# Credits

- This presentation leans heavily on other people's work and graphics
- All credits are available in the **speaker notes** which you should consult to find out who made all these great movies and images
- **Thank you so much Wikipedia Commons in particular!**

https://berthub.eu/revdna/

# Online questions!
https://webchat.oftc.net/?channels=why2025-andromeda

# IRC: oftc.net, channel: #why2025-andromeda

The **ISME** Journal
*Multidisciplinary Journal of Microbial Ecology*

⊙ Login   🛒 Cart

Search [    ] go  Advanced search

Journal home > Archive > Original Articles > Full text

**Journal home**
**Advance online publication**
  L. About AOP
**Current issue**
**Archive**
**Focuses**
**Browse by subject**
**Press releases**

≡⟋Online submission
For authors
For referees
Contact editorial office
About the journal
  L. Editors and Editorial Board
About the society
For librarians
Subscribe

## Original Article

# Density-dependent adaptive resistance allows swimming bacteria to colonize an antibiotic gradient

Felix J H Hol[1], Bert Hubert[1], Cees Dekker[1] and Juan E Keymer[1,2,3]

[1]Department of Bionanoscience, Kavli Institute of Nanoscience, Delft University of Technology, Delft, The Netherlands
[2]Department of Ecology, Faculty of Biological Sciences, P. Catholic University of Chile, Santiago, Chile
[3]Institute of Physics, Faculty of Physics, P. Catholic University of Chile, Santiago, Chile

Correspondence: FJH Hol or JE Keymer, Department of Bionanoscience, Kavli Institute of Nanoscience, Delft University of Technology, Lorentzweg 1, Delft 2628CJ, The Netherlands. E-mail: f.j.h.hol@tudelft.nl or jkeymer@uc.cl

Received 5 January 2015; Revised 9 April 2015; Accepted 19 May 2015
Advance online publication 3 July 2015

## Abstract                                                    ▲ Top

**During antibiotic treatment, antibiotic concentration gradients develop. Little is know regarding the effects of antibiotic gradients on populations of nonresistant bacteria. Using a microfluidic device, we show that high-**

**T̃U**Delft Delft University of Technology

https://www.nature.com/articles/ismej2015107

# scientific **data**

# SkewDB, a comprehensive database of GC and 10 other skews for over 30,000 chromosomes and plasmids

Bert Hubert ✉

Download PDF

**Sections**    Figures

Abstract

Background & Summary

Methods

Data Records

*SkewDB*

https://www.nature.com/articles/s41597-022-01179-8

"Imagine a flashy spaceship lands in your backyard. The door opens and you are invited to investigate everything to see what you can learn. The technology is clearly millions of years beyond what we can make.

This is biology."



"Kindly take us to your President!"

Although…

© NASA (oddly enough)

- Let's study DNA the way we study random binary blobs

- Highlight many cool DNA things

- I want YOU to Join in!

- DNA: Millions, billions of nucleotides or "bases":
  - A, C, G, T
- Organized in chromosomes & genes
- Absolutely **atom for atom** universal across all life
- >4 billion years old

ARCHEIA

ANIMALS, PLANTS, FUNGI, **US**

BACTERIA

2 billion years ago

4 billion years ago

First water: 4.2 billion

Age of earth: 4.5 billion

https://giphy.com/gifs/sea-vents-hydrothermal-1bTEQnjArFBy8

# Basics: DNA

A

C

G

T

"00"

"01"

"10"

"11"

**DNA is very much like tape**

Sometimes circular tape - no beginning, no end!

No addressing! No alignment!

It is a **nucleotide stream** which can be compared to a **bitstream**

It IS however **content addressable**!

