

# Counting kmers for fun and profit

Will Trimble  
Argonne National Laboratory

9 February 2016



# Counting is good for you.

- Counting kmers – how
- Counting – what you can see
- Visualizations and case studies
- Error correction options

“Why is kmer error correction better than quality-score error correction?”



# Definition

kmers are *overlapping, fixed-length oligonucleotide sequences*

AGCTTTCAATTCTGACTGCAACGGGCAATATGTCTCTGTGGATTAAAAAAAGAGT

AGCTTTCAATTCTGACTGCAA  
GCTTTCAATTCTGACTGCAAC  
CTTTCAATTCTGACTGCAACG  
TTTCATTCTGACTGCAACGG  
TTTCATTCTGACTGCAACGGG  
TTCATTCTGACTGCAACGGGC  
TCATTCTGACTGCAACGGGCA  
CATTCTGACTGCAACGGGCAA  
ATTCTGACTGCAACGGGCAAT  
TTCTGACTGCAACGGGCAATA  
TCTGACTGCAACGGGCAATAT  
.....



# Definition

|                  |     |
|------------------|-----|
| ATCGCGAAAAGTCCC  | 2   |
| AAAAAAAAAAAAAA   | 459 |
| AAAAAAAAAAAAAC   | 71  |
| AAAATAAAAAAAATA  | 1   |
| AAAAAAAAAAAAAG   | 36  |
| ACATGAAAAACAACT  | 1   |
| AAAAAAAAAAAAAT   | 23  |
| AAAAAAAAAAAAACA  | 95  |
| GTAGGAAAAGCCCCAC | 1   |
| AAAAAAAAAAAAACC  | 7   |
| AAAAAAAAAAAAACG  | 8   |
| AAAAAAAAAAAAACT  | 9   |
| AAAAAAAAAAAAAGA  | 36  |
| AACAAGAAAAACAAA  | 1   |
| AAAAAAAAAAAAAGC  | 10  |
| AAATAAAAAAAATAG  | 1   |
| AACAGAAAAACACG   | 1   |
| AAAAAAAAAAAAAGG  | 2   |
| AAAAAAAAAAAAAGT  | 6   |

kmers are *overlapping, fixed-length oligonucleotide sequences*

Just a tool—but a super useful one.

Can count for any sequence dataset

Makes sense for **genomes / samples / haplotypes / genes** if you want

Programming / plotting puzzle

<http://angus.readthedocs.org/en/2016/automation.html>

# Short kmers

| C         | 1179554 | AA         | 337870 | AAAAAAA       | 711 | AAAAAAAAAA    | 0  |
|-----------|---------|------------|--------|---------------|-----|---------------|----|
| G         | 1176923 | AC         | 256662 | AAAAAAC       | 722 | AAAAAAAAC     | 1  |
| A         | 1142228 | AG         | 237877 | AAAAAAG       | 815 | AAAAAAAAG     | 2  |
| T         | 1140970 | AT         | 309819 | AAAAAAT       | 941 | AAAAAAAAT     | 4  |
|           |         | CA         | 325149 | AAAAACA       | 791 | AAAAAAAACA    | 12 |
|           |         | CC         | 271673 | AAAAACC       | 629 | AAAAAAAACC    | 11 |
|           |         | CG         | 346670 | AAAAACG       | 732 | AAAAAAAACG    | 4  |
|           |         | CT         | 236061 | AAAAACT       | 607 | AAAAAAAACT    | 5  |
|           |         | GA         | 267247 | AAAAAGA       | 665 | AAAAAAAAGA    | 10 |
|           |         | GC         | 383931 | AAAAAGC       | 819 | AAAAAAAAGC    | 23 |
|           |         | GG         | 270137 | AAAAAGG       | 522 | AAAAAAAAGG    | 4  |
|           |         | GT         | 255608 | AAAAAGT       | 481 | AAAAAAAAGT    | 4  |
|           |         | TA         | 211961 | AAAAATA       | 805 | AAAAAAAATA    | 13 |
|           |         | TC         | 267288 | AAAAATC       | 770 | AAAAAAAATC    | 12 |
|           |         | TG         | 322239 | AAAAATG       | 799 | AAAAAAAATG    | 11 |
|           |         | TT         | 339482 | AAAAATT       | 665 | AAAAAAAATT    | 7  |
|           |         |            |        | AAAACAA       | 522 | AAAAAAAACAA   | 13 |
|           |         |            |        | AAAACAC       | 344 | AAAAAAAACAC   | 7  |
|           |         |            |        | AAAACAG       | 657 | AAAAAAAACAG   | 18 |
| k=1       |         | k=2        |        | k=7           |     | k=10          |    |
| $4^1 = 4$ |         | $4^2 = 16$ |        | $4^7 = 16384$ |     | $4^{10} = 1M$ |    |



# Long kmers

AAAAAAAAAAAAA  
AAAAAAAAAAAAA  
AAAAAAAAAAAAAC  
AAAAAAAAAAAAAG  
  
...  
AAAAAAAAAACCTCTTTTTT  
AAAAAAAAACCTGAAAAAAA  
AAAAAAAAACCTGAAAAAAC  
  
...  
AAAAAAAAAGAAAGGTAACG  
AAAAAAAAAGAAAGGTAACT  
AAAAAAAAAGAAAGGTAAGA  
  
...  
AAAAAAAAAGCCAGCACCCC  
AAAAAAAAAGCCAGCACCCG  
AAAAAAAAAGCCAGCACCCCT

$$k=20$$

$10^{11}$   
pigeons



$10^{12}$   
holes

## Short kmers

Tables of all  $4^k$  symbols can be computed by novice programmers.

Digests are compact

Individual short kmers are not interesting; only aggregate patterns matter.

$$1 \leq k < 16$$

## Long kmers

Storing  $4^k > 10^{12}$  doesn't work very well, so digital representations are either sparse or probabilistic.

Digests are bulky

Because of sparseness, **individual kmers can be interesting.** (kmers map.)

$$k > 20$$

# Genome kmers vs. data kmers

When sequencing a genome, some of your reads will perfectly represent the genome, and some will contain errors.

But if you sequence a genome to 100x, you will see the **same genome sequence over and over again**, but most errors you will only see once.

# kmer table -> kmer spectrum.

| kmer              | number of occurrences | number of occurrences | number of kmers |
|-------------------|-----------------------|-----------------------|-----------------|
| ATCGCGAAAAGTCCC   | 2                     |                       |                 |
| AAAAAAAAAAAAAA    | 459                   |                       |                 |
| AAAAAAAAAAAAAAAC  | 71                    | 1                     | 79317651        |
| AAAATAAAAAAAATA   | 1                     | 2                     | 3726721         |
| AAAAAAAAAAAAAAAG  | 36                    | 3                     | 952261          |
| ACATGAAAAACAAC    | 1                     | 4                     | 365963          |
| AAAAAAAAAAAAAAAT  | 23                    | 5                     | 173303          |
| AAAAAAAAAAAAAAACA | 95                    | 6                     | 92876           |
| GTAGGAAAAGCCCAC   | 1                     | ...                   |                 |
| AAAAAAAAAAAAAAACC | 7                     | 100                   | 1276            |
| AAAAAAAAAAAAAAACG | 8                     | 101                   | 1166            |
| AAAAAAAAAAAAAAACT | 9                     | ...                   |                 |
| AAAAAAAAAAAAAAAGA | 36                    | 3268                  | 1               |
| AACAAGAAAAACAAA   | 1                     | 3370                  | 1               |
| AAAAAAAAAAAAAAAGC | 10                    |                       |                 |
| AAATAAAAAAAATAG   | 1                     |                       |                 |

kmer table -> kmer spectrum.

| kmer              | number of occurrences | number of occurrences | number of kmers |
|-------------------|-----------------------|-----------------------|-----------------|
| ATCGCGAAAAGTCCC   | 2                     |                       |                 |
| AAAAAAAAAAAAAA    | 459                   |                       | 79317651        |
| AAAAAAAAAAAAAAAC  | 71                    | 1                     | 3726721         |
| AAAATAAAAAAAATA   | 1                     | 2                     | 952261          |
| AAAAAAAAAAAAAAAG  | 36                    | 3                     | 365963          |
| ACATGAAAAACAACT   | 1                     | 4                     | 173303          |
| AAAAAAAAAAAAAAAT  | 23                    | 5                     | 92876           |
| AAAAAAAAAAAAAAACA | 95                    | 6                     |                 |
| GTAGGAAAAGCCCCAC  | 1                     | ...                   |                 |
| A                 | 6                     | 100                   | 1276            |
| A                 |                       | 101                   | 1166            |
| A                 |                       | ...                   |                 |
| A                 |                       | 3268                  | 1               |
| A                 |                       | 3370                  | 1               |
| A                 |                       |                       |                 |
| AAA'TAAAAAAAT'AG  | 1                     |                       |                 |

Yes, it's the histogram of the histogram.

