

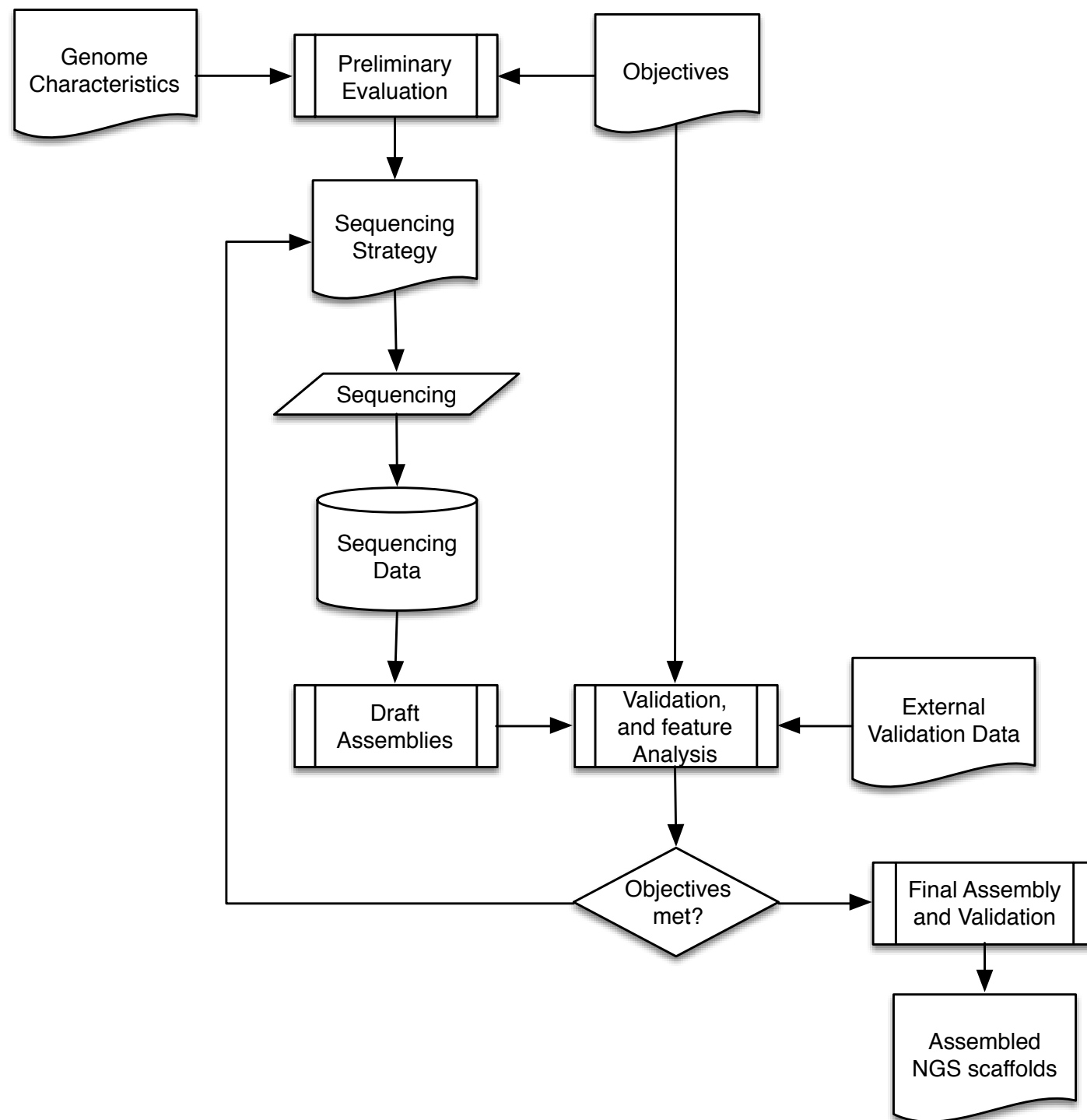
4 - Genome Assembly and Validation (Concepts)

