

# Computational identification of within individual contamination for more sensitive somatic mutation profiling

**Sangwoo Kim**, Kyowon Jeong, Kunal Bhutani, Hayan Lee, and Vineet  
Bafna

Department of Computer Science and Engineering, UCSD

[sak042@cs.ucsd.edu](mailto:sak042@cs.ucsd.edu), [vbafna@cs.ucsd.edu](mailto:vbafna@cs.ucsd.edu)

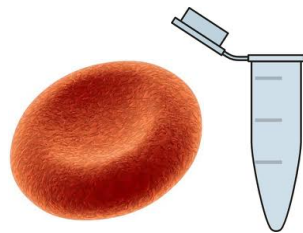
# Contents

- Background
  - somatic mutation calling
  - **The Problem:** How to consider within individual contamination in somatic mutation calling
- Virmid:
  - **The solution:**
    - Estimate of sample heterogeneity ( $\alpha$ )
    - Use  $\alpha$  to call variants
  - Test & Validation
- Application:
  - Test on HME data set
  - New SNPs and their functions
- Summary and Conclusions

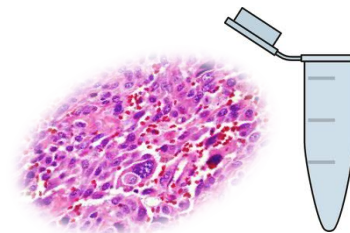
Variant calling and Within individual contamination

# **BACKGROUND**

# finding somatic mutations



sample from  
non-disease site



sample from  
disease site

reference sequence (e.g. hg19)



TGATGATT
TGATGATT
TGATGATT
TGATGATT
TGATGATT
TGATGATT
TGATGATT
TGATGATT

ACTCCATG
ACGCCATG
ACTCCCTG
ACGCCCTG
ACTCCATG
ACGCCATG
ACTCCCTG
ACGCCCTG

- **UnifiedGenotyper**
- **VarScan**
- **SomaticSniper**
- ...

# Using a mixed model

**BIOINFORMATICS ORIGINAL PAPER** Vol. 28 no. 7 2012, pages 907–913  
doi:10.1093/bioinformatics/bts053

Sequence analysis Advance Access publication January 27, 2012

**JointSNVMix: a probabilistic model for accurate detection of somatic mutations in normal/tumour paired next-generation sequencing data**

Andrew Roth<sup>1</sup>, Jiarui Ding<sup>1,3</sup>, Ryan Morin<sup>2</sup>, Anamaria Crisan<sup>1</sup>, Gavin Ha<sup>1</sup>, Ryan Giuliany<sup>1</sup>, Ali Bashashati<sup>1</sup>, Martin Hirst<sup>2</sup>, Gulisa Turashvili<sup>1</sup>, Arusha Oloumi<sup>1</sup>, Marco A. Marra<sup>2</sup>, Samuel Aparicio<sup>1,4</sup> and Sohrab P. Shah<sup>1,3,4,\*</sup>

<sup>1</sup>Department of Molecular Oncology, BC Cancer Agency, <sup>2</sup>Canada's Michael Smith Genome Sciences Centre, <sup>3</sup>Department of Computer Science and <sup>4</sup>Department of Pathology, University of British Columbia, Vancouver, BC, Canada

Associate Editor: Alex Bateman

**ABSTRACT**

**Motivation:** Identification of somatic single nucleotide variants (SNVs) in tumour genomes is a necessary step in defining the mutational landscape of cancer. From a clinical perspective, seen a number of studies exploring the mutational landscapes of various cancer subtypes. NGS investigations into prostate (Berger *et al.*, 2011), breast (Ding *et al.*, 2010; Shah *et al.*, 2009a), ovarian (Liang *et al.*, 2010; Shah *et al.*, 2009b; Wenzel *et al.*, 2010).

**BIOINFORMATICS ORIGINAL PAPER** Vol. 28 no. 14 2012, pages 1811–1817  
doi:10.1093/bioinformatics/bts271

Genome analysis Advance Access publication May 10, 2012

**Strelka: accurate somatic small-variant calling from sequenced tumor-normal sample pairs**

Christopher T. Saunders<sup>1,\*</sup>, Wendy S. Wong<sup>2</sup>, Sajani Swamy<sup>1</sup>, Jennifer Becq<sup>2</sup>, Lisa J. Murray<sup>2</sup> and R. Keira Cheetham<sup>2</sup>

<sup>1</sup>Illumina, Inc., 5200 Illumina Way, San Diego, CA 92122, USA and <sup>2</sup>Illumina Cambridge Ltd., Chesterford Research Park, Little Chesterford, Essex CB10 1XL, UK

Associate Editor: Michael Brudno

**ABSTRACT**

**Motivation:** Whole genome and exome sequencing of matched tumor-normal sample pairs is becoming routine in cancer research. The consequent increased demand for somatic variant analysis of paired samples requires methods specialized to model this problem calling methods, particularly for SNVs and small indels where the number of somatic variants can easily overwhelm manual review. An additional challenge for somatic variant calling on matched tumor-normal samples is robust handling of impurity and copy-number variation in the tumor sample, ideally without requiring external

normal  $p(AA), p(AB), p(BB)$   
disease  $p(AA), p(AB), p(BB)$

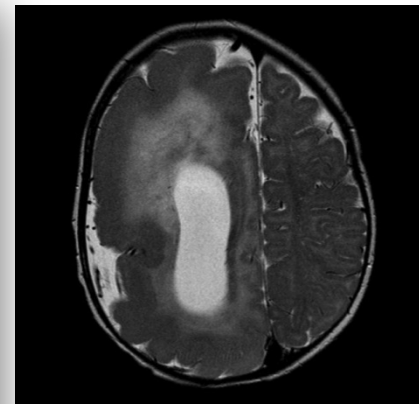


disease  
■ AA&AB&BB  
normal  
■ ■ A@1 &p↓12 &p↓13 @  
AB@B&p↓22 &p↓23 @p↓31  
&p↓32 &p↓33 )

**G: joint genotype probability matrix**

## *De novo* somatic mutations in components of the PI3K-AKT3-mTOR pathway cause hemimegalencephaly

Jeong Ho Lee<sup>1,2,9</sup>, My Huynh<sup>3,4</sup>, Jennifer L Silhavy<sup>1,2</sup>, Sangwoo Kim<sup>5</sup>, Tracy Dixon-Salazar<sup>1,2</sup>, Andrew Heiberg<sup>1,2</sup>, Eric Scott<sup>1,2</sup>, Vineet Bafna<sup>5</sup>, Kiley J Hill<sup>1,2</sup>, Adrienne Collazo<sup>1,2</sup>, Vincent Funari<sup>6,7</sup>, Carsten Russ<sup>8</sup>, Stacey B Gabriel<sup>8</sup>, Gary W Mathern<sup>3,4,10</sup> & Joseph G Gleeson<sup>1,2,10</sup>



**Focal cortical dysplasia**

## Universal noninvasive detection of solid organ transplant rejection

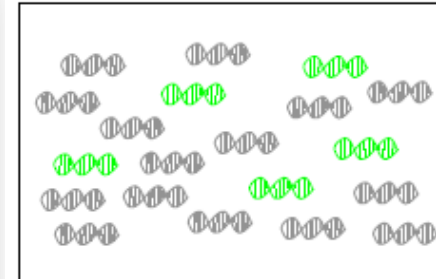
Thomas M. Snyder<sup>a,b</sup>, Kiran K. Khush<sup>c</sup>, Hannah A. Valentine<sup>c,1</sup>, and Stephen R. Quake<sup>a,b,1</sup>

