

Peng Zhang<sup>1</sup>, Hua Ling<sup>1</sup>, Elizabeth Pugh<sup>1</sup>, Kurt Hetrick<sup>1</sup>, Dane Witmer<sup>1</sup>, Nara Sobreira<sup>2</sup>, David Valle<sup>2</sup>, Kimberly Doheny<sup>1</sup>

<sup>1</sup>Center for Inherited Disease Research, Institute of Genetic Medicine, The Johns Hopkins School of Medicine

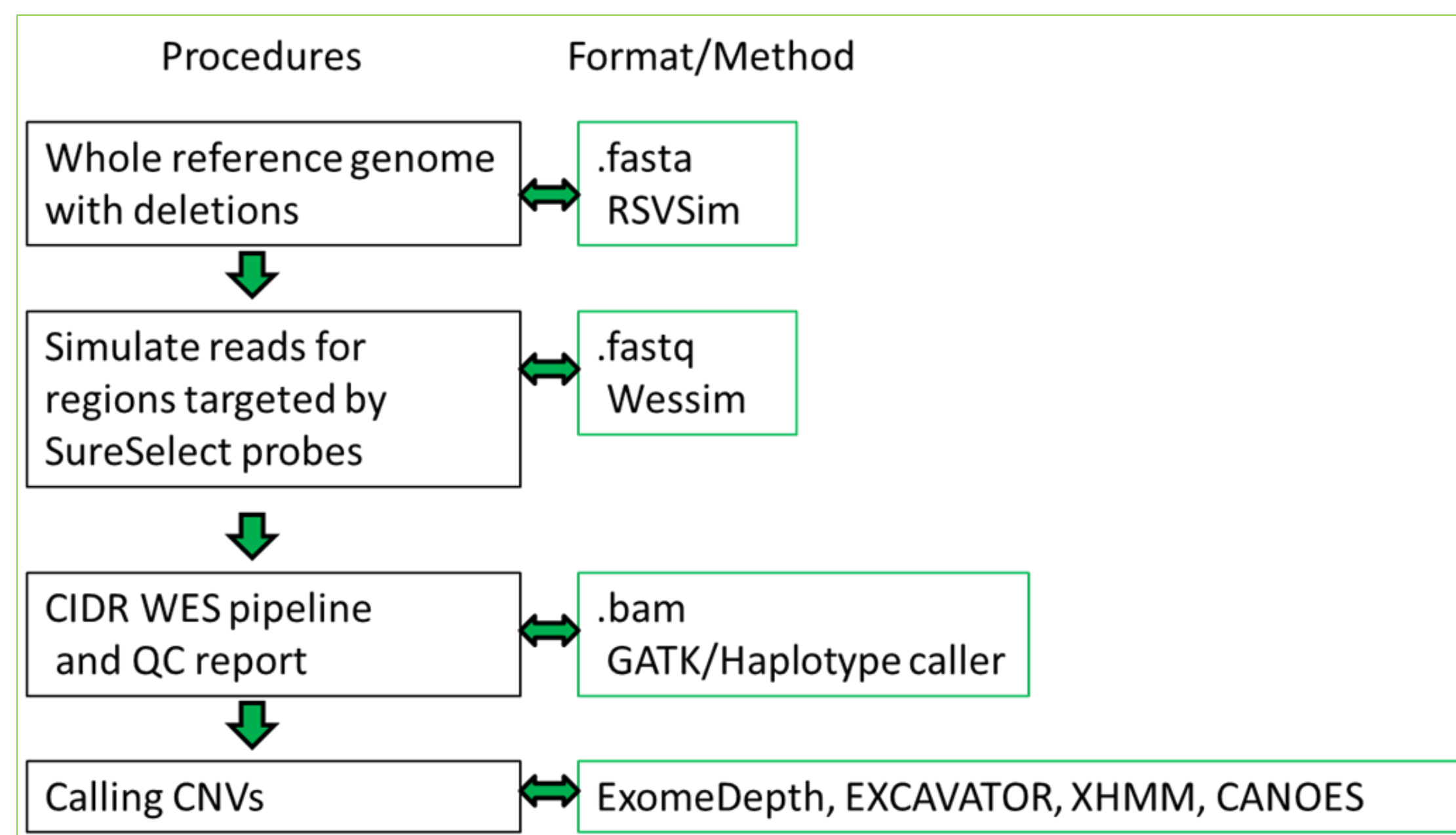
<sup>2</sup>Institute of Genetic Medicine, Johns Hopkins University School of Medicine

## Introduction

As part of the Baylor Hopkins Center for Mendelian Genomics (BHCMG) (<http://www.mendelian.org/>), CIDR has performed whole exome sequencing (WES) for over 1000 samples to discover the genetic basis for Mendelian conditions. However, calling Copy Number Variants (CNVs) from WES is a challenge because of the sparseness of the target regions and the non-uniform distribution of reads across genome. As a result, the concordance across different methods is usually low, and it's often hard to evaluate the CNV calls from real sequencing data. We present here a simulation framework for evaluating the performance of ExomeDepth, EXCAVATOR, XHMM and CANOES. We simulated WES data with CNVs of various sizes and frequencies. We then evaluate the sensitivity and specificity of each CNV calling programs under different scenarios.

