

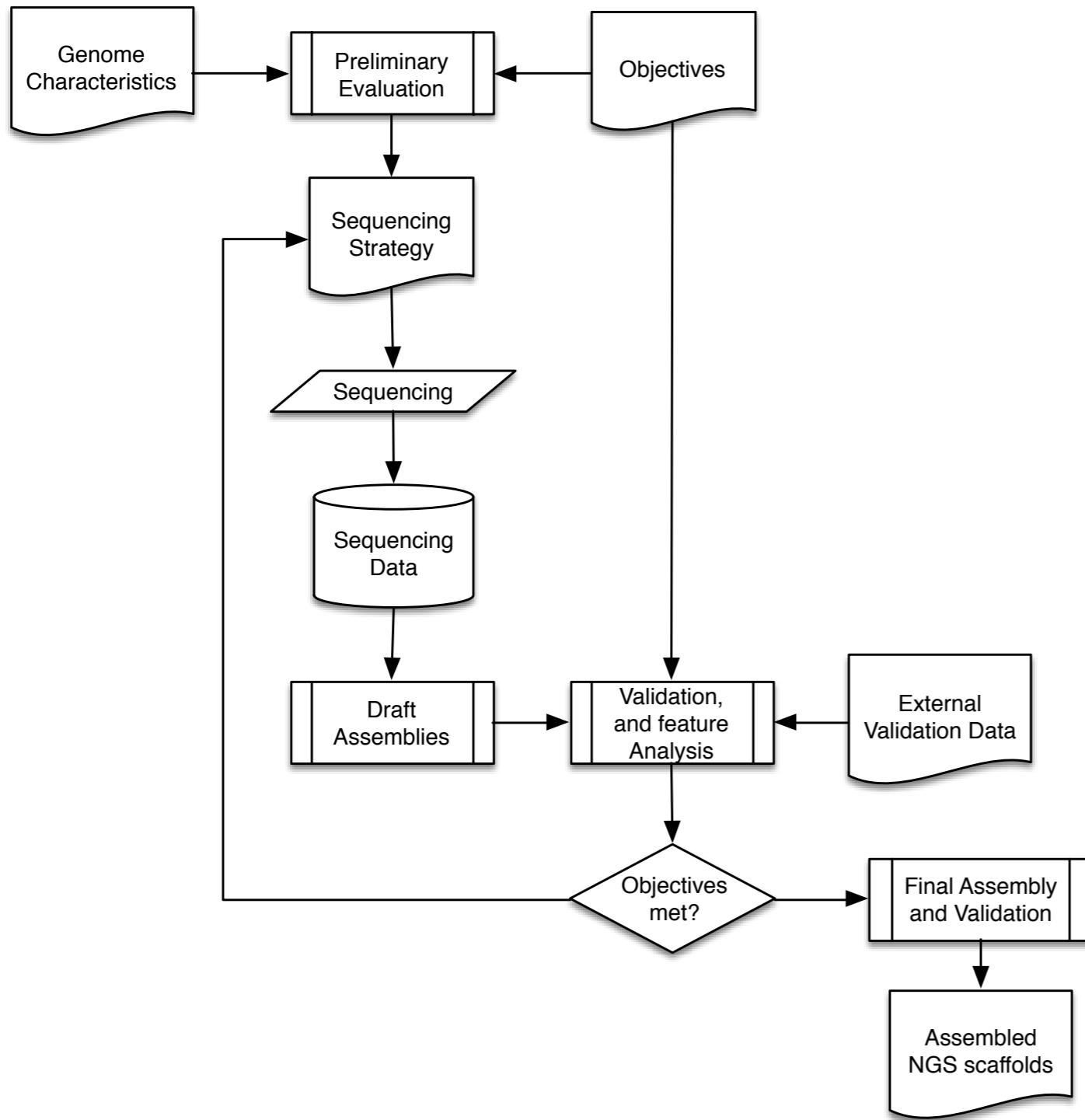
4 - Genome Assembly and Validation (Concepts)

Wednesday afternoon

*Bernardo J. Clavijo
Richard Smith-Unna
Gonzalo Garcia / Jon Wright*



Assembly project workflow | Prior Knowledge

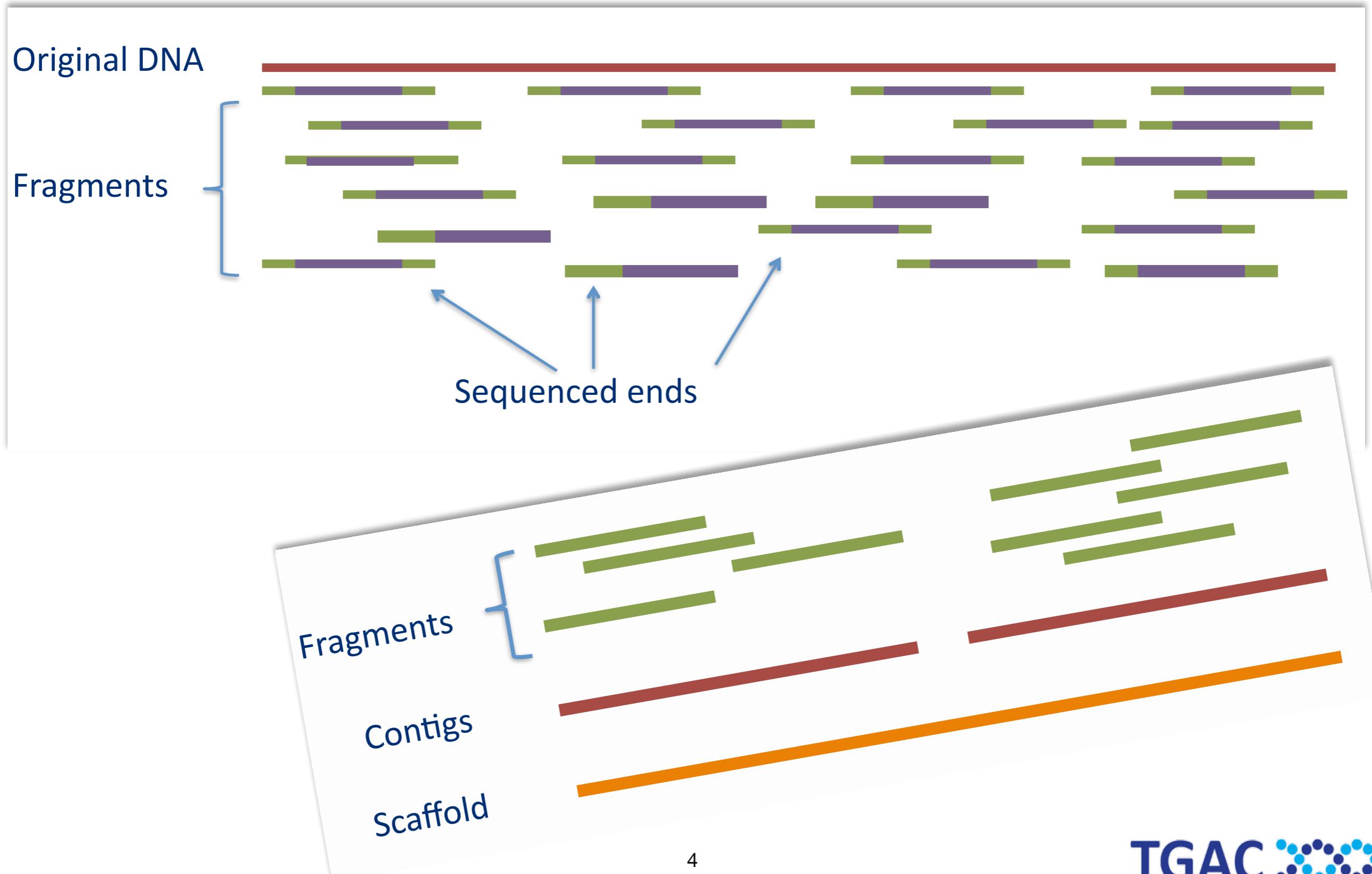


- Kariotype: Genome size, Ploidy
- Heterozygosity
- GC content
- Contaminants / Symbionts
- Data Sets:
 - Close relatives
 - Genes / ESTs / RNAseq / Markers
- Mitocondria
- Chloroplast

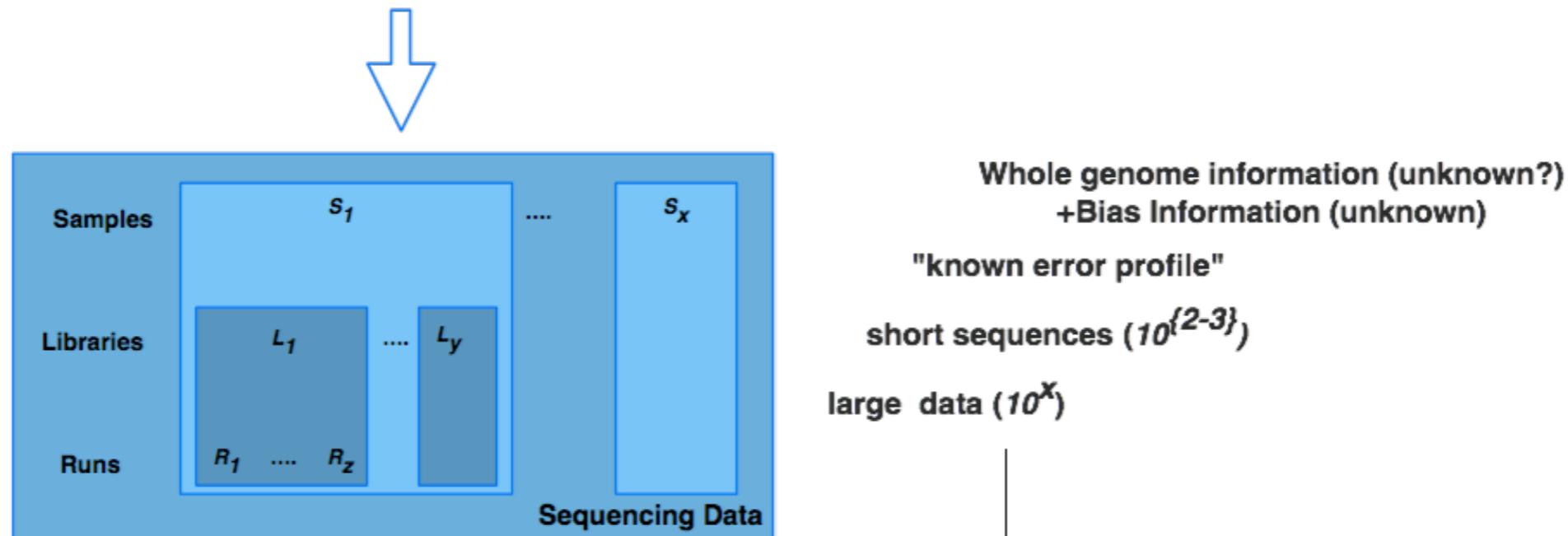
Experiment design (you choose the data!)

