

De novo genome assembly versus mapping to a reference genome

Beat Wolf

PhD. Student in Computer Science

University of Würzburg, Germany

University of Applied Sciences Western Switzerland

beat.wolf@hefr.ch



- Genetic variations
- De novo sequence assembly
- Reference based mapping/alignment
- Variant calling
- Comparison
- Conclusion

What are variants?

- Difference between a sample (patient) DNA and a reference (another sample or a population consensus)
- Sum of all variations in a patient determine his genotype and phenotype

- Small variations (< 50bp)
 - SNV (Single nucleotide variation)

Reference

A

T

Sample

- Indel (insertion/deletion)

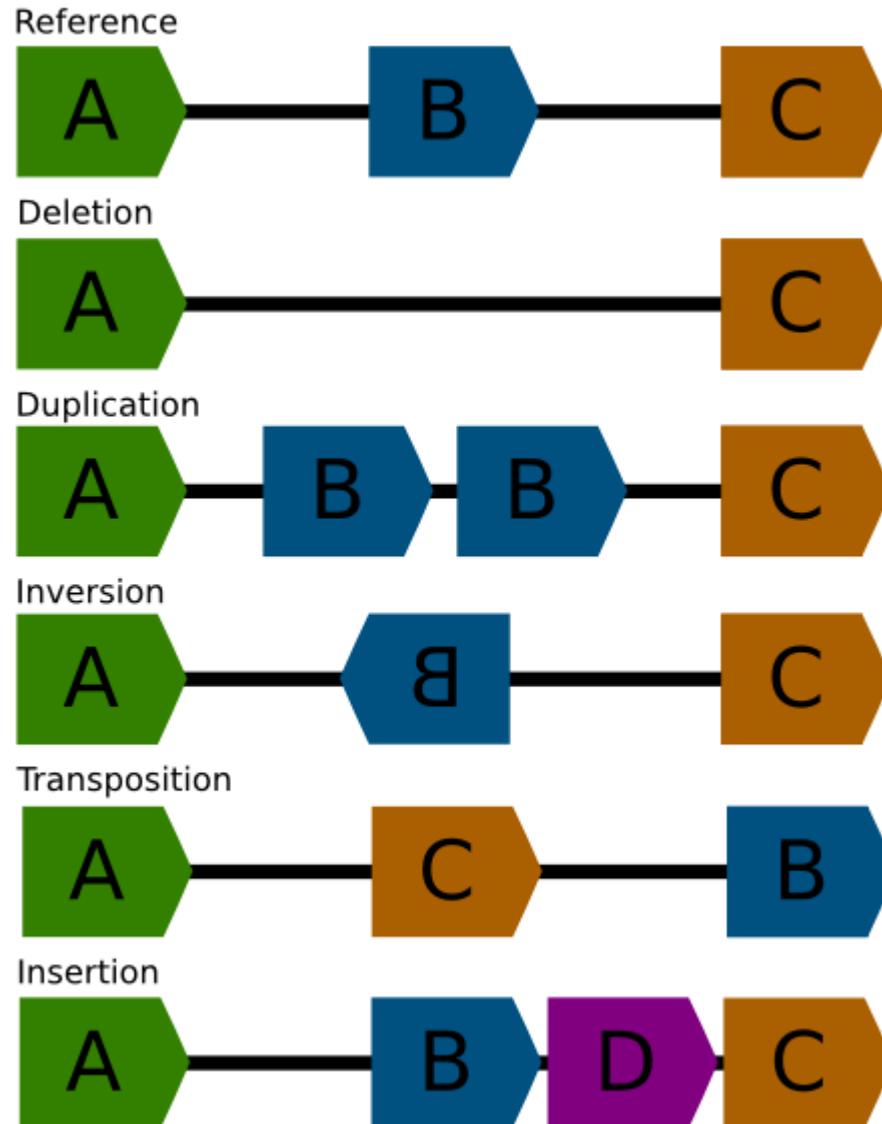
Reference

T A G

T - G

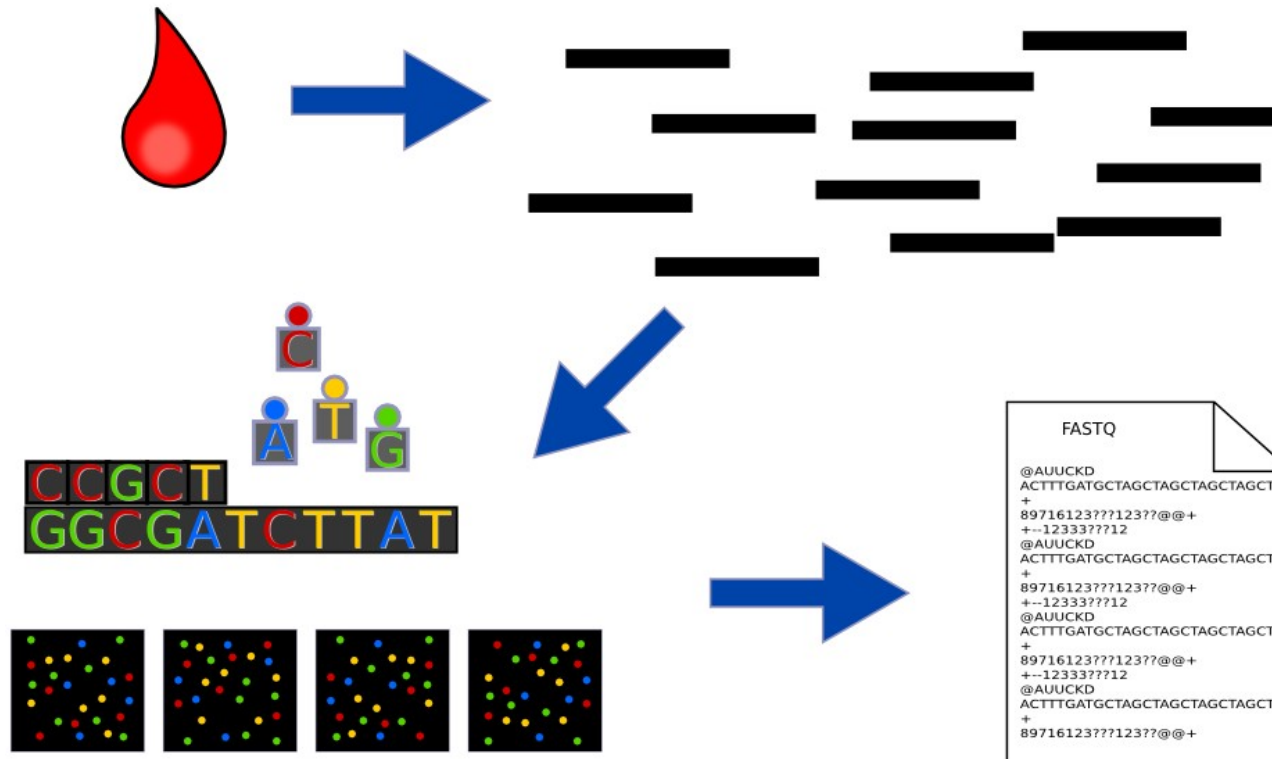
Sample

Structural variations

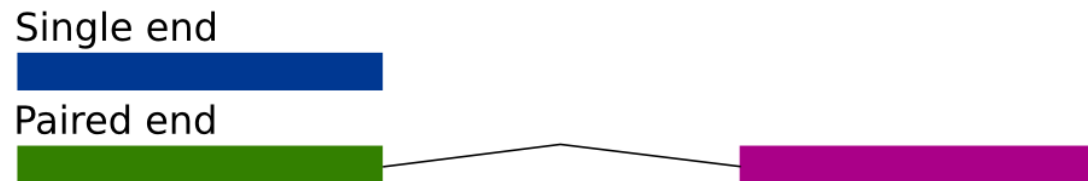


Sequencing technologies

- Sequencing produces small overlapping sequences



- Difference read lengths, 36 – 10'000bp (150-500bp is typical)
- Different sequencing technologies produce different data



