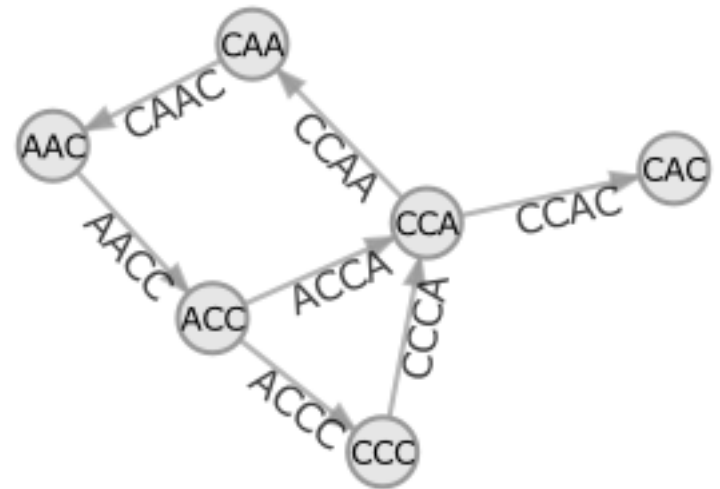


## *De novo* transcriptome assembly

How it works, limitations and how to tell if your assembly is any good

# *De novo*

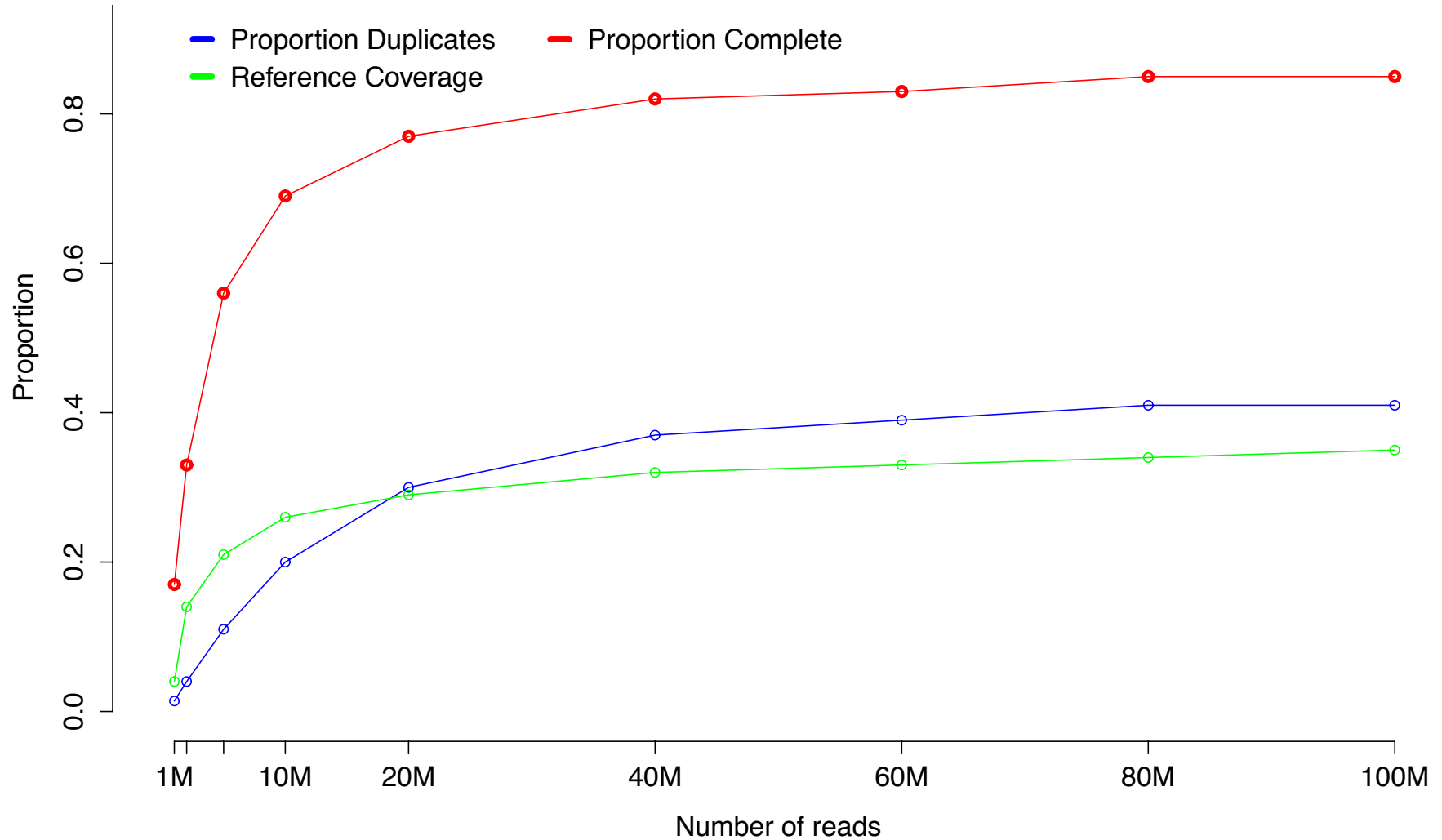
- Latin expression meaning "from the beginning," "afresh," "anew," "beginning again."
- De novo, a term for any method that makes predictions about biological features using only a computational model without extrinsic comparison to existing data
- No reference genome, then you must assemble your reads into genes



