Assembly and Alignment

Workshop on Comparative Genomics

İnanç Birol, Shaun Jackman

Smithsonian Institution – 5 October 2011





Assembly and Alignment

Workshop on Comparative Genomics

İnanç Birol, Shaun Jackman

Smithsonian Institution – 5 October 2011





Moore's Law





Growth of Knowledge



Next Generation Sequencing

 Illumina sequencing throughput at GSC Cost of sequencing human genome







- Old paradigm:
 - long and non-uniform reads (800bp 1000bp)



- Old paradigm:
 - long and non-uniform reads (800bp 1000bp)
 - overlap; overlay; consensus



- New paradigm:
 - short and uniform reads (50bp 150bp)





- New paradigm:
 - short and uniform reads (50bp 150bp)
 - overlap; overlay; consensus V



- New paradigm:
 - short and uniform reads (50bp 150bp)
 - , de Bruijn graphs



- New paradigm:
 - long range information through read pairs
 - graph theoretic approaches

Assembly By Short Sequencing

Resource	IFFE InfoVis 2009
ABySS: A parallel assembler for short read sequence data	ABySS-Explorer: Visualizing Genome Sequence Assemblies
Jared T. Simpson, ¹ Kim Wong, Shaun D. Jackman, Jacqueline E. Schein, Steven J.M. Jones, and Inanç Birol ² Genome Sciences Centre, British Columbia Cancer Agency, Vancouver, British Columbia VSZ 4E6, Canada Widespread adoption of massively parallel deoryribonucleic add (DNA) sequencing instruments has prompted the resent development of de novo short read assembly algorithms. A common shortcoming of the available took is their in ability to difficuent to the novo short read assembly algorithms. A common shortcoming of the available took is their in ability to difficuent for the novo short read assembly algorithms. A common shortcoming of the available took is their in ability to difficuent to the novo short read assembly algorithms.	Cydney B. Nielsen, Shaun D. Jackman, Inanç Birol, and Steven J.M. Jones
Jackman and Birol Genome Biology 2010, 11:202 http://genomebiology.com/2010/11/1/202	 ABySS-Explorer employs a novel graph representation enabling biologists to examine the global structure of a genome ence assembly. Charles Content and Content a
Centron Infector Sequent Sequent Sequent Sequent Sequent Sequent Sequent Sequent Shaun D Jackman and Inanç Birol [®]	sequencing BRIEF COMMUNICA
Original Paper De novo Transcriptome Assembly with ABySS Inanç Birol ^{1,*} , Shaun D Jackman ¹ , Cydney Nielsen ¹ , Jenny Q Qian ¹ , Richard Varhol ¹ , Stazyk ¹ , Ryan D Morin ¹ , Yongjun Zhao ¹ , Martin Hirst ¹ , Jacqueline E Schein ¹ , Doug J man ³ , Joseph M Connors ² , Randy D Gascoyne ² , Marco A Marra ¹ and Steven JM Jone ¹ Genome Sciences Centre, 100-570 W 7th Avenue, Vancouver BC V5Z 486, Canada, www.bcga.ca ² British Columbia Cancer Agency, 600 West 10th Avenue, Vancouver, BC V5Z 4E6, Canada, www.bccancer Associate Editor: Dr. Alex Bateman	Are hundred and by the set of the





• SE Assembly: k-mer extension on a de Bruijn graph



 PE Assembly:search for unambiguous contig merging along paths



 Scaffolding: search for unambiguous linkage across distant contigs





naturenews

news archive

News



Stories by subject

Genome sequence

Genome assembly

Blogs linking to this

Add to Connotea

Add to Facebook

A DE LE MELLER

Add to Digg

Bioinformatics

This article

article

elsewhere

Genome sequencing

Genetics

Technology

Genomics

nature news home

Published online <u>23 March 2011</u> | *Natur*e **471**, 425 (2011) | doi:10.1038/471425a

Genome builders face the competition

specials

Three independent projects seek to contrast approaches in preparation for routine analysis of genetic data.

