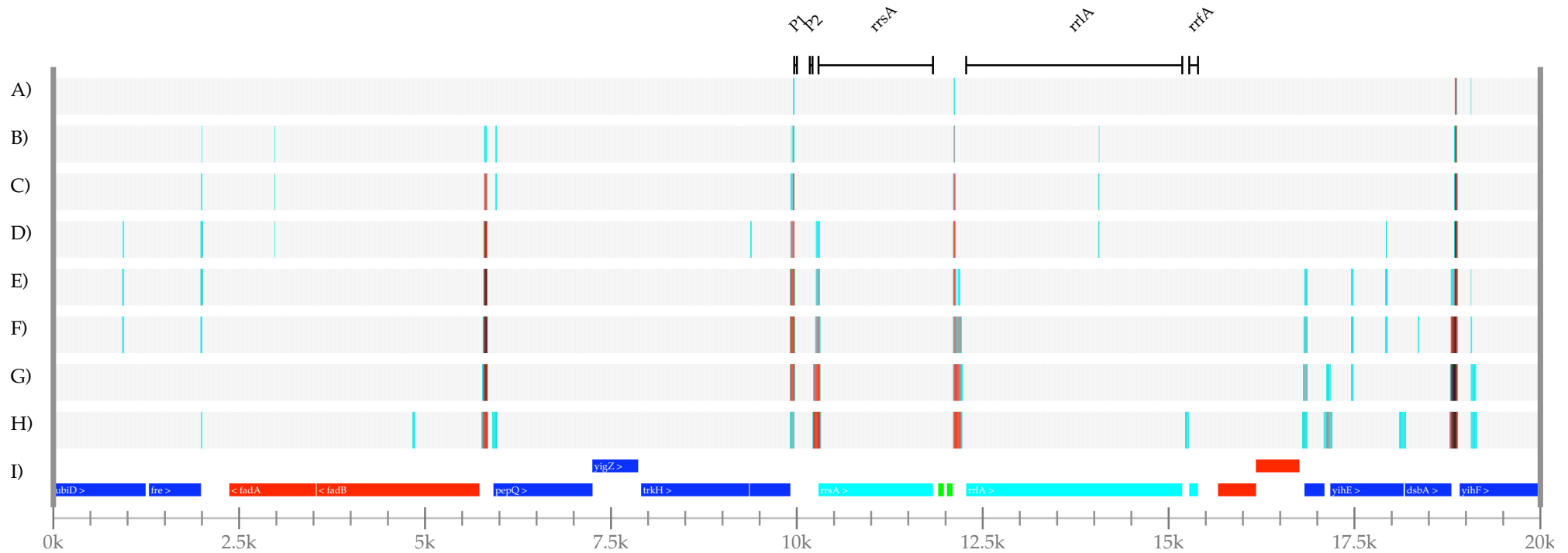




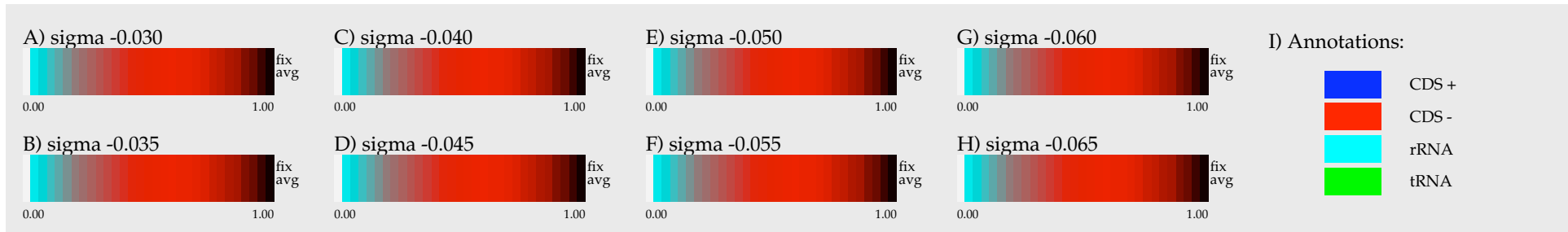
<http://www.cbs.dtu.dk/>  
 Center for Biological Sequence Analysis

