

Transcriptome and isoform reconstruction with short reads

Tangled up in reads

Enabler for Life Sciences

Topics of this lecture

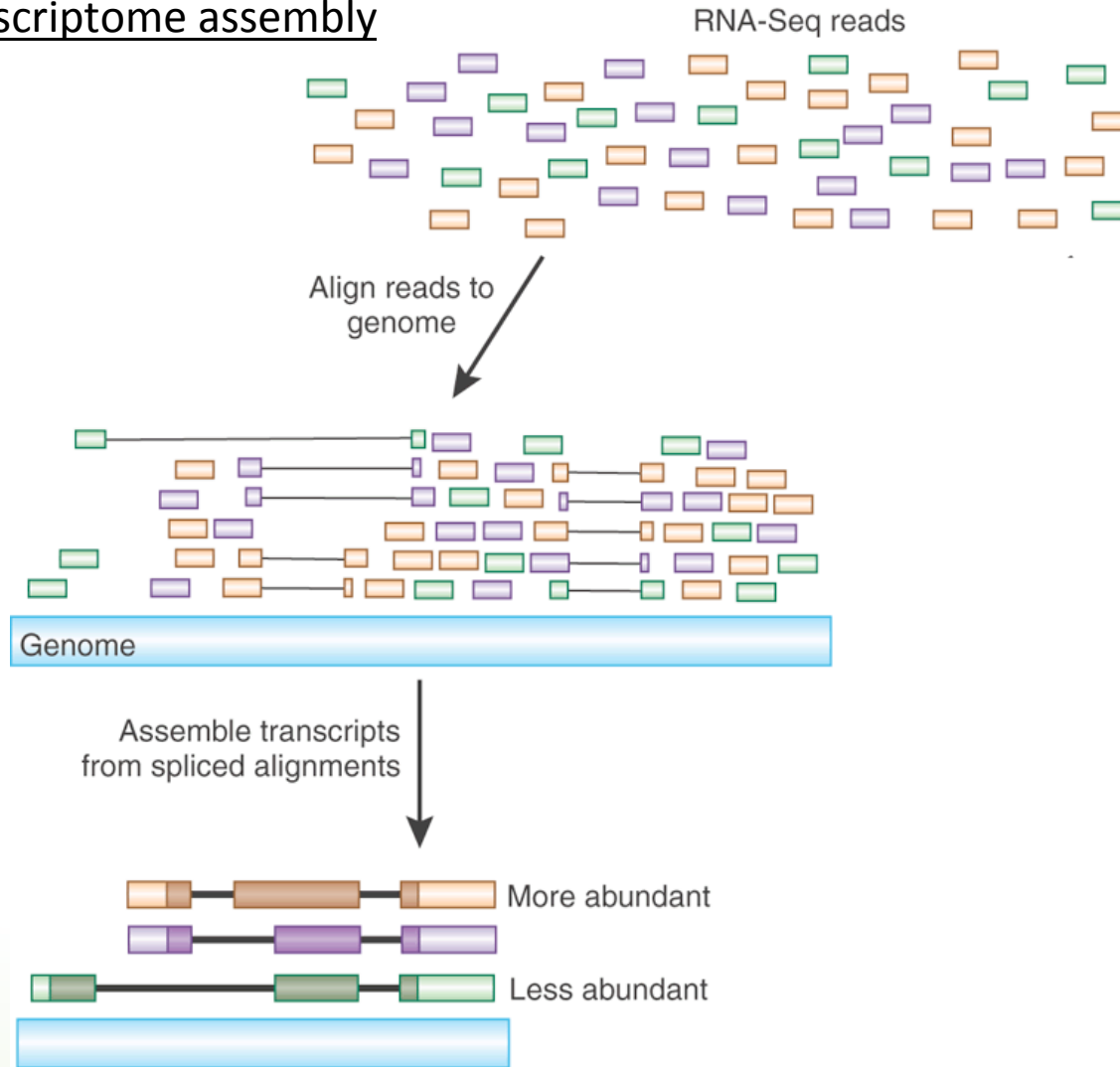
- Mapping-based reconstruction methods
 - Case study: The domestic dog
- De-novo reconstruction method
 - Trinity

Transcriptome assembly



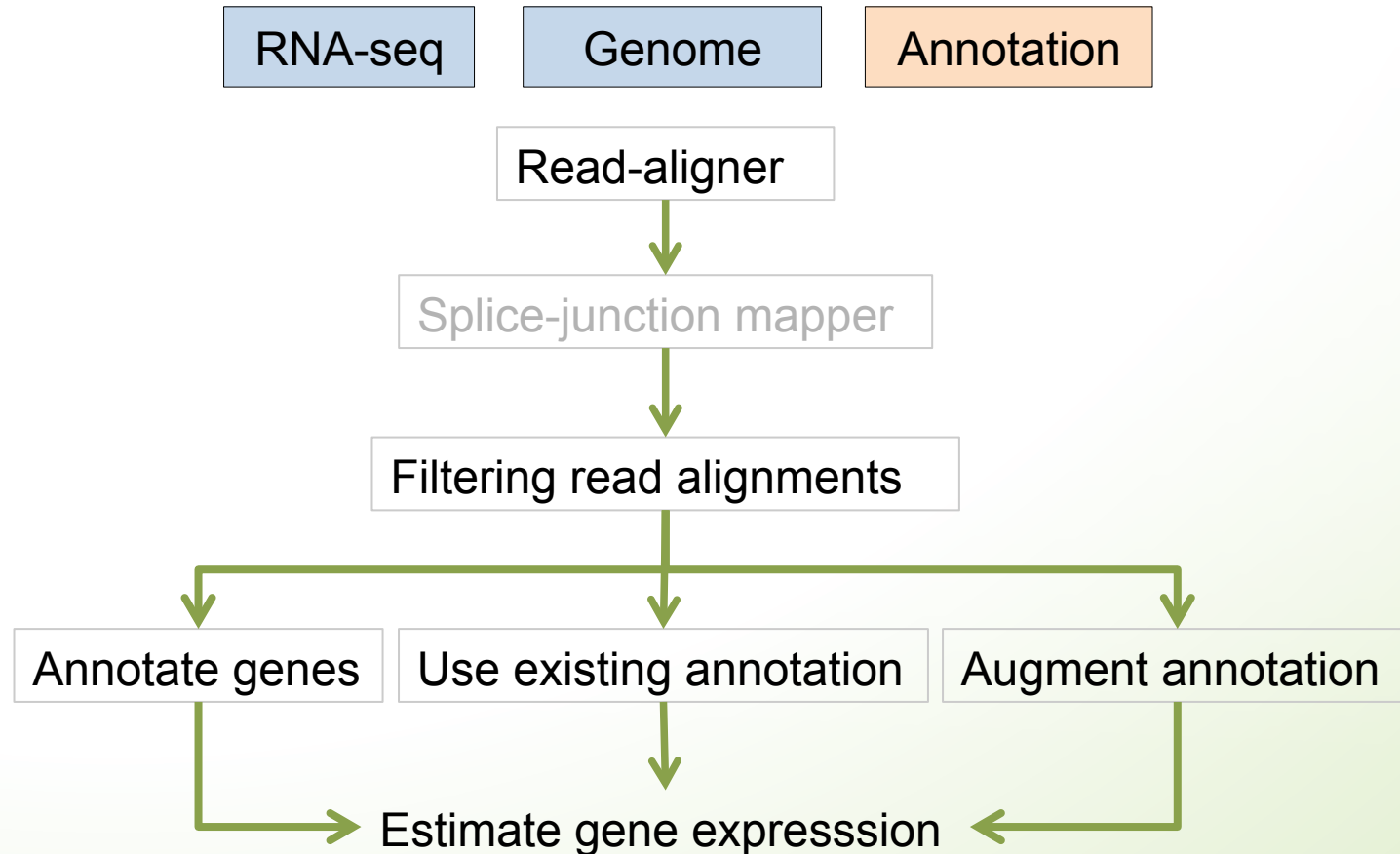
Haas and Zody, Nature Biotechnology 28, 421–423 (2010)

Transcriptome assembly



Haas and Zody, Nature Biotechnology 28, 421–423 (2010)

Mapping-based transcriptome reconstruction



Case study: The transcriptome of the domestic dog



Case study: The transcriptome of the domestic dog

Has shared an environment with humans for > 5000 years

- > Exposed to many of the same environ. influences

Affected by many of the same diseases as man

- > Cancer
- > Heart disease

Extensive breeding and selection

- > Many dog breeds are prone to certain diseases
- > Long haplotypes ideal for association studies



Question: what genes are located in my region of interest?

Requires a high quality genome...and detailed annotation!

Case study: The transcriptome of the domestic dog

Recently, the Broad institute released an updated build, canFam3.1

85 Mb of additional sequence integrated

99.8% of euchromatic portion of genome covered, high quality

Recovered 100s of GC-rich promoter regions

Now approaches level of quality/completion of mouse or human

> the annotation...not so much.

Case study: The transcriptome of the domestic dog

strong discrepancy between well-annotated human genome and dog. Why?

