



# Геном

Андрей Владимирович  
Алексеевский  
[aba@belozersky.msu.ru](mailto:aba@belozersky.msu.ru)

# На примере вируса SARS-CoV-2

>NC\_045512.2 Wuhan seafood market pneumonia virus isolate Wuhan-Hu-1, complete genome

```
ATTAAAGGTTTATACCTTCCCAGGTAACAAACCAACCAACTTTCGATCTCTTGTAGATCTGTTCTCTAAA  
CGAACTTTAAAATCTGTGTGGCTGTCACCTCGGCTGCATGCTTAGTGCACCTCACGCAGTATAATTAATAAC  
TAATTA CTGTCGTTGACAGGACACGAGTAACTCGTCTATCTTCTGCAGGCTGCTTACGGTTTTCGTCCGTG  
TTGCAGCCGATCATCAGCACATCTAGGTTTCGTCCGGGTGTGACCGAAAGGTAAGATGGAGAGCCTTGTC  
CCTGGTTTTCAACGAGAAAACACACGTCCAACCTCAGTTTGCCTGTTTTACAGGTTTCGCGACGTGCTCGTAC  
GTGGCTTTGGAGACTCCGTGGAGGAGGTCTTATCAGAGGCACGTCAACATCTTAAAGATGGCACTTGTGG  
CTTAGTAGAAGTTGAAAAAGGCGTTTTTGCCTCAACTTGAACAGCCCTATGTGTTTCATCAAACGTTCCGGAT  
GCTCGAACTGCACCTCATGGTCATGTTATGGTTGAGCTGGTAGCAGA ACTCGAAGGCATTCAGTACGGTC  
GTAGTGGTGAGACACTTGGTGTCCTTGTCCCTCATGTGGGCGAAATACCAGTGGCTTACCGCAAGGTTCT  
TCTTCGTAAGAACGGTAATAAAGGAGCTGGTGGCCATAGTTACGGCGCCGATCTAAAGTCATTTGACTTA  
GGCGACGAGCTTGGCACTGATCCTTATGAAGATTTTCAAGAAA ACTGGAACACTAAACATAGCAGTGGTG  
TTACCCGTGAACTCATGCGTGAGCTTAACGGAGGGGCATACACTCGCTATGTCGATAACA ACTTCTGTGG  
CCCTGATGGCTACCCTCTTGAGTGCATTAAGACCTTCTAGCACGTGCTGGTAAAGCTTCATGCACTTTG  
TCCGAACA ACTGGACTTTATTGACACTAAGAGGGGTGTATACTGCTGCCGTGAACATGAGCATGAAATTG  
CTTGGTACACGGAACGTTCTGAAAAGAGCTATGAATTGCAGACACCTTTTGAAATTAAATTGGCAAAGAA  
ATTTGACACCTTCAATGGGGAATGTCCAAATTTTGTATTTCCCTTAAATTCATAATCAAGACTATTCAA  
CCAAGGGTTGAAAAGAAAAGCTTGATGGCTTTATGGGTAGAATTCGATCTGTCTATCCAGTTGCGTCCAC
```

