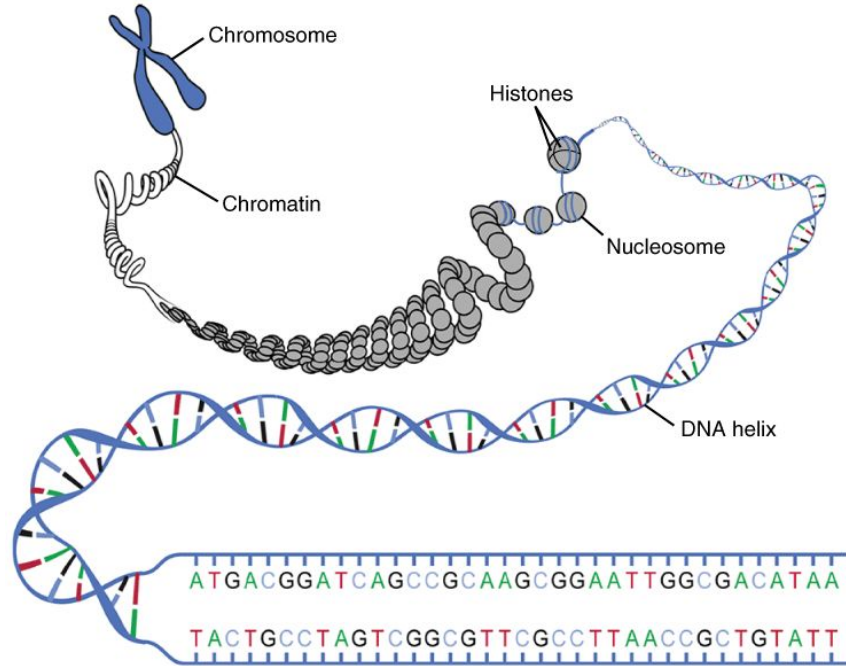
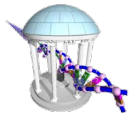


Comp 555 - BioAlgorithms - Spring 2021



Finding Patterns in DNA



Login to Course Website

1) Login to your Comp555 account

Logged in as: *guest* [Log in](#)

mcmillan@unc.edu

Home Research Courses Publications

Announcements

- **January 9:** First class meeting in SN014. See you there

Course Description

2) Your username is your UNC ONYEN and password is your PID

Username:

Password:

Login



Next Steps

3) Once you are logged in, press "Course" and then a "Setup" button should appear. Press "Setup" and you should see something like:

Comp555S20 Problem Sets and Exams:

Comp555S20 Exercises:

Exercises:

leehart has submitted 1 of 0 exercises

Exercise01:

<https://forms.g1e/f6y85beL8Hw5zoF47>

Your Profile

Username: leehart

First Name: Lee

Last Name: Hart

Email:

Institution:

New Password:

Verify Password:

4) (BTW, you can also change your password here if you want).



For those without a login...

- Go back to the login page, and click "registered"

Username:

Password:

Login

No password is required to logon as "guest"
You must be **registered** to have full access or modify content.

- Then enter the following information:
- Once registered a screen will indicate you've been verified; then click "Course" and "Setup" as before.

Username: MUST be your ONYEN

First Name:

Last Name: Your UNC email

Email:

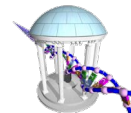
Institution:

Password:

Verify Password:

Register

You've seen a genome... let's scale up

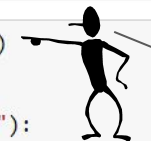


```
In [15]: def loadFasta(filename):
        """ Parses a classically formatted and possibly
            compressed FASTA file into a list of headers
            and fragment sequences for each sequence contained"""
        if filename.endswith(".gz"):
            fp = gzip.open(filename, 'r')
        else:
            fp = open(filename, 'r')
        # split at headers
        data = fp.read().split('>')
        fp.close()
        # ignore whatever appears before the 1st header
        data.pop(0)
        headers = []
        sequences = []
        for sequence in data:
            lines = sequence.split('\n')
            headers.append(lines.pop(0))
            # add an extra "+" to make string "1-referenced"
            sequences.append('+ ' + ''.join(lines))
        return (headers, sequences)
```

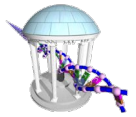


This is the same FASTA
format file parser from
last lecture

```
In [19]: header, seq = loadFasta("GCA_000001405.15_GRCh38_genomic.fna")
        print(len(header), "sequences")
        for i in range(len(header)):
            if header[i].startswith("CM") or header[i].startswith("J0"):
                print(header[i])
                print(len(seq[i])-1, "bases", seq[i][:30], "...", seq[i][-30:])
```

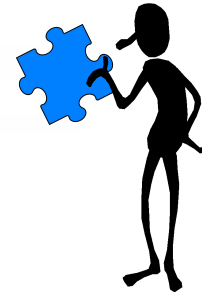
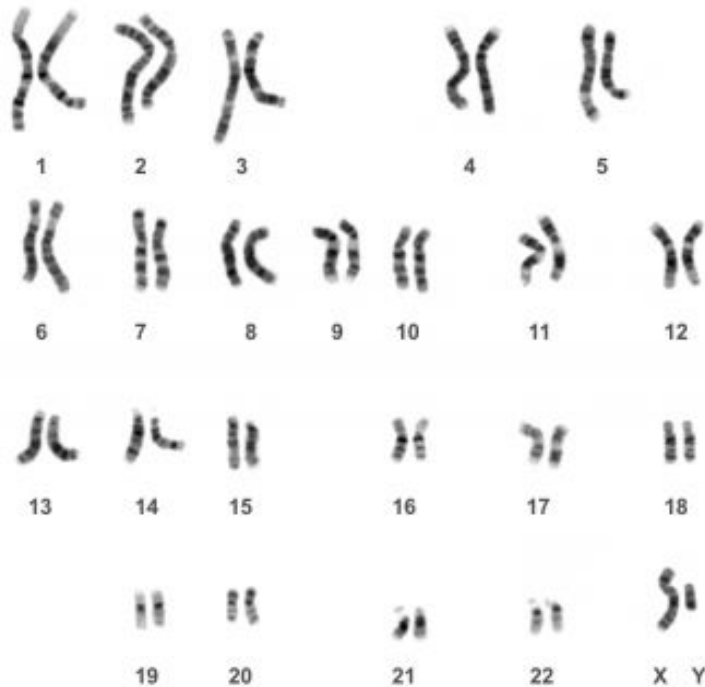