<sup>a</sup>The Howard Hughes Medical Institute and <sup>b</sup>Departments of Applied Physics and Bioengineering, Stanford University, Stanford, CA 94305; and <sup>c</sup>Division of Cardiovascular Medicine, Stanford University School of Medicine, Stanford, CA 94305

Edited\* by Leonard A. Herzenberg, Stanford University, Stanford, CA, and approved February 24, 2011 (received for review September 15, 2010)

It is challenging to monitor the health of transplanted organs, particularly with respect to rejection by the host immune system. Because transplanted organs have genomes that are distinct from the recipient's genome, we used high throughput shotgun sequencing to develop a universal noninvasive approach to monitoring organ health. We analyzed cell-free DNA circulating in the blood of heart

donor-specific chromosome Y has been detected in recipient urine and plasma (16–19). To date, most measurements of cell-free DNA in organ transplantation have been limited to the special case of women who receive male organs, which has prevented the widespread use of cell-free DNA as a diagnostic tool, because female recipients of male donor organs represent less than a quarter of all transplant procedures. HLA markers can be quantified to identify



**Cell free DNA in recipients' blood**

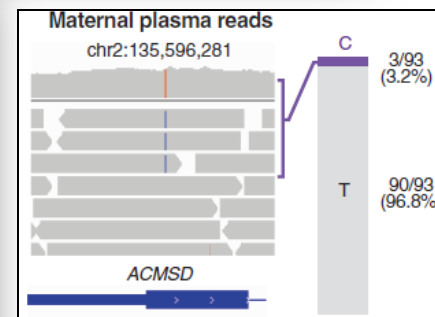
### RESEARCH ARTICLE

#### GENOMICS

## Noninvasive Whole-Genome Sequencing of a Human Fetus

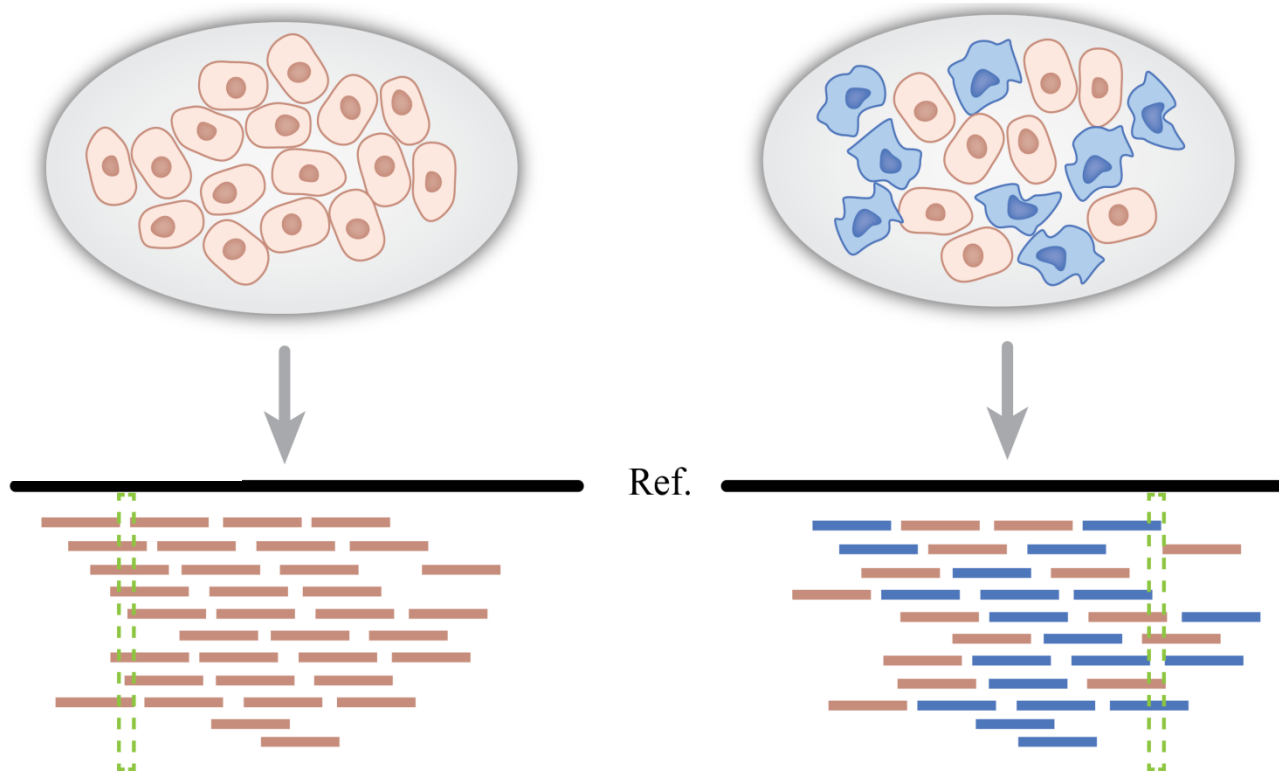
Jacob O. Kitzman,<sup>1\*</sup> Matthew W. Snyder,<sup>1</sup> Mario Ventura,<sup>1,2</sup> Alexandra P. Lewis,<sup>1</sup> Ruolan Qiu,<sup>1</sup> LaVone E. Simmons,<sup>3</sup> Hilary S. Gammill,<sup>3,4</sup> Craig E. Rubens,<sup>5,6</sup> Donna A. Santillan,<sup>7</sup> Jeffrey C. Murray,<sup>8</sup> Holly K. Tabor,<sup>5,9</sup> Michael J. Bamshad,<sup>1,5</sup> Evan E. Eichler,<sup>1,10</sup> Jay Shendure<sup>1\*</sup>

Analysis of cell-free fetal DNA in maternal plasma holds promise for the development of noninvasive prenatal genetic diagnostics. Previous studies have been restricted to detection of fetal trisomies, to specific paternally inherited mutations, or to genotyping common polymorphisms using material obtained invasively, for example, through chorionic villus sampling. Here, we combine genome sequencing of two parents, genome-wide maternal haplotyping, and deep sequencing of maternal plasma DNA to noninvasively determine the genome sequence of a human fetus at 18.5 weeks of gestation. Inheritance was predicted at  $2.8 \times 10^6$  parental heterozygous sites



