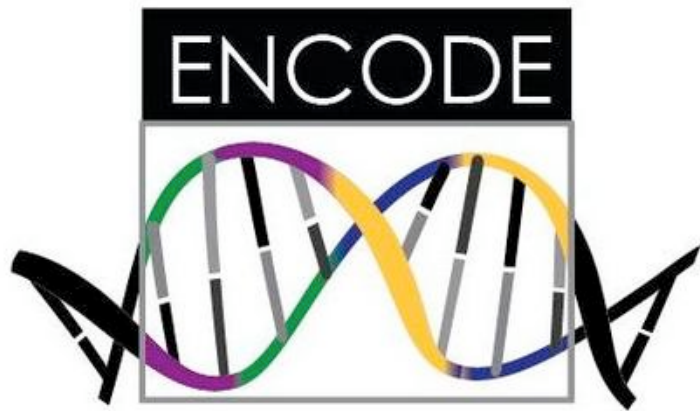


# Personalized Phased Diploid Genomes of the EN-TEEx Samples

Michael Schatz, Fritz Sedlazeck, Han Fang, Maria Nattestad, Ruibang Luo, Srividya Ramakrishnan, Charlotte Darby, Philipp Rescheneder, Alex Dobin, Carrie Davis, Ashwin Prakash, Anna Vlasova, Alessandra Breschi, Roderic Guigo, Tom Gingeras

Feb 15, 2017  
AGBT Informatics





Catalog “functional elements” through large-scale RNA-seq, ChIP-seq and other assays in many tissue types

***No DNA sequencing,  
Mostly cell lines***



Analyze how genomic variations impacts expression in many tissue types in many people

***No regulatory assays,  
Mostly unphased SNP analysis***

BIOSAMPLE

ASSAY

841 results

Clear Filters

ChIP-seq  
RNA-seq  
RAMPAGE  
small RNA-seq  
ATAC-seq  
microRNA counts  
microRNA-seq  
DNase-seq  
genotyping HTS  
eCLIP

| tissue                        | ChIP-seq | RNA-seq | RAMPAGE | small RNA-seq | ATAC-seq | microRNA counts | microRNA-seq | DNase-seq | genotyping HTS | eCLIP |
|-------------------------------|----------|---------|---------|---------------|----------|-----------------|--------------|-----------|----------------|-------|
| transverse colon              | 26       | 4       | 4       | 4             | 4        |                 |              | 3         | 8              | 4     |
| sigmoid colon                 | 31       | 4       | 4       | 4             | 4        |                 |              |           |                | 2     |
| body of pancreas              | 32       | 2       | 2       | 2             | 2        | 2               | 2            | 2         |                |       |
| adrenal gland                 | 24       | 4       | 4       | 2             | 2        | 2               | 2            | 3         |                | 2     |
| thyroid gland                 | 26       | 4       | 4       | 1             | 3        | 2               | 2            | 3         |                |       |
| gastrocnemius medialis        | 29       | 4       | 3       |               | 3        | 2               | 2            | 1         |                |       |
| stomach                       | 29       | 4       | 4       | 4             | 2        |                 |              |           |                |       |
| upper lobe of left lung       | 23       | 4       | 4       | 4             |          |                 |              |           |                |       |
| gastroesophageal sphincter    | 19       | 4       | 4       | 4             | 2        |                 |              |           |                |       |
| breast epithelium             | 23       | 3       | 2       |               | 3        |                 |              |           |                |       |
| spleen                        | 18       | 4       | 4       | 4             | 1        |                 |              |           |                |       |
| esophagus squamous epithelium | 16       | 4       | 4       | 4             |          |                 |              | 1         |                |       |
| esophagus muscularis mucosa   | 16       | 4       | 4       | 3             |          |                 |              |           |                |       |
| Peyer's patch                 | 16       | 4       | 4       | 2             |          |                 |              |           |                |       |
| suprapubic skin               | 16       | 4       | 4       | 4             |          |                 |              |           |                |       |
| tibial nerve                  | 15       | 4       | 3       | 1             |          |                 |              |           |                |       |
| heart left ventricle          | 11       | 2       | 2       |               | 2        | 2               | 2            |           |                |       |
| lower leg skin                | 11       | 4       | 2       |               |          | 2               | 2            |           |                |       |
| omental fat pad               | 11       | 4       | 4       | 1             | 1        |                 |              |           |                |       |
| subcutaneous adipose tissue   | 12       | 4       | 3       | 1             | 1        |                 |              |           |                |       |
| vagina                        | 13       | 1       | 1       | 1             |          | 2               | 2            |           |                |       |
| prostate gland                | 8        | 2       | 2       | 2             |          | 2               | 2            | 1         |                |       |
| ovary                         | 9        | 2       | 2       | 1             | 1        | 1               | 1            | 1         |                |       |
| right lobe of liver           | 8        | 2       | 2       | 2             | 1        | 1               | 1            | 1         |                |       |
| thoracic aorta                | 15       | 2       | 1       |               |          |                 |              |           |                |       |
| uterus                        | 8        | 2       | 1       | 1             |          | 2               | 2            | 2         |                |       |
| testis                        | 6        | 2       | 2       | 2             |          | 2               | 2            | 1         |                |       |
| ascending aorta               | 15       | 1       |         |               |          |                 |              |           |                |       |
| right atrium auricular region | 6        | 2       | 2       |               |          | 2               | 2            |           |                |       |
| tibial artery                 | 5        |         |         |               | 1        |                 |              |           |                |       |
| coronary artery               | 5        |         |         |               |          |                 |              |           |                |       |

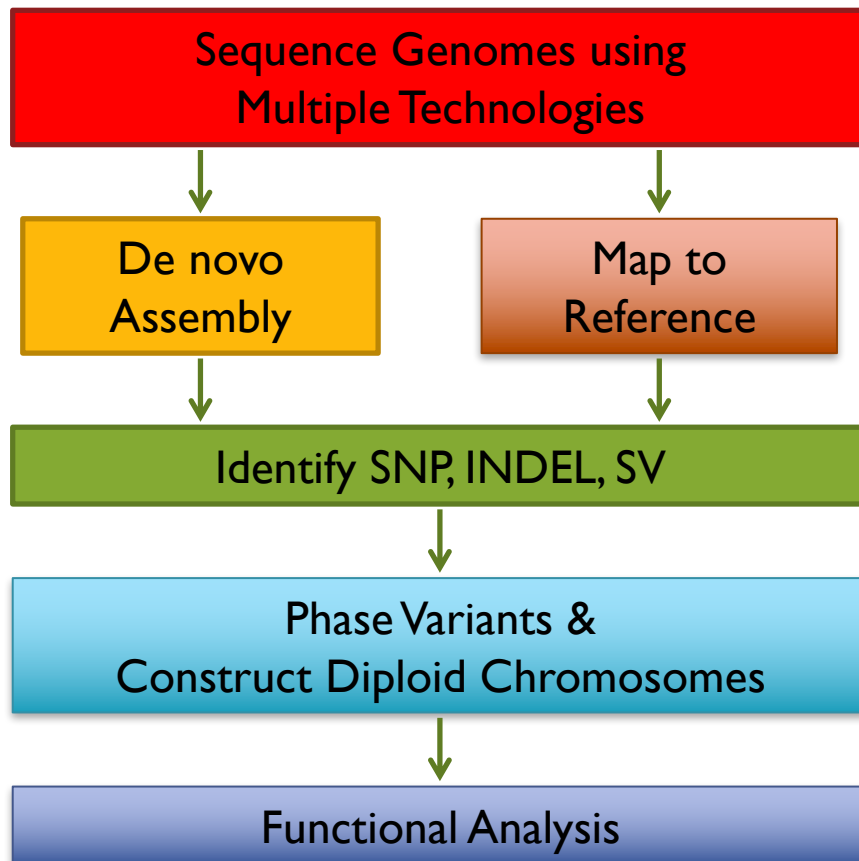
# EN-TEEx: Expression & Regulation Analysis of Personalized Genomes

|                        | ENC-001 | ENC-002 | ENC-003                    | ENC-004                    |
|------------------------|---------|---------|----------------------------|----------------------------|
| <b>Age</b>             | 37      | 54      | 53                         | 51                         |
| <b>Sex</b>             | Male    | Male    | Female                     | Female                     |
| <b>Cause of Death</b>  | Anoxia  | Anoxia  | Cerebral Vascular Accident | Cerebral Vascular Accident |
| <b>Total Libraries</b> | 319     | 299     | 488                        | 299                        |