- > largely homology-based
- > almost no isoform information
- > only few dog-specific gene annotations

Majority of loci likely incomplete, many dog-specific genes probably missing

Case study: The transcriptome of the domestic dog

10 tissues at great depth (> 20 million reads)

blood, brain, heart, kidney, liver, lung, muscle, ovary, skin, testes

Stranded paired-end libraries

Poly-A selected: default approach, recovers mostly protein-coding genes

DSN prep: Targets all RNAs, but normalizes library to avoid strong biases

An improved canine genome and a comprehensive catalogue of coding genes and non-coding transcripts. Hoepfner MP et al. PLoS One 2014 Mar 13;9(3):e91172

Mapping-based transcriptome reconstruction

Align reads with Tophat/Bowtie

Reconstruct transcripts with
Cufflinks

Reconcile de-novo annotation
with reference

Annotate novel transcripts

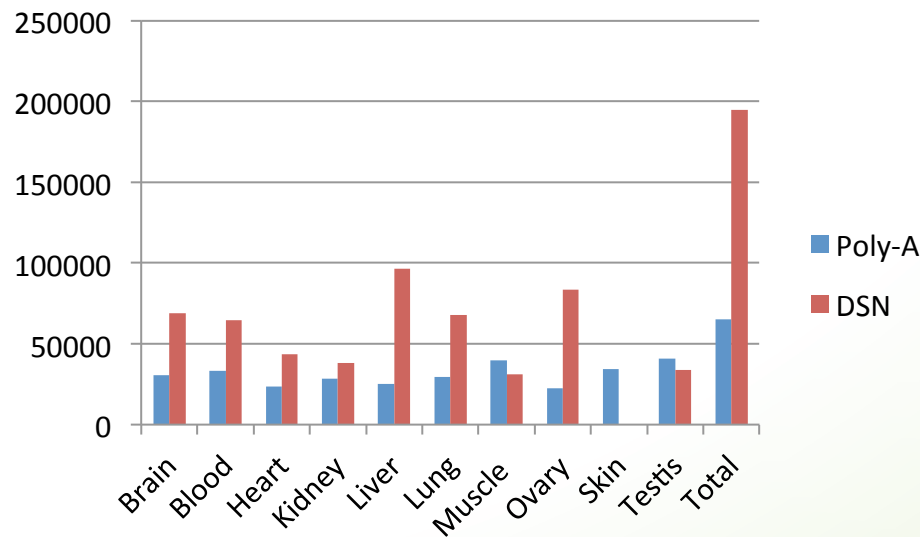
Quantify

Mapping-based transcriptome reconstruction



Case study: The transcriptome of the domestic dog

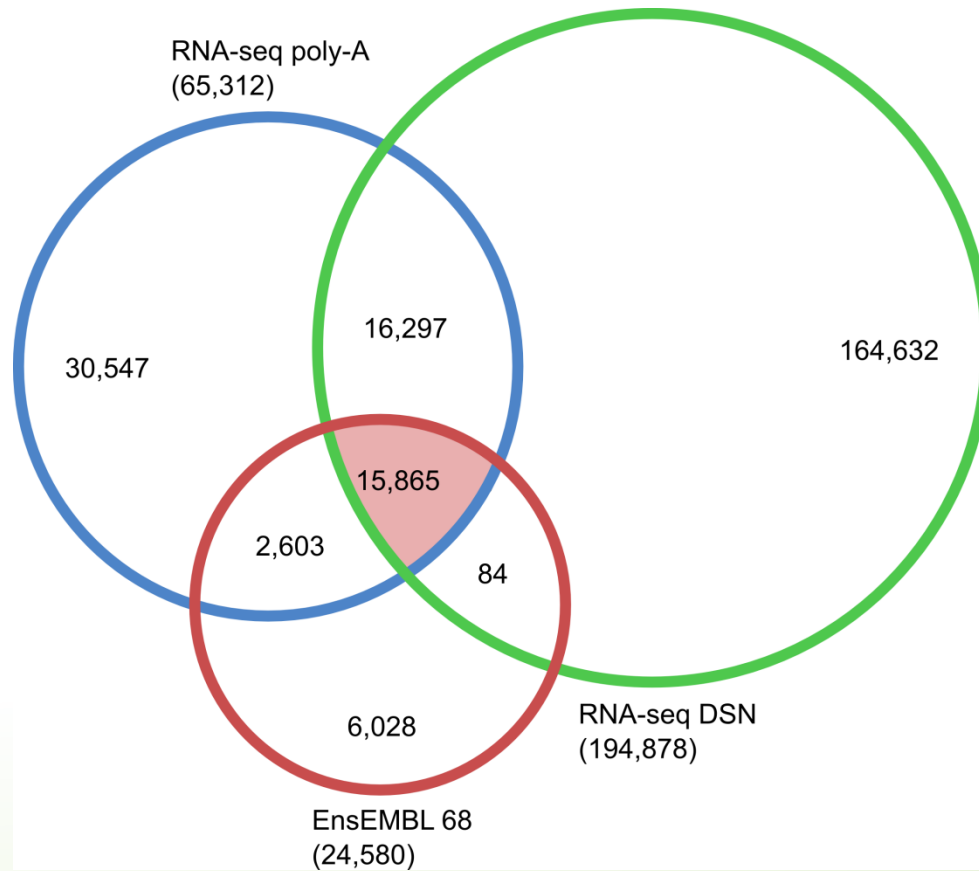
Transcript reconstruction using cufflinks for both libraries



DSN recovers more transcripts than polyA

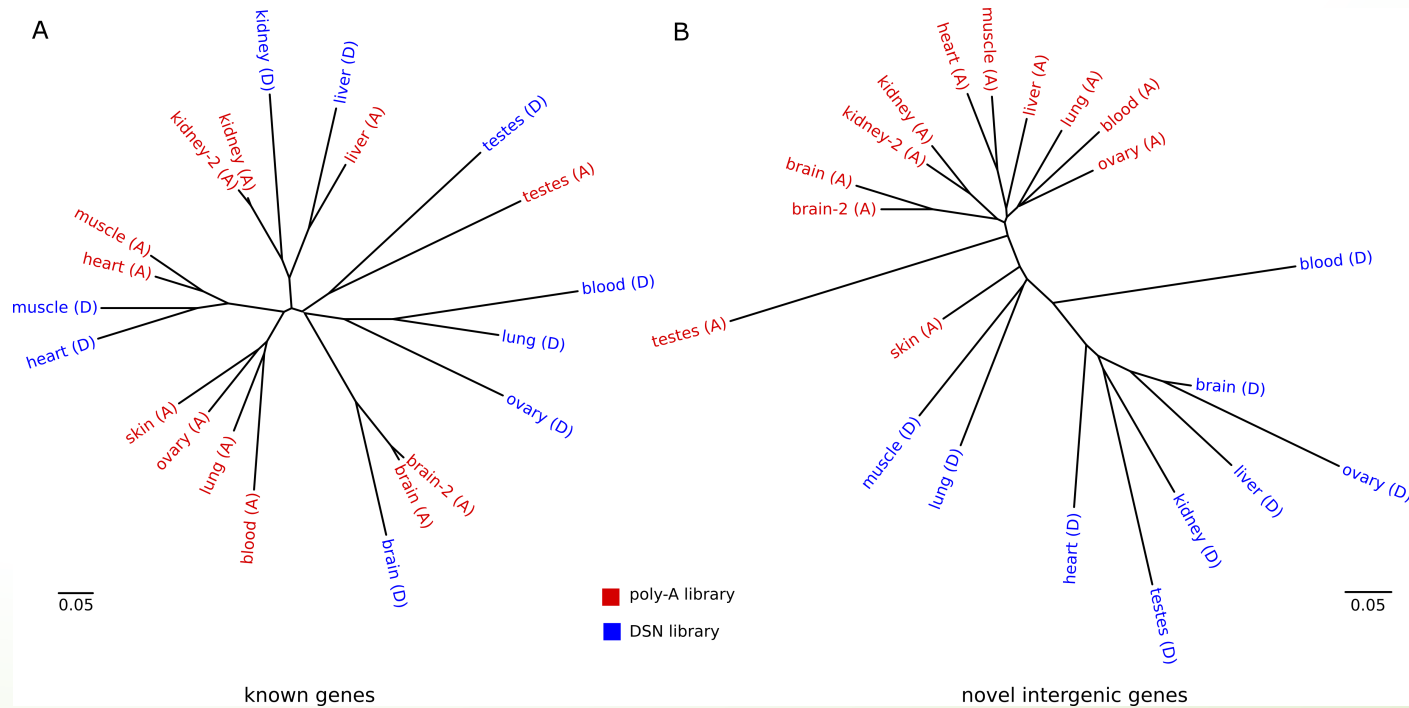
Transcriptional diversity is highest in testes