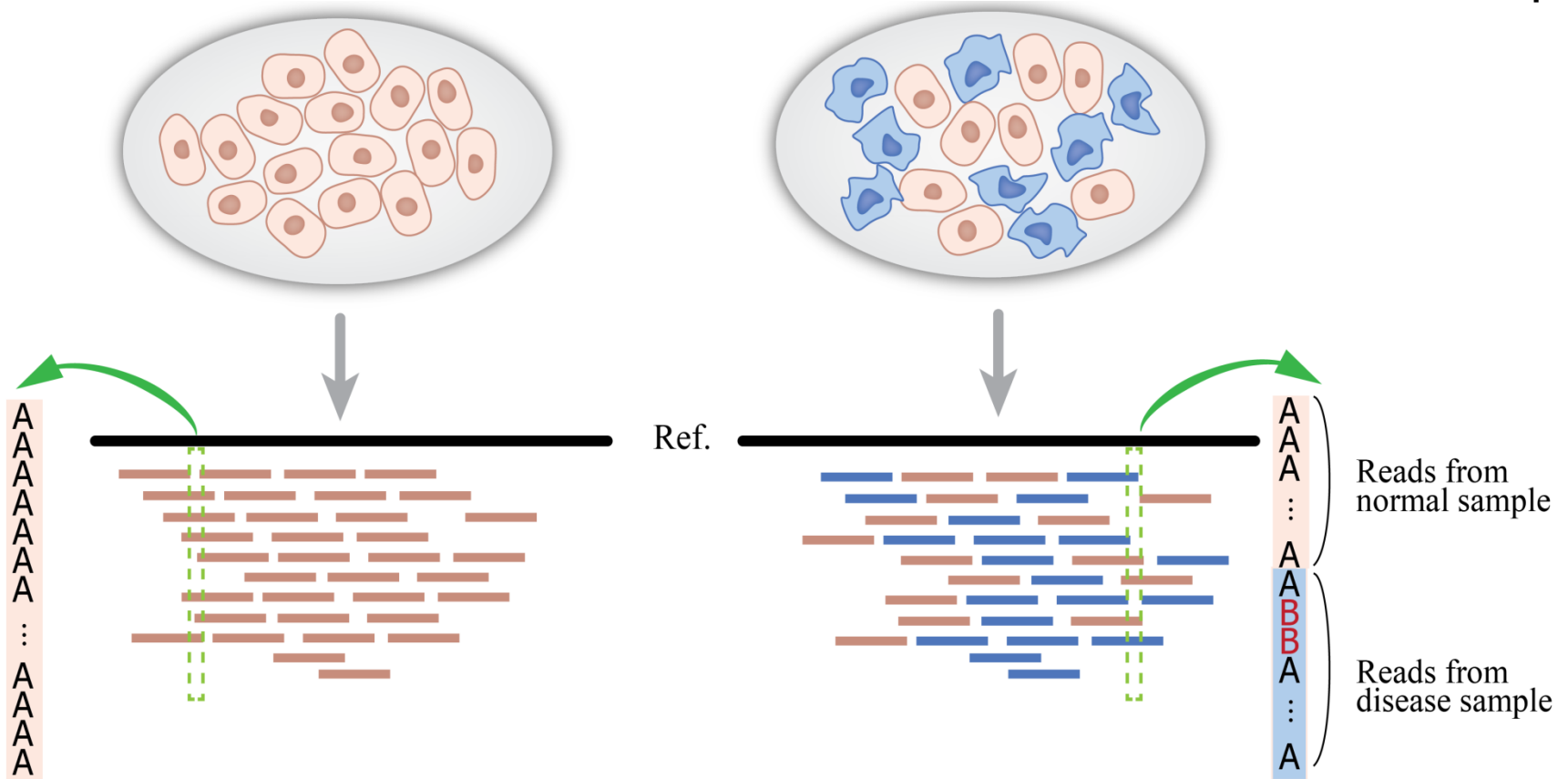
**Fetus genome in maternal blood**

# Sample heterogeneity



# Loss of somatic mutation calling

$\alpha$  proportion of **control** sample in the mixed disease sample

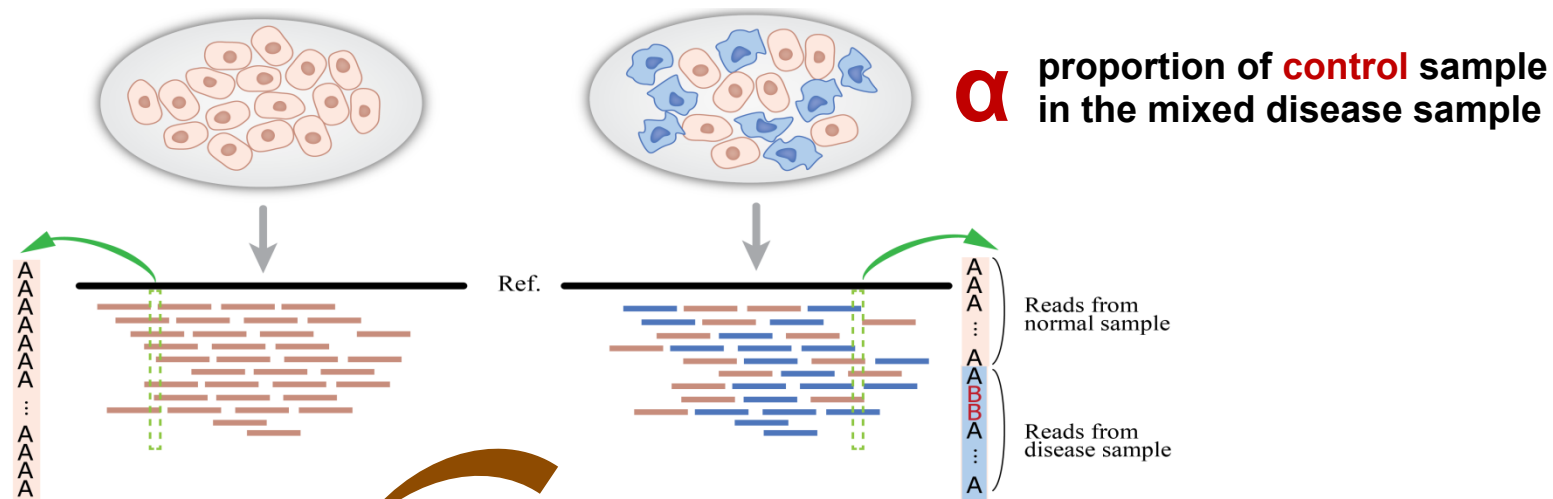


$$\begin{aligned} \alpha=0 & \quad g(AA,AB) \rightarrow \dots ABBBABBAAAB \dots \rightarrow \sim 50\% \text{ BAF} \\ \alpha=0.5 & \quad g(AA,AB) \rightarrow \dots AAAAAABBAAAB \dots \rightarrow \sim 25\% \text{ BAF} \end{aligned}$$