# Transcriptome vs Genome Assembly

- Transcriptome is easier
  - Smaller total volume of bases
  - Less low complexity regions
- Transcriptome is harder
  - Alternative splicing
  - Expression variability between tissues/cells
    - Difficult/impossible to fully sample all transcripts
    - Exponentially distributed coverage levels
    - Not the same in every cell
- Because of the very unique properties of transcriptome assemblies, it is important to use an assembler meant for transcriptomes (not genomes)

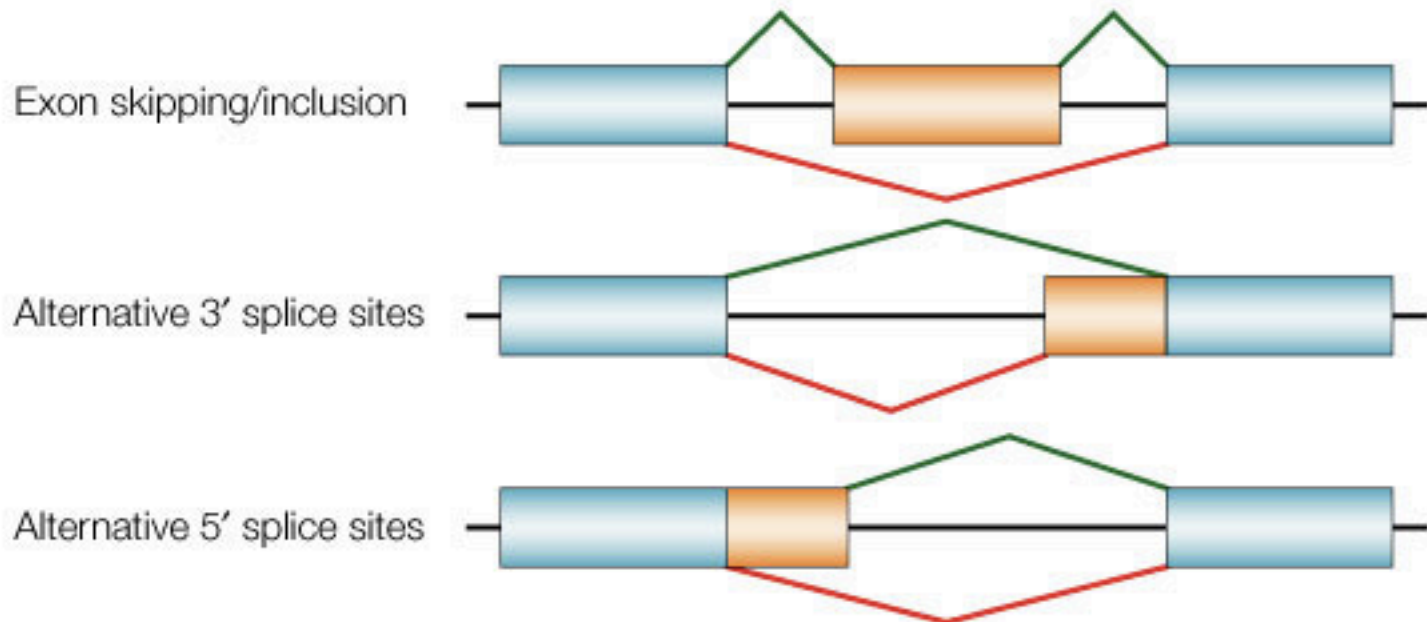
# How many reads?



