Long Read Sequencing Technology

- Algorithms and its applications -

Hayan Lee@Schatz Lab May 4, 2015 Proposal



Outline

- Background
 - Long read sequencing technology
- The limitations of short read mapping illustrated by the Genome Mappability Score (GMS)
- The Resurgence of reference quality genome (3Cs)
 - The next version of Lander-Waterman Statistics (Contiguity)
 - Historical human genome quality by gene block analysis (Completeness)
 - The effectiveness of long reads in de novo assembly (Correctness)
- Sugarcane de novo genome assembly challenges
 - The effectiveness of accurate long reads in de novo assembly especially for highly heterozygous aneuploid genome
 - Pure long read de novo assembly, combine with accurate long reads and erroneous long reads
- Contributions

Background

- BAC-by-BAC + Sanger Era (1995 to 2007)
 - Very high quality reference genomes for human, mouse, worm, fly, rice, Arabidopsis and a select few other high value species.
 - Contig sizes in the megabases, but costs in the 10s to 100s of millions of dollars

• Next-Gen Era (2007 to current)

- Costs dropped, but genome quality suffered
- Genome finishing almost completely abandoned; "exon-sized" contigs
- These low quality draft sequences are (1) missing important sequences, (2) lack context to discover regulatory elements or evolutionary patterns, and (3) contain many errors

• Third-Gen Era (current)

- New biotechnologies (single molecule, chromatin assays, etc) and new algorithms (MHAP, LACHESIS, etc) are leading to a *Resurgence of Reference Quality Genomes*
- De novo assemblies of human and other large genomes with contig sizes over 1Mbp.

CSH Cold Spring Harbor Laboratory

Third-Gen Sequencing Technology

Long Read Sequencing: De novo assembly, SV analysis, phasing



Long Span Sequencing: Chromosome Scaffolding, SV analysis, phasing



Outline

- Background
 - Long read sequencing technology and algorithms
- The limitations of short read mapping illustrated by the Genome Mappability Score (GMS)
- The Resurgence of reference quality genome (3Cs)
 - The next version of Lander-Waterman Statistics (Contiguity)
 - Historical human genome quality by gene block analysis (Completeness)
 - The effectiveness of long reads in de novo assembly (Correctness)
- Sugarcane de novo genome assembly challenge
 - The effectiveness of accurate long reads in de novo assembly especially for highly heterozygous aneuploidy genome
 - Pure long read de novo assembly, combine with accurate long reads and erroneous long reads
- Contributions

Short read mapping (Resequencing)



- Discovering genome variations
- Investigating the relationship between variations and phenotypes
- Profiling epigenetic activations and inactivations
- Measuring transcription rates

old Spring Harbor Laboratory







Read Quality Score – MAQ

Sensitivity of Read Mapping Score



Challenges

- There is inherent uncertainty to mapping
- Read quality score is very sensitive to a minute change
- Base quality score is useful only inside a single read
- Read quality score is assigned to each read not a position of a genome, thus provides only local view
- However, there is no tool to measure the reliability of mapped reads to the reference genome in a global perspective.

It does not consider all possible reads We need more stable "GPS" for a genome



Genome Mappability Score (GMS)



- x is a reference
- z is a read
- I is read length



Genome Mappability Analyzer (GMA)



GMS vs. MAQ Sensitivity of Read Mapping Score

Comparison GMS vs MAQ (Read length: 100bp, error rate: 1%, Paired-end)











Variation Accuracy Simulator (VAS)



 Simulation of resequencing experiments to measure the accuracy of variation detection



Genomic Dark Matter



 Unlike false negatives in high GMS region that can be discovered in high coverage (>=20-fold), false negatives in low GMS regions cannot be discovered, because variation calling program will not use poorly mapped reads