Case study: The transcriptome of the domestic dog

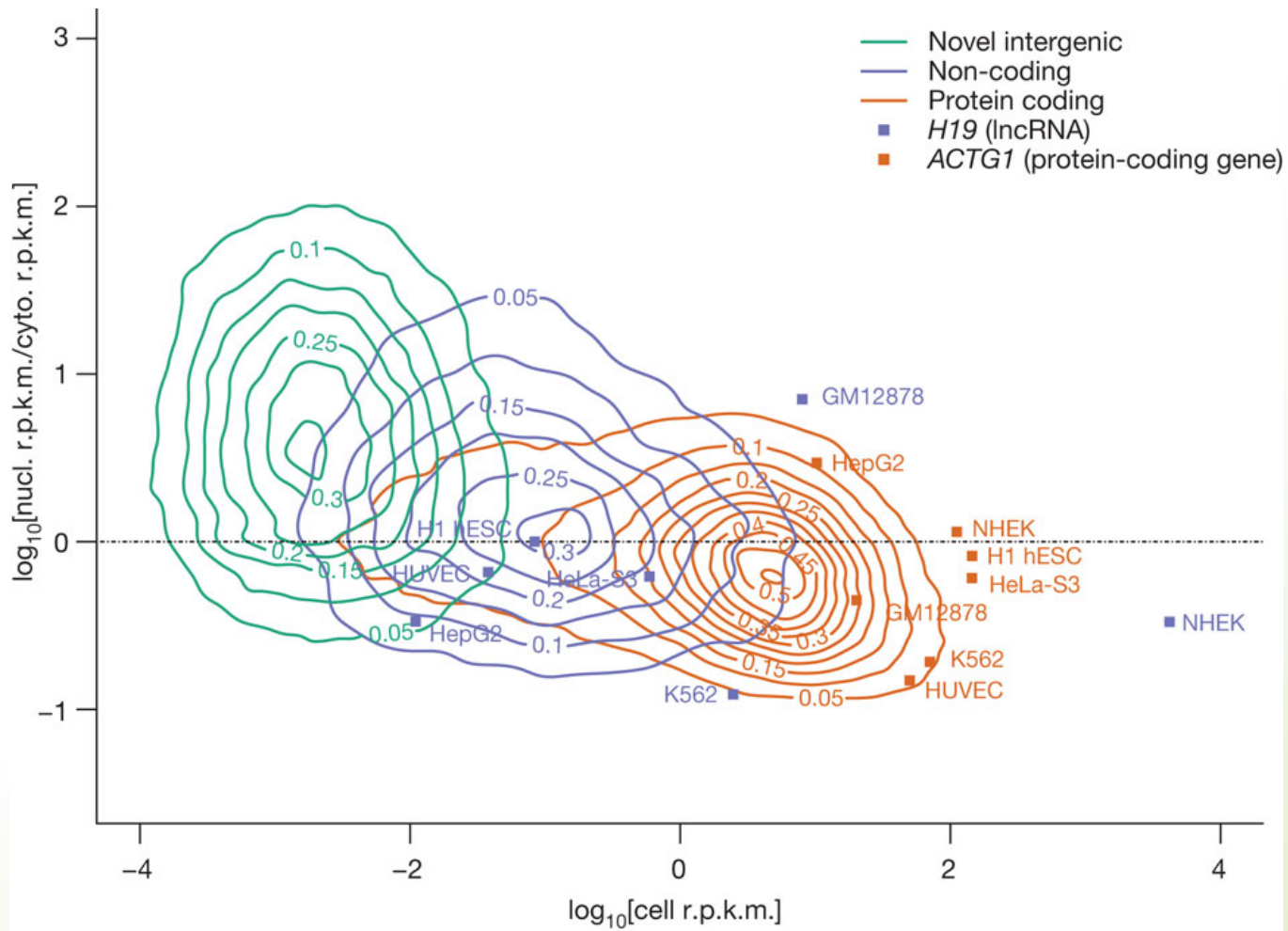


Case study: The transcriptome of the domestic dog

Transcript reconstruction using cufflinks for both libraries

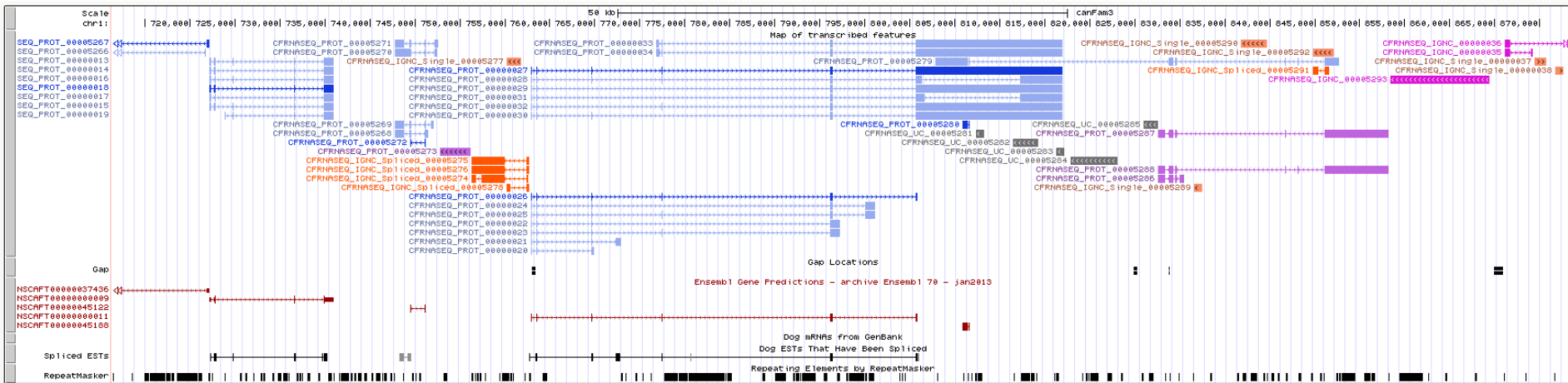


RNA flavors



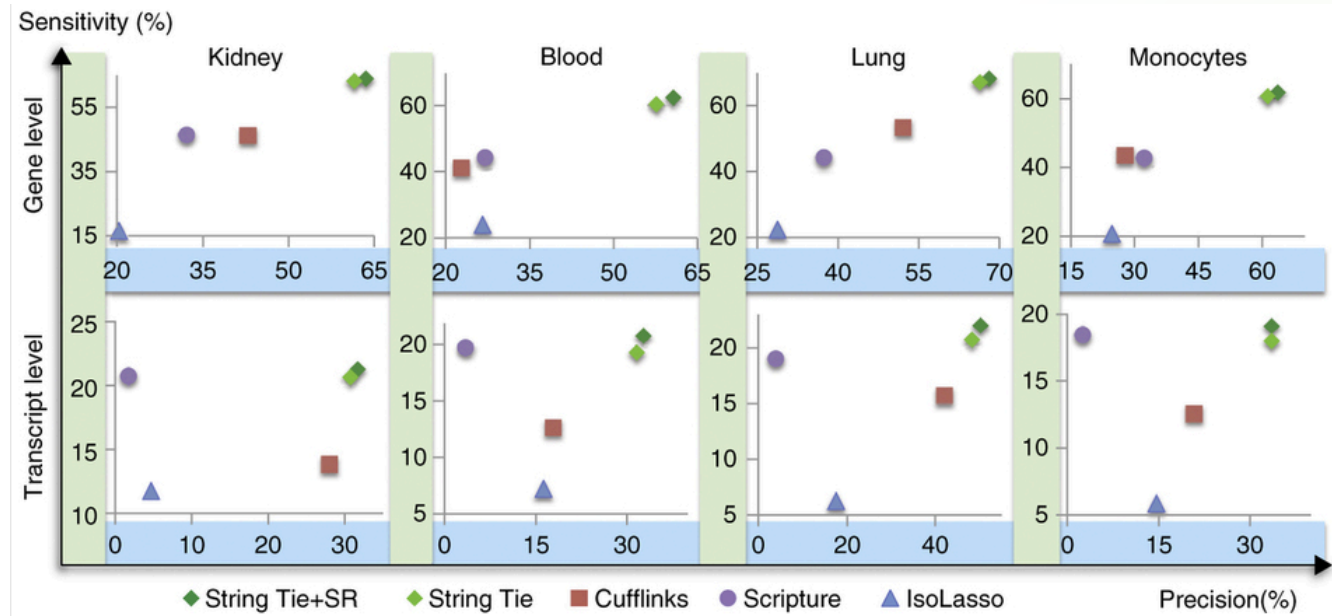
Case study: The transcriptome of the domestic dog