And different kinds of errors

- Substitutions (Base replaced by other)
- Homopolymers (3 or more repeated bases)
 - AAAAA might be read as AAAA or AAAAAA
- Insertion (Non existent base has been read)
- Deletion (Base has been skipped)
- Duplication (cloned sequences during PCR)
- Somatic cells sequenced

- Standardized output format: **FASTQ**
 - Contains the read sequence and a quality for every base

@9WV6Z:791:946

GCTCTTCCGATCTATGGATGCACCAAGATATATGACCCTGTCTGTGGGACTGATGGAA

+

7743992220342217743992220342217743992220342217743992220342

http://en.wikipedia.org/wiki/FASTQ_format



- The problem:
 - Recreate the original patient genome from the sequenced reads
 - For which we don't know where they came from and are noisy
- Solution:
 - Recreate the genome with no prior knowledge using de novo sequence assembly
 - Recreate the genome using prior knowledge with reference based alignment/mapping

De novo sequence assembly

- Ideal approach
- Recreate **original genome** sequence through overlapping sequenced reads

T G A C A A G C
A A G C G T T A
C G T T A C A G
T T A C A G C G

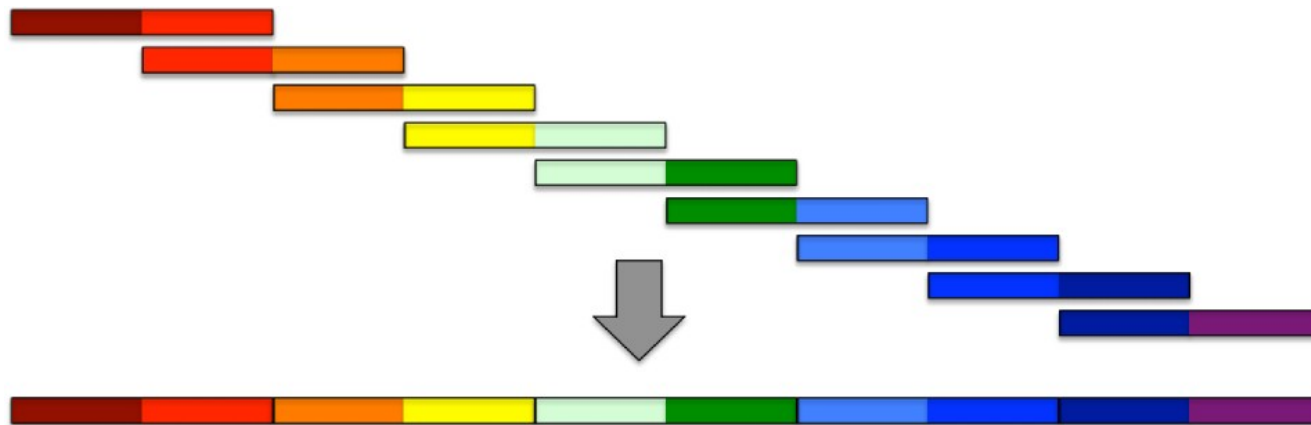
- Construct assembly graph from overlapping reads

...AGCCTAGGGATGCGCGACACGT

GGATGCGCGACACGT CGCATATCCGGTTTGGTCAACCTCGGACGGAC

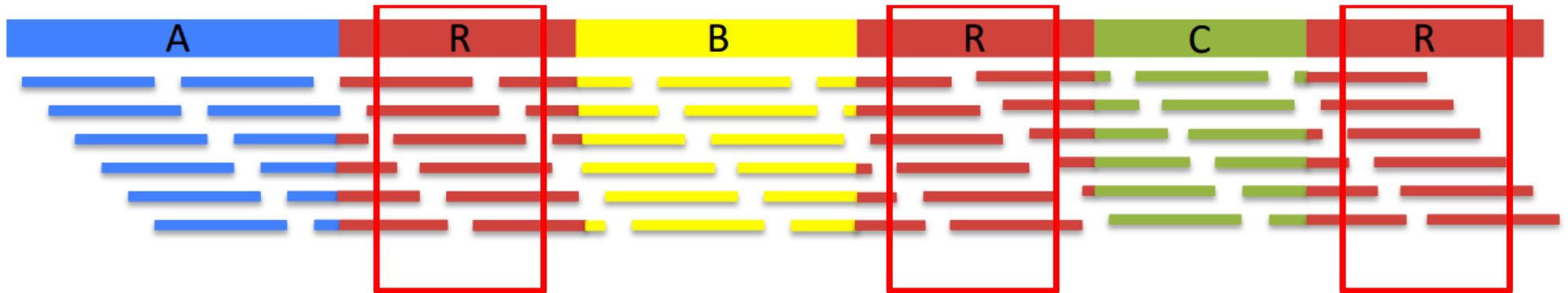
CAACCTCGGACGGACCTCAGCGAA...

- Simplify assembly graph



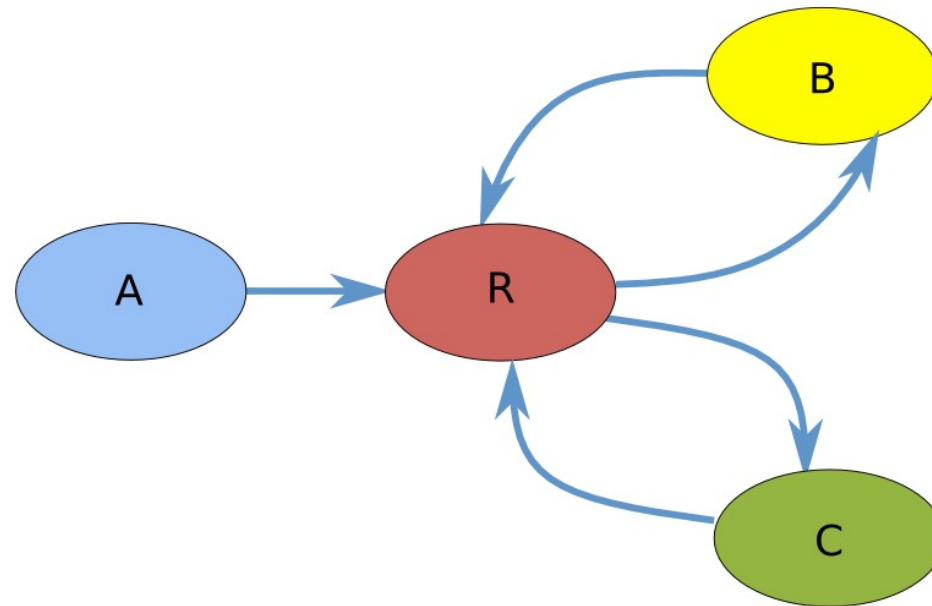
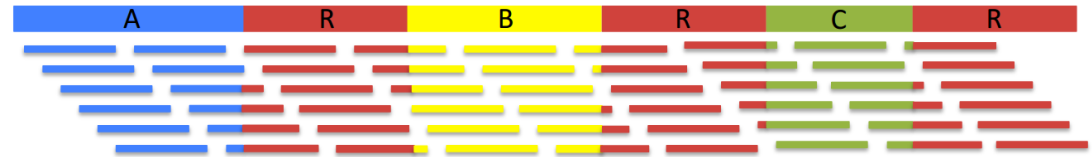
Modified from: De novo assembly of complex genomes using single molecule sequencing, Michael Schatz

- Genome with repeated regions



Modified from: De novo assembly of complex genomes using single molecule sequencing, Michael Schatz