- **Know your biological question.**
- Plan your data processing (from an information perspective).
- Decide on conditions and biological/technical replicas.
- Decide on technologies and coverages:
 - How will the typical bias affect your experiment?
 - Is the coverage enough? Significant results?

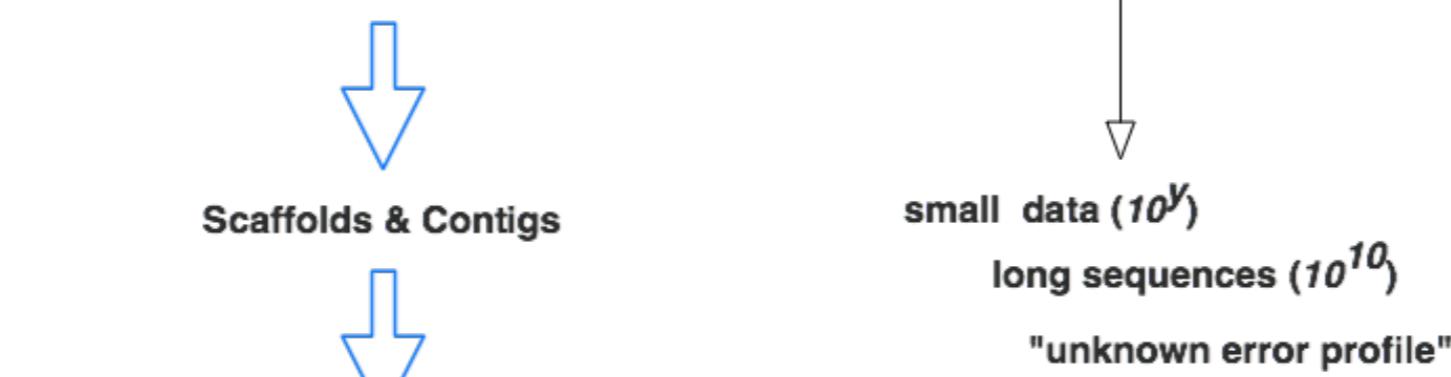
The genome assembly problem (WGS)



Planning and "informed guesses"



Assemble and Scaffold



Validate and release

Whole genome information (known?)
+Bias Information (unknown, reduced?)

The assembly is just a probabilistic model of a genome, condensing the information from the experimental evidence.

All the information is already present in the experimental results.

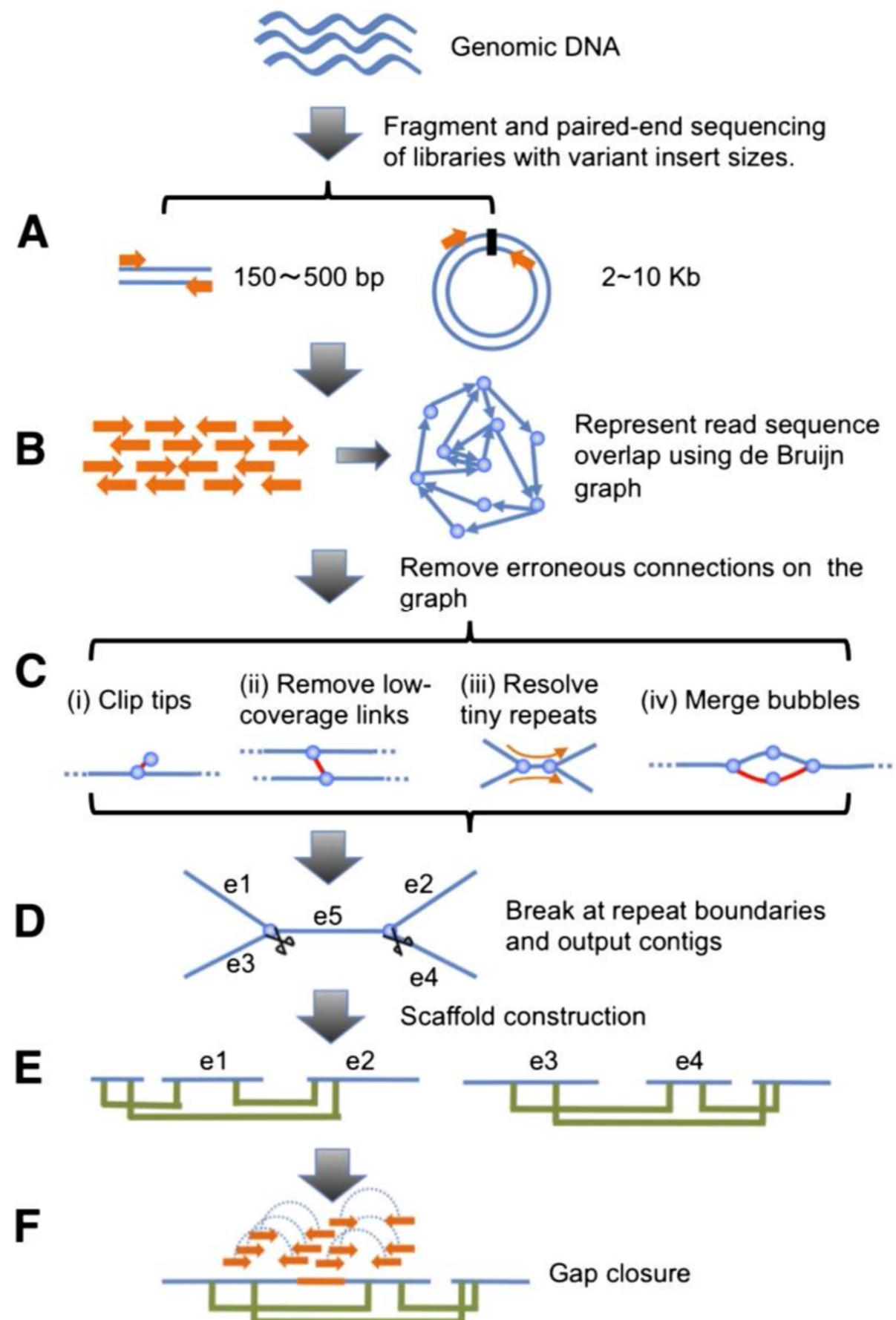
A correct assembly has:

The right *motifs*,
the correct number of times,
in correct order and position.

None of which is assessed by length stats.

A modern assembler

Using SOAPdenovo2 as an example



Fragment and paired end sequencing of libraries with variant insert sizes.

A



B



Remove erroneous connections on the graph

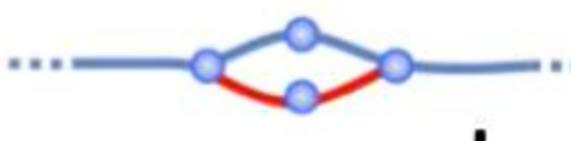
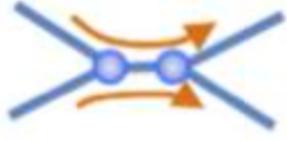
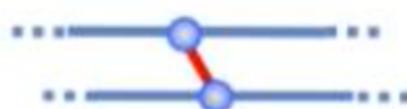
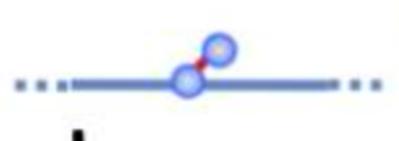
C