http://onlinelibrary.wiley.com/doi/10.1002/smll.201400265/full

Hepatitis B - 800 bytes (unpleasant)

```
ctccactgccttccaccaagctctgcaggatcccaaagtcaggggtctgtattttcctgctggtggctccagttcaggaacagtaaaccctgctccgaatattgcctctcacatc
tcgtcaatctccgcgaggactggggaccctgtgacgaacatggagaacatcacatcaggattcctaggacccctgctcgtgttacaggcggggtttttcttgttgacaagaatcc
tcacaataccgcagagtctagactcgtggtggacttctctcaattttctagggggatcacccgtgtgtcttggccaaaattcgcagtccccaacctccaatcactcaccaacctc
ctgtcctccaatttgtcctggttatcgctggatgtgtctgcggcgtttttatcatattcctcttcatcctgctgctatgcctcatcttcttattggttcttctggattatcaagt
atgttgcccgtttgtcctctaattccaggatcaacaacaaccagtacgggaccatgcaaaacctgcacgactcctgctcaaggcaactctatgtttccctcatgttgctgtacaa
aacctacggatggaaattgcacctgtattcccatcccatcgtcctgggctttcgcaaaatacctatgggagtgggcctcagtccgtttctcttggctcagtttactagtgccatt
tgttcagtggttcgtagggctttcccccactgtttggctttcagctatatggatgatgtggtattggggggccaagtctgtacagcatcgtgagtccctttataccgctgttacca
atttttcttttgtctctgggtatacatttaaaccctaacaaaacaaaaagatggggttattccctaaacttcatgggttacataattggaagttggggaactttgccacaggatca
tattgtacaaaagatcaaacactgttttagaaaacttcctgttaacaggcctattgattggaaagtatgtcaaagaattgtgggtcttttgggctttgctgctccatttacacaa
tgtggatatcctgccttaatgcctctgtatgcatgtatacaagctaaacaggctttcactttctcgccaacttacaaggcctttctaagtaaacagtacatgaacctttaccccg
ttgctcggcaacggcctggtctgtgccaagtgtttgctgacgcaacccccactggctggggcttggccataggccatcagcgcatgcgtggaacctttgtggctcctctgccgat
ccatactgcggaactcctagccgcttgttttgctcgcagccggtctggagcaaagctcatcggaactgacaattctgtcgtcctctcgcggaaataatacatcgtttccatggct
gctaggctgtactgccaactggatccttcgcgggacgtcctttgtttacgtcccgtcggcgctgaatcccgcggacgacccctctcggggccgcttgggactctctcgtcccctt
ctccgtctgccgttccagccgaccacggggcgcacctctctttacgcggtctccccgtctgtgccttctcatctgccggtccgtgtgcacttcgcttcacctctgcacgttgcat
ggagaccaccgtgaacgcccatcagatcctgcccaaggtcttacataagaggactcttggacttccagcaatgtcaacgaccgaccttgaggcctacttcaaagactgtgtgttt
aaggactgggaggagctggggggaggagattaggttaaaggtctttgtattaggaggctgtaggcataaattggtctgcgcaccaacaccatgcaacttttttcacctctgcctaat
catctcttgtacatgtcccactgttcaagcctccaagctgtgccttgggtggctttggggcatggacattgacccttataaagaatttggagctactgtggagttactctcgttt
ttgccttctgacttctttccttccgtcagagatctcctagacaccgcctcagctctgtatcgagaagccttagagtctcctgagcattgctcacctcaccatactgcactcaggc
aagccattctctgctgggggaattgatgactctagctacctgggtgggtaataatttggaagatccagcatccagggatctagtagtcaattatgttaatactaacatgggttt
aaagatcaggcaactattgtggtttcatatatcttgccttacttttggaagagagactgtacttgaatatttggtctctttcggagtgtggattcgcactcctccagcctataga
ccaccaaatgccccctatcttatcaacacttccggaaactactgttgttagacgacgggaccgaggcaggtcccctagaagaagaactccctcgcctcgcaaacgcagatctcaat
cgccgcgtcgcagaagatctcaatctcgggaatctcaatgttagtattccttggactcataaggtgggaaactttacggggctttattcctctacagtacctatctttaatcctg
aatggcaaactccttccttttcctaagattcatttacaagaggacattattaataggtgtcaacaatttgtgggccctctcactgtaaatgaaaagagaagattgaaattaattat
gcctgctagattctatcctacccacactaaatattttcccttagacaaaggaattaaaccttattatccagatcaggtagttaatcattacttccaaaccagacattatttacat
actctttggaaggctggtattctatataagagggaaaccacacgtagcgcatcattttgcgggtcaccatattcttgggaacaagagctacagcatgggaggttggtcatcaaaa
cctcgcaaaggcatggggacgaatctttctgttcccaaccctctgggattctttcccgatcatcagttggaccctgcattcggagccaactcaaacaatccagattgggacttca
accccatcaaggaccactggccagcagccaaccaggtaggagtgggagcattcgggccagggctcacccctccacacggcggtattttggggtggagccctcaggctcagggcat
attgaccacagtgtcaacaattcctcctcctgcctccaccaatcggcagtcaggaaggcagcctactcccatctctccacctctaagagacagtcatcctcaggccatgcagtgg
aa
```
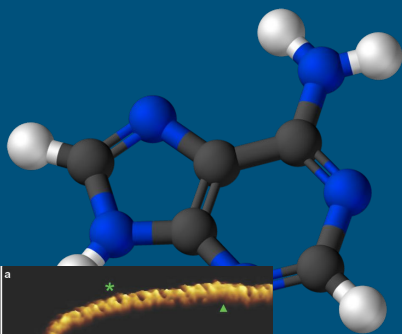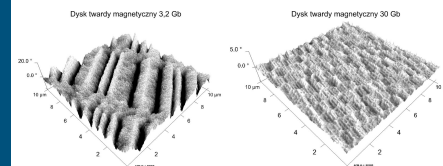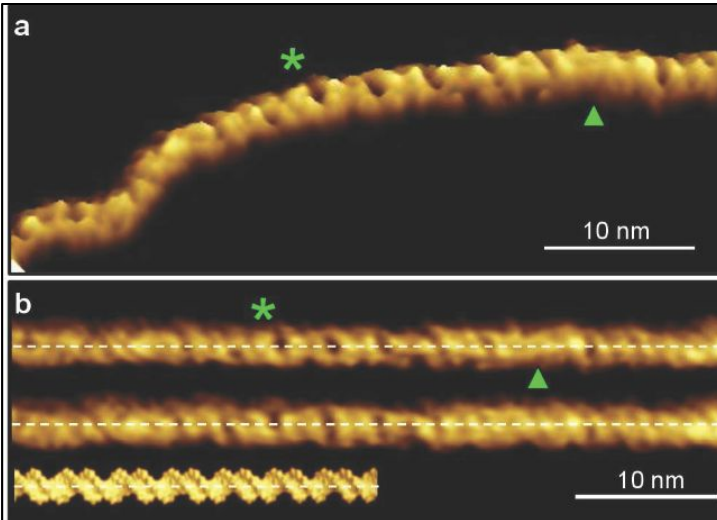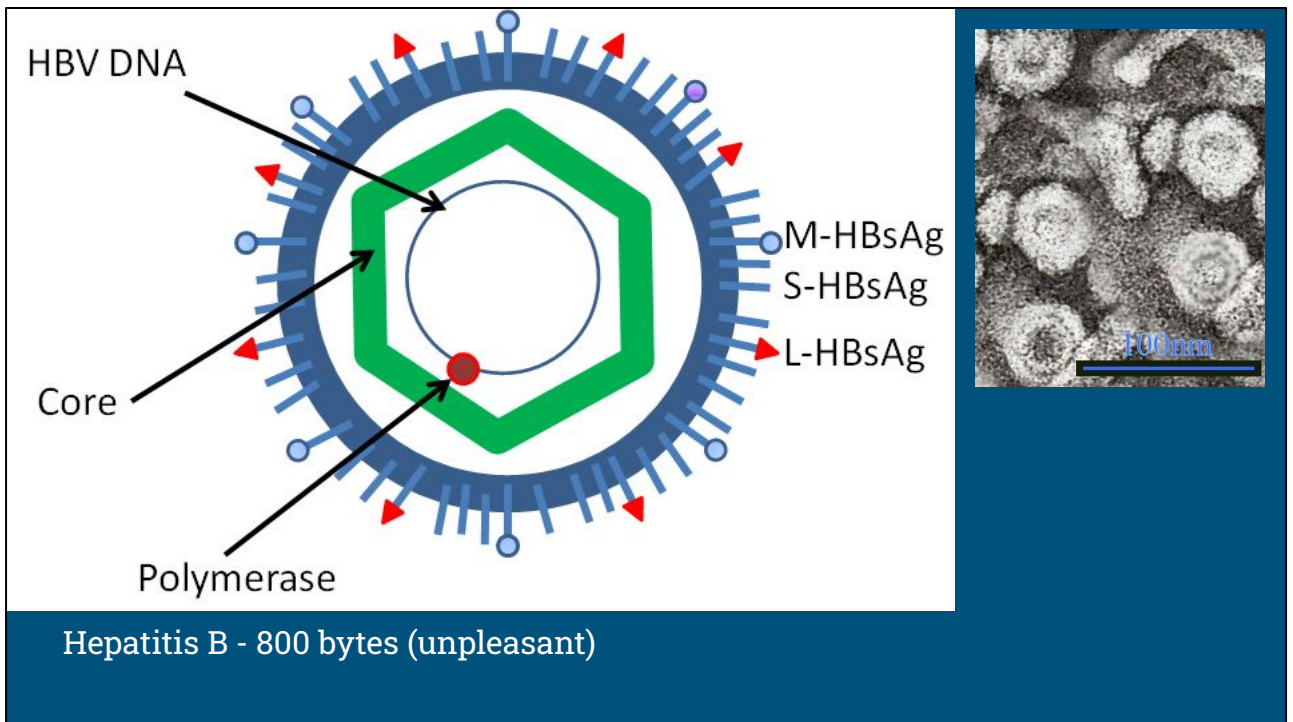
All of Hepatitis-B - 800 bytes

But is life REALLY digital?

Histogram of 49177 BACTERIAL chromosome sizes
Plus select eukaryotes

"Small company Life"

"Enterprise Life" - various departments and procedures

Yeast

Histogram of 49177 BACTERIAL chromosome sizes
Plus select eukaryotes

Histogram of 49177 BACTERIAL chromosome sizes
Plus select eukaryotes

Legend:
- Yeast
- Genlisea aurea (flesh eating plant)
- Pufferfish

Y-axis: Count
X-axis: Chromosome size (megaBYTES)

Most compact animal. Genome almost completely used for proteins.

Histogram of 49177 BACTERIAL chromosome sizes
Plus select eukaryotes

Legend:
- Yeast
- Genlisea aurea (flesh eating plant)
- Pufferfish
- Human beings
- Marbled lungfish

33 gigabytes!!

Log graph!

Count

Chromosome size (megaBYTES)

Unimaginably over the top bloated megacorp stuff!

NCBI HOME    LITERATURE    HEALTH    GENOMES    GENES    PROTEINS    CHEMICALS    POPULAR RESOURCES ▼

All Databases ⌄    | Search NCBI |    🔍 Search

# Genomes

NCBI's Genome resources include information on large-scale genomics projects, genome sequences and assemblies, and mapped annotations, such as variations, markers and data from epigenomics studies.