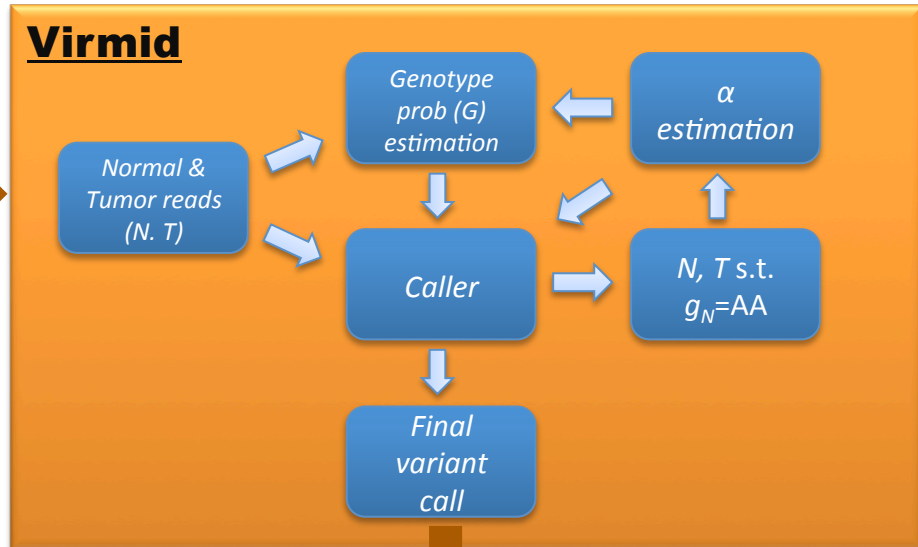


Virtual Microdissection for mixed disease sample

**VIRMID**

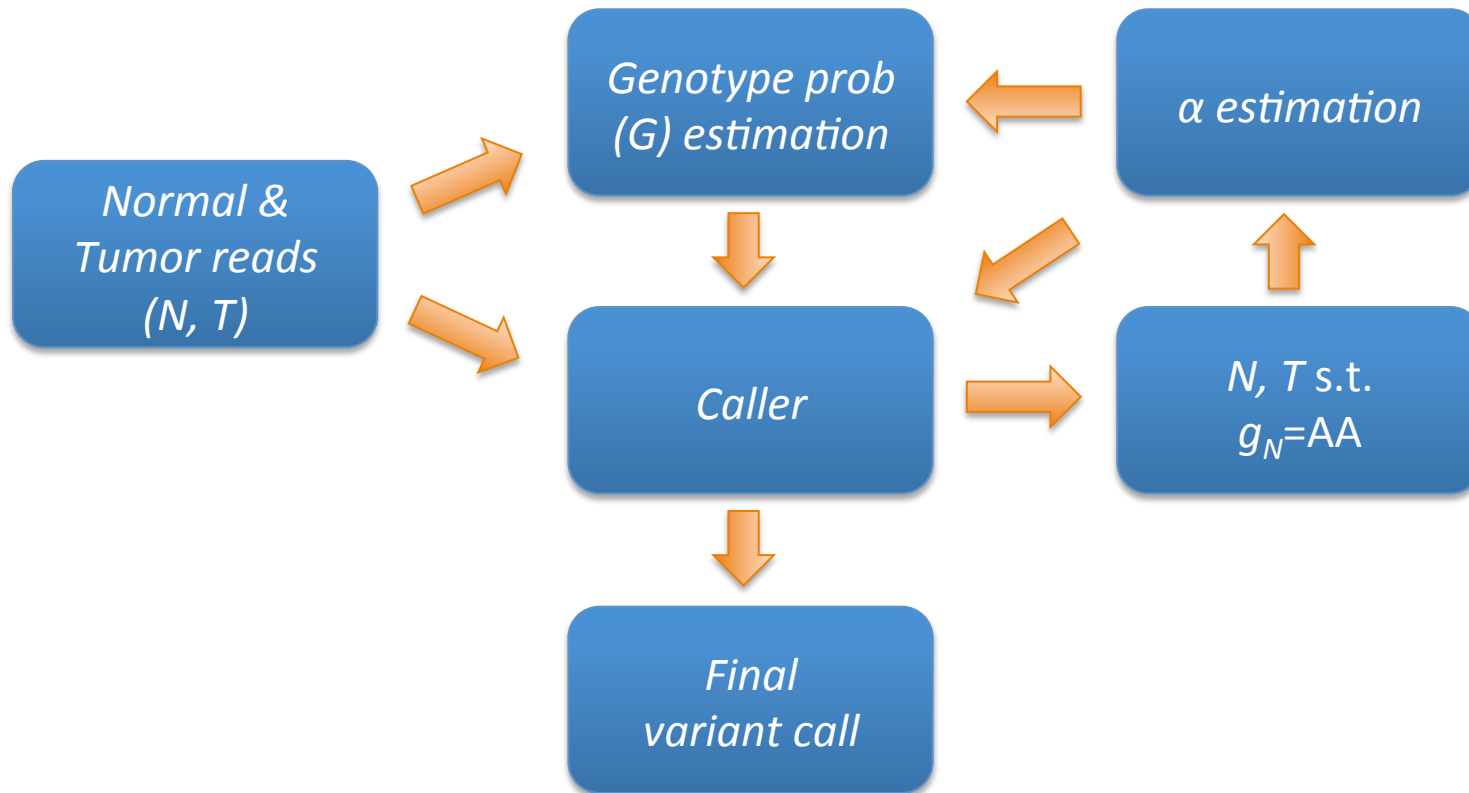


1. **read alignment of disease & normal**
2. **reference sequence**



1. **estimated  $\alpha$**
2. **a list of somatic (germline) mutations**

# Virmid Model



# Likelihood function

- Likelihood function

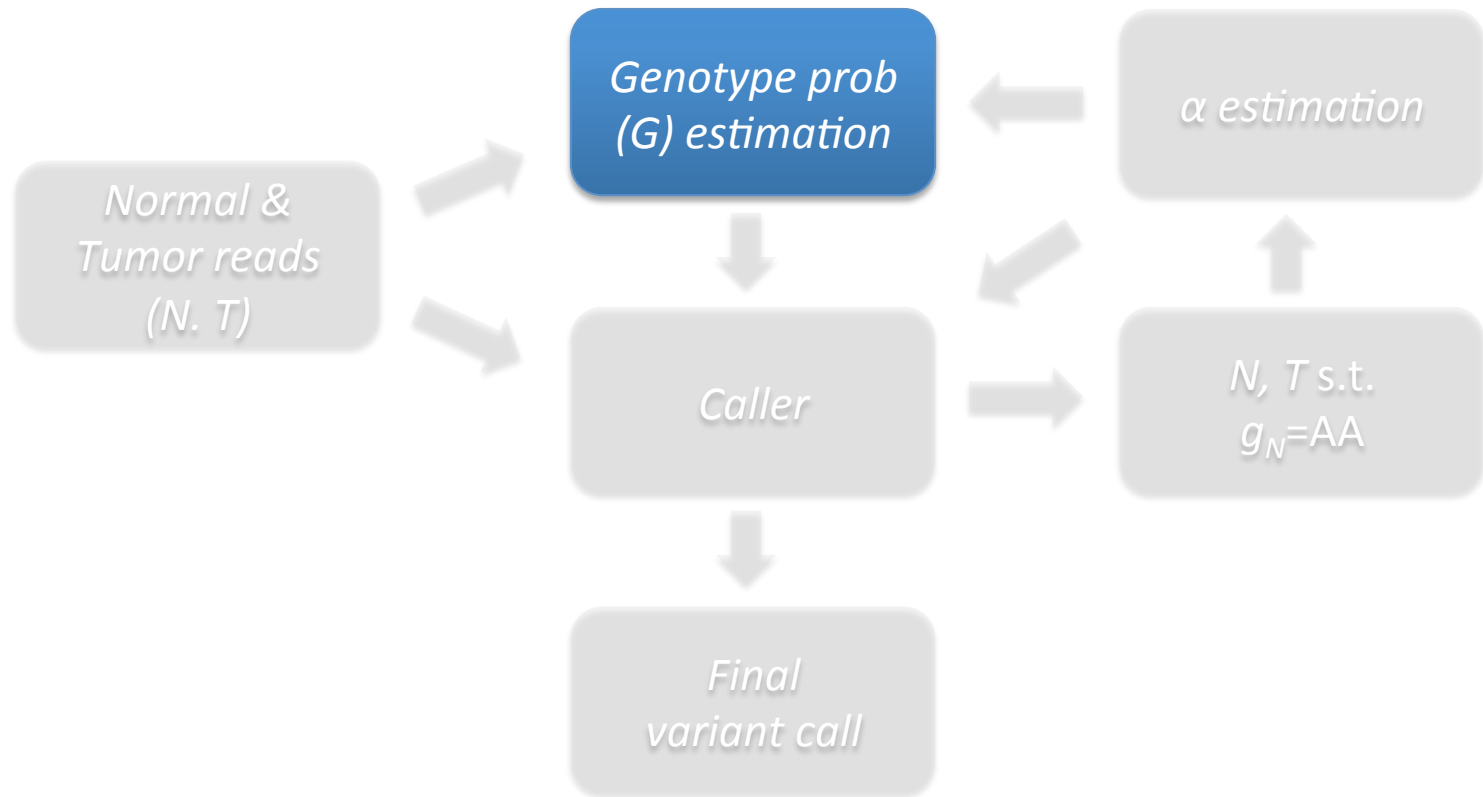
$$L(\alpha, G | N, T) = P_{\alpha, G}(N, T)$$