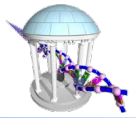
This is a recent version
of the human genome.
But, we're only going to
look at a subset of it.



Missing Puzzle Pieces

There are still missing, partially assembled, pieces, that we don't yet know where they are placed.





What its sequence looks like

As with SARS-CoV-2, we can get some insights into a genome by examining its k-mer distributions. But before we start, let's look consider the genome's size?

- In total, there are 3,272,116,950 base pairs in the primary (forward) sequence
- Many of these are unknown, and are indicated by 'N'
- There are also a small number of ambiguous bases indicated using a standard called UIPAC

Chromosome lengths

Total lengths Ungapped lengths N50s Gaps Counts

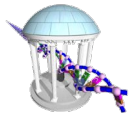
Chromosome lengths are calculated by summing the length of the placed scaffolds and estimated gaps.

Chromosome	Total length (bp)	GenBank accession	RefSeq accession
1	248,956,422	CM000663.2	NC_000001.11
2	242,193,529	CM000664.2	NC_000002.12
3	198,295,559	CM000665.2	NC_000003.12
4	190,214,555	CM000666.2	NC_000004.12
5	181,538,259	CM000667.2	
6	170,805,979	CM000668.2	
7	159,345,973	CM000669.2	
8	145,138,636	CM000670.2	
9	138,394,717	CM000671.2	
10	133,797,422	CM000672.2	
11	135,086,622	CM000673.2	
12	133,275,309	CM000674.2	
13	114,364,328	CM000675.2	
14	107,043,718	CM000676.2	
15	101,991,189	CM000677.2	
16	90,338,345	CM000678.2	
17	83,257,441	CM000679.2	
18	80,373,285	CM000680.2	
19	58,617,616	CM000681.2	
20	64,444,167	CM000682.2	
21	46,709,983	CM000683.2	
22	50,818,468	CM000684.2	
X	156,040,895	CM000685.2	NC_000023.11
Y	57,227,415	CM000686.2	NC_000024.10

IUPAC degenerate base symbols^[2]

Description	Symbol	Bases represented				Complementary bases ^[4]
		No.	A	C	G	
Adenine	A	1	A			T
Cytosine	C	1		C		G
Guanine	G	1			G	C
Thymine	T	1				A
Uracil	U	1				A
Weak	W	2	A			T
Strong	S	2		C	G	
Amino	M	2	A	C		
Keto	K	2			G	T
Purine	R	2	A	G		
Pyrimidine	Y	2		C		T
Not A ^[b]	B	3		C	G	T
Not C ^[b]	D	3	A		G	T
Not G ^[b]	H	3	A	C		T
Not T ^[b]	V	3	A	C	G	
Any one base	N	4	A	C	G	T
Zero	Z	0				

a. ^a I.e., here, read the represented bases in reverse.
b. ^{a b c d} Represented by the letter following (excluding U).



Let's reformat our sequences

It's a little annoying to load a series of sequences from FASTA files over and over again. Especially when we will mostly deal with a subset, and of those we will consider them one at a time.

So I decided to write out each sequence as a single string to its own file.

```
In [ ]: header, seq = loadFasta("data/GCA_000001405.15_GRCh38_genomic.fna")
print(len(header), "sequences")
for i in range(len(header)):
    if header[i].startswith("CM") or header[i].startswith("J0"):
        start = header[i].find('chromosome ')
        chromo = header[i][start+11:header[i].find(',')] if (start >= 0) else "MT"
        with open("data/Chr%s.seq" % chromo, 'w') as fp:
            fp.write(seq[i])
```

You might want to wait and do this later.

A quick helpful function



- DNA is actually two sequences, a primary and reverse-complement version
- Genomes report only one (the primary one), and the reverse complement version can be derived from it.
- When we consider k-mers, in most cases, we don't care which of the sequences they come from
- Here's a simple function that maps back and forth

```
In [ ]: def revComp(dnaSeq):  
        return ''.join(['A':'T','C':'G','G':'C','T':'A'][base] for base in reversed(dnaSeq))
```



We didn't consider the reverse complement sequence of our viral genome because it was an RNA genome.

- Here's an example:

```
In [4]: print(revComp("GAGACAT"))  
        print(revComp("ATGTCTC"))
```

```
ATGTCTC  
GAGACAT
```



Let's consider some k-mer statistics

For what value of k?

```
In [ ]: import time

chromo = [str(i) for i in xrange(1,23)] + ['X', 'Y', 'MT']

kmerCount = {}
K = 11
L = 0
for contig in chromo:
    tick = time.time()
    with open("Chr%s.seq" % contig, 'r') as fp:
        seq = fp.read()
        for i in xrange(1, len(seq)-K+1):
            kmer = seq[i:i+K]
            for base in "RYSWKMBDQVHVN":
                if (base in kmer):
                    break
            else:
                kmerCount[kmer] = kmerCount.get(kmer, 0) + 1
                kmer = revComp(kmer)
                kmerCount[kmer] = kmerCount.get(kmer, 0) + 1
    tock = time.time()
    print(contig, len(seq)-1, len(kmerCount), "%6.2f secs" % (tock - tick))
    tick = tock
    L += len(seq) - 1
print(L, len(kmerCount))
```

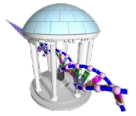
A "for-else" statement. Have you seen one of those before?