- Graph generation



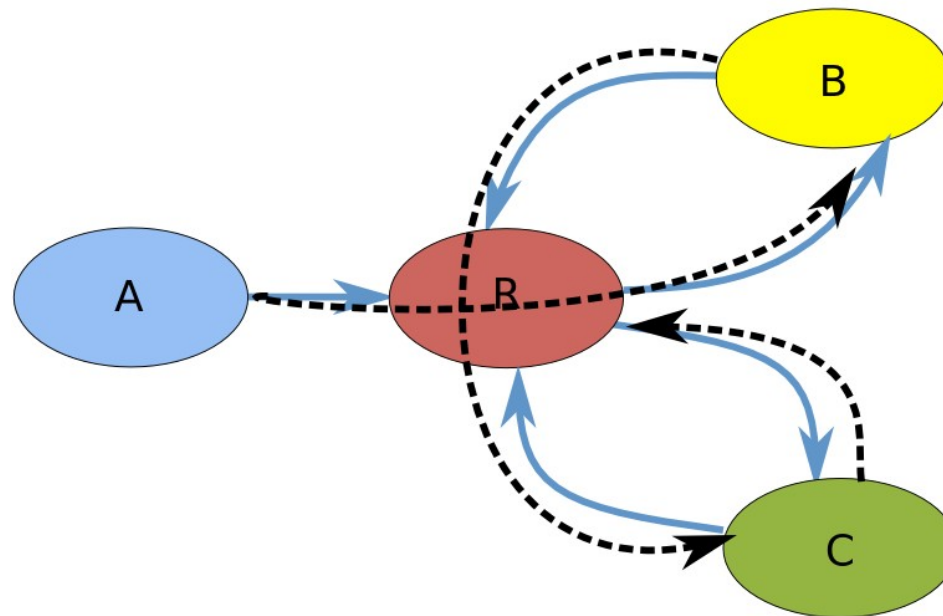
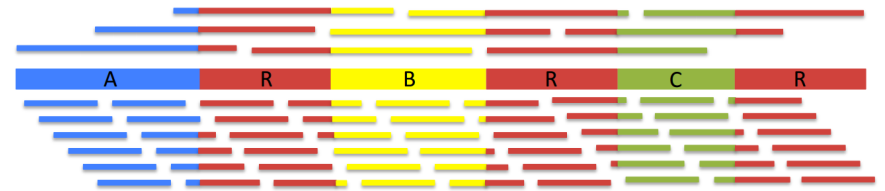
Modified from: De novo assembly of complex genomes using single molecule sequencing, Michael Schatz

- Double sequencing, once with short and once with long reads (or paired end)



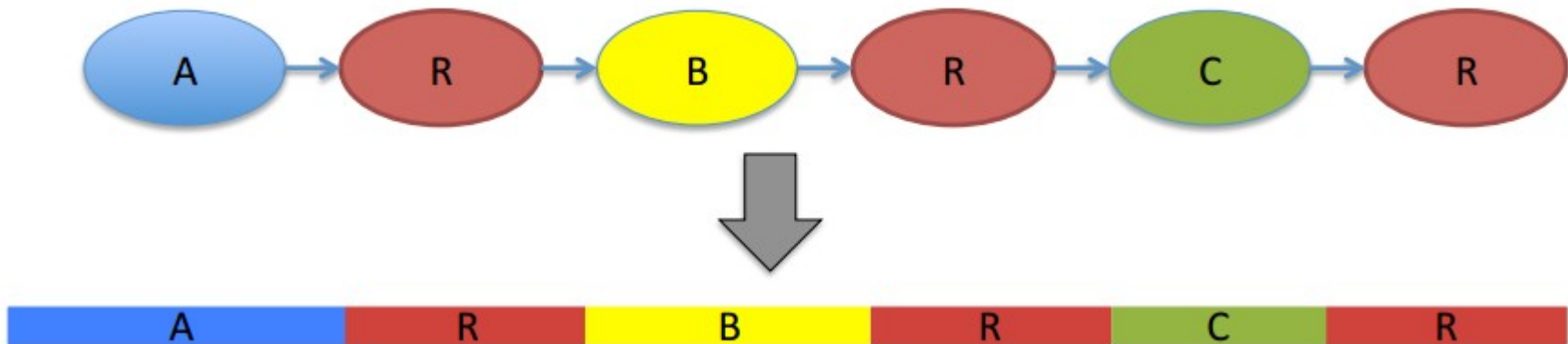
Modified from: De novo assembly of complex genomes using single molecule sequencing, Michael Schatz

- Finding the correct path through the graph with:
 - Longer reads
 - Paired end reads



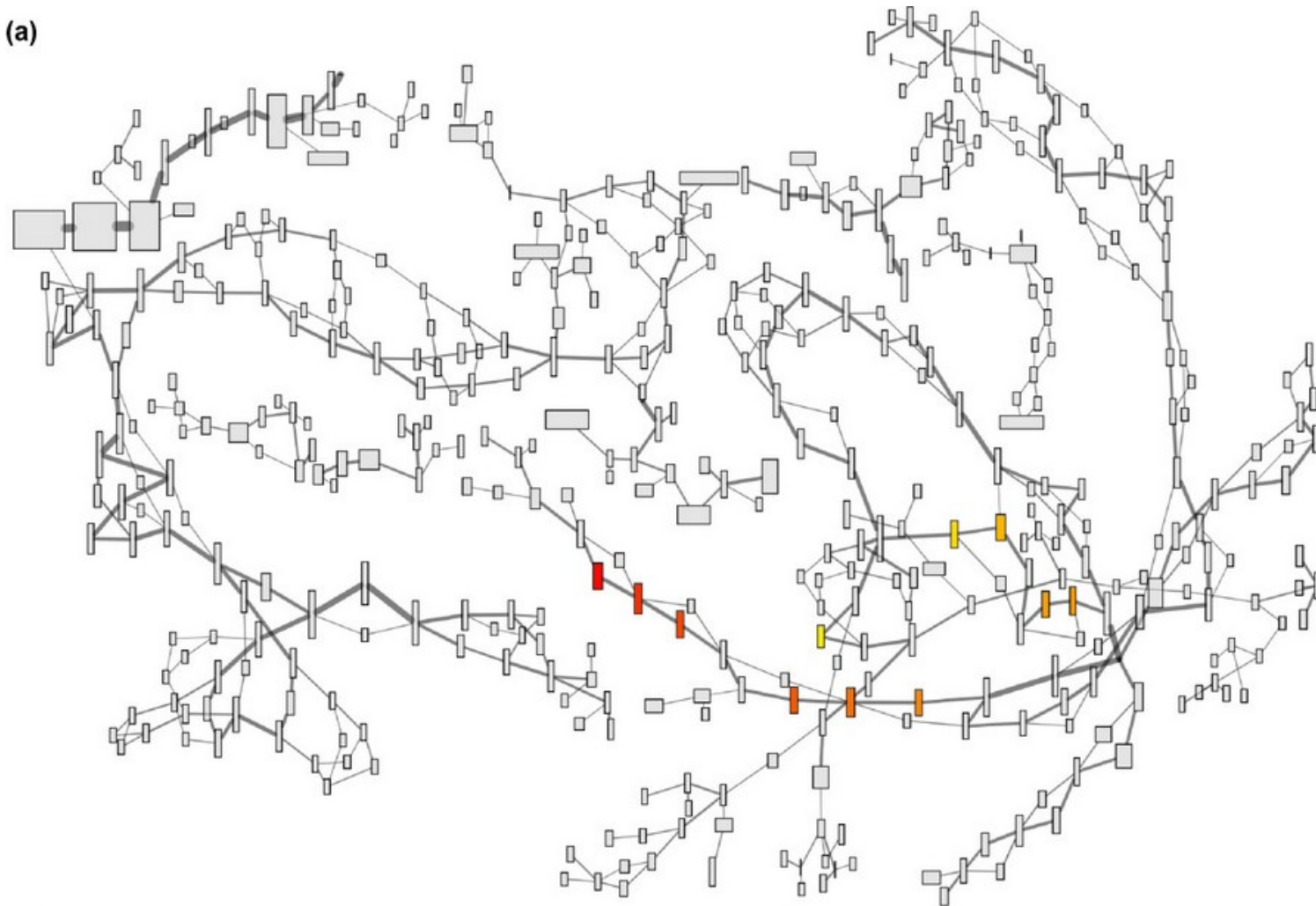
Modified from: De novo assembly of complex genomes using single molecule sequencing, Michael Schatz