CSH Cold Spring Harbor Laboratory

BIOMFORMATICS ORIGINAL PAPER Vol. 28 no. 16 2012, pages 2097-2106 doi:10.1093/bioinformatics/bts330 Monte analysis Advance Access publication June 4, 2012 Genomic dark matter: the reliability of short read mapping ilustrated by the genome mappability score Advance Access publication June 4, 2012 Hayan Lee^{1,4} and Michael C. Schatz^{1,2} Advance Access publications Center for Quantitive Biology, Cold Spring Harbor Laboratory, Cold Spring Harbor, NY, USA and ²Simons Center for Quantitive Biology, Cold Spring Harbor Laboratory, Cold Spring Harbor, NY, USA and ²Simons Center for Quantitive Biology, Cold Spring Harbor Laboratory, Cold Spring Harbor, NY, USA and Paines Center for Quantitive Biology, Cold Spring Harbor Laboratory, Cold Spring Harbor, NY, USA and Paines Center for Quantitive Biology, Cold Spring Harbor Laboratory, Cold Spring Harbor, NY, USA and Paines Center for Quantitive Biology, Cold Spring Harbor Laboratory, Cold Spring Harbor, NY, USA and Paines Center for Quantitive Biology, Cold Spring Harbor Laboratory, Cold Spring Harbor, NY, USA and Paines Center for Quantitive Biology, Cold Spring Harbor Laboratory, Cold Spring Harbor, NY, USA

ABSTRACT

Motivation: Genome resequencing and short read mapping are two of the primary tools of genomics and are used for many important applications. The current state-of-the-art in mapping uses the quality values and mapping quality scores to evaluate the reliability of the mapping. These attributes, however, are assigned to individual reads and do not directly measure the problematic repeats across the genome. Here, we present the Genome Mappability Score (GMS) as a novel measure of the complexity of resequencing a genome. The GMS is a weighted probability that any read could be unambiguously mapped to a given position and thus measures the overall composition of the genome itself.

Results: We have developed the Genome Mappability Analyzer to compute the GMS of every position in a genome. It leverages the parallelism of cloud computing to analyze large genomes, and enabled us to identify the 5–14% of the human, mouse, fly and yeast genomes that are difficult to analyze with short reads. We examined the accuracy of the widely used BWA/SAMtools polymorphism discovery pipeline in the context of the GMS, and found discovery errors are dominated by false negatives, especially in regions with poor GMS. These errors are fundamental to the mapping process and cannot be overcome by increasing coverage. As such, the GMS should be considered in every resequencing project to pinpoint the 'dark matter' of the genome, including of known clinically relevant uvariations in these regions. sequencing, including several large projects to sequence thousands of human genomes and exomes, such as the (1000 Genomes Project Consortium, 2010) or (International Cancer Genome Consortium, 2010). Other projects, such as (ENCODE Project Consortium, 2004) and (modENCODE Consortium, 2010), are extensively using resequencing and read mapping to discover novel genes and binding sites.

The output of current DNA sequencing instruments consists of billions of short, 25–200 bp sequences of DNA called reads, with an overall per base error rate around 1–2% (Bentley *et al.*, 2008). In the case of whole genome resequencing, these short reads will originate from random locations in the genome, but nevertheless, entire genomes can be accurately studied by oversampling the genome, and then aligning or 'mapping' each read to the reference genome to computationally identify where it originated. Once the entire collection of reads has been mapped, variations in the sample can be identified by the pileup of reads that significantly disagree from the reference genome (Fig. 1).

The leading short read mapping algorithms, including BWA (Li and Durbin, 2009), Bowtie (Langmead *et al.*, 2009), and SOAP (Li *et al.*, 2009b), all try to identify the best mapping position for each read that minimizes the number of differences between the read and the genome, i.e. the edit distance of the nucleotide strings, possibly weighted by base quality value. This is made practical through sophisticated indexing schemes such as the Burrows-Wheeler



Cited by

Kim et al. Genome Biology 2013, 14:R90 http://genomebiology.com/2013/14/8/R90



METHOD

Open Access

Virmid: accurate detection of somatic mutations with sample impurity inference

Sanguese Kim^{1*†}, Kyowon Jeong^{2†}, Kunal Bhutani¹, Jeong Ho Lee^{3,6}, Anand Patel¹, Eric Scott³, Hojung Nam⁴, Havan Lee⁵, Joseph G Gleeson³ and Vineet Bafna^{1*}

mut

Abstract

Detection of somatic variation using sequence from disease-con many cases including cancer, however, it is hard to isolate pure mutation analysis by disrupting overall allele frequencies. Here, a determines the level of impurity in the sample, and uses it for it tests on simulated and real sequencing data from breast cancer of our model. A software implementation of our method is avai

Background

Identifying mutations relevant to a specific phenotype is relat one of the primary goals in sequence analysis. With the prol advent of massively parallel sequencing technologies, we exo can produce an immense amount of genomic informapote tion to estimate the landscape of sequence variations. schi However, the error rates for base-call and read alignment still remain much higher than the empirical frecove quencies of single nucleotide variations (SNVs) and de imp novo mutations [1]. Many statistical methods have been exar proposed to strengthen mutation discovery in the pretain sence of confounding errors [2-4]. acq