(i) Clip tips

(ii) Remove low-coverage links

(iii) Resolve tiny repeats

(iv) Merge bubbles



D



Break at repeat boundaries and output contigs

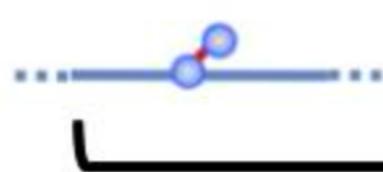




Remove erroneous connections on the graph

C

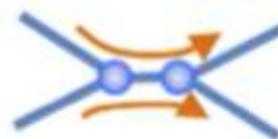
(i) Clip tips



(ii) Remove low-coverage links



(iii) Resolve tiny repeats

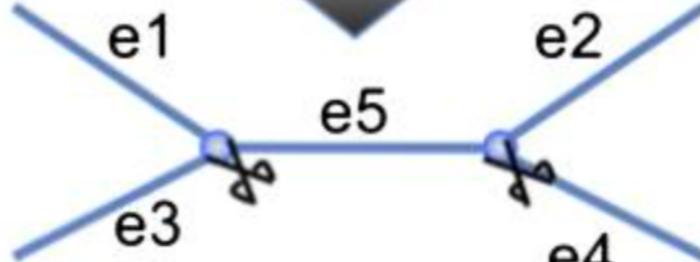


(iv) Merge bubbles



D

Break at repeat boundaries
and output contigs



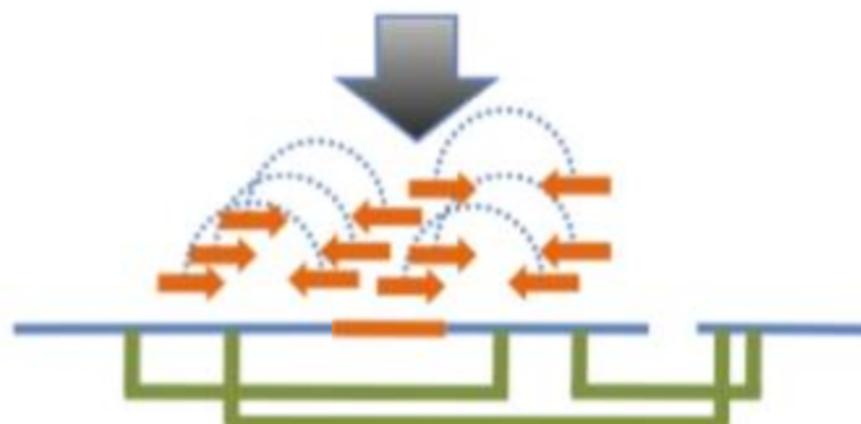
Scaffold construction

E



F

Gap closure



Assembly validation

The power of the kmer spectra...

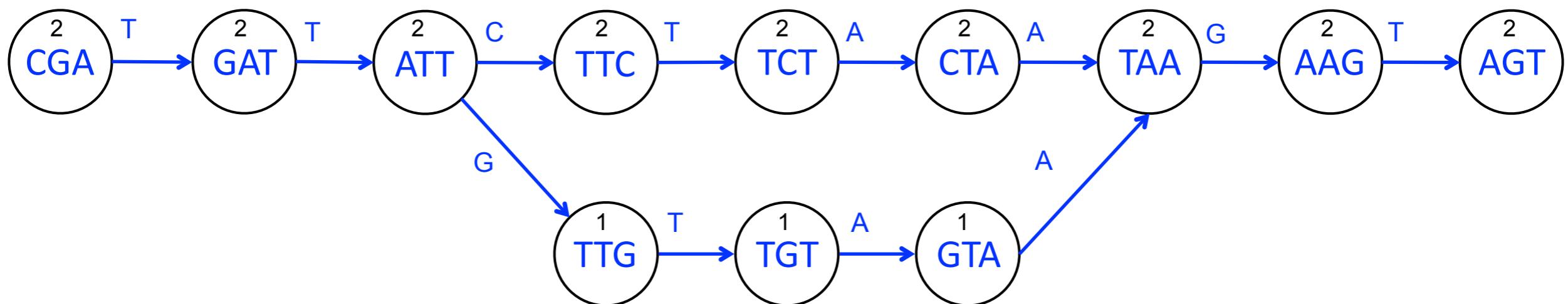
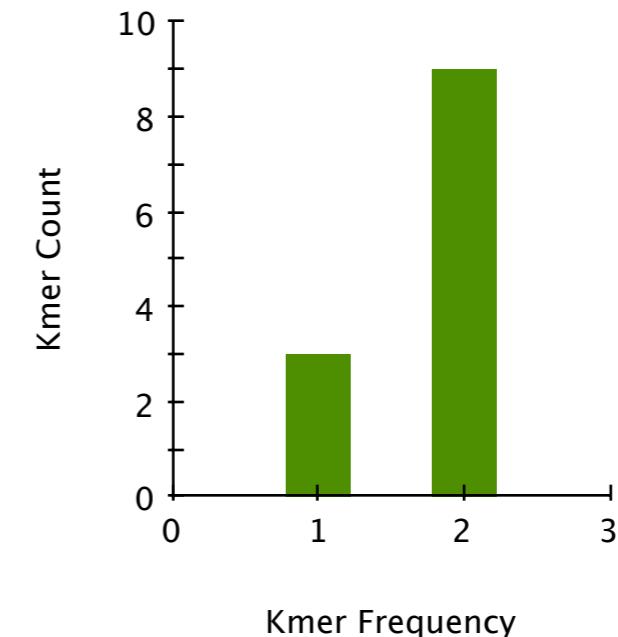


The right *motifs*,
the correct number of times,
in correct order and position.

Counting kmers

```
>seq1  
TTCTAAAGT  
>seq2  
CGATTCTA
```

```
>seq3  
CGATTGTAAGT
```



The kmer spectra

spectra | 'spektrə |

plural form of **SPECTRUM**.

spectrum | 'spektrəm |

noun (pl. **spectra** | -trə | or **spectrums**)

1 a band of colours, as seen in a rainbow, produced by separation of the components of light by their different degrees of refraction according to wavelength.

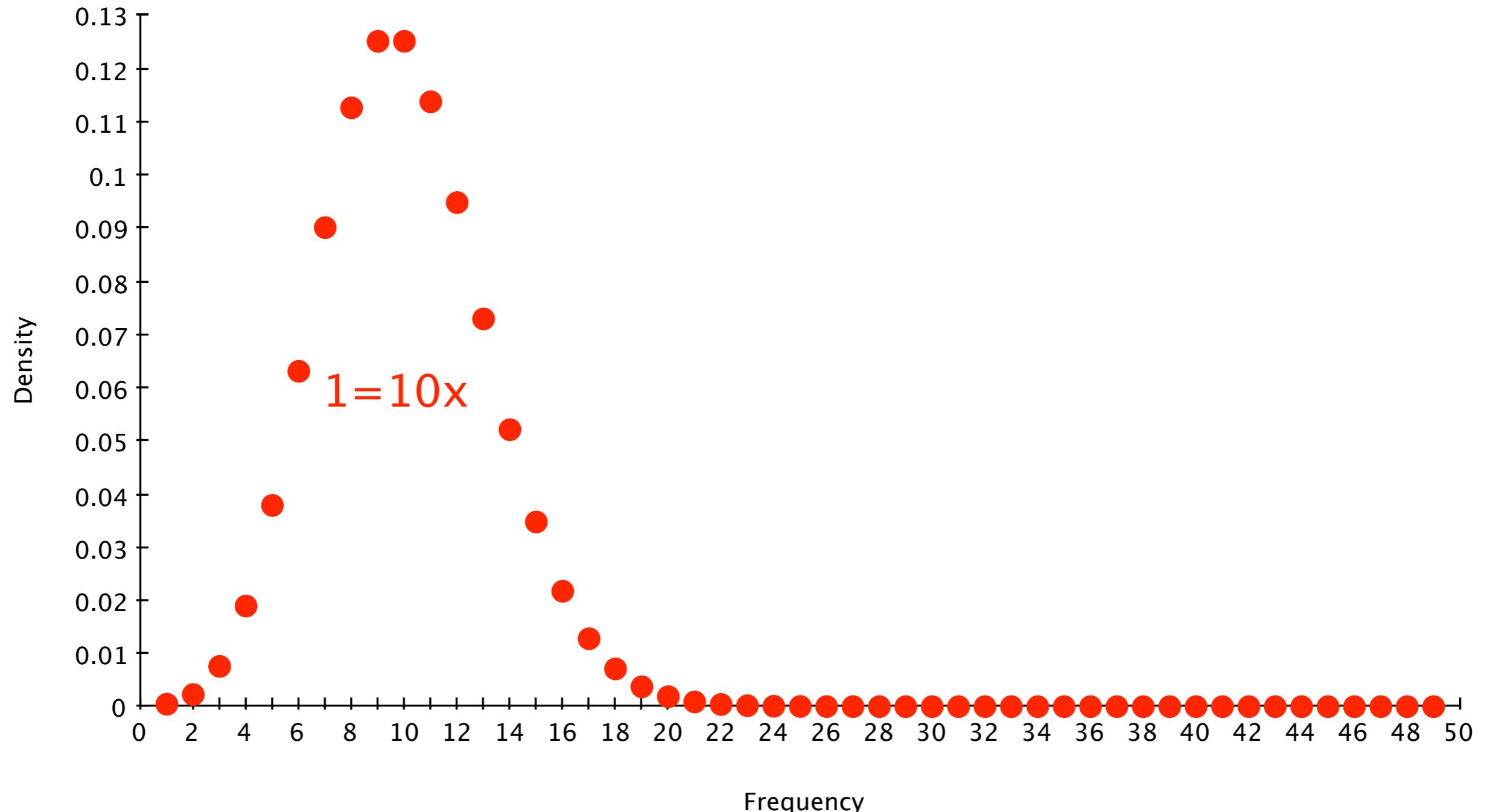
- (**the spectrum**) the entire range of wavelengths of electromagnetic radiation.
- a characteristic series of frequencies of electromagnetic radiation emitted or absorbed by a substance.
- the components of a sound or other phenomenon arranged according to such characteristics as frequency, charge, and energy.

2 used to classify something in terms of its position on a scale between two extreme points: *the left or the right of the political spectrum*.