# *E. coli*\_K12\_MG1655 *rrsA*

20,000 bp

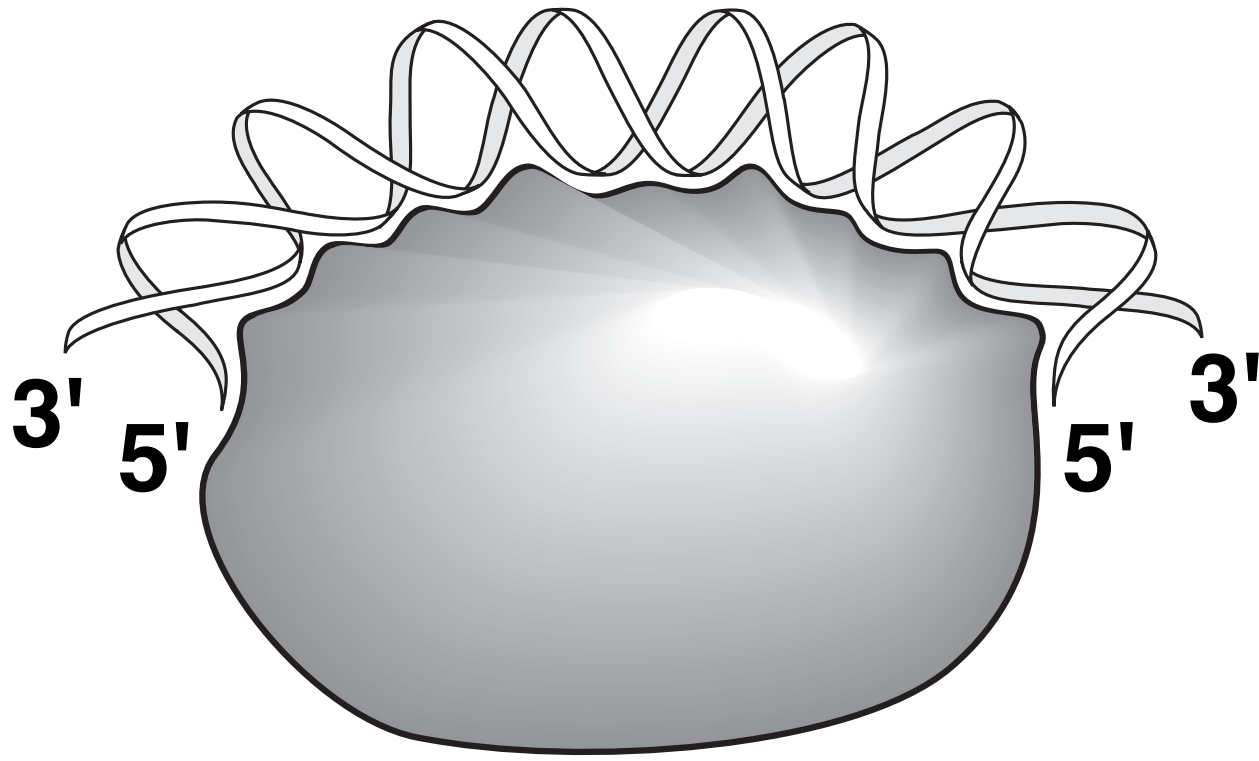


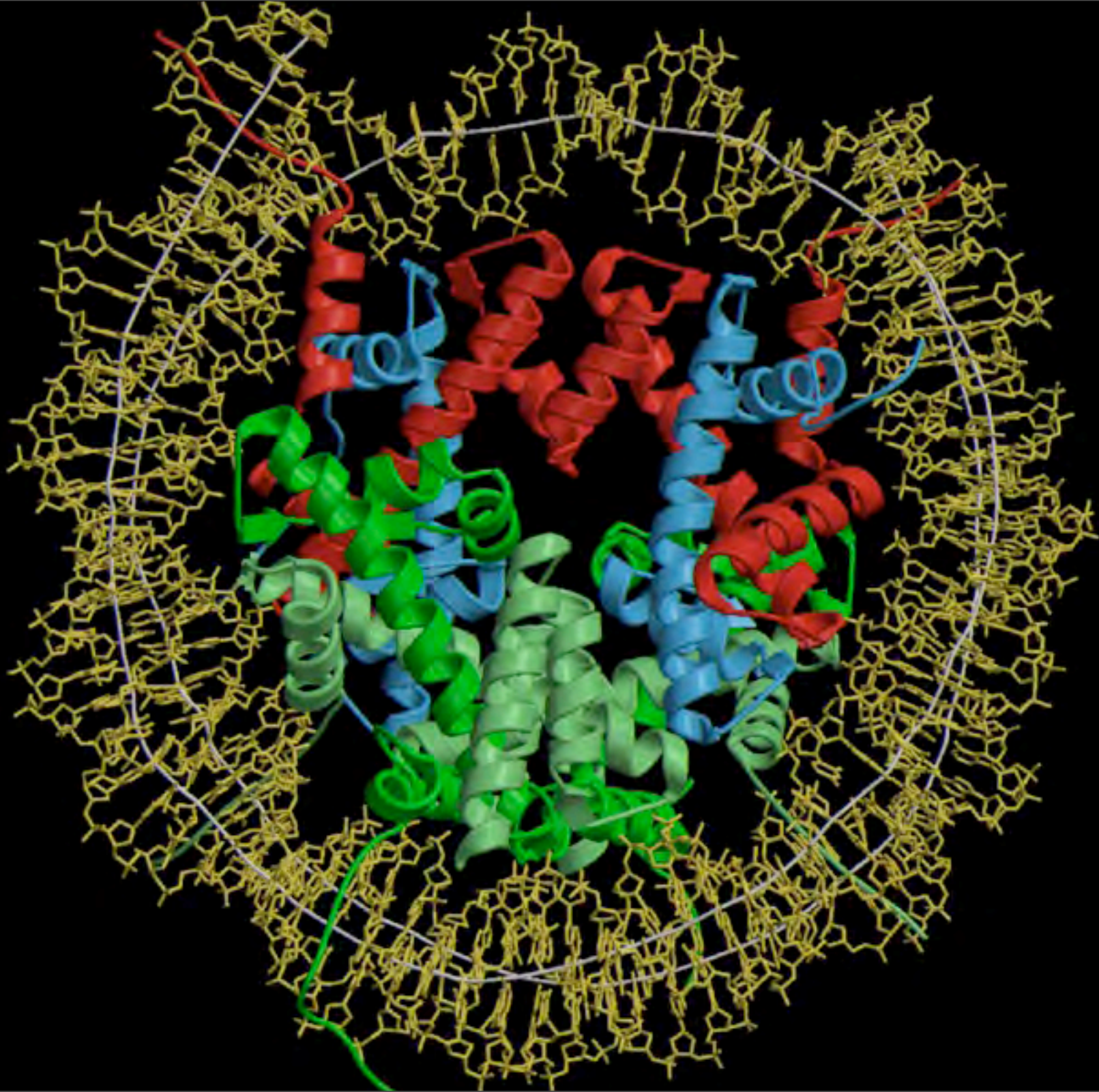
Resolution: 9

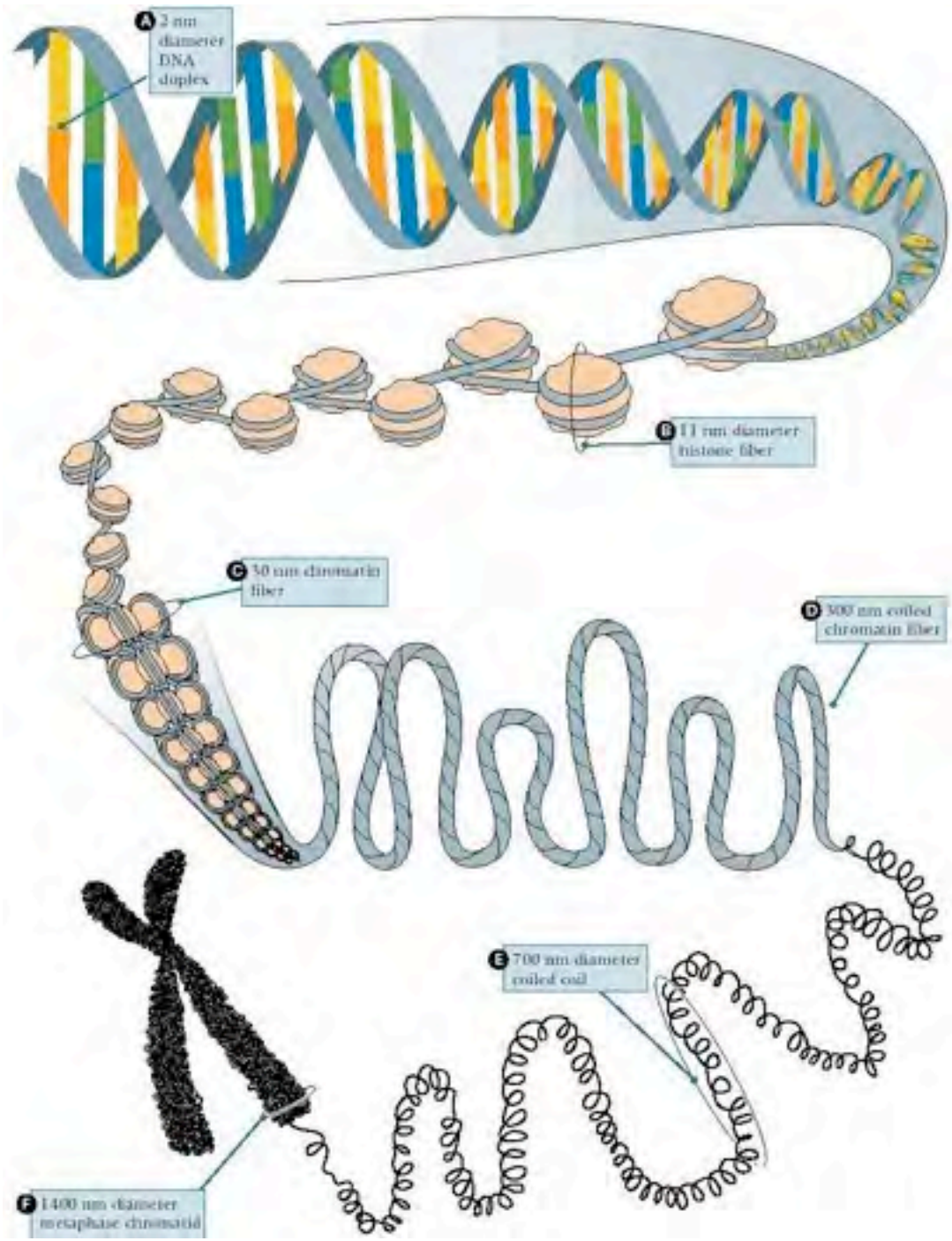


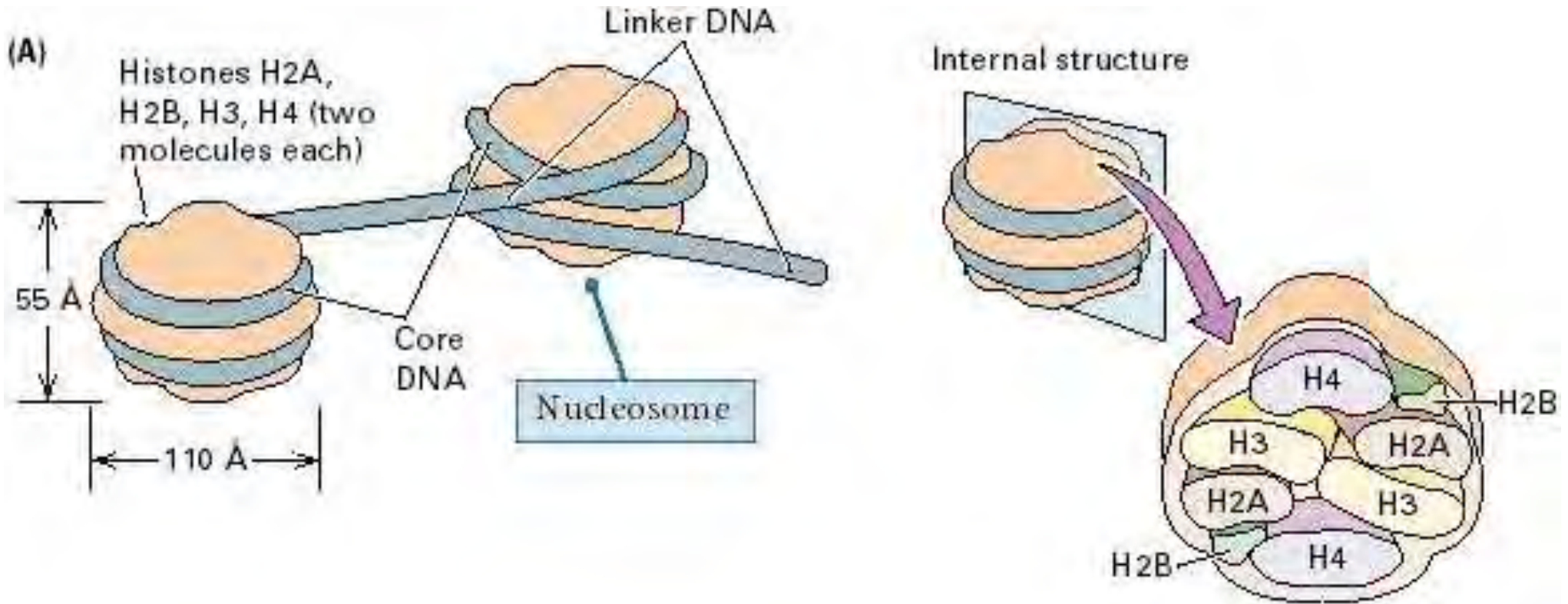


# Pause



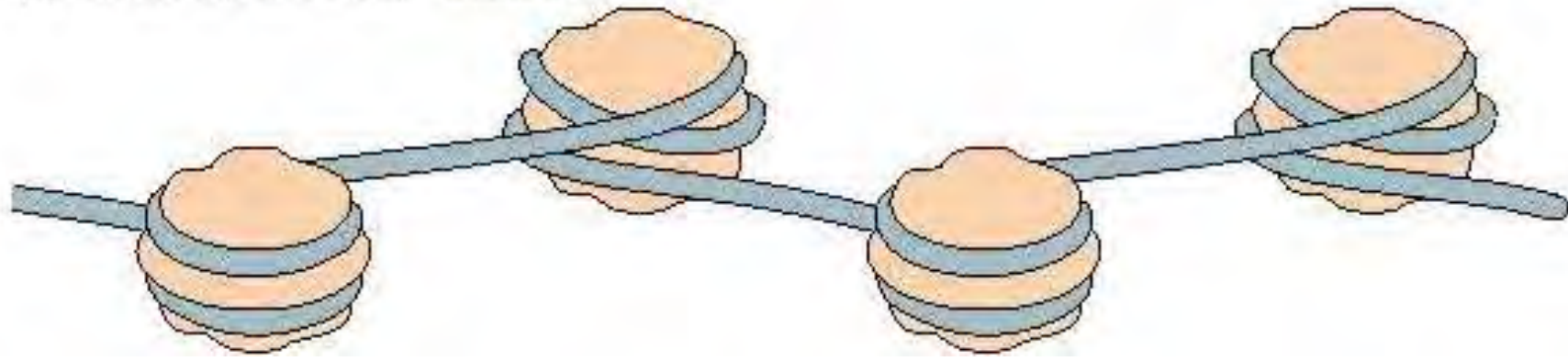




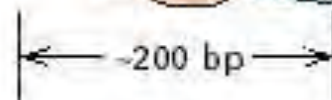
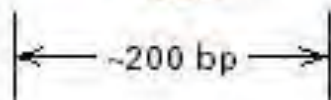
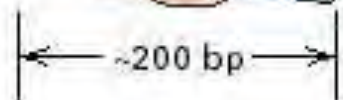
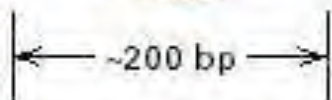
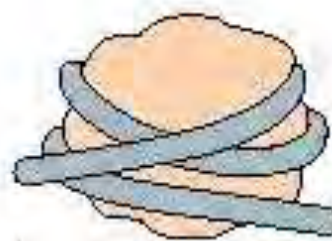




**(B) Chromatin in form of "beads on a string"**



**Light nuclease digestion**

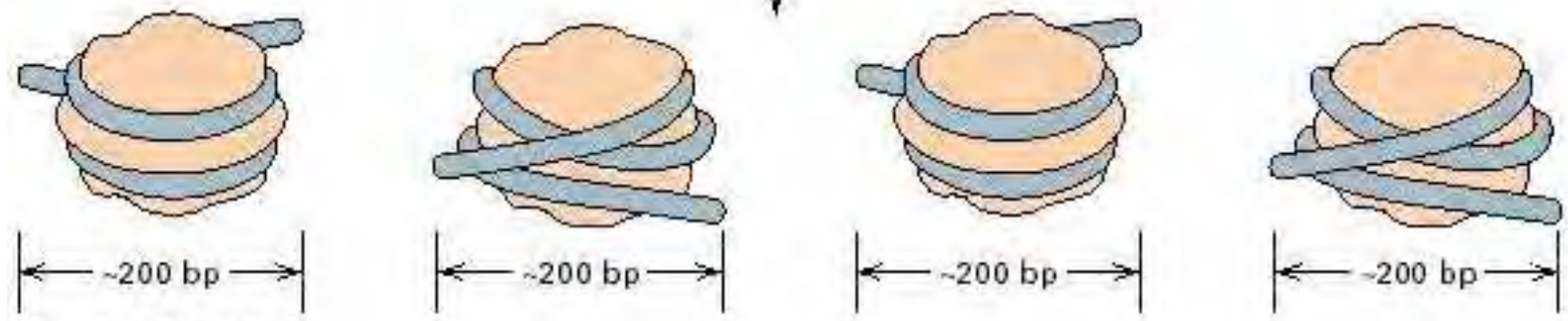