Finding somatic mutations is one particular type of variant calling, which constitutes an essential step of clinical genotyping. Unlike the procedures used for germ line mutation discovery, the availability of a matched control complexity in the second s

LETTER

doi:10.1038/nature13907

Resolving the complexity of the human genome using single-molecule sequencing

Mark J. P. Chaisson¹, John Huddleston^{1,2}, Megan Y. Dennis¹, Peter H. Sudmant¹, Maika Malig¹, Fereydoun Hormozdiari¹, Francesca Antonacci³, Urvashi Surti⁴, Richard Sandstrom¹, Matthew Boitano⁵, Jane M. Landolin⁵, John A. Stamatoyannopoulos¹, Michael W. Hunkapiller⁵, Jonas Korlach⁵ & Evan E. Eichler^{1,2}

The human genome is arguably the most complete mammalian reference assembly¹⁻³, yet more than 160 euchromatic gaps remain⁴⁻⁶ and aspects of its structural variation remain poorly understood ten years after its completion⁷⁻⁹. To identify missing sequence and genetic variation, here we sequence and analyse a haploid human genome (CHM1) using single-molecule, real-time DNA sequencing¹⁰. We dose or extend 55% of the remaining interstitial gaps in the human GRCh37 reference genome—78% of which carried long runs of degenerate short tandem repeats, often several kilobases in length, embedded within (G+C)-rich genomic regions. We resolve the complete sequence of 26,079 euchromatic structural variants at the base-pair level, including inversions, complex insertions and long tracts of tandem repeats. Most have not been previously reported, with the greatest increases in sensitivity occurring for events less than 5 kilobases in size. Com-

for recruiting additional sequence reads for assembly (Supplementary Information). Using this approach, we closed 50 gaps and extended into 40 others (60 boundaries), adding 398 kb and 721 kb of novel sequence to the genome, respectively (Supplementary Table 4). The closed gaps in the human genome were enriched for simple repeats, long tandem repeats, and high (G+C) content (Fig. 1) but also included novel exons (Supplementary Table 20) and putative regulatory sequences based on DNase I hypersensitivity and chromatin immunoprecipitation followed by high-throughput DNA sequencing (ChIP-seq) analysis (Supplementary Information). We identified a significant 15-fold enrichment of short tandem repeats (STRs) when compared to a random sample (P < 0.00001) (Fig. 1a). A total of 78% (39 out of 50) of the closed gap sequences were composed of 10% or more of STRs. The STRs were frequently embedded

Outline

- Background
 - Long read sequencing technology and algorithms
- The limitations of short read mapping illustrated by the Genome Mappability Score (GMS)
- The Resurgence of reference quality genome (3Cs)
 - The next version of Lander-Waterman Statistics (Contiguity)
 - Historical human genome quality by gene block analysis (Completeness)
 - The effectiveness of long reads in de novo assembly (Correctness)
- Sugarcane de novo genome assembly challenge
 - The effectiveness of accurate long reads in de novo assembly especially for highly heterozygous aneuploidy genome
 - Pure long read de novo assembly, combine with accurate long reads and erroneous long reads
- Contributions



2. Construct assembly graph from overlapping reads

...AGCCTAG<mark>GGATGCGCGACACG</mark>T

CAACCTCGGACGGACCTCAGCGAA...

3. Simplify assembly graph

4. Detangle graph with long reads, mates, and other links



Many Genomes Are Sequenced... Many Questions Are Raised... But...

- How long should the read length be?
- What coverage should be used?

Given the read length and coverage,

- How long are contigs? <- Contiguity prediction
- How many contigs?
- How many reads are in each contigs?
- How big are the gaps?

Cold Spring Harbor Laboratory

Lander-Waterman Statistics

GENOMICS 2, 231-239 (1988)

Genomic Mapping by Fingerprinting Random Clones: A Mathematical Analysis

ERIC S. LANDER + AND MICHAEL S. WATERMAN

*Whitehead Institute for Biomedical Research, 9 Cambridge Center, Cambridge, Massachusetts 02142; †Harvard University, Cambridge, Massachusetts 02138; and ‡Departments of Mathematics and Molecular Biology, University of Southern California, Los Angeles, California 90089

Received January 13, 1988; revised March 31, 1988

Results from physical mapping projects have recently been reported for the genomes of *Escherichia coli*, *Saccharomyces cerevisiae*, and *Caenorhabditis elegans*, and similar projects are currently being planned for other organisms. In such projects, the physical map is assembled by first "fingerprinting" a large number of clones chosen at random from a recombinant library and then inferring overlaps between clones with sufficiently similar fingerprints. available region of up to several megabases and of studying its properties. In addition, the overlapping clones comprising the physical map would constitute the logical substrate for efforts to sequence an organism's genome.