- Sequenced the genomes for 2 male and 2 female samples using transverse colon tissue
- Large number of ChIP-seq, RNA-seq, ATAC-seq, DNase-seq, and other functional datasets available in dozens of tissues

<http://encodeproject.org>

# Assembling and Analyzing Personal Genomes

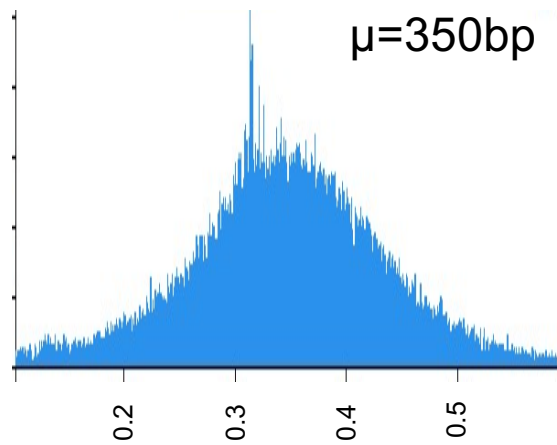


## *Goals*

1. What are the most effective biotechnologies for sequencing?
2. What do we learn from a personalized genome instead of the reference?
3. Can we use the genomic variants as natural perturbations of the encyclopedia elements?

# Genomic Sequencing Data

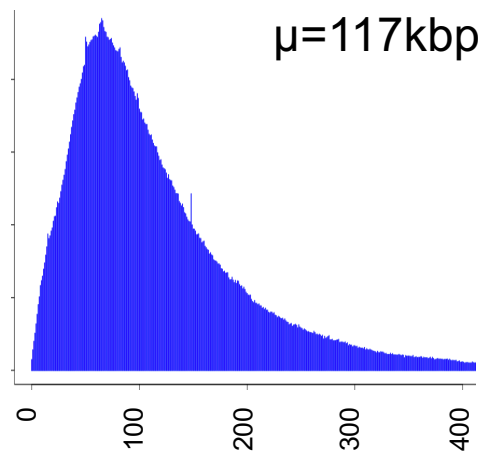
## Illumina



Fragment Length (kbp)

**60x Paired End**  
*All 4 samples*

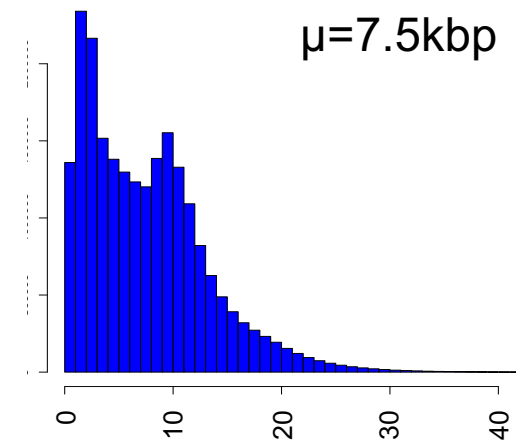
## 10X Genomics



Molecule Length (kbp)

**35x Linked Reads**  
*All 4 samples*

## PacBio

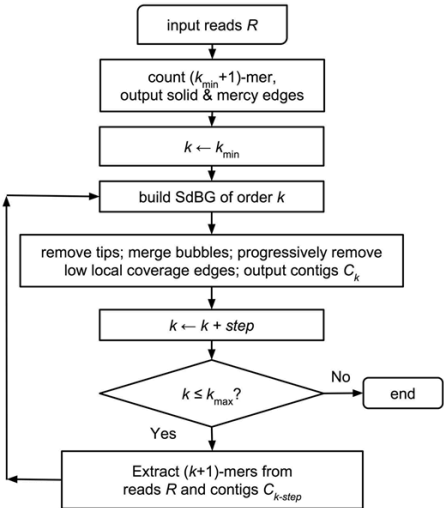


Read Length (kbp)

**55x Long Reads**  
*\*Only ENC-002*

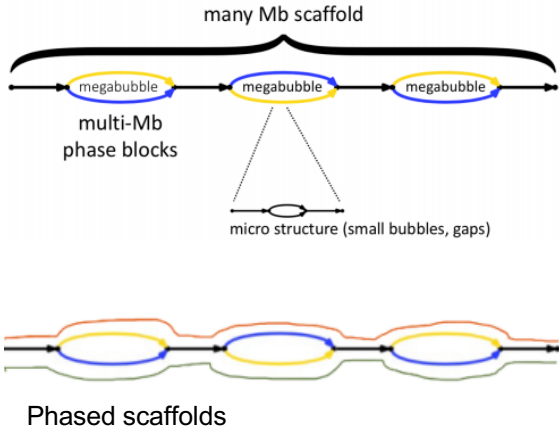
# De Novo Assembly

## Illumina



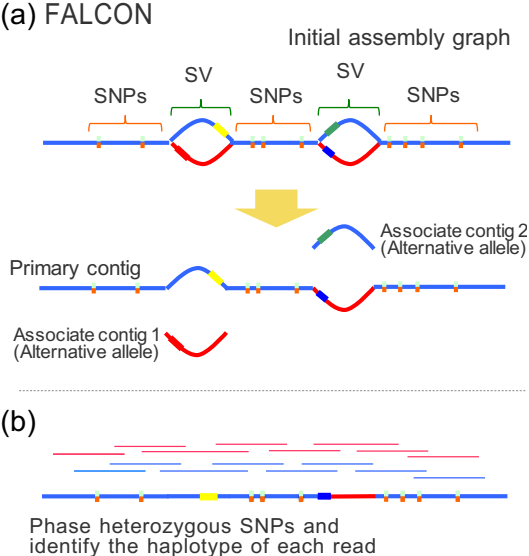
**MegaHit**  
(Li et al, 2015)