Average length of  
DNA fragment ~200 bp  
per nucleosome



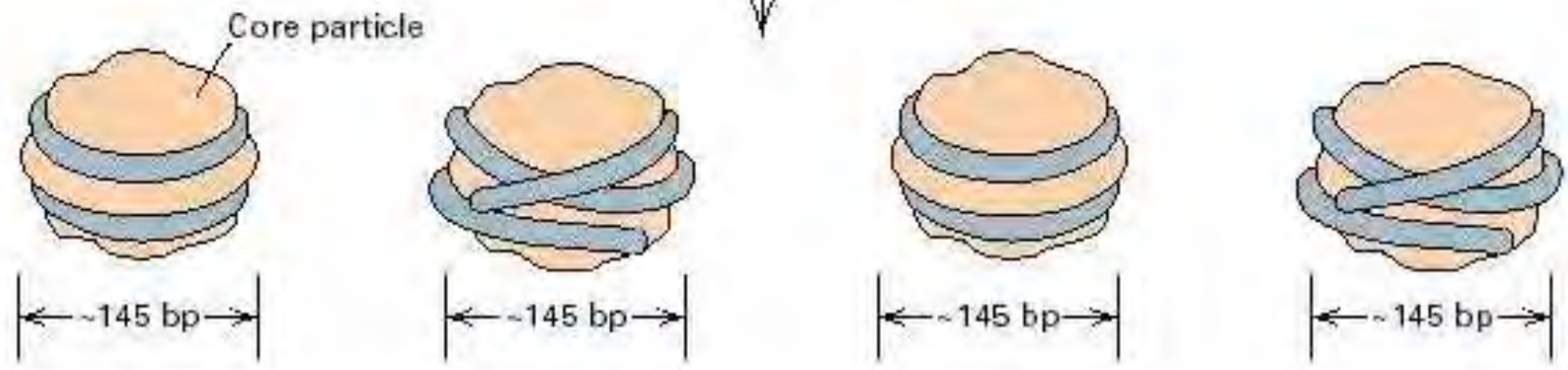


Light nuclease digestion



Average length of DNA fragment ~200 bp per nucleosome

Heavier nuclease digestion



Average length of core DNA ~145 bp per core particle

# "Travers" trinucleotide scale:



Satchwell, S.C., Drew, H.R., and Travers, A.A., "Sequence periodicities in chicken nucleosome core DNA", *J. Mol. Biol.*, **191**:659-675, (1986).

NOTE: we use a (slight) modification, in which the absolute value (magnitude) of the values is used to reflect trinucleotides which tend to exclude nucleosomes.