Recently, three pioneering efforts have investigated the feasibility of assembling physical maps by means of "fingerprinting" randomly chosen clones. The fingerprints consisted of information about restriction

Lander-Waterman Statistics



In practice, it's useful only in low coverage (3-5x) but becomes nonsensical in high coverage.



HG19 Genome Assembly Performance by Lander-Waterman Statistics



Empirical Data-driven Approach

- We selected 26 species across tree of life and exhaustively analyzed their assemblies using simulated reads for 4 different length (6 for HG19) and 4 different coverage per species
- For the extra long reads, we fixed the Celera Assembler(CA) to support reads up to 0.5Mbp



26 Species Across Tree of Life

	Model	ID	Genome Size
	Organism		
	M.jannaschii	1	1,664,970
	C.hydrogenoformans	2	2,401,520
	E.coli	3	4,639,675
	Y.pestis	4	4,653,728
	B.anthracis	5	5,227,293
	A.minum	6	8,248,144
	yeast	7	12,157,105
	Y.lipolytica	8	20,502,981
	slime mold	9	34,338,145
	Red bread mold	10	41,037,538
	sea squirt	- 11	78,296,155
	roundworm	12	100,272,276
	green alga	13	112,305,447
	arahidonsis	14	119,667,750
	fruitfly	15	130,450,100
	peach	16	227,252,106
	rice	17	370,792,118
	poplar	18	417,640,243
	tomato	19	781,666,411
	soybean	20	973,344,380
	turkey	21	1,061,998,909
	zebra fish	22	1,412,464,843
	lizard	23	1,799,126,364
	com	24	2,066,432,718
	mouse	25	2,654,895,218
-	human	26	3,095,693,983









HG19 Genome Assembly Performance



Why?

Lander-Waterman Statistics

- Assumptions!!!
- If genome is a random sequence, it will work
- It works only in low coverage 3-5x
- It works for small genomes (< yeast)

Our Approach

- Stop assuming that we cannot guarantee!!!
- We tried to assume as least as possible.
- Instead of building on top of assumptions, we let the model learn from the data
- Empirical data-driven approach



& CSH

Repeats in Rice



Our Goal





Assembly Challenge (1) Read Length

- Read length is very important
- A matter of technology
- The longer is the better
- Quality was important but can be corrected
 - PacBio produces long reads, but low quality (~15% error rate)
 - Error correction pipeline are developed
 - Errors are corrected very accurately up to 99%

CSH Cold Spring Harbor Laboratory

- Assembly Challenge (1) - Read Length



CSH Cold Spring Harbor Laboratory

Assembly Challenge (2) Coverage

- A matter of money
- Using perfect reads, assembly performance increased for most genomes : Lower bound
- Using real reads, overall performance line will shift to the higher coverage
- The higher is the better (?)
- But still it suggests that there would be a threshold that can maximize your return on investment (ROI)

Assembly Challenge (2)



Assembly Challenge (3) **Repeats**

- Genome is not a random sequence
- Repeat hurts genome assembly performance
- Isolating the impact of repeats is not trivial
- Quantifying repeat characteristics is not trivial as well





Stony Brook University

Dept. of Computer Science

Longest Repeat Size and Genome Size


Assembly Challenge (4) Genome Size

- Increase the assembly complexity
- Make a hard problem harder.