De novo sequence assembly



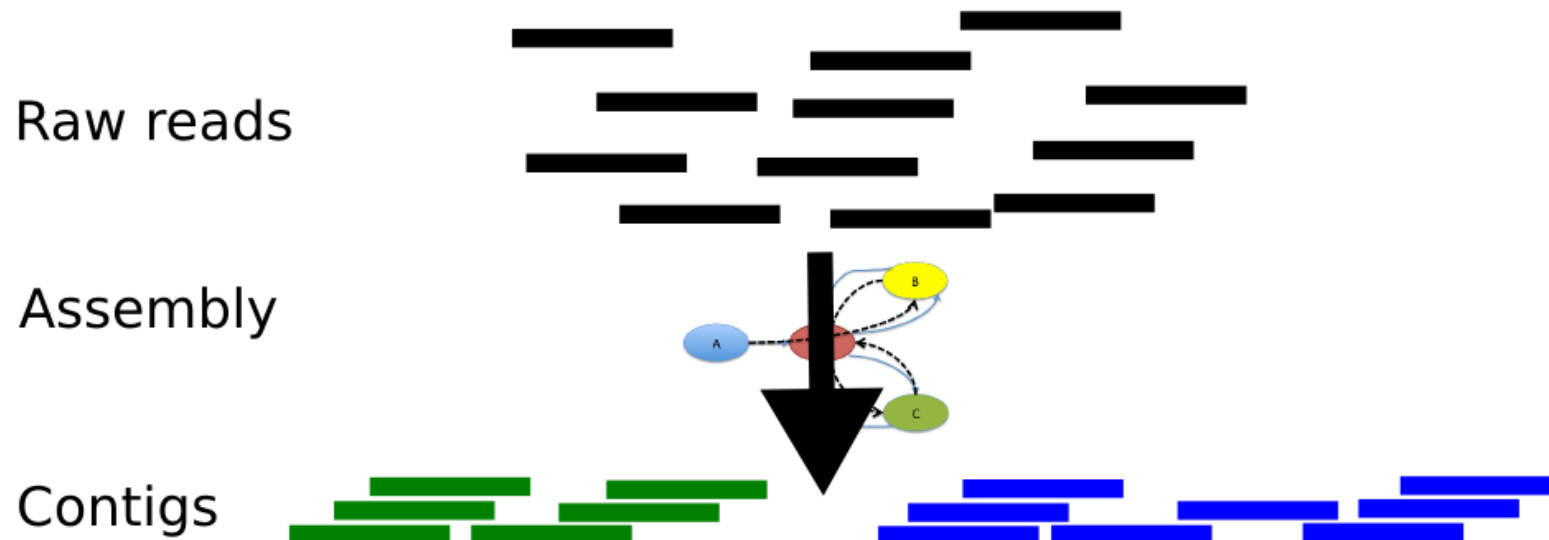
Modified from: De novo assembly of complex genomes using single molecule sequencing, Michael Schatz

(a)

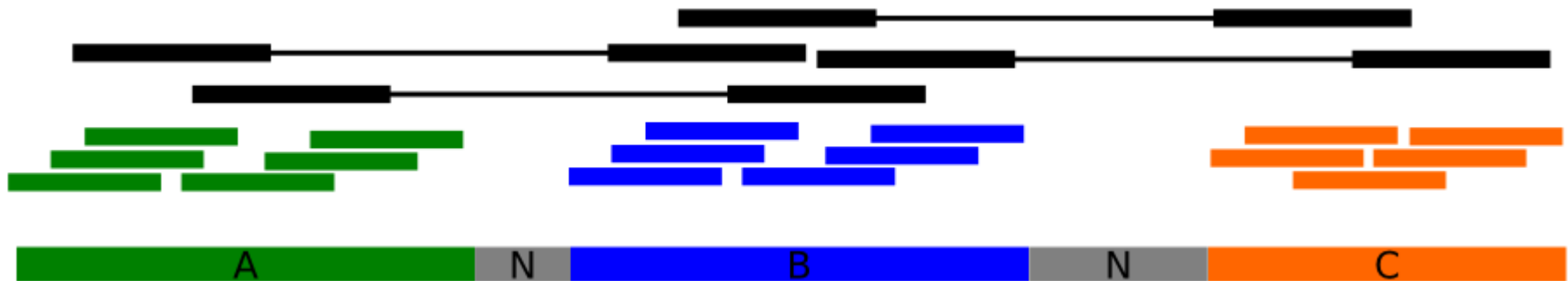


Modified from: EMIRGE: reconstruction of full-length ribosomal genes from microbial community short read sequencing data, Miller et al.

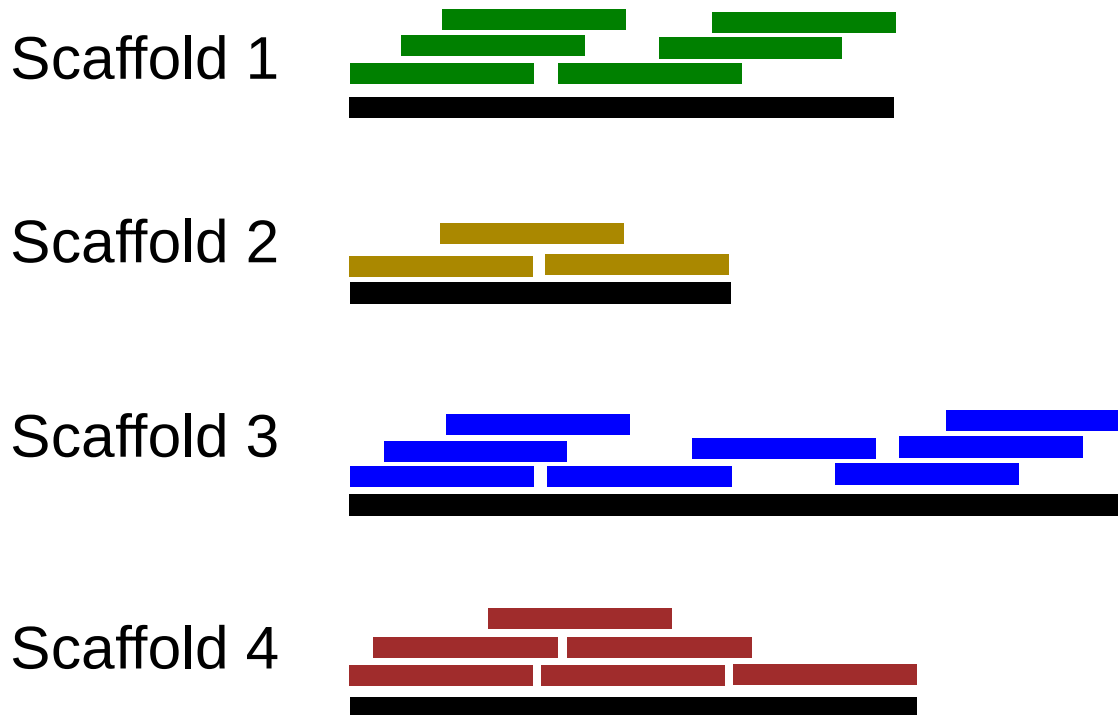
- Overlapping reads are assembled into groups, so called contigs



- Scaffolding
 - Using paired end information, contigs can be put in the right order



- Final result, a list of scaffolds
 - In an ideal world of the size of a chromosome, molecule, mtDNA etc.