## How to

Submit sequence data to NCBI

Download a complete genome

Convert feature coordinates between genomic assemblies

Find an interactive view of a genomic annotation

more...

| **Genome Sequences** | **Functional Genomics** | **Variation Resources** | **Additional Tools** |
|---|---|---|---|
| **Genome** | **GEO DataSets** | **dbSNP** | Genome Data Viewer |
| information about organisms' genomes | functional genomics study data | catalog of short genetic variations | displays data tracks in an interactive genome browser |
| **Assembly** | GEO2R | **dbVar** | Genome Workbench |
| genomic assembly statistics | identifies differentially expressed genes | genome structural variation studies | displays and analyzes sequence data |

# [https://skewdb.org/](https://skewdb.org/) - [https://skewdb.org/view/](https://skewdb.org/view/) - CSV FILES!

## SkewDB, a comprehensive database of GC and 10 other skews for over ~~30,000~~ chromosomes and plasmids

50,000

Bert Hubert ✉

## Abstract

GC skew denotes the relative excess of G nucleotides over C nucleotides on the leading versus the lagging replication strand of eubacteria. While the effect is small, typically around 2.5%, it is robust and pervasive. GC skew and the analogous TA skew are a localized deviation from Chargaff's second parity rule, which states that G and C, and T and A occur with (mostly) equal frequency even within a strand. Different bacterial phyla show different kinds of skew, and differing relations between TA and GC skew. This article introduces an open access database ([https://skewdb.org](https://skewdb.org)) of GC and 10 other skews for over 30,000 chromosomes and plasmids. Further details like codon bias, strand bias, strand lengths and taxonomic data are also included. The *SkewDB* can be used to generate or verify hypotheses. Since the origins of both the second parity rule and GC skew itself are not yet satisfactorily explained, such a database may enhance our understanding of prokaryotic DNA.

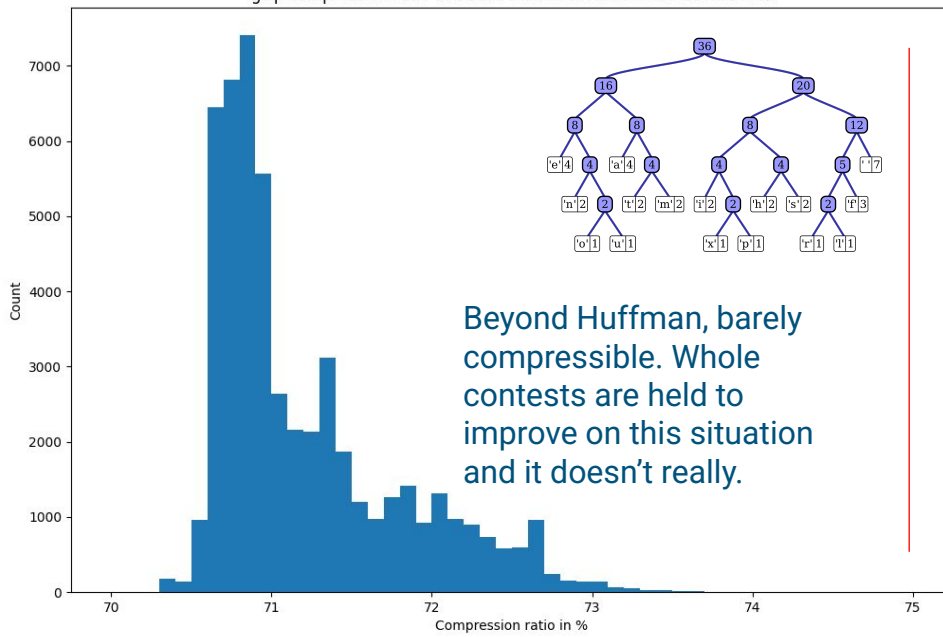Let's dive into the >50,000 binaries!
(in ASCII)

```
>NZ_CP150338.1 Sorangium sp. So ce388 chromosome, complete genome
ATGACGATTCCGAAACACGAGCCCCGCGAAGTCTTCGATCGGGCGATCGAGCATACGCGGGCCCTTTCTCCCGCAACTTT
TGATCAGTGGTTTGGGGGAGTTCAGTTCGATGACCTGACCGACGGCGTGCTCACGCTGCGAGTCCAGAACGAGTTCGTCC
TCGAGTGGGTCAGGGACAATTTCCTGCCCGCGCTGACCGACAAGATCCGCGAGATCACGGGCTGGTCGGTCCAGGTGGCG
TGGACGGTGGATCAGCACCTTGAGTCGCCGATCGCGCAGCGCGTCGAGCTCACCCCGGTTCGCCCGCGGGCGCTCGTGGT
GCGTCCGACGAGCACGGCGCCGACGCCTGCGCCCCGCCCCGAGCAGCCGCTGCGGCGCGTCTCGCCGATGCCCGACGACC
TCAACCCGAAGCACACCTTCGCGAGCTTCGTCGTGGGCCCGTCGAACCAGCTGGCGCACGCGGCCGCGATCGCGGCCGCG
GGCGGCGGGGGTCGCCGGTACAACCCGCTCTTCATCTGCGGCGGAACGGGGCTCGGCAAGACCCACCTGATGCACGCGAT
CGCGCATCGCGTCTTCGAGGGCAGGCCGGACGCGCGGATCATCTACGTCTCGGCGGAGAAGTTCACGAACGACTTCATCA
CGGCCATCCAGCACCACCGGATGGACGACTTCCGCACGAGGTACAGGTCGAGCTGCGACGTGCTGCTCGTCGACGACATC
CAGTTCCTGGCCGGGCGCGAGCAGACGCAGGAGGAGTTCTTCCACACCTTCAACGCGCTCCACACGCTCGACCGGCAGAT
CGTGGTGACGAGCGACAAGTACCCCCAGAACCTCGAGCGCATGGAGGAGCGCCTCGTCTCGCGTTTCTCGTGGGGGCTCG
TCGCCGACATCCAGGTGCCGGAGCTGGAGACGCGCGTCGCGATCGTCCGGAACAAGGCGGCGCTCGAGGGCATCGCGCTC
ACGGACGACGTGGCGCTCTACCTCGCGCAGATGGTCCGCTCGAACGTCCGCGAGCTCGAGGGGACGCTGATCCGGCTTGC
GGCCAAGAGCTCGCTCACGGGCCGCCCCGTGGATCTCCCGTTCGCGCGCGCCGAGATCACGGCCACGTCGCCGCCGCGGG
```

Beyond Huffman, barely compressible. Whole contests are held to improve on this situation and it doesn't really.

```
>NC_004337.2 Shigella flexneri 2a str. 301 chromosome, complete genome
AGCTTTTCATTCTGACTGCAACGGGCAATATGTCTCTGTGTGGATTAAAAAAAGAGTGTCTGATAGCAGCTTCTGAACTG
GTTACCTGCCGTGAGTAAATTAAAATTTTATTGACTTAGGTCACTAAATACTTTAACCAATATAGGCATAGCGCACAGAC
...
GACGAAATGCTGAACCAGGGCTGAAGCGTTGGTTTCTTTCCACNNNNNNCAGCTTCAGCCGTTATTGGTGCGATTGTCGT
GCTATTTATCTACAGGAAGATTAAAAGTTAACGCTTAAATTGCACAAAGGCTGCACACAGGCAGCCTTTGCTATTTTTTA
GAGGTGACGTTACCACCATCCCAGCTCAGTGCCGAGTACGACAATAATCGCCACACCCATCATAATAATCGTTGTTTTTT
TCATCTTTATGCTCCCGGGGCAGCATAGCAGCCAATAAAAAACCATGCTAAAAATACGACCGCTGTAATGAAGATTACAA
CCGGAAAAATAATACCGATTCGCATGTTATCACCAATCAAAAGGATGTGAATACGCCTCAGGAATGTAGCGCTGGATGCG
```

"2bit" formats would achieve 75%
compression, and allow better
further compression, but the field
doesn't care.