Assembly Challenge (4) Genome Size



Assembly Challenge (4)

Genome Size



Challenges for Prediction

- Sample size is small
- Quality is not guaranteed
- Predictive Power
- Overfitting

Support Vector Regression (SVR) Cross Validation



Stony Brook University



The resurgence of reference genome qaultiy

Lee, H, Gurtowski, J, Yoo, S, Marcus, S, McCombie, WR, Schatz MC et al. (2015) In preparation

Predictive Power

- Average of residual is 15%
- We can predict the new genome assembly performance in 15% of error residual boundary
- Genome size, read length and coverage used explicitly
- Repeats are included implicitly



Web Service for Contiguity Prediction

	See Genome Assembly Performance Prediction - Mozilla Firefox
	Genome Assembly Perf ×
	Genome Assembly Performance Prediction
	Http://qb.cshl.edu/asm_model/predict.html
	Given genome size, we internally set read lengths and coverages for you. With 3 features, our model predicts the expected performance of assembly. Performance is defined as follows:
7	Performance(%) = N50 of assembly / N50 of chromosome segments Genome size : 1000123456 Submit
	Assembly Prediction of Genome Size 1000123456 By Coverage
	88
	20 0 0 0 0 0 0 0 0 0 0 0 0 0

Reference Genome Quality



Contiguity



de novo human genome assembly

What happens as we sequence the human genome with longer reads?

- Red: Sizes of the chromosome arms of HG19 from largest to shortest
- Green: Results of our assemblies using progressively longer and longer reads
- Orange: Results of Allpaths/Illumina assemblies

Lengths selected to represent the biotechnologies:

- mean1: ~Moleculo
- mean2: ~PacBio/ONT
- mean16: ~10x / Chromatin
- mean32: ~Optical mapping (log-normal with increasing means)

Validated by MHAP

Add results

CSH Spring hinRviv	HOME ABOUT SUBMIT ALERTS / RSS
Laboratory	Search Q
THE PREPRINT SERVER FOR BIOLOGY	Advanced Search
New Perula	Previous
New Results	C Previous
Assembling Large Genomes with Single-Molecule Sequencing	and Locality Posted August 14, 2014.
Sensitive Hashing	Download PDF Share
Konstantin Berlin , Sergey Koren , Chen-Shan Chin , James Drake , Jane M Landolin ,	Adam M Phillippy Semail Citation Tools
doi: http://dx.doi.org/10.1101/008003	
Abstract Info/History Metrics Data Supplements	Preview PDF Tweet 167 F Like 13 2
	Subject Area
Abstract	Bioinformatics
We report reference-grade de novo assemblies of four model organ	sms and the
human genome from single-molecule, real-time (SMRT) sequencing.	Long-read SMRT Subject Areas
sequencing is routinely used to finish microbial genomes, but the a	vailable assembly
methods have not scaled well to larger genomes. Here we introduce	the MinHash All Articles
Alignment Process (MHAP) for efficient overlapping of noisy, long re	ads using Animal Behavior and Cognition
probabilistic, locality-sensitive hashing. Together with Celera Assem	bler, MHAP was Biochemistry
used to reconstruct the genomes of Escherichia coli, Saccharomyce	cerevisiae, Bioengineering
Arabidopsis thaliana, Drosophila melanogaster, and human from hi	h-coverage SMRT Bioinformatics
sequencing. The resulting assemblies include fully resolved chromo	some arms and Biophysics

١

Developmental Biology

Ecology

Preprint

CSH Spring Harbor Laboratory biology	HOME AE	BOUT SUBMIT ALERTS / RSS	
New Results		Previous	
Error correction and assembly complexity of single molecule s	equencing reads.	Posted June 18, 2014.	
Hayan Lee , James Gurtowski , Shinjae Yoo , Shoshana Marcus , W. Richard McCombi loi: http://dx.doi.org/10.1101/006395	ie , Michael Schatz	Download PDF Email	
Abstract Info/History Metrics Data Supplements	Preview PDF	Tweet 61	
Abstract		Subject Area	
Third generation single molecule sequencing technology is poised to	o revolutionize	Bioinformatics	
genomics by enabling the sequencing of long, individual molecules	of DNA and RNA.		
These technologies now routinely produce reads exceeding 5,000 b	asepairs, and can	Subject Areas	
achieve reads as long as 50,000 basepairs. Here we evaluate the lim	its of single	011.0	
molecule sequencing by assessing the impact of long read sequenci	ng in the assembly	All Articles	
of the human genome and 25 other important genomes across the t	tree of life. From	Animal Behavior and C	
this, we develop a new data-driven model using support vector regre	ession that can	Biochemistry	
accurately predict assembly performance. We also present a novel h	ybrid error	Bioengineering	
correction algorithm for long PacBio sequencing reads that uses pre	-assembled	Bioinformatics	
Illumina sequences for the error correction. We apply it several prok	aryotic and	Biophysics	
eukaryotic genomes, and show it can achieve near-perfect assemblie	es of small	Cancer Biology	
genomes (< 100Mbp) and substantially improved assemblies of large	er ones. All source	Cell Biology	
code and the assembly model are available open-source.		Developmental Biolog	

Completeness

Human Reference Genome Quality by gene block analysis



Completeness

Human Reference Genome Quality by gene block analysis



Completeness

Human Reference Genome Quality by gene block analysis



Larger contigs and scaffolds empowers analysis at every possible level.