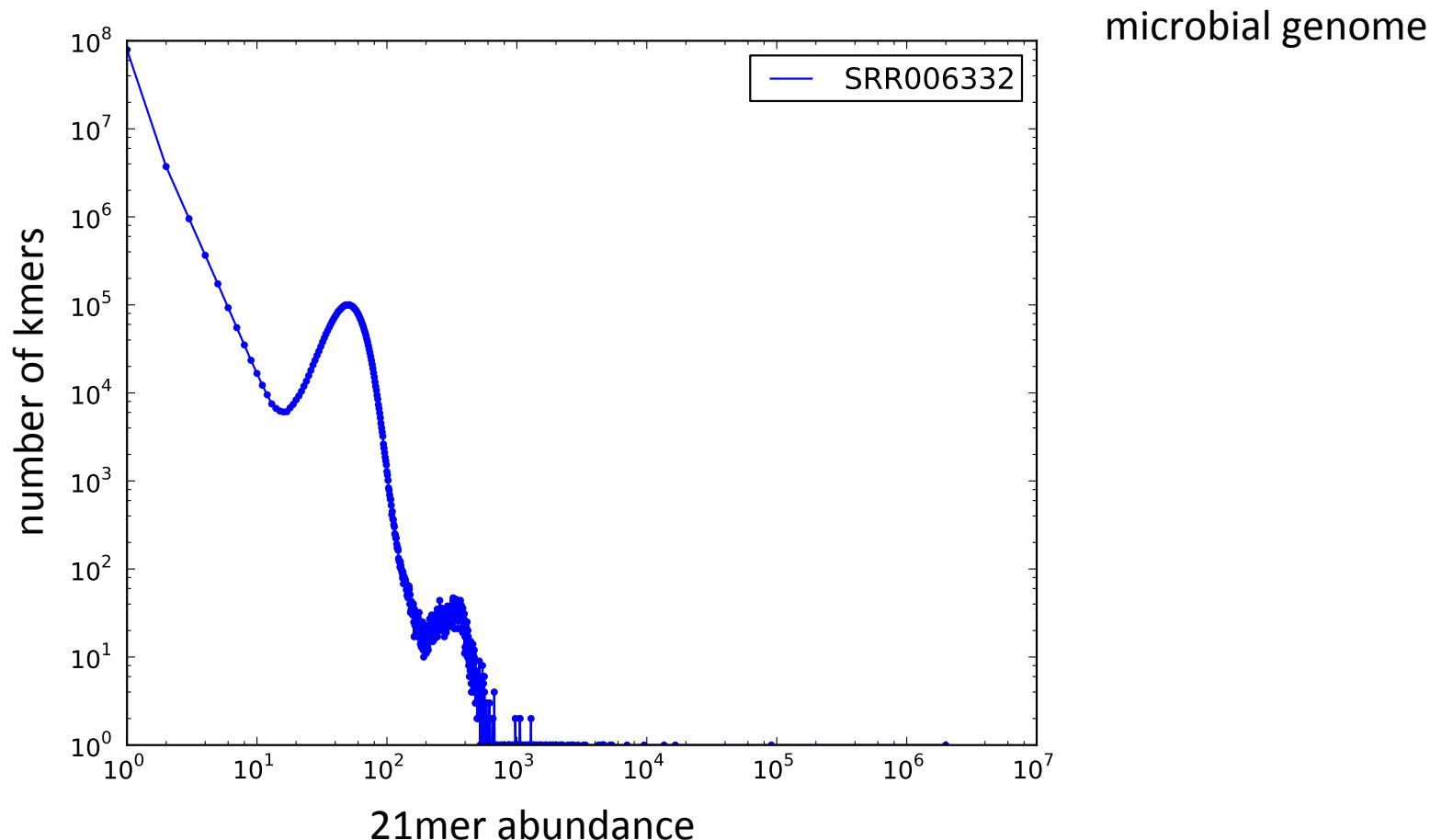
And it tells us something about pigeon family size – do the same kmers keep showing up over and over or not?



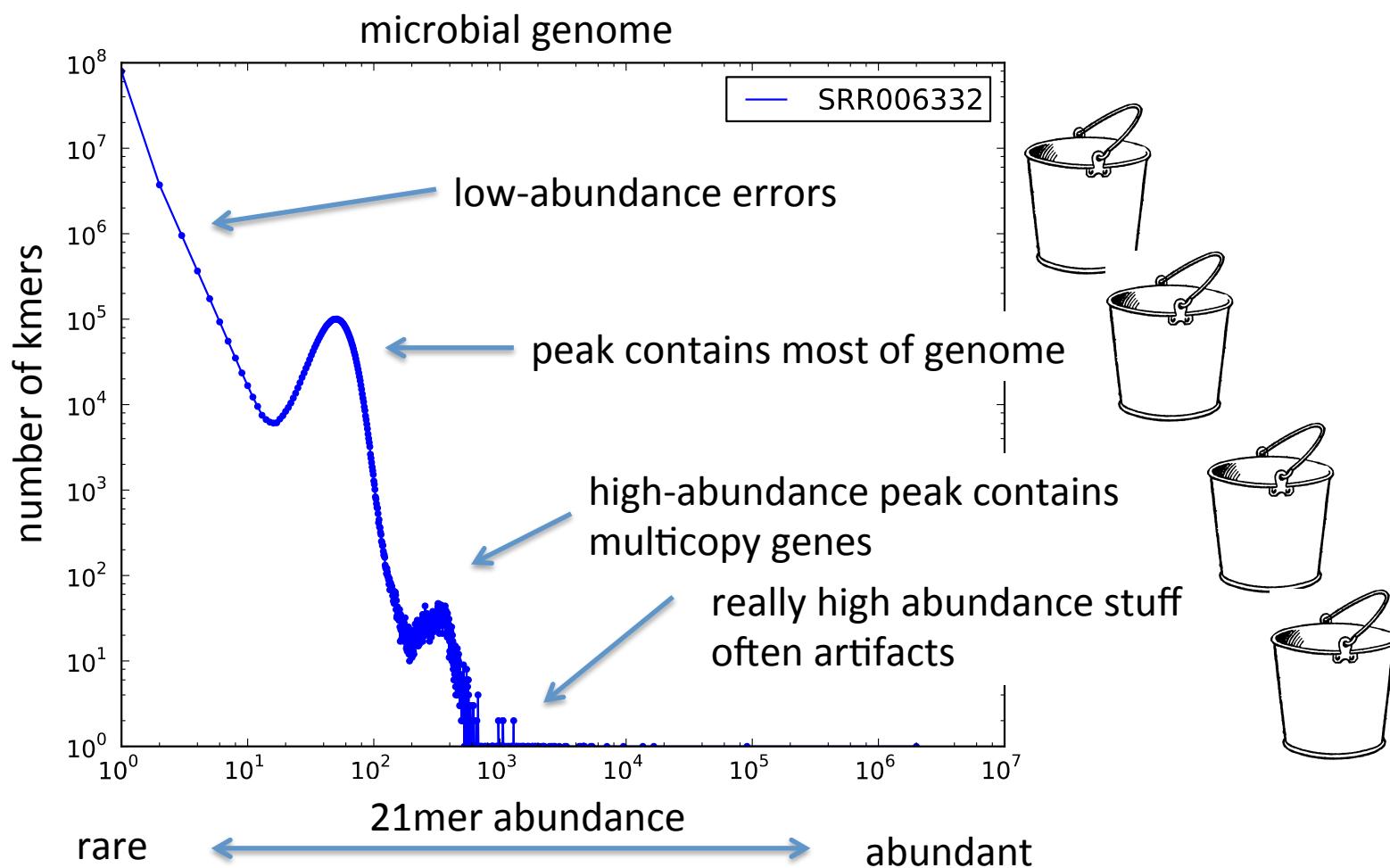
5 million holes  
>50 pigeons

100 million holes  
1 pigeon each

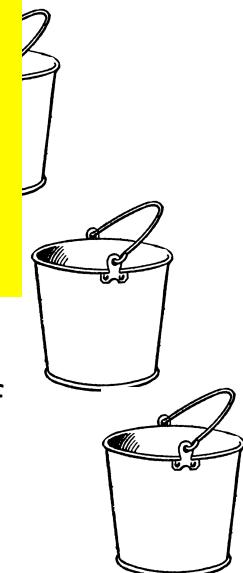
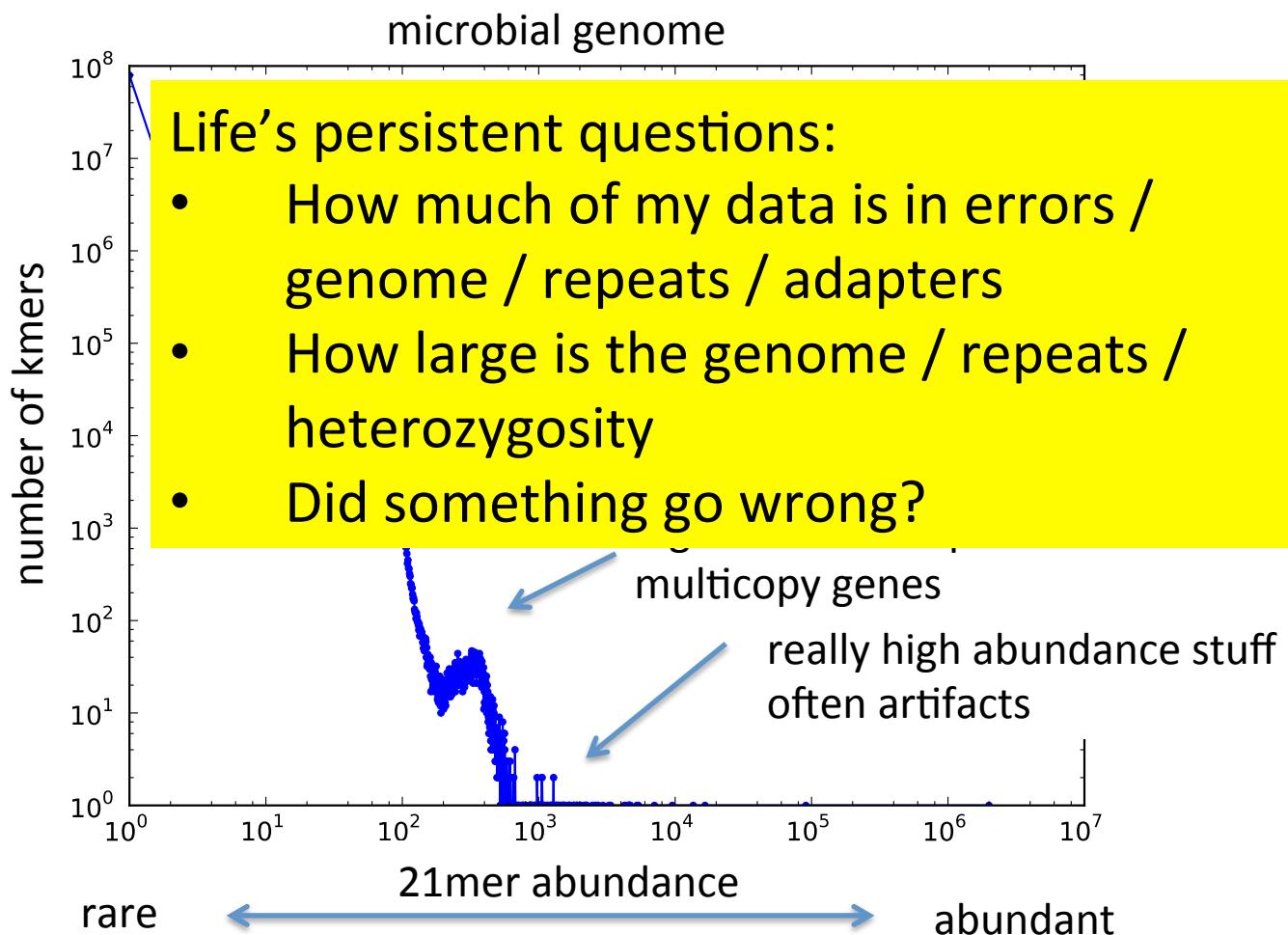
# The kmer spectrum.