Baldi, P., Brunak, S., Chauvin, Y., and Krogh, A., "Naturally occurring nucleosome positioning signals in human exons and introns",

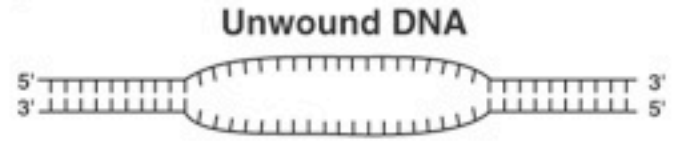
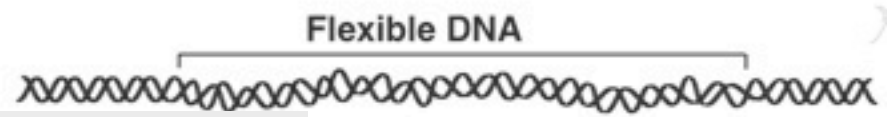
*J. Mol. Biol.*, **263**:503-510, (1996).

trinuc.	% Out
AAT/ATT	-0.280
AAA/TTT	-0.274
CCA/TGG	-0.246
AAC/GTT	-0.205
ACT/AGT	-0.183
CCG/CGG	-0.136
ATC/GAT	-0.110
AAG/CTT	-0.081
CGC/GCG	-0.077
AGG/CCT	-0.057
GAA/TTC	-0.037
ACG/CGT	-0.033
ACC/GGT	-0.032
GAC/GTC	-0.013
CCC/GGG	-0.012
ACA/TGT	-0.006
CGA/TCG	-0.003
GGA/TCC	0.013
CAA/TTG	0.015
AGC/GCT	0.017
GTA/TAC	0.025
AGA/TCT	0.027
CTC/GAG	0.031
CAC/GTG	0.040
TAA/TTA	0.068
GCA/TGC	0.076
CTA/TAG	0.090
GCC/GGC	0.107
ATG/CAT	0.134
CAG/CTG	0.175
ATA/TAT	0.182
TCA/TGA	0.194

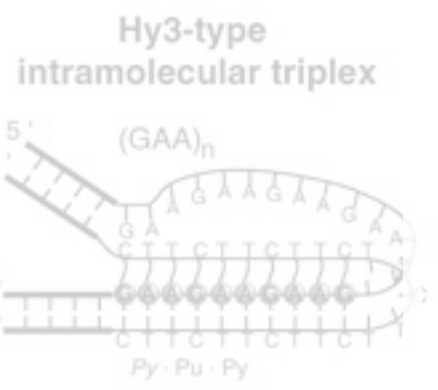
High position pref.

Low position pref.

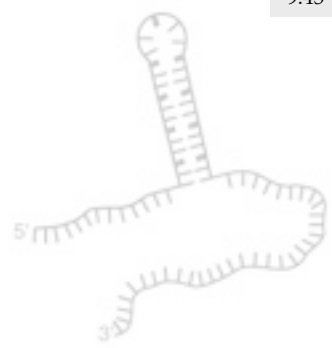
High position pref.



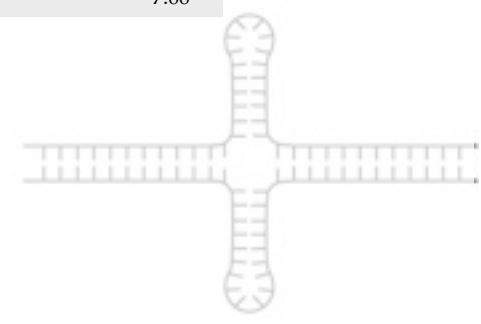
CGG quadruplex



Hairpin structure



Y-junction structure



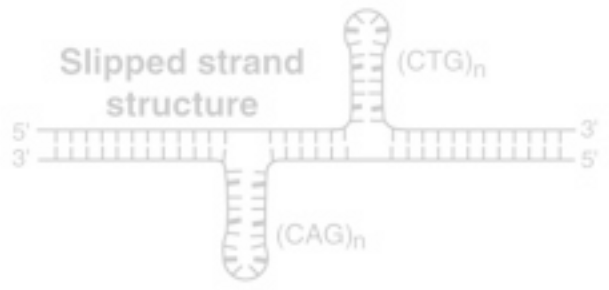
Anti Parallel DNA



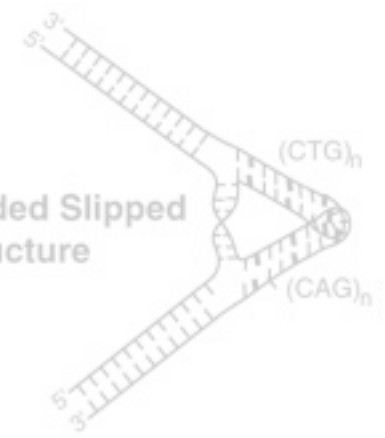
Parallel DNA



Slipped strand structure

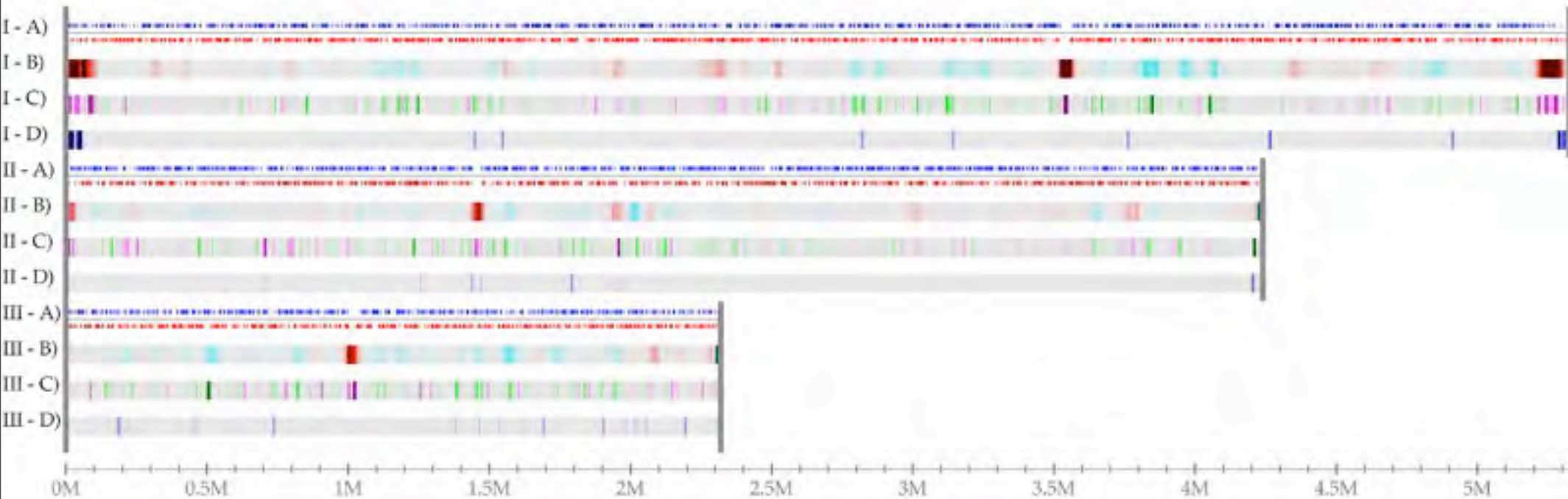


Folded Slipped Structure



# Schizosaccharomyces pombe

All Three Chromosomes 11,896,623 bp total



Resolution: 1066

A) Annotations:



I - B) Percent AT



D) Watson Repeats



C) Position Preference



Center for Biological Sequence Analysis  
<http://www.cbs.dtu.dk/>

GENOME ATLAS

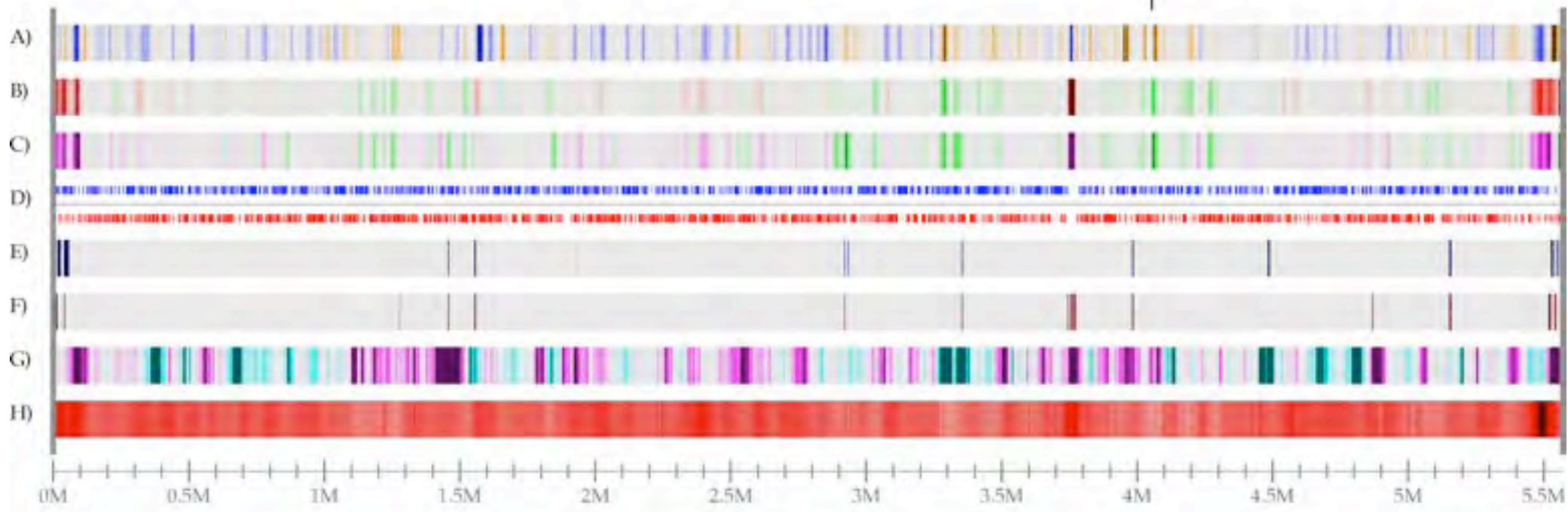
36

Nature, 415:871-880, (2002).

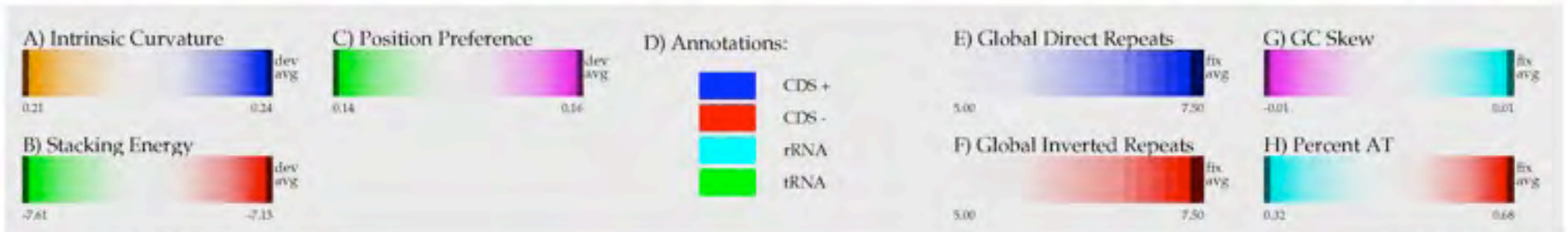
# Schizosaccharomyces pombe

strain 972, chromosome I 5,570,797 bp

SPAC23A1.07



Resolution: 2229



Center for Biological Sequence Analysis  
<http://www.cbs.dtu.dk/>

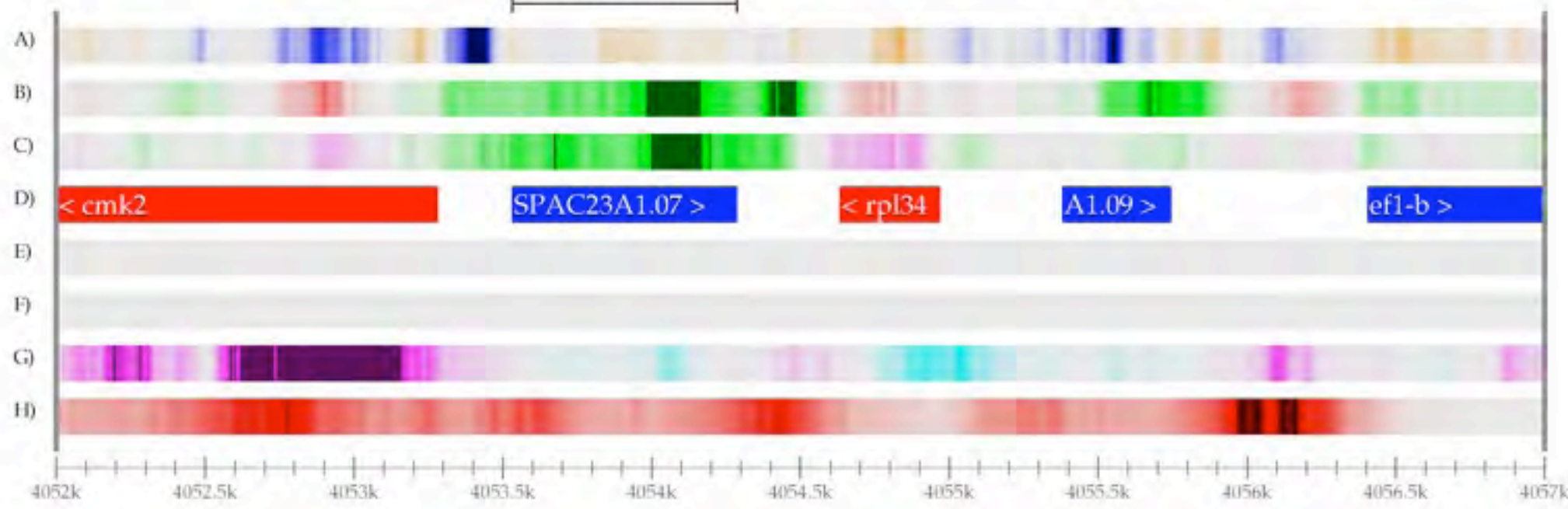
GENOME ATLAS

Range:  
4052000  
..  
4057000

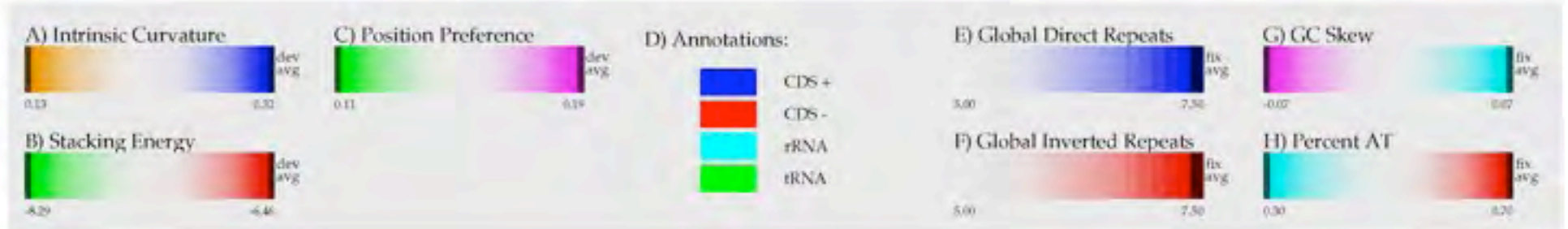
# Schizosaccharomyces pombe

strain 972, chromosome I 5,570,797 bp

SPAC23A1.07



Resolution: 2



Center for Biological Sequence Analysis  
<http://www.cbs.dtu.dk/>

GENOME ATLAS

### Global Direct Repeats

#### Direct repeats

5' - AGAGAGAGATCTAGACTAGAGAGAGAGATACCCTAGAGAGAGAGAGAAAG - 3'  
 3' - TCTCTCTCTAGATCTGATCTCCCTCTCTCTATGGGATCTCTCTCTCTCTTTC - 5'

#### Local direct simple repeats

5' - AGGCAAAGTCGCCATGCTGCCATCTGGAGGCAAAGTCGCCATGCTGCCATCTGG - 3'  
 3' - TCCGTTTCAGCGGTACGACGGTAGACGTCCGTTTCAGCGGTACGACGGTAGACC - 5'

#### Local direct repeats

5' - AGGCAAAGTCGCCATGCTGCCATCTGGTGCCAGGCAAAGTCGCCATGCTGCCATCTGG - 3'  
 3' - TCCGTTTCAGCGGTACGACGGTAGACGACGGGTCCGTTTCAGCGGTACGACGGTAGACC - 5'

#### Local direct repeats with spacer

### Global Inverted Repeats

#### Inverted repeats

5' - AGGCAAAGTCGCCATGCTGCCATCTGGCCAGATGGCAGCATGGCGACTTTGCCT - 3'  
 3' - TCCGTTTCAGCGGTACGACGGTAGACCGGTCTACCGTCTACCGCTGAAACGGA - 5'