- SNPs (~10k clinically relevant)
- Genes
- Regulatory elements
- Synteny blocks
- Chromosome structure

Correctness Summary in HG19

-N50 misleading

HG19	(major) misassembly	major) breaks) (major)
	False Positive	False Negative
	Increase N50 (falsely lengthen contiguity)	Decrease N50 (shorten contiguity)
	Mislead us in biological meaning	Negatively impact on downstream research
Mean1	209	4069
Mean2	70	462
Mean4	49	296
Mean8	33	197
Mean16	9	42
Mean32	7	5

Misassembly



HG19.m4.c20.misassemble

HG19.m8.c20.misassemble

Misassembly Analysis in HG19



Misassembly Analysis in HG19



Long read sequencing technology helps to reduce both misassembly and breaks thus increase correctness of de novo genome assembly

Summary & Recommendations

Reference quality genome assembly is here

- Use the longest possible reads and spans for the best assembly
- Coverage and algorithmics overcome most random errors

Megabase N50 improves the analysis in every dimension

- Better resolution of genes and flanking regulatory regions
- Better resolution of transposons and other complex sequences
- Better resolution of chromosome organization

Need to develop methods to jointly analyze multiple highquality references at once



Related Work

BIOINFORMATICS ORIGINAL PAPER

Genome analysis

Advance Access publication November 13, 2014

Vol. 30 no. 24 2014, pages 3476-3483

doi:10.1093/bioinformatics/btu756

SplitMEM: a graphical algorithm for pan-genome analysis with suffix skips

Shoshana Marcus¹, Hayan Lee^{1,2} and Michael C. Schatz^{1,2,*} ¹Simons Center for Quantitative Biology, Cold Spring Harbor Laboratory, Cold Spring Harbor, NY 11724, USA and ²Department of Computer Science, Stony Brook University, Stony Brook, NY 11794, USA

Associate Editor: Gunnar Ratsch

ABSTRACT

Motivation: Genomics is expanding from a single reference per species paradigm into a more comprehensive pan-genome approach that analyzes multiple individuals together. A compressed de Bruijn graph is a sophisticated data structure for representing the genomes of entire populations. It robustly encodes shared segments, simple single-nucleotide polymorphisms and complex structural variations far beyond what can be represented in a collection of linear sequences alone.

Results: We explore deep topological relationships between suffix trees and compressed de Bruijn graphs and introduce an algorithm, splitMEM, that directly constructs the compressed de Bruijn graph in time and space linear to the total number of genomes for a given maximum genome size. We introduce *suffix skips* to traverse several suffix links simultaneously and use them to efficiently decompose maximal exact matches into graph nodes. We demonstrate the utility of splitMEM by analyzing the nine-strain pan-genome of *Bacillus anthracis* and up to 62 strains of *Escherichia coli*, revealing their core-genome properties.

Availability and implementation: Source code and documentation available open-source http://splitmem.sourceforge.net. Contact: mschatz@cshl.edu

Supplementary information: Supplementary data are available at ⁵⁶ Bioinformatics online.

resequencing projects, gene discovery and numerous other important applications. However, reference genomes also suffer in that they represent a single individual or a mosaic of individuals as a single linear sequence, making them an incomplete catalog of all the known genes, variants and other variable elements in a population. Especially in the case of structural and other largescale variations, this creates an analysis gap when modeling the role of complex variations or gene flow in the population. For the human genome, for example, multiple auxiliary databases including dbSNP, dbVAR, DGV and several others must be separately queried through several different interfaces to access the population-wide status of a variant (MacDonald *et al.*, 2014).