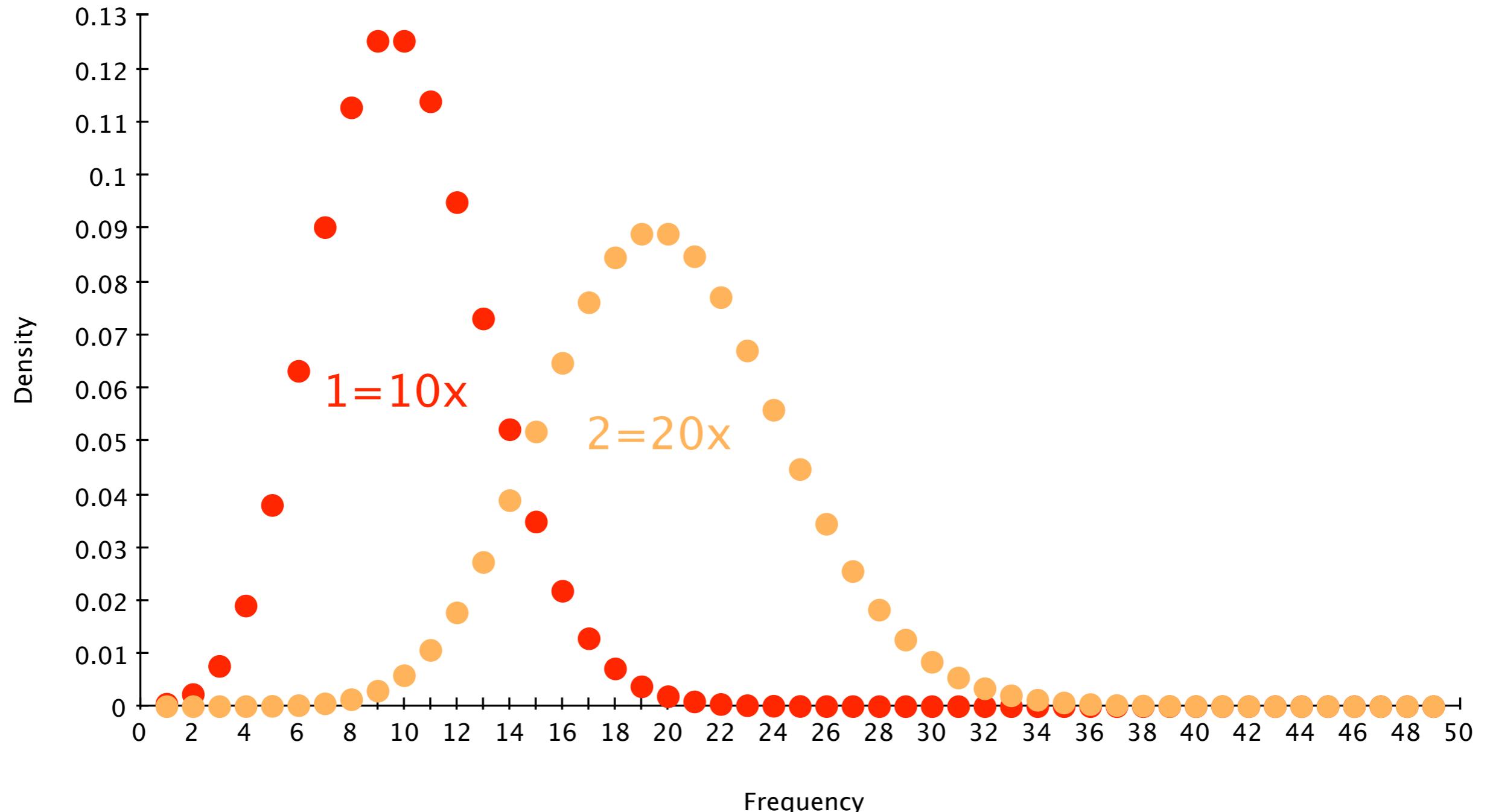
- a wide range: *self-help books are covering a broader and broader spectrum*.

ORIGIN early 17th cent. (in the sense ‘spectre’): from Latin, literally ‘**image, apparition**’, from **specere** ‘to look’.