## 10X Genomics



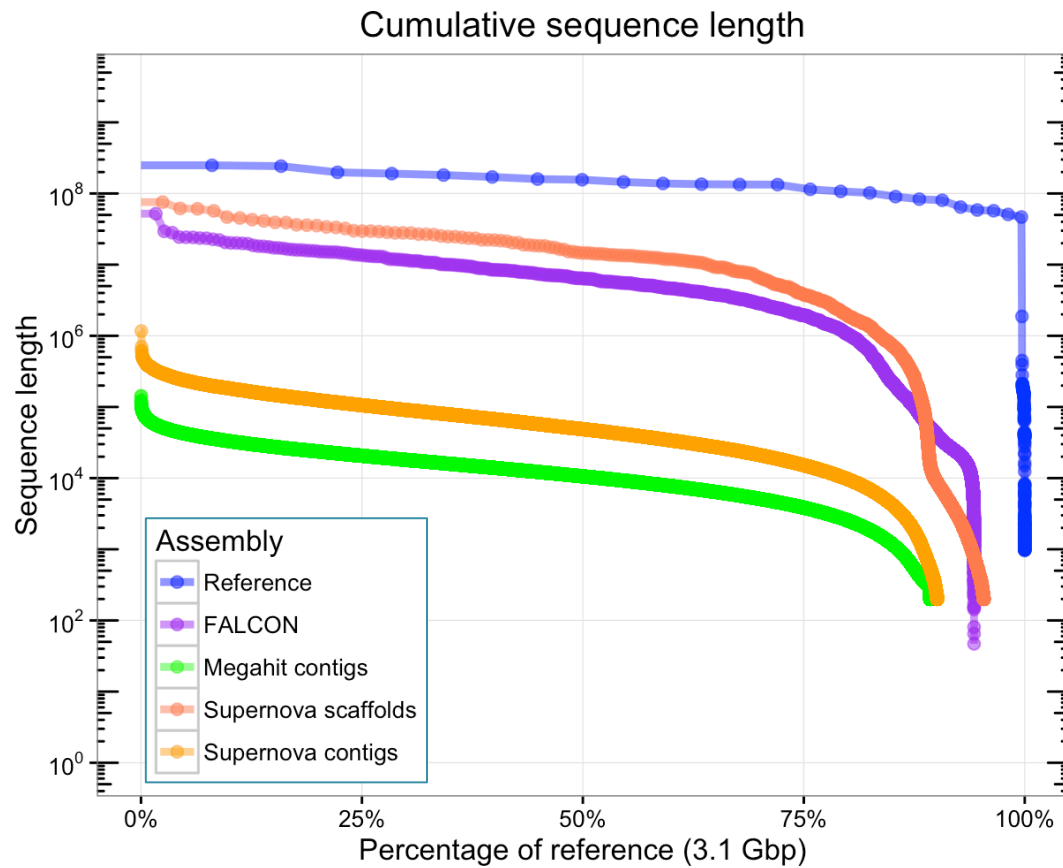
**SuperNova**  
(Weisenfeld et al, 2016)

## PacBio



**FALCON-unzip**  
(Chin et al, 2016)

# Assembly Contiguity



## GRC38 Reference

- Includes alt sequences

## 10X Genomics/SuperNova

- 21 Mbp scaffold N50
- 162 Mbp in scaffold gaps

## PacBio/Falcon-unzip

- 7.0 Mbp contig N50

## 10X Genomics/Supernova

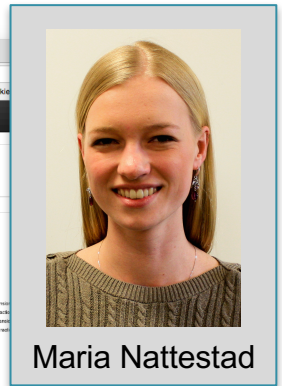
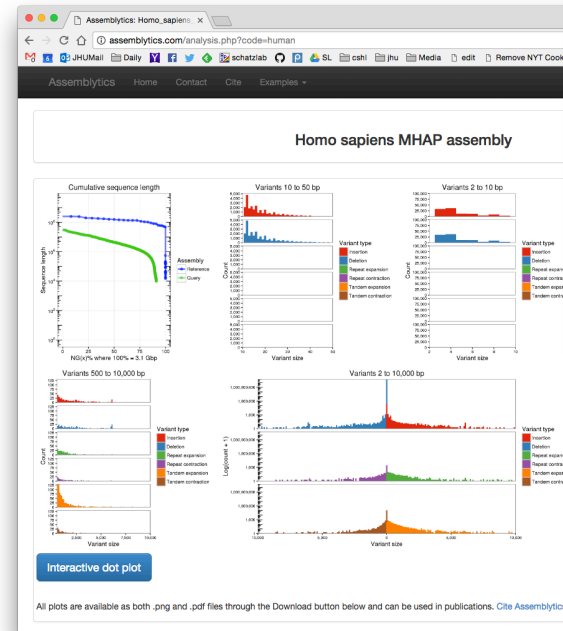
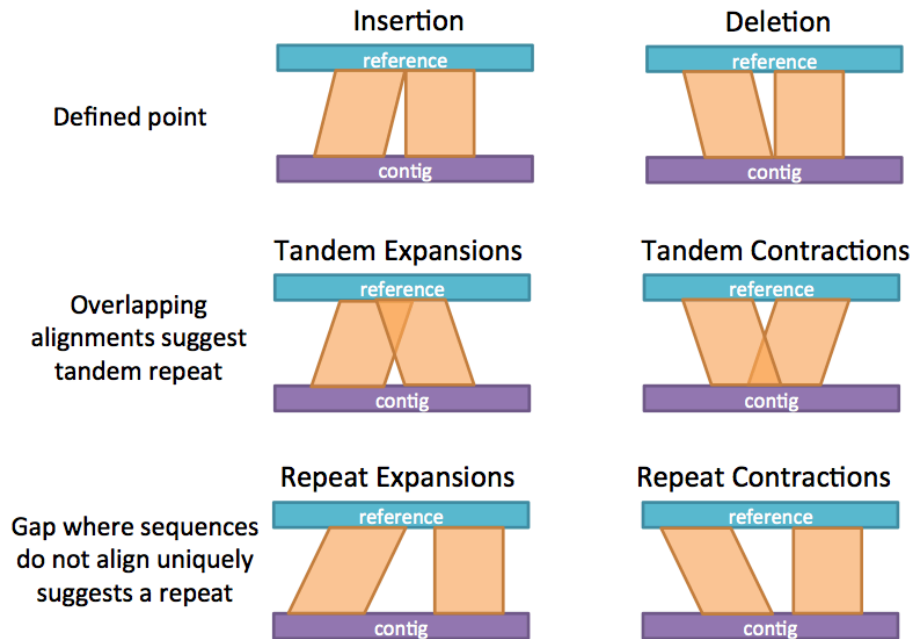
- 50 kbp contig N50

## Illumina/MegaHit

- 13 kbp contig N50

# Assemblytics: Assembly-Based Variant-Caller

<http://assemblytics.com>



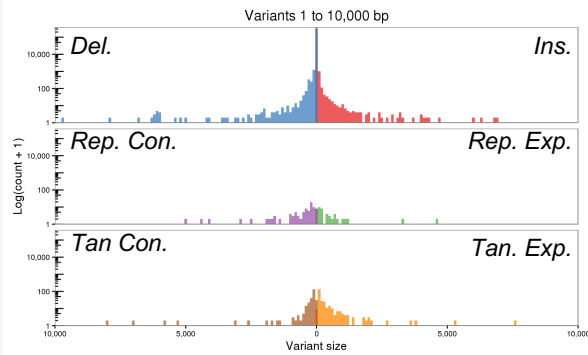
**Assemblytics: a web analytics tool for the detection of variants from an assembly**

Nattestad, M, Schatz, MC (2016) *Bioinformatics* doi: 10.1093/bioinformatics/btw369



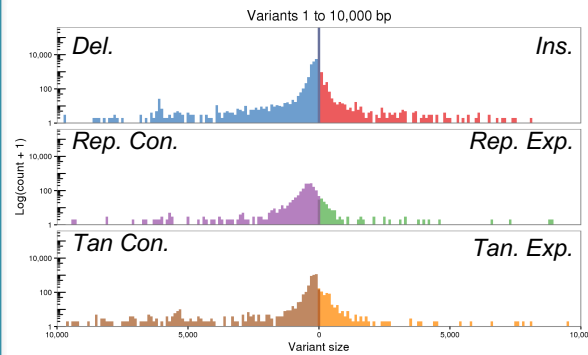
# Structural Variations vs GRCh38

## Illumina



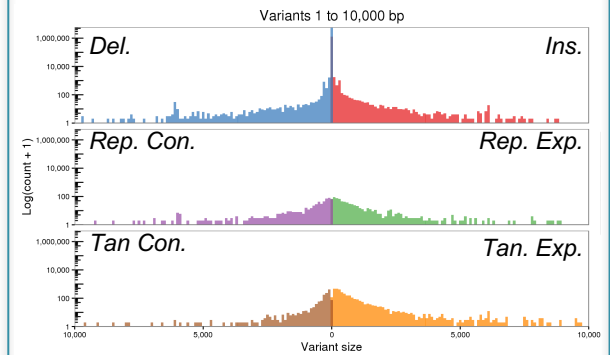
**SVs (50bp – 10kb)  
Count: 3,997  
Bases: 1.11 Mbp**

