# Under the Hood of Alignment Algorithms for NGS Researchers

April 16, 2014
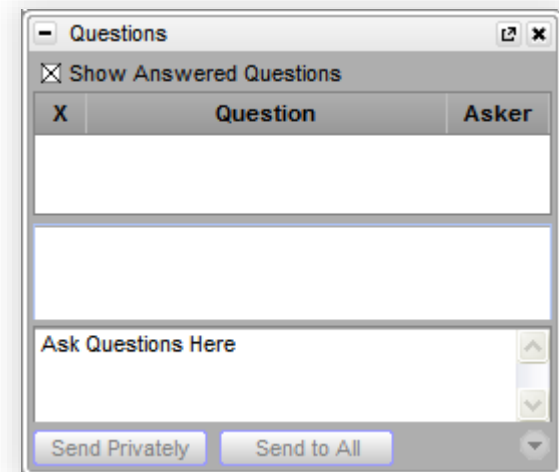
Gabe Rudy
VP of Product Development

Golden Helix

GOLDEN HELIX
*Accelerating the Quest for Significance*™

# Questions during the presentation

Use the Questions pane in your GoToWebinar window

# My Background

- **Golden Helix**
  - Founded in 1998
  - Genetic association software
  - Analytic services
  - Hundreds of users worldwide
  - Over 800 customer citations in scientific journals

- **Products I Build with My Team**
  - **SNP & Variation Suite (SVS)**
    - SNP, CNV, NGS tertiary analysis
    - Import and deal with all flavors of upstream data
  - **GenomeBrowse**
    - Visualization of everything with genomic coordinates. All standardized file formats.
  - **RNA-Seq Pipeline**
    - Expression profiling bioinformatics





GOLDEN HELIX
*Accelerating the Quest for Significance™*

# Agenda

| 1 | Alignment 101 |
| 2 | A Brief History of Time |
| 3 | Know Your CIGAR |
| 4 | It's All about the Variants |
| 5 | Q&A |

GOLDEN HELIX
Accelerating the Quest for Significance™

# Analytics and Sequencing

| Primary Analysis | | data |
|---|---|---|

**Alignment**

**Local Realignment**

**Variant Calling**

**Tertiary Analysis**

**"Sense Making"**

- Population structure analysis
- Genome browser-driven exploratory analysis
- Phenotypic association testing

Alignment 101

# Agenda

**1** Alignment 101

**2** **A Brief History of Time**

**3** Know Your CIGAR

**4** It's All about the Variants

**5** Q&A

GOLDEN HELIX
*Accelerating the Quest for Significance™*

# Types of Alignment

- Multiple Sequence Alignment

- Phylogenic analysis

- Database Search (BLAST)

- **Pairwise Alignment**
  - Local vs Global
  - Dynamic Programing vs Word Based

# Pairwise Alignment with Dynamic Programming

- **Needleman-Wunsch (1970)**

  - Dynamic programming optimal alignment of two sequences globally

  - O(n*m) space and time

  - Weighting function critical to define

    - Penalty matrix for mismatches

    - Penalty for gaps open and extensions (insertions, deletions)

- **Smith-Waterman (1981)**

  - NW based piecewise (local) alignment

  - Many optimizations, still O(n*m)

Needleman-Wunsch

match = 1     mismatch = -1     gap = -1

|   |   | G | C | A | T | G | C | U |
|---|---|---|---|---|---|---|---|---|
|   | 0 | -1 | -2 | -3 | -4 | -5 | -6 | -7 |
| G | -1 | 1 | 0 | -1 | -2 | -3 | -4 | -5 |
| A | -2 | 0 | 0 | 1 | 0 | -1 | -2 | -3 |
| T | -3 | -1 | -1 | 0 | 2 | 1 | 0 | -1 |
| T | -4 | -2 | -2 | -1 | 1 | 1 | 0 | -1 |
| A | -5 | -3 | -3 | -1 | 0 | 0 | 0 | -1 |
| C | -6 | -4 | -2 | -2 | -1 | -1 | 1 | 0 |
| A | -7 | -5 | -3 | -1 | -2 | -2 | 0 | 0 |

# Pairwise Alignment with Word Methods

- Heuristics methods based on finding matching k-tuples

- Significantly more efficient when the majority of sequence will not match (database search, reference-based alignment)

- FASTA (1985), BLAST (1990) designed for large DNA/Protein searches

- New class of problem emerged with high-throughput sequencing

# Alignment Versus Assembly

- **Assembly**

  - Orders of magnitude slower and memory intensive than alignment

  - Potentially compare every read with each other $O(n^2)$

  - Steps:

    - Merge overlapping reads into a de Bruijn graph

    - Simply the graph iteratively, construct contigs

    - Detangle with orthogonal tech (long reads, mates, optical mapping)

  - "Draft" genomes from short reads, have ~1kb sized contigs

- **Alignment**

  - Requires a finished genome for your species (draft genomes possible, but of limited utility)

  - Precompute an index of the reference genome (can be costly as you do it once)

  - Each short read uses the index to find its best placements (potentially multiple)
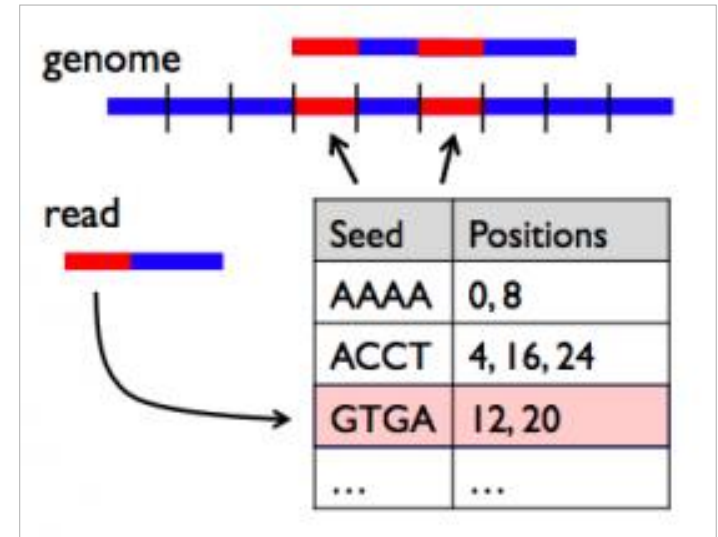


GOLDEN HELIX
Accelerating the Quest for Significance™

# Hash Based Alignment Algorithms

- **Hash based**
  - Pick k-mer size, build lookup of every k-mer in the reference to its positions
  - ~16GB of RAM required for hg19
  - **Seed-and-extend strategy**
  - **Popular tools:**
    - BLAST: tunable for different uses
    - MAQ (2008): Heng Li, et al
    - NovaAlign: Slower, but very accurate
    - Isaac (2013): High mem, but fast
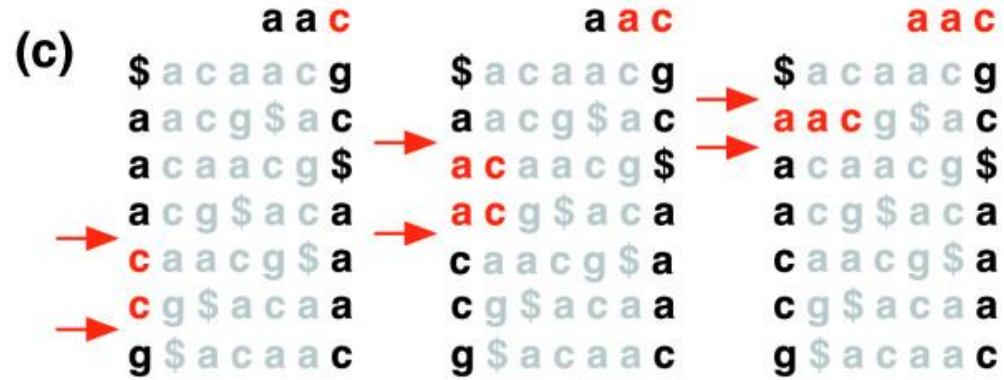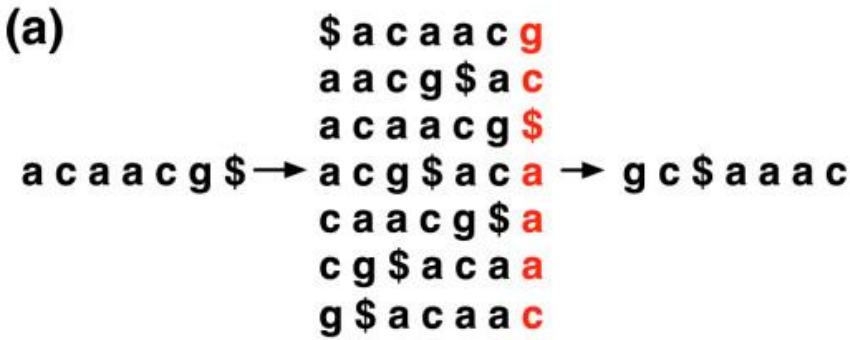    - MOSAIK (2014): Hash clustering+SW

# Burrows-Wheeler Transform

- BWT is a reversible permutation of characters that can be used for fast substring-searching when used with an index



GTTTGCGAA^TGTC$A

^TAGTCGAGGCTTTA$

1. All possible rotations

^TAGTCGAGGCTTTA$
TAGTCGAGGCTTTA$^
AGTCGAGGCTTTA$^T
GTCGAGGCTTTA$^TA
TCGAGGCTTTA$^TAG
CGAGGCTTTA$^TAGT
GAGGCTTTA$^TAGTC
AGGCTTTA$^TAGTCG
GGCTTTA$^TAGTCGA
GCTTTA$^TAGTCGAG
CTTTA$^TAGTCGAGG
TTTA$^TAGTCGAGGC
TTA$^TAGTCGAGGCT
TA$^TAGTCGAGGCTT
A$^TAGTCGAGGCTTT
$^TAGTCGAGGCTTTA

2. Sort

AGGCTTTA$^TAGTCG
AGTCGAGGCTTTA$^T
A$^TAGTCGAGGCTTT
CGAGGCTTTA$^TAGT
CTTTA$^TAGTCGAGG
GAGGCTTTA$^TAGTC
GCTTTA$^TAGTCGAG
GGCTTTA$^TAGTCGA
GTCGAGGCTTTA$^TA
TAGTCGAGGCTTTA$^
TA$^TAGTCGAGGCTT
TCGAGGCTTTA$^TAG
TTA$^TAGTCGAGGCT
TTTA$^TAGTCGAGGC
^TAGTCGAGGCTTTA$
$^TAGTCGAGGCTTTA

3. Select final column

^TAGTCGAGGCTTTAGATCCGATGAGGCTTTAGAGACAG$   Genomic sequence

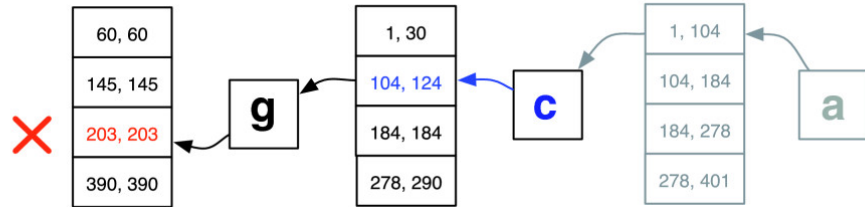GGTTGGTCGGATTCGGAATCACGGAAAATT^AGATTCC$G   Transform

# BWT Based Algorithms
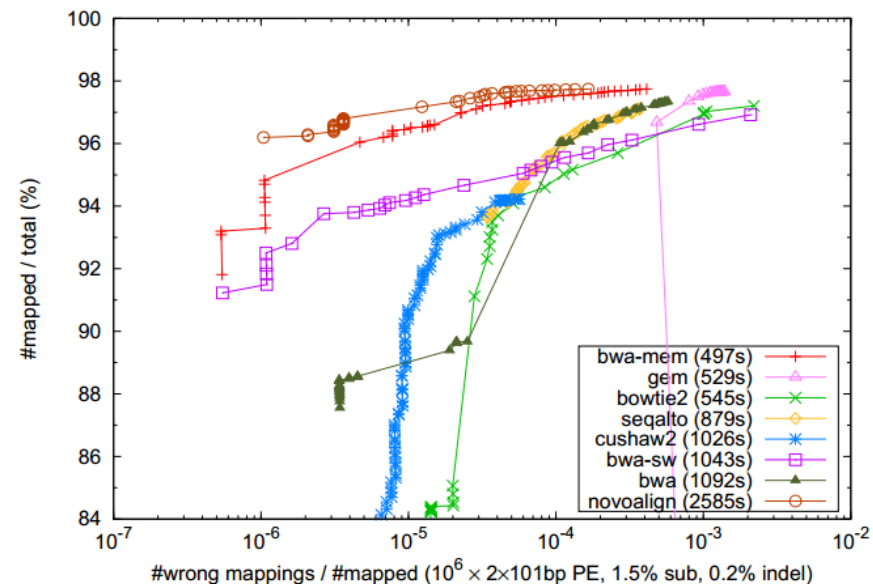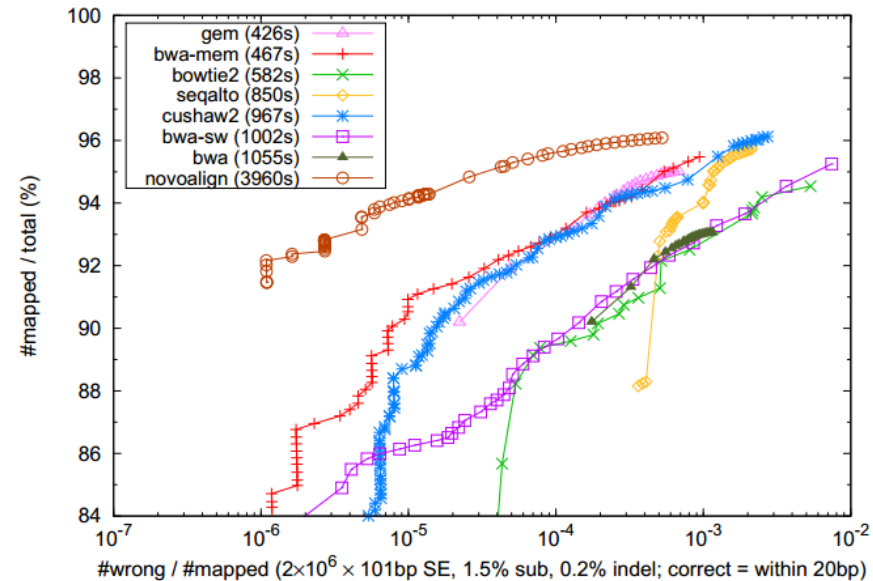
- **Compute a FM index of the reference**
  - Requires only ~1.5GB to hold in RAM of hg19

- **Requires a back-tracking algorithm to account for mismatches and gaps**

- **Designed for speed**
  - BowTie, BWA, SOAP2 (2009)
  - Order of magnitude less RAM and Time

- **More recent algorithms are often a hybrid:**
  - BWA-SW (2010)
  - Bowtie2 (2012)
  - BWA-MEM (2013)

# BWA and Friends

- **BWA (backtrack) - 2009**
  - Very mature, handles short reads up to 100bp

- **BWA-SW (Smith Waterman) - 2010**

- **BWA-MEM (Max Exact Matches) - 2013**
  - >70bp read length recommended, but up to 1Mbp
  - Seed and extend with SW
  - Allowable error rate adjust with sequence length
  - Finds larger gaps
  - Faster! Generally supersedes BWA-SW

# Algorithm Comparison

| | BWA | BWA-SW | BWA-MEM | Bowtie | Bowtie2 | NovaAlign | MOSAIK | Isaac | Tmap |
|---|---|---|---|---|---|---|---|---|---|
| Affiliation | | Heng Li | | U of Maryland | | Novacraft | Boston College | Illumina | Ion Torrent |
| First Published | 2009 | 2010 | 2013 | 2009 | 2012 | - | 2014 | 2013 | - |
| Read Length | <100 | 70bp-1Mbp | | <100 | >50 | | | | |
| Gapped Alignments | | | | No | | | | | |
| Trimming | | | | No | | | | | |
| Error Rates Allowed | Low | High | Med | Low | Med | Med | High | Low | Med |
| Chim Reads | No | Yes | **Yes** | No | Opt | Opt | **Yes** | No | No |
| Mem Usage | Med | Med | Med | Low | Low | Low | Med | **High** | Med |
| Speed | Med | Med | Fast | Fast | Fast | **Slow** | Fast | Fast | Fast |

# Agenda

| 1 | Alignment 101 |
| 2 | A Brief History of Time |
| 3 | **Know Your CIGAR** |
| 4 | It's All about the Variants |
| 5 | Q&A |

GOLDEN HELIX
*Accelerating the Quest for Significance™*

# SAM/BAM

- **Spec defined by bwa/samtools author** Heng Li, aka Li H, aka lh3.

- **SAM is text version** (easy for any program to output)

- **BAM is binary/compressed version** with indexing support

- Alignment encoded in CIGAR code of matches, insertions, deletions, gaps and clipping

- **Can have any custom flags set** by alignment tool (mix of standard and custom two-letter tags)

```
@HD VN:1.3 SO:coordinate
@SQ SN:ref LN:45
r001 163 ref  7 30 8M2I4M1D3M = 37  39 TTAGATAAAGGATACTG *
r002   0 ref  9 30 3S6M1P1I4M *  0   0 AAAAGATAAGGATA    *
r003   0 ref  9 30 5H6M        *  0   0 AGCTAA    *   NM:i:1
r004   0 ref 16 30 6M14N5M     *  0   0 ATAGCTTCAGC       *
r003  16 ref 29 30 6H5M        *  0   0 TAGGC     *   NM:i:0
r001  83 ref 37 30 9M          =  7 -39 CAGCGCCAT         *
```

## Key Fields

- Chr, position
- **Mapping quality**
- **CIGAR**
- Name/position of mate
- Total template length
- Sequence
- Per-Base Quality Scores

# CIGAR String

| Op | Description | Used |
|---|---|---|
| M | alignment match (can be sequence match or mismatch) | by default |
| I | insertion to the reference | InDel |
| D | deletion from the reference | InDel |
| N | skipped region from the reference | spanning intron |
| S | soft clipping (not-aligned) | per-base quality drops to threshold to trim read |
| H | hard clipping (not-present in reference) | chimeric reads, breakpoints, end of seq |
| P | padding (silent deletion from padded reference) | multiple sequence alignments (not common) |
| = | sequence match | when compared to ref |
| X | sequence mismatch | when compared to ref |



My Exome (BWA/FreeBayes)   Father_Realigned

Pile-up

| Cigar Ops | M4 D1 M72 |
| Adjusted Cigar Ops | =4 D1 =7 X1 =18 X1 =44 X1 |

Different Alignment Outcomes

# Agenda

| 1 | Alignment 101 |
|---|---|

| 2 | A Brief History of Time |
|---|---|

| 3 | Know Your CIGAR |
|---|---|

| 4 | **It's All about the Variants** |
|---|---|

| 5 | Q&A |
|---|---|

# Mapping and Calling Variants on the Human Genome

- **Classes of Confounders:**
  - Issues with the **Reference Assembly**:
    - Sequence under-represented (exact match not in human reference, so get poor match)
    - Tiling issues creating artificial splices
  - **Repeated** regions and Low **Mapping Quality** Regions:
    - Over 50% of the genome is repetitive
    - Low sequence "complexity" or "information density" means short reads cannot uniquely map. "Mappability"
  - Interference with larger classes of variation: **Structural Variation**
    - Calling genotypes of SNPs/short-InDels in a deletion
    - Inversion/Translocation/CNV break points
  - Disagreement in **Representing Complex Variants**

# Responsibility of the Alignment Algorithm?

- **Placing reads in the right part of the genome**

- **Providing accurate mapping quality scores**
  - Often need to empirically train an aligner to produce Gaussian spectrum of scores

- **Providing the best data to variant callers**

- **What Variant Callers expect?**
  - Multi-mapped read status (often filtered out by MQ=0)
  - Mate-pair mapping information
  - Just "localizing" the read?
  - Consistently described gapped alignments?

# GRCh38 – Here Now, but still Waiting

- **A better human reference**
  - Revised Cambridge Reference Sequence (rCRS) MT
  - Has centromere models
  - ~2000 incorrect alleles fixed
  - ~100 assembly gaps updated

- **No Gene Annotations**
  - RefSeqGene - Feb 2014
  - Ensembl Q4 2014

- **No Variant Annotations**
  - Re-align 1000 Genomes and NHLBI 6500?
  - dbSNP?

## My Exome

Bar chart showing variant counts for GRCh37 (total 331,824) and GRCh38 (total 319,442), broken down into snps, mnps, indels, and complex categories. Y-axis ranges from 270000 to 340000.

| | GRCh37 | GRCh38 |
|---|---|---|
| Ts/Tv | 2.06558 | 2.10171 |

# InDel Alignment: Watch for ambiguities

# InDel Alignment: Watch for repeats and read ends

# MNP vs Allelic Primitives

# Genome In a Bottle

- NIST Sponsored, Community Effort

- NA12878 cell line, sequenced many platforms, read lengths and sample preps

- Create "ensembl" variant call set

- Create many making regions

  - Regions not able to make consensus call

  - Repeat and low-complexity regions

  - SV in NA12878

- Variants, BED and alignment data available

Justin Zook

Genome in a Bottle Consortium

# Resources

- **GCAT**

- **Benchmarks**
  - Alignment
  - Variant Calling
  - GIAB Truth Set
  - Various bench samples

- **Interactive filtering**

Some GIAB Examples

# Questions?

Use the Questions pane in your GoToWebinar window