>NZ_CP043307.1 Acinetobacter johnsonii strain Acsw19
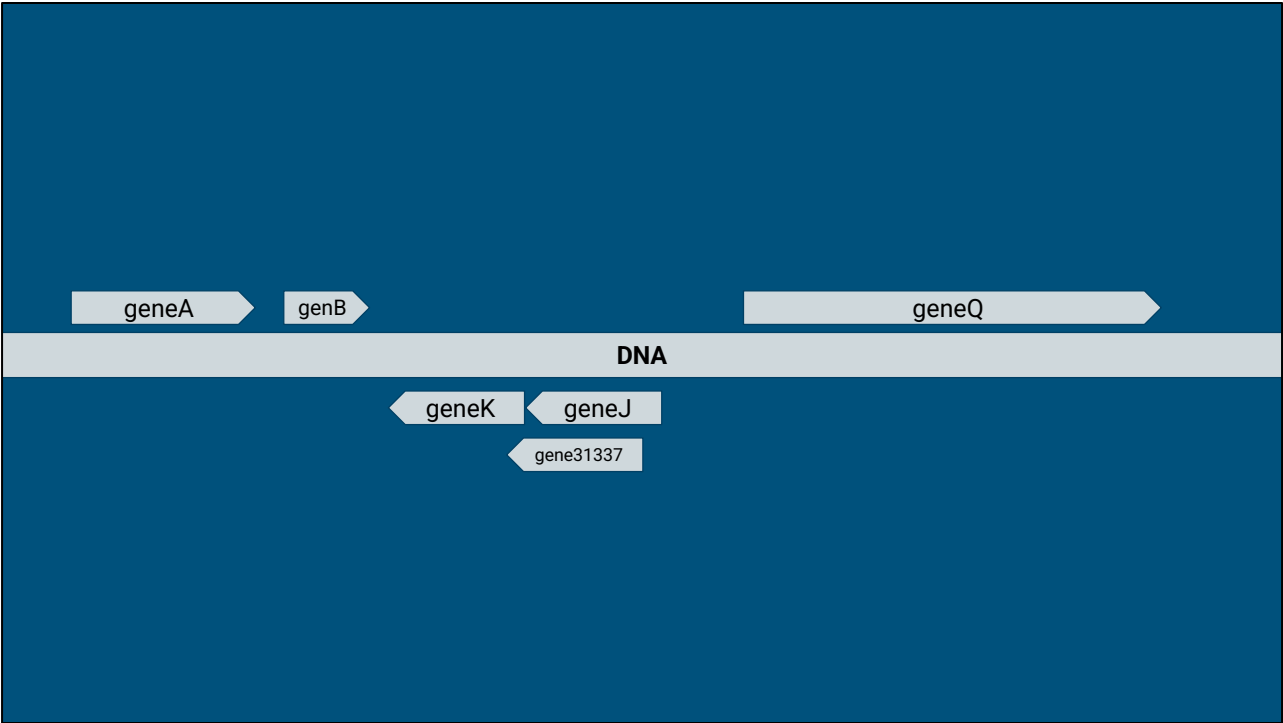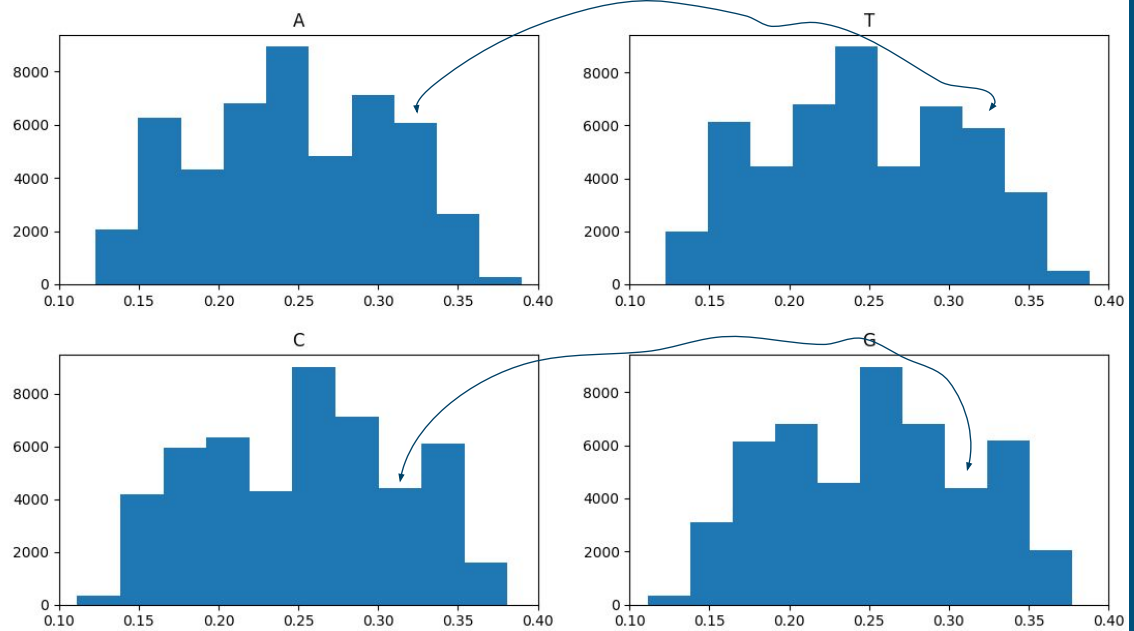ATGCTTTGGACGGACTGCTTAACTCGCTTGCGACAAGAGCTCTCTGGGAATGTCTTTACAATGT

```
A    C
|    |
T    G
```

**ATGCTTTGGACGGACTGCTTAACTCGCTTGCGACAAGAGCTCTCTGGGAATGTCTTTACAATGT**
|||||||||||||||||||||||||||||||||||||*|||||||||||||||||||||||||||||||||||||
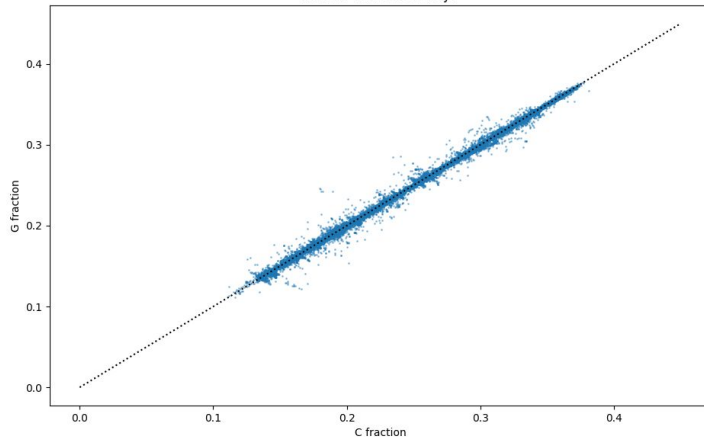**TACGAAACCTGCCTGACGAATTGAGCGAACGCTGTTCTCGAGAGACCCTTACAGAAATGTTACA**

>NZ_CP043307.1.**rev** Acinetobacter johnsonii strain Acsw19
ACATTGTAAAGACATTCCCAGAGAGCTCTTGTCGCAAGCGAGTTAAGCAGTCCGTCCAAAGCAT

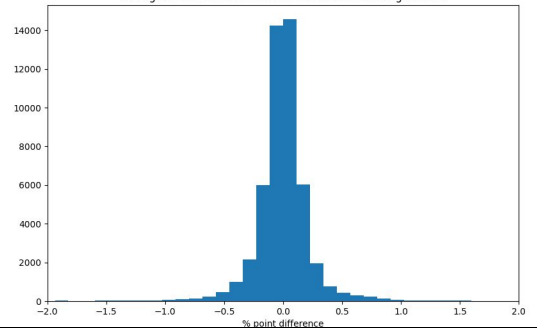Histogram of nucleotide fraction for 49366 chromosomes

For 49366 bacterial genomes
And no one knows why!