opinion

features

news blog

Stories by keywords Erika Check Hayden

Sequencing DNA on an industrial scale is no longer difficult: the challenge is in assembling a full genome from the multitude of short, overlapping snippets that second-generation sequencing machines





08 July 2011

Related stories

nature journal

etly embrace whole-genome 2010 es beyond DNA sequence uencing: the third generation 009

At a meeting last week at the University of California, Santa Cruz, three winners emerged: ALLPATHS-LG, developed by the Broad Institute in Cambridge, Massachusetts; ABySS, developed at the British Columbia Cancer Agency's Genome Sciences Centre in Vancouver, Canada; and SOAPdenovo, developed by the Beijing Genomics Institute. But, Korf notes, "it's not just the software, it's how people are running it" that determines the quality of each assembly.



Assembly Problem

TCGATCGATTTTCGGCCTAA read1 ATTTTCGGCCTAATATTAGG read2

...GCATCGATCGATTTTCGGCCTAATATTAGGCCGATAATCGACGATC...

A <u>partial</u> and <u>unambiguous</u> read-to-read alignment extends the length of sequence information

- First stage of an assembly algorithm is to find such alignments
- Assembly algorithms differ in the way they find and use these alignments



Greedy Assembly

- Find two reads with the largest overlap
- Merge them
 Repeat until no more

Pro: fast

Con: prone to misassembly

• Assumes largest overlaps are unambiguous



Overlap Overlay Consensus

• Overlap

Find all pairs of sequences that overlap

- Overlay (a.k.a. Layout) Remove redundant and weak overlaps
- Consensus

Merge pairs of sequences that overlap unambiguously Build a consensus sequence from all reads overlaid in a region



Find Overlapping Reads

- Naïve algorithm: make all binary comparisons Untenable when too many reads ARACHNE
 - $-O(n^2)$
 - RAM
 - CPU
- Ferragina-Manzini index
 - Apply Burrows-Wheeler transform
 - Small memory footprint SGA
- Build an overlap graph

ARACHNE CAP3 Celera assembler MIRA Newbler Phred/Phrap

Forget About Overlapping Reads!

- Shred reads to a uniform length k
- Build a special overlap graph: de Bruijn Graph



Euler Velvet ABySS SOAPdenovo ALLPATHS

...<u>GCATCGATCGATTTTCGGCCTAATATTAGGCCGATAATCGACGATC</u>...



De Bruijn Graph

- Load *k*-mers in memory
 - 2x4 possible extension of every k-mer
- Check if there is a "next" k-mer
 - -O(n) algorithm
 - ...<u>GACATTGC</u>... seq1 ...<u>GACATTAT</u>... seq2



Memory Concerns

- Human genome has over 2 billion unique *k*-mers
- If we represent every *k-mer* using, say 50 bytes
 we require over 100 GB RAM
 just to represent *k*-mers

Solution #1: Clustering reads

Curtain (w/ Velvet) Phusion (w/ Phrap)

Solution #2: Distributed computing

ABySS SOAPdenovo ALLPATH-LG



Partitioning Read Space

Distribute sub-reads and reverse-complements over nodes





Graph Generation

- A given k-mer can have up to 8 extensions
- Each node announces the list of k-mers that it has to the nodes that hold their possible extensions
- Each node records if there are any extensions of the k-mers that it stores



• This forms adjacency information for *k*-mers over a distributed de Bruijn graph



Trimming

- Data would have experimental noise
- de Bruijn graph would have false branches
- Some read errors are filtered by removing such branches
- Trimming prevents the later assembly step to come to a premature end because of read errors





Bubble Popping

- Repeat read errors and single nucleotide allelic differences would cause "bubbles" of length 2k-1
- Bubbles are popped by removing either of those branches
- Complex bubbles can form when multiple bubbles intersect
 - Bubble popping step either reduces the bubble orders by one
 - Or creates dead branches
- Popped bubbles are recorded in a log file to study potential allelic differences







Assembly - SET

- Remaining de Bruijn graph is analyzed for contig extension ambiguities
- If there is a multiplicity in the inbound or outbound contig extensions, then contig growth is terminated