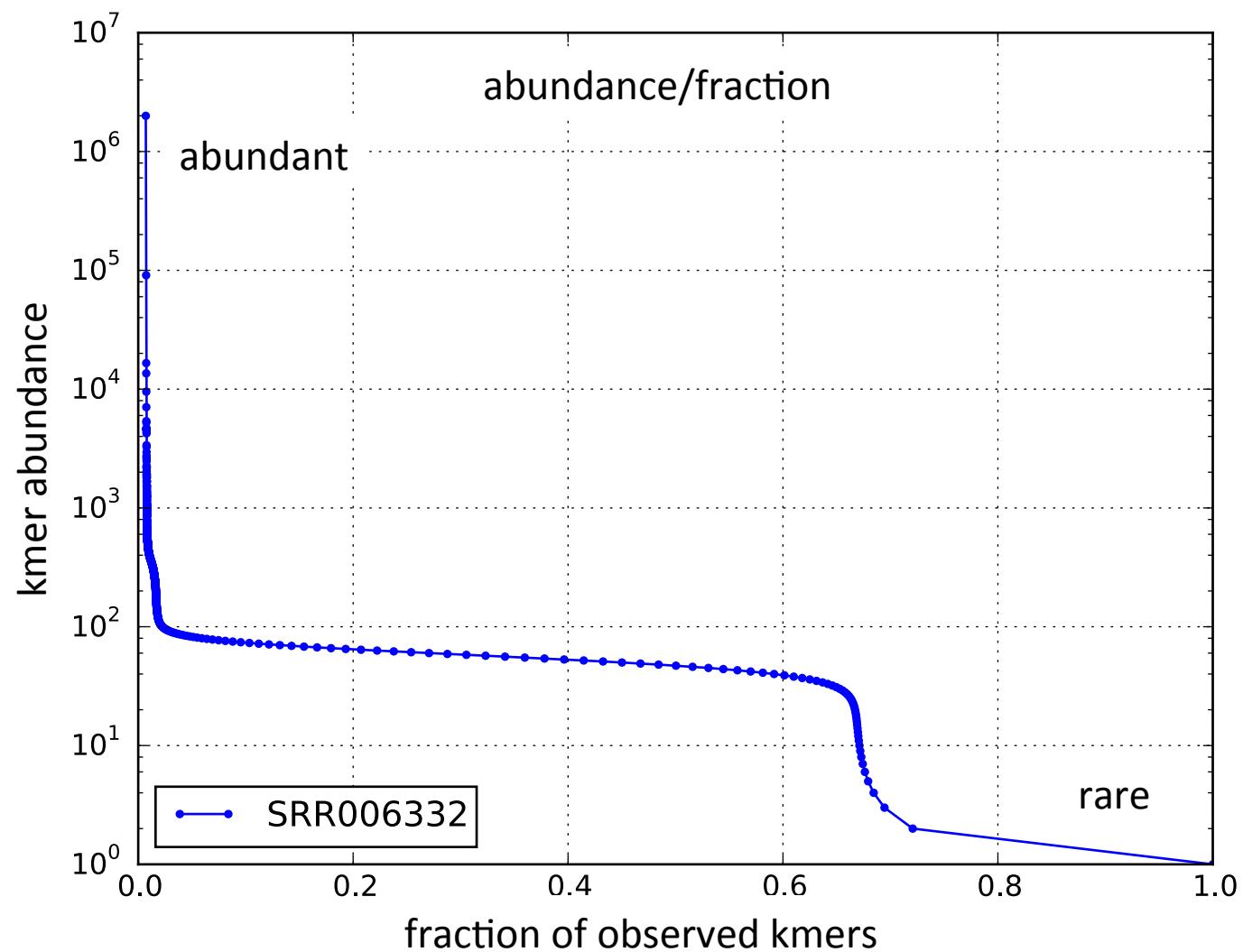
# The kmer spectrum.



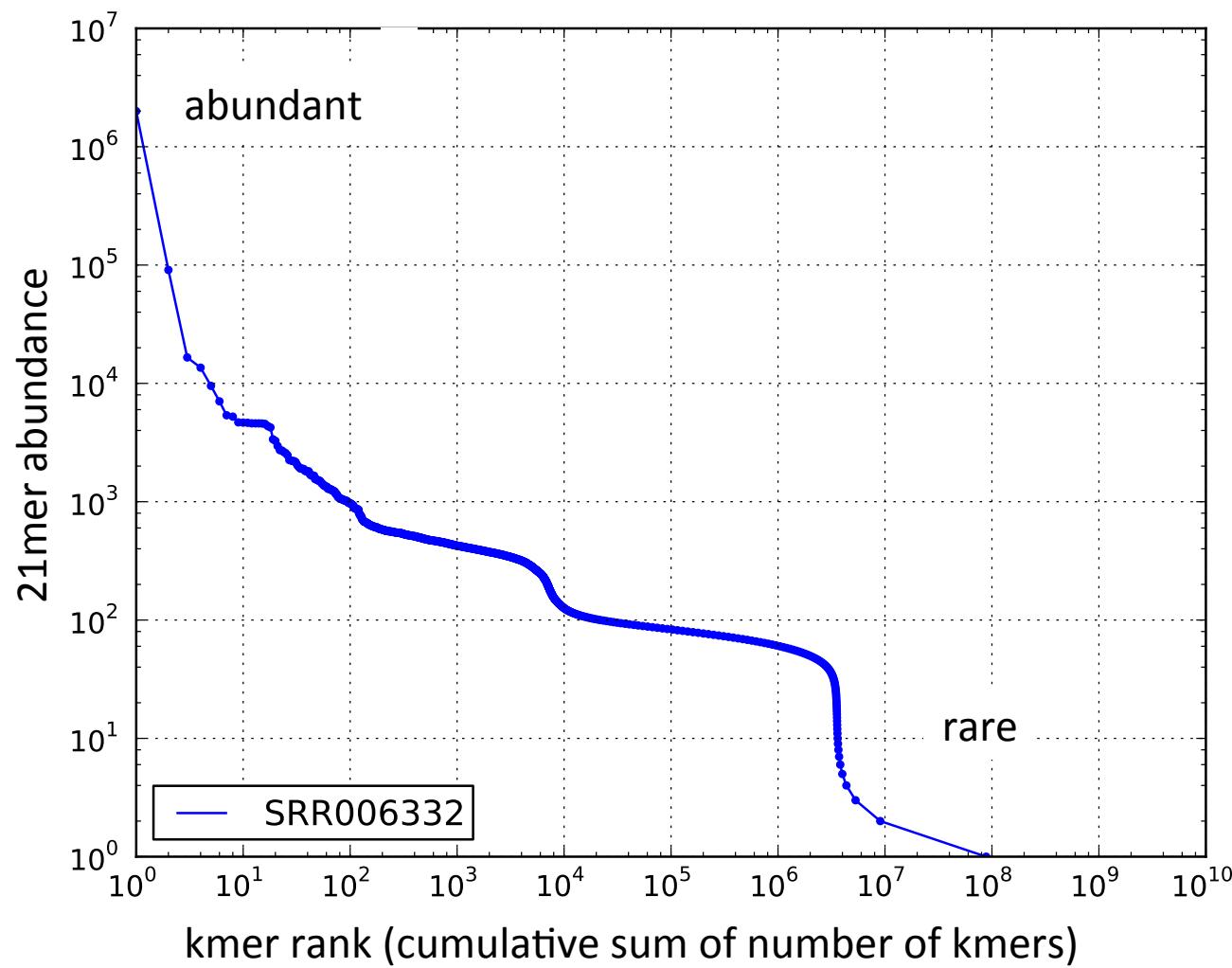
# The kmer spectrum.