The 'reference-centric' approach in genomics has been established largely because of technological and budgetary concerns. Especially in the case of mammalian-sized genomes, it remains prohibitively expensive and technically challenging to assemble each sample into a complete genome *de novo*, making it substantially cheaper and more accessible to analyze a new sample relative to an established reference. However, for some species, especially medically or otherwise biologically important microbial genomes, multiple genomes of the same species are available. In the current version of National Center for Biotechnology Information (NCBI) GenBank, 296 of the 1471 bacterial species listed have at least two strains present, including 9 strains of

Pan-Genome Alignment & Assembly



Time to start considering problems for which N complete genomes is the input to study the "pan-genome" •Available today for many microbial species, near future for higher eukaryotes



Pan-genome colored de Bruijn graph

- Encodes all the sequence relationships between the genomes
- How well conserved is a given sequence?
- What are the pan-genome network properties?

SplitMEM: A graphical algorithm for pan-genome analysis with suffix skips Marcus, S, Lee, H, Schatz MC (2014) *Bioinformatics.* doi: 10.1093/bioinformatics/btu756

Outline

- Background
 - Long read sequencing technology and algorithms
- The limitations of short read mapping illustrated by the Genome Mappability Score (GMS)
- The Resurgence of reference genome quality (3Cs)
 - The next version of Lander-Waterman Statistics (Contiguity)
 - Historical human genome quality by gene block analysis (Completeness)
 - The effectiveness of long reads in de novo assembly (Correctness)
- Sugarcane de novo genome assembly challenge
 - The effectiveness of accurate long reads in de novo assembly especially for highly heterozygous aneuploidy genome
 - Pure long read de novo assembly, combine with Moleculo and PacBio reads.
- Contributions

Sugarcane for food and biofuel

• Food

- By 2050, the world's population will grow by 50%, thus another
 2.5 billion people will need to eat!
- Rapidly rising oil prices, adverse weather conditions, speculation in agricultural markets are causing more demand

• Biofuel

- By 2050, global energy needs will double as will carbon dioxide emission
- Low-carbon solution
- Sugarcane ethanol is a clean, renewable fuel that produces on average 90 percent less carbon dioxide emission than oil and can be an important tool in the fight against climate change.



S. officinarum

(Contribute to sweetness)

Sugarcane

A hybrid sugarcane cultivar SP80-3280

S. spontaneum

(Contribute to robustness)

- S.spontaneum x S.officinarum
- A century ago....
- Saccharum genus
 - S. spontaneum (2n=40-128, x=8)
 - S. officinarum (2n=8x=80)



- Monoploid genome is about 1Gbp
- 8-12 copies per chromosome
- In total, 100-130 chromosomes
- Total size is about 10Gbp

H Cold Spring Harbor Laboratory

Why is sugarcane assembly harder? (1)

Polyploidy/Aneuploidy

 10% of the chromosomes are inherited in their entirety from *S. spontaneum*, 80% are inherited entirely from *S. officinarum*

Large scale recombination

 10% is the result of recombination between chromosomes from the two ancestral species, a few being double recombinants



(source) http://ars.elscdn.com/content/image/1-s2.0-S1369526602002340-gr1.jpg Simons Center for Quantitative Biology

Four Important Questions in Sugarcane

Scaffold polyploidy/aneuploidy genome

 How do we connect contigs/cluster contigs per chromosome/fill gaps among contigs?

^J Phasing haplotypes

Not solved in diploid genome yet

赵 Heterozygosity

- How do we measure heterozygosity in polyploidy/aneuploidy genome?
- How do we quantify alleles and get ratio?
- Inference of polyploidy/aneuploidy estimation
 - How do we infer the number of copies per chromosome in aneuploidy genome, especially in the large scale of recombination?

Margarido GRA, Heckerman D (2015) ConPADE: Genome Assembly Ploidy Estimation from Next-Generation Sequencing Data. PLoS Comput Biol 11(4): e1004229. doi: 10.1371/journal.pcbi.1004229

Assembly Complexity by Repeats





Long Reads is the solution!!!



Assembly Complexity by Heterozygosity







Assembly Complexity by Polyploidy





CSH Cold Spring Harbor Laboratory

Dept. of Computer Science

Moleculo Reads

- (1) The DNA is sheared into fragments of about 10Kbp
- (2) Sheared fragments are then diluted
- (3) and placed into 384 wells, at about 3,000 fragments per well.
- (4) Within each well, fragments are amplified through long-range PCR, cut into short fragments and barcoded
- (5) before finally being pooled together and sequenced.
- (6) Sequenced short reads are aligned and mapped back to their original well using the barcode adapters.
- (7) Within each well, reads are grouped into fragments, which are assembled to long reads.