Wednesday afternoon

Bernardo J. Clavijo
Richard Smith-Unna
Gonzalo Garcia



Assembly project workflow | Prior Knowledge

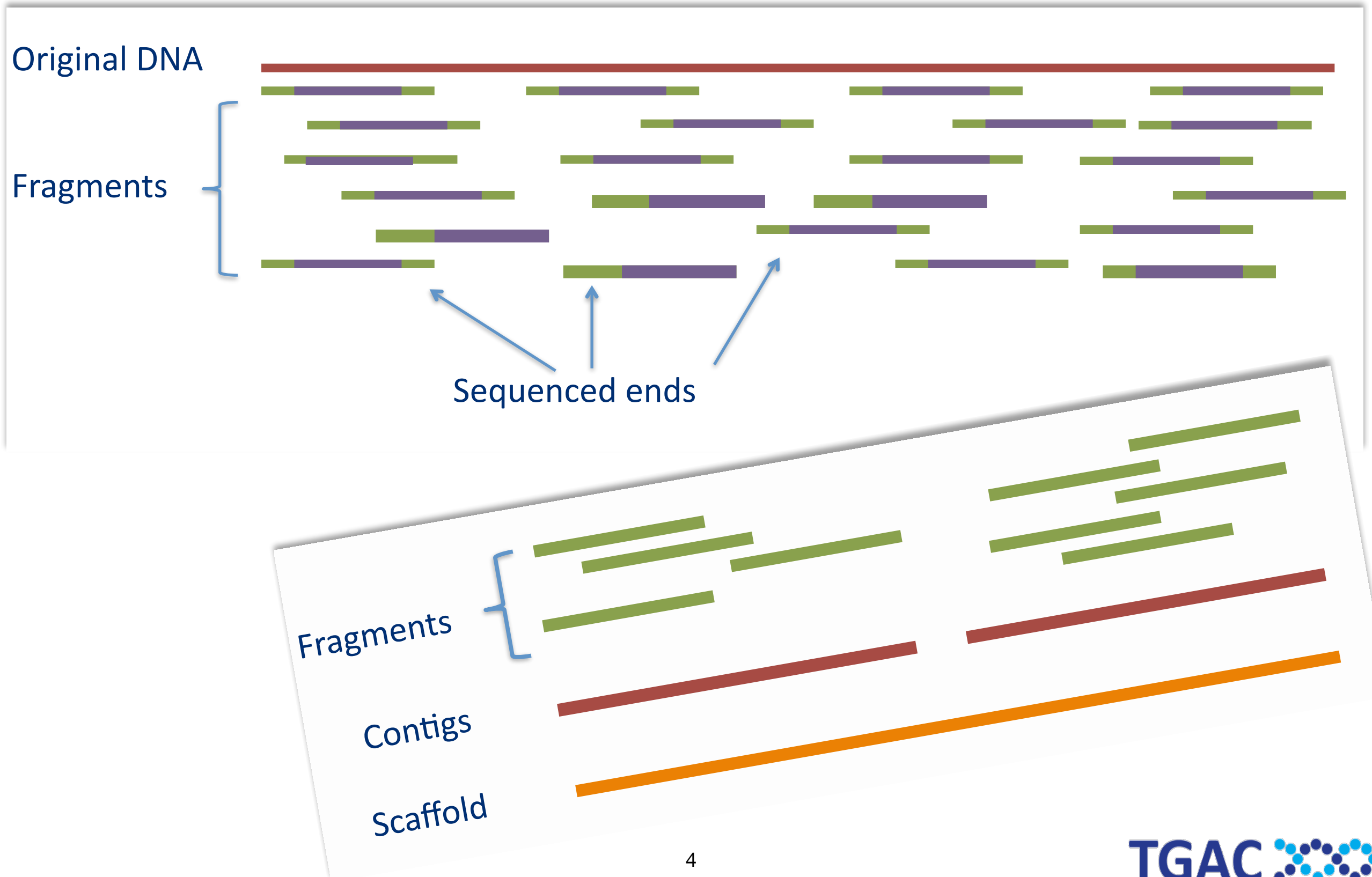


- Kariotype: Genome size, Ploidy
- Heterozygosity
- GC content
- Contaminants / Symbionts
- Data Sets:
 - Close relatives
 - Genes / ESTs / RNAseq / Markers
- Mitochondria
- Chloroplast