 SET assembly step then concatenates the remaining connected nodes in the di-graph, creating independent contigs that overlap by no more than k-1 bases



Assembly - PET

- After SET assembly, reads are aligned to contigs
- Using reads that hit the same contig, empirical fragment size distribution(s) is (are) calculated
 - calculated Using reads that hit multiple contigs, inter-contig distances are inferred with a maximum likelihood estimator
- Contigs with coherent and unambiguous distances are joined





Adjacency Graph

- SET assembly result as a graph
 - Nodes: contigs
 - Edges: overlaps (k-1 bp)







Adjacency Graph

- SET assembly result as a graph
 - Nodes: overlaps (k-1 bp)
 - Edges: contigs







Assembly As a Hairball





ABySS-Explorer: Visualizing Genome Sequence Assemblies

Cydney B. Nielsen, Shaun D. Jackman, Inanç Birol, and Steven J.M. Jones





Paired End Tag Information





Paired End Contigs



Blue gradient: SET contig path in PET assembly Orange: selected SET contig



ABySS-Explorer GUI



Single-end contig id: 5- (32 bp; 61 kmer cov)



Paired-end contig id: 1829+

Single-end contig members: 498+ 1077-1621-1696-239+1438+291+1638-80+1181+1007+1045+1487-1679-612+129+152+1591+792+891-739-636-523+344+ 714-396-1275-158+434+1518+1437-1360-1028-1315-1251-923+112+407+1724+753+1411-155-121-1680+1431-1510+1395+789+1724+407+112+923+ 1254+1136-605-1435-897-21-1145-1165-1626+401+194+1369-1576-1515+401+194+1369-1576-1515+401+194+1178+349-714-396-1275-158+433+ 1517+516+71+1709+1026-1644-1751+1557+1415-1283-1361-5-452+1266+1748-

Statistics Display



Nxx plot and N50



- The N50 is the weighted median of contig sizes
- The N50 summarizes a single point on the Nxx plot
- Better assemblies are further to the right
k-mer Coverage Histogram

- Counts the number of occurrences of each k-mer
- Useful for
 - estimating the genome size
 - measuring mean coverage
 - library quality control





ABySS — Canada's Michael Smith Genome Sciences Centre - Mozilla Firefox					
<u>F</u> ile <u>E</u> dit <u>V</u> iew Hi <u>s</u> tory <u>B</u> ookmark	s <u>T</u> ools <u>H</u> elp	5 ¹			
🗢 🔿 🗸 🖗 🖉 👘 http	://www.bcgsc.ca/platform/bioinfo/software/abyss	☆ 🗸 Google 🍭			
BC Cancer Ag CARE + RESEARCH An agency of the Provincial Health	Log in Search Site Only in current section				
Home Platforms You are here: Home → Platfor	Projects Data Training Services Faculty	Careers About Project Resources			
	Accomply Dy Short Servences of de neve	- Delegere			
Bioinformatics Bioinformatics Licenses	parallel, paired-end sequence assembler	Support Contact address			
GSC Software Centre	Current release				
PASsiT		Project owner: Shaun Jackman			
Adapter Trimming for Small RNA Sequencing	ABySS 1.3.0 Released Sep 09, 2011	Subscribe to updates for this project			
Spark	Mate-pair data can be used to scaffold contigs. Specify your	Assembly Algorithm			
TASR	mate-pair libraries using the `mp' parameter of abyss-pe.	1) Persisioning Read Space 2) Adjacency Generation			
XpressAlign: FPGA Short Read Aligner	More about this release Get ABySS for all platforms (498 kB) Source	Respondence de la construir de			
Anchor	Get ABvSS for Linux (1.2 MB)	3) Trimming 4) Bubble Popping			
BLISS	Debian package (amd64)	A this work from the second for			
MiRNA Profiling		Successful and instrument of the second seco			
ORegAnno: Open Regulatory Annotation	Project Description	Change generative filt loss of the second seco			
SNVMix	sequence assembler that is designed for	Constantiante antegran la constantiante a nstantiante a constantiante a constantiante a constantiante a c			
SliderII	short reads. The single-processor				
ABySS	version is useful for assembling genomes up to 100 Mbases in size. The parallel version is implemented using MPI and is capable.				
Releases	of assembling larger genomes.	Screenshot — click to enlarge			
ABySS-Explorer	To assemble transcriptome data, see Trans-ABySS.				
Chinook	Publications	•			
Done					