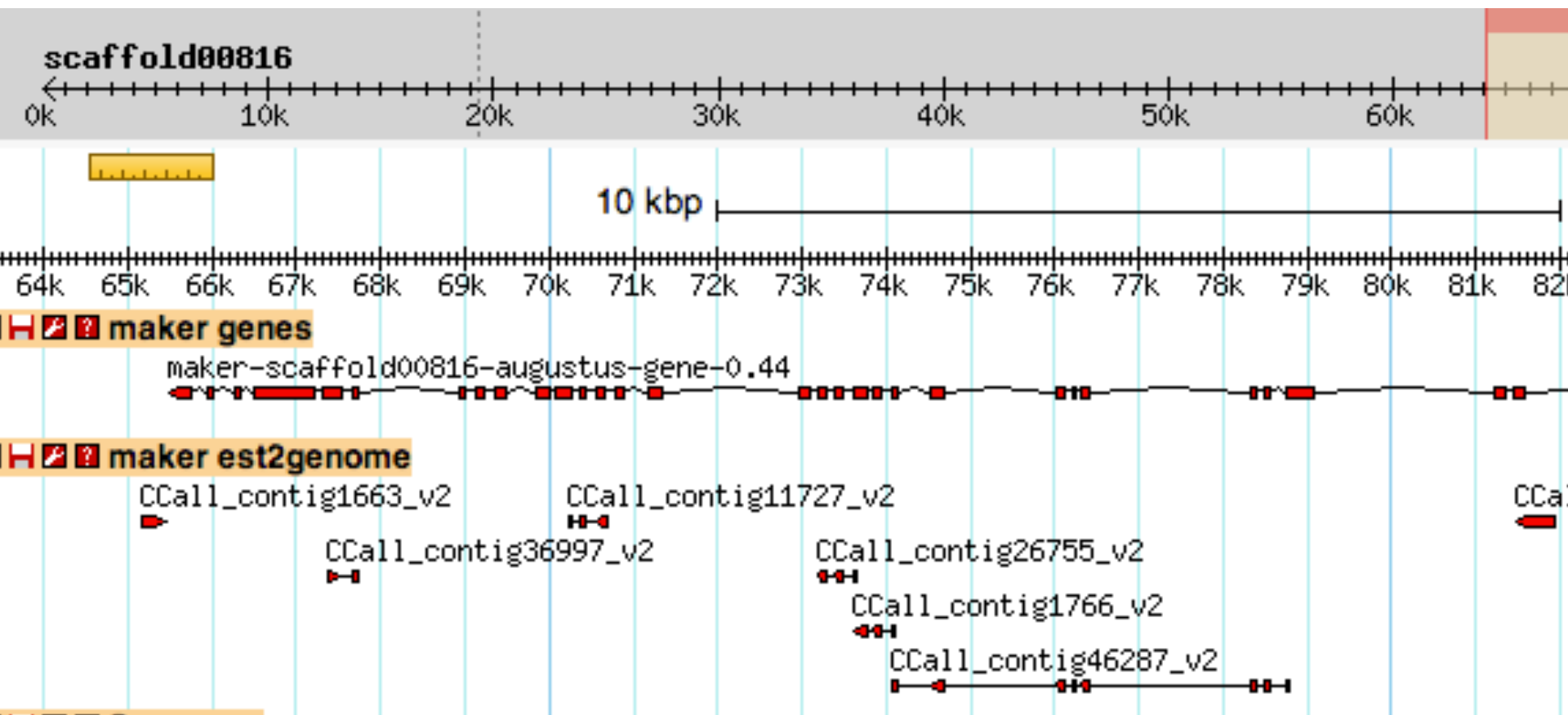
How many individuals? How many  
tissues?

# Problems with *de novo* assemblies

- Results
  - Highly fragmented assemblies
  - Chimeras (can be biological, experimental or computational)
  - Paralogs, alleles and alternative splicing variants combined or fragmented