Augmented annotation and transcript classification

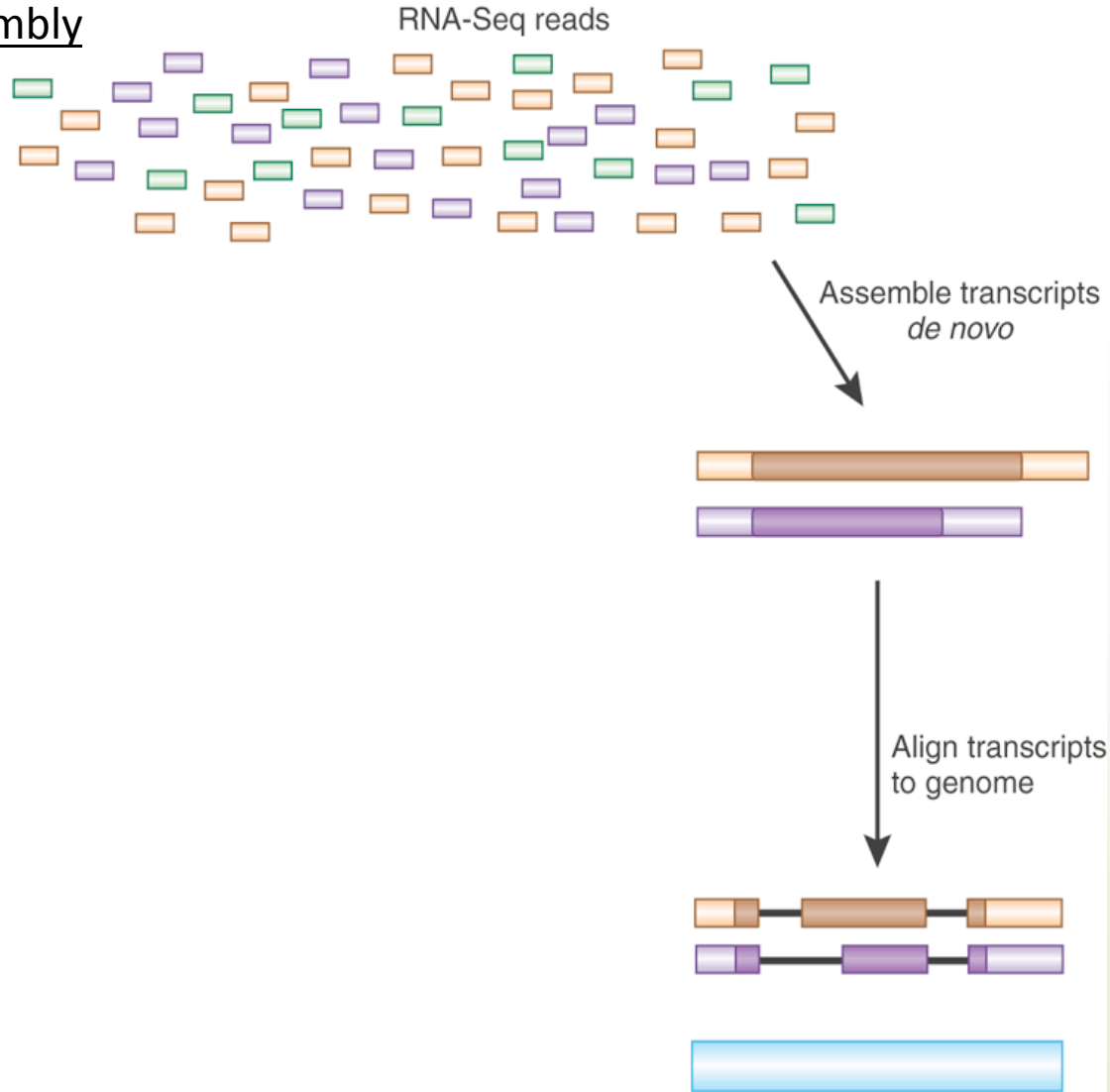


Several softwares

- Cufflinks
- Scripture
- Ballgown
- StringTie



Transcriptome assembly



Haas and Zody, Nature Biotechnology 28, 421-423 (2010)

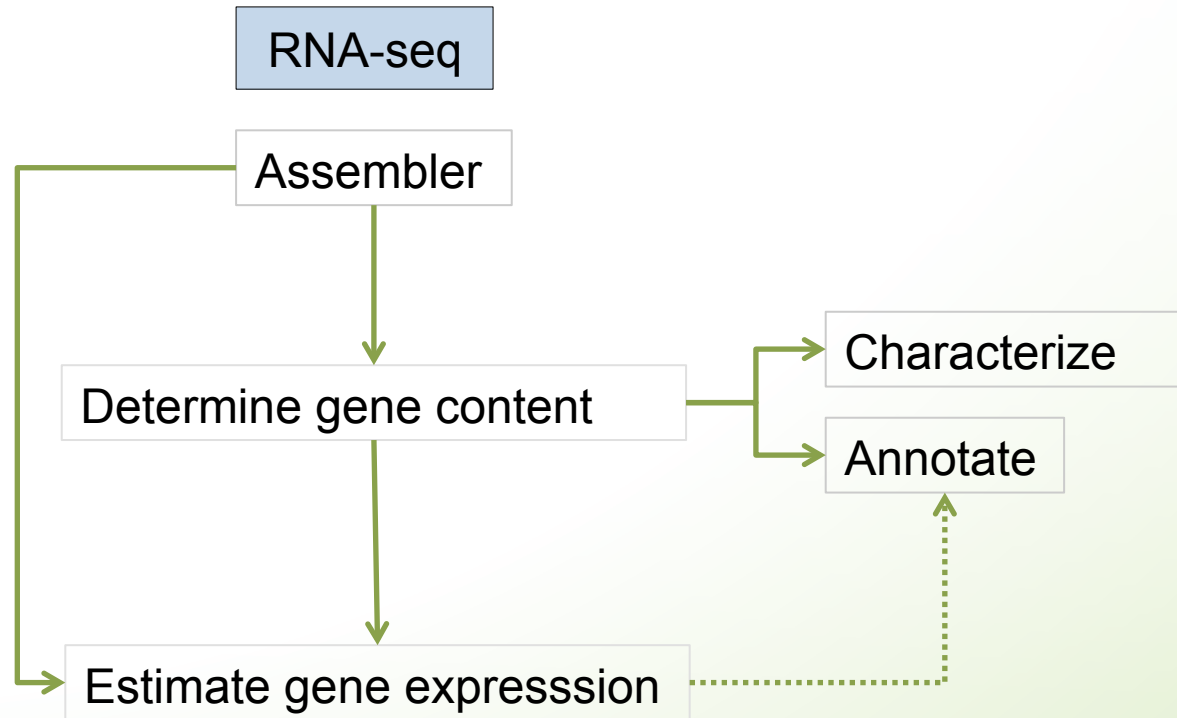
De-novo transcriptome assembly

For the majority of species, there are no comprehensive genome sequences...

Transcriptomics can inform a broad range of questions without reference

→ De-novo transcriptome assembly from extracted RNA

De-novo transcriptome reconstruction



De-novo transcriptome assembly

Manfred Grabherr

Brian Haas

Moran Yassour

Kerstin Lindblad-Toh

Aviv Regev

Nir Friedman

David Eccles

Alexie Papanicolaou

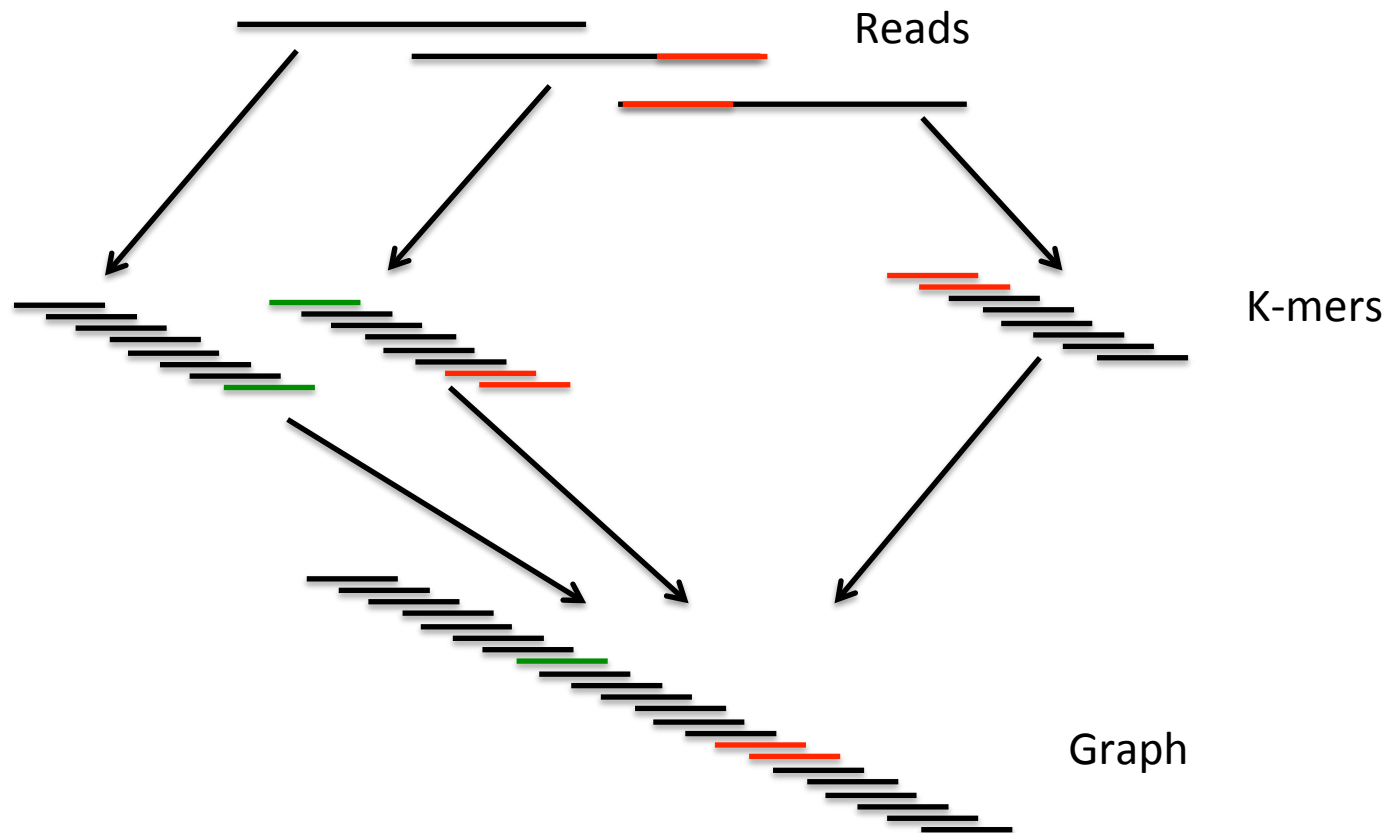
Michael Ott

...