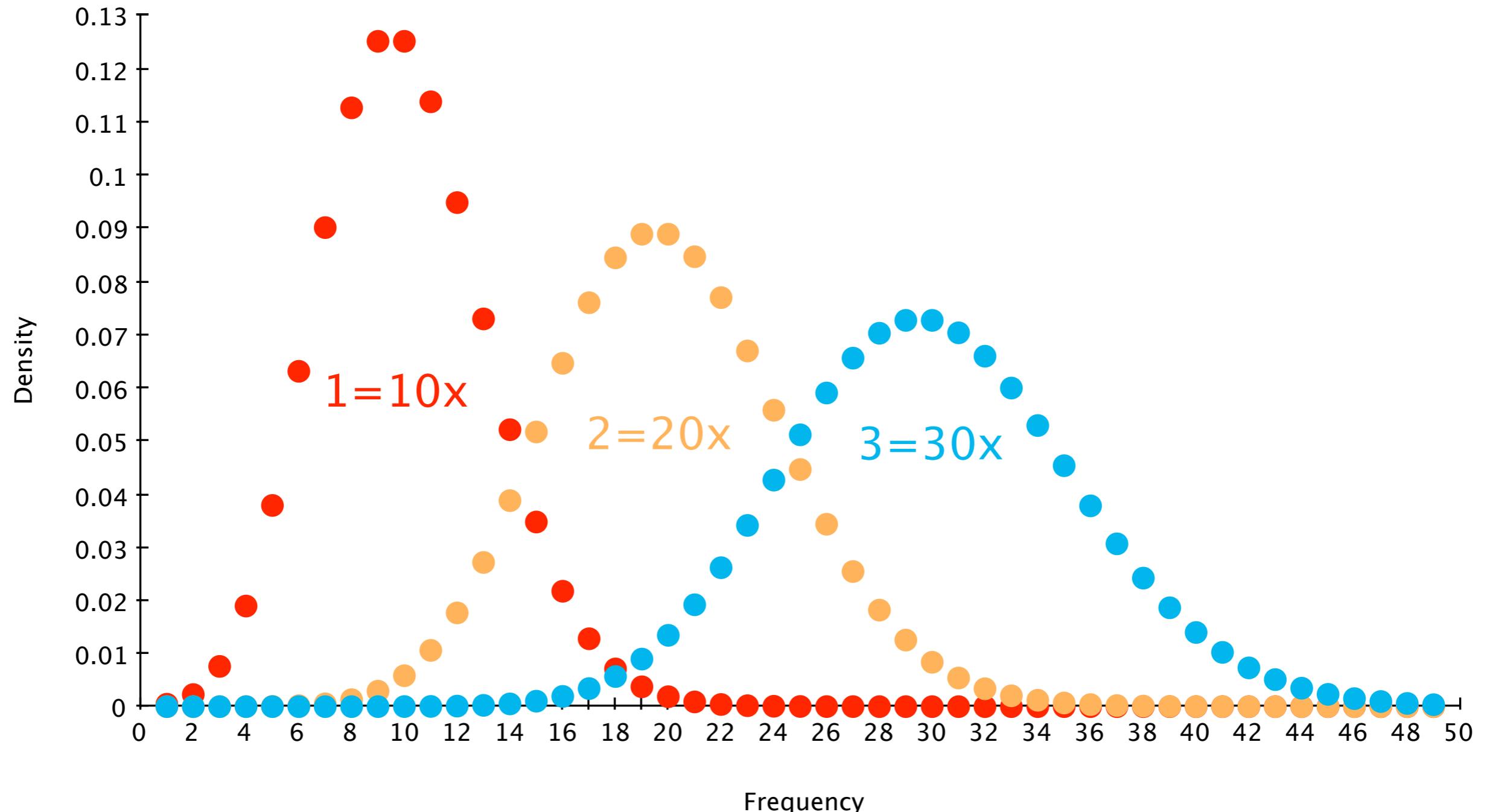
The kmer spectra's components



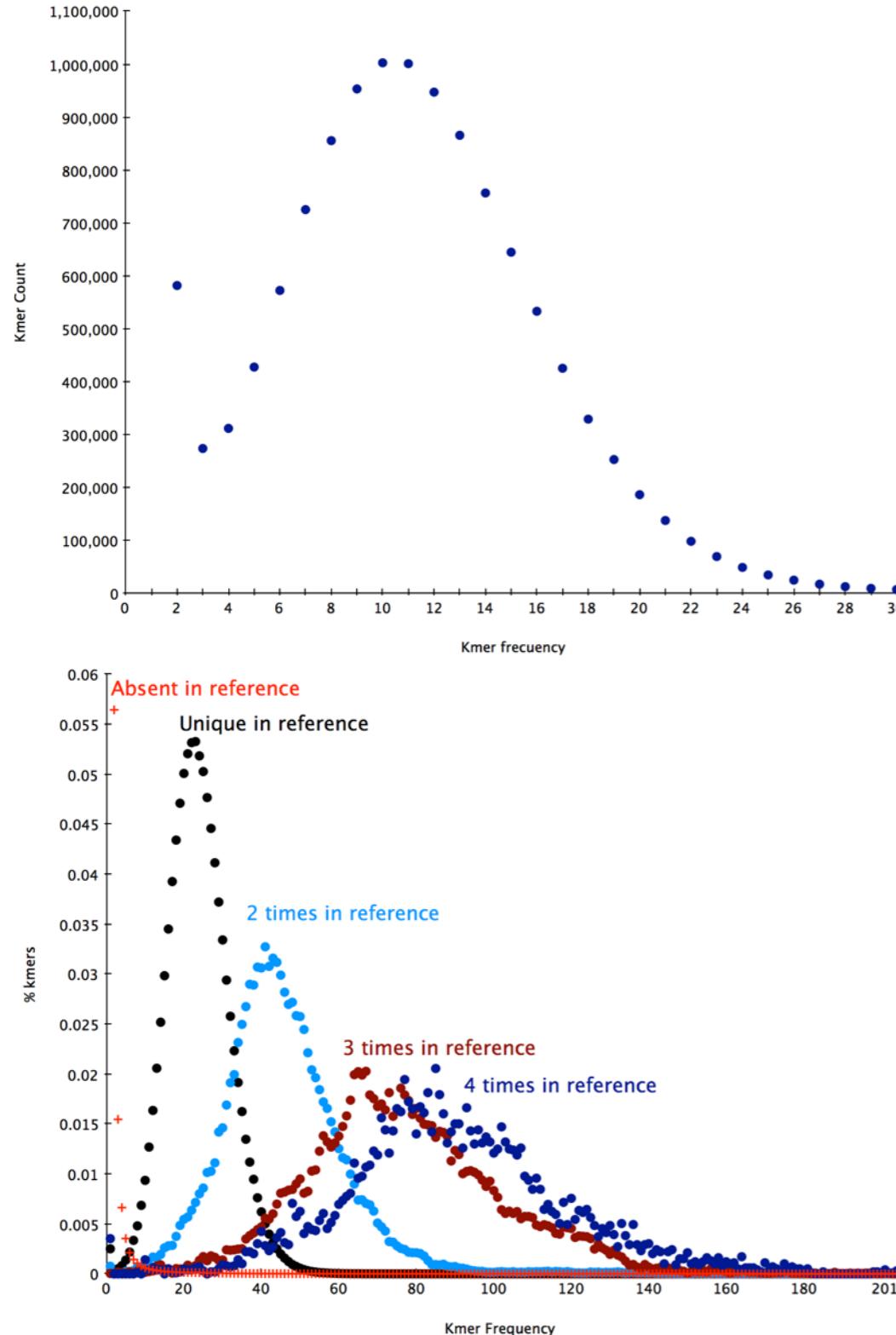
The kmer spectra's components



The kmer spectra's components



The kmer spectrum... and its dissection.



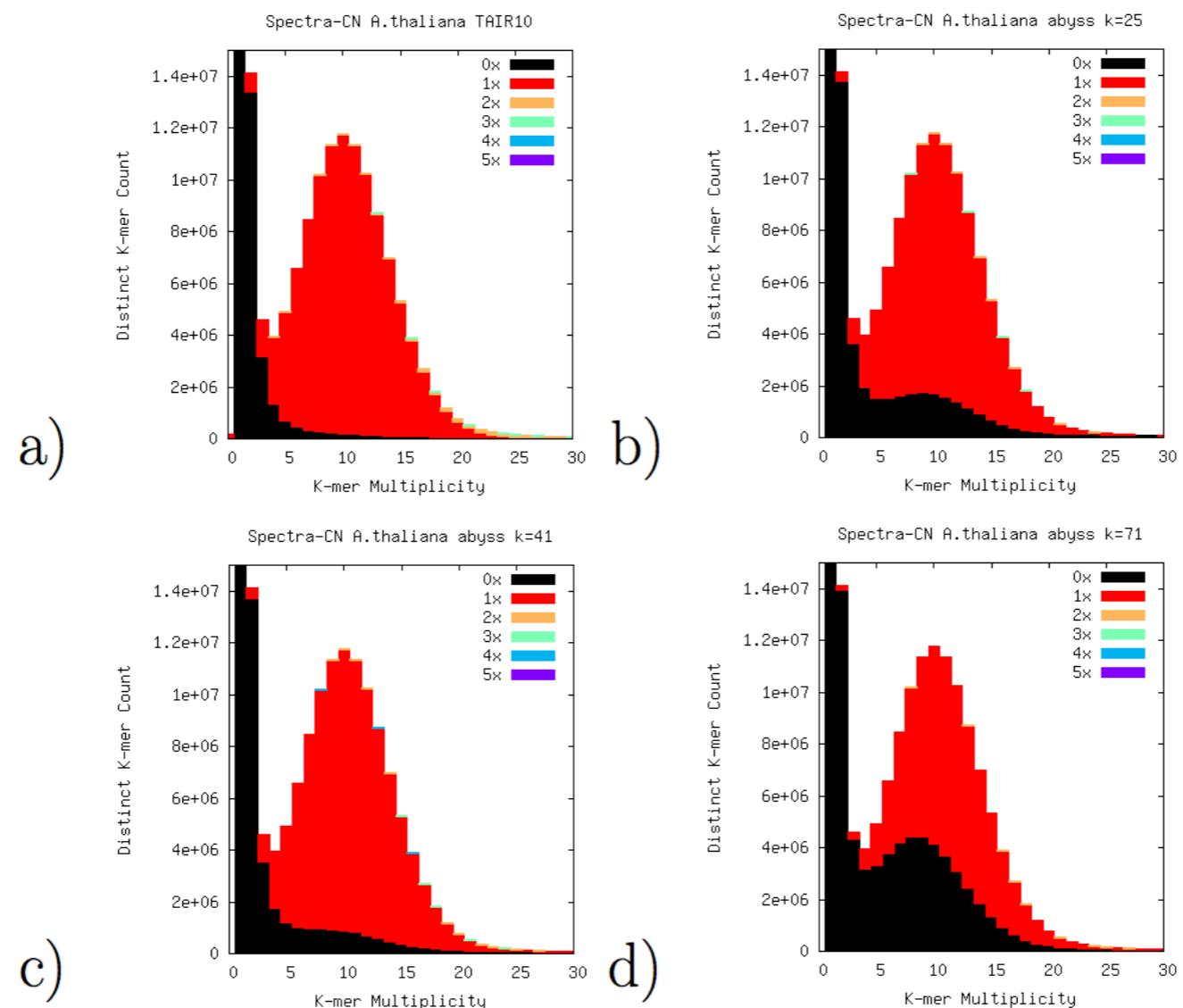
- We typically use KAT to kmer-count.
- You can “read”:
 - Kmer coverage.
 - Genome size.
 - Errors vs. Good kmers.
- Comparing different spectrum (KAT):
 - Is a reference free library assessment.
 - Runs fast.
 - Gives at least a better vs. worse result.



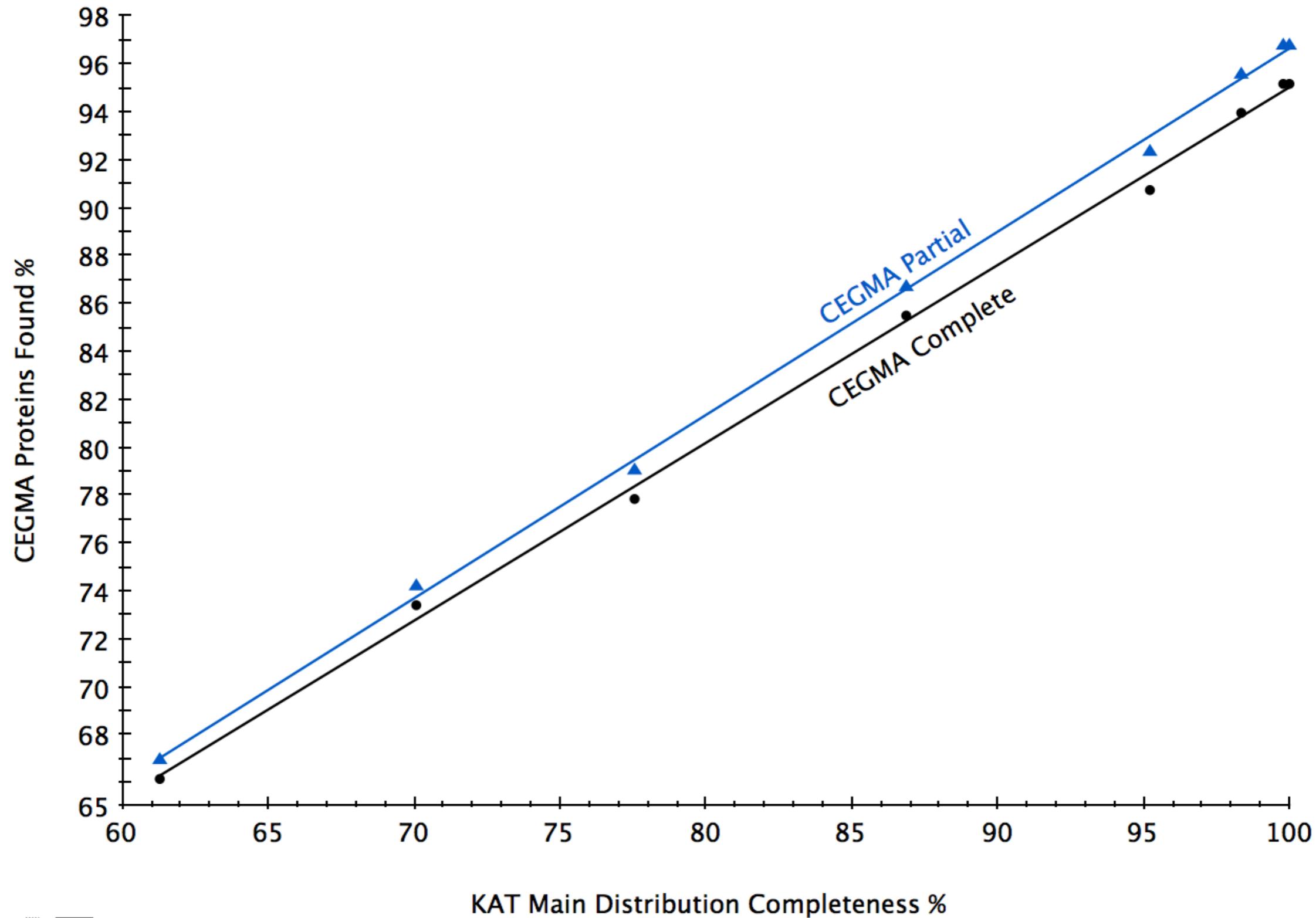
If an assembly is correct, then the original reads should be a plausible sequencing set for the resulting genome model.