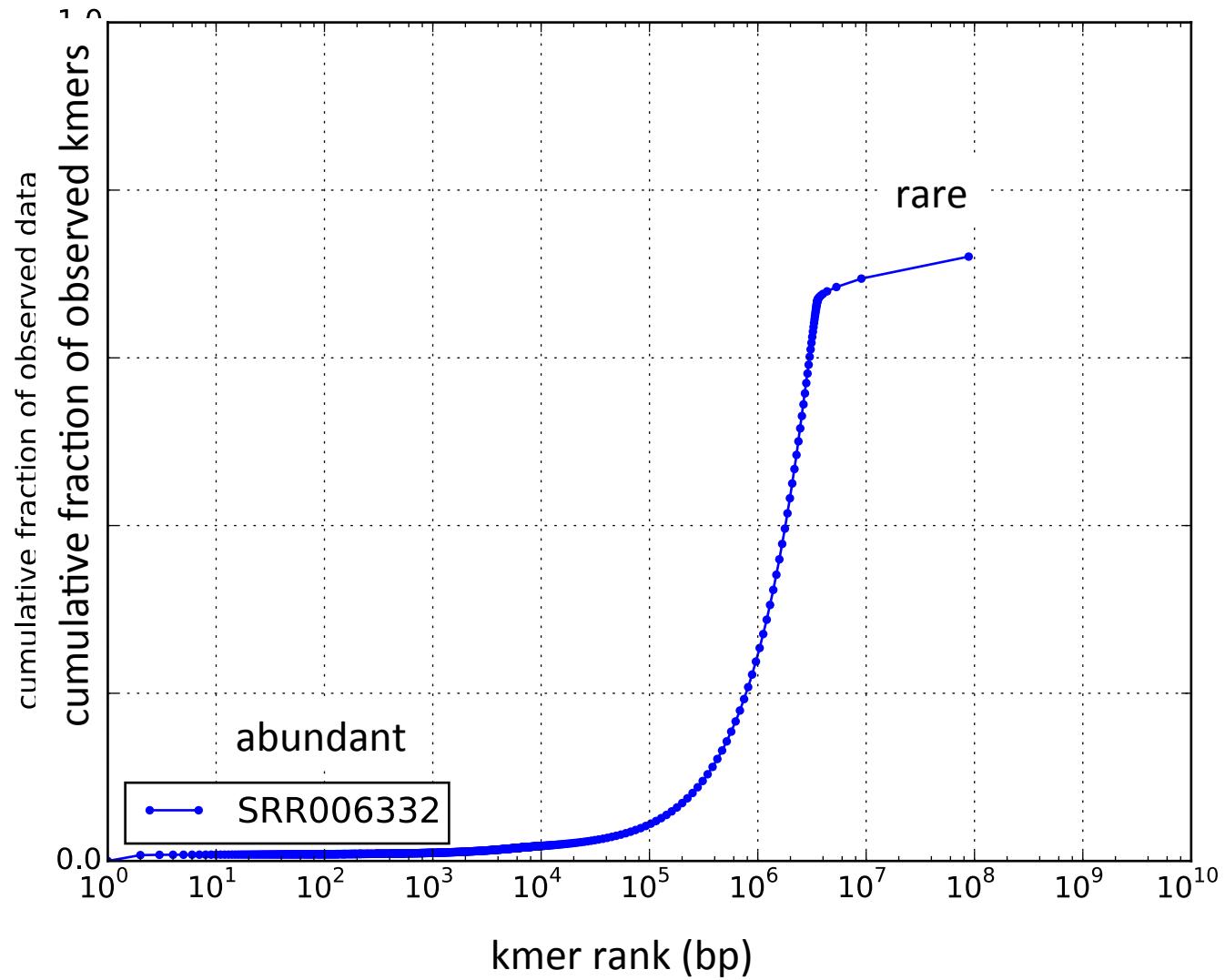
# Abundance /fraction



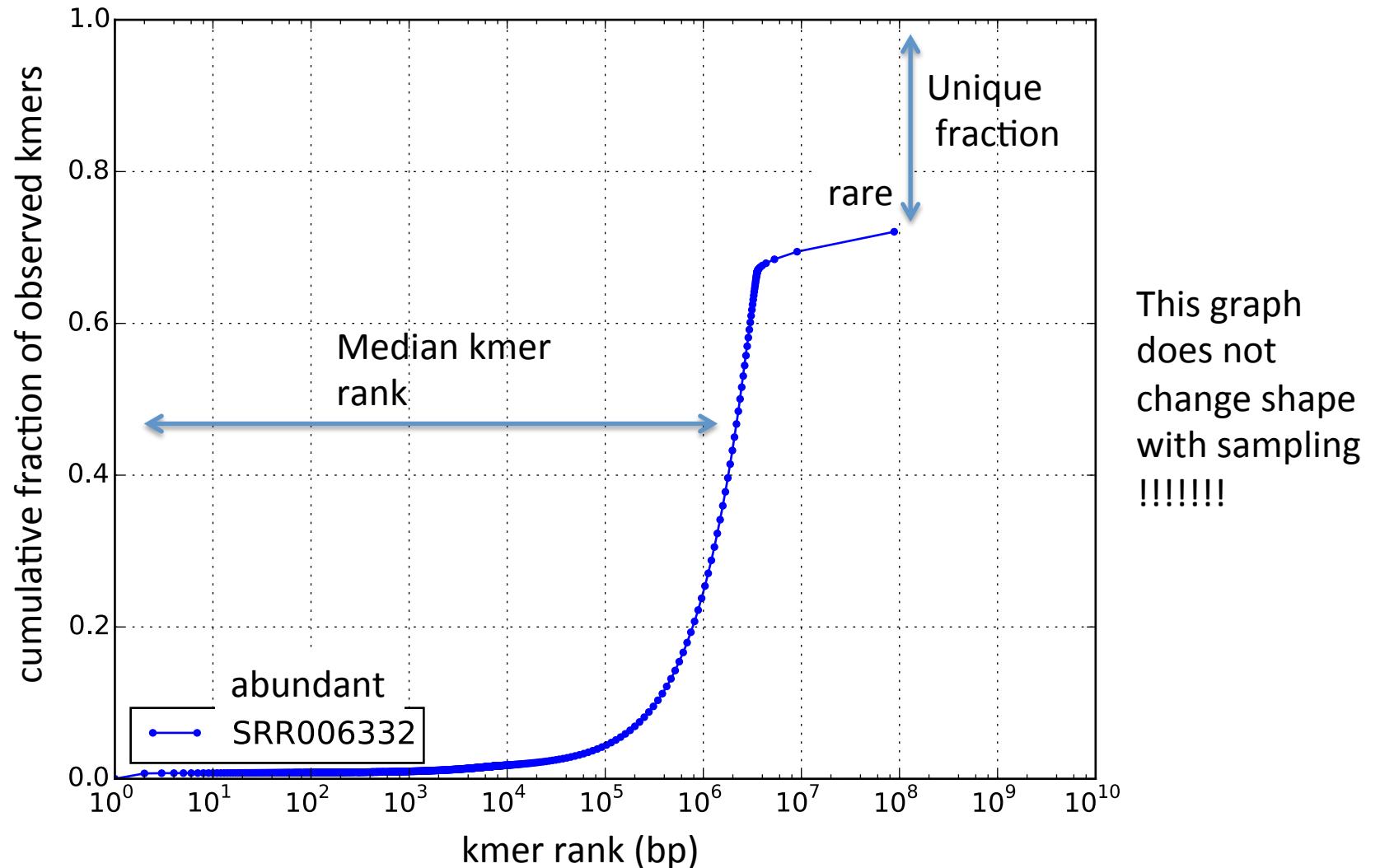
# Cumulative ranked kmer spectrum

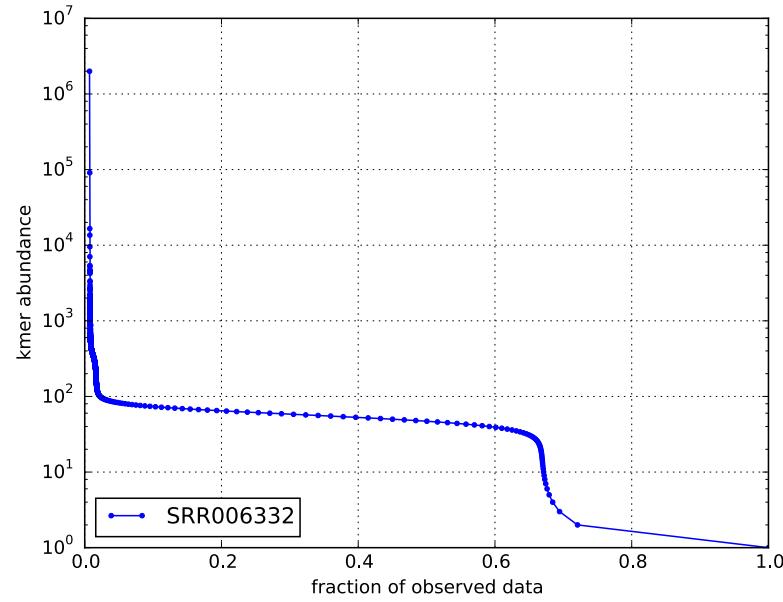
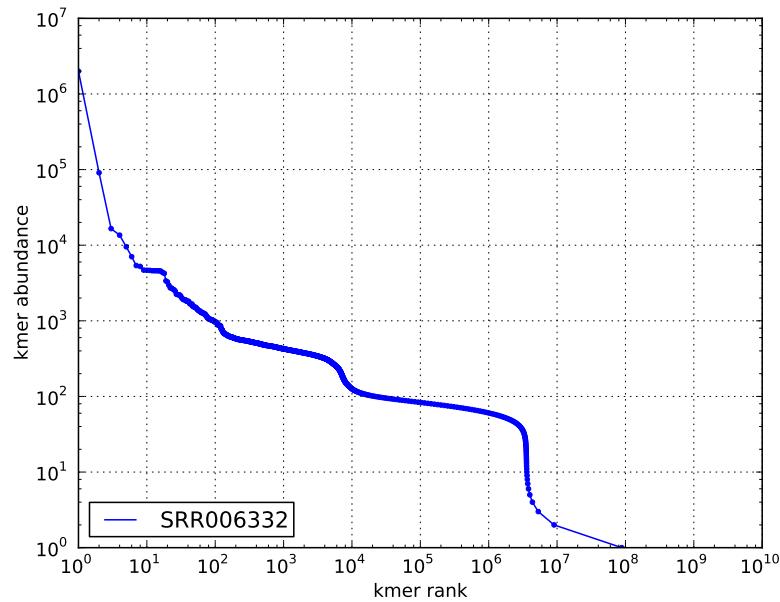
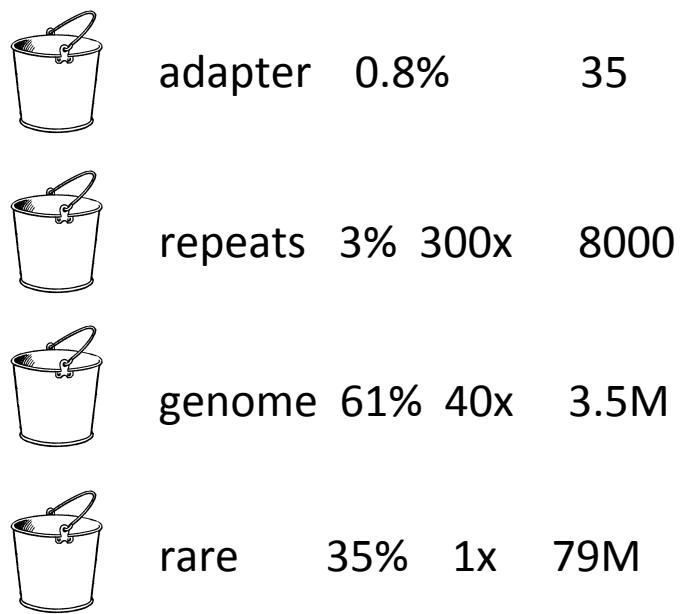
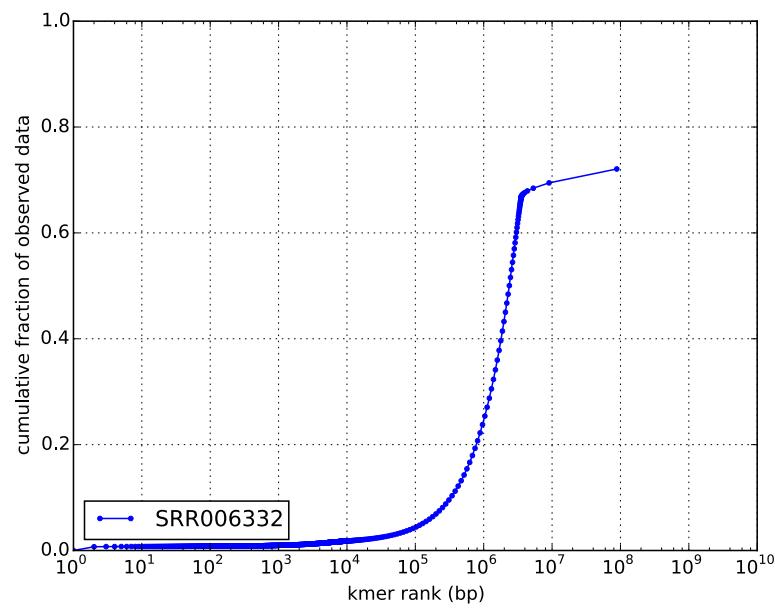