The k-mer

- K consecutive nucleotides



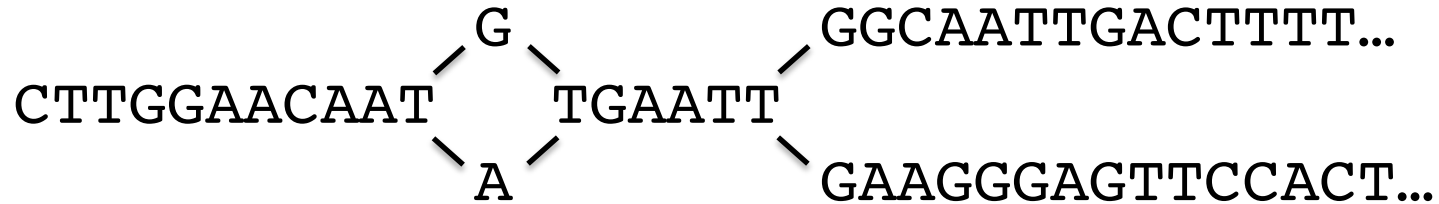
The de Bruijn Graph

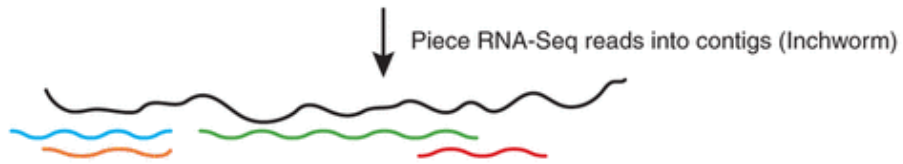
- Graph of overlapping sequences
- Intended for cryptology
- Fixed length element: k

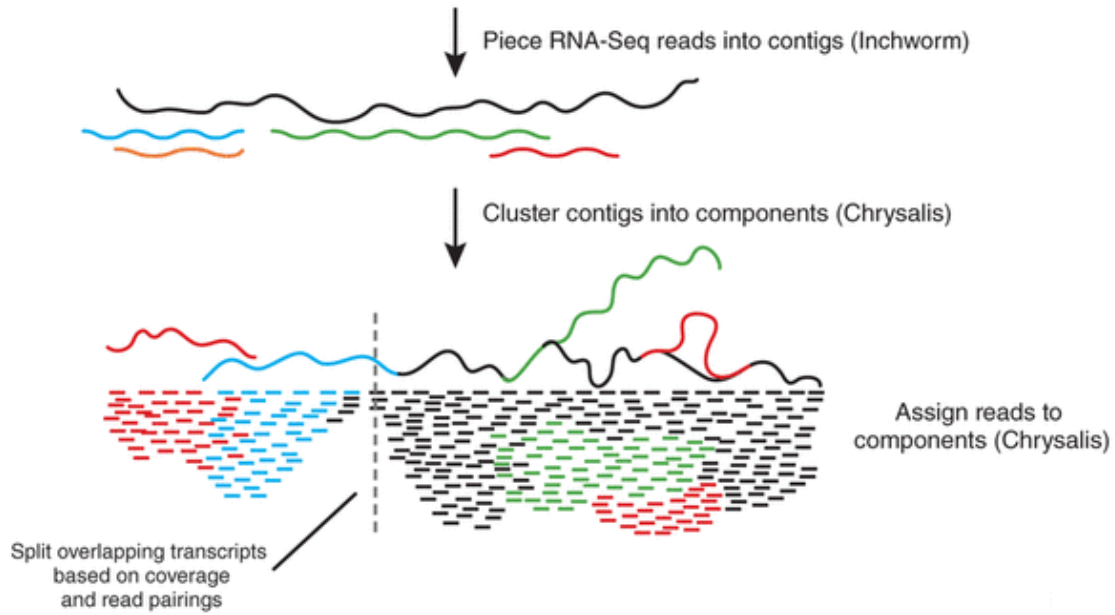
CTTGGAA
 TTGGAAC
 TGGAACA
 GGAACAA
 GAACAAT

The de Bruijn Graph

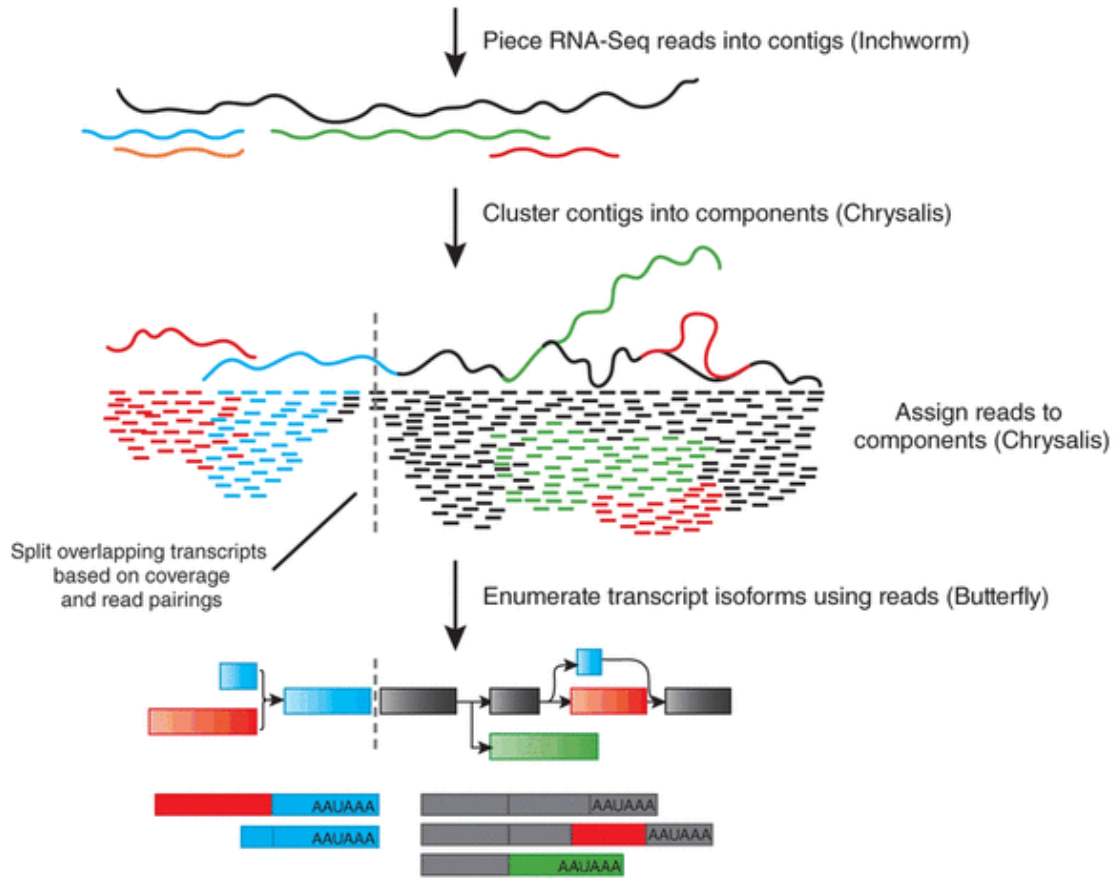
- Graph has “nodes” and “edges”