ABySS

		////////==×			
<u>F</u> ile <u>E</u> dit <u>V</u> iew <u>G</u> o	<u>B</u> ookmarks <u>H</u> elp				
Back Forward		🤣 🙀 Reload Home	Computer Search		
👔 🕻 🔯 ibirol tr	abyss-1.3.0			🔍 100% 🔍	View as Icons 😽
Places → ×					
🐼 ibirol 😻 Desktop	ABYSS	AdjList	Align	Assembly	Common
🗇 File System		B			
	Consensus	DAssembler	DataLayer	DistanceEst	FMIndex
	Graph	KAligner	Мар	MergePaths	Overlap
4	Parallel	ParseAligns	PathOverlap	PopBubbles	Scaffold
	SimpleGraph	bin	dialign	doc	kmerprint
	COPYRIGHT	ChangeLog	LICENSE	Makefile.am	Makefile.in
				2-	
	README	aclocal.m4	config.h.in	configure	configure.ac
	2-			2-	
	depcomp	doxygen.conf	install-sh	missing	
39 items, Free space: 673	3.1 GB				





Assembly Operations

- SET contig building: de Bruijn
 - k-mer overlap information
- SET error removal: adjacency
- PET contig merging: adjacency & linkage
 PET alignments
- PET/MPET scaffolding: adjacency & linkage
 - PET/MPET alignments
- Gap closure and contig extensions: read overlap
 - PET alignments



1. Erode low-coverage tips

-e, --erode=COVERAGE

erode bases at the ends of blunt contigs with coverage less than this threshold

-E, --erode-strand=COVERAGE

erode bases at the ends of blunt contigs with coverage less than this threshold on either strand





2. Trim tips

-t, --trim-length=TRIM_LENGTH

maximum length of dangling edges to trim





3. Remove low coverage contigs

-c, --coverage=COVERAGE

remove contigs with mean k-mer coverage less than this threshold





4. Pop bubbles

-b, --bubbles=N

pop bubbles shorter than N bp (default: 3*k)

-b0, --no-bubbles

do not pop bubbles





1. Resolve forks





2. Trim tips





3. Remove repeats





4. Remove transitive edges





5. Trim tips





6. Pop bubbles





7. Remove weak edges





Running ABySS

- Assemble the paired-end reads in the file reads.fa
 - > abyss-pe name=ecoli k=32 n=10
 in=reads.fa
- Assemble the paired-end reads in the files reads_1.fa and reads_2.fa:

> abyss-pe name=ecoli k=32 n=10
in='reads_1.fa reads_2.fa'



Running ABySS in Parallel

- Run ABySS using eight threads
 - > abyss-pe np=8 name=ecoli k=32
 n=10 in='reads_1.fa reads_2.fa'
- ABySS uses MPI, the Message Passing Interface. OpenMPI is an open-source implementation of MPI



Running Parallel Jobs on a Cluster

- Run ABySS on a cluster using 8 threads
 - > qsub -pe openmpi 8 -N ecoli abyss-pe np=8 name=ecoli k=32 n=10 in='reads_1.fa reads_2.fa'
- abyss-pe uses the environment variables *JOB_NAME* and *NSLOTS* passed to it by SGE as the default values for *name* and *np*



Running for Multiple k Values

- Assemble every 8th *k* from 32 to 96
 - > qsub -pe openmpi 8 -N ecoli
 -t 32-96:8 abyss-pe k=32 n=10
 in='reads_1.fa reads_2.fa'
- abyss-pe uses the environment variable
 SGE_TASK_ID passed to it by SGE as the default value for k