# Ranked kmers consumed



# Ranked kmers consumed



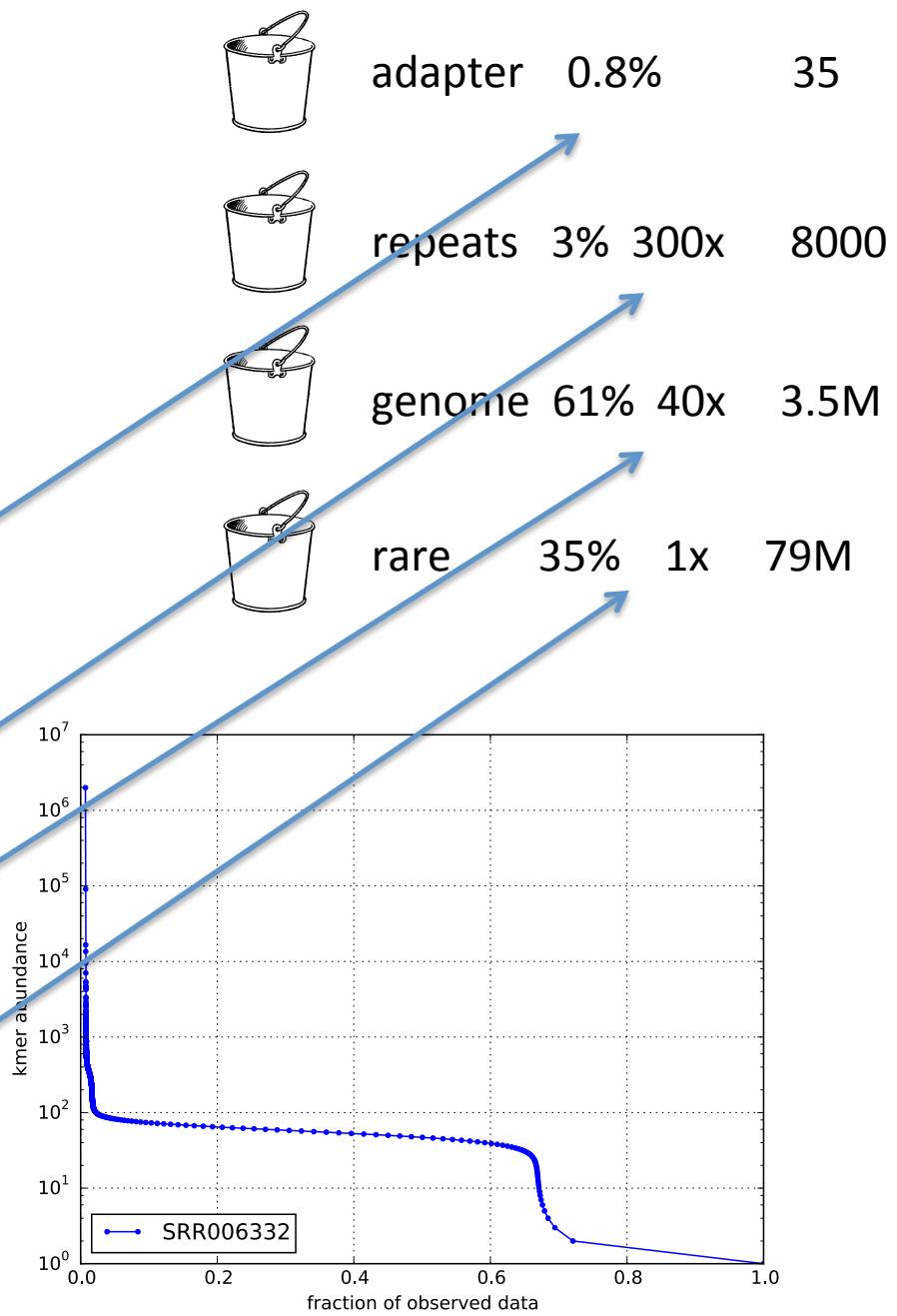
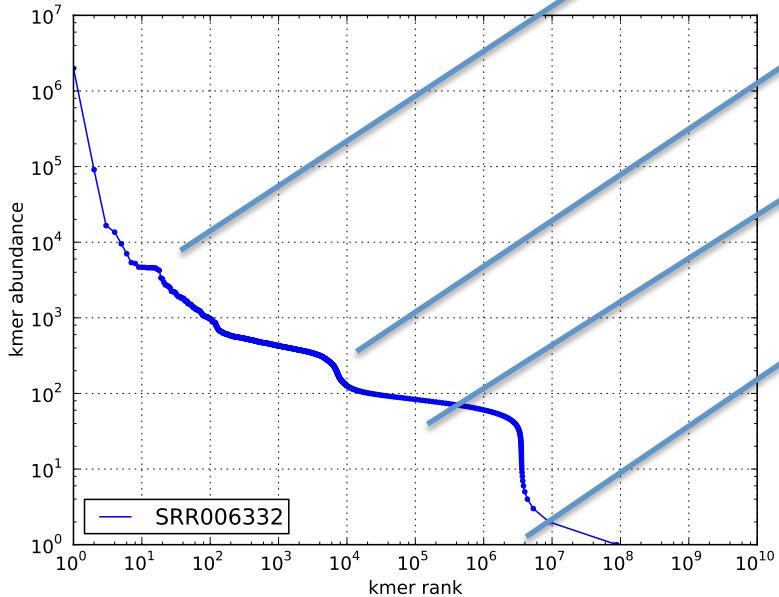
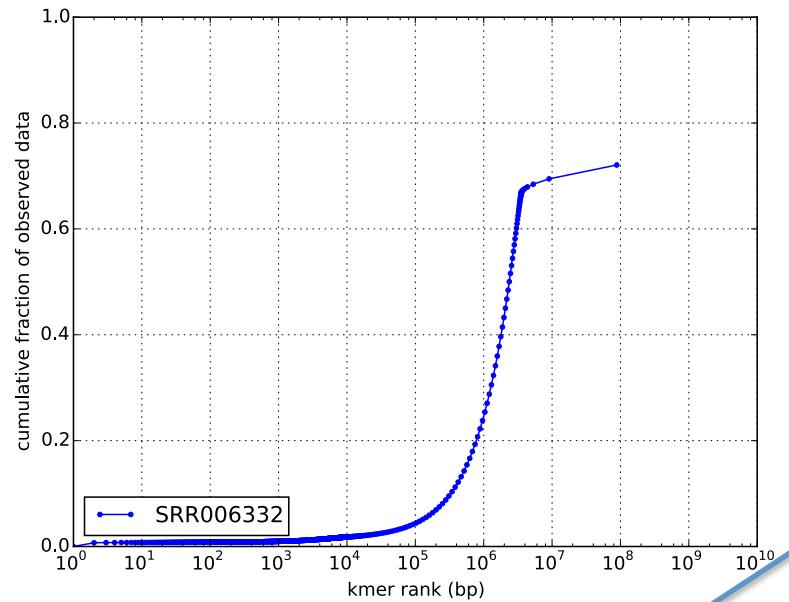


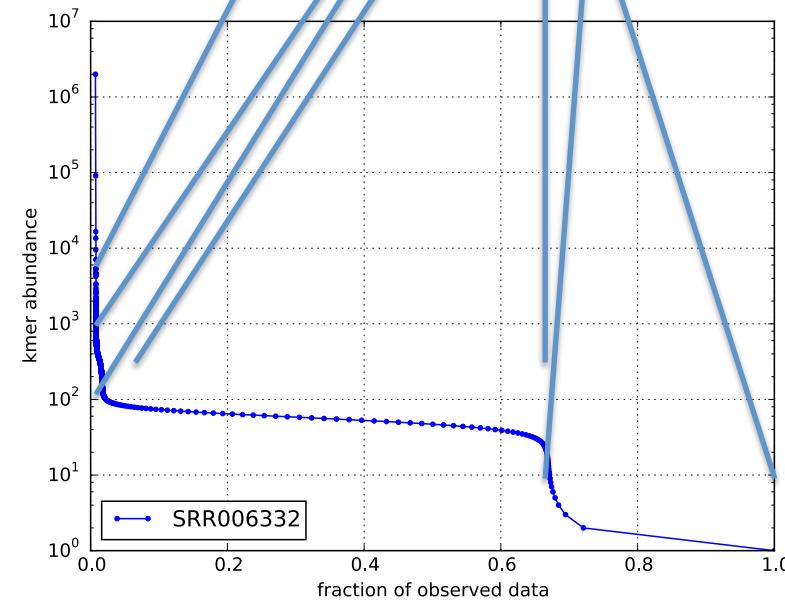
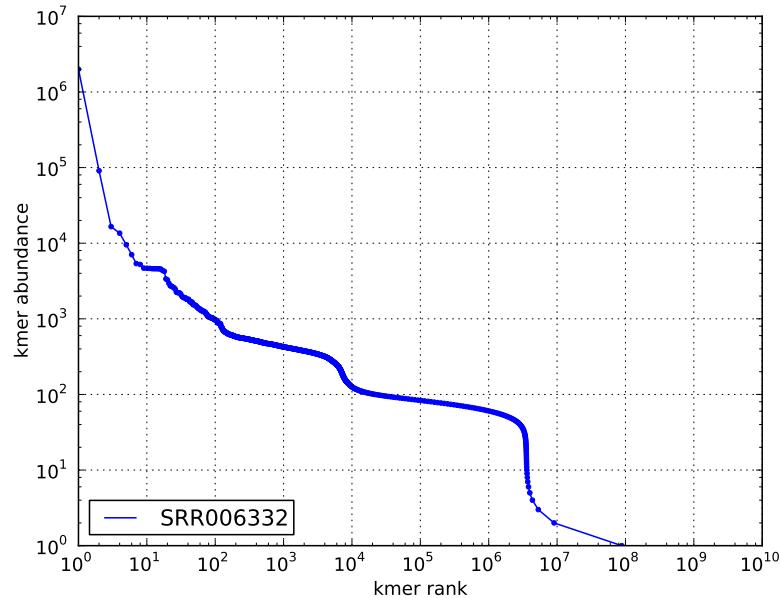
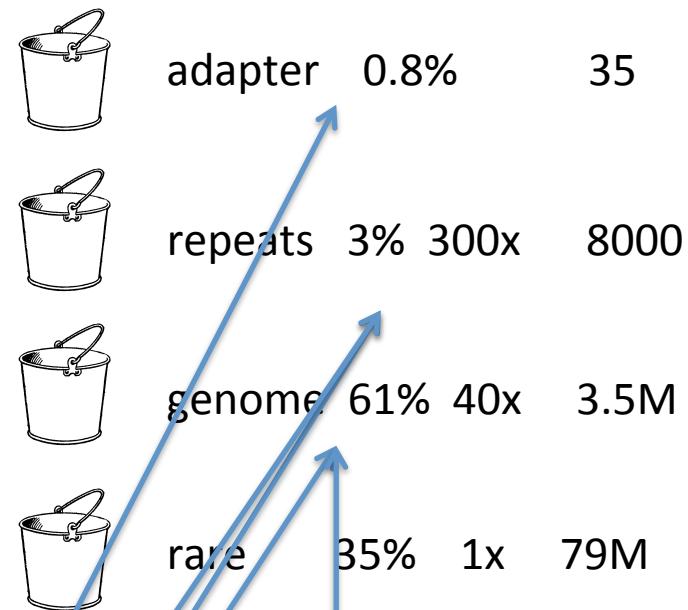
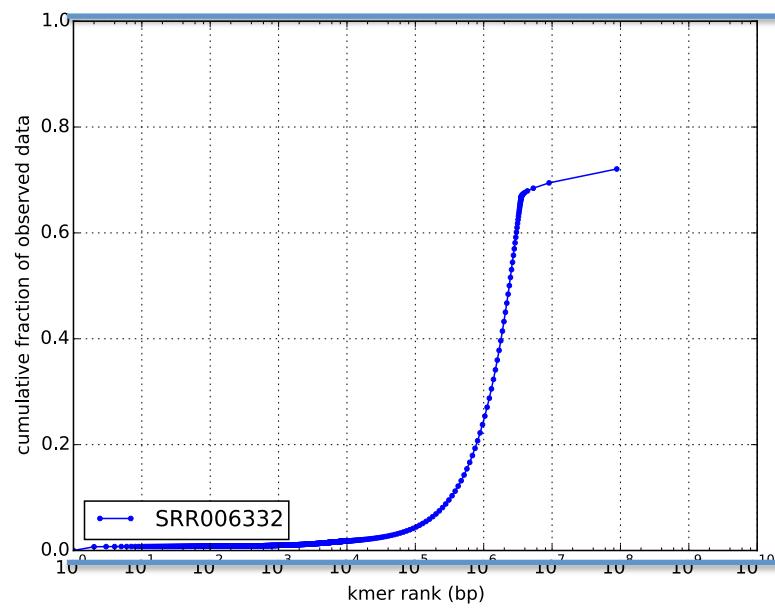
Yes, all these graphs are necessary

