

Detection and analysis of structural variants in personal genomes

> Mark B Gerstein Yale

> > Slides at

Lectures.GersteinLab.org (See Last Slide for References & More Info.)

## **Plummeting Cost of Sequencing...**



**2** - Lectures.GersteinLab.org (e) '09

## <u>... has led to the era of</u> <u>Personal Genomics</u>

- Resequencing of individuals' genomes
- Now for a few, eventually for many

Easy (!) except for the genome's complex,
 <u>**REPEAT</u>-CON<u><b>TA**</u>INING STRUCTURE</u>
</u>





### 4. Local Reassembly

Breakpointer: Segmentation of Array Signal as precursor to Read Depth





7 - Lectures.GersteinLab.org (c) 00



- To get highest resolution on breakpoints need to smooth & segment the signal
- BreakPtr: prediction of breakpoints, dosage and crosshybridization using a system based on Hidden Markov Models



Korbel\*, Urban\* et al., PNAS (2007)

http://breakptr.gersteinlab.org

*BreakPtr* statistically integrates array signal and DNA sequence signatures (using a discrete-valued bivariate HMM)



Korbel\*, Urban\* et al., PNAS (2007)

### <u>'Active' approach for breakpoint identification: initial scoring</u> with preliminary model, targeted validation (with sequencing), retraining, and rescoring



CNV breakpoints sequenced in ~10 cases following BreakPtr analysis;

#### Median resolution <300 bp

No improvement in accuracy with higher resolution (9nt tiling)

HMM optimized iteratively (using Expectation Maximization, EM) Korbel\*, Urban\* *et al.*, PNAS (2007)

## MSB: Read-Depth Segmentation



### <u>Mean-shift-based</u> (MSB) Segmentation: no explicit model

- For each bin attraction (meanshift) vector points in the direction of bins with most similar RD signal
- No prior assumptions about number, sizes, haplotype, frequency and density of CNV regions
- Not Model-based (e.g. like HMM) with global optimization, distr. assumption & parms. (e.g. num. of segments).
- Achieves discontinuity-preserving smoothing
- Derived from image-processing applications



[Wang et al. Gen. Res ('09) 19:106]

RD signal

## Some Intuition on how MSB works: Non-Parametric Density Estimation

Assumption : The data points are sampled from an underlying PDF MSB determines grad. of this function



[ Adapted from S Ullman et al. "Advanced Topics in Computer Vision," www.wisdom.weizmann.ac.il/~vision/courses/2004\_2 ]



[ Adapted from S Ullman et al. "Advanced Topics in Computer Vision," www.wisdom.weizmann.ac.il/~vision/courses/2004\_2 ]



[ Adapted from S Ullman et al. "Advanced Topics in Computer Vision," www.wisdom.weizmann.ac.il/~vision/courses/2004\_2 ]

**15** - Lectures.GersteinLab.org<sub>00</sub>



[ Adapted from S Ullman et al. "Advanced Topics in Computer Vision," www.wisdom.weizmann.ac.il/~vision/courses/2004\_2 ]

60, (c) - Lectures.GersteinLab.org 6 H



[ Adapted from S Ullman et al. "Advanced Topics in Computer Vision," www.wisdom.weizmann.ac.il/~vision/courses/2004\_2 ]

**17** - Lectures.GersteinLab.org

60, (c)



[ Adapted from S Ullman et al. "Advanced Topics in Computer Vision," www.wisdom.weizmann.ac.il/~vision/courses/2004\_2 ]

60, (c) - Lectures.GersteinLab.org 8 H



[ Adapted from S Ullman et al. "Advanced Topics in Computer Vision," www.wisdom.weizmann.ac.il/~vision/courses/2004\_2 ]



[ Adapted from S Ullman et al. "Advanced Topics in Computer Vision," www.wisdom.weizmann.ac.il/~vision/courses/2004\_2 ]

**20** - Lectures.GersteinLab.org



**21** - Lectures.GersteinLab.org  $_{
m 0,\infty}$ 



The blue data points were traversed by the windows towards the mode

[ Adapted from S Ullman et al. "Advanced Topics in Computer Vision," www.wisdom.weizmann.ac.il/~vision/courses/2004\_2 ]