## 10X Genomics



**SVs (50bp – 10kbp)  
Count: 18,025  
Bases: 6.13 Mbp**

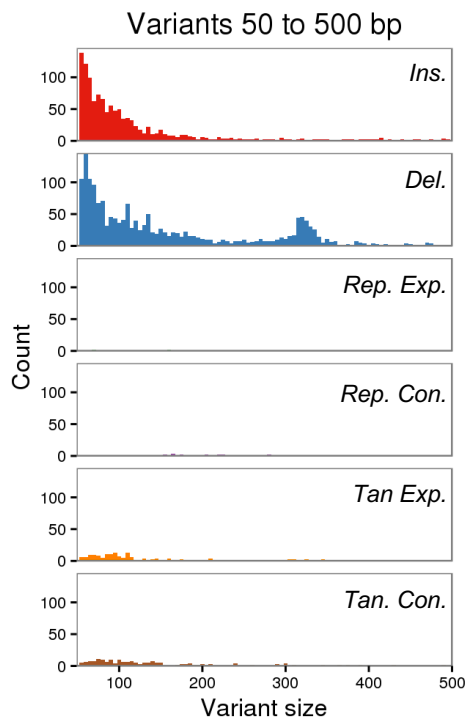
## PacBio



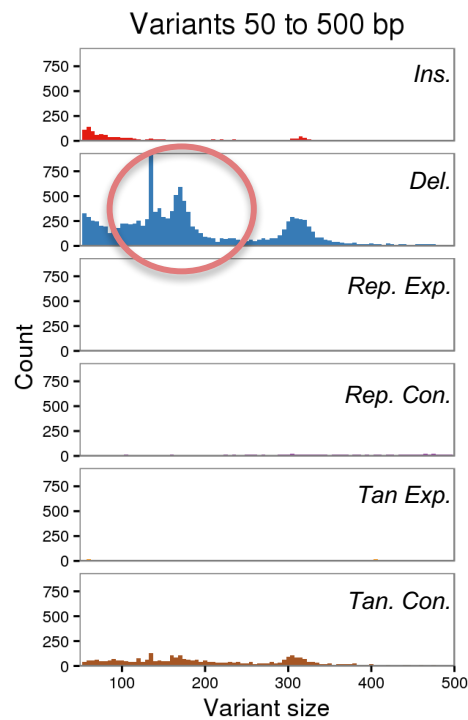
**SVs (50bp-10kbp)  
Count: 12,965  
Bases: 8.13 Mbp**

# Missing Insertions from Short and Linked Read?

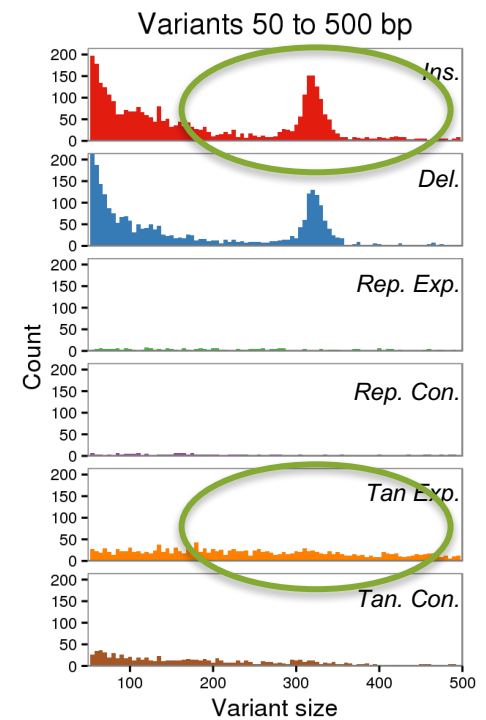
## Illumina



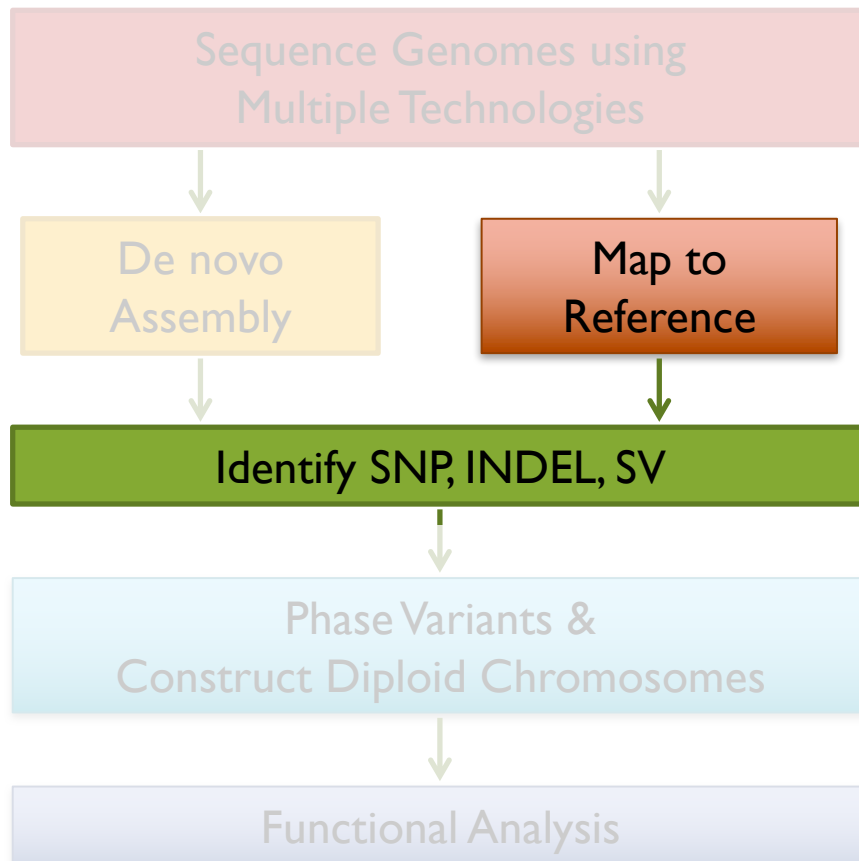
## 10X Genomics



## PacBio



# Assembling and Analyzing Personal Genomes

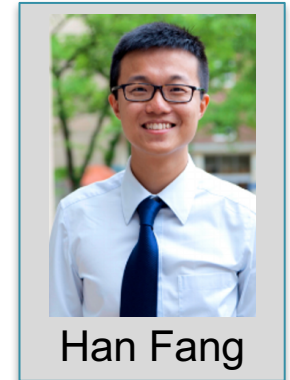


## *Goals*

1. What are the most effective biotechnologies for a sequencing?
2. What do we learn from an personalized genome instead of the reference?
3. Can we use the genomic variants as natural perturbations of the encyclopedia elements?