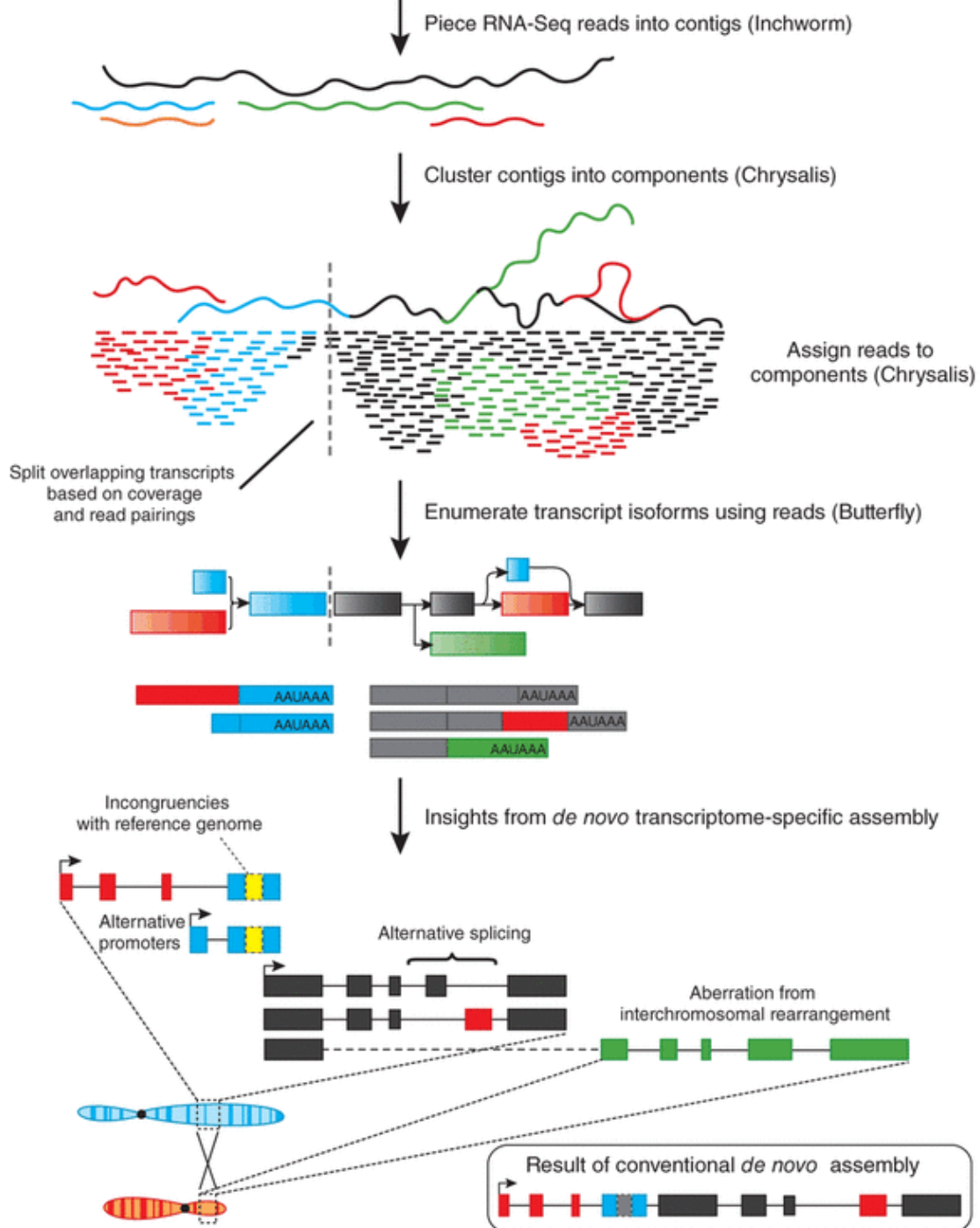






Iyer MK, Chinnaiyan AM (2011)
Nature Biotechnology **29**, 599–600





Iyer MK, Chinnaiyan AM (2011)
Nature Biotechnology **29**, 599–600

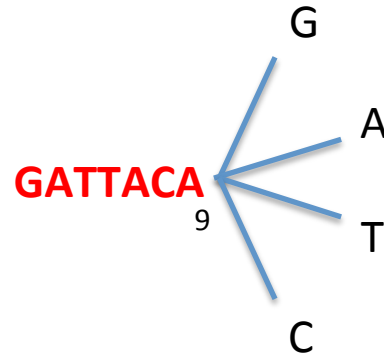


Inchworm Algorithm

Decompose all reads into overlapping Kmers (25-mers)

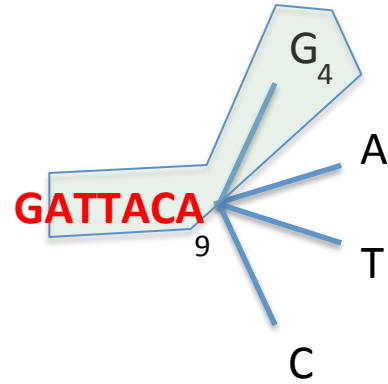
Identify seed kmer as most abundant Kmer, ignoring low-complexity kmers.

Extend kmer at 3' end, guided by coverage.