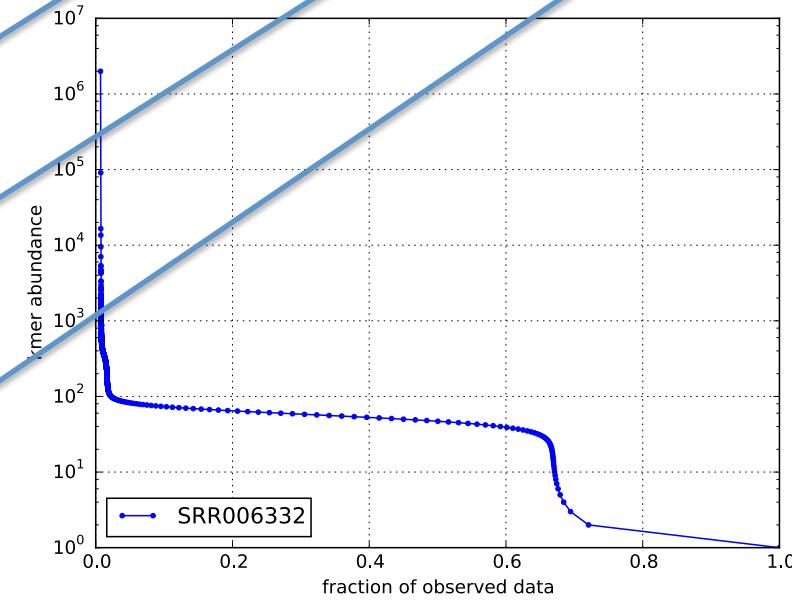
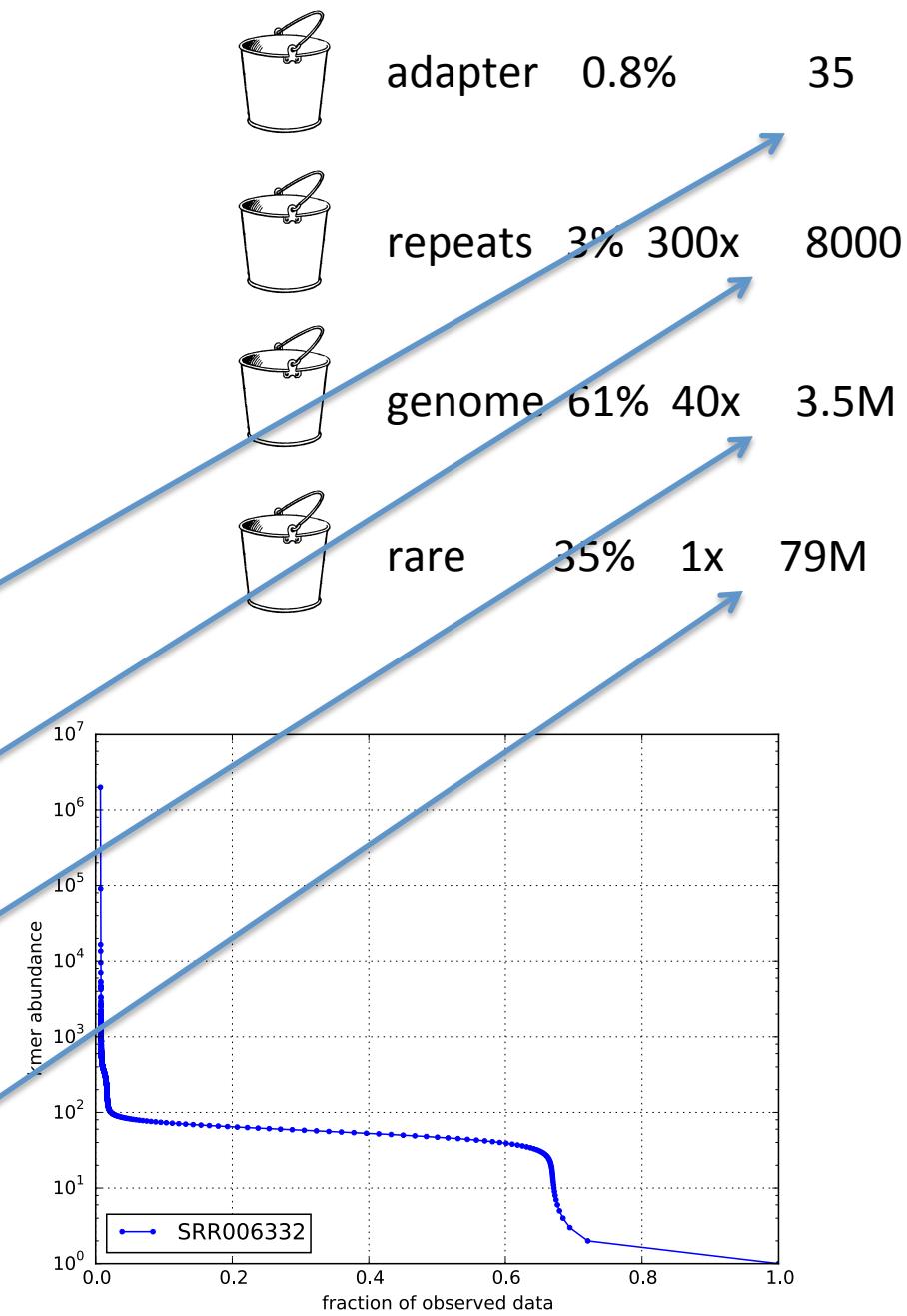
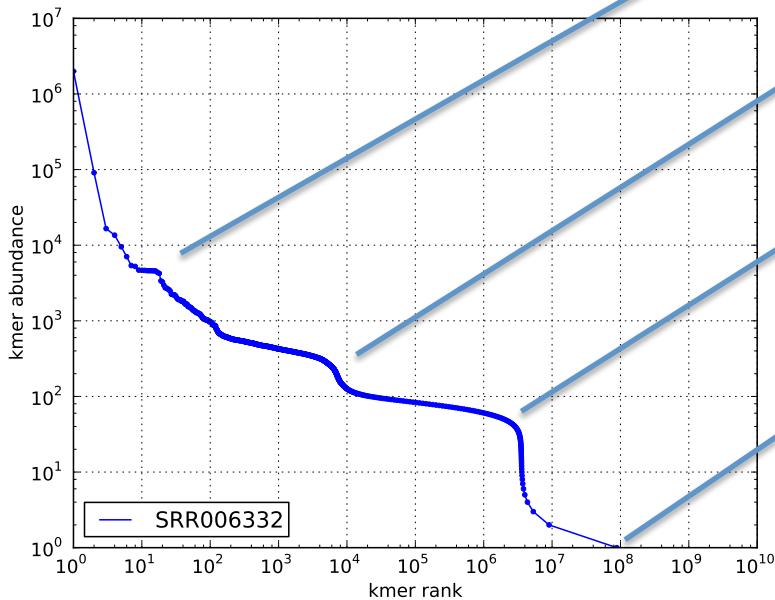
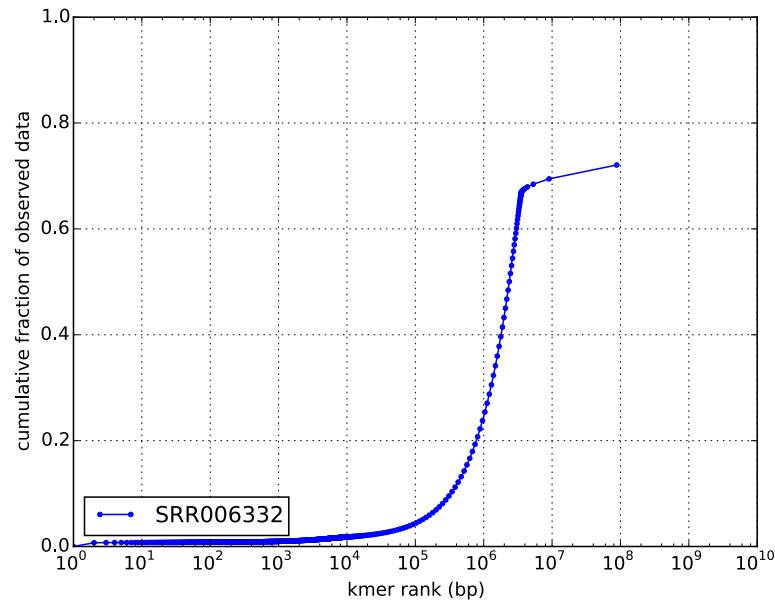
Genome size  $\sim$  number of kmers (M)

depth  $\sim$  abundance (C)