"The second Chargaff rule holds that both A% ≈ T% and G% ≈ C% are valid for **each** of the two DNA strands"

"The basis for this rule is still under investigation" (!!!)

Histogram of C fraction minus G fraction for 49366 genomes

https://en.wikipedia.org/wiki/Chargaff%27s_rules

# But wait, it gets weirder

Complementary DNA:

C -> G
A -> T

Reverse & complement:

**CTA** -> GAT -> **TAG**

**ATG** -> TAC -> CAT

- Do the complement thing
- Reverse the string

"Reverse Complement" (RC)

```
ATGCTTTGGACGGACTGCTTAACTCGCTTGCGACAAGAGCTCTCTGGGAATGTCTTTACAATGT
||||||||||||||||||||||||||||||||||*|||||||||||||||||||||||||||||
TACGAAACCTGCCTGACGAATTGAGCGAACGCTGTTCTCGAGAGACCCTTACAGAAATGTTACA
```

Count of trigram/codon versus count of its reverse complement
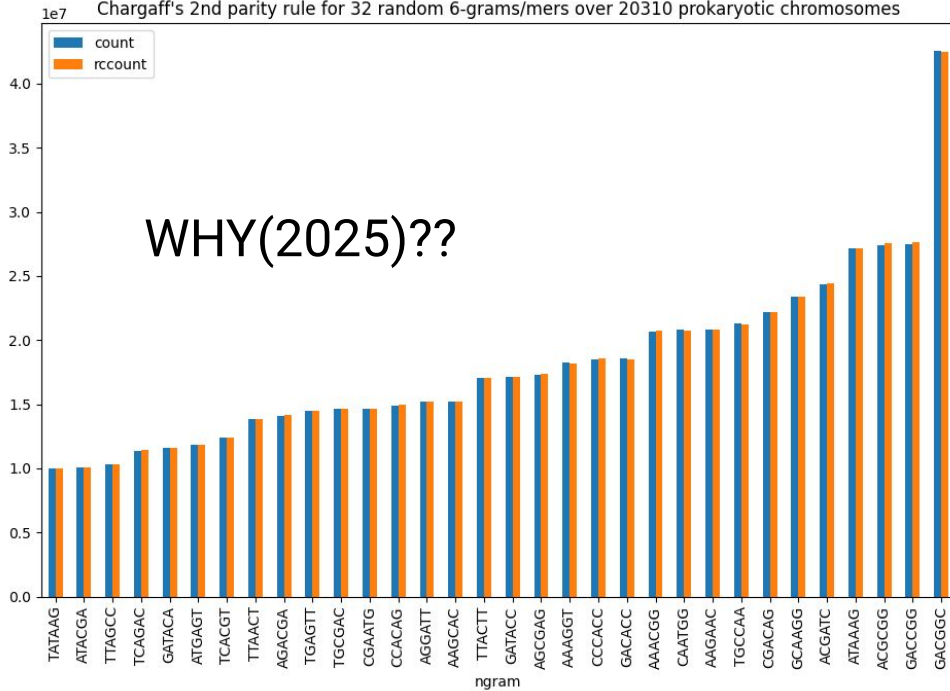In a single bacterial chromosome

Occurrence of ngrams in DNA compared to their reverse complement
In a single bacterial chromosome

Count of 6-mer versus count of its reverse complement
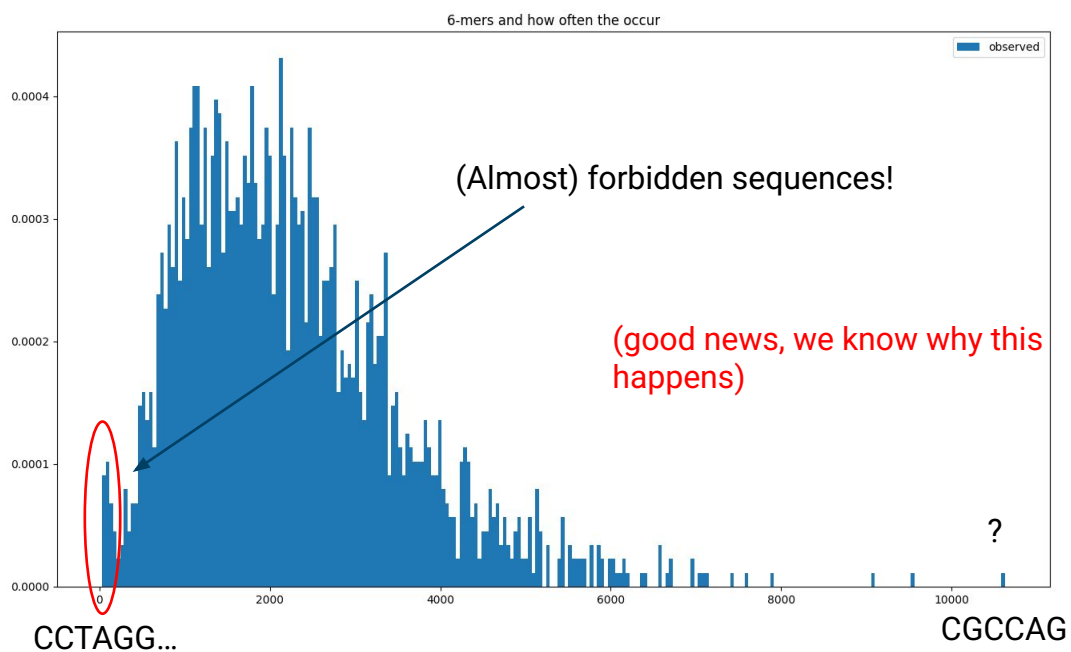In a single bacterial chromosome

But wait, now
for longer
sequences,
like CCGATT

WHY(2025)??

Chargaff's 2nd parity rule for 32 random 6-grams/mers over 20310 prokaryotic chromosomes

WHY(2025)??

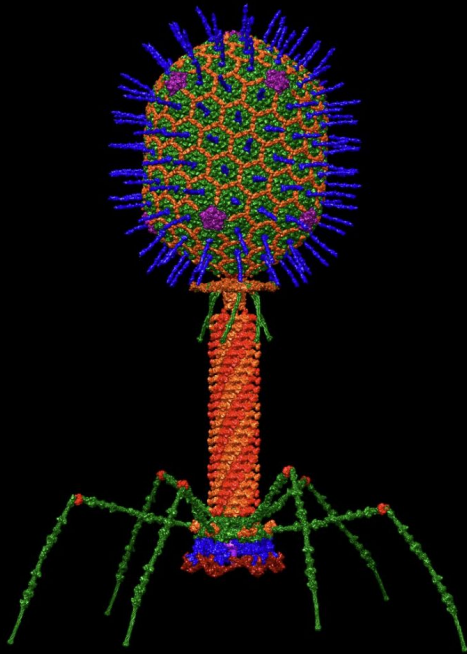We need to definitively find out why this is so, and if it means something!