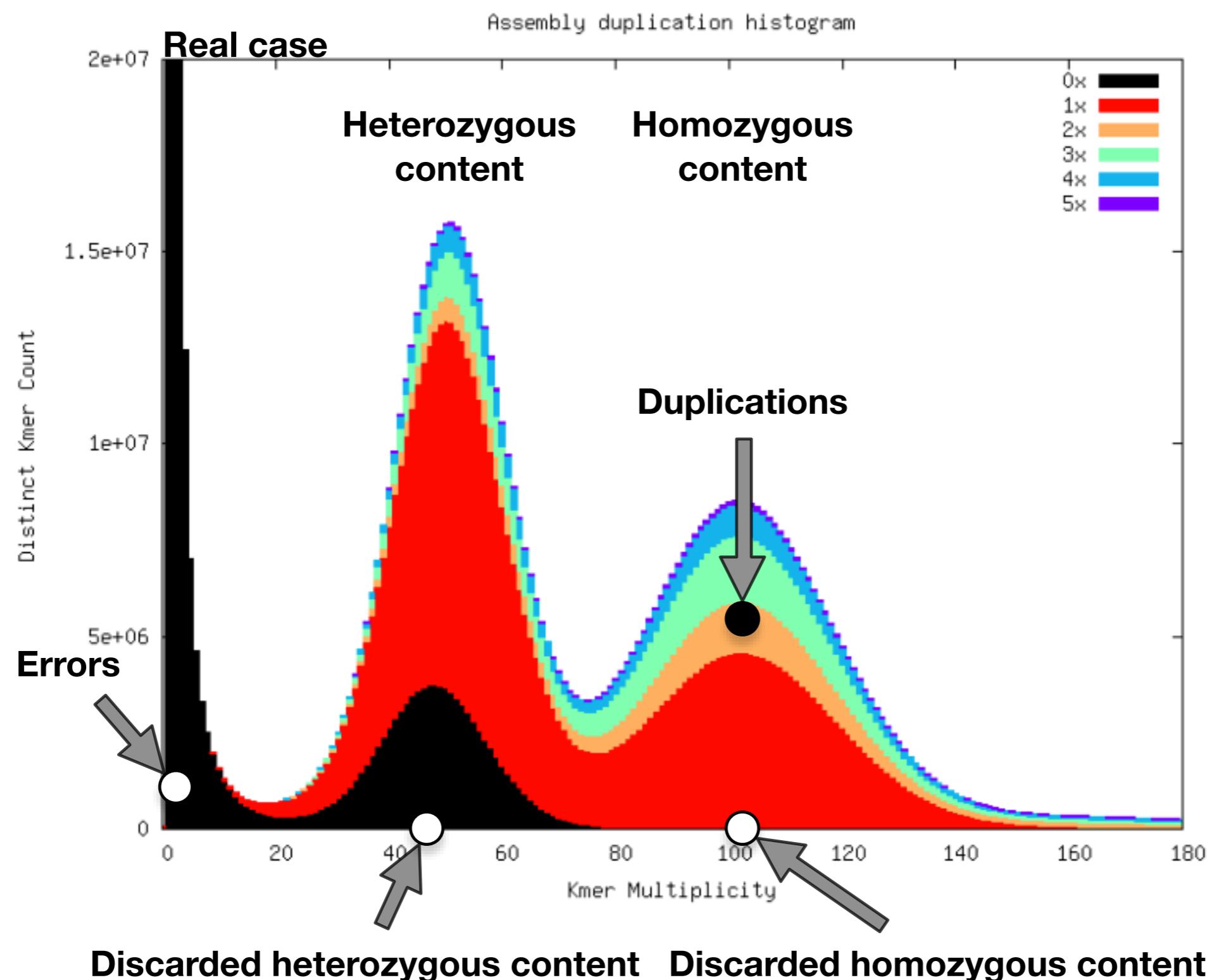
Checking content inclusion using KAT

- Just compare the frequency of kmers in the assembly to the reads spectrum.



KAT vs. CEGMA





Assembly validation

Using biological knowledge to figure out what are...

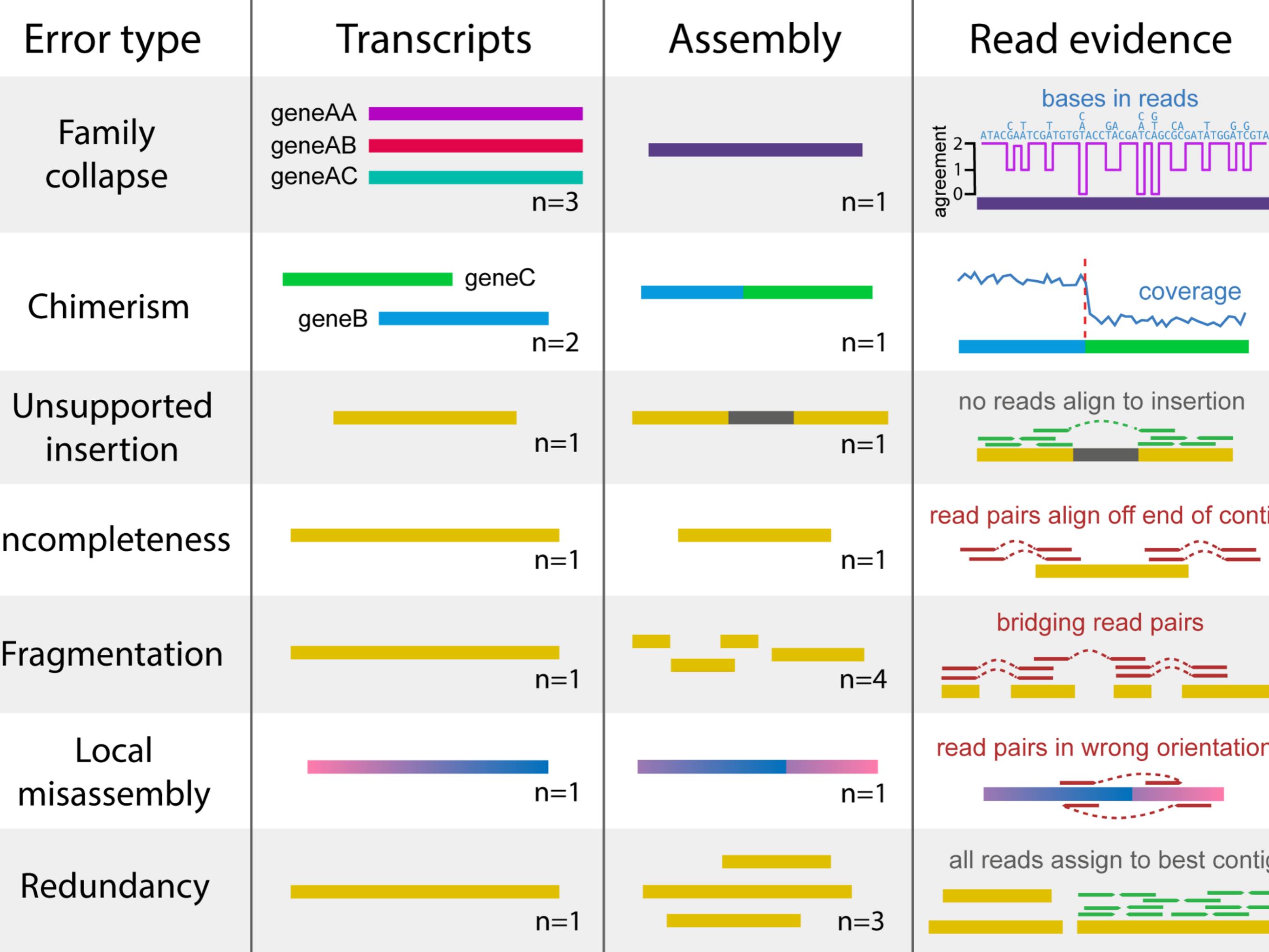
The right *motifs*,
the correct number of times,
in correct order and position.

Direct experimental evidence: the reads

ACTGACTGCCTGTGTGTGTGTGTGTGTGGACTGTTAAA
ACTGACTGC GACTGTTAAA

structure sequence

The right *motifs*,
the correct number of times,
in correct order and position.



Direct experimental evidence: other evidence

- Genome size, ploidy
- GC content
- Symbionts
- Plastids
- ESTs, cDNAs, peptides, genome walking

The right *motifs*,
the correct number of times,
in correct order and position.

Indirect experimental evidence: genomes in general

- Genes! They have structure
- Repeats
- Chromosome macrostructure
 - (circular?, number, telomeres, ...)

The right *motifs*,
the correct number of times,
in correct order and position.