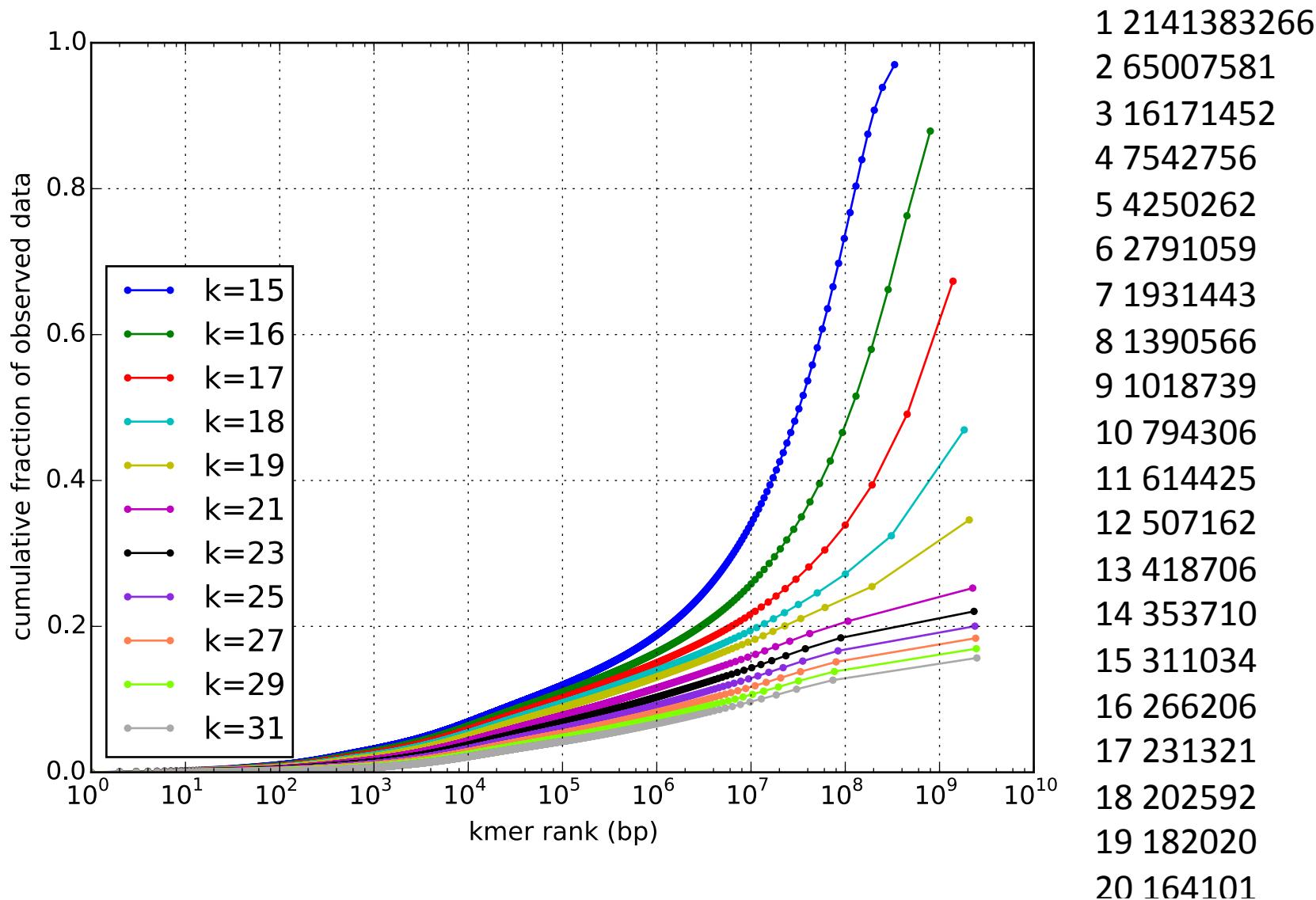
DATA FRACTION (F)







# Check: Human reference genome



# kmer statistical summaries

- H0 kmer richness (VERY BAD)
  - H1 Shannon entropy (BAD)
  - H2 Reymi entropy / Simpson index (GOOD)
- 
- observation-weighted coverage (BAD)
  - observation-weighted size (BAD)
  - observation-median coverage (GOOD)
  - observation-median rank (GOOD)
  - fraction in top 100 kmers (USEFUL)
  - fraction unique (OK but requires size correction)

# kmer statistical summaries

- H0 kmer richness
  - Most of these give answers which vary so strongly with sampling depth as to be unusable.
    - Observation-weighted fraction-of-data metrics behave fairly well. Fractions of the data with particular properties are stable with respect to sampling.
  - (VERY BAD)
  - (BAD)
  - x (GOOD)
  - (BAD)
  - (BAD)
  - (GOOD)
  - (GOOD)
  - EFUL)
- fraction unique (OK but requires size correction)

# kmer statistical summaries

- H0 kmer richness (VERY BAD)
- H1 Shannon entropy (BAD)
- H2 Reymi entropy / Simpson index (GOOD)
  
- observation-weighted coverage (BAD)
- observation-weighted size (BAD)
- **observation-median coverage (GOOD) C50**
- **observation-median rank (GOOD) M50**
- **fraction in top 100 kmers (USEFUL) F100**
- **fraction unique (OK but requires size correction)**

# Generalities from the kmer counting mines

- Many datasets have as much as 5-45% of the sequence yield in **adapters**.
- **FEW DATASETS have well-separated abundance peaks** (of the sort metavelvet was engineered to find)
- Diverse datasets have a **featureless, geometric relationship** between kmer rank and kmer abundance.
- **Shannon entropy is oversensitive to errors.** \*

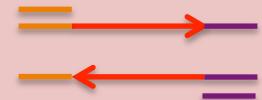
# Three approaches to error correction:

- Filtering (Q-values, tile quality)
- Pair-joining
- Error correction

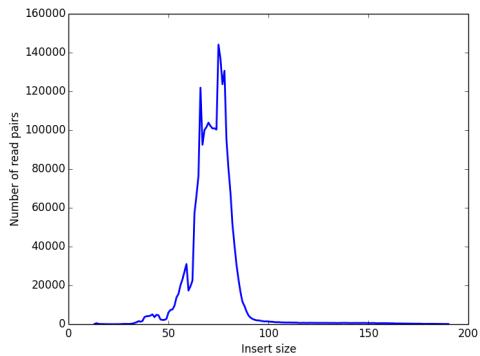
# Choosing insert size



$0 < \text{insert size} < L$   
eats into adapters



Sequencing adapter  
contamination

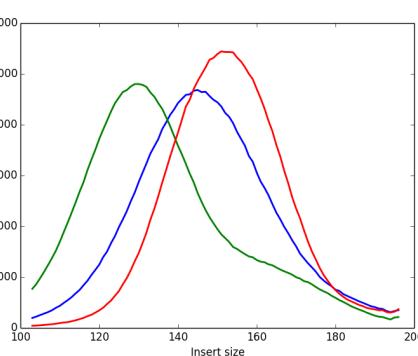


60% of sequencing  
yield wasted on  
adapters; half of the  
remainder is  
redundant.

overlapping reads  
 $L < \text{insert size} < 2L$



Error-corrected  
overlapping reads

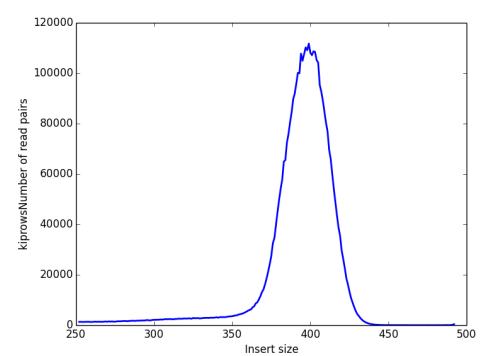


25% of sequencing  
yield consumed in  
error-correcting  
overlaps. 75% of  
basepairs in reads  
 $1.5-1.8L$

unsequenced insert  
 $2L < \text{insert size}$



Unsequenced insert  
No error correction



None of the  
sequencing is lost,  
but past 200bp the  
uncorrected quality is  
very poor.

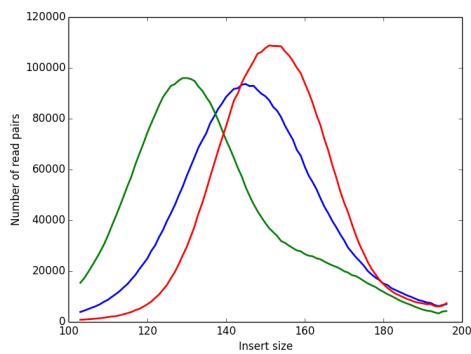
# Overlapping reads

**overlapping reads**

$$L < \text{insert size} < 2L$$



Error-corrected  
overlapping reads



- Careful size selection results in 90+ % of reads with overlapping ends.
- **Read joining is easy, cheap.**
- Produces pile of joined reads and two piles of unjoined reads (and a whole pile of headache keeping read pairs together)
- Error correction when reads differ.
- Overlap operation benignly reduces total coverage
- Size selection may disturb mixture proportions (“bias”)
- What to do with the qualitatively different unoverlapped pairs?
- Keep them: possible size – length effect
- Discard them: possible size—depth effect