This is similar to our kmer counter from last lecture, except every k-mer is considered twice. Once as it appears, and once as its reverse complement.

We'll also skip over any k-mer with an 'N' or with one of those strange IUPAC bases.

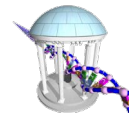
DON'T RUN IT!



```
'1', 248956422, 4133410, '1388.88 secs'  
'2', 242193529, 4175438, '1364.77 secs'  
'3', 198295559, 4184312, '1064.28 secs'  
'4', 190214555, 4188228, '1155.32 secs'  
'5', 181538259, 4190446, '1031.80 secs'  
'6', 170805979, 4191700, '930.77 secs'  
'7', 159345973, 4192490, '1025.24 secs'  
'8', 145138636, 4192908, '931.29 secs'  
'9', 138394717, 4193190, '751.69 secs'  
'10', 133797422, 4193464, '827.02 secs'  
'11', 135086622, 4193648, '780.05 secs'  
'12', 133275309, 4193788, '724.03 secs'  
'13', 114364328, 4193862, '635.01 secs'  
'14', 107043718, 4193926, '534.36 secs'  
'15', 101991189, 4193988, '530.65 secs'  
'16', 90338345, 4194048, '478.42 secs'  
'17', 83257441, 4194088, '447.81 secs'  
'18', 80373285, 4194110, '448.31 secs'  
'19', 58617616, 4194136, '345.26 secs'  
'20', 64444167, 4194158, '363.14 secs'  
'21', 46709983, 4194166, '251.95 secs'  
'22', 50818468, 4194182, '253.91 secs'  
'X', 156040895, 4194200, '866.71 secs'  
'Y', 57227415, 4194200, '189.44 secs'  
'MT', 16569, 4194200, '0.10 secs'  
3088286401, 4194200
```

- It takes a while to run (there are actually faster ways to do this!)
- 1000 secs is around 16 minutes
- And we still don't see every possible 11-mer

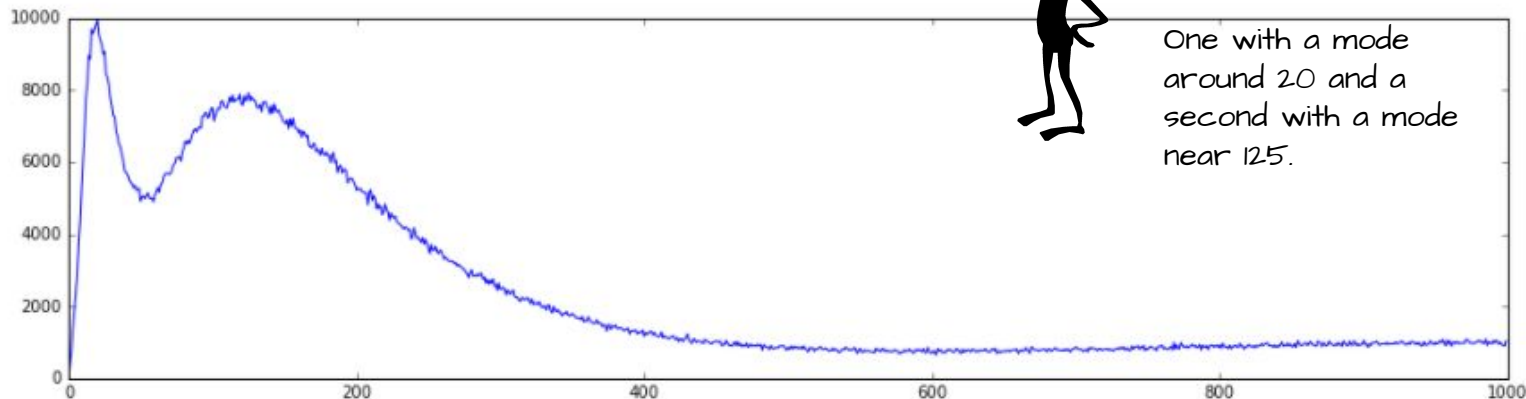
What does the distribution look like?



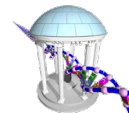
```
In [14]: import matplotlib
import matplotlib.pyplot as plot
%matplotlib inline

# Compute a histogram of kmerCount (i.e. how many kmers appear 1 time, 2 times, 3 times ...)
maxcount = 1000
hist = [0 for i in range(maxcount)]
for kmer in kmerCount:
    count = kmerCount[kmer]
    if (count < maxcount):
        hist[count] += 1

fig = plot.figure(figsize=(16,4))
plot.plot([i for i in range(maxcount)], hist)
plot.show()
```



What could we learn from other values of k



- Our genome includes every possible 11-mer
- How large should k be so that we'd expect most k -mers to be unique?
- Recall the genome has 3,272,116,950 bases

There are 4,194,304 11-mers

There are 67,108,864 13-mers

There are 1,073,741,824 15-mers

There are 17,179,869,184 17-mers

There are 274,877,906,944 19-mers

There are 4,398,046,511,104 21-mers

There are 70,368,744,177,664 23-mers

There are 1,125,899,906,842,624 25-mers

There are 18,014,398,509,481,984 27-mers

There are 288,230,376,151,711,744 29-mers

There are 4,611,686,018,427,387,904 31-mers

There are 73,786,976,294,838,206,464 33-mers

There are 1,180,591,620,717,411,303,424 35-mers

There are 18,889,465,931,478,580,854,784 37-mers

There are 302,231,454,903,657,293,676,544 39-mers

There are 4,835,703,278,458,516,698,824,704 41-mers

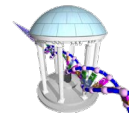
There are 77,371,252,455,336,267,181,195,264 43-mers

There are 1,237,940,039,285,380,274,899,124,224 45-mers



What's with all the odd numbers?