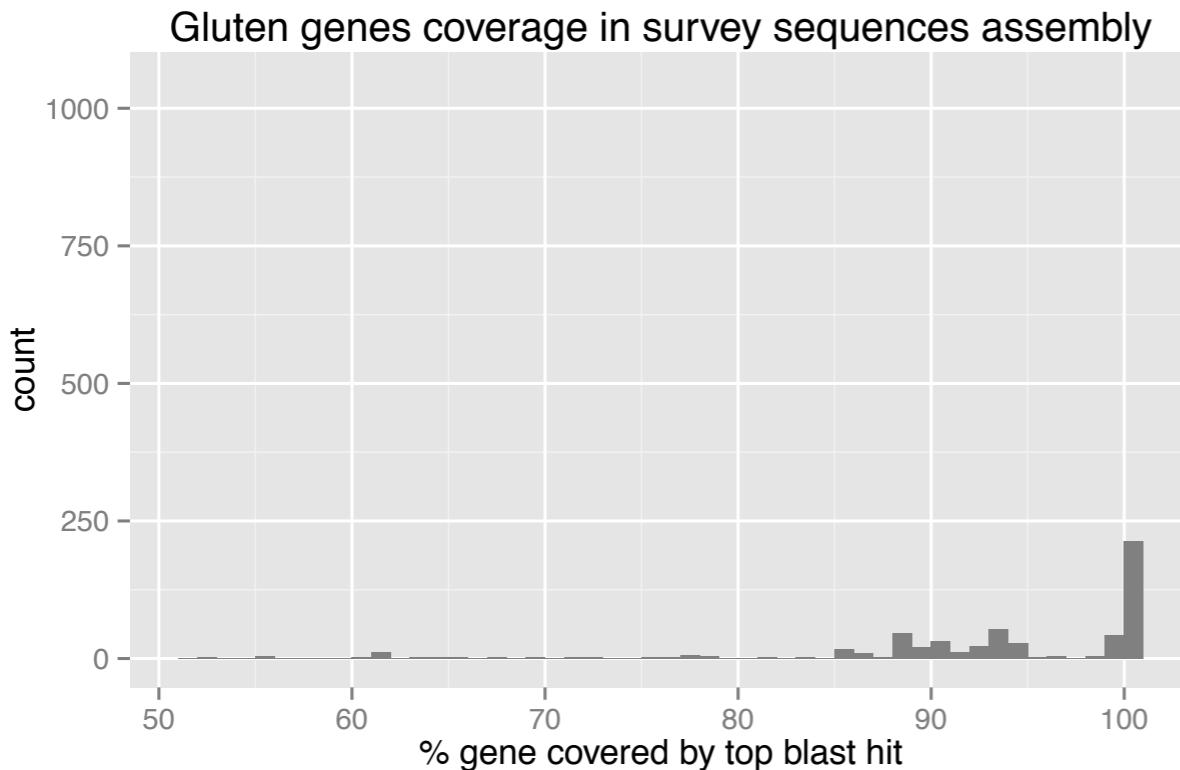
Indirect experimental evidence: other species

- Close relatives: proteins, transcripts, genomes
- Distant relatives: single-copy genes,
phylogeny, HGT

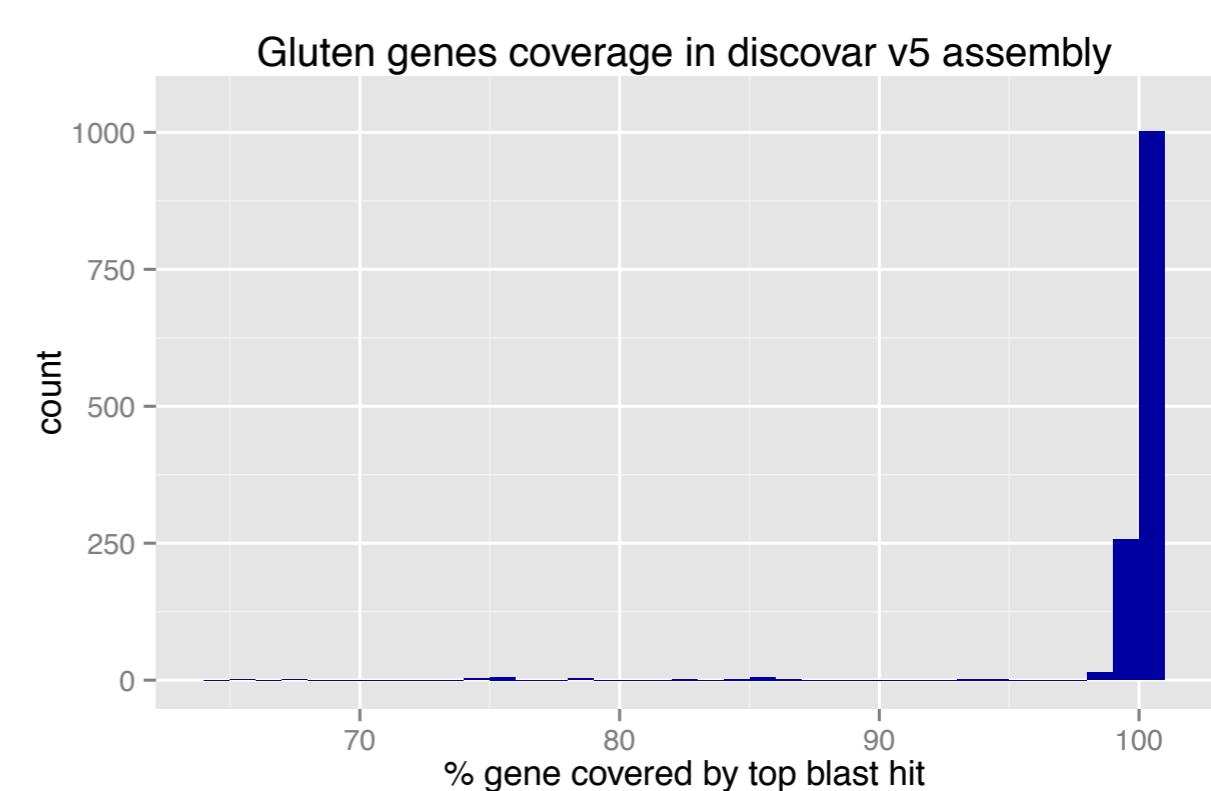
The right *motifs*,
the correct number of times,
in correct order and position.

Examples from the wheat genome

```
>gi|146261042|gb|AB014771.1| HMW glutenin subunit 1By16 [Triticum aestivum]
MAKRLVLFATVVITLVALTAAEGEASRQOCERELQESSLEACRPVVDQOLAGRLPWSTGLCMRCCQQLR
DVSAKCRPVAVSVVRQYEQIVVPPKGGSFYYPGETTPLQQLOVIFWGTSSQTVQGYYPSVSSPQQGPYY
PGQASPOOPGQGQOPGKWDQELGQGQOGYYPTSLHDSGQGQOGYYPSSLQOPGQGQOTGQGQOGYYPTSLQ
QPGQGQIGQGQGQGYYPTSPQHPGQRQPGQGQIGQEQQLGQGRQIGQGQSGQGQGQGYYPTSPQQLGQ
GQOPGQWQOQSGQGQOGYYPTSQQOPGQGQOGQYPASQQOPGQGQOGQYPASQQOPGQGQOGQYPASQQOP
AQGQGQGQYPASQQOPGQGQGHYLASQQOPGQGQRHYPASIQQPGQGQGHYTASIQQPGQGQGHYPAS
SLQVGQGQIGQLGQROOPGQGQTRQGQOLEQGQOPGQGQTRQGQOLEQGQOPGQGQTRQGQOLEQ
GQOPGQGQGYYPTSPQQSGQGQOPGQGQPGQGQGYYSTSLQOPGQGQGQGHYPASLQOPGQGQHPGQRO
QPGQGQOPKQGQGQPGQGQGYYPTSSQQPGQGKQLGQGQGYYPTSPQQPGQGQGQPGQGQGHCPTSPQQ
TGQAAQOPGQGQIGQVQEPGQGQGYYPISLQSGQGQSGQGQGHQOLGQGQSGQGQGQGYDNPYH
VNTEQQTASPKVAKVQQPATQLPIMCRMEEGGDALSTQ
```



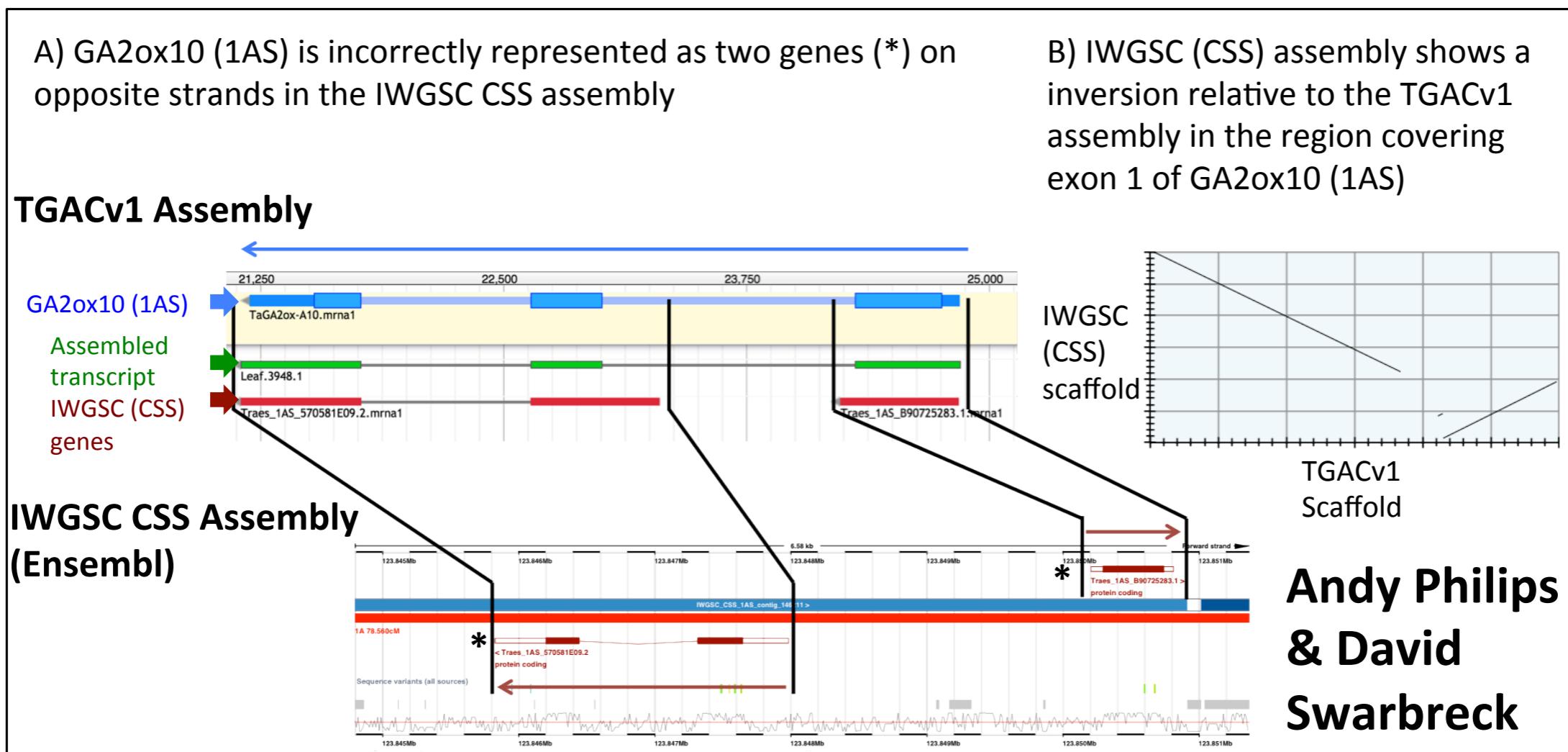
~1000 gluten gene classes are represented:
a-, b-, g- gliadins and high
and low molecular weights glutenins