It works up to *10* characters at least.

| | ngram | rcngram | count | rccount | procdiff | totcount |
|---|---|---|---|---|---|---|
| 839 | CCTAGG | CCTAGG | 16 | 16 | 0.00000 | 32 |
| 1004 | CTAGGA | TCCTAG | 17 | 21 | 23.52940 | 38 |
| 837 | CCTAGA | TCTAGG | 23 | 19 | -17.39130 | 42 |
| 997 | CTAGAC | GTCTAG | 19 | 27 | 42.10530 | 46 |
| 1006 | CTAGGG | CCCTAG | 28 | 37 | 32.14290 | 65 |
| ... | ... | ... | ... | ... | ... | ... |
| 679 | CAGCGC | GCGCTG | 3947 | 3629 | -8.05675 | 7576 |
| 736 | CCAGCA | TGCTGG | 4072 | 3804 | -6.58153 | 7876 |
| 738 | CCAGCG | CGCTGG | 4598 | 4489 | -2.37060 | 9087 |
| 1291 | GCCAGC | GCTGGC | 4815 | 4750 | -1.34995 | 9565 |
| 913 | CGCCAG | CTGGCG | 5372 | 5253 | -2.21519 | 10625 |

6-mers and how often the occur

(Almost) forbidden sequences!

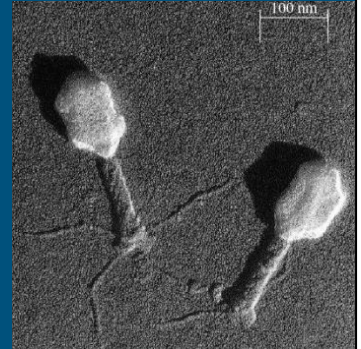(good news, we know why this happens)

?

CCTAGG…                                    CGCCAG

Bacteria have viruses too..

T4

Dr. Victor Padilla-Sanchez, PhD https://www.drvictorpadillasanchez.com

https://en.wikipedia.org/wiki/Bacteriophage#/media/File:PhageExterior.svg
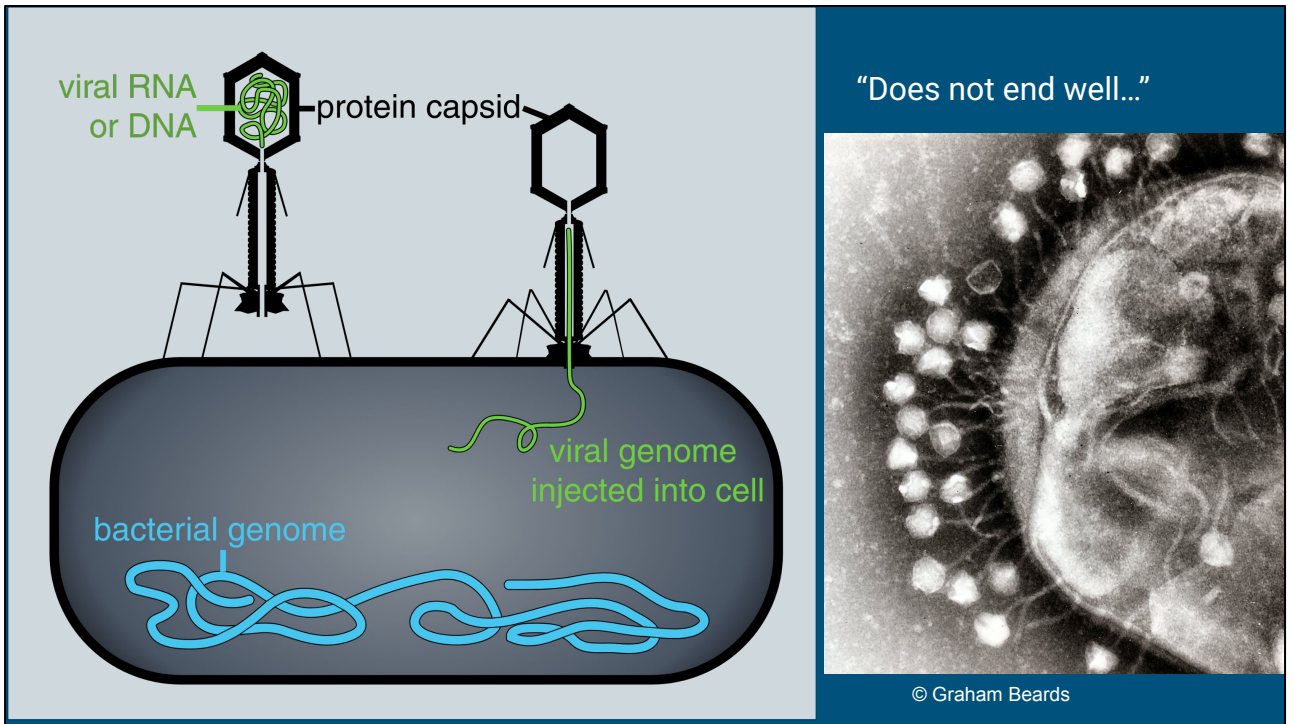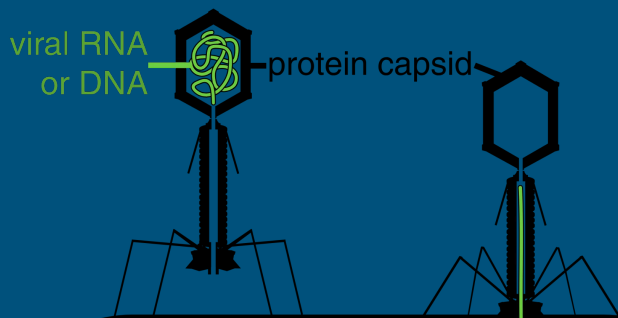http://stl-bjb.ac-dijon.fr/spip.php?article32
https://en.wikipedia.org/wiki/Escherichia_virus_T4
https://commons.wikimedia.org/wiki/File:Bacteriophage_T4_Structural_Model_at_Atomic_Resolution.tif?page=1

Viral RNA or DNA — protein capsid

viral genome injected into cell

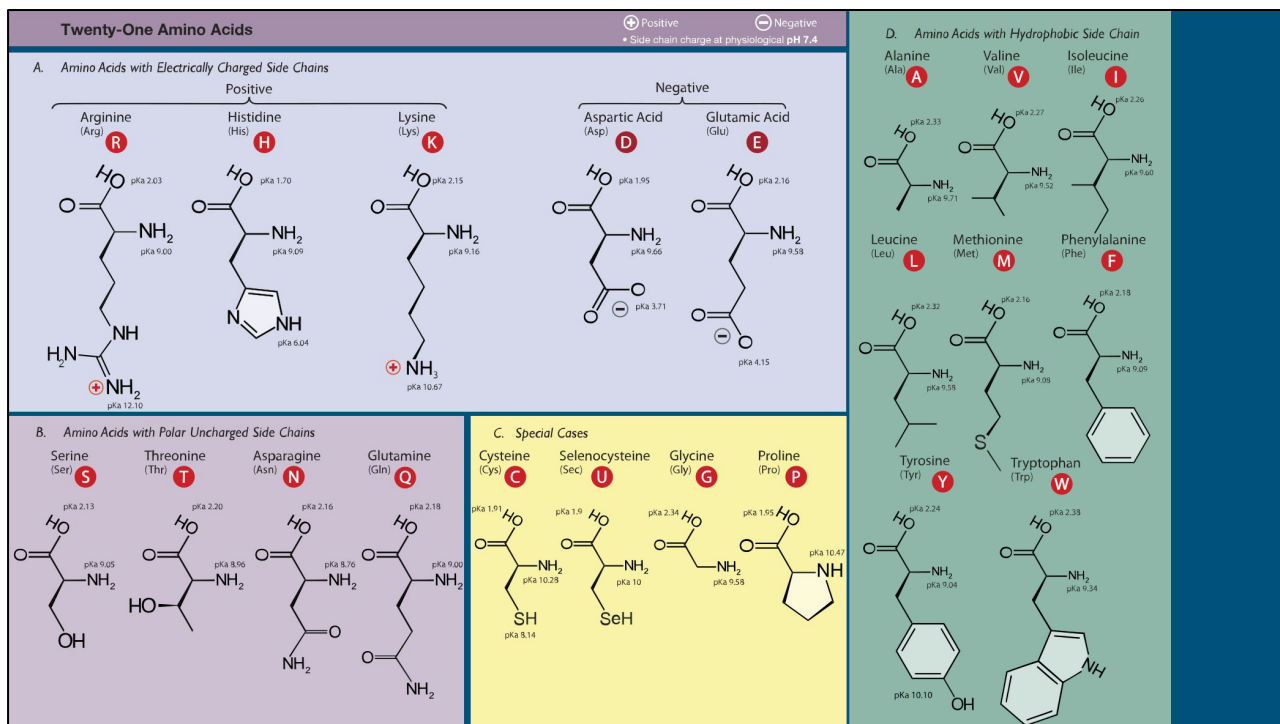bacterial genome

"Does not end well…"