While I was bored last night...



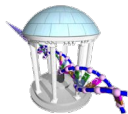
I broke the human genome into non-overlapping 45-mers, and counted how many time each appears in the genome...



TGGCCGAATAGGAACAGCTCCGGTCTACAGCTCCCAGCGTGAGCG | ACGCAGAAGACGGTGATTTCTGCATTTCCATCTGAGGTACCGGGT | TCATCTCACTAGGGAGTGCCAGACAGTGGGCGCAGGCCAGTGTGT | GTGCGCACCGTGCACGAGCCGAAGCAGGGCGAGGCATTGCCTCAC |



CTGGGAAGCGCAAGGGGTCAGGGAGTTCCTTTCCGAGTCAAAGA | AAGGGGTGACGGTGCACCTGGAAAATCGGGTCACTCCACCCGA | ATATTGCGCTTTTCAGACCGGCTTAAGAAACGGCGCACCCAGAGA | CTATATCCCACACCTGGCTCGGAGGGTCTACGCCACGGAATCT | . . .



Most places look like this...

Chromosome 4

```

79935885 1 TCCAGCTGTTGCATAGCTTTGTTAAAGAGTGACACTTAGGCTAAT
79935930 1 GTACTCTAAGGAAATGACTCCGCTCCAGTGGAAATCTCTTCTG
79935975 1 AAACAATAAATGCCTGTTCCAACAAAAGAGCACCCTAAACTATGA
79936020 1 TTCCATTCCAAGTGTAAAAAATGGAAATTAATAGGATTTAGC
79936065 1 AAGTGTACAACCTCTAGCCAGGAGTCAATATTTCTAATTTTGAGA
79936110 1 TATTATTCATAGTCTCAATGCAGAGAGTTACTACACATTATTTT
79936155 1 TACTATCAGTACAATACCACCTTTTTAAAAGGGTTCAGATGTTTTA
79936200 1 ATCACTATTACACAAGTACTACAAAATGATATAATTAGTTGCATTC
79936245 1 TTATTTGCAGAATATTTAATGTATCTCTATTACAGATAAATTTTTA
79936290 1 AATGACAAAATGCTATTTAACTGTCTTATTTTTCAGACCTCCCTGTC
79936335 1 ATCAGAGCTTAGTACTCTCTTTTCAAACCATTACTATTTCTCTG
79936380 1 CACAACCTAGCATAAACTGGTTGCACTGCATTTAGAAGCTCTGCT
79936425 1 TAGGTACTTCTATGCAAGTTTTCTCTTCTATCTCAACAACAAA
79936470 1 CCTAAATGAACAATCAAAACATTCTCACATTTTTCTAATTTCCCA
79936515 1 AAAGTTTTCTTTTTCTTTCAAGGTAATTTATGTCCTCATAAAAGCT
79936560 1 TAGGTATAATCTGTATGGGAGAAATAGTAAAATATTTCAACAG
79936605 1 TATTTCAAGTGTGCTCAGGGACTGGGAAATAATGCAAAAAGAAATA
79936650 1 AAAAAAATCCCTTATTAACAGTAAAAAGGGAAGAAAGAAAGACTAA
79936695 1 CACATGAAATAATTAGAAAACAATTAATTTAAAAAATTAAGTGCACA
79936740 1 ATATTTAACTAATTTCACTTTAAGTCTGATTATTAATTTGCACTA
79936785 1 TGGGTTATAATCACTTTTTTTTGTGCTTATAATCACTCCACCTAG
79936830 1 TACAATGGTACTACACTAGTAGTACACACTAAAGGTACTACACT
79936875 1 CTACCAGTACTAACCAATGGTAAAACTACACTTATATTTTCT
79936920 1 ATTTTTTTCTACTCTGTATAATGTAGGCAGAAAAGTCAAGCAAG
79936965 1 CAATGGATAAATGATTATGTATTCTTCAATCATTTTAAAGCATT
79937010 1 TTTTCAATTTAATCTTTGTTGCAAAAAGAAAAATGATTAATAAT
79937055 1 TTTTACTTTTAAATAACATTAAGCATTAACTTTATACAGTTTTTA
79937100 1 AAACACTCATAAACTTAATTAAGCTTTTAAATTCAGTTATAAATAG
79937145 1 GACTCATTCACTGTATTTCTCAACAGTAGCATTAAAAAACCAGGT
79937190 1 GCCTATTTTCATATTTCTTAATGAAGCAATGTGCTAGCAATAGGAAA
79937235 1 CCTCAAAAGATTCACATTTGGCTCAACTAAGTTCTCTGAAAATTA
79937280 1 ATTACATATAATCATTTAAAAACAGCACAATAATGAGCAGAAAGAAA
79937325 1 AAAAAATTTTGAAGAATGTTGTAATATCCATAAATGTTTAGGC
79937370 1 TAGTTTTGGCTGGTTTCTGATTAACTGCATTTTGGACATATCTTCAT
79937415 1 TGAAGATTTCACTGTAACTACTCACAGAAAGCTTTTTATCTGCA
79937460 1 AGTGACTTTTTGTGCCACTTGTCTTGGGCCACTTTTTCCAACCTCT
79937505 1 AATTTGCAATTTGTATCTACCTCGAGAGAGTACTGTCTATCAAG
79937550 1 GTATATAGTACCATACTCAAAACAGATTTGTTCCGTTATCAAACT
79937595 1 AGAAAAATAAATAATCATAAAATGTTATGTGTCACTAACCAAGGTA
79937640 1 CAACCTGAATGCTTATGTATATATTGAGCATCAATATGTACCCA

```

As you'd expect, most 45-mers are unique.