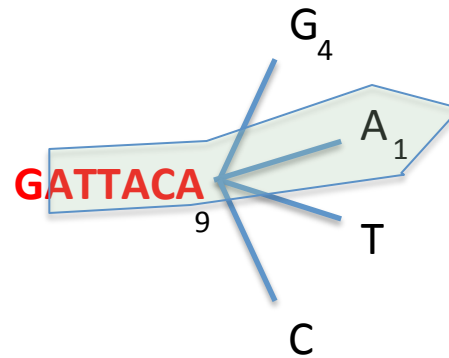


Inchworm Algorithm



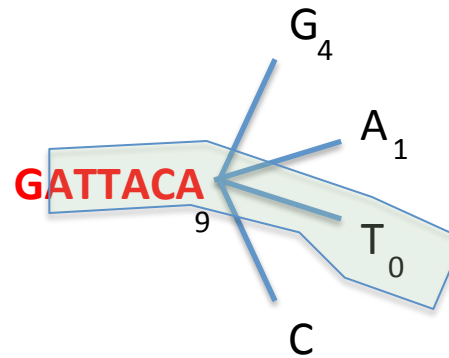


Inchworm Algorithm



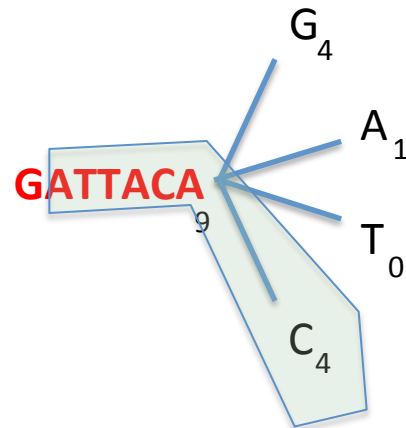


Inchworm Algorithm



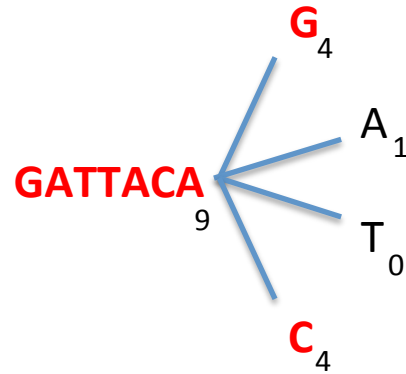


Inchworm Algorithm



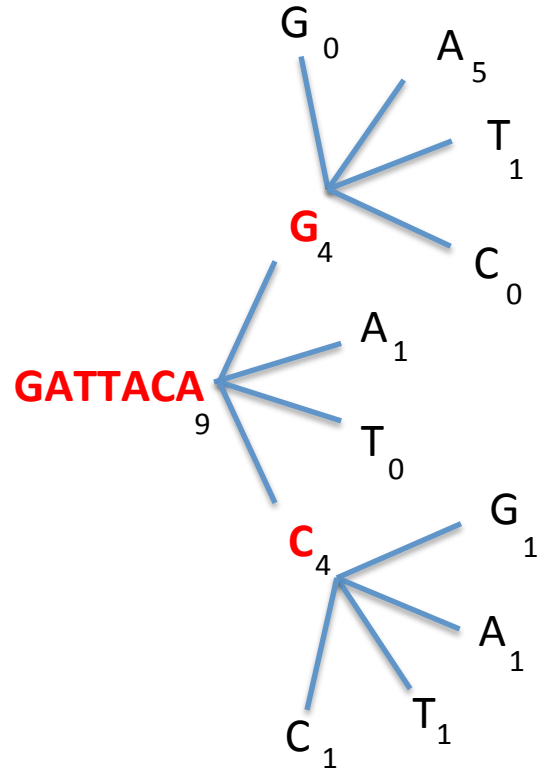


Inchworm Algorithm



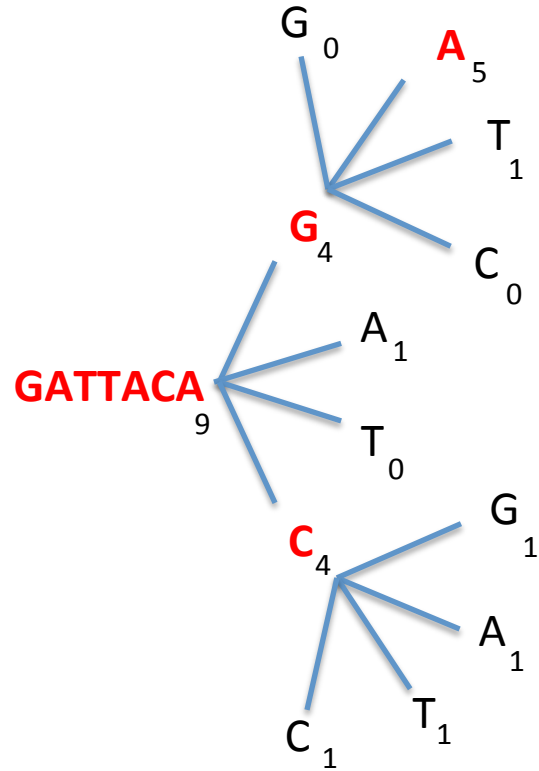


Inchworm Algorithm



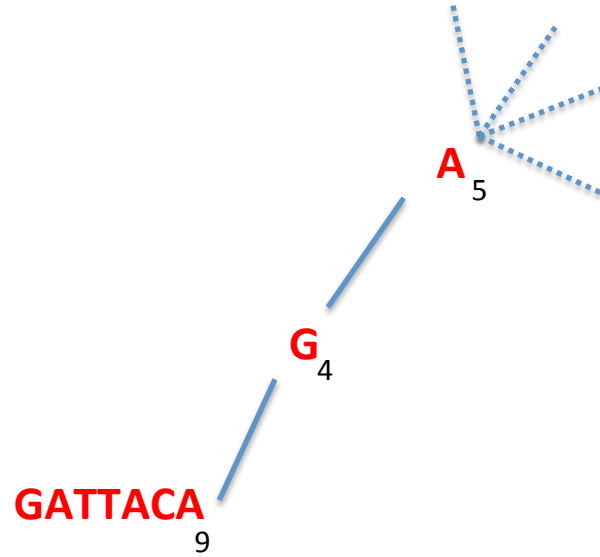


Inchworm Algorithm



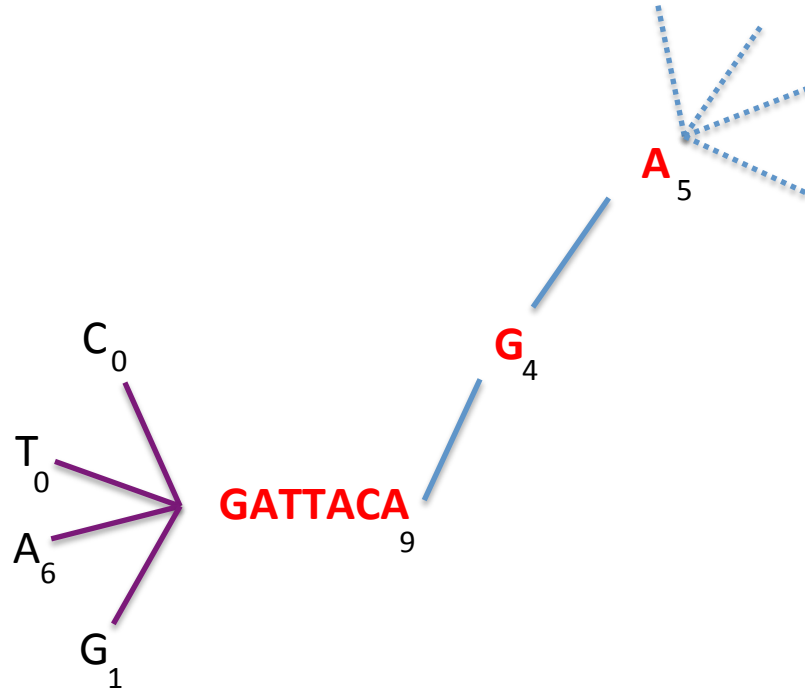


Inchworm Algorithm



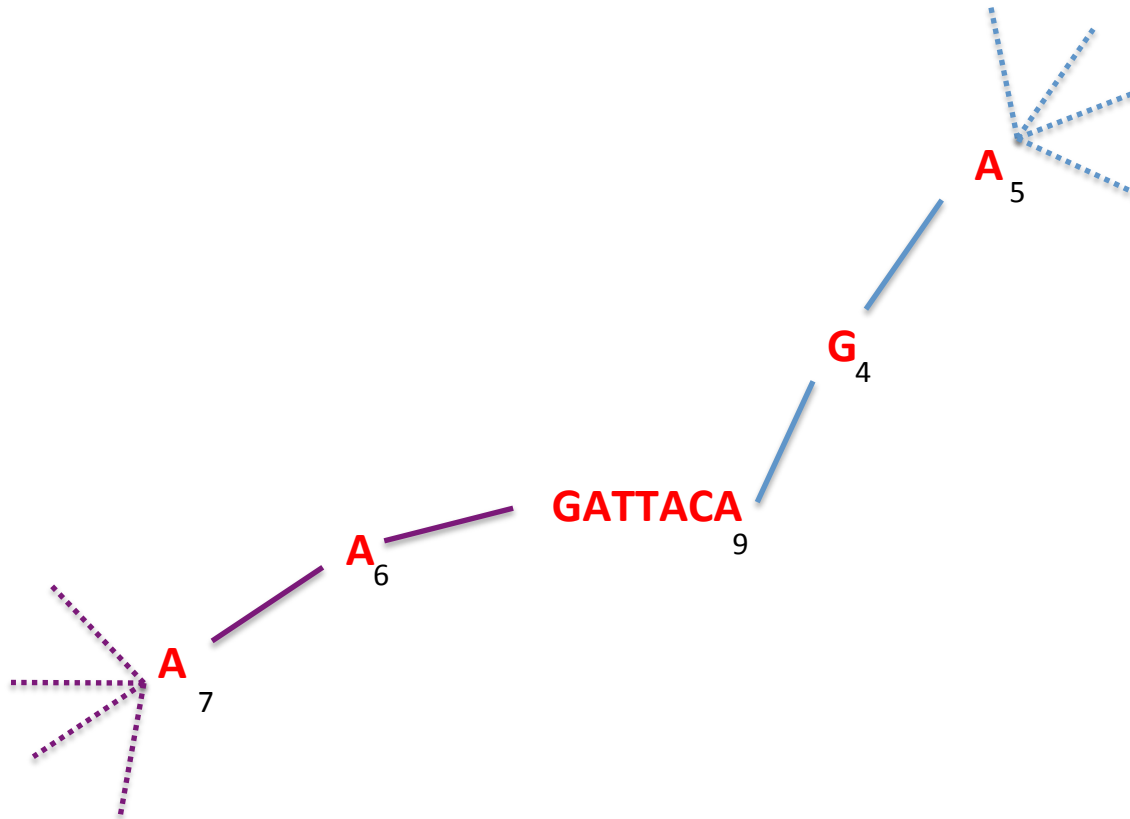


Inchworm Algorithm





Inchworm Algorithm

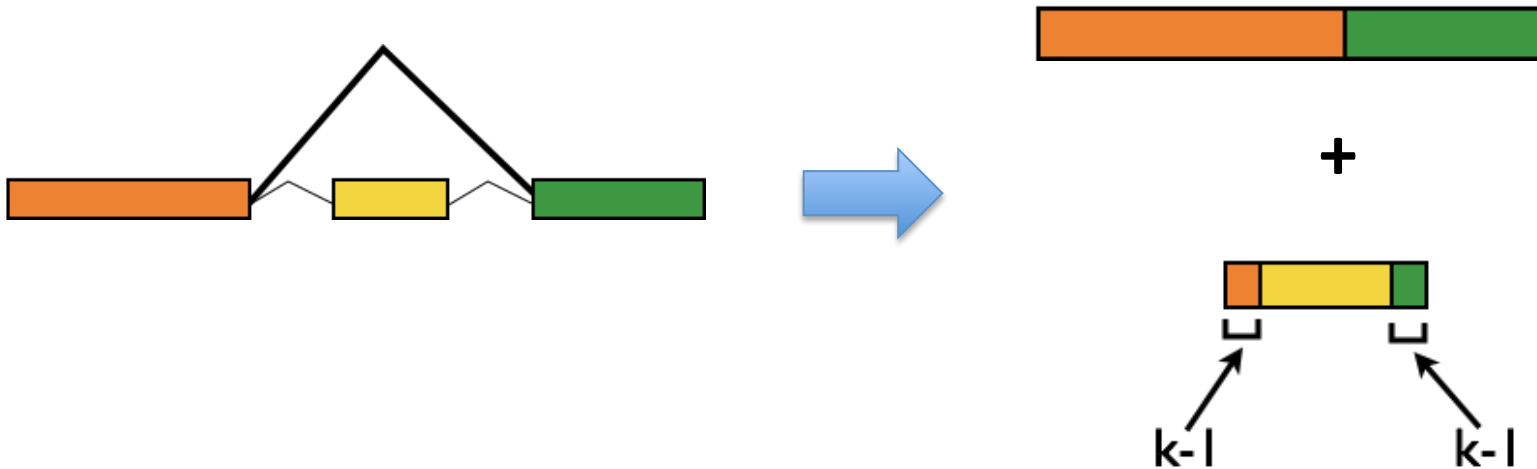


Report contig: **....AAGATTACAGA....**

Remove assembled kmers from catalog, then repeat the entire process.



Inchworm Contigs from Alt-Spliced Transcripts => Minimal lossless representation of data



Chrysalis

>a121:len=5845



>a122:len=2560



>a123:len=4443



>a124:len=48



>a125:len=8876



>a126:len=68



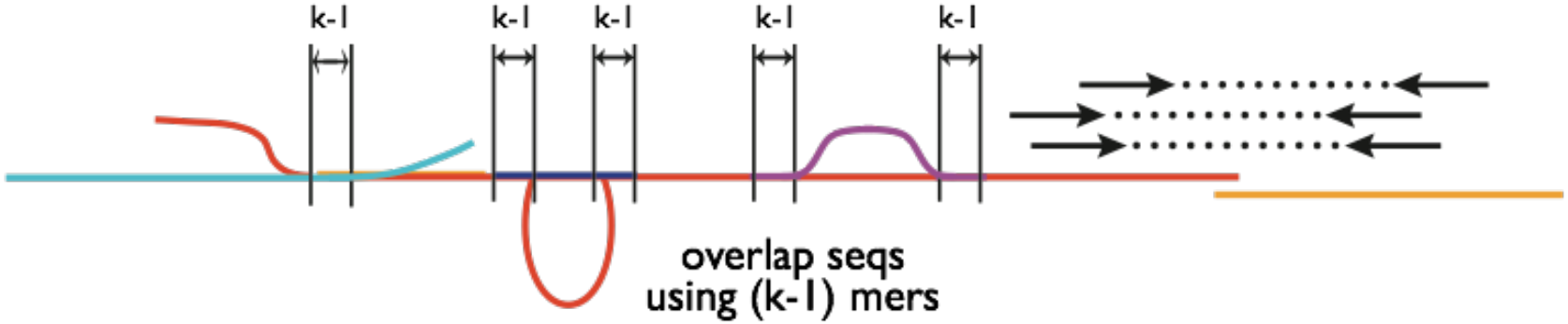
Integrate isoforms
via k-1 overlaps

Chrysalis

>a121:len=5845
>a122:len=2560
>a123:len=4443
>a124:len=48
>a125:len=8876
>a126:len=68