© Graham Beards

# AvrII

Avrll has been reformulated with Recombinant Albumin (rAlbumin) beginning with Lot #10128047. **Learn more.**

We are excited to announce that all reaction buffers are now BSA-free. NEB began switching our BSA-containing reaction buffers in April 2021 to buffers containing **Recombinant Albumin** (rAlbumin) for restriction enzymes and some DNA modifying enzymes. Find more details at **www.neb.com/BSA-free**.

```
5´...C CTAGG...3´
3´...GGATC C...5´
```

**Isoschizomers | Single Letter Code | Pronunciation:**

- Time-Saver™ qualified for digestion in 5-15 minutes

- 100% activity in rCutSmart™ Buffer (over 210 enzymes are available in the same buffer) simplifying double digests

- Supplied with 1 vial of Gel Loading Dye, Purple (6X)

- Restriction Enzyme Cut Site: C/CTAGG

Modified from https://commons.wikimedia.org/wiki/File:Amino_Acids.svg

| 1st base | | | | | | | | | 3rd base |
|---|---|---|---|---|---|---|---|---|---|
| | 2nd base | | | | | | | | |
| | T | | C | | A | | G | | |
| T | TTT | (Phe/F) Phenylalanine (np) | TCT | (Ser/S) Serine (p) | TAT | (Tyr/Y) Tyrosine (p) | TGT | (Cys/C) Cysteine (p) | T |
| | TTC | | TCC | | TAC | | TGC | | C |
| | TTA | (Leu/L) Leucine (np) | TCA | | TAA | Stop (Ochre) *[note 2] | TGA | Stop (Opal) *[note 2] | A |
| | TTG ⇒ | | TCG | | TAG | Stop (Amber) *[note 2] | TGG | (Trp/W) Tryptophan (np) | G |
| C | CTT | (Leu/L) Leucine (np) | CCT | (Pro/P) Proline (np) | CAT | (His/H) Histidine (b) | CGT | (Arg/R) Arginine (b) | T |
| | CTC | | CCC | | CAC | | CGC | | C |
| | CTA | | CCA | | CAA | (Gln/Q) Glutamine (p) | CGA | | A |
| | CTG | | CCG | | CAG | | CGG | | G |
| A | ATT | (Ile/I) Isoleucine (np) | ACT | (Thr/T) Threonine (p) | AAT | (Asn/N) Asparagine (p) | AGT | (Ser/S) Serine (p) | T |
| | ATC | | ACC | | AAC | | AGC | | C |
| | ATA | | ACA | | AAA | (Lys/K) Lysine (b) | AGA | (Arg/R) Arginine (b) | A |
| | ATG ⇒ | (Met/M) Methionine (np) | ACG | | AAG | | AGG | | G |
| G | GTT | (Val/V) Valine (np) | GCT | (Ala/A) Alanine (np) | GAT | (Asp/D) Aspartic acid (a) | GGT | (Gly/G) Glycine (np) | T |
| | GTC | | GCC | | GAC | | GGC | | C |
| | GTA | | GCA | | GAA | (Glu/E) Glutamic acid (a) | GGA | | A |
| | GTG ⇒ | | GCG | | GAG | | GGG | | G |

Multi-billion year old table!

Multiple codons for same amino acids

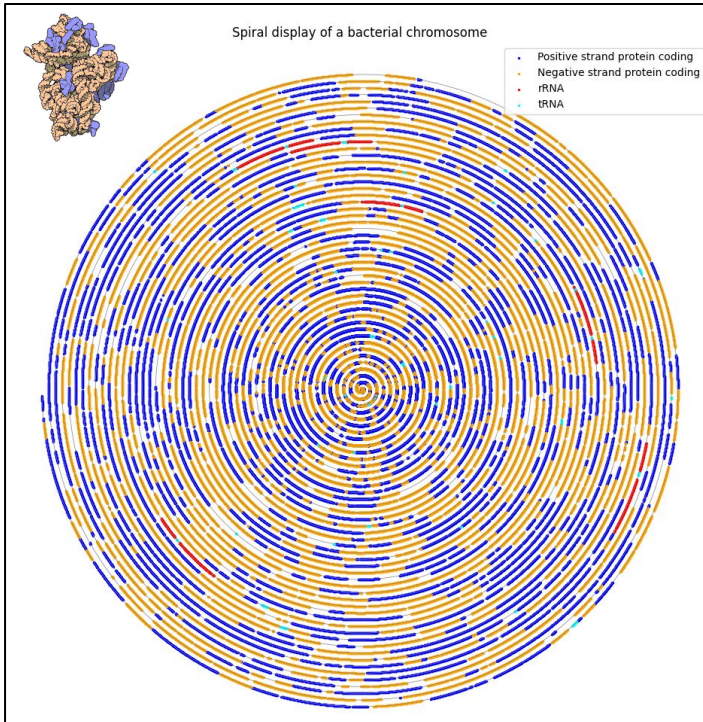This allows for **dialects** and shaping DNA

Where do genes begin and end?

HNGGGGGG!!!

FASTA

GFF

DNA layer
RNA layer
Amino acid layer
..
RNA layer
DNA layer

RNA

Amino Acids

Amino acids code

Promoter        -35 .......        -10                    Shine-Dalgarno    ATG...TGA

"TTGACA"        "TATAAT"                    "AGGAGG"

Spiral display of a bacterial chromosome

Legend:
- Positive strand protein coding
- Negative strand protein coding
- rRNA
- tRNA

Escherichia coli str. K-12 substr. W3110

Topologically this chromosome is actually a circle.

Spiral shape however is better for an overview

The blank areas are **not genes**. Partially we know what this is. Partially not!