But, occasionally, we run into a series that are repeated all over the genome.

And, they aren't trivial repetitive sequences.

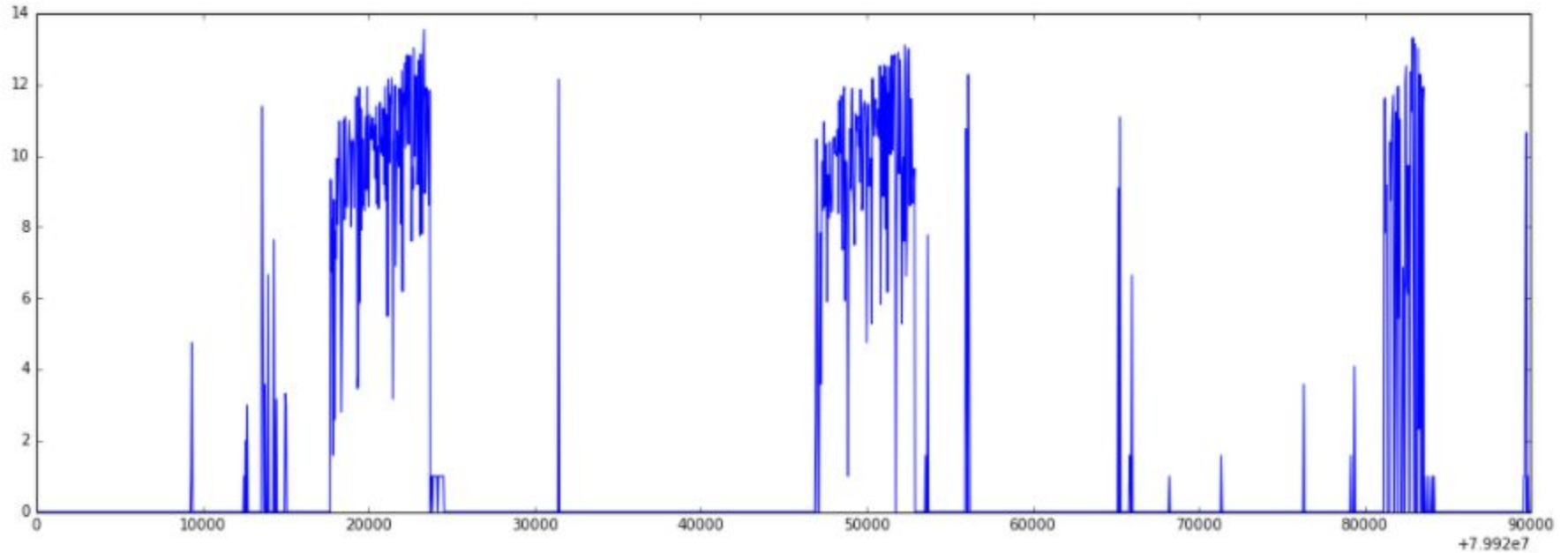
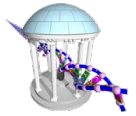
Chromosome 4

```

79937415 1 TGAAAGTTTCACTGTAACATACTCACAGAAAGCTTTTTATCTGCA
79937460 1 AGTGACTTTTTGTGCCACTTGTCTGGGGCACTTTTTCCCAACTCT
79937505 1 AATTTGCAATTTGTATCTACCTGAGAGAGGACTGTCTATCAGG
79937550 1 GTATATAGTACCATACTCAACAGATTTGTTCCGTTATCTAACT
79937595 1 AGAAAATAAATAATCATAAAATGTTATGTGTCACTAAACAAGGTAA
79937640 1 CAACCTGAATGCTTATGTATATATTGAGCATCAATATGTACCCA
79937685 1 GCCTGTGATAGTGTTTTTAAAACCCCTAAGAGAGGAGCCAAGA
79937730 645 TGCGCGAATAGGAAACAGCTCCGGCTCACAGCTCCCAAGCTGGAGG
79937775 108 ACGCAGAAAGACGGTGATTTCTGCATTTCCATCTGAGGTCACGGGT
79937820 295 TCATCTCACTAGGAGTGCCAGACAGTGGGCGCAGGCAAGTGTGT
79937865 3 GTGCGCACCTGTGCAGGAGCGAAGCAGGGCAGGCAATGCTCTCAC
79937910 436 CTGGGAAGCGCAAGGGGTGAGGAGTTCCTTTCCGAGTCAAGA
79937955 6 AAGGGGTGACGGTCCGACTGGAAAATCGGGTCACTCCCAACCGCA
79938000 233 ATATTTGCGCTTTTTCAGACGGCTTAAGAAACGGCCACCACGAGA
79938045 138 CTATATCCCAACACTGGCTGGAGGGTCTTACGCCCAAGGAACTTC
79938090 973 CGCTGAATGCTAGCACAGCACTGTAGATCAAACTGCAAGGCGCC
79938135 270 AACGAGGCTGGGGAGGGGCGCCGCCATTGCCAGGCTTGCTTA
79938180 546 GGTAAACAAGACAGCCGGGAAGCTGCAACTGGGTGGAGCCCAACA