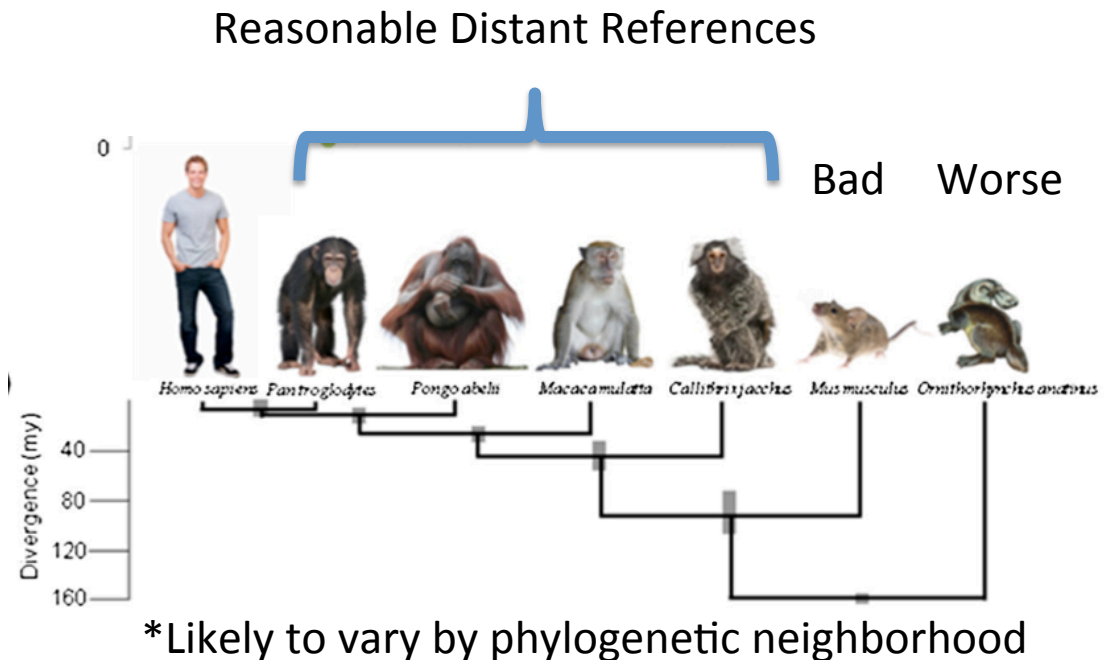


# Chestnut



# Avoid an entirely *de novo* assembly - Use a distant or close relative

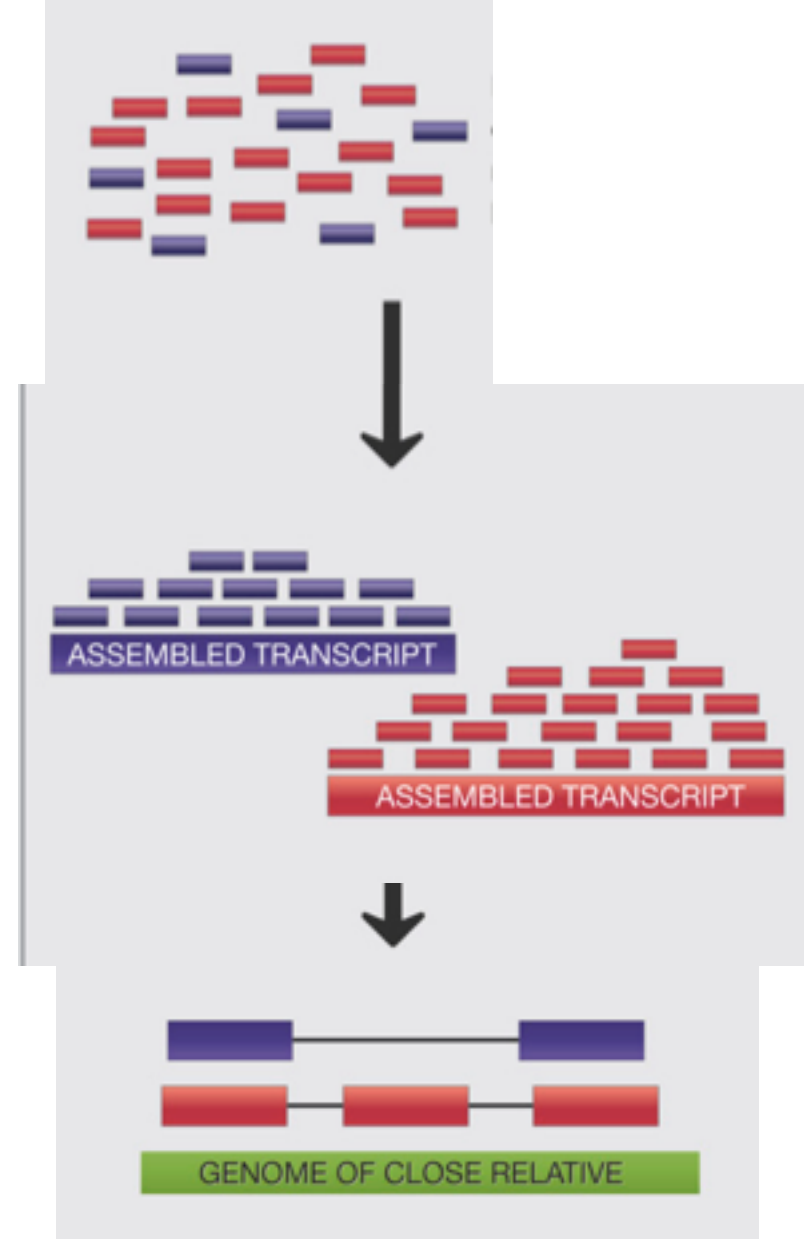
- Is there a close relative with a sequenced genome?
- How close is close enough?
  - Align then assemble
  - Assemble then align