Biologists study genes.. "Looking for your keys where the light is"

**But us nerds could take a look at the data!**

```
Usage:  prodigal [-a trans_file] [-c] [-d nuc_file] [-f output_type]
                 [-g tr_table] [-h] [-i input_file] [-m] [-n] [-o output_file]
                 [-p mode] [-q] [-s start_file] [-t training_file] [-v]
```
So standard that Debian ships it! (Also, congrats on Trixie!)
```
         -a:  Write protein translations to the selected file.
         -c:  Closed ends.  Do not allow genes to run off edges.
         -d:  Write nucleotide sequences of genes to the selected file.
         -f:  Select output format (gbk, gff, or sco).  Default is gbk.
         -g:  Specify a translation table to use (default 11).
         -h:  Print help menu and exit.
         -i:  Specify FASTA/Genbank input file (default reads from stdin).
         -m:  Treat runs of N as masked sequence; don't build genes across them.
         -n:  Bypass Shine-Dalgarno trainer and force a full motif scan.
         -o:  Specify output file (default writes to stdout).
         -p:  Select procedure (single or meta).  Default is single.
         -q:  Run quietly (suppress normal stderr output).
         -s:  Write all potential genes (with scores) to the selected file.
         -t:  Write a training file (if none exists); otherwise, read and use
              the specified training file.
         -v:  Print version number and exit.
```

```
ahu@xeon:~/skewdb/skewdb-articles/antonie2$ prodigal -i in -o genes
--------------------------------------
PRODIGAL v2.6.3 [February, 2016]
Univ of Tenn / Oak Ridge National Lab
Doug Hyatt, Loren Hauser, et al.
--------------------------------------
Request:  Single Genome, Phase:  Training
Reading in the sequence(s) to train...4646332 bp seq created, 50.80 pct GC
Locating all potential starts and stops...241263 nodes
Looking for GC bias in different frames...frame bias scores: 1.54 0.18 1.27
Building initial set of genes to train from...done!
Creating coding model and scoring nodes...done!
Examining upstream regions and training starts...done!
--------------------------------------
Request:  Single Genome, Phase:  Gene Finding
Finding genes in sequence #1 (4646332 bp)...done!
```

# Realloc zero, fixes undefined behaviour #119

🎋 Open   **berthubert** wants to merge 2 commits into `hyattpd:GoogleImport` from `berthubert:realloc-zero` ⎘

💬 Conversation  0    ⊶ Commits  2    ⊟ Checks  0    ⊞ Files changed  1

**berthubert** commented 4 days ago                                                    ...

Valgrind saw Prodigal access uninitialized memory. This was because resizing the nodes array using realloc did not zero the newly allocated memory. This could have had consequences for generated gene predictions if you were unlucky.

This PR makes sure that the new memory is zeroed as well. In addition, an error message if /dev/stdin was not available was fixed to not crash.

⤴ **berthubert** added 2 commits 4 days ago

⊶  🔲 fix error report if /dev/stdin could not be opened                           d4e3d76

⊶  🔲 when resizing the nodes array, the newly allocated space was not zero… ...   4e3b87b

```
485  491
486  492            /* Reallocate memory if this is the biggest sequence we've seen */
487  493            if(slen > max_slen && slen > STT_NOD*8) {
488       -              nodes = (struct _node *)realloc(nodes, (int)(slen/8)*sizeof(struct _node));
     494  +              size_t newnodesize = (int)(slen/8)*sizeof(struct _node);
     495  +              nodes = (struct _node *)realloc(nodes, newnodesize);
489  496                if(nodes == NULL) {
490  497                  fprintf(stderr, "Realloc failed on nodes\n\n");
491  498                  exit(11);
492  499                }
     500  +              memset( ((char*) &nodes[0]) + nodesize, 0, newnodesize-nodesize);
     501  +              nodesize = newnodesize;
493  502                max_slen = slen;
494  503            }
495  504
```

The state of bioinformatics software is… not great (second time this happened to me)

More bacterial anti-viral defenses.

This war has been raging for 2 billion years at least!

Maybe we could learn..

A multi-generational immune system.

A forensic record of previously survived viruses!

**Can we see it?**

**People looked at this for 20 years w/o knowing what it ws**

By James atmos - Own work, CC BY-SA 3.0, https://commons.wikimedia.org/w/index.php?curid=7821536

# CRISPR: "The origin of the spacer sequences remains unknown" - 2002

Signature:
"AGTTTCCGTATCTCCGGATTTATAAAGCTGA"

Why does the CRISPR system not destroy itself?

GC SKEW!
https://skewdb.org/view

GC and TA skew in random bacterial chromosomes

**And no one really knows why!**

Data for 7970 Firmicute chromosomes

Lots of data available to test hypotheses

I also have a few

Have a go with the data from skewdb.org!

Spiral display of a bacterial chromosome

- Positive strand protein coding
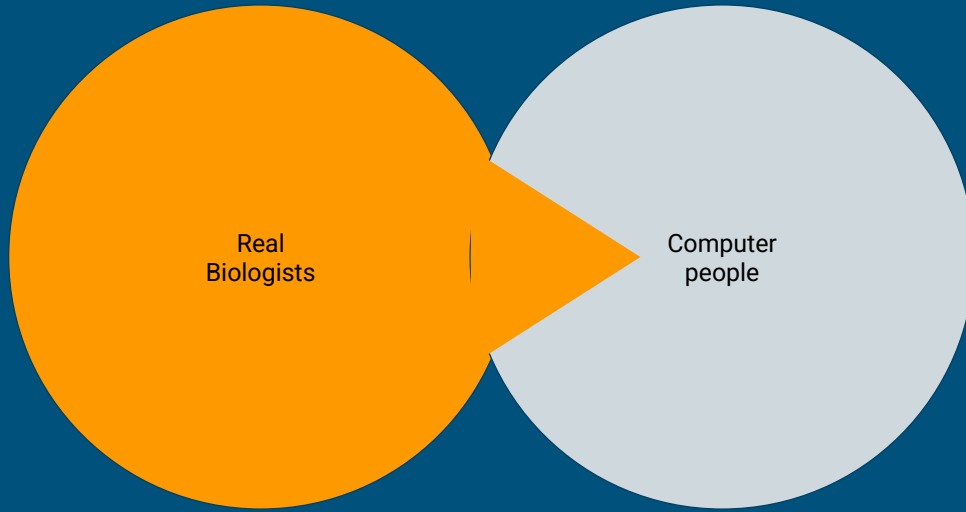- Negative strand protein coding
- rRNA
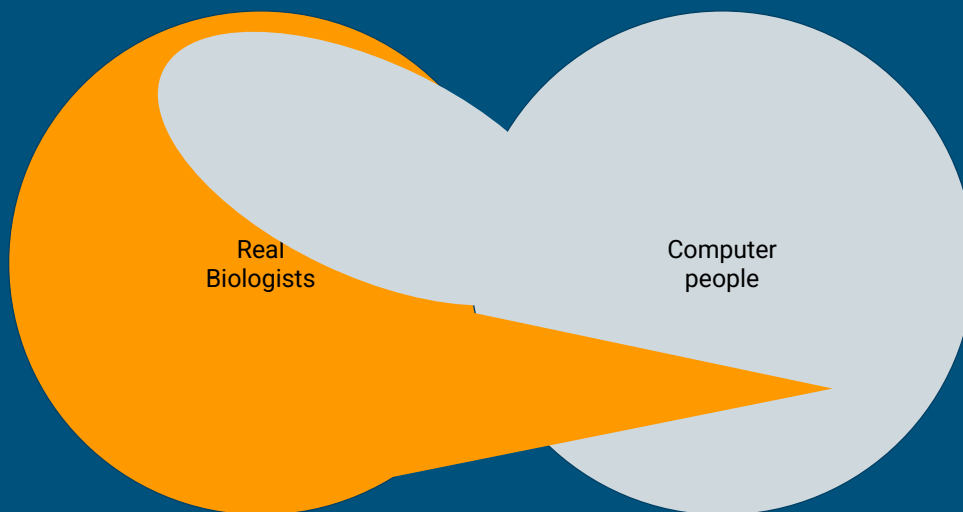- tRNA

Escherichia coli str. K-12 substr. W3110

Look at the white spaces

Very interesting things could be hiding there

**And you could help find out what it is!**

Why does GC-skew exist?
Why even Chargaff's 2nd rule?
What is in the bacterial whitespace?
Why is so much bioinformatics tooling so terrible?

**GOOD LUCK!**

# Come to the afterparty!

Tomorrow, Monday, 2025-08-11 15:00–15:50, Cassiopeia

Interactive session, featuring all the skipped slides and lots of room for questions and answers!



https://berthub.eu/revdna/