$$= \prod_i \sum_{g_N, g_T} P_{\alpha, G}(g_N, g_T) \cdot \prod_j P_{\alpha, G}(N_j^i | g_N) \cdot \prod_k P_{\alpha, G}(T_k^i | g_N, g_T)$$

Joint genotype probability  
 Probability that we observe the normal reads given the genotype of normal sample  
 Probability that we observe the tumor reads given the genotype of normal/tumor sample

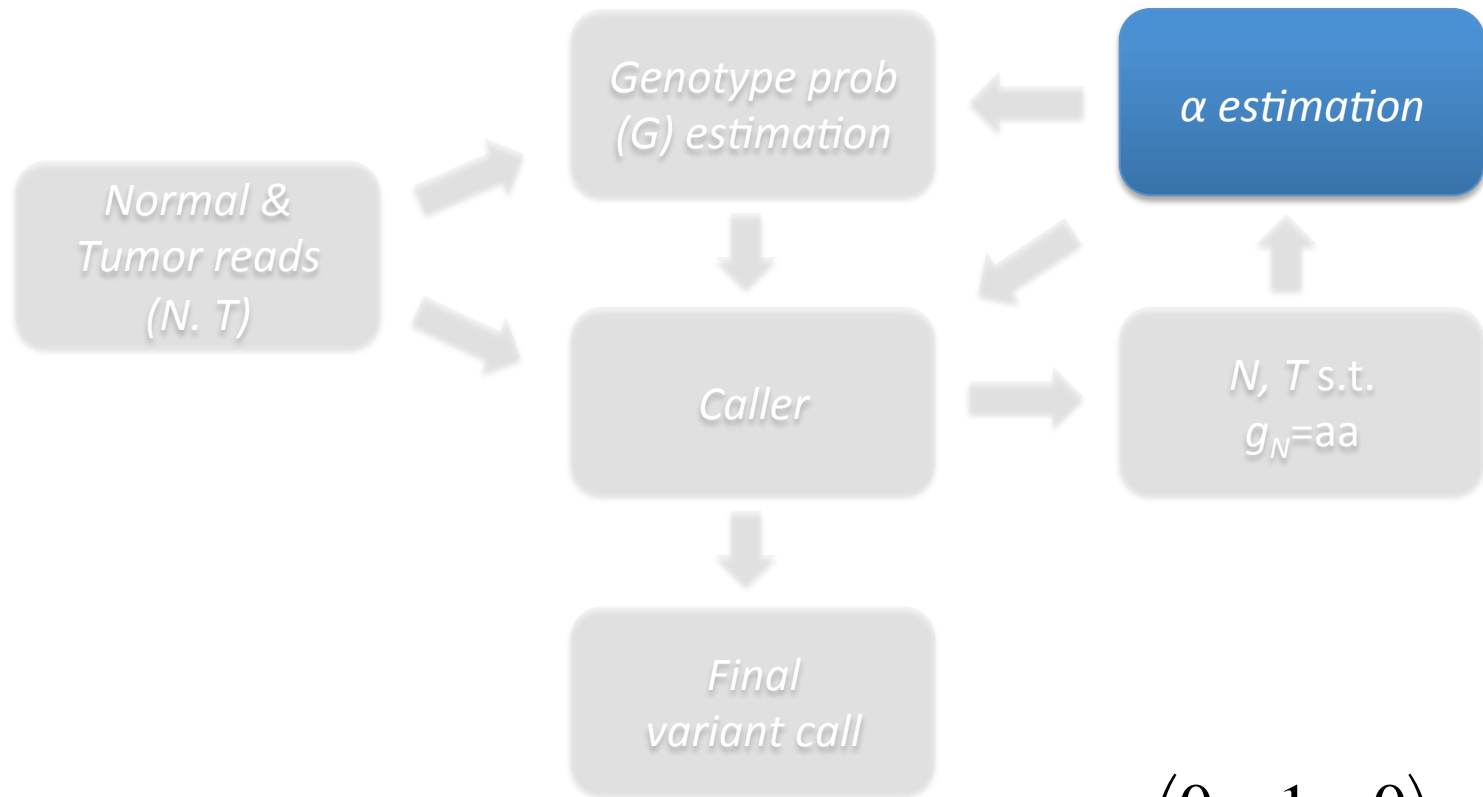
- Assumptions
  - Reads at different positions are independent
  - If genotypes are given, reads in a positions are independent
- Considers
  - Read error rate
  - Mapping error rate
  - Other biases