# kmer statistical summaries

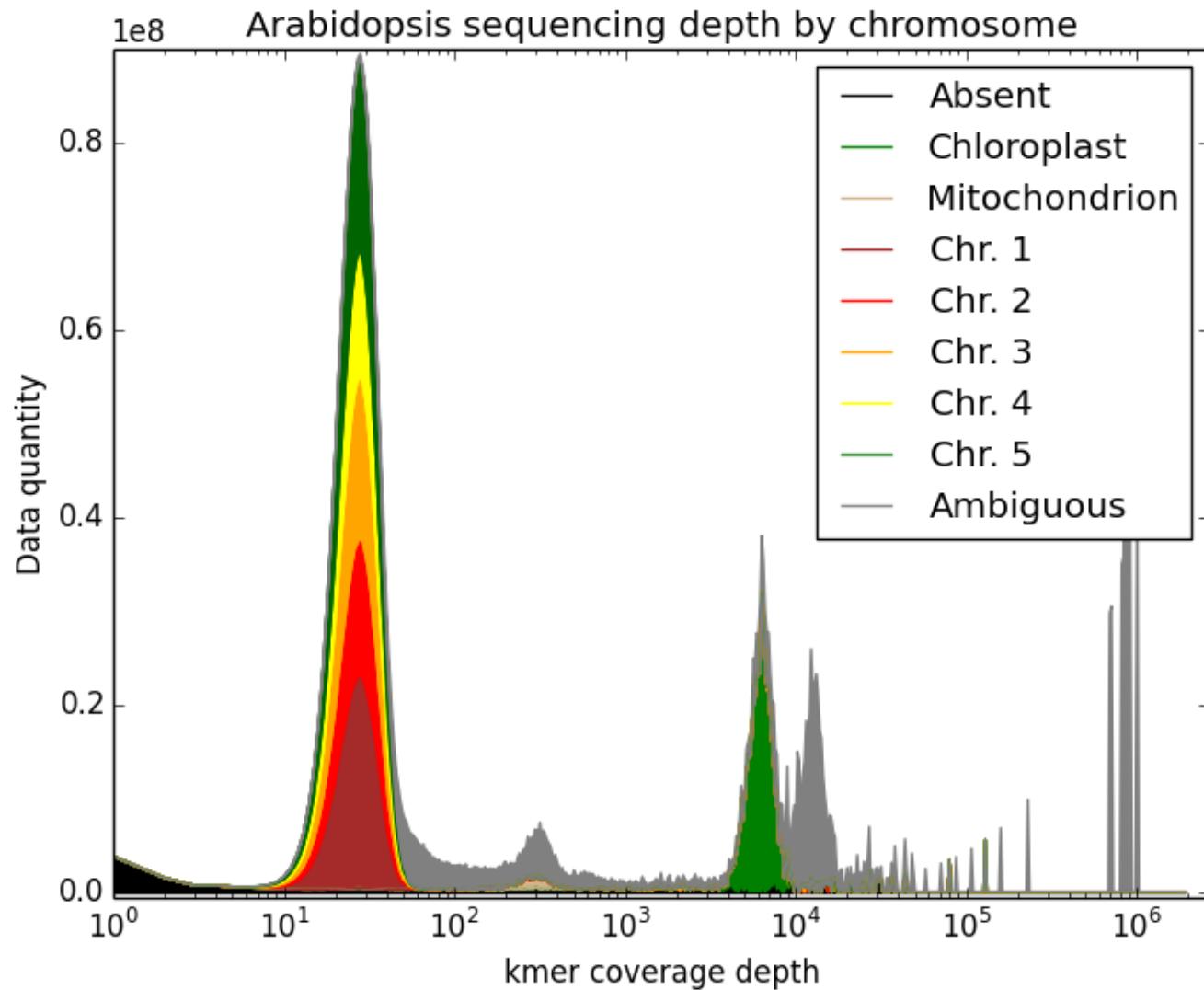
**observation-median coverage (C50)**

**observation-median rank (M50)**

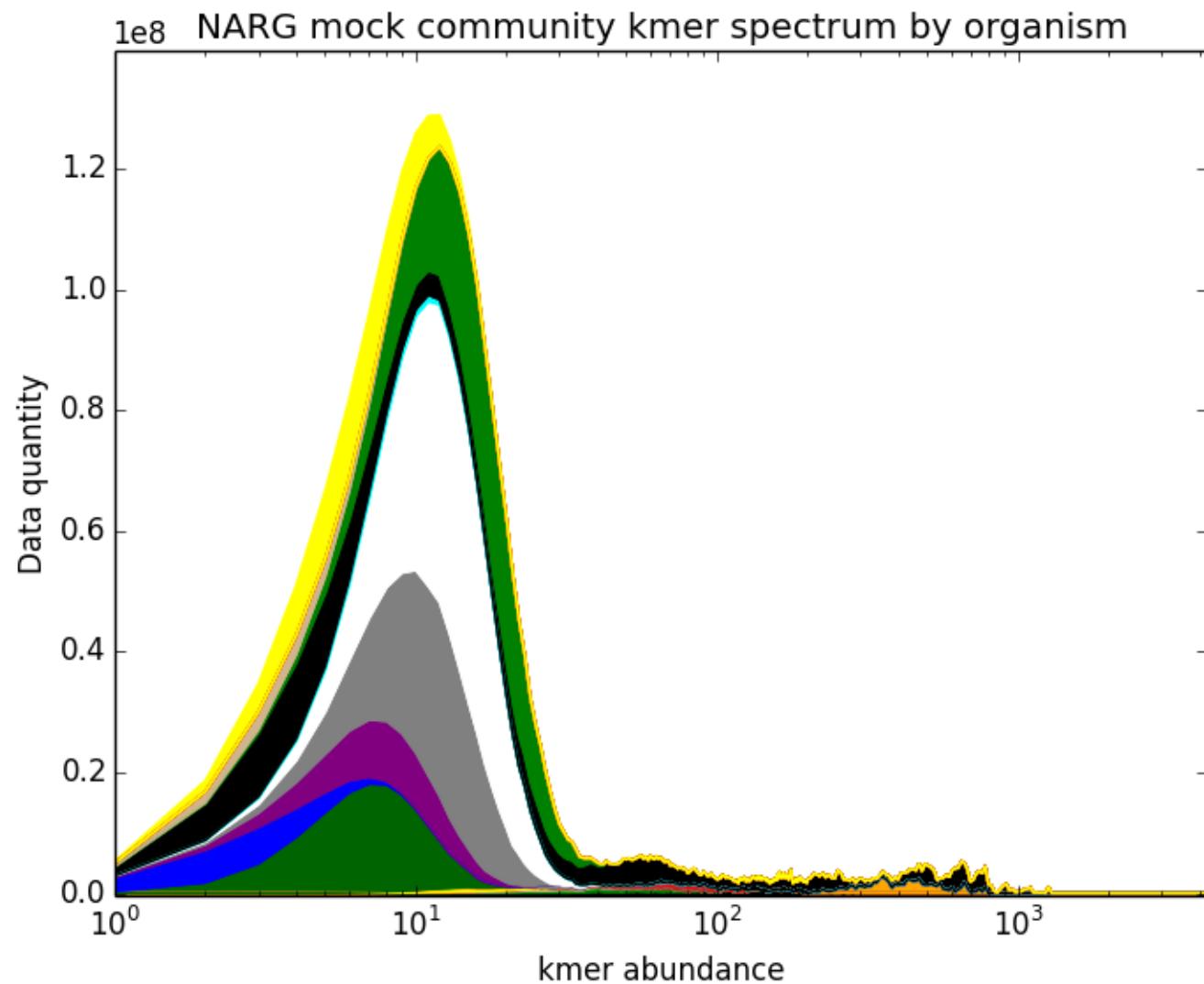
**fraction in top 100 kmers (F100)**

**fraction unique**

# Straightforward abundance distribution: plant



# Complex abundance distribution: mock community



# The double-counting quirk

|                 |     |
|-----------------|-----|
| ATCGCGAAAAGTCCC | 2   |
| AAAAAAAAAAAAAA  | 459 |
| AAAAAAAAAAAAAC  | 71  |
| AAAATAAAAAAATA  | 1   |
| AAAAAAAAAAAAAG  | 36  |
| ACATGAAAAACAAC  | 1   |
| AAAAAAAAAAAAAT  | 23  |
| AAAAAAAAAAAAACA | 95  |
| GTAGGAAAAGCCCAC | 1   |
| AAAAAAAAAAAAACC | 7   |
| AAAAAAAAAAAAACG | 8   |
| AAAAAAAAAAAAACT | 9   |
| AAAAAAAAAAAAAGA | 36  |
| AACAAGAAAAACAAA | 1   |
| AAAAAAAAAAAAAGC | 10  |
| AAATAAAAAAAATAG | 1   |
| AACAGAAAAACACG  | 1   |
| AAAAAAAAAAAAAGG | 2   |
| AAAAAAAAAAAAAGT | 6   |

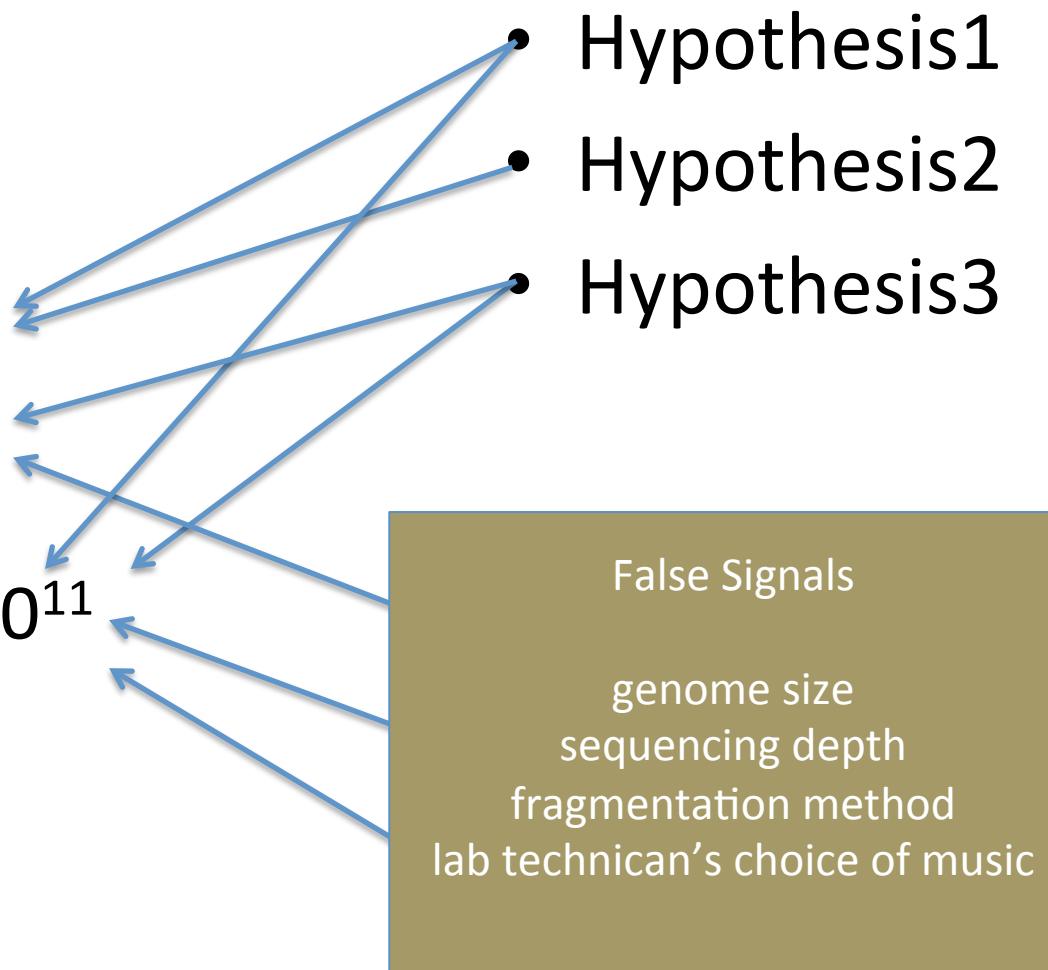
There is a sense in which  
**AAAAAAAAAAAAAA**  
and  
**TTTTTTTTTTTT**  
are the same.

- 2 approaches:
- double-count everything
  - only count once



# We don't care about the left side!

- Observation 1
- Observation 2
- ...
- Observation  $1+10^{11}$



# Generic experiment

- Observation 1
- Observation 2
- ...
- Observation  $1+10^{11}$