- What is needed for a good assembly?
 - High coverage
 - High read lengths
 - Good read quality
- Current sequencing technologies do not have all three
 - Illumina, good quality reads, but short
 - PacBio, very long reads, but low quality

- Combined sequencing technologies assembly
 - High quality contigs created with short reads
 - Scaffolding of those contigs with long reads



- Double sequencing means
 - High infrastructure requirements
 - High costs

- Field of assemblers is constantly evolving
 - Competitions like Assemblathon 1 + 2 exist
<https://genome10k.soe.ucsc.edu/assemblathon>
- The results vary greatly depending on datatype and species to be assembled
- High memory and computational complexity

- Short list of assemblers
 - ALLPATHS-LG
 - Meraculous
 - Ray
- Software used by winners of Assemblathon 2:

SeqPrep, KmerFreq, Quake, BWA, Newbler, ALLPATHS-LG, Atlas-Link, Atlas-GapFill, Phrap, CrossMatch, Velvet, BLAST, and BLASR
- Creating a high quality assembly is complicated

- Human Genome project
 - Produced the first „complete“ human genome
- Human genome reference consortium
 - Constantly improves the reference
 - GRCh38 released at the end of 2013



Reference based alignment

- A previously assembled genome is used as a reference
- Sequenced reads are independently aligned against this reference sequence
- Every read is placed at its most likely position
- Unlike sequence assembly, no synergies between reads exist

- Naive approach:
 - Evaluate every location on the reference

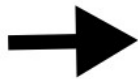
Reference

ACTGA TGAAT ACTGA

Reads

ACTGA

TGAAT



- Too slow for billions of reads on a big reference

- Speed up with the creation of a reference index



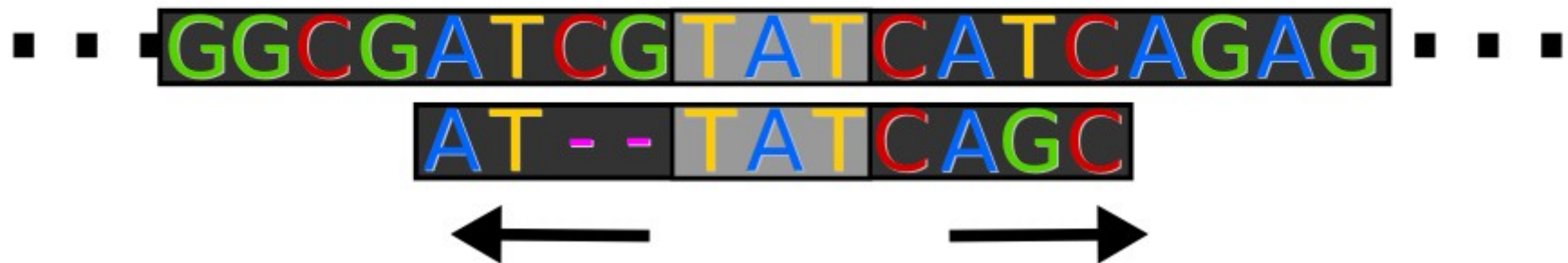
Index

TGA	1		
ACG	2	5	8
TTC	3	7	
CTG	4		
ATT	6		

- Fast lookup table for subsequences in reference

- Determine optimal alignment for the best candidate positions
- Insertions and deletions increase the complexity of the alignment

Seed



- Most common technique, dynamic programming
- Smith-Watherman, Gotoh etc. are common algorithms

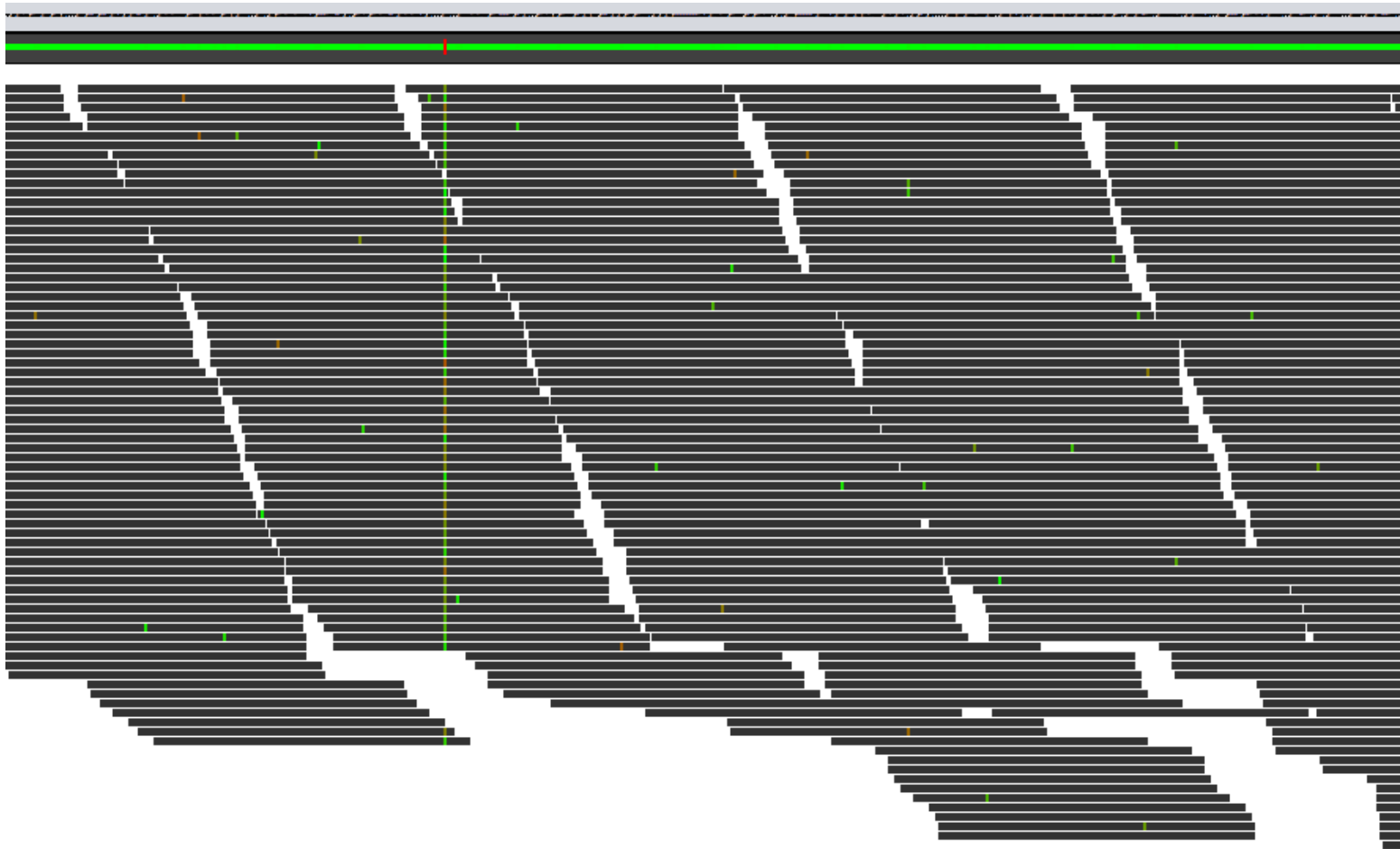
$$H = \begin{pmatrix} - & A & C & A & C & A & C & T & A \\ - & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ A & 0 & 2 & 1 & 2 & 1 & 2 & 1 & 0 & 2 \\ G & 0 & 1 & 1 & 1 & 1 & 1 & 1 & 0 & 1 \\ C & 0 & 0 & 3 & 2 & 3 & 2 & 3 & 2 & 1 \\ A & 0 & 2 & 2 & 5 & 4 & 5 & 4 & 3 & 4 \\ C & 0 & 1 & 4 & 4 & 7 & 6 & 7 & 6 & 5 \\ A & 0 & 2 & 3 & 6 & 6 & 9 & 8 & 7 & 8 \\ C & 0 & 1 & 4 & 5 & 8 & 8 & 11 & 10 & 9 \\ A & 0 & 2 & 3 & 6 & 7 & 10 & 10 & 10 & 12 \end{pmatrix}$$