# Estimation of genotype probabilities



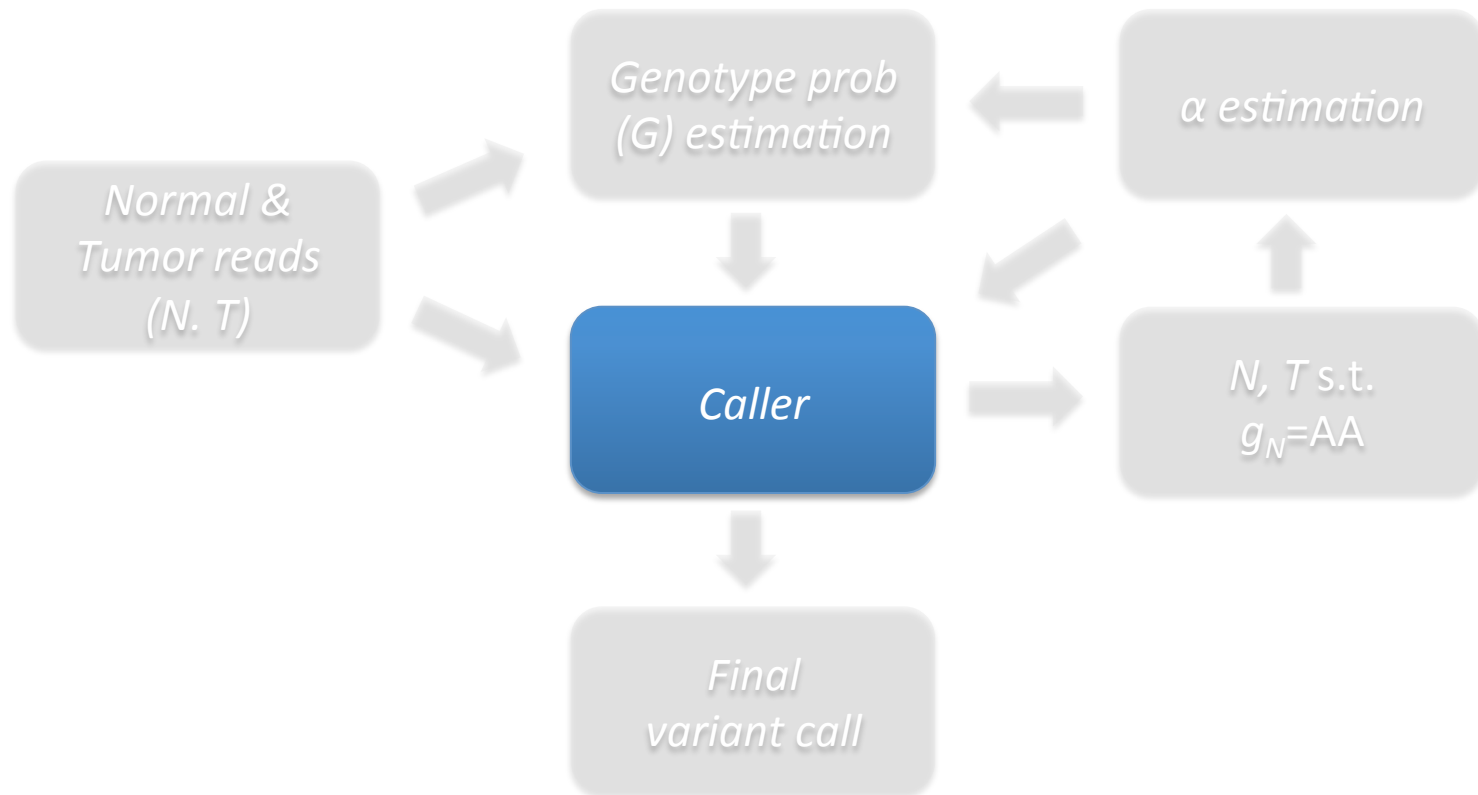
$$\hat{G} = \arg \max_G L(\hat{\alpha}, G \mid N, T)$$

# Estimation of $\alpha$



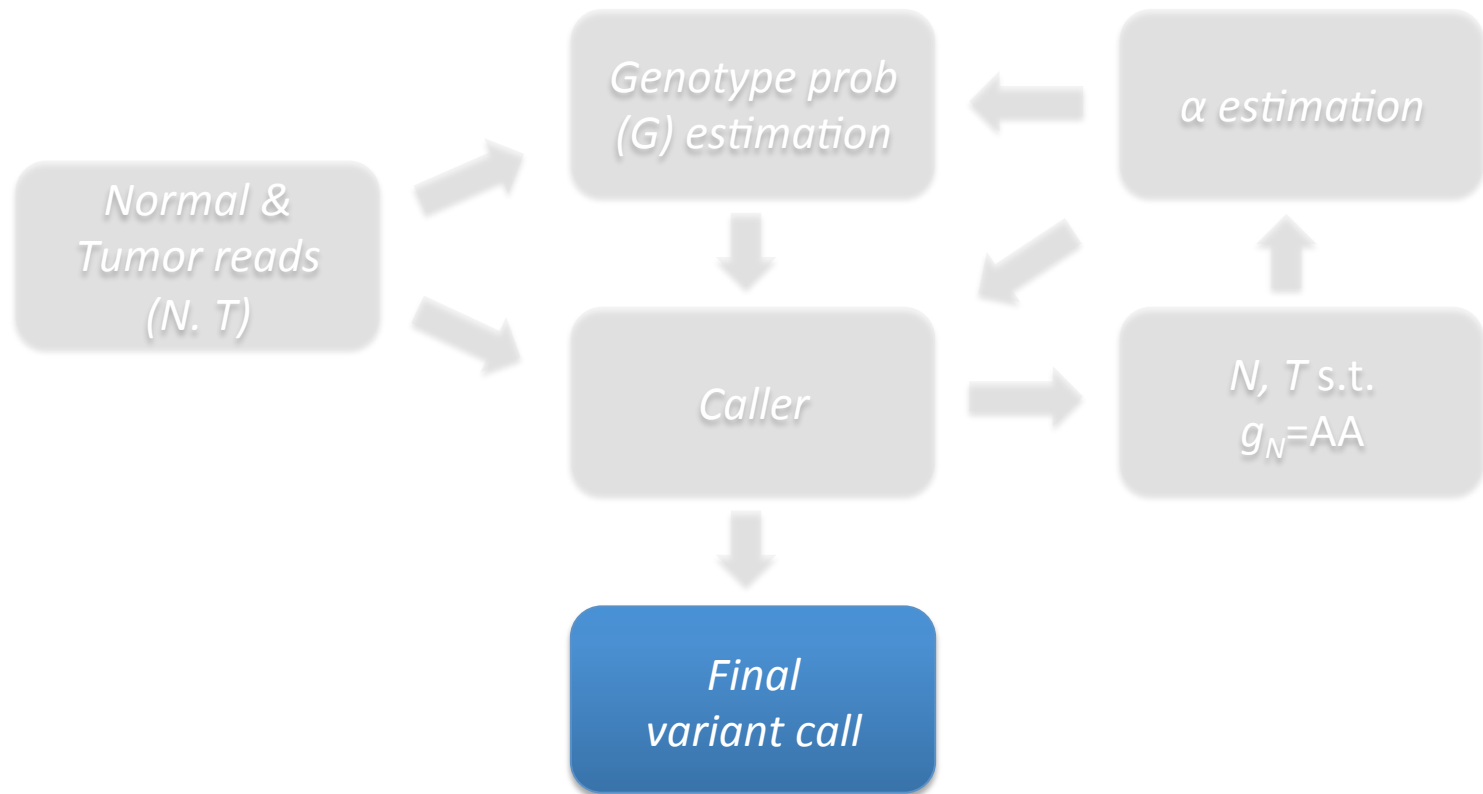
$$\hat{\alpha} = \arg \max_{\alpha} L(\alpha, \bar{G} \mid N, T) \quad \text{where} \quad \bar{G} \approx \begin{pmatrix} 0 & 1 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{pmatrix}$$

# Variant Caller



- Use of MAP (Maximum a posteriori probability estimate)
- Use estimated  $G$  as a priori distribution
- Estimate allele frequencies

# Final Variant Call



- Output final variants after  $\alpha$  converged
- Call if the probability of being somatic mutation is higher than not