79938225 2005 CAGCTCAAGGAGGCTGCTGCTCTGTAGGCTCACTCTGGGG
79938270 2010 GCAGGGCACAGACAACAAAAAGACAGCAGTAACCTCTGCAGACT
79938315 405 TAAGTGTCCCTGTCTGACAGCTTTGAAGAGAGCAGTGGTCTCC
79938360 7 AGCAGCGAGCTGGAGATCTGAGAACGGGCGAGCAGACTGCCTCT
79938405 73 CAAGTGGTCCCTGACTCTGACCCCGAGCAGCTTAAGTGGAG
79938450 590 GCACCCCGCAGCAGGGGCACACTGACACTCACAGGCGAGGAT
79938495 2051 TCCAACAGACTGCACTGAGGGTCTGTCTGTAGAGGAAAC
79938540 295 TAACAACAGAAAAGCAGACTCAACCGAAAACCATCTGTACATC
79938585 2174 ACCATCATCAAGACCAAAAGTAGATAAAACCAAAAGATGGGA
79938630 501 AAAAAACAGAACAGAAAACCTGAAACTCTAAAACCGCAGGCGCT
79938675 377 CTCTCCTCAAAGGACGCAAGTTCTCACCAGCAACAGAAACAAA
79938720 382 GCTGGATGGAGAATGATTTTGTGAGGCTGAGAGAAGAAAGGCTCA
79938765 955 GACGATCAAAATTAAGTCTGAGCTACGGGAGGACATCAAAACCAAG
79938810 1453 GCAAGGAAATGAAAACCTTGAAGAAAATTTAGAAGAAATGATAA
79938855 2046 CTAGAATAACCAATACAGAGAAGTCTTAAAGGAGCTGTAGGAG
79938900 1271 TGAAAACCAAGGCTCGAGAACTACGTGAAGAATGCAAGAGCCTCA
79938945 258 GGAGCGATGCGATCAACTGGAAAGAAAGGATACAGCGATGGAAG
79938990 1313 ATGAAATGAATGAAATGAAGCGAGAAAGGAAAGTTTGAAGAAA
79939035 1398 GAATAAAAAGAAATGAGCAAAGCTCCAAGAAATATGGGACTATG
79939080 1315 TGAAAAGACAAAATCACTGCTGATTTGTTGTGTGCAAGGATGATG
79939125 371 TGAAGAAATGGAACCAAGTTGGAAAAACACTCTGACGGATATTTCC
79939170 1155 AGGAGAATTTCCCAACTAGCAAGGCAAGGCAAGCTTCAAGATTC