### **GC** bias correction



 $RD_{corrected} = \overline{RD}_{global} RD/\overline{RD}_{GC}$ 

## Example of Application of MSB to RD data



### MSB works well on array data too





wavelet

CLAC

GLAD

[Wang et al. Gen. Res ('09) 19:106]



Looking for Aberrantly Placed Paired Ends





PEMer: Detecting Structural Variants from Discordant Paired Ends in Massive Sequencing

> [Korbel et al., Science ('07); Korbel et al., GenomeBiol. ('09)]



**28** - Lectures.GersteinLab.org (a) w

Parameterize Error Models <u>through</u> Simulation

Reconstruction efficiency at different coverage **Deletion** size Reconstruction efficiency at 5x coverage by 2.5 kb inserts 1000 3 2000 11 49 3000 4000 80 5000 91 92 6000 88 10000 Total 414 False positives 5



[Korbel et al., GenomeBiol. ('09)] **29** - Lectures.GersteinLab.org 🤲

## **Reconstruction of heterozygous** <u>insertions</u>

5x coverage by 2.5 kb inserts		5x coverage by 10 kb inserts	
Insertion size	Reconstruction efficiency	Insertion size	Reconstruction efficiency
250	0	1000	8
500	1	2000	42
750	2	3000	72
1000	1	4000	69
1250	8	5000	61
1500	3	6000	55
1750	3	7000	37
2000	1	8000	23
2250	1	9000	4
2500	0	10000	1
2750	0		
3000	0		
False positives	4		4

Better coverage and fewer reads allow to relax cutoff on outlier lengths and reconstruct more insertions

[Korbel et al., GenomeBiol. ('09)]

**30** - Lect

## Local Reassembly



# Simple Local Assembly: iterative contig extension



G Iterative contig elongation with the best supported extension

### **Optimal integration of sequencing technologies:** Local Reassembly of large novel insertions

Given a fixed budget, what are the sequencing coverage A, B and C that can achieve the maximum reconstruction rate (on average/worst-case)? Maybe a few long reads can bootstrap reconstruction process.



### **Optimal integration of sequencing technologies:** *Need Efficient Simulation*

Different combinations of technologies (i.e. read lenghs) very expensive to actually test.

Also computationally expensive to simulate.

(Each round of whole-genome assembly takes >100 CPU hrs; thus, simulation exploring 1K possibilities takes 100K CPU hr)

**C** Simplification of the simulation to the insertion region only



### **Optimal integration of sequencing technologies:** Efficient Simulation Toolbox using Mappability Maps



### Optimal integration of sequencing technologies: Simulation shows power of PEs

Simulation results w/ shotgun & paired-end reads on the same ~10Kb insertion



### **Optimal integration of sequencing technologies:** Simulation shows combination better than single technology



# Analyzing Repeated Blocks in the Genome (SDs & CNVs)



### SEGMENTAL DUPLCATIONS AND COPY NUMBER VARIANTS ARE RELATED PHENOMENA AND HAVE BEEN CREATED BY SEVERAL DIFFERENT MECHANISMS



#### NAHR

(Non-allelic homologous recombination)

Flanking repeat (e.g. Alu, LINE...)



#### NHEJ

(Non-homologous-endjoining)

No (flanking) repeats. In some cases <4bp microhomologies

## PERFORM LARGE SCALE CORRELATION ANALYSIS TO DETECT REPEAT SIGNATURES OF SDs AND CNVs



#### OLDER SDs ARE MUCH MORE LIKELY TO BE FORMED BY ALU ELEMENTS



- The co-localization of Alu elements with SDs is highly significant.
- Older SDs have a much higher association with Alus than younger SDs.
- Hence it is likely, that Alu elements were more active in mediating NAHR in the past (consistent with the Alu burst)

## FOCUSSING ON SDS: SDS CAN PROPAGATE THEMSELVES, WHICH LEADS TO A POWER-LAW DISTRIBUTION



#### Corollary

- SDs can mediate NAHR and lead to the formation of CNVs
- CNVs can become fixed and then be SDs
- Such mechanisms ("preferential attachment") are well studied in physics and should leads a very skewed ("power-law") distribution of SDs.