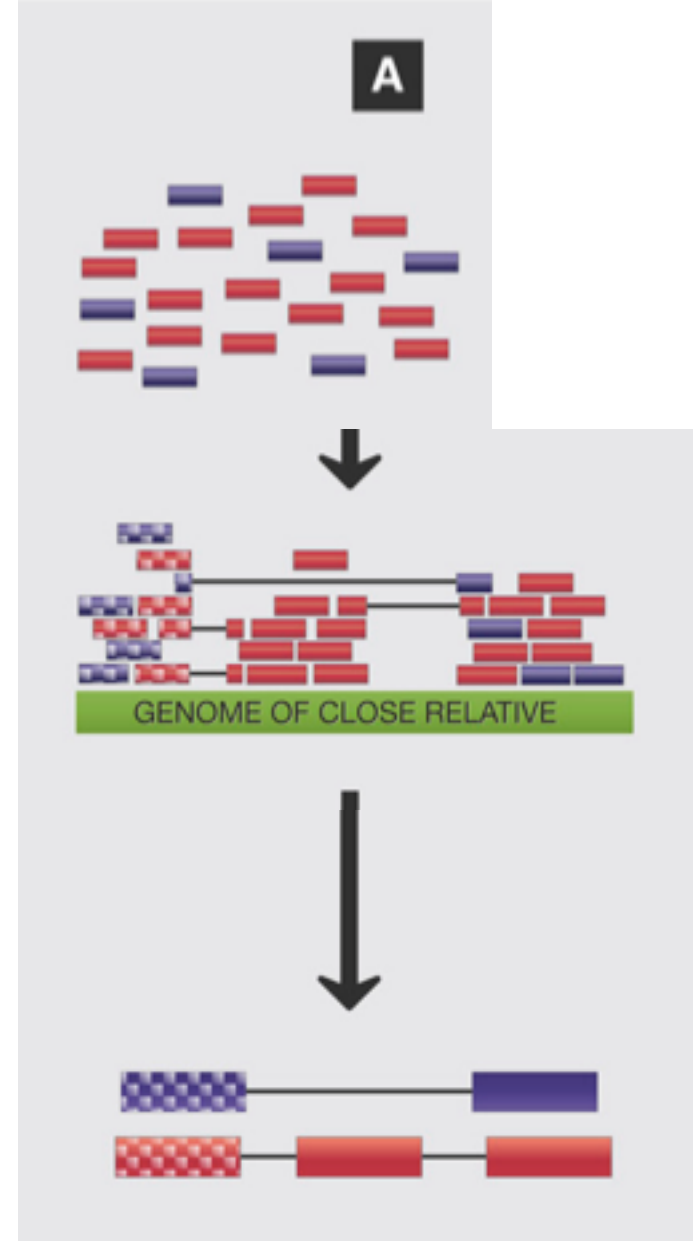
# Assemble then align

- First, assemble
- Next, align to a close relative
- Main Problems:
  - Fragmented assemblies – gene pieces are scattered in a different consensus pieces
  - More difficult to sort out gene family members
- Main Advantages:
  - Alignment to a close relative can identify exon/exon boundaries (sort out alternative splicing)
  - Less bias – can discover novel gene sequences



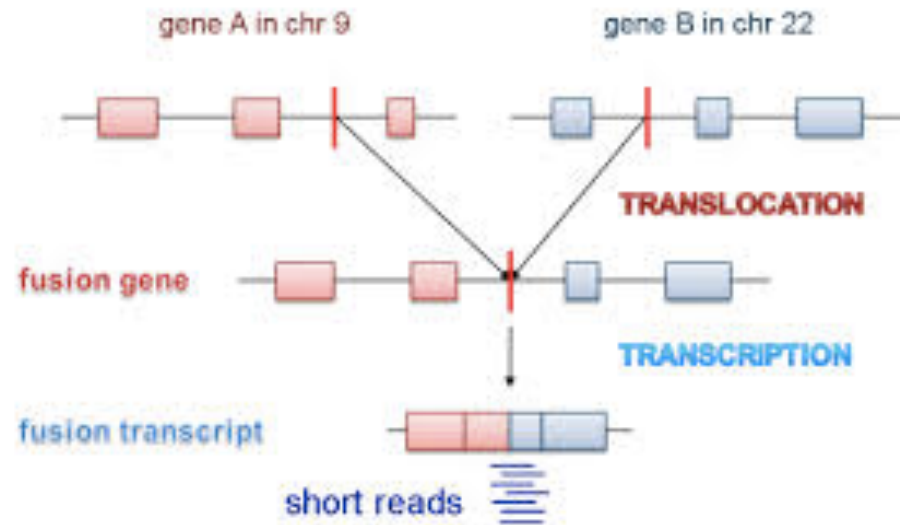
# Align then assemble

- First, map reads to (distant) reference
- Next, do local assemblies for each gene
- Main Problems
  - Read alignment may be poor due to lack of sequence similarity
  - Gene family expansion/contraction
- Main Advantage
  - Transcript assembly is less likely to be fragmented
  - Even where it is fragmented, you can identify all the fragments that originate from a single locus