$$T = \begin{pmatrix} - & A & C & A & C & A & C & T & A \\ - & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ A & 0 & \swarrow & \leftarrow & \swarrow & \leftarrow & \swarrow & \leftarrow & \swarrow \\ G & 0 & \uparrow & \swarrow & \uparrow & \swarrow & \uparrow & \swarrow & \uparrow \\ C & 0 & \uparrow & \swarrow & \leftarrow & \swarrow & \leftarrow & \swarrow & \leftarrow \\ A & 0 & \swarrow & \uparrow & \swarrow & \leftarrow & \swarrow & \leftarrow & \swarrow \\ C & 0 & \uparrow & \swarrow & \uparrow & \swarrow & \leftarrow & \swarrow & \leftarrow \\ A & 0 & \swarrow & \uparrow & \swarrow & \uparrow & \swarrow & \leftarrow & \swarrow \\ C & 0 & \uparrow & \swarrow & \uparrow & \swarrow & \uparrow & \swarrow & \leftarrow \\ A & 0 & \swarrow & \uparrow & \swarrow & \uparrow & \swarrow & \uparrow & \swarrow \end{pmatrix}$$

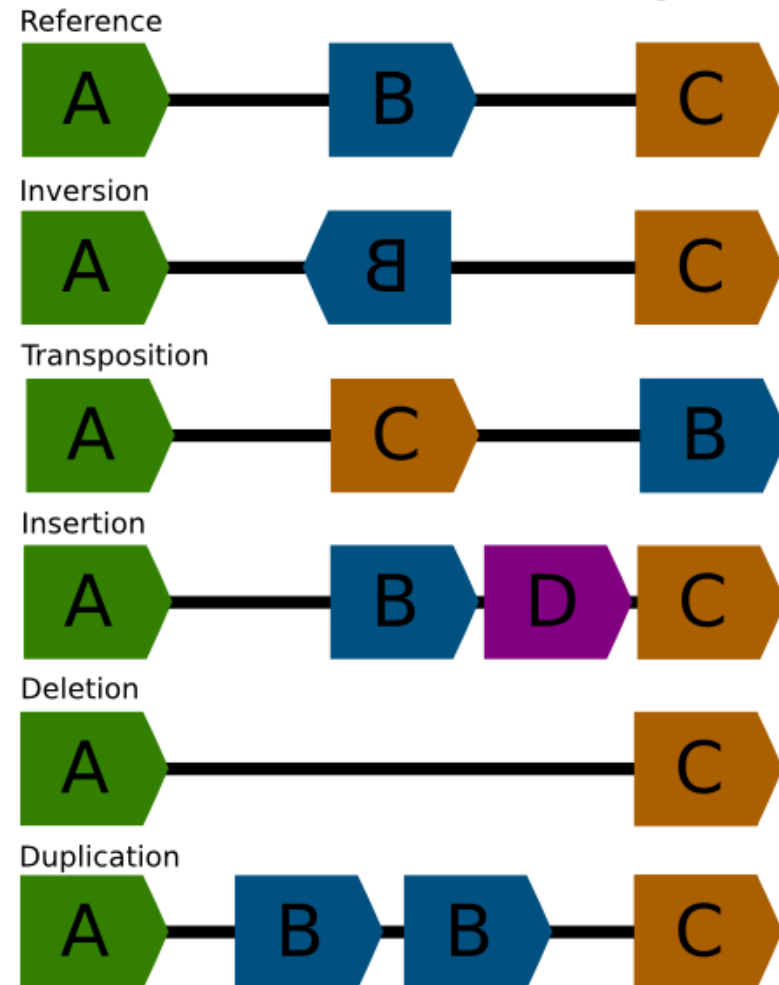
http://en.wikipedia.org/wiki/Smith-Waterman_algorithm



- Final result, an alignment file (BAM)



- Regions very different from reference sequence
 - Structural variations
 - Except for deletions and duplications



- Reference which contains duplicate regions
- Different strategies exist if multiple positions are equally valid:
 - Ignore read
 - Place at multiple positions
 - Choose one location at random
 - Place at first position
 - Etc.