[Kim et al. Gen. Res. (submitted, '08), arxiv.org/abs/0709.4200v1]

### FOCUSSING ON SDS: SDs COLOCALIZE WITH EACH OTHER



#### Corollary

- SDs can mediate NAHR and lead to the formation of CNVs
- CNVs can become fixed and then be SDs
- SDs of similar age should co-localize better with each other:



#### Pseudogenes & CNV/SDs (whole genome, not just encode pilot)



[Kim et al. Gen. Res. (submitted, '08), arxiv.org/abs/0709.4200v1 ]

**CNVs ARE LESS** 

#### **ASSOCIATIONS ARE DIFFERENT FOR SDs AND CNVs**

#### **ASSOCIATED WITH** SD association with repeats **SDs THAN THE GENERAL SD TREND** 0.27 CNV 0.21 0.094 Association 0.07 with SDs Alu Microsatellite Pseudogenes LINE 0.31 (<0.001 (<0.001) (0.046) 0.001 0.11 **CNV** association with repeats 0.0739 0.048 0.0466 0.0006 >99% SDs\* CNVs Microsatellite Pseudogenes LINE Alu < 0.001 0.92 0.046 0.001

[Kim et al. Gen. Res. ('08), arxiv.org/abs/0709.4200v1 ]



AFTER THE ALU BURST, THE **IMPORTANCE OF ALU ELEMENTS FOR GENOME** REARRANGEMENT DECLINED RAPIDLY

- About 40 million years ago there was a burst in retrotransposon activity
- The majority of Alu elements stem from that time
- This, in turn, led to rapid genome rearrangement via NAHR
- The resulting SDs, could create more SDs, but with Alu activity decaying, their creation slowed





### Identifying Structural Variants in Human Population

- BreakPtr
  - Ø Model-based segmentation using bivariate HMM
- MSB
  - Mean-shift segmentation approach following grad. of PDF
  - Equally applied to aCGH and depth of coverage of short reads

- PEMer
  - Detecting Variants from discordantly placed pairedends
  - Simulation to paramaterize statistical model
- ReSeqSim
  - Efficiently simulating assembly of a representative variant
  - Shows that best reconstruction has a combination of long, med. and short reads

### Analysis of Duplication in the Genome: SVs and SDs

- Large-scale analysis of existing CNVs & SDs in human genome
- SDs assoc. with Alu, pseudogenes and older SDs
- CNVs assoc. other repeats (microsat.) and not as much with SDs
- Suggestion: Alu burst 40 MYA triggered much NAHR rearrangement, then dupl. feed on itself in hotspots but now dying down and NAHR assoc. with other repeats and CNVs also from NHEJ

YK Lam J Du J Korbel L Wang P Kim A Abyzov M Snyder

X Mu, D Greenbaum, A Urban, P Cayting, J Rozowsky, R Bjornson, S Weissman, Z Zhang, S Balasubramanian

GenomeTECH.gersteinlab.org



## **More Information on this Talk**

**SUBJECT:** GenomeAssembly

DESCRIPTION: Banbury meeting on Structural Variation in the Human Genome, Lloyd Harbor, NY 2009.11.17, 9:00-9:20; [I:BANBURYCNV] (Adaption of GenomeAssembly talk, building on [I:BIBM].)

(Works equally well on mac or PC. Paper references in the talk were mostly from Papers.GersteinLab.org. The above topic list can be easily cross-referenced against this website. Each topic abbrev. which is starred is actually a papers "ID" on the site. For instance, the topic **pubnet**\* can be looked up at <a href="http://papers.gersteinlab.org/papers/pubnet">http://papers.gersteinlab.org/papers/pubnet</a> )

**PERMISSIONS**: This Presentation is copyright Mark Gerstein, Yale University, 2008. Please read permissions statement at http://www.gersteinlab.org/misc/permissions.html . Feel free to use images in the talk with PROPER acknowledgement (via citation to relevant papers or link to gersteinlab.org).

<u>PHOTOS & IMAGES</u>. For thoughts on the source and permissions of many of the photos and clipped images in this presentation see http://streams.gerstein.info . In particular, many of the images have particular EXIF tags, such as **kwpotppt** , that can be easily queried from flickr, viz: http://www.flickr.com/photos/mbgmbg/tags/kwpotppt .