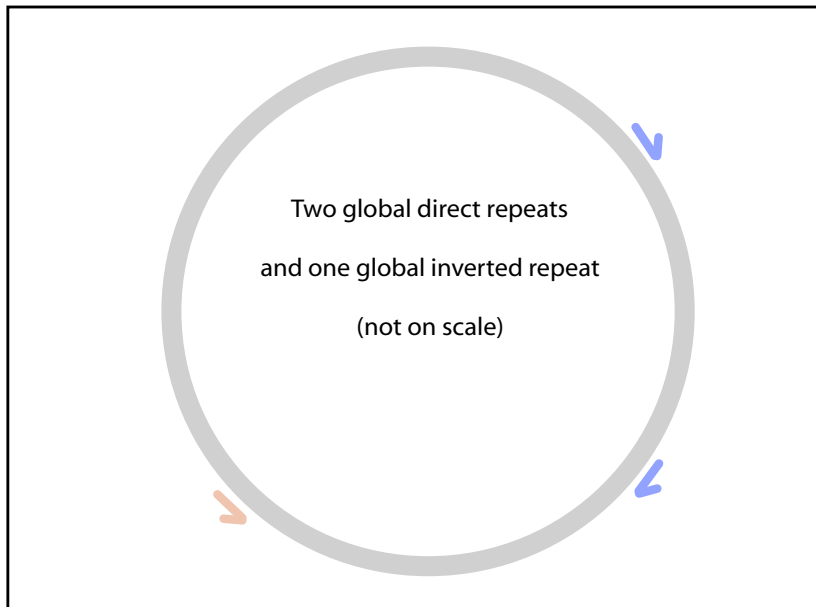
#### Local inverted repeat

5' - AGGCAAAGTCGCCATGCTGCCATCTGGTGCCAGATGGCAGCATGGCGACTTTGCCT - 3'  
 3' - TCCGTTTCAGCGGTACGACGGTAGACCGGTCTACCGTCTACCGCTGAAACGGA - 5'

#### Inverted repeat with spacer

5' - AGGCACGTGAGAATGAATCCTGCTTTAGATCTAAAAGGCAAAGCTTTGCCT - 3'  
 3' - TCCGTGCCTCTTACTTAAGGACGAAATCTAGATTTTCCGTTTCGAAACGGT - 5'

#### Palindromic repeats



#### Mirror and Everted repeats

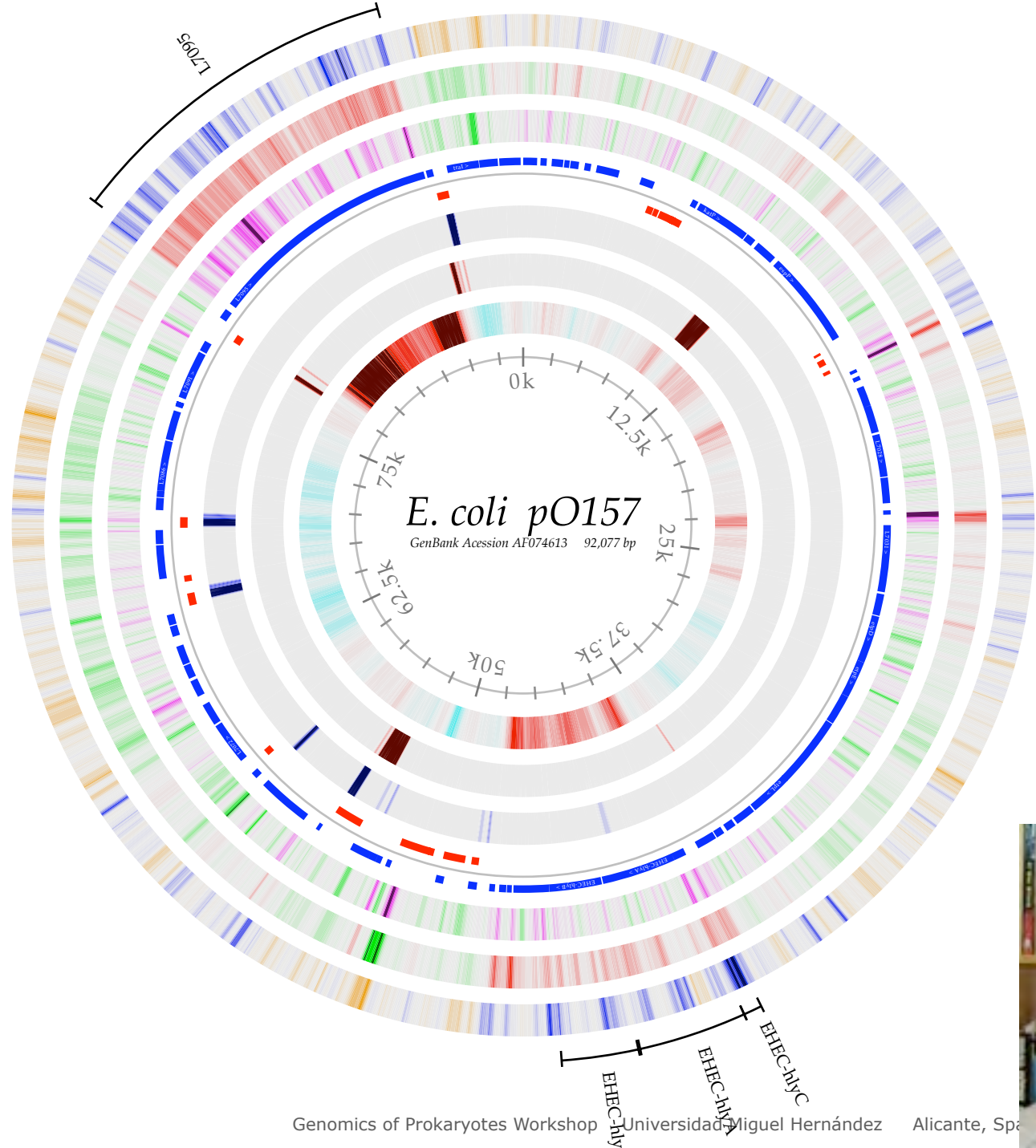
5' - AGGCAAAGTCGCCATGCTGCCATCTGGGGTCTACCGTCTACCGCTGAAACGGA - 3'  
 3' - TCCGTTTCAGCGGTACGACGGTAGACCCAGATGGCAGCATGGCGACTTTGCCT - 5'

#### Mirror repeat

5' - AGGCAAAGTCGCCATGCTGCCATCTGGTCCGTTTCAGCGGTACGACGGTAGACC - 3'  
 3' - TCCGTTTCAGCGGTACGACGGTAGACCAGGCAAAGTCGCCATGCTGCCATCTGG - 5'