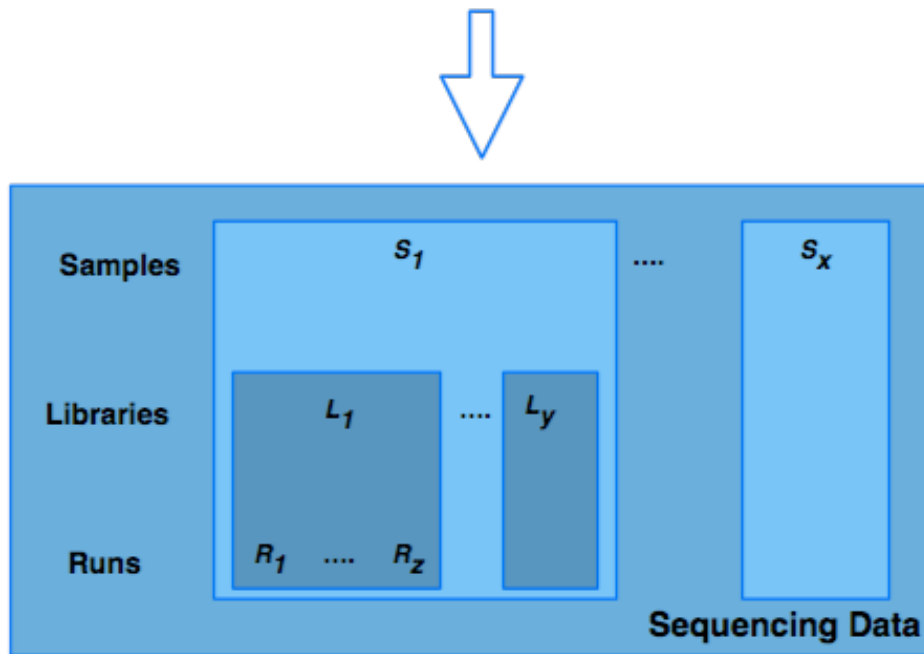
Experiment design (you choose the data!)

- **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)