CCAAGGGTTGAAAAGAAAAAGCTTGATGGCTTTATGGGTAGAATTCGATCTGTCTATCCAGTTGCGTCAC  
CAAATGAATGCAACCAAATGTGCCTTTCAACTCTCATGAAGTGTGATCATTGTGGTGAACTTCATGGCA  
GACGGGCGATTTTGTAAAGCCACTTGCGAATTTTGTGGCACTGAGAATTTGACTAAAGAAGGTGCCACT  
ACTTGTGGTTACTTACCCCAAATGCTGTTGTTAAAATTTATTGTCCAGCATGTCACAATTCAGAAGTAG  
GACCTGAGCATAGTCTTGCCGAATACCATAATGAATCTGGCTTGAAAACCATTCTTCGTAAGGGTGGTTCG  
CACTATTGCCTTTGGAGGCTGTGTGTTCTCTTATGTTGGTTGCCATAACAAGTGTGCCTATTGGGTTCCA  
CGTGCTAGCGCTAACATAGGTTGTAACCATAACAGGTGTTGTTGGAGAAGGTCCGAAGGTCTTAATGACA  
ACCTTCTTGAAATACTCCAAAAAGAGAAAGTCAACATCAATATTGTTGGTGACTTTAAACTTAATGAAGA  
GATCGCCATTATTTTGGCATCTTTTTCTGCTTCCACAAGTGCTTTTGTGGAACTGTGAAAGGTTTGGAT  
TATAAAGCATTCAAACAAATTGTTGAATCCTGTGGTAATTTTAAAGTTACAAAAGGAAAAGCTAAAAAAG  
GTGCCTGGAATATTGGTGAACAGAAATCAATACTGAGTCCTCTTTATGCATTTGCATCAGAGGCTGCTCG  
TGTTGTACGATCAATTTTCTCCCGCACTCTTGAAACTGCTCAAATTTCTGTGCGTGTTTTACAGAAGGCC  
GCTATAACAATACTAGATGGAATTTCACAGTATTCACTGAGACTCATTGATGCTATGATGTTCCACATCTG  
ATTTGGCTACTAACAATCTAGTTGTAATGGCCTACATTACAGGTGGTGTGTTGTTTCCAGTTGACTTCGCAGTG  
GCTAACTAACATCTTTGGCACTGTTTATGAAAACTCAAACCCGTCCTTGATTGGCTTGAAGAGAAGTTT  
AAGGAAGGTGTAGAGTTTCTTAGAGACGGTTGGGAAATTGTTAAATTTATCTCAACCTGTGCTTGTGAAA  
TTGTCGGTGGACAAATTGTCACCTGTGCAAAGGAAATTAAGGAGAGTGTTTCCAGACATTCTTTAAGCTTGT  
AAATAAATTTTTGGCTTTGTGTGCTGACTCTATCATTATTGGTGGAGCTAAACTTAAAGCCTTGAATTTA  
GGTGAACATTTGTCACGCACTCAAAGGGATTGTACAGAAAGTGTGTTAAATCCAGAGAAGAACTGGCC  
TACTCATGCCTCTAAAAGCCCCAAAAGAAATTATCTTCTTAGAGGGAGAAACACTTCCACAGAAGTGTT  
AACAGAGGAAGTTGTCTTGAAACTGGTGAATTTACAACCATTAGAACAACCTACTAGTGAAGCTGTTGAA  
GCTCCATTGGTTGGTACACCAGTTTGTATTAACGGGCTTATGTTGCTCGAAATCAAAGACACAGAAAAGT  
ACTGTGCCCTTGCACCTAATATGATGGTAACAACAATACCTTCACACTCAAAGGCGGTGCACCAACAAA  
GGTACTTTTTGGTGATGACACTGTGATAGAAGTGCAAGGTTACAAGAGTGTGAATATCACTTTTGAACCTT  
GATGAAAGGATTGATAAAGTACTTAATGAGAAGTGCTCTGCCTATACAGTTGAACTCGGTACAGAAGTAA  
ATGAGTTCGCCTGTGTTGTGGCAGATGCTGTCATAAAAACCTTTGCAACCAGTATCTGAATTACTTACACC  
ACTGGGCATTGATTTAGATGAGTGGAGTATGGCTACATACTACTTATTTGATGAGTCTGGTGAGTTTAAA<sup>3</sup>

И еще 27 страниц такого текста ...

Всего 29 903 букв.

Геном содержит информацию:  
инструкцию для клеток организма хозяина  
(человека) как размножить вирус  
SARS-CoV-2. При этом хозяин заболевает!!!

Информация: Текст и читатель

# Что такое информация?

Армянское радио

- «Правда ли, что Иштоян выиграл в лотерею машину?»
- «**Правда.** Но не Иштоян, а Петросян, не машину, а швейную машинку; не в лотерею, а в карты; и не выиграл, а проиграл»