## Materials and methods

**Figure 1:** Software and procedures in simulating WES data.



### Simulation design:

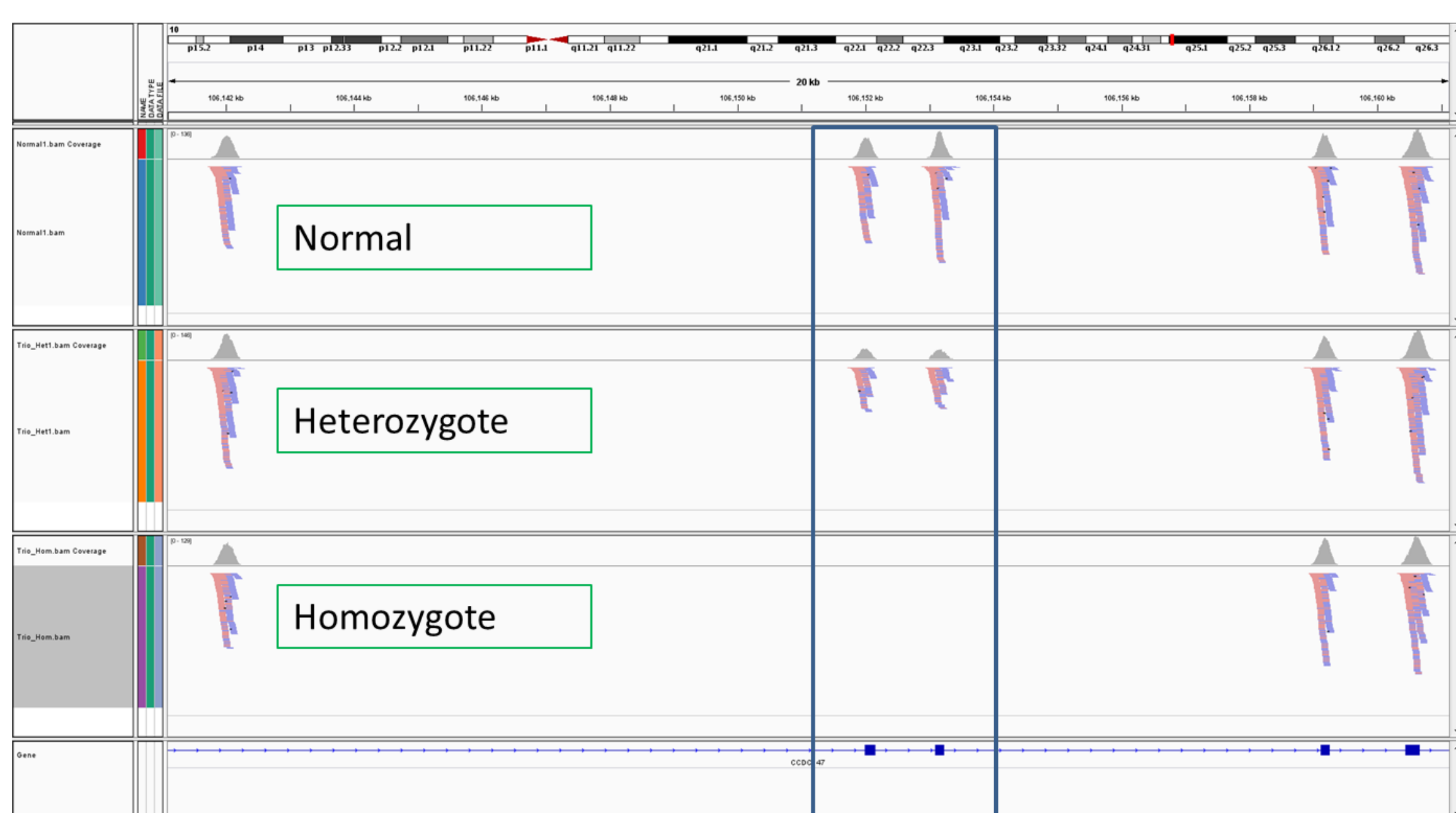
We simulated a total of 50 samples that carry each of the five sets of deletions in the whole genome reference sequence including:

- 1): **Singleton**, 1 sample with this set of homozygous deletions.
  - 2): **Doubleton**, 2 samples, both are homozygous deletions.
  - 3): **Trio**, 1 sample with homozygous deletions, 2 with heterozygous deletions.
  - 4): **Pct10**, 1 sample with homozygous deletions, 8 samples with heterozygous deletions.
  - 5): **Pct20**, 2 samples with homozygous deletions, 16 samples with heterozygous deletions.
- The rest of 17 samples are normal samples with no deletions.