```


Zooming out





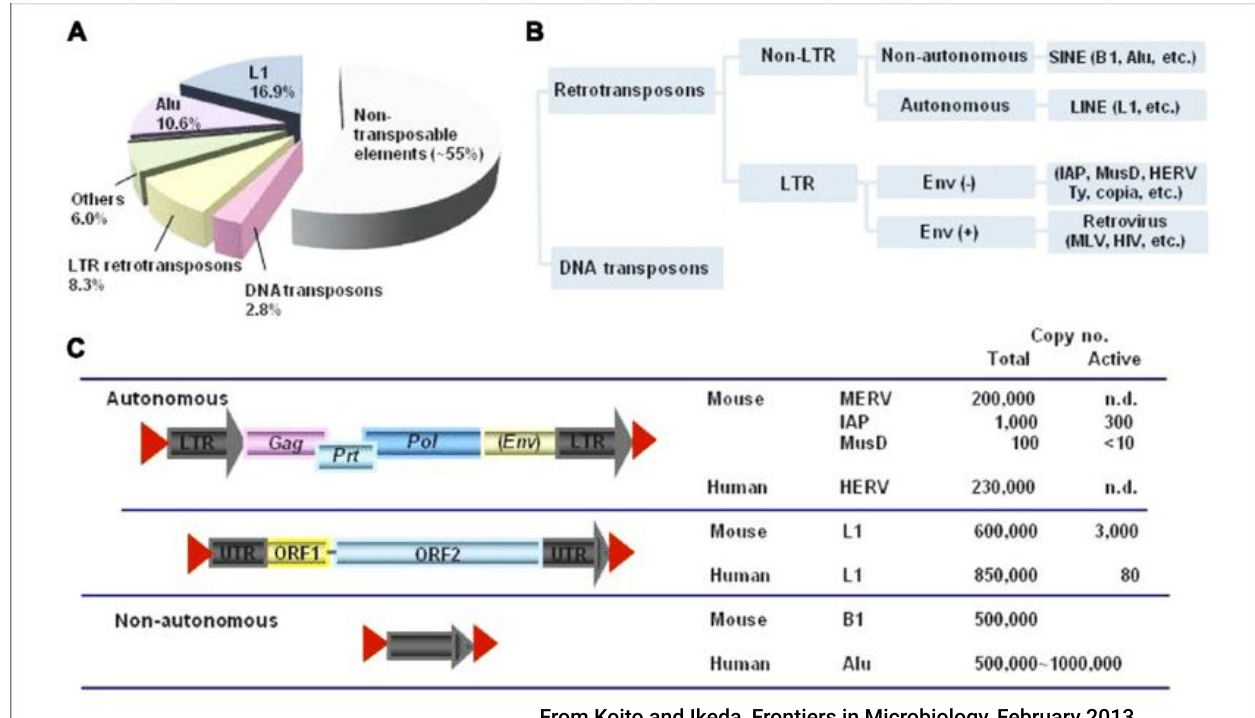
Repeated regions of our genome

Our genome is full of copies... most are due to Transposable Elements

About 45%

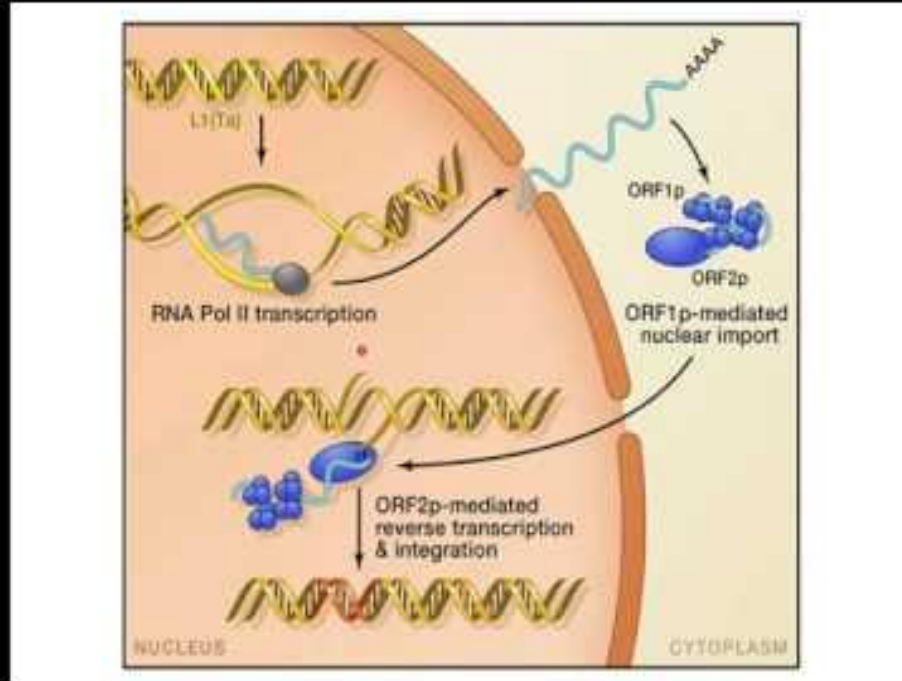
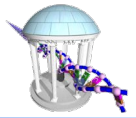
Cut-and-paste
DNA transposons

Copy-and-paste
Retrotransposons



From Koito and Ikeda, *Frontiers in Microbiology*, February 2013

TEs are everywhere



Let's find all copies of one of our repeats



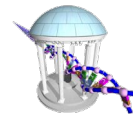
```
In [5]: chromo = [str(i) for i in xrange(1,23)] + ['X', 'Y', 'MT']

target = "AGCACGCAGCTGGAGATCTGAGAACGGGCAGACAGACTGCCTCCT" # was 7 times
revtar = revComp(target)

for contig in chromo:
    with open("Chr%s.seq" % contig, 'r') as fp:
        seq = fp.read()
        start = 0
        while True:
            i = seq.find(target, start)
            if (i > 0):
                print(contig, i, "+")
                start = i + 1
            else:
                break
        start = 0
        while True:
            i = seq.find(revtar, start)
            if (i > 0):
                print(contig, i, "-")
                start = i + 1
            else:
                break
```

```
2 169249269 +
4 79938360 +
5 156067276 -
8 72880901 -
11 93137285 +
11 93421619 +
16 33957922 -
```

And look around where we found them



Here's one of the copies of
"AGCACGCAGCTGGAGATCTGAGAACGGGCAGACAGACTGCCTCCT"

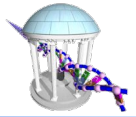
Not surprisingly, it is surrounded by repeated 45-mers that
are similar to the ones on chromosome 4.

So, with what we now know let's go hunting for a particular type
of Transposable Element.