- Example situation
 - 2 duplicate regions, one with a heterozygote variant



Alignment problems

- Map to first position

CTACTAGCGCAT ————— CTACTAGCGCAT

CTACTAGCGCAT
CTACTAGCGCAT
CTACTAGCGCAT
CTACTAGCGCAT
CTAC**G**AGCGCAT
CTAC**G**AGCGCAT
CTACTAGCGCAT
CTACTAGCGCAT

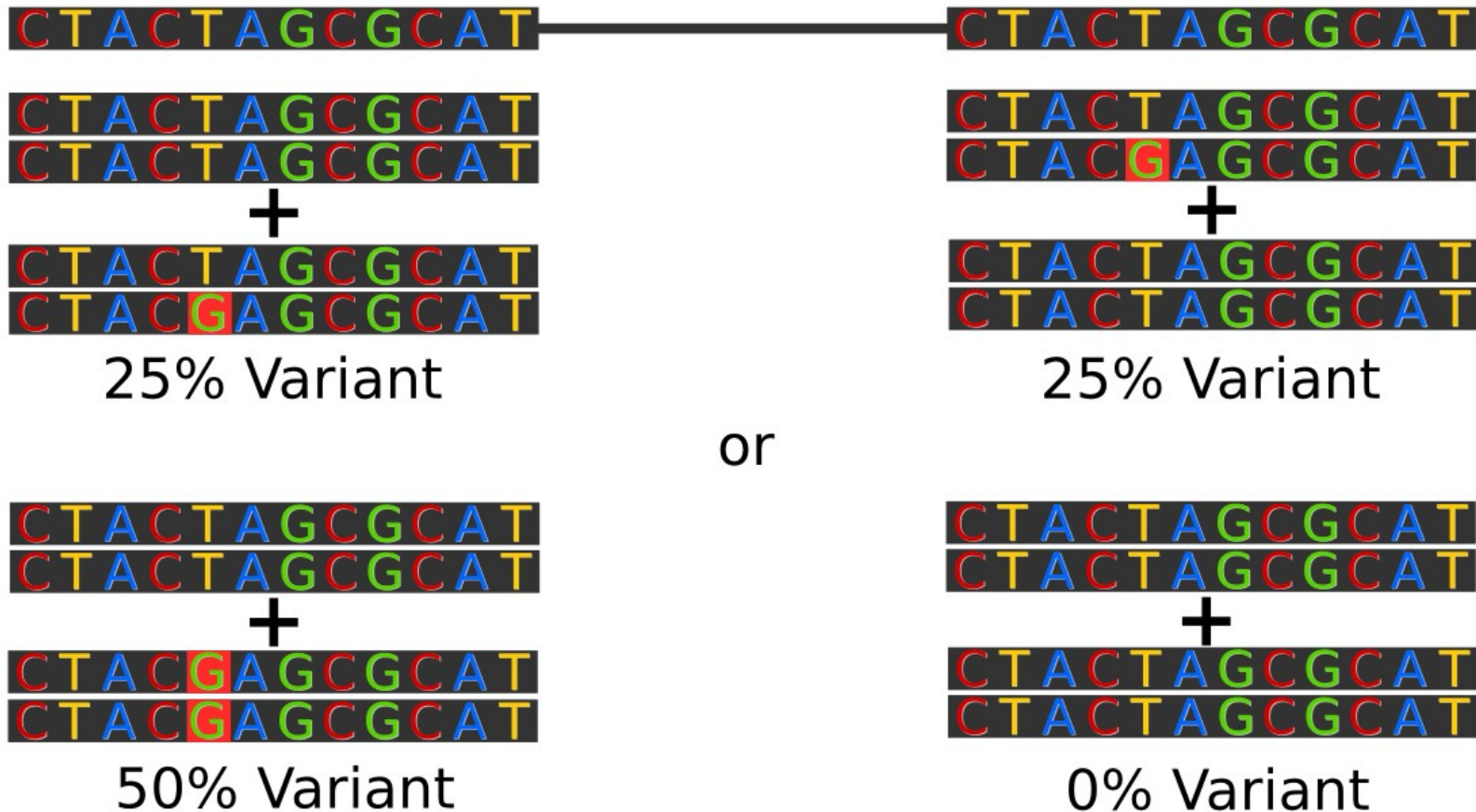
25% Variant

no data



Alignment problems

- Map to random position



Alignment problems

- To dustbin

CTACTAGCGCAT ————— CTACTAGCGCAT

deletion

deletion

- Sequences that are not aligned can be recovered in the dustbin
 - Sequences with no matching place on reference
 - Sequences with multiple possible alignments
- Several strategies exist to handle them
 - De novo assembly
 - Realigning with a different aligner
 - Etc.
- Important information can often be found there

- Popular aligners
 - Bowtie 1 + 2 (<http://bowtie-bio.sourceforge.net/>)
 - BWA (<http://bio-bwa.sourceforge.net/>)
 - BLAST (<http://blast.ncbi.nlm.nih.gov/>)
- Different strengths for each
 - Read length
 - Paired end
 - Indels

A survey of sequence alignment algorithms for next-generation sequencing. Heng Li & Nils Homer, 2010



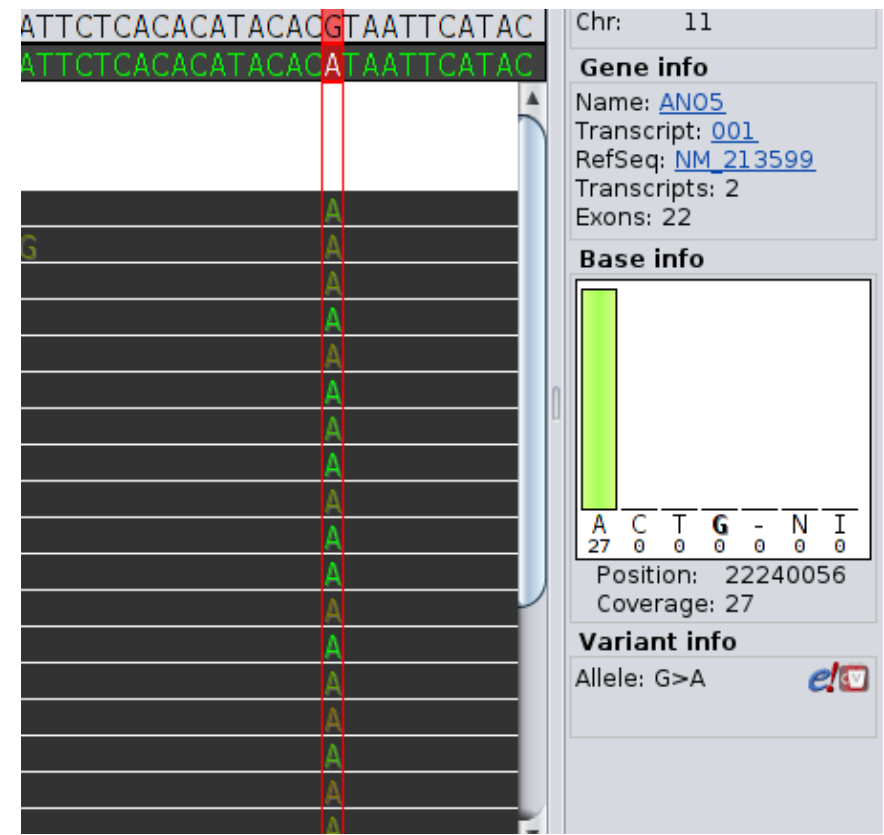
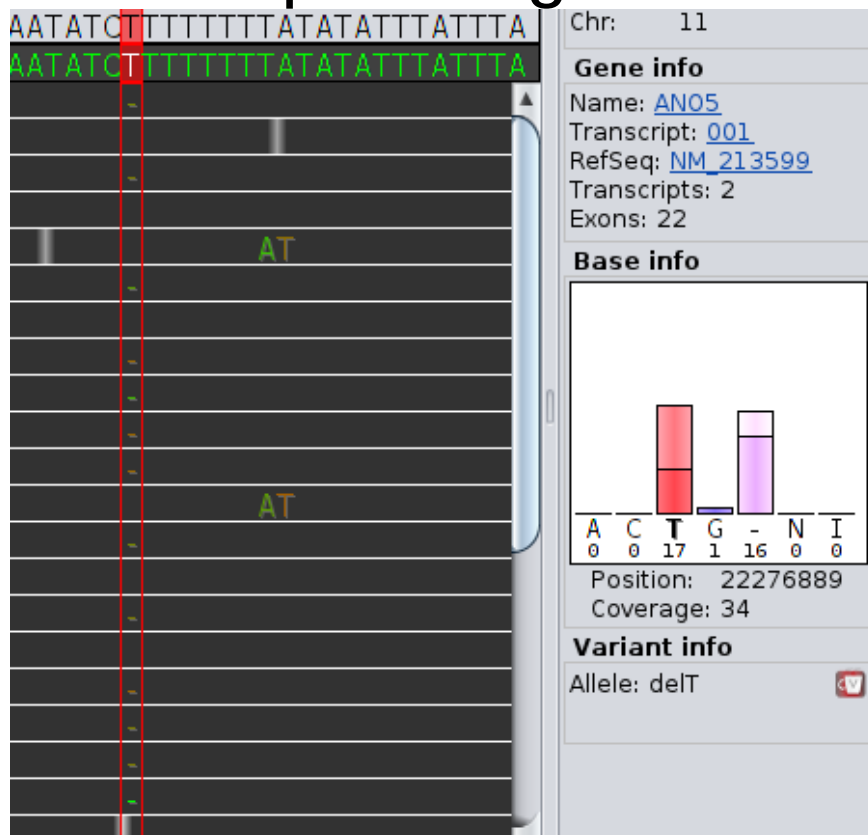
- Hybrid methods
 - Assemble contigs that are aligned back against the reference, many popular aligners can be used for this



- Reference aided assembly

- Difference in underlying data (alignment vs assembly) require different strategies for variant calling
 - Reference based variant calling
 - Patient comparison of de novo assembly
- Hybrid methods exist to combine both approaches
 - Alignment of contigs against reference
 - Local de novo re-assembly