# Small Variant Analysis

- Mapped Illumina PE reads using BWA-MEM (Li, 2013)
- Identified 3.7M SNPs using GATK (Van der Auwera *et al.* 2013)
- Identified 700k indels using Scalpel (Fang *et al.*, 2016)
- Annovar (Wang *et al.*, 2010) characterization of variants

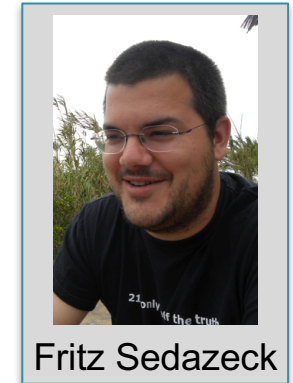


|                                    | ENC-001  | ENC-002  | ENC-003  | ENC-004  |
|------------------------------------|----------|----------|----------|----------|
| <b><i>Synonymous SNP</i></b>       | 12,007   | 12,249   | 12,524   | 12,172   |
| <b><i>Non Syn. SNP</i></b>         | 11,507   | 11,816   | 12,078   | 12,009   |
|                                    |          |          |          |          |
| <b><i>Frameshift Indel</i></b>     | 304      | 344      | 344      | 322      |
| <b><i>Stop Gain + Loss</i></b>     | 113 / 27 | 120 / 33 | 136 / 25 | 135 / 28 |
| <b><i>Splicing SNP + Indel</i></b> | 109 / 41 | 102 / 43 | 117 / 55 | 111 / 50 |

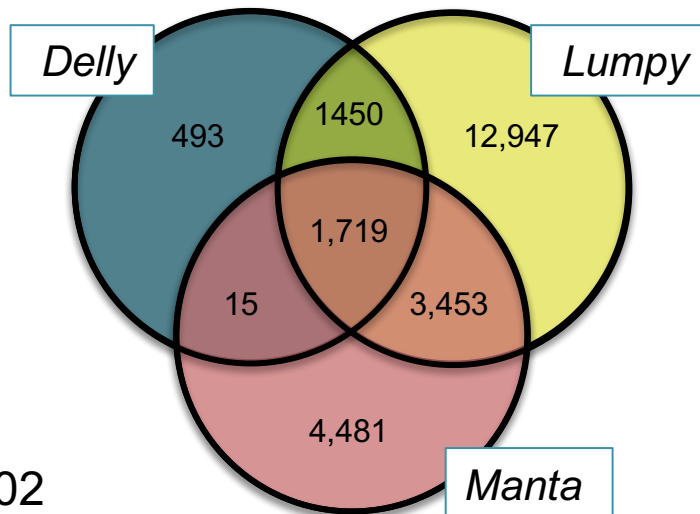
***Establishes a catalog of variation, heterozygous positions informative for phasing***

# Consensus SV Analysis with SURVIVOR

<https://github.com/fritzsedlazeck/SURVIVOR>



- Analyzed the Illumina PE sequence data using 3 different algorithms that use split-reads and discordant pairs to identify SVs: Manta (Chen *et al.*, 2015), Delly (Rausch *et al.*, 2012), and Lumpy (Layer *et al.*, 2014)
- Use SURVIVOR (Jeffares *et al.*, 2016) to improve accuracy by excluding variants identified by only 1 method



| Type           | SURVIVOR2    |
|----------------|--------------|
| Deletions      | 3,692        |
| Duplications   | 1,144        |
| Insertions     | 253          |
| Inversions     | 602          |
| Translocations | 676          |
| <b>All</b>     | <b>6,367</b> |

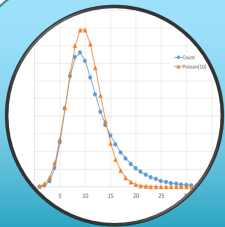
ENC-002

# 10X LongRanger Variant Calls

- Recent versions of LongRanger report SVs in addition to phasing
- Two classes: large\_sv.vcf (>30kbp) and dels.vcf (40bp to 30kbp)



|         | Large SV | DEL | DUP | INV | Small Del. | Del. Span | Del. Mean |
|---------|----------|-----|-----|-----|------------|-----------|-----------|
| ENC-001 | 96       | 21  | 7   | 7   | 4,022      | 3.77 Mbp  | 937 bp    |
| ENC-002 | 194      | 36  | 6   | 6   | 3,796      | 3.21 Mbp  | 852 bp    |
| ENC-003 | 96       | 32  | 2   | 8   | 4,055      | 3.74 Mbp  | 927 bp    |
| ENC-004 | 103      | 33  | 1   | 4   | 4,294      | 3.43 Mbp  | 805 bp    |



## **LRSim: Linked Read Simulator**

Lead Author: Ruibang Luo  
<https://github.com/aquaskyline/LRSIM>  
bioRxiv: <https://doi.org/10.1101/103549>



## **TopSorter: 10X SV Analysis**

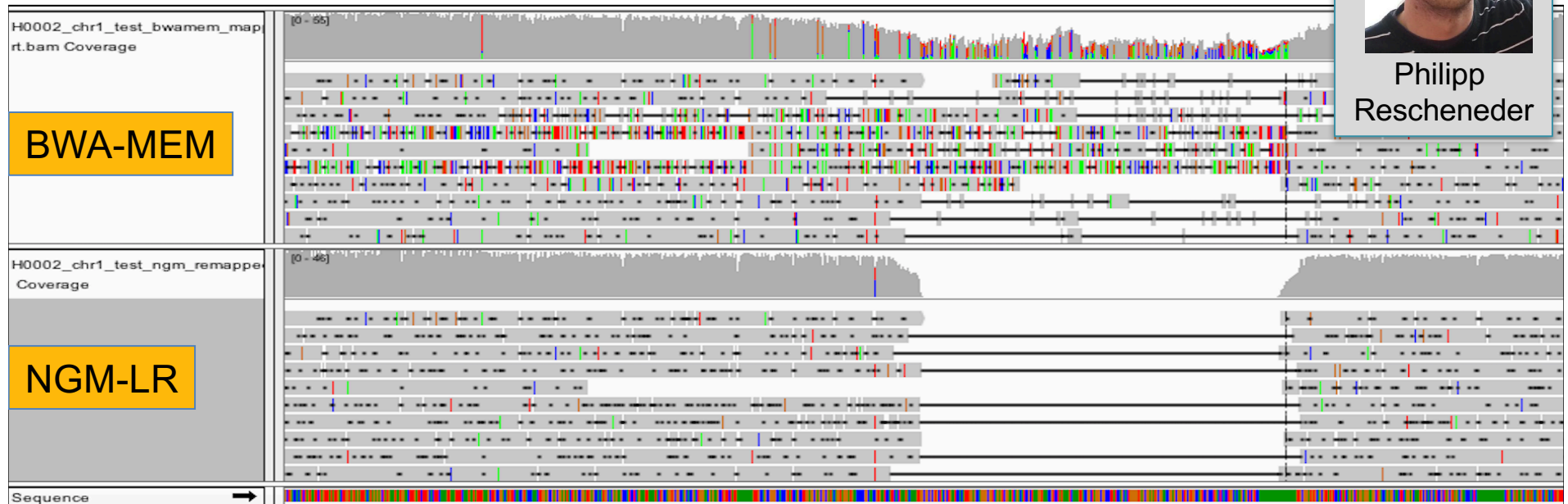
Lead Author: Han Fang  
<https://github.com/hanfng/TopSorter>  
Preprint: Coming soon

# NGM-LR + Sniffles: PacBio SV Analysis Tools

<https://github.com/philres/ngmlr> & <https://github.com/fritzsedlazeck/Sniffles>