Integrate isoforms
via $k-1$ overlaps



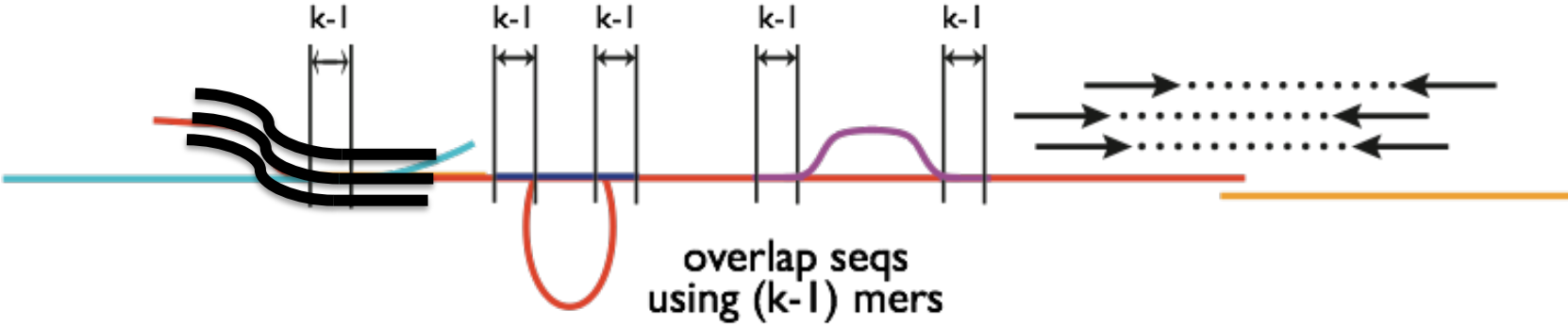
overlap seqs
using $(k-1)$ mers

Chrysalis

```
>a121:len=5845  
_____  
>a122:len=2560  
_____  
>a123:len=4443  
_____  
>a124:len=48  
_____  
>a125:len=8876  
_____  
>a126:len=68  
_____
```



Integrate isoforms
via $k-1$ overlaps
Verify via "welds"



Chrysalis

>a121:len=5845

>a122:len=2560

>a123:len=4443

>a124:len=48

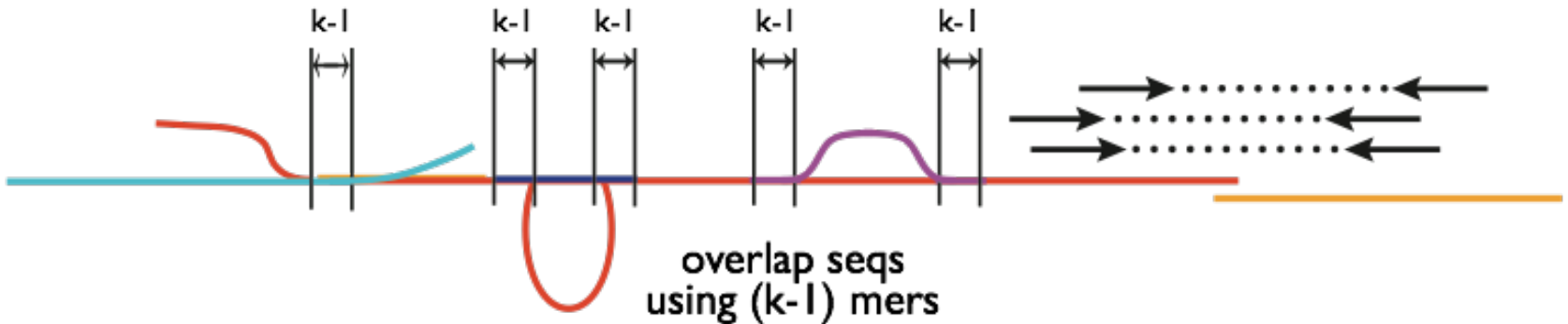
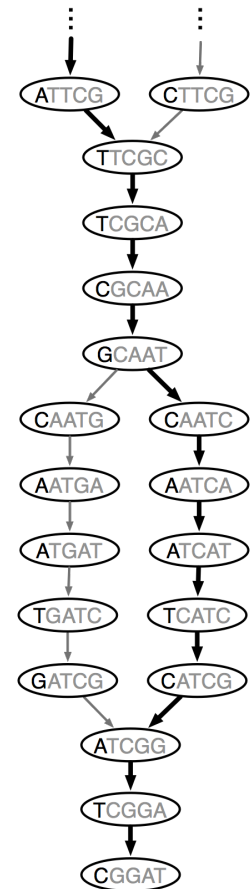
>a125:len=8876

>a126:len=68

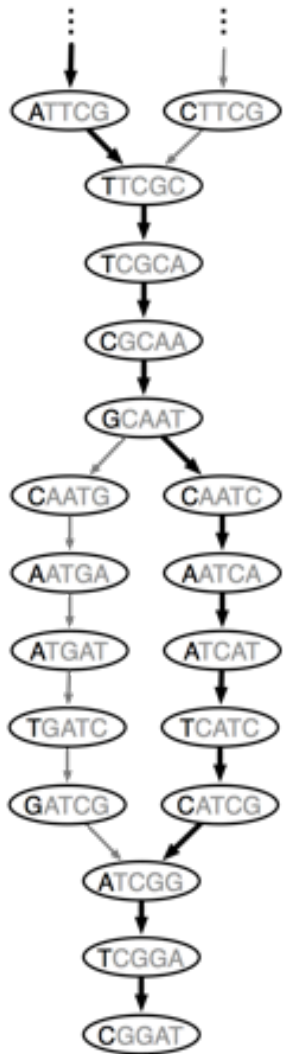


Integrate isoforms
via $k-1$ overlaps
Verify via "welds"

Build de Bruijn Graphs
(ideally, one per gene)

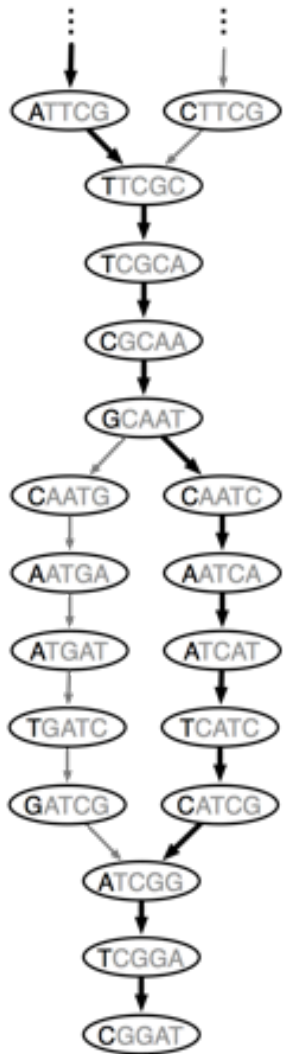


Butterfly



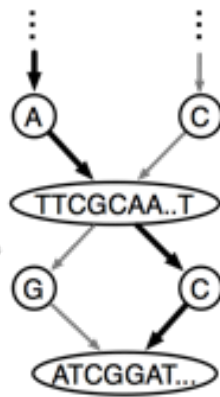
de Bruijn
graph

Butterfly



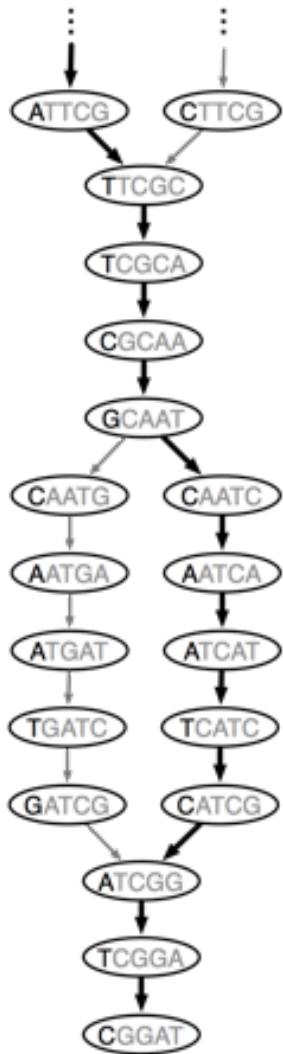
de Bruijn graph

compacting



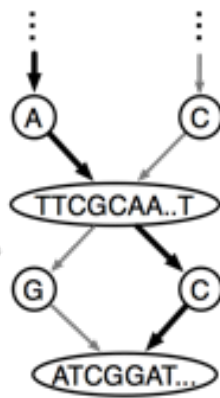
compact graph

Butterfly



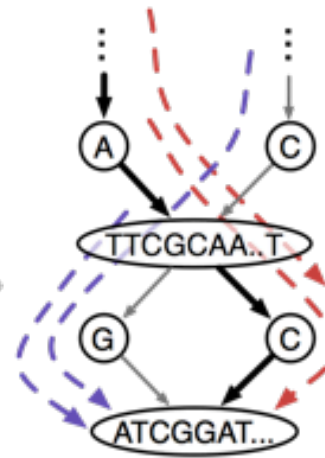
de Bruijn graph

compacting



compact graph

finding paths



compact graph with reads

extracting sequences

..CTTCGCAA..TGATCGGAT...
..ATTGCAA..TCATCGGAT...

sequences

Completeness and coverage as function of read counts