#### Everted repeat

# Introducing the "Genome Atlas"



**Intrinsic Curvature**  
0.08 0.30

**Stacking Energy**  
-9.52 -6.41

**Position Preference**  
0.11 0.17

**Annotations:**  
■ CDS +  
■ CDS -

**Watson Repeats**  
5.00 7.50

**Crick Repeats**  
5.00 7.50

**Percent AT**  
0.30 0.70

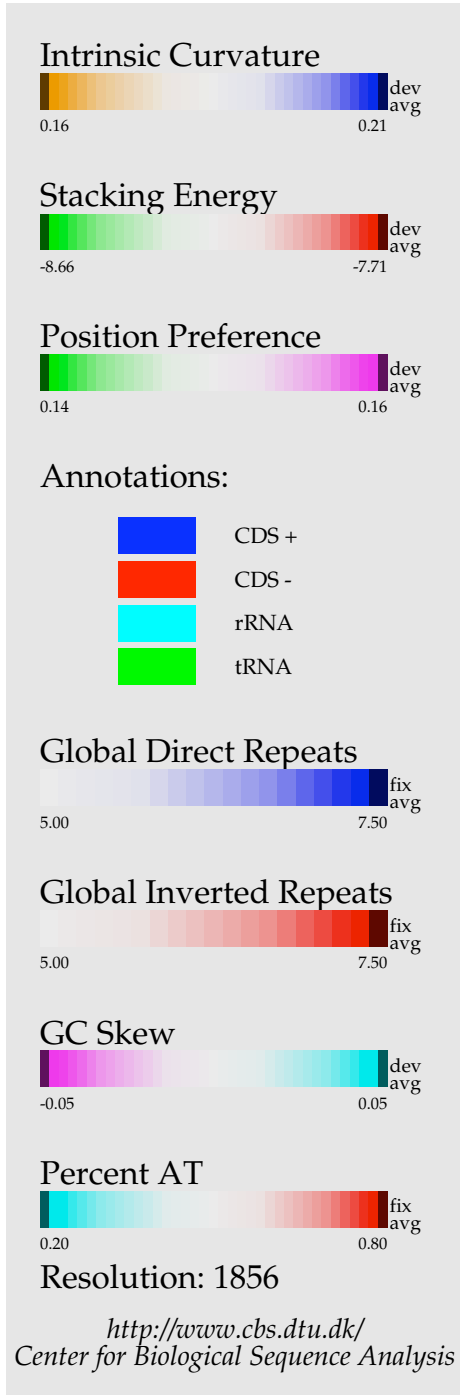
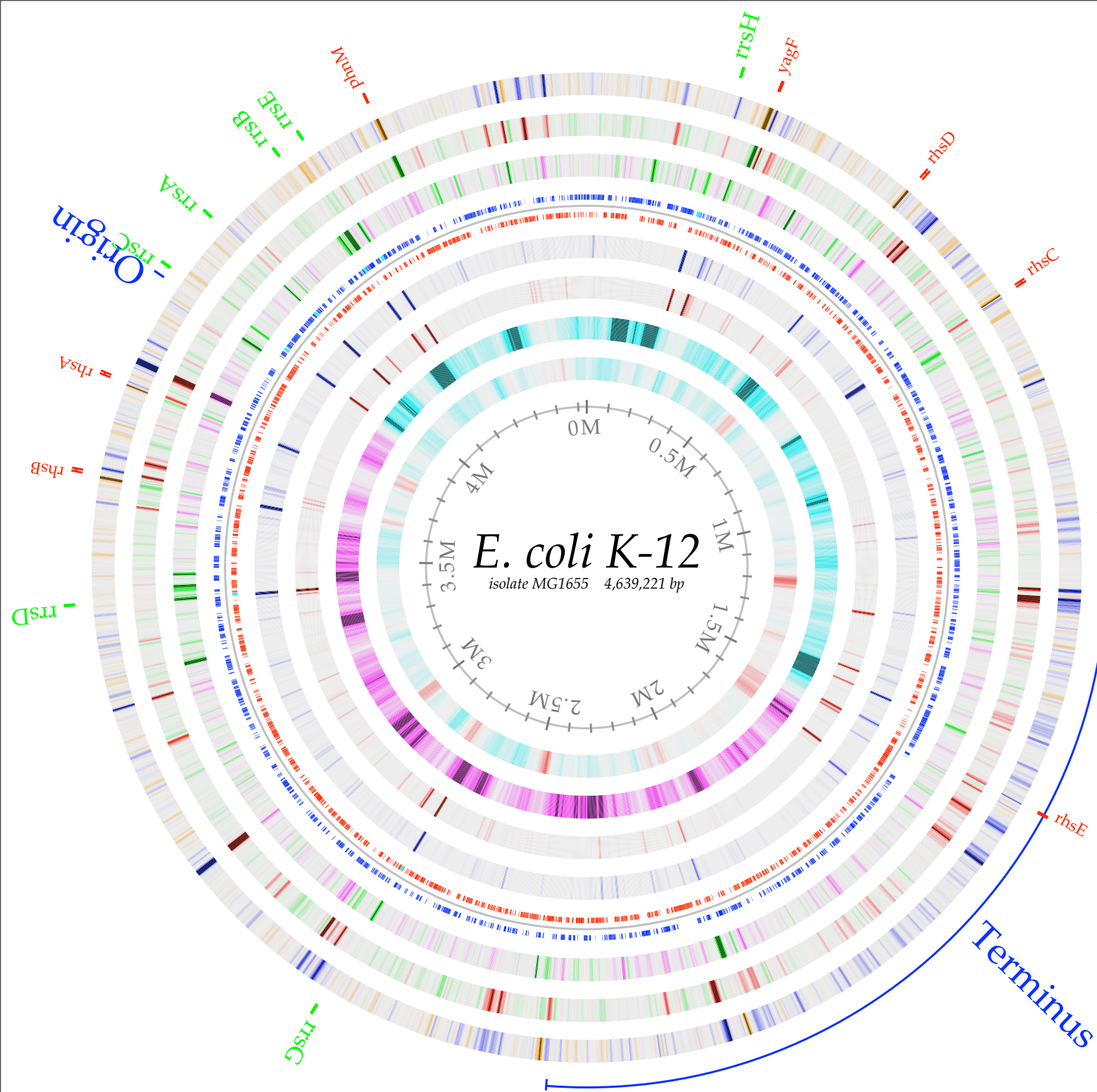
**Resolution: 37**

Friis, Jensen, and Ussery  
[Genetica, 108:47-51, \(2000\).](#)  
<http://www.cbs.dtu.dk/>  
 Center for Biological Sequence Analysis

**GENOME ATLAS**







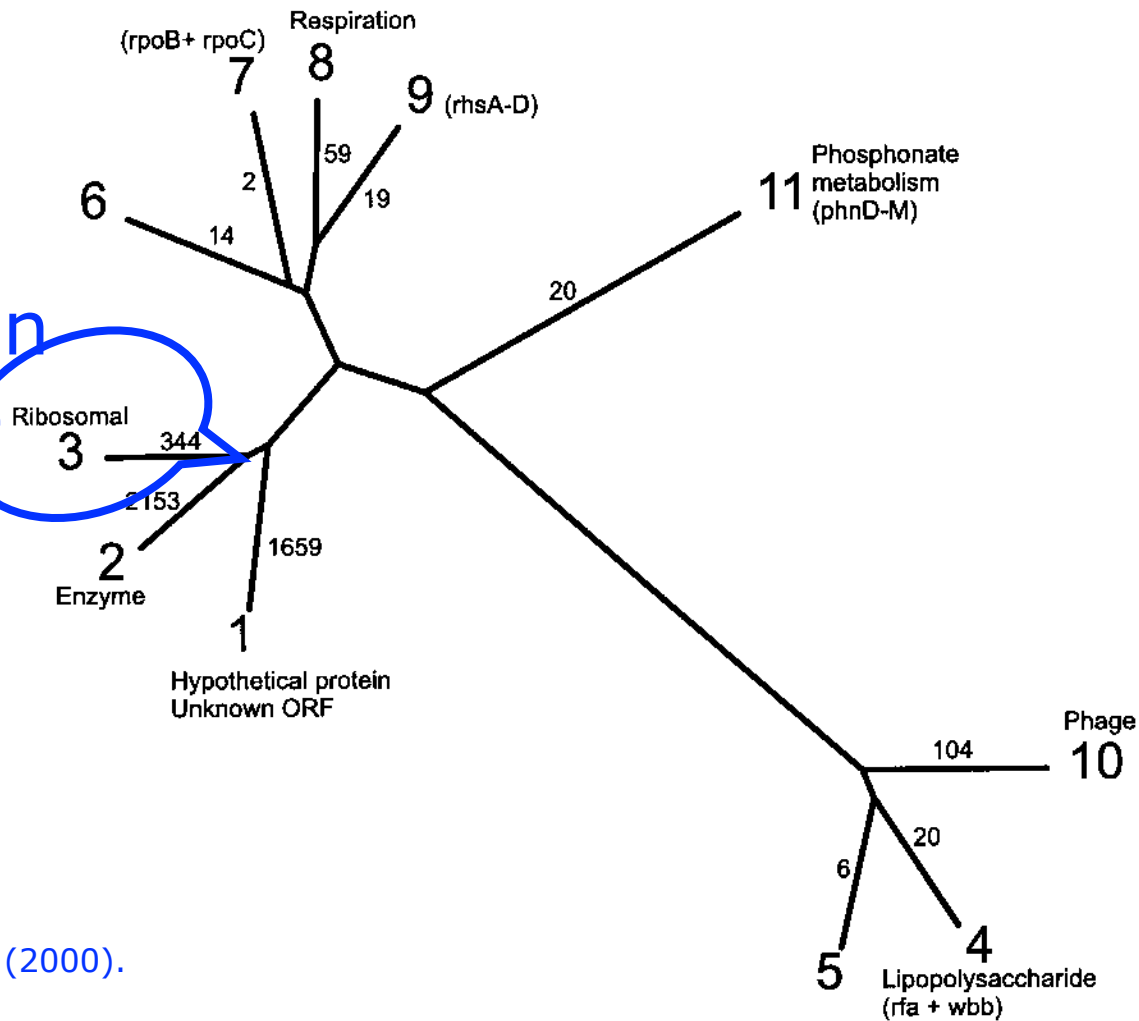
GENOME ATLAS

[Res. Microbiol.](#) **150**:773-777 (1999).  
[J. Mol. Biol.](#) **299**:907-930, (2000).