Philipp  
Rescheneder



## Improved SV Variant Detection with long reads

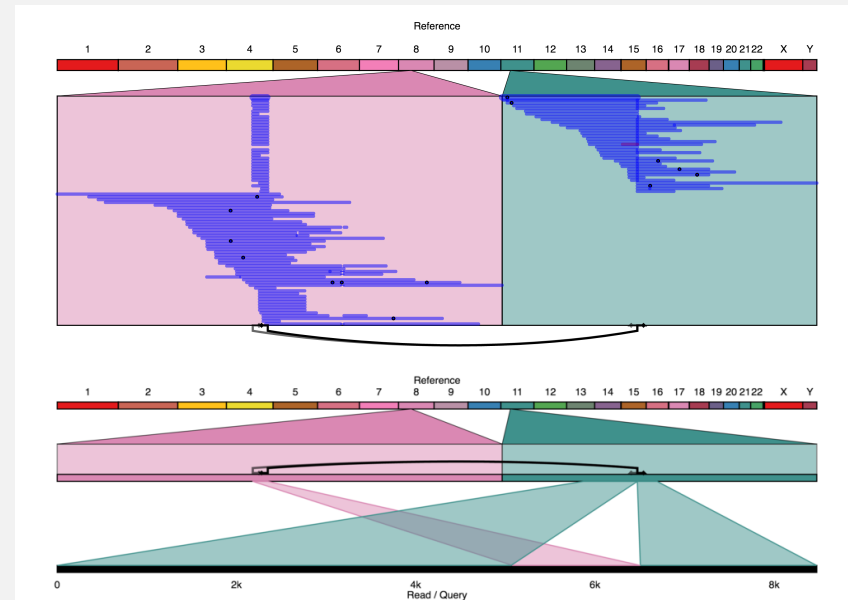
1. **NGM-LR**: Improve mapping of noisy long reads: improved seeding, convex gap scoring
2. **Sniffles**: Integrates evidence from split-reads, alignment fidelity, breakpoint concordance

# Sniffles PacBio Variant Calls

## Sniffles calls

|                | All SVs<br>(50bp+) | Large SVs<br>(10kbp+) |
|----------------|--------------------|-----------------------|
| Deletions      | 7,389              | 164                   |
| Duplications   | 1,284              | 139                   |
| Insertions     | 8,382              | 4                     |
| Inversions     | 229                | 116                   |
| Translocations | 170                | 170                   |
| <b>All</b>     | <b>17,454</b>      | <b>593</b>            |

## Translocation in Ribbon



**Ribbon: Visualizing complex genome alignments and structural variation**

Nattestad et al. (2016) *bioRxiv* doi: <http://dx.doi.org/10.1101/082123>



# Structural Variations Concordance

|            |        |        |       |        |       |       |
|------------|--------|--------|-------|--------|-------|-------|
| Sniffles   | 17,139 |        |       |        |       |       |
| Falcon     | 7,857  | 12,241 |       |        |       |       |
| LongRanger | 2,823  | 1,946  | 3,785 |        |       |       |
| SuperNova  | 3,394  | 2,837  | 1,486 | 18,862 |       |       |
| SURVIVOR2  | 3,291  | 2,163  | 2,274 | 1,646  | 6,631 |       |
| MegaHit    | 1,858  | 1,529  | 569   | 1,378  | 687   | 3,855 |

# Structural Variations Concordance

|            |          |        |            |           |           |         |
|------------|----------|--------|------------|-----------|-----------|---------|
| Sniffles   | 17,139   |        |            |           |           |         |
| Falcon     | 7,857    | 12,241 |            |           |           |         |
| LongRanger | 2,823    | 1,946  | 3,785      |           |           |         |
| SuperNova  | 3,394    | 2,837  | 1,486      | 18,862    |           |         |
| SURVIVOR2  | 3,291    | 2,163  | 2,274      | 1,646     | 6,631     |         |
| MegaHit    | 1,858    | 1,529  | 569        | 1,378     | 687       | 3,855   |
|            | Sniffles | Falcon | LongRanger | SuperNova | SURVIVOR2 | MegaHit |

## **Main Diagonal**

- Calls per tool

# Structural Variations Concordance

|            |          |        |            |           |           |         |              |
|------------|----------|--------|------------|-----------|-----------|---------|--------------|
| Sniffles   | 17,139   |        |            |           |           |         | PacBio       |
| Falcon     | 7,857    | 12,241 |            |           |           |         |              |
| LongRanger | 2,823    | 1,946  | 3,785      |           |           |         | 10X Genomics |
| SuperNova  | 3,394    | 2,837  | 1,486      | 18,862    |           |         |              |
| SURVIVOR2  | 3,291    | 2,163  | 2,274      | 1,646     | 6,631     |         | Illumina     |
| MegaHit    | 1,858    | 1,529  | 569        | 1,378     | 687       | 3,855   |              |
|            | Sniffles | Falcon | LongRanger | SuperNova | SURVIVOR2 | MegaHit |              |

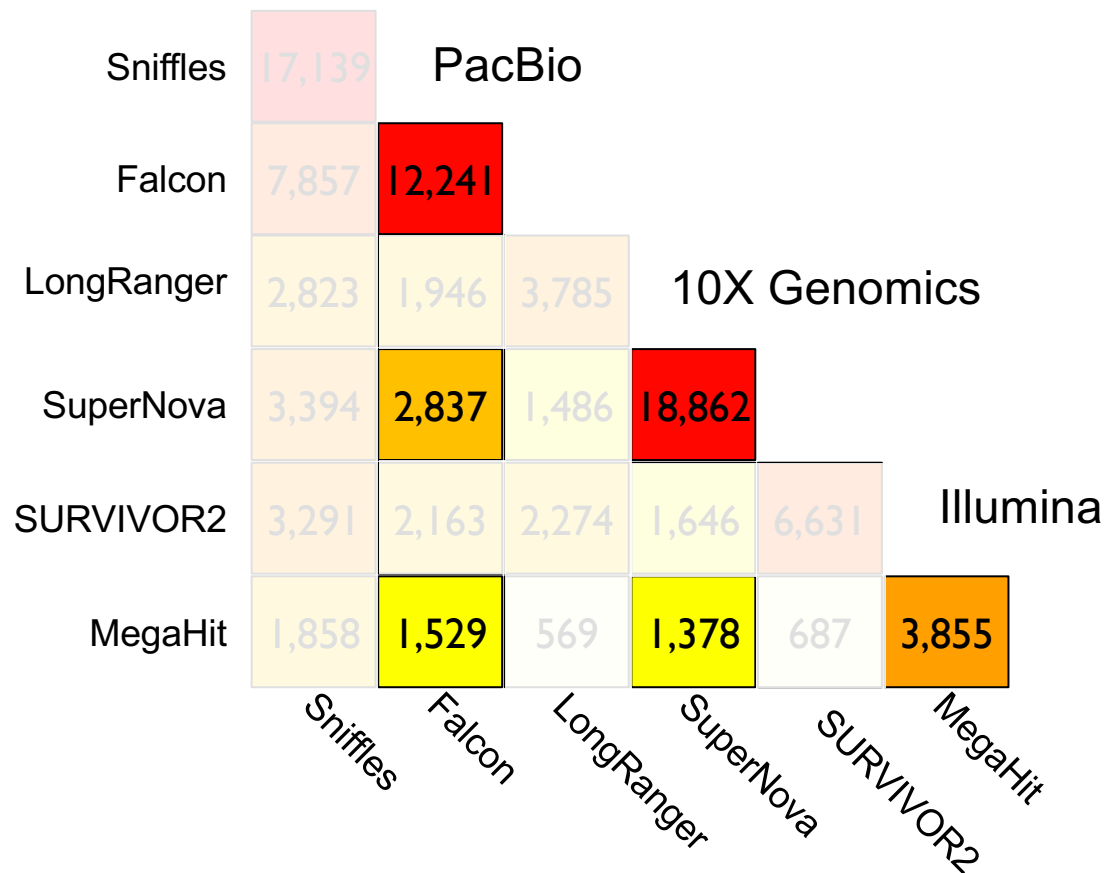
## **Main Diagonal**

- Calls per tool

## **Outer triplets**

- Concordance by Technology

# Structural Variations Concordance



## Main Diagonal

- Calls per tool

## Outer triplets

- Concordance by Technology

## Inner triplets

- Concordance by Assembly

# Structural Variations Concordance

|            |          |        |            |           |           |         |              |
|------------|----------|--------|------------|-----------|-----------|---------|--------------|
| Sniffles   | 17,139   |        |            |           |           |         | PacBio       |
| Falcon     | 7,857    | 12,241 |            |           |           |         |              |
| LongRanger | 2,823    | 1,946  | 3,785      |           |           |         | 10X Genomics |
| SuperNova  | 3,394    | 2,837  | 1,486      | 18,862    |           |         |              |
| SURVIVOR2  | 3,291    | 2,163  | 2,274      | 1,646     | 6,631     |         | Illumina     |
| MegaHit    | 1,858    | 1,529  | 569        | 1,378     | 687       | 3,855   |              |
|            | Sniffles | Falcon | LongRanger | SuperNova | SURVIVOR2 | MegaHit |              |