Assembling Multiple Libraries

> abyss-pe name=ecoli k=32 n=10 lib='pe200 pe500' pe200='pe200_1.fa pe200_2.fa' pe500='pe500_1.fa pe500_2.fa'



Assembling a Mix of PET and SET

> abyss-pe name=ecoli k=32 n=10 lib='pe200 pe500' pe200='pe200_1.fa pe200_2.fa' pe500='pe500_1.fa pe500_2.fa' se='long.fa'



Parameters of ABySS

- name: name of the assembly
- lib: name of the libraries (one or more)
- se: paths of the single-end read files
- \${lib}: paths of the read files for that library
- Example

> abyss-pe name=ecoli k=32 n=10
lib='pe200 pe500'
pe200='pe200_1.fa pe200_2.fa'
pe500='pe500_1.fa pe500_2.fa'
se='long.fa'



Parameters of ABySS (SET)

- k: the size of a k-mer
- q: quality trimming removes low-quality bases from the ends of reads
- e and c: coverage-threshold parameters
 - e: erosion removes bases from the ends of contigs
 - c: coverage threshold removes entire contigs
- **p**: the minimum identity for bubble popping



Parameters of ABySS (PET)

- s: the minimum size of a seed contig
- n: the number of pairs required to join two contigs
- Example

> abyss-pe name=ecoli k=64 q=3 p=0.9 s=100 n=10 lib='pe200 pe500' pe200='pe200_1.fa pe200_2.fa' pe500='pe500_1.fa pe500_2.fa' se='long.fa'



Optimizing k

- Assemble every 8th k from 32 to 96
 Nine assemblies: 32 40 48 56 64 72 80 88 96
- Find the peak
- Assemble every 2nd k around the peak
 For example, if the peak were at k=64...
 Eight assemblies: 56 58 60 62 66 68 70 72
- SGE:
 - > qsub -t 32-96:8 qsub-abyss.sh
 - > qsub -t 56-72:2 qsub-abyss.sh



Output Files of ABySS

\${name}-contigs.fa

The final contigs in FASTA format

- \${name}-bubbles.fa
 The equal-length variant sequences (FASTA)
- \${name}-indel.fa
 The different-length variant sequences
 (FASTA)
- \${name}-contigs.dot
 The contig overlap graph in Graphviz format



Intermediate Output Files of ABySS

- .adj: contig overlap graph in ABySS adj format
- .dist: estimates of the distance between contigs in ABySS dist format
- .path: lists of contigs to be merged
- .hist: fragment-size histogram of a library
- coverage.hist: k-mer coverage histogram



Case Study

Mountain Pine Beetle Genome Assembly



Mountain Pine Beetle Genome

Assembly statistics

	contigs	scaffolds
n	1,128,463	1,103,221
n:500bp	33,591	11,657
n:N50	4,324	82
N50 (bp)	11,220	541,443
Max (bp)	276,135	3,583,207
Reconstruction (Gb)	201.9	200.4







Assembly As a Hairball

- ABySS v1.2.7
 - PET/MPET information disambiguates short contig extensions









• Contig 4 is (eventually) followed by Contig 7





Fragment size distribution



Fragment size distribution



Biotin Read-Through





Illumina's




Triage of MPET Reads



Information:

- Distances from contig ends
- Base mismatches on read ends
- Inferred contig orientations

Triage of MPET Reads





Scaffolding





Anchor

• Scrubbing "homozygous" variations



Indel (2,935)

SNPs (19,715)



www.bcgsc.ca

Anchor

- Local directional assembly
 - scaffold gap filling

(10,499 of 63,986)



extension



(20,213 of 53,487)



Quality Assessment

Alignment of 81,047,980 reads

	Before Anchor	After Anchor	Change
Mapped	65,624,456 (80.97%)	66,949,341 (82.60%)	+ 1,324,885
Paired	43,207,118 (53.31%)	44,732,320 (55.19%)	+ 1,525,202
Single-end	9,536,178 (11.77%)	8,846,977 (10.92%)	-689,201

Gene alignments

	2,180 ESTs		248 Conserved Genes		
	Complete	Partial	Complete	Partial	
Contigs	968	1169	212	18	
Scaffolds	1,481	619	228	5	