# DNA Structural Atlas for Escherichia coli

low position preference

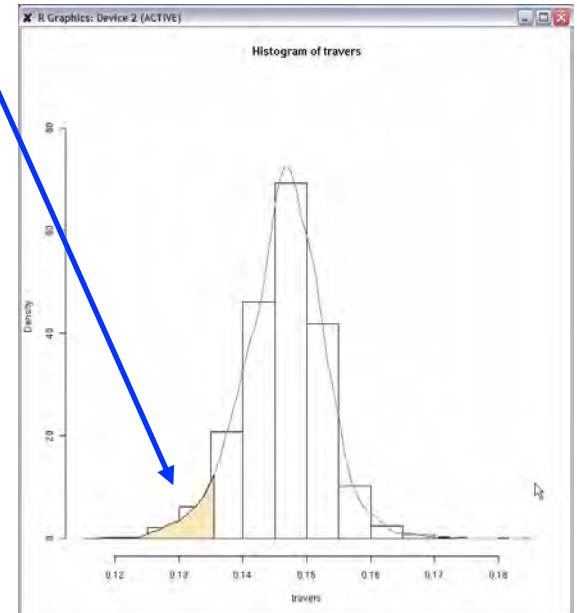
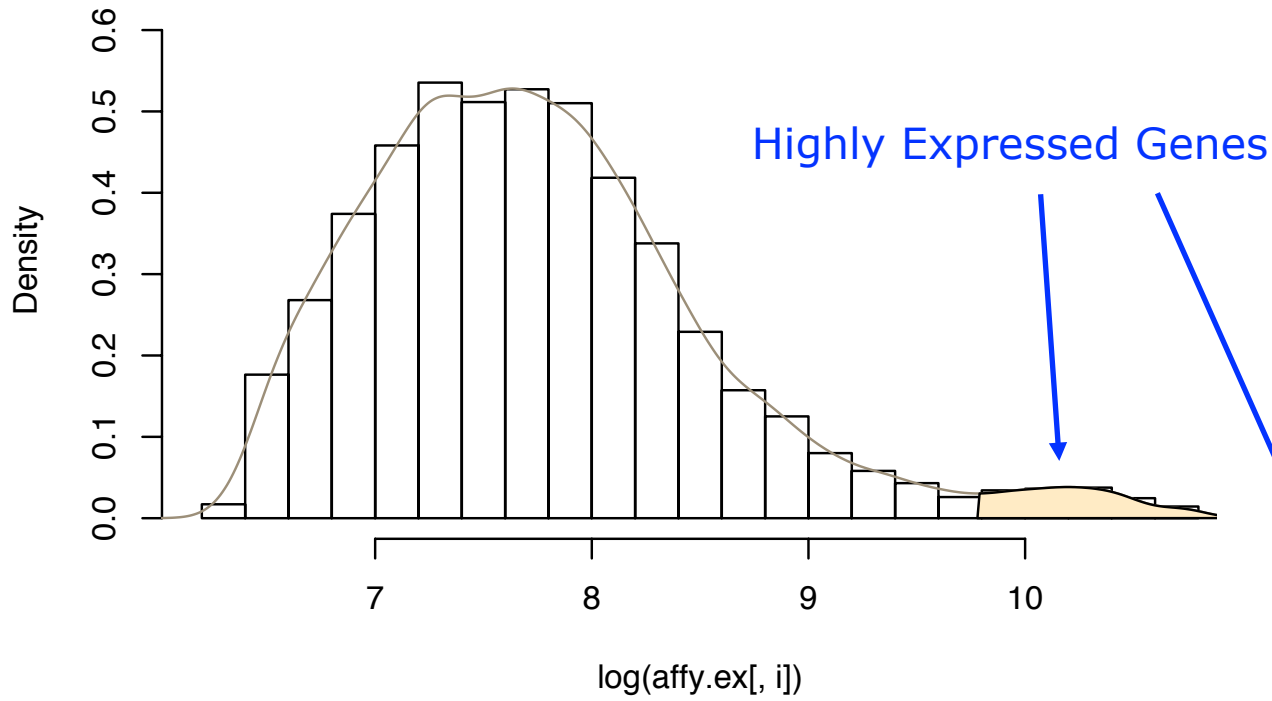


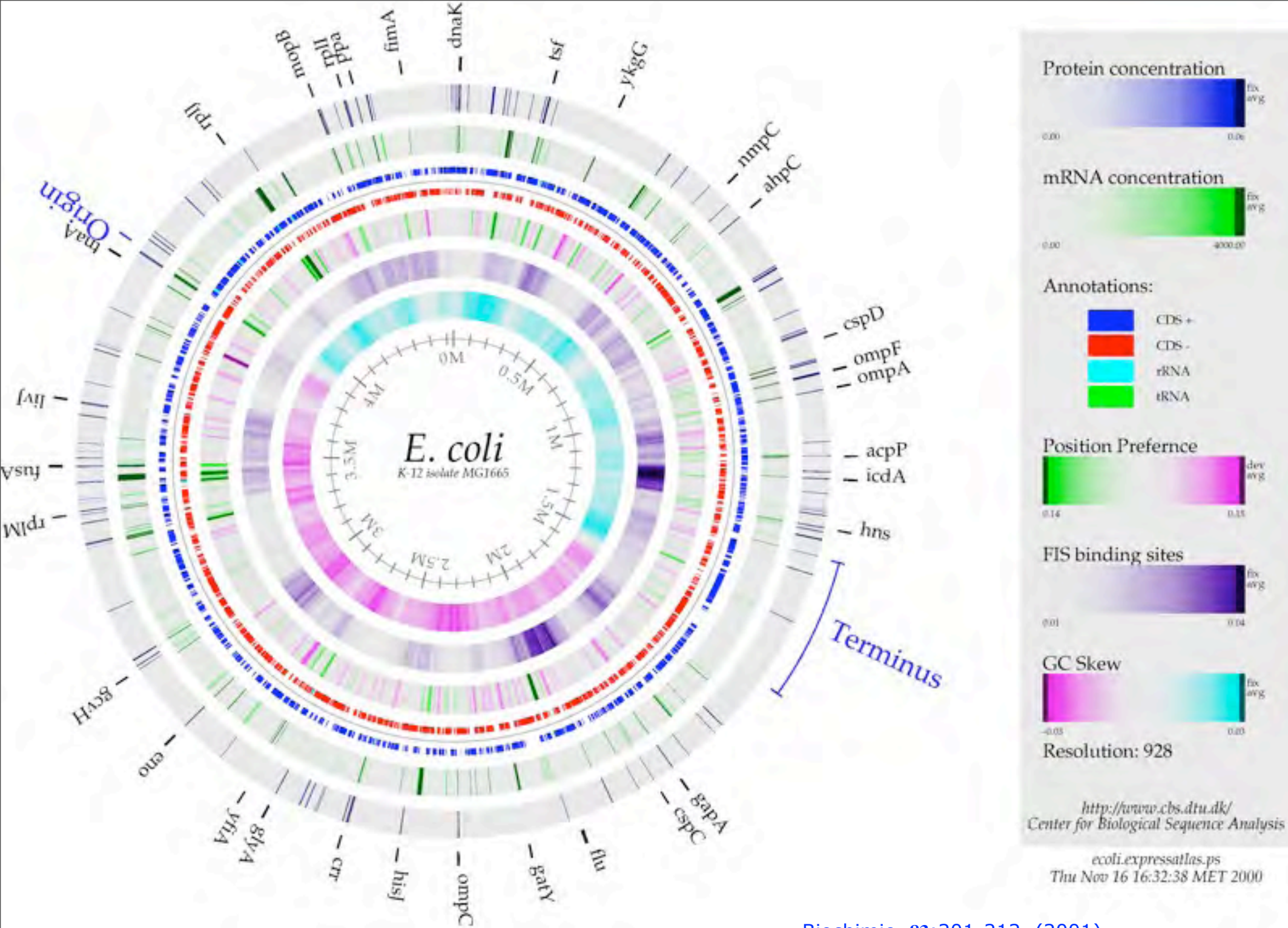
[J. Mol. Biol.](#), 299:907-930, (2000).

**Figure 7.** Structural cluster analysis. Distance tree showing the relative location of 11 gene clusters based on average structural measures. The number of genes in



### Histogram of $\log(\text{affy.ex}[, i])$





Biochimie, 83:201-212, (2001).



Monday, 12 December, 2011