## **Main Diagonal**

- Calls per tool

## **Outer triplets**

- Concordance by Technology

## **Inner triplets**

- Concordance by Assembly
- Concordance by Mappers

# Structural Variations Concordance

|            |          |        |            |              |           |          |
|------------|----------|--------|------------|--------------|-----------|----------|
| Sniffles   | 17,139   | PacBio |            |              |           |          |
| Falcon     | 7,857    | 12,241 |            |              |           |          |
| LongRanger | 2,823    | 1,946  | 3,785      | 10X Genomics |           |          |
| SuperNova  | 3,394    | 2,837  | 1,486      | 18,862       |           |          |
| SURVIVOR2  | 3,291    | 2,163  | 2,274      | 1,646        | 6,631     | Illumina |
| MegaHit    | 1,858    | 1,529  | 569        | 1,378        | 687       | 3,855    |
|            | Sniffles | Falcon | LongRanger | SuperNova    | SURVIVOR2 | MegaHit  |

## **Main Diagonal**

- Calls per tool

## **Outer triplets**

- Concordance by Technology

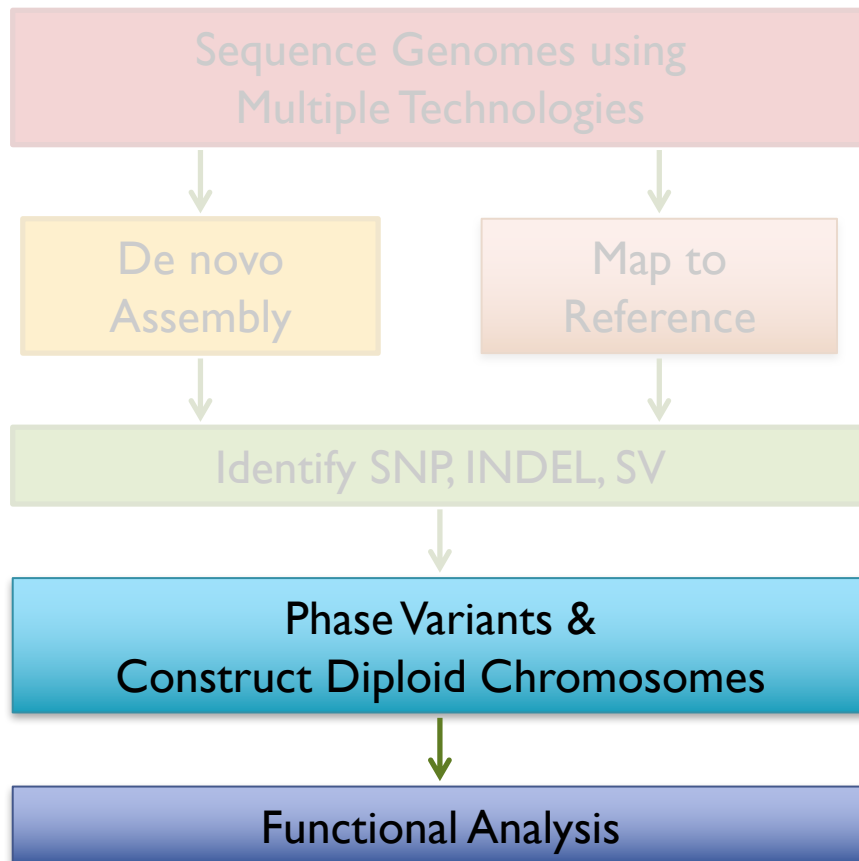
## **Inner triplets**

- Concordance by Assembly
- Concordance by Mappers

## **Overall:**

- We need multiple technologies and approaches

# Assembling and Analyzing Personal Genomes

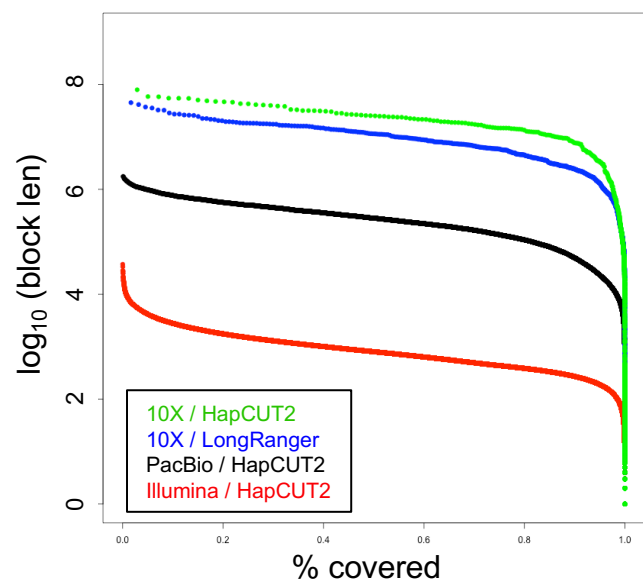


## *Goals*

1. What are the most effective biotechnologies for sequencing?
2. What do we learn from a personalized genome instead of the reference?
3. Can we use the genomic variants as natural perturbations of the encyclopedia elements?

# Phasing Results

- Phasing attempts to link together variants that came from the same molecule
- Long reads & fragments are needed to link distant heterozygous sites



**10X + HapCUT2**

**N50: 25.1 Mbp**



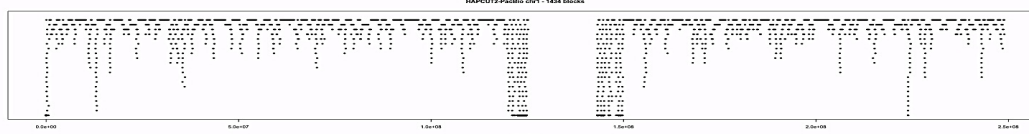
**10X + LongRanger**

**N50: 11.3 Mbp**



**PacBio + HapCUT2**

**N50: 283kbp**

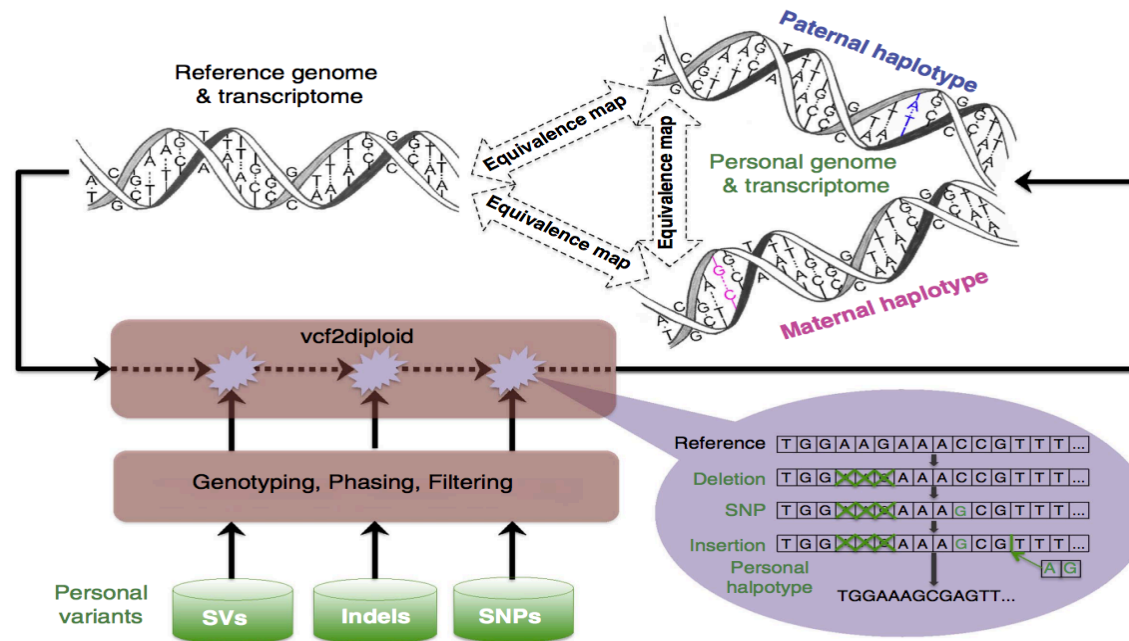


**HapCUT2: robust and accurate haplotype assembly for diverse sequencing technologies**

Edge, P, Bafna, V, Bansal, V (2016) *Genome Research*. doi: 10.1101/gr.213462.116