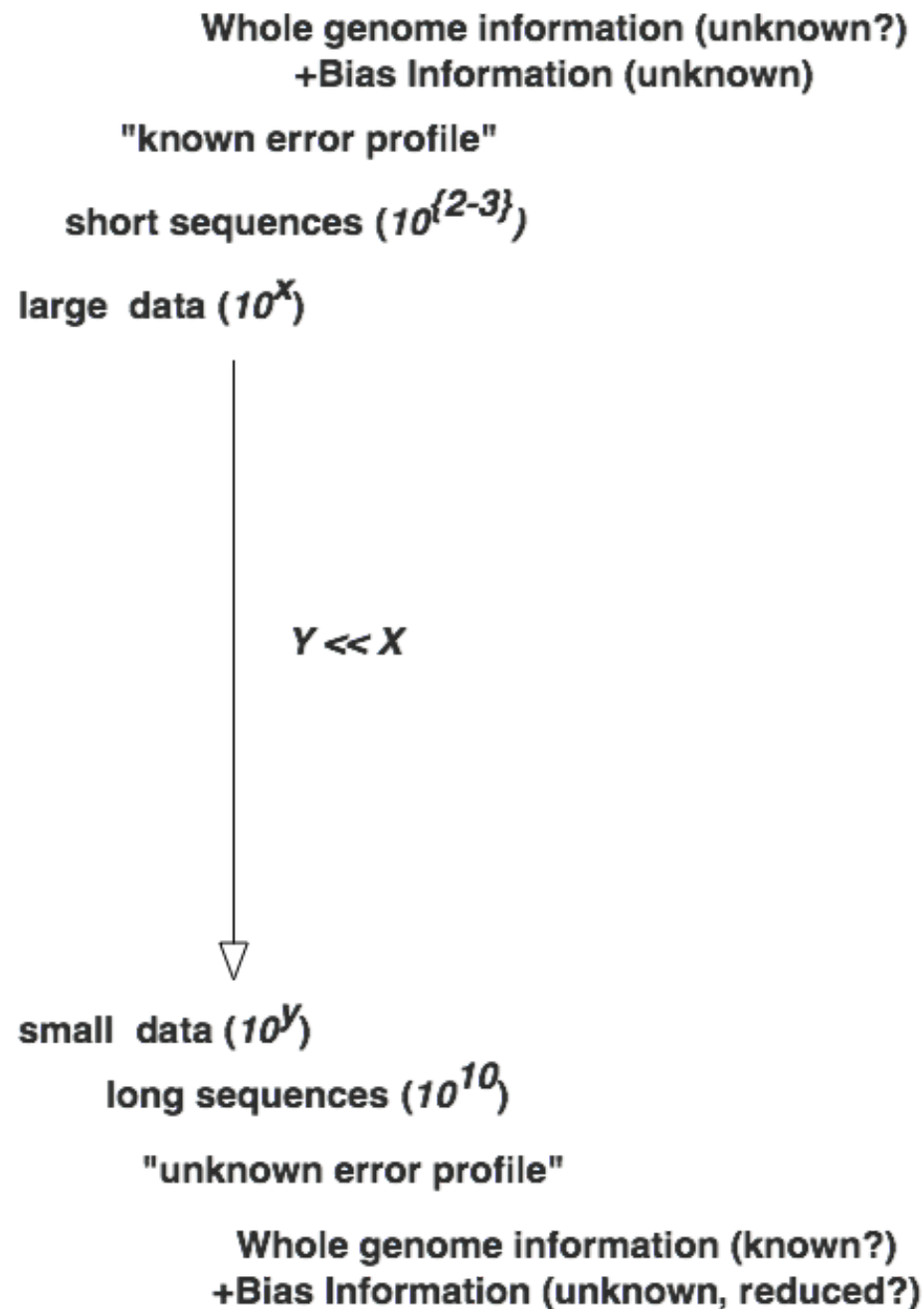
Planning and "informed guesses"



Assemble and Scaffold

Scaffolds & Contigs

Validate and release



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.