*«Сколько информации в этом сообщении?»  
И.М.Гельфанд*

# Анекдоты нам устраивают и вирусы



Вирус друг человека? К.Северинов



Гены и сигналы

## **2. ЧТО ЗАПИСАНО В ГЕНОМЕ?**

# Что видим своими глазами

>NC\_045512.2 Wuhan seafood market pneumonia virus isolate Wuhan-Hu-1, complete genome **SARS-CoV-2**

```
ATTAAAGGTTTATACCTTCCCAGGTAACAAACCAACCAACTTTCGATCTCTTGTAGATCTGTTCTCTAAA  
CGAACTTTAAAATCTGTGTGGCTGTCACTCGGCTGCATGCTTAGTGCACCTCACGCAGTATAATTAATAAC  
TAATTA CTGTCGTTGACAGGACACGAGTAACTCGTCTATCTTCTGCAGGCTGCTTACGGTTTTCGTCCGTG  
TTGCAGCCGATCATCAGCACATCTAGGTTTTCGTCCGGGTGTGACCGAAAGGTAAGATGGAGAGCCTTGTC  
CCTGGTTTTCAACGAGAAAACACACGTCCAACCTCAGTTTGCCTGTTTTTACAGGTTTCGCGACGTGCTCGTAC  
GTGGCTTTGGAGACTCCGTGGAGGAGGTCTTATCAGAGGCACGTCAACATCTTAAAGATGGCACTTGTGG  
CTTAGTAGAAGTTGAAAAGGCGTTTTTGCCTCAACTTGAACAGCCCTATGTGTTTCATCAAACGTTCCGGAT  
GCTCGAACTGCACCTCATGGTCATGTTATGGTTGAGCTGGTAGCAGAACTCGAAGGCATTCAGTACGGTC  
GTAGTGGTGAGACACTTGGTGTCCTTGTCCTCATGTGGGCGAAATACCAGTGGCTTACCGCAAGGTTCT  
TCTTCGTAAGAACGGTAATAAAGGAGCTGGTGGCCATAGTTACGGCGCCGATCTAAAGTCATTTGACTTA  
GGCGACGAGCTTGGCACTGATCCTTATGAAGATTTTCAAGAAAACCTGGAACACTAAACATAGCAGTGGTG  
TTACCCGTGAACTCATGCGTGAGCTTAACGGAGGGGCATACACTCGCTATGTGATAACAACCTTCTGTGG  
CCCTGATGGCTACCCTCTTGAGTGCATTAAGACCTTCTAGCACGTGCTGGTAAAGCTTCATGCACTTTG  
TCCGAACAACCTGGACTTTATTGACACTAAGAGGGGTGTATACTGCTGCCGTGAACATGAGCATGAAATTG  
CTTGGTACACGGAACGTTCTGAAAAGAGCTATGAATTGCAGACACCTTTTGAATTAATTTGGCAAAGAA  
ATTTGACACCTTCAATGGGGAATGTCCAAATTTTGTATTTCCCTTAAATTCATAATCAAGACTATTCAA  
CCAAGGGTTGAAAAGAAAAGCTTGATGGCTTTATGGGTAGAATTCGATCTGTCTATCCAGTTGCGTCAC
```



# Что видим своими глазами

- В геноме четыре буквы А, Т, G, С  
(понятно)
- Буквы идут неупорядоченно,  
похоже на случайную последовательность

# Лингвистический анализ текста

- Правда ли, что в файле в последовательности генома нет других букв?
- Частоты букв
- Часто и редко встречающиеся слова
- Равномерность частоты букв и слов вдоль текста

Эти вопросы изучаются и имеют биологически смысл! Примеры наблюдений:

- $\#C \approx \#G$ ,  $\#T \approx \#A$  ( $\#$  = число)
- Слов CG *мало* в определенных геномах
- Слов TA *мало* во всех геномах
- В некоторых геномах  $\#C > \#G$  в одной части и  $\#G > \#C$  в другой части («GC skew»)

## “Много, нормально, или мало”

- Чтобы ответить надо знать сколько - нормально: сколько изучаемых слов ожидается, если предположить, что никакой причины, влияющей на число слов нет – чистая случайность!
  - Можно предположить, что буквы А, Т, G, С в геноме имеют одинаковую частоту  $\frac{1}{4}$  и проверить так ли это в вашем геноме
  - Можно использовать наблюдаемые в вашем геноме частоты букв и вычислить сколько слов ТА ожидается в вашем геноме, если соседние буквы встречаются случайно и независимо друг от друга и сравнить с наблюдаемым в геноме числом слов ТА
    - Подсказка: как вычисляется вероятность (частота) двух независимых событий, если вероятности каждого из них известны? Например, события (1) увидеть ворону по дороге в МГУ для сдачи зачёта и (2) получить зачёт похоже независимы.

# Всерьёз думают о живых вакцинах, основанных на вирусах с увеличенным числом CG или TA!

Interestingly, most mammalian RNA viruses have low frequencies of CpGs ([45,46](#)). Furthermore, viruses with high CpG frequencies may be more recognizable by pathogen innate immune sensors ([47–50](#)).

Attenuation of the classical oral poliovirus vaccine is based on very few point mutations, which can revert to virulence after a few rounds of viral replication ([144](#)). These pioneering results obtained with recoded polioviruses suggest that codon-usage in recoded viruses may be much more stable than most RNA virus point mutants, and could possibly enable the development of live attenuated RNA virus vaccines with superior genetic stability.

Martinez et al., 2019, NAR

**КОНЕЦ ПРЕЗЕНТАЦИИ**