- Reference based variant calling
 - Compare aligned reads with reference

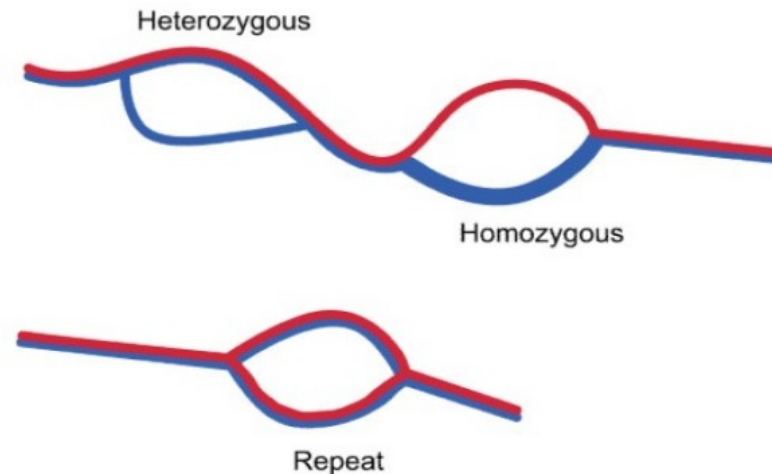


- Common reference based variant callers:
 - GATK
 - Samtools
 - FreeBayes

- Works very well for (in non repeat regions):
 - SNVs
 - Small indels

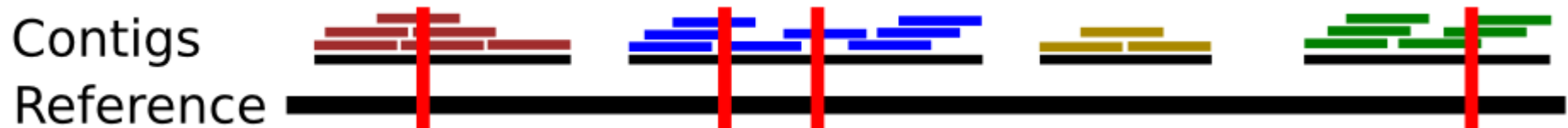
- De novo assembly
 - Either compare two patients
 - Useful for large structural variation detection
 - Can not be used to annotate variations with public databases
 - Or realign contigs against reference
 - Useful to annotate variants
 - Might lose information for the unaligned contigs

- Cortex
 - Colored de Bruijn graph based variant calling



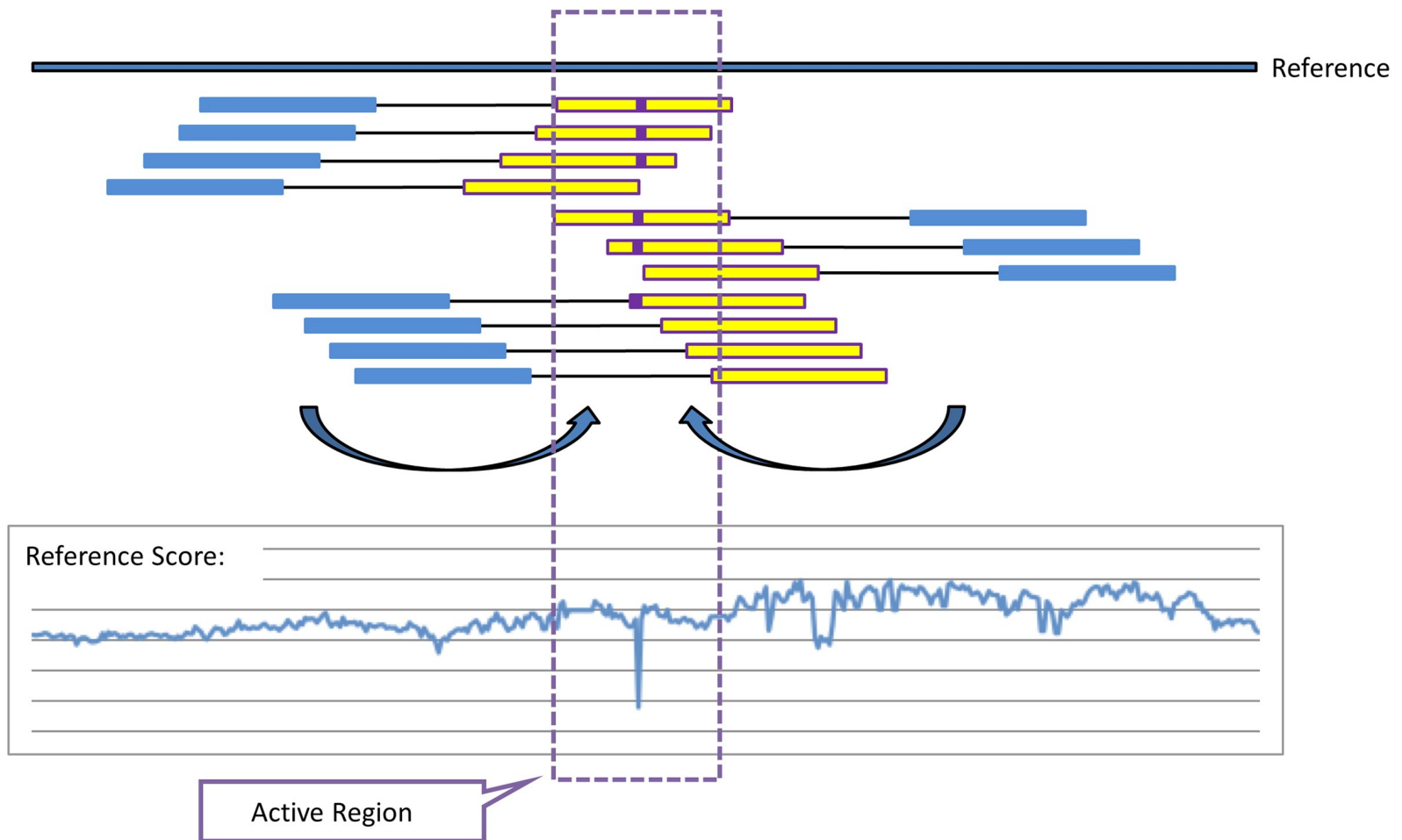
- Works well for
 - Structural variations detection

- Contig alignment against reference
 - Using aligners such as BWA
 - Uses standard reference alignment tools for variant detection
 - Helpful to „increase read size“ for better alignment
 - Variant detection is done using standard variant calling tools



- Local de novo assembly
 - Used by the Complete Genomics variant caller
- Every read around a variant is de novo assembled
- Contig is realigned back against the reference
- Final variant calling is done

Variant calling



Computational Techniques for Human Genome Resequencing Using Mated Gapped Reads, Paolo Carnevali et al., 2012

Variant calling

- Local de novo realignment allows for bigger features to be found than with traditional reference based variant calling
- Faster than complete assembly

- Reference based alignment
 - Good for SNV, small indels
 - Limited by read length for feature detection
 - Works for deletions and duplications (CNVs)
 - Using coverage information
 - Alignments are done “quickly“
 - Very good at hiding raw data limitations
 - The alignment does not necessarily correspond to the original sequence
 - Requires a reference that is close to the sequenced data

- De novo assembly
 - Assemblies try to recreate the original sequence
 - Good for structural variations
 - Good for completely new sequences not present in the reference
 - Slow and high infrastructure requirements
 - Very bad at hiding raw data limitations

- Unless necessary, stick with reference based alignment
 - Easier to use
 - More tools to work with the results
 - Easier annotation and comparison
 - Current standard in diagnostics
 - Can still benefit from de novo alignment through local de novo realignment
 - Analyze dustbin if results are inconclusive

- Transcriptomics, similar problematic to DNaseq
 - If small variations and gene expression analysis is done, alignment against reference is used
 - If unknown transcripts/genes are searched, de novo assembly is used
 - Used to detect transcripts with new introns, changed splice sites
 - Is able to handle RNA editing much better than alignment
 - Different underlying data (single strand, non uniform coverage, many small contigs)

- Reference based alignment is the current standard in diagnostics
- Assemblies can be used if reference based alignment is not conclusive
- Assembly will become much more important in the future when sequencing technologies are improved



Thank you for your attention

beat.wolf@hefr.ch

Further resources

Next Generation Variant Calling:

<http://blog.goldenhelix.com/?p=1434>

De novo alignment:

<http://schatzlab.cshl.edu/presentations/>

Structural variation in two human genomes mapped at single-nucleotide resolution by whole genome de novo assembly:

<http://www.nature.com/nbt/journal/v29/n8/abs/nbt.1904.html>