# Validation

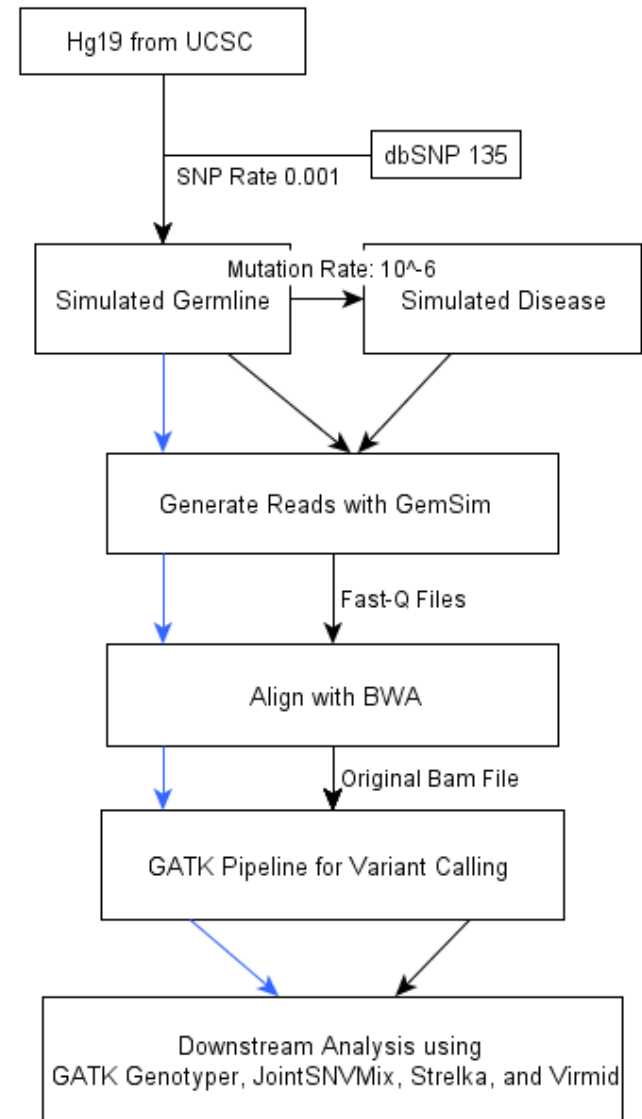
**simulated mutation & mixture in hg19 Chr 1**

**tested  $\alpha$  values:**

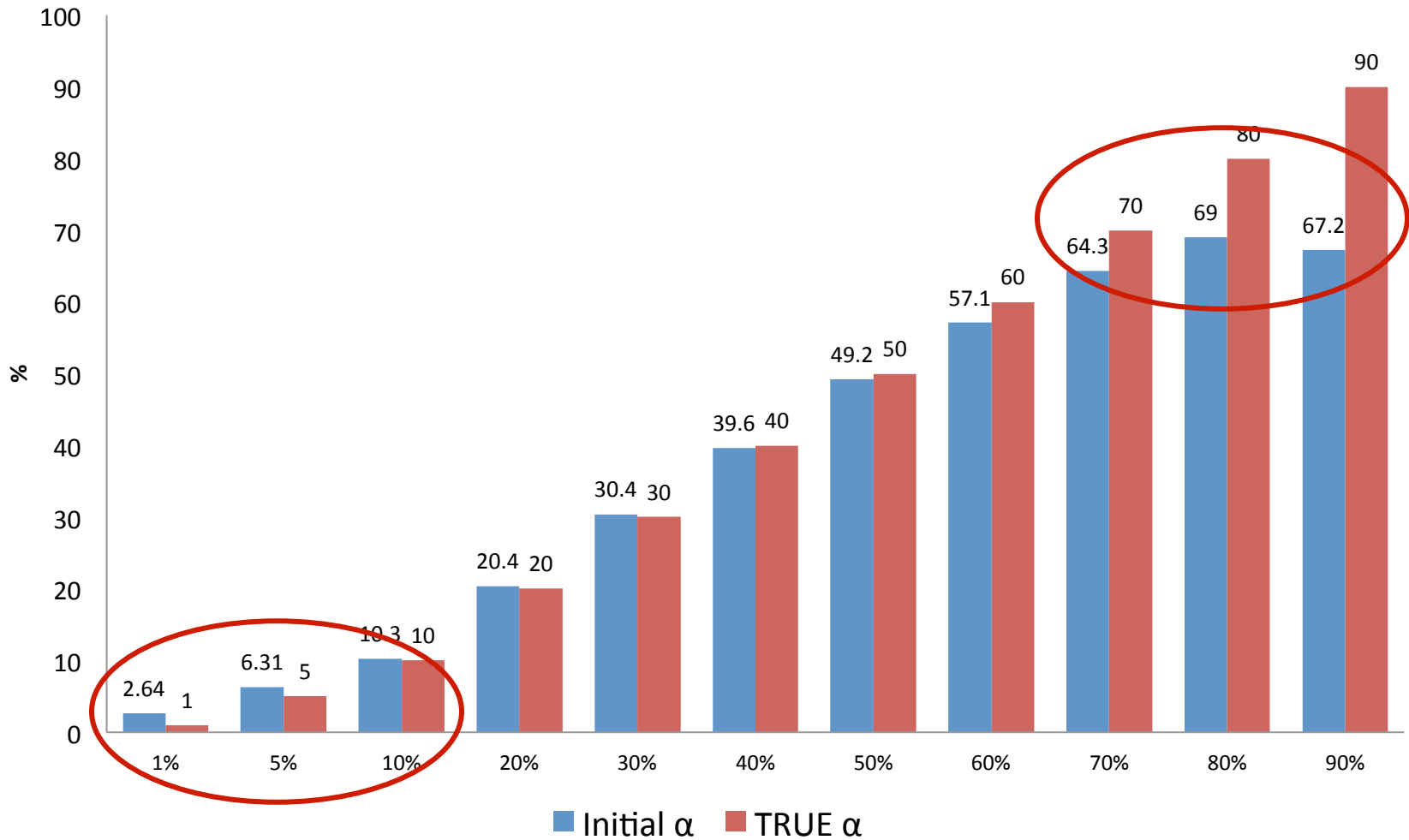
**1%, 5%, 10%, 20%, 30%, 40%,  
50%, 60%, 70%, 80%, 90%**

**# true somatic (germline) mutations:**

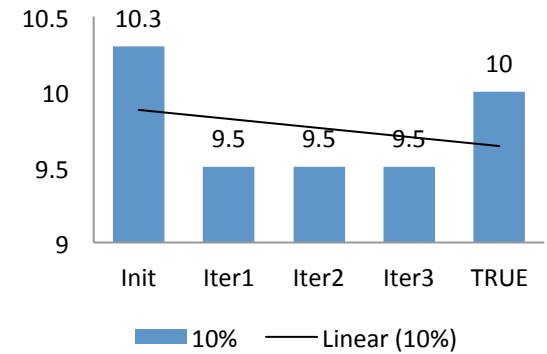
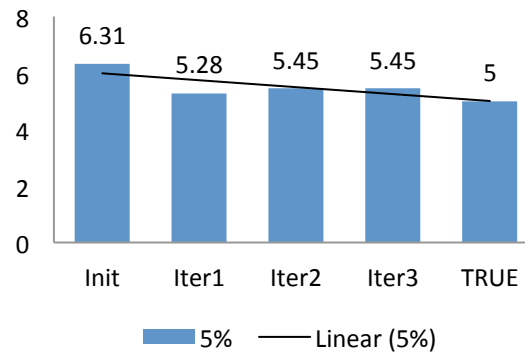
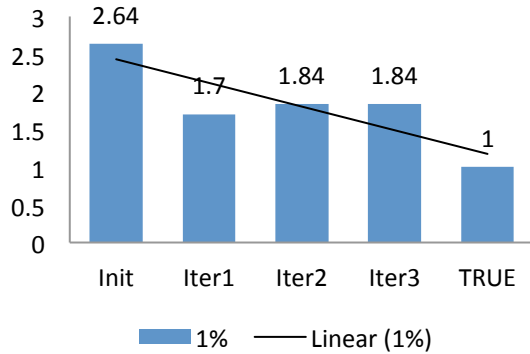
**2226 (228,018)**