Final Hairball

- ABySS v1.2.7
 - Read pairs and inferred distances allow for scaffolding

	contigs	scaffolds
n	1,128,463	1,103,221
n:500bp	33,591	11,657
n:N50	4,324	82
N50 (bp)	11,220	541,443
Max (bp)	276,135	3,583,207
Reconstruction (Gb)	201.9	200.4





Date	ABySS Version	Data	n:500	N50	Max	Sum
August 2009	1.0.11	3x GAiix	81,431	1,526	20,755	107.3e6
November 2009	1.0.15	+2x GAiix	104,958	2,333	55,845	195.8e6
February 2010	1.1.1	+4x GAiix	157,081	2,790	136,637	346.3e6
July 2010	1.2.0	+2x GAiix	146,313	3,354	129,008	376.2e6
November 2010	1.2.4	+1x GAiix +1x GAiix (MPET)	100,690	4,474	294,323	268.8e6
May 2011	1.2.7		18,660	108,158	1,908,773	201.4e6
July 2011	1.2.7	+ 1x HiSeq +1x HiSeq (MPET)	11,657	541,443	3,583,207	200.4e6
August 2011	1.2.7		11,523	561,847	3,746,698	206.5e6



Future Work

- Clean up the chaff
 - Place short contigs on Anchored scaffolds
 - Annotate repeat elements



Transcriptome Assembly



Transcriptome Sequencing

- RNA-seq protocol
- Brings information on how a genome "acts"
 - Expression levels
 - Allelic expression
 - Present isoforms
 - Gene fusions
 - Other transcriptional events
 - Post-transcriptional RNA editing





Rodrigo Goya

Transcriptome Assembly



Transcriptome assembly is different from genome assembly

- − varying coverage levels
 ⇒ varying expression levels
- split assembly paths
 ⇒ isoforms/splice variants
- small contig sizes
 ⇒ small product sizes

What Overlap to Choose?





What Overlap to Choose?

- Selection of parameter k depends on read coverage depth
- Expression levels vary over 5 orders of magnitude



Selection of k





Assembly Merging







Multi-k Assembly

We capture a wide range of expression levels

- Gray: all transcripts with a read alignment
- Blue: at least 80% of a transcript in a single contig
- Red: at least 80% of a transcript is reconstructed





Trans-ABySS



A versatile tool for

- Transcript reconstruction
- Gene identification
- InDel and SNV discovery
- Chimeric transcript discovery
 - Gene fusions
 - Trans-splicing
- Expression analysis



Transcriptome Assembly



De novo assembly based on ABySS

Reference-based assembly based on TopHat alignments

[Trapnell et al., 2010; Guttman et al., 2010; Trapnell et al., 2009]



 δ_{32}^{00}

+ chimeric transcripts

Detecting Fusions



- Conventionally detected through identifying translocations in genomes
- Assembled transcriptome contigs span multiple genes
- Break points (usually) correspond to exon boundaries
- Break points are supported by
 - Spanning reads
 - Read pairs linking regions



Lucas Swanson, Readman Chiu and Gordon Robertson

Detecting Partial Tandem Duplications



- One or more exons get repeated in their entirety
- Usually coexist with the wild-type
- PTD events are manifested in a particular contig type
 - A short contig with 50/50 split alignment
- Break points are supported by
 - Spanning reads
 - Read pairs in opposite orientation

 δ_{94}^{00}

Lucas Swanson, Readman Chiu and Gordon Robertson

Detecting Internal Tandem Duplications



- Tandem duplications internal to exons
- Contig alignments result in
 - Query gaps
 - Contiguous target blocks
- Read support on break point(s)
- Aberrant read pair distances



Lucas Swanson, Readman Chiu and Gordon Robertson

Performance

- Compared to mapping-based analysis tools Trans-ABySS constructs
 - as many transcripts
 - with better sensitivity and specificity





[Trapnell et al., 2010; Guttman et al., 2010; Trapnell et al., 2009]