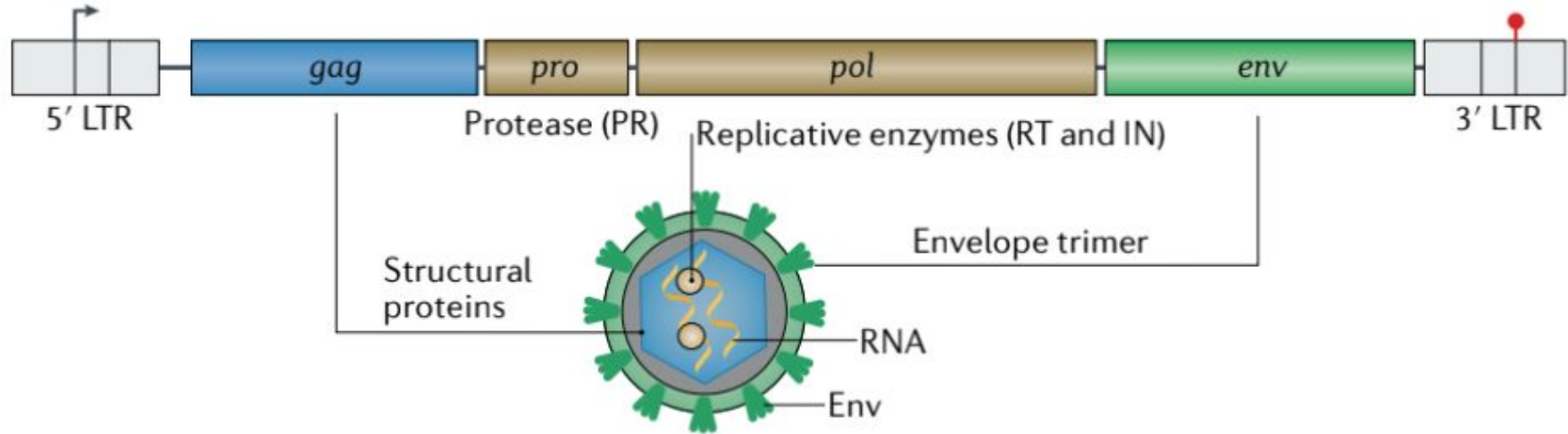
One of Viral origin! Yes, our genomes have parasites.

```
Chromosome 11
93136365 1 TTGCAGAGAAGTAGGAATGCTTTTACACTGTCGGTGGGAATGTAA
93136410 1 ATTAGGTCAACTATTGTGGAAGACAGTGTGCAATTCCTCAAAGAT
93136455 5 CTAGAACCCAGAAATACGATTTTGACCCAGCAATCCCATTTACTGGGT
93136500 1 AAATACCCAAAAGAATATAAATCATTCTATTATAGAGATACATGC
93136545 1 ATGAGTATGTTTCAAGTGCAGCACACTTCAAAATAGCAAAAGACATGG
93136590 1 AATCAACACAACAGTCCCATCAATGATAGACTAAAAGAAAACGTGGT
93136635 1 ACATGGGGGAGGAGCCAAAGATGGCCGAATAGGAACAGCTCCGGTC
93136680 93 TACAGCTCCCAGCGTGTGAGCGACGCAGAAAGACGGTGTATTCGCAT
93136725 1334 TTCCATCTGAGGTACCGGTTCACTCAGGAGTGGCCAGACA
93136770 169 GTGGGCGCAGGCCAGTGTGTGCGCACCGTGTGCGGAGCCGAAGC
93136815 449 AGGGCGAGGCATTGCTCACCCTGGGAAGCGCAAGGGGTGAGGAG
93136860 3 TTCCCTTTCTGAGTCAAAGAAAAGGGGTGACGGTGCACCTGGAAA
93136905 437 ATCGGGTCACTCCACCCGAATATTGCGCTTTTTCAGACCCGGCTTA
93136950 55 AGAAAACGGCGCACCCAGAGACTATATCCACACCTGGCTCGGAGG
93136995 393 GTCCTACGCCACCGAACTCTCGCTGATTGCTAGCACAGCAGTCTG
93137040 270 AGATCAAAGTGAAGGGGCAACGAGGCTGGGGGAGGGGCGCCCG
93137085 581 CCATTCGCCAGGCTTGTAGTAAACAAAGCAGCCGGGAAGCTC
93137130 1565 GAAGTGGGTGGAGCCACCCAGCTCAAGGAGGCTGCTGCCTC
93137175 1908 TGTAGGCTCCACCTCTGGGGGAGGGCACAGACAAAACAAAAGAC
93137220 406 AGCAGTAACCTCTGAGACTTAAGTGTCCCTGTCTGACAGCTTTG
93137265 1104 AAGAGAGCAGTGGTTCTCCAGCACGCTGGAGATCTGAGAAC
93137310 8 GGGCAGACAGACTGCCTCCTCAAGTGGTCCCTGACTCTGACCC
93137355 817 CCGAGCAGCCTAACTGGGAGGCCACCCCCAGCAGGGGACACTGA
93137400 81 CACCTCACATGGCAGGGTATTCCAACAGACCTGACGCTGAGGGTC
93137445 329 CTGTCTGTTAGAAGGAAAATAACAACAGAAAGGACATCTACAC
93137490 526 CGAAAACCCATCTGTACATCACCATCATCAAAGACCAAAAGTAGA
93137535 1308 TAAAACCACAAGATGGGGAAAAACAGAACAGAAAACCTGGAAA
93137580 446 CTCTAAAACGCAGAGCGCTCTCCTCCTCAAAGGAAACGCAGTTC
93137625 246 CTCACCAAGCAACAGAACAAAGCTGGATGGAGAATGATTTTGACGA
93137670 886 GCTGAGAGAAGAAGGCTTACAGACGATCAAACTACTCTGAGCTACG
93137715 1553 GGAGGACATTCAAAACCAAAGGCAAGAAAGTTGAAAACCTTTGAAAA
93137760 1512 AAATTTAGAAAGATGTATAACTAGAAATAACCAATACAGAGAAAGTG
93137805 1272 CTTAAAGGAGCTGATGGAGCTGAAAACCAAGGCTCGAGAAGTACG
93137850 1265 TGAAGAATGCAGAAAGCCTCAGGAGCCGATGCGATCAACTGGGAAGA
93137895 775 AAGGGTATCAGCAATGGAAGATGAAATGAATGAAATGAAGCGAGA
93137940 1267 AGGGAAAGTTTAGAGAAAAAGAAATAAAAAGAAATGAGCAAAAGCCT
93137985 115 CCAAGAAAATATGGGACTATGTGAAAAGACCAAACTACTAGCTTGAC
93138030 367 TGGTGTACCTGAAAGTGTGTGGAGAATGGAACCAAGTTGGAAAA
93138075 2615 CACTCTGAGGATATTATCCAGGAGAACTCCCAATCTAGCAAG
93138120 955 GCAGGCCAACGTTTCAGATTTCAGGAAAATACAGAGAAGCCACAAAG
```

An Endogenous Retrovirus

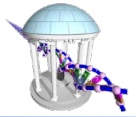


These TEs include many genes related to those in viral sequences.



The LTRs are identical sequences that enclose the a virus-like genome.

Next Time



We will develop a strategy to find LTR-like sequences in a genome.