Read length distribution in Moleculo



CSH Cold Spring Harbor Laboratory

Choose the right data and the right method

DATA	 Hiseq 2000 PE (2x100bp) 575Gbp 600x of haploid genome Roche454 9x of haploid genome [min=20 max=1,168] Mean=332bp 	Moleculo 19Gbp 19x of haploid genome [min=1,500 max=22,904] Mean = 4,930bp
Algorithm	SOAPdenovo (De Bruijn Graph)	Celera Assembler (Overlap Graph)
RESULT	Max contig = 21,564 bp NG50= 823 bp Coverage= 0.86x	Max contig = 467,567 bp NG50= 41,394 bp Coverage= 3.59x # of contigs = 450K

CSH Cold Spring Harbor Laboratory

CEGMA

• CEGs - Korf Lab in UC. Davis selected 248 core eukaryotic genes

Statistics of the completeness

	Prots	%Completeness	Total	Average	%Ortho
Complete	219	88.31	827	3.78	89.04
Partial	242	97.58	1083	4.48	95.45

- Gene prediction aided by sorghum gene model
 - In progess...
 - 39k sorghum genes were found in sugarcane contigs at least partially



NP-Hard Hairball of Sugarcane

Vertices are contigs Edges are linking information Edges are reliable linking information from 120 Gbp 10K jumping library

of vertices : 81,552

of edges : 82,269

Average degree of a node : 1

of connected components = 17,919 Average number of vertices per CC= 2.54 The biggest CC has 25 vertices

Benefits of Long Read Scaffolding



- Read Length is increasing, the cost is decreasing
- Very informative whether it has high error rate or not
- More repeats resolved
- Better scaffolding solution than long jumping library
- We don't have to approximate insert size by MLE or so.
- It's much better to fill gaps with some base information rather than just NNNNNN.

Dept. of Computer Science

Prototype for scaffolding



Simulate heterozygous 1. polyploidy genome

> - 10 copies with 5% of difference from original chromosome

2. Simulate Moleculo reads from polyploidy genome

- Read length distribution follows exactly real molecule read distribution

3. Simulate PacBio reads from polyploidy genome

> - Simulate P6-C4, the lastest PacBio chemistry

- 4. Run Celera Assembler(CA) to assemble contigs with Moleculo reads
- Run LRScf to scaffold the 5. contigs with PacBio reads
Preliminary Results

 \checkmark

×

- Moleculo-based contigs from CA
 - Around 700 contigs
- Long Read Scaffolding
 - Align PacBio reads to all contigs
 - Find PacBio reads that link between two contigs
 - Around 1600 alignments out of 40K PacBio Reads

Sugarcane Scaffolding Challenges

- How to represent aneuploidy genome?
- How to screen out false positive link information?
 - # Weakly connected components
 - # Strongly connected components 61
 - True value 5 < 10 < 61
- How to assemble PacBio reads across gaps?

• How to extend contigs with PacBio reads?



Contributions

- The limitations of short read mapping illustrated by the Genome Mappability Score (GMS)
 - A new metric that measure reliability per position of a genome
 - Cloud computing pipeline for efficient computation for big genomes
 - Analysis of biological importance in variation discover low/high GMS region
- The Resurgence of reference genome quality (3Cs)
 - Provide the data-driven model, a.k.a. the next version of Lander-Waterman
 Statistics to predict contiguity of de novo genome assembly project
 - Analysis of completeness and correctness in historical human genome assembly
- Sugarcane de novo genome assembly challenge
 - Showed the effectiveness of accurate long reads in de novo assembly especially for highly heterozygous aneuploidy genome
 - <u>NG50</u> contig length improved 50 times
 - The longest contig extended 25 times to half million bp
 - Pure long read de novo assembly for both contigs and scaffolding



Committee



The State University of New York

Steven Skiena Rob Patro Michael Schatz



Acknowledgements



Schatz Lab Michael Schatz Fritz Sedlazeck James Gurtowski Sri Ramakrishnan Han fang Maria Nattestad Rob Aboukhalil Tyler Garvin Mohammad Amin Shoshana Marcus

McCombie Lab Dick McCombie Sara Goodwin



Shinjae Yoo

Microsoft[®] Research

Ravi Pandya Bob Davidson David Heckerman



University of São Paulo Gabriel Rodrigues Alves Margarido Jonas W. Gaiarsa

Carolina G. Lembke Marie-Anne Van Sluys Glaucia M. Souza



The State University of New York



Thank You Q & A