# *De novo* transcriptome assemblies

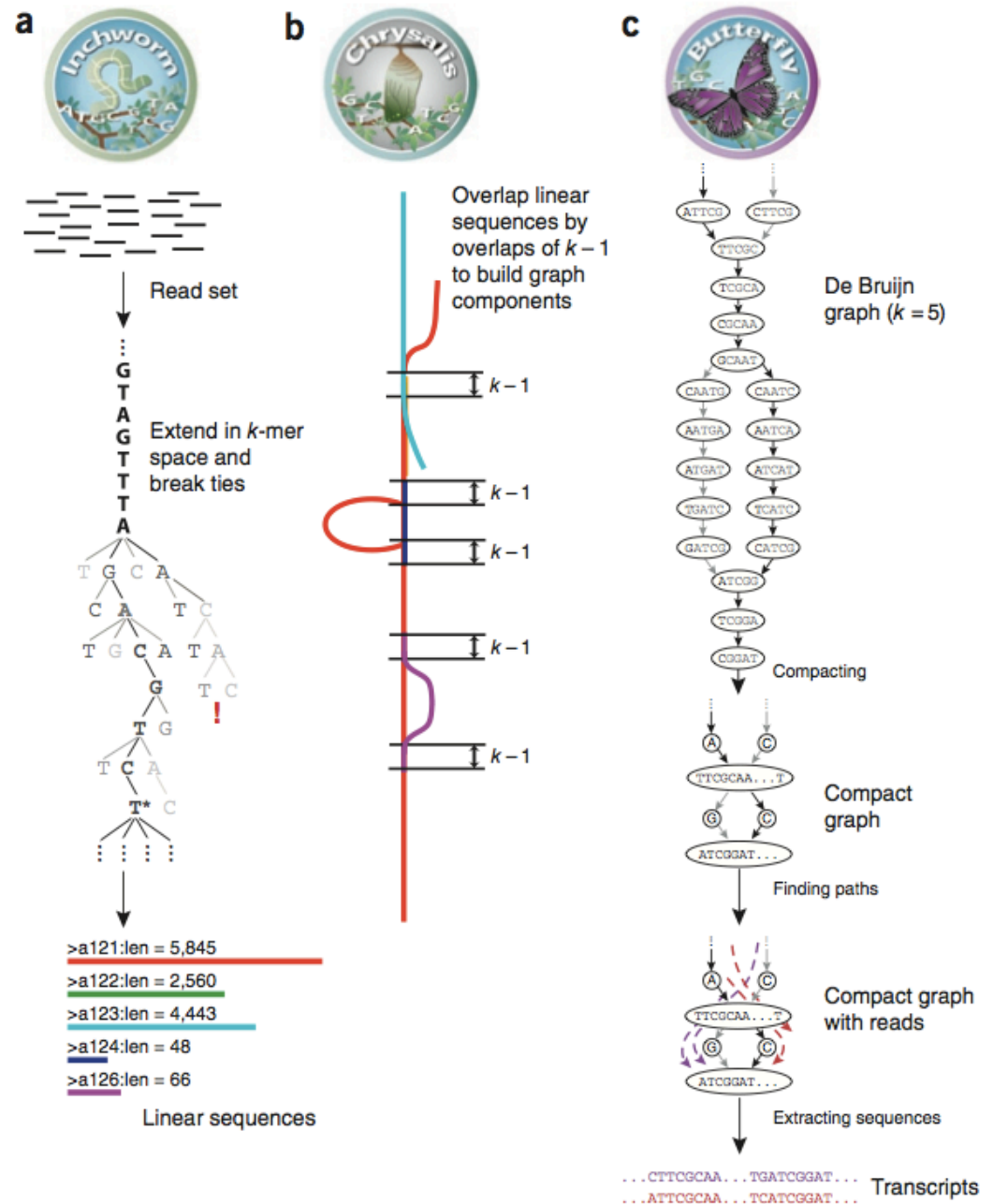
- Completely reference free
- What are they useful for?
  - Transcriptome characterization
  - Enabling proteomics experiments
  - Candidate gene discovery
  - Marker discovery/development
  - Cancer or other tissues where fusion events are important
  - Metatranscriptomics – surveying microbiota