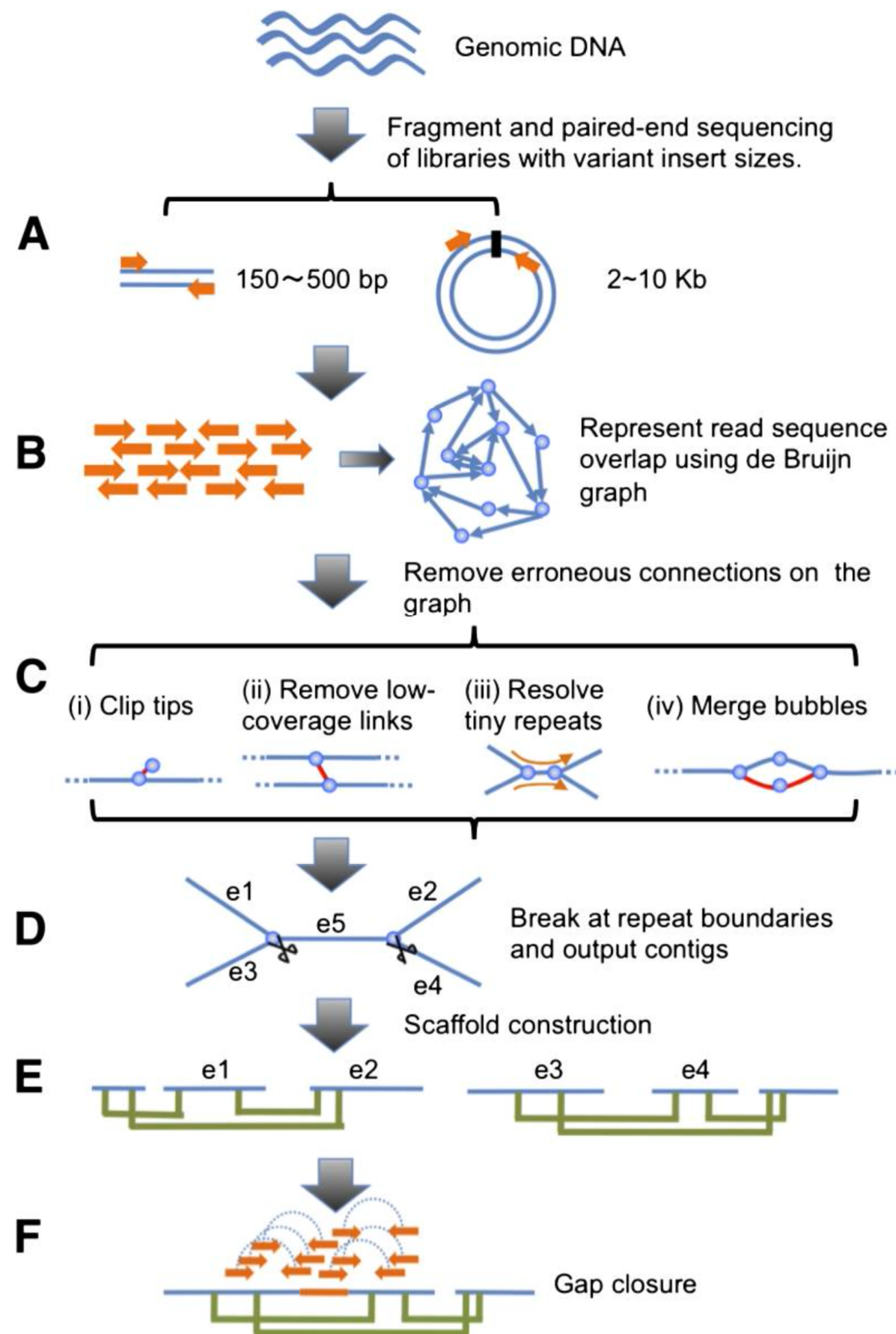
A correct assembly has:

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

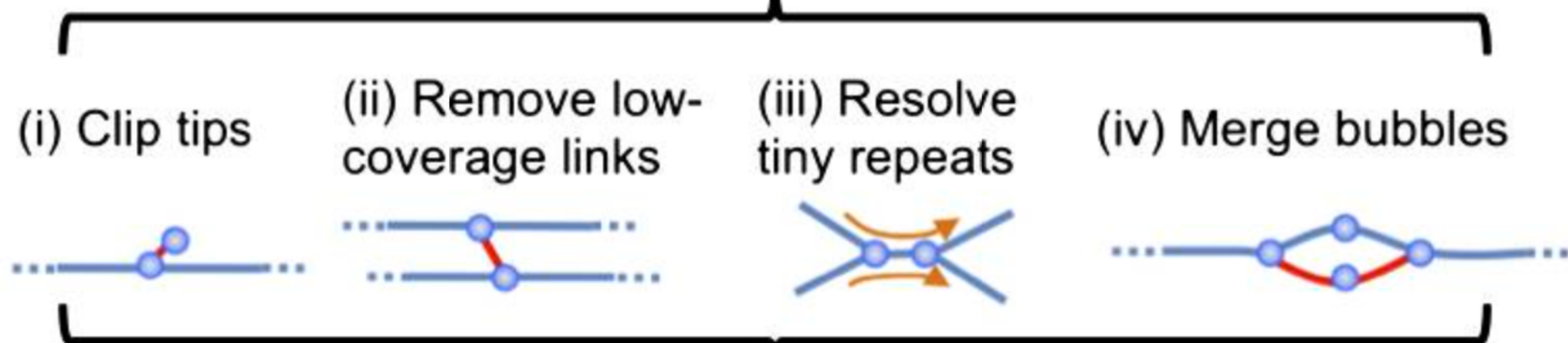


B

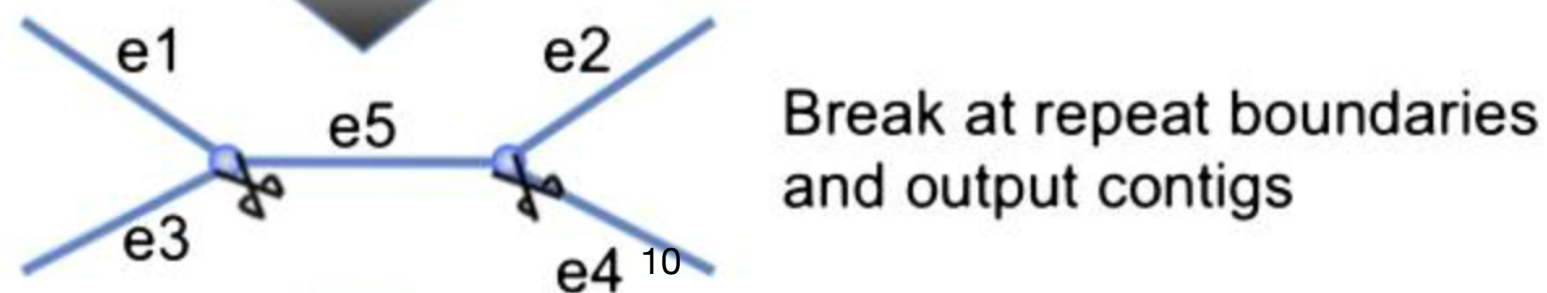


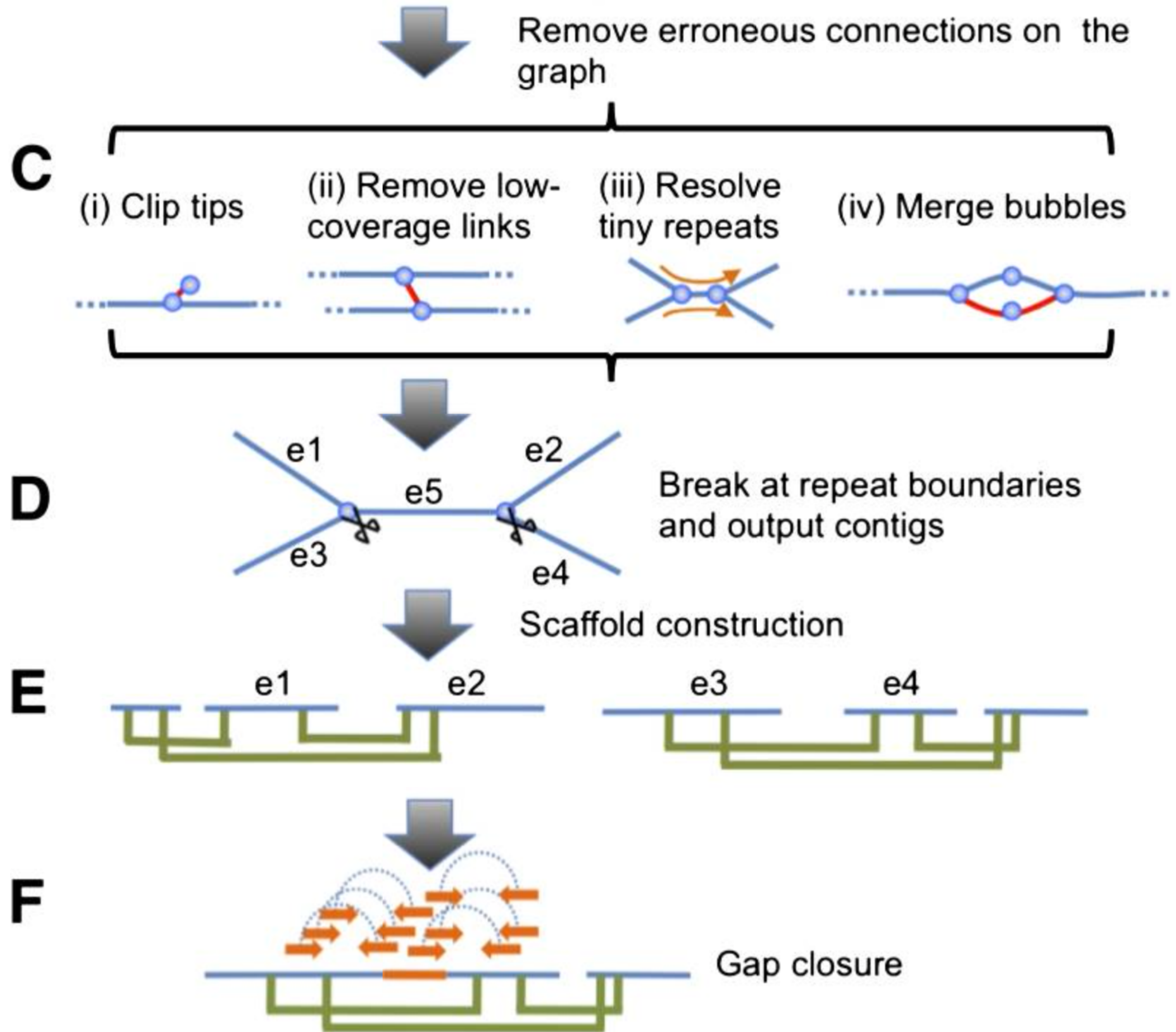
Remove erroneous connections on the graph

C



D





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



ACTGACTGCCTGTGTGTGTGTGTGTGTGTGTGGACTGTTAAA



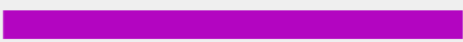


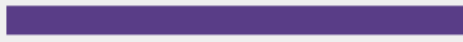




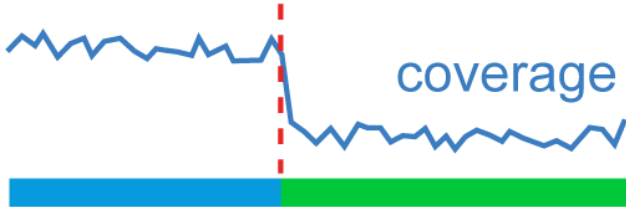
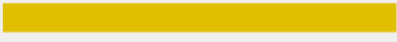
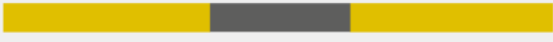










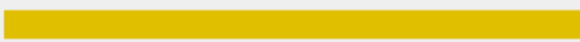


ACTGACTGC



GACTGTTAAA

structure
sequence

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

Error type	Transcripts	Assembly	Read evidence
Family collapse	geneAA  geneAB  geneAC  n=3	 n=1	
Chimerism	 geneC geneB  n=2	 n=1	
Unsupported insertion	 n=1	 n=1	no reads align to insertion 
Incompleteness	 n=1	 n=1	read pairs align off end of contig 
Fragmentation	 n=1	 n=4	bridging read pairs 
Local misassembly	 n=1	 n=1	read pairs in wrong orientation 
Redundancy	 n=1	 n=3	all reads assign to best contig 

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.**

Questions?