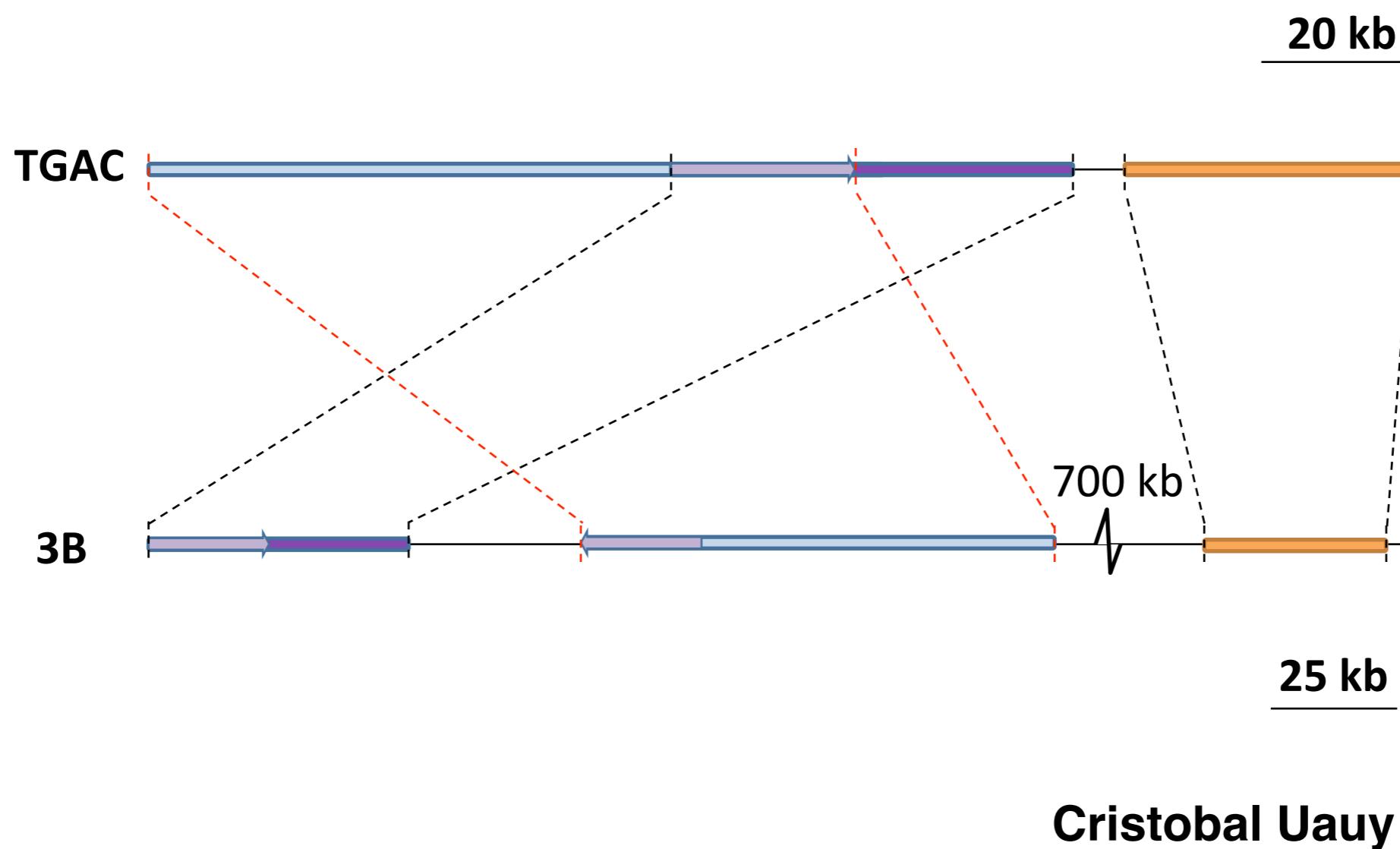
Ksenia Krasileva

Examples from the wheat genome (II)

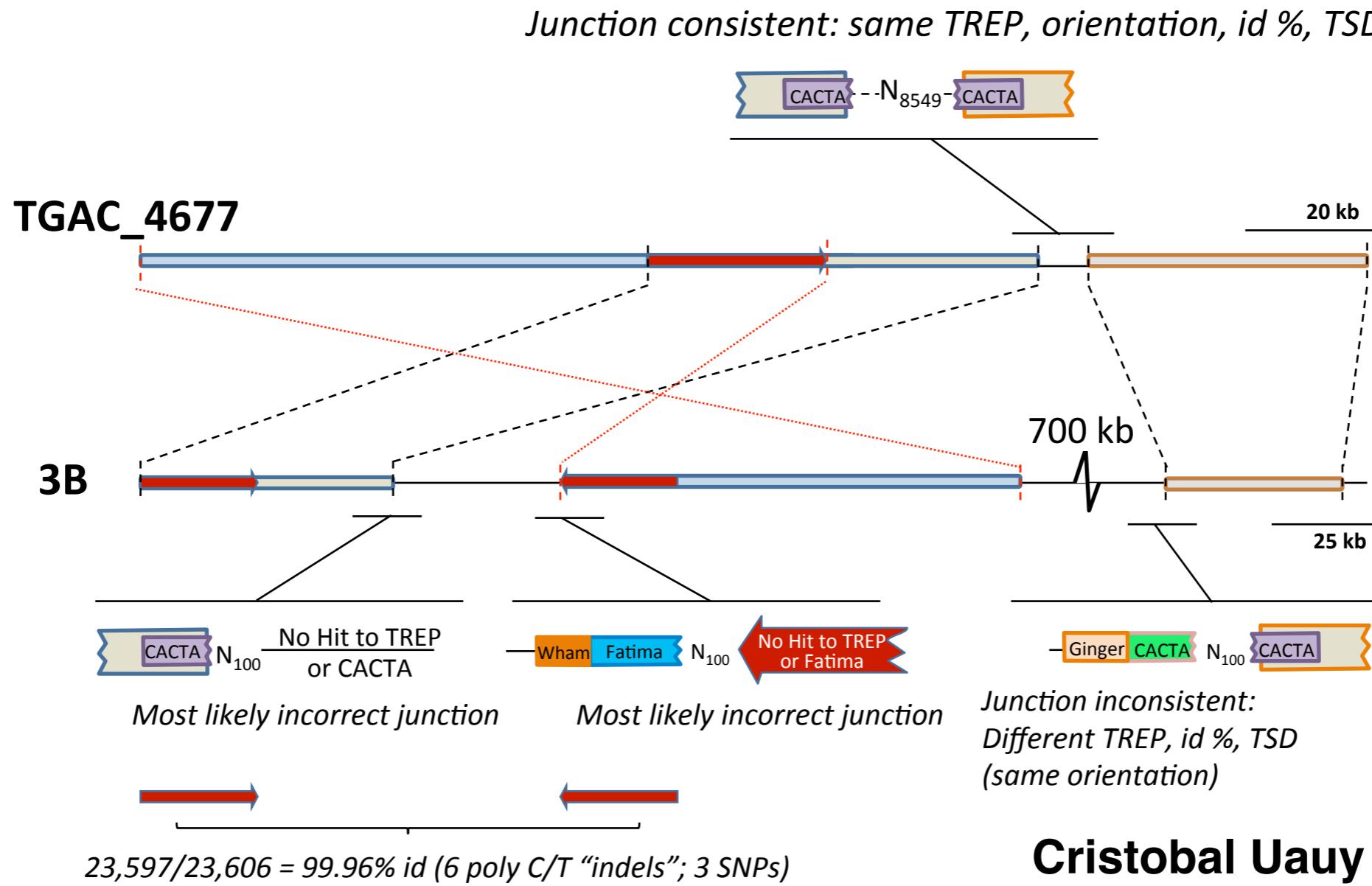
- Gibberellin (GA) pathway plays a central role in plant development.
- 84 genes in the pathway (including homoeologues).
 - **25 genes (30%)** are found as full-length sequences in **CSS**
 - **79 of them (94%)** are full-length in **TGACv1**



Examples from the wheat genome (III)



Examples from the wheat genome (III)



Questions?