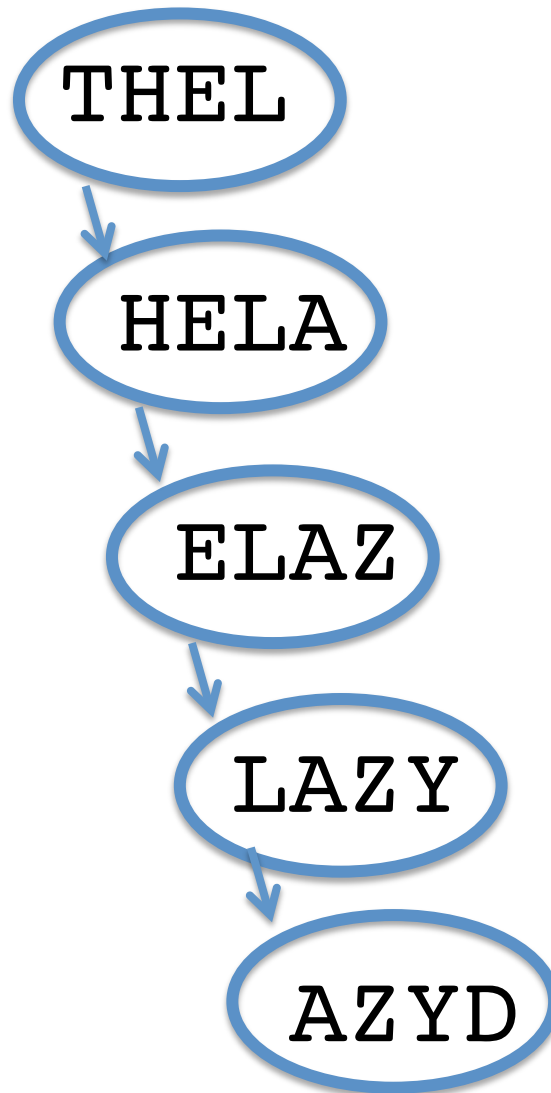
# Trinity strategy

## Three stages

1. Inchworm
2. Craylis
3. Butterfly



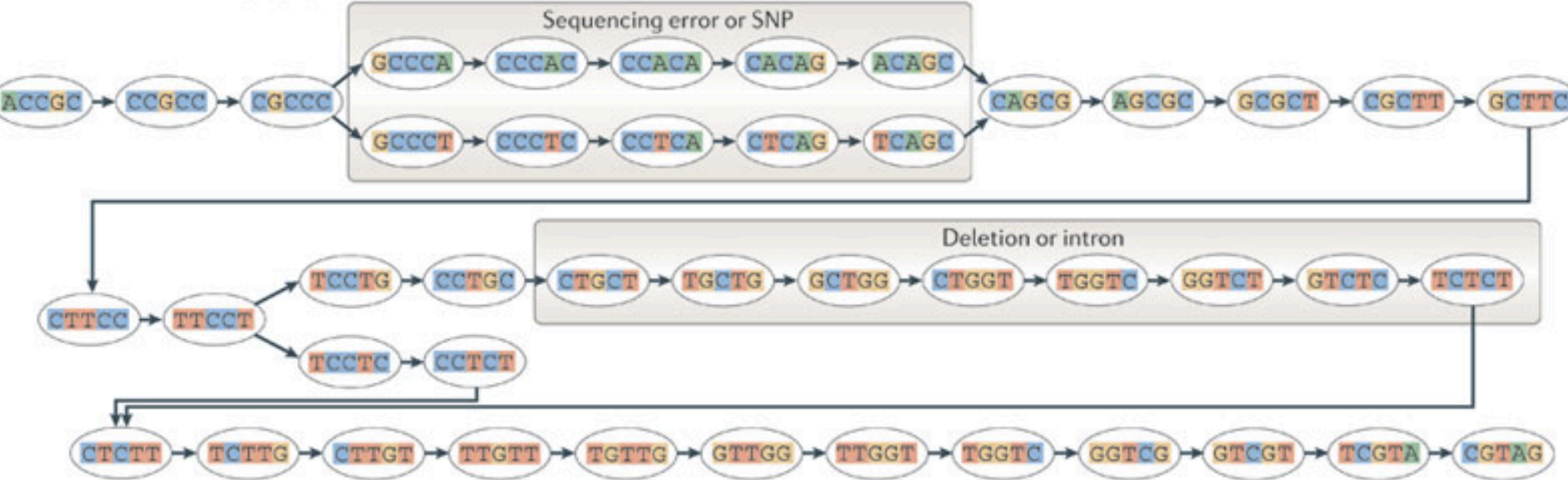
THELAZYBROWND OG



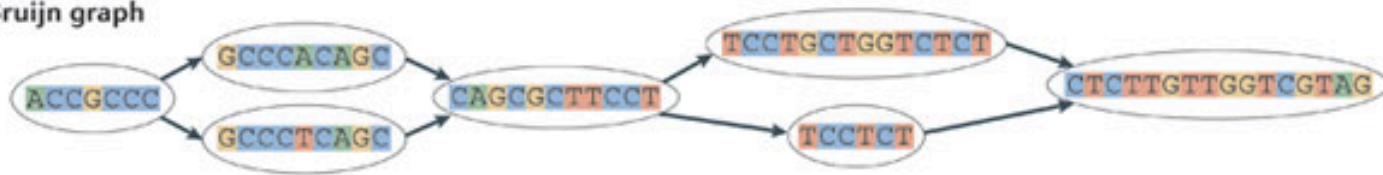
## De Bruijn Graph

Directed acyclic graph  
of kmers

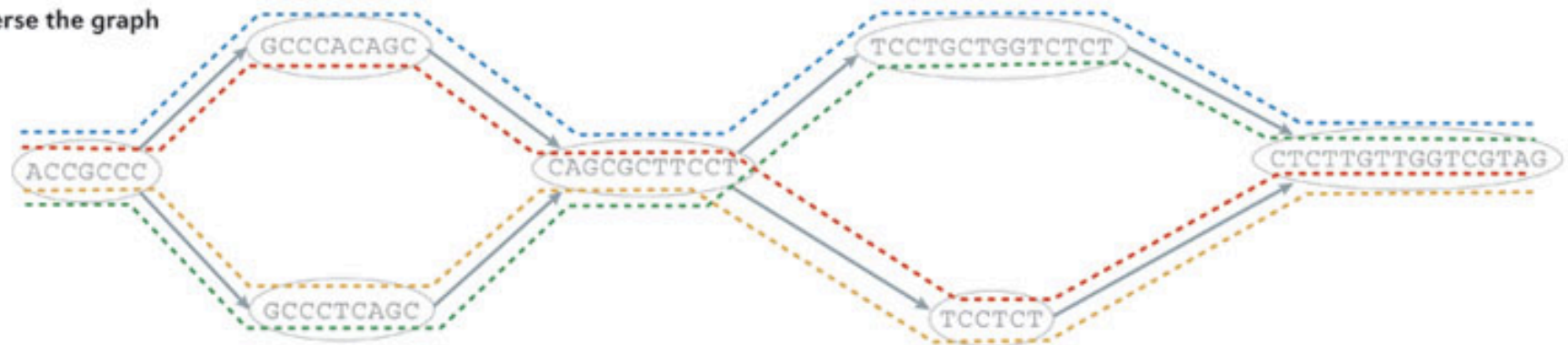
Generate the De Bruijn graph



**c Collapse the De Bruijn graph**



**d Traverse the graph**

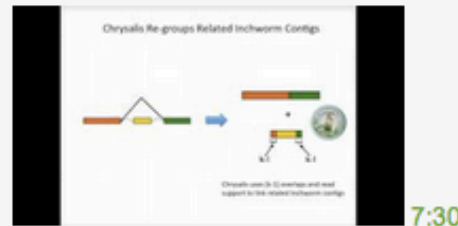


**e Assembled isoforms**

..... ACCGCCACAGCGCTTCCTGCTGGTCTCTTGTTGGTCGTAG  
 ..... ACCGCCACAGCGCTTCCT ..... CTTGTTGGTCGTAG  
 ..... ACCGCCCTCAGCGCTTCCT ..... CTTGTTGGTCGTAG  
 ..... ACCGCCCTCAGCGCTTCCTGCTGGTCTCTTGTTGGTCGTAG

## A Collection of new RNA-Seq Videos from The Broad Institute

Posted by: RNA-Seq Blog Administrator In Presentations ⌚ October 10, 2013 👁 1,134 Views



7:30



5:20

## Videos!

<http://www.rna-seqblog.com/a-collection-of-new-rna-seq-videos-from-the-broad-institute/>

## BroadE: Trinity – How it works

## BroadE: The General Approach to De novo RNA-Seq Assembly Using De Bruijn Graphs



5:38



3:00

## BroadE: Introduction to De Novo RNA-Seq Assembly using Trinity

## BroadE: Strand-specific RNA-Seq is Preferred



# Trinity output – deciphering the naming

- An example Fasta entry for one of the transcripts is formatted like so:

>c115\_g5\_i1 len=247 path=[31015:0-148  
23018:149-246]

Component –  
a collection of  
contigs that are  
likely to be  
derived from  
alternative splice  
forms or closely  
related paralogs

Gene – best  
guess at an  
individual locus

Isoform –  
alternative splicing  
events and alleles

These divisions are guesses only!

## II. Improving the assembly and checking quality

# How to figure out if your assembly is good

- BUSCO
  - Benchmarking Universal Single-Copy Orthologs
  - based on evolutionarily informed expectations of gene content from near-universal single-copy orthologs selected from OrthoDB.
  - Use to assess completeness of transcriptome
  - <http://busco.ezlab.org/>



# How to figure out if your assembly is good

- Map reads back and see what % are captured in the assembly
- Transrate
  - analyses a transcriptome assembly in three key ways:
    - by inspecting the contig sequences
    - by mapping reads to the contigs and inspecting the alignments
    - by aligning the contigs against proteins or transcripts from a related species and inspecting the alignments



### III. De novo transcriptome sequencing – after assembly

# ORF Finding - TransDecoder

- Searches all frames for ORFs, start codons and stop codons
- Maximizes length and log-likelihood score of ORF
- a single transcript can report multiple ORFs (allowing for operons, chimeras, etc).