# AlleleSeq: Constructing the Personal Genomes



(J Rozowsky et al, 2011)

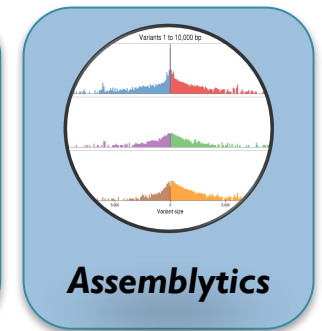
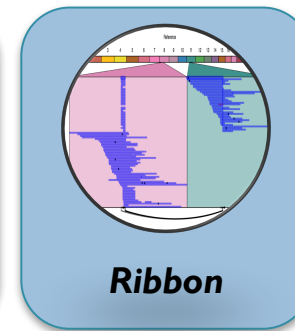
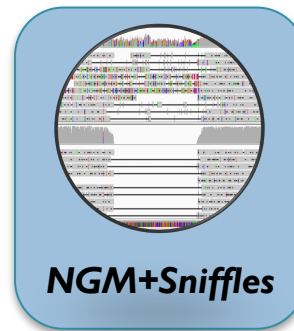
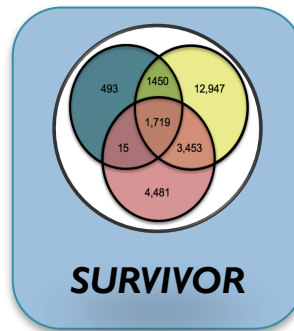
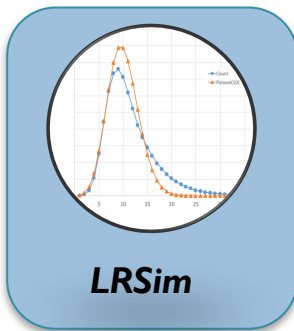
***AlleleSeq/vcf2diploid inserts phased variants from a VCF file into the reference genome to create a pair of phased chromosome fasta files***



# Personalized Phased Diploid Genomes



- **Multiple sequencing technologies & approaches needed**
  - *PacBio*: Best Resolution of SVs
  - *10X*: Best Resolution of Phasing
  - *De novo*: Best Resolution of smaller events
  - *Mapping*: Best resolution of larger events
- **We have just begun to explore the universe of variants that can be detected**
  - Tens of thousands of SVs per person`
  - Thousands of genes, thousands of regulatory elements impacted per person



<http://schatz-lab.org>

# Acknowledgements

## **Schatz Lab**

Charlotte Darby  
Han Fang  
Sam Kovaka  
Ruibang Luo  
Maria Nattestad  
Srividya  
Ramakrishnan  
Philipp  
Rescheneder  
Fritz Sedlazeck

## **Gingeras Lab**

Carrie Davis  
Alex Dobin  
Ashwin Prakash

## **McCombie Lab**

Sara Goodwin

## **Guigo Lab**

Alessandra Breschi  
Anna Vlasova

## **ENCODE Partners**

Berstein Lab  
Gerstein Lab  
Myers Lab  
Ren Lab  
Snyder Lab  
Stam Lab  
Wold Lab

+ All ENCODE  
Members



National Human  
Genome Research  
Institute



ALFRED P. SLOAN  
FOUNDATION



PACIFIC  
BIOSCIENCES®

10x GENOMICS®



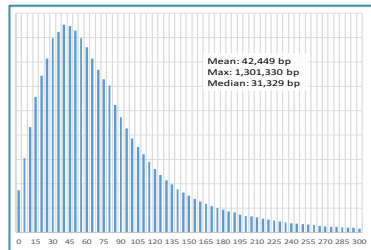
**Now hiring postdocs!**  
<http://schatz-lab.org/apply>

# Thank you

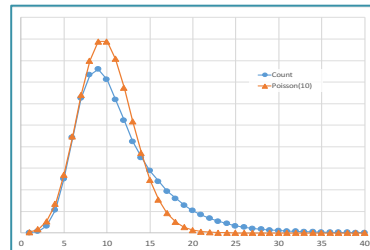
<http://schatz-lab.org>  
[@mike\\_schatz](#)

# LRSim: Linked Read Simulator

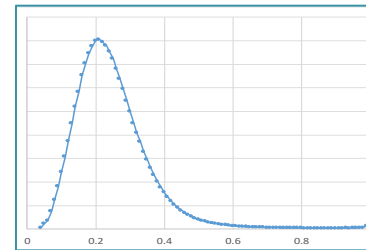
<https://github.com/aquaskyline/LRSIM>



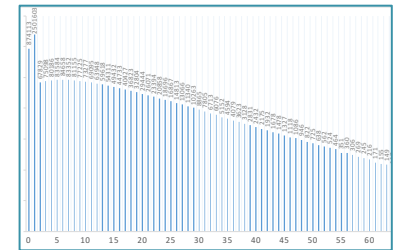
Molecule size ( $f$ )



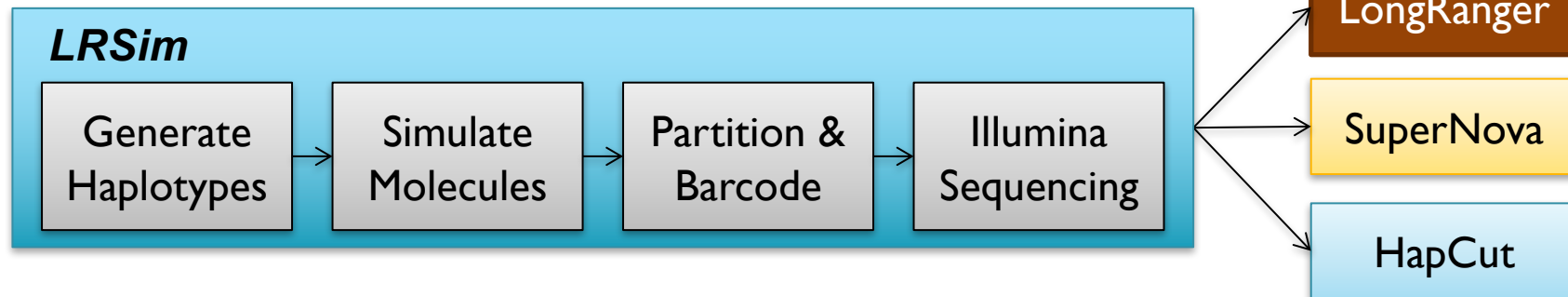
# of molecules per partition ( $m$ )



Coverage per molecule ( $c$ )



# of partitions ( $t$ )



**LRSim: a Linked Reads Simulator generating insights for better genome partitioning**

Luo, R et al MC (2017) bioRxiv doi: <https://doi.org/10.1101/103549>