**Size:** 50 bp to 50 Kb

**Mean target coverage:** 110X

**Figure 2:** Example of simulated WES data with two-exon deletion shown in IGV (Chr10: 106151892-106157746).



### Running the CNV calling methods:

ExomeDepth can select reference control from a set of sample based on pairwise correlation of read counts. We ran EXCAVATOR using the first five samples ExomeDepth picked for controls. For XHMM, we added 250 real WES samples from BHCMG to reach an appropriate sample size. For CANOES, we ran the 50 simulated samples alone.

## Results

**Table 1:** Detection of homozygous deletions.

		Singleton (n = 1066)	Doubleton (n = 1193)	Trio (n = 1143)	Pct10 (n = 1123)	Pct20 (n = 1105)
<b>ExomeDepth</b>	sensitivity	<b>0.977</b>	<b>0.873</b>	<b>0.899</b>	<b>0.963</b>	<b>0.896</b>
		1041/1066	1042/1193	1027/1143	1081/1123	990/1105
	specificity	0.974	0.951	0.993	0.982	0.996
		931/956	962/1012	943/950	998/1016	938/942
<b>EXCAVATOR</b>	sensitivity	0.374	0.372	0.368	0.391	0.381
		399/1066	444/1193	421/1143	439/1123	421/1105
	specificity	<b>0.993</b>	<b>0.982</b>	<b>0.997</b>	<b>1.0</b>	<b>1.0</b>
		289/291	320/326	311/312	298/298	292/292
<b>XHMM</b>	sensitivity	0.443	0.425	0.430	0.295	0.02
		472/1066	507/1193	492/1143	331/1123	23/1105
	specificity	0.930	0.937	0.936	0.894	0.351
		413/444	417/445	409/437	270/302	20/57
<b>CANOES</b>	sensitivity	0.90	0	0.524	0.902	0.033
		959/1066		599/1143	1013/1123	36/1105
	specificity	0.952	0	0.987	0.976	0.939
		893/938		520/527	950/973	31/33

**Table 2:** Detection of heterozygous deletions.

		Trio (n = 1143)	Pct10 (n = 1123)	Pct20 (n = 1105)	Trio_Het (n = 1143)	Pct10_Het (n = 1123)	Pct20 (n = 1105)
<b>ExomeDepth</b>	sensitivity	0.381	0.001	0	0.787	0.816	0.807
		435/1143	1/1123	0/1105	899/1143	916/1123	892/1105
	specificity	0.995	0.50	0	0.964	0.959	0.963
		412/414	1/2	0/1	832/863	861/898	839/871
<b>EXCAVATOR</b>	sensitivity	0.01	0	0	0.267	0.265	0.257
		12/1143			305/1143	298/1123	284/1105
	specificity	0.833	0	0	0.968	0.976	0.990
		10/12			210/217	201/206	197/199
<b>XHMM</b>	sensitivity	0.226	0.034	0			
		258/1143	38/1123				
	specificity	0.936	0.944	0			
		225/243	34/36				
<b>CANOES</b>	sensitivity	0	0	0			
	specificity	0	0/3	0/17			

**Table 3:** Better sensitivity and specificity for calling heterozygous deletions by selecting normal samples as controls.

## Summary

### For homozygous deletions:

- ExomeDepth has the best sensitivity and higher specificity in general; EXCAVATOR has the best specificity but not high sensitivity.
- XHMM have high specificity only for less common deletions.
- CANOES can have high sensitivity and specificity for rare deletions, but may not be stable.

### For heterozygous deletions:

- All the methods have lower sensitivity for detecting heterozygous deletions, while ExomeDepth generated the best results.
- High specificity was only observed for rare heterozygous deletions.

### Discussions:

- ExomeDepth produced the best performance overall in this simulated dataset but we are aware the simulation setting is in favor of ExomeDepth, as there are a lot of small deletions and EXCAVATOR is designed to detect larger CNVs. In addition, EXCAVATOR can differentiate between heterozygous and homozygous deletions. And we added 250 real WES samples in order to run XHMM, which may introduce noise to its CNV calling.
- XHMM and CANOES are more affected by the frequency of deletions in the dataset than ExomeDepth and EXCAVATOR.
- **Reference selection:** Using normal samples as reference significantly increase the sensitivity and specificity calling heterozygous deletions, but not for homozygous deletions.
- **Still need to check how the methods affected by CNV sizes.**