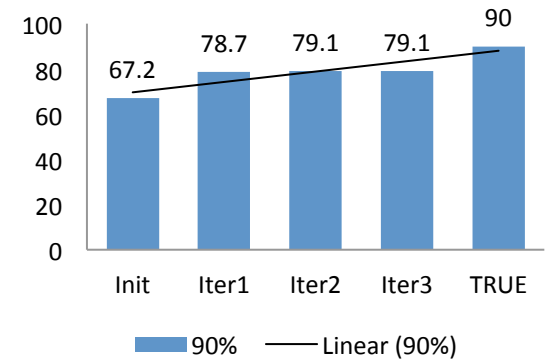
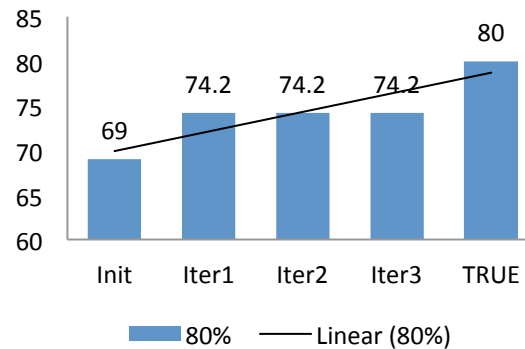
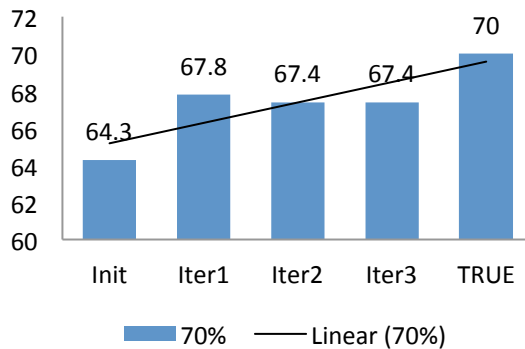
# Estimation of $\alpha$



# Reducing biases

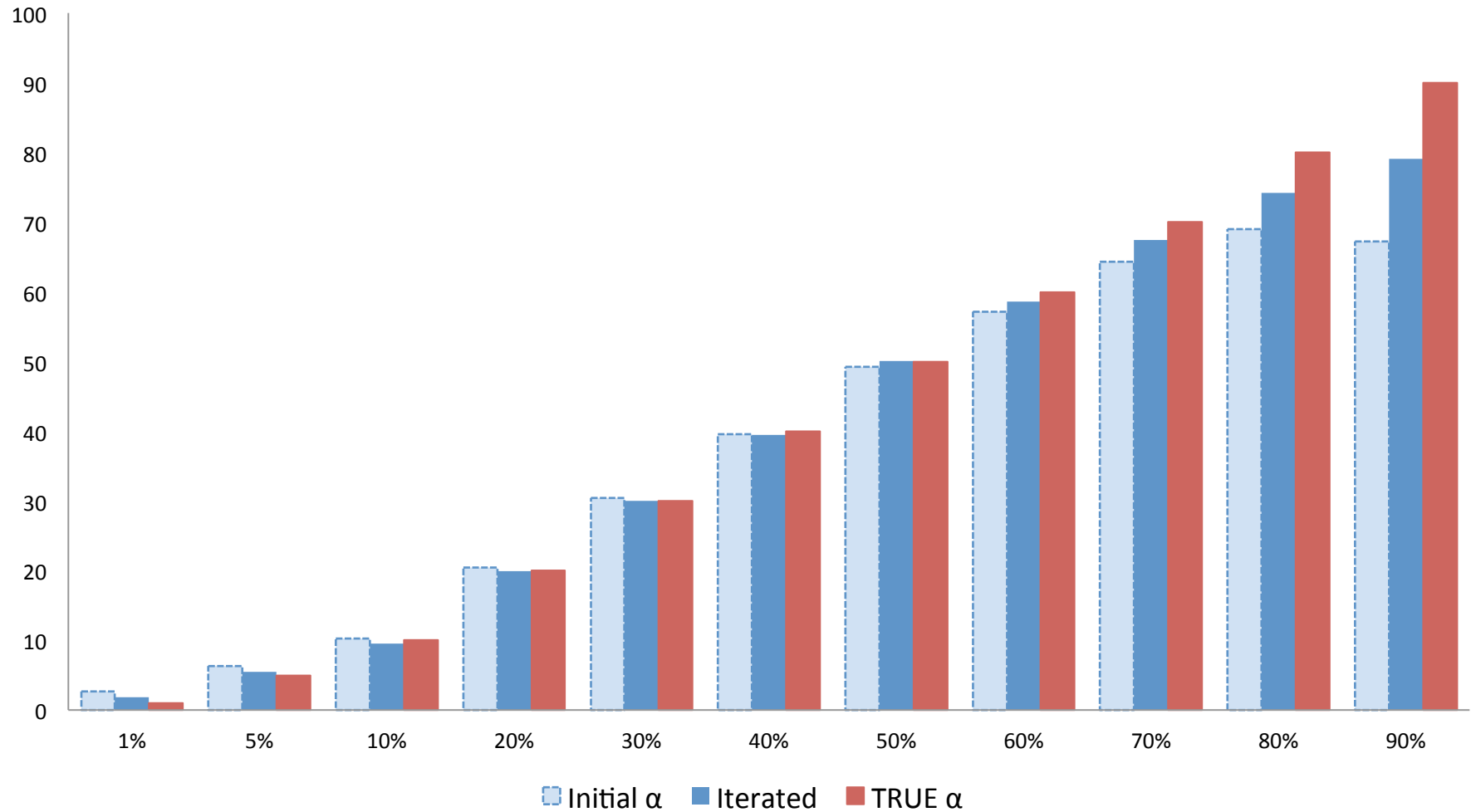


## *Changes in low $\alpha$*



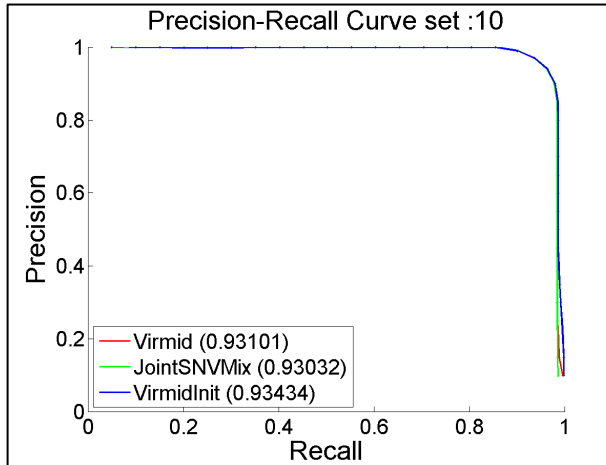
## *Changes in high $\alpha$*

# $\alpha$ after iteration