Sequence Alignment



W http://en.wikipedia.org/wiki/Sequence_alignment_software

Short-Read Sequence Alignment

Name	Description
BFAST	Explicit time and accuracy tradeoff with a prior accuracy estimation, supported by indexing the reference sequences. Optimally compresses indexes. Can handle billions of short reads. Can handle insertions, deletions, SNPs, and color errors (can map ABI SOLiD color space reads). Performs a full Smith Waterman alignment.
BLASTN	BLAST's nucleotide alignment program, slow and not accurate for short reads, and uses a sequence database (EST, sanger sequence) rather than a reference genome.
BLAT	Made by Jim Kent. Can handle one mismatch in initial alignment step.
Bowtie	Uses a Burrows-Wheeler transform to create a permanent, reusable index of the genome; 1.3 GB memory footprint for human genome. Aligns more than 25 million Illumina reads in 1 CPU hour. Supports Maq-like and SOAP-like alignment policies (can be run from inside Geneious Server).
BWA	Uses a Burrows-Wheeler transform to create an index of the genome. It's a bit slower than bowtie but allows indels in alignment (can be run from inside Geneious Server).
CASHX	Quantify and manage large quantities of short-read sequence data. CASHX pipeline contains a set of tools that can be used together or as independent modules on their own. This algorithm is very accurate for perfect hits to a reference genome.
CUDA-EC	Short-read alignment error correction using GPUs.
drFAST	Read mapping alignment software that implements cache obliviousness to minimize main/cache memory transfers like mrFAST and mrsFAST, however designed for the SOLiD sequencing platform (color space reads). It also returns all possible map locations for improved structural variation discovery.
ELAND	Implemented by Illumina. Includes ungapped alignment with a finite read length.
GNUMAP	Accurately performs gapped alignment of sequence data obtained from next-generation sequencing machines (specifically that of Solexa/Illumina) back to a genome of any size. Includes adaptor trimming, SNP calling and Bisulfite sequence analysis.
GEM	High-quality alignment engine (exhaustive mapping, that is 100% of sensitivity, for any number of substitutions; 1 non-exhaustive indel). Several standalone applications (mapper, split mapper, mappability, and other) provided.
GMAP and GSNAP	Robust, fast, short-read alignment. GMAP: longer reads, with multiple indels and splices (see entry above under Genomics analysis); GSNAP: shorter reads, with a single indel or up to two splices per read. Useful for digital gene expression, SNP and indel genotyping. Developed by Thomas Wu at Genentech. Used by the National Center for Genome Resources (NCGR) in Alpheus.
Geneious Assembler	Fast, accurate overlap assembler with the ability to handle any combination of sequencing technology, read length, any pairing orientations, with any spacer size for the pairing, with or without a reference genome.
LAST	
MAQ	Ungapped alignment that takes into account quality scores for each base (can be run from inside Geneious Server).
mrFAST and	Gapped (mrFAST) and ungapped (mrsFAST) alignment software that implements cache obliviousness to minimize main/cache memory transfers. They are designed for the Illumina sequencing platform and they can return all possible man locations for improved structural variation discovery.

Sequence alignment

- Global
- Local
- Glocal



Global alignment

- Base-by-base alignment of one sequence to another allowing for both mismatches and gaps
- Example:

AGAGTGCTGCCGCC AGATGTACTGCGCC

- Alignment:
 AGA-GTGCTGCCGCC
 ||| || ||| |||
 AGATGTACTGC-GCC
- 12 matches of 15 bp = 80% identity



Local Alignment

- Given two sequences, find a matching substring from each of those two sequences
- Example:

AGATGTGCTGCCGCC TTTGTACTGAAA AGATGTGCTGCCGCC ||| ||| TTTGTACTGAAA

• 6 matches of 7 bp = 86% identity



Glocal Alignment

- Given a query sequence and a reference sequence, identify a substring of the reference sequence that matches the entirety of the query sequence
- Example:

Reference: AGATGTGCTGCCGCCACGT Query: TTTGTACTGAAA ACGTAGATGTGCTGCCGCCACGT ||| ||| TTTGTACTGAAA

• 6 matches of 12 bp = 50% identity



Criteria for Choosing an Aligner

- Global, local or glocal alignment
- Aligning short sequences to long sequences
- Aligning long sequences to long sequences
- Handling small gaps (insertions and deletions)
- Handling large gaps (introns)
- Handling split alignments (chimera)
- Speed and ease of use



Popular Alignment Software

Short reads

- BWA
- GSNAP
- Bowtie
 - TopHat
- SOAP

Long sequence

- BWA-SW
- GMAP
- BLAT
- BLAST
- Exonerate
- MUMmer



Seed and Extend

- For large sequences, an exhaustive alignment is very slow
- Many aligners start by finding perfect or near perfect matches to seeds
- The seeding strategy has a large effect on the sensitivity of the aligner
 - e.g. BLAT requires two perfect nearby 11-mer matches



Memory Use

Hashing

- Load a representation of all the reads and/or the reference into memory
 - GSNAP
 - SOAP
 - mr/mrsFAST
 - KAligner

Burrows-Wheeler Transformation (Ferragina-Manzini, indexing)

- Compress reads and/or the reference before loading
 - BWA
 - Bowtie
 - Abyss-map



Hashing





BW Transform




Inverse BW Transform



TGCACT



Summary

De Novo Assembly Problem



Old paradigm:

long and non-uniform reads (800bp - 1000bp)

De Bruijn Graph

- Load *k*-mers in memory
 - 2x4 possible extension of every k-mer
- Check if there is a "next" k-mer
 - O(n) algorithm
 -GACATTGC... seq1

19



Adjacency Graph

- SET assembly result as a graph
 - Nodes: overlaps (k-1 bp)
 - Edges: contigs





ABySS-Explorer GUI



Summary

Assembly As a Hairball

- ABySS v1.2.7
- PET/MPET information disambiguates short contig extensions Node connectivity* ∖out 4 5 6. in 530 109 9814 1817 456 72 1817 1074 238 31 1882 13 238 126 530 456 72 31 13 10 109 For contigs $\geq 2 \text{ kb}$ 66
 - **Biotin Read-Through**



Scaffold Graph Operations

1. Resolve forks



Anchor

Local directional assembly

 scaffold gap filling

76



Summary

Trans-ABySS



A versatile tool for

- Transcript reconstruction
- Gene identification
- InDel and SNV discovery
- Chimeric transcript
 discovery
 - Gene fusions
 - Trans-splicing
- Expression analysis

Detecting Fusions



Lucas Swanson, Readman Chiu and Gordon Robertson

- Conventionally detected through identifying translocations in genomes
- Assembled transcriptome contigs span multiple genes
- Break points (usually) correspond to exon boundaries
- Break points are supported by
- Spanning reads
 - Read pairs linking regions

BW Transform

Rotate	Sort	Index
TGCACT\$	\$TGCAOT	т
GCACT\$T	ACT\$TGC	С
CACT\$TG	CACT\$TG	G
ACT\$TGC	CT\$TGCA	А
CT\$TGCA	GCACTST	т
T\$T <mark>GCAC</mark>	TSTGCAC	С
\$TGCACT	TGCACT\$	\$

www.bcgsc.ca -> software ABySS, Trans-ABySS, ABySS-Explorer, Anchor ABySS and Trans-ABySS Google Groups ibirol@bcgsc.ca sjackman@bcgsc.ca





ACKNOWLEDGEMENTS

<u>ABySS Team:</u> Shaun Jackman Tony Raymond Rod Docking

Cydney Nielsen

<u>Beetle Project:</u> Joerg Bohlmann Chris Keeling

GenomeCanada

<u>Trans-ABySS Team:</u> Readman Chiu Karen Mungall Gordon Robertson Ka Ming Nip Jenny Qian Rong She Lucas Swanson

Nancy Liao Greg Taylor Simon Chan Diana Palmquist ONAL INSTITUTE

GenomeBritishColumbia

GSC:

Sequencing Team

Library Core

Steven Jones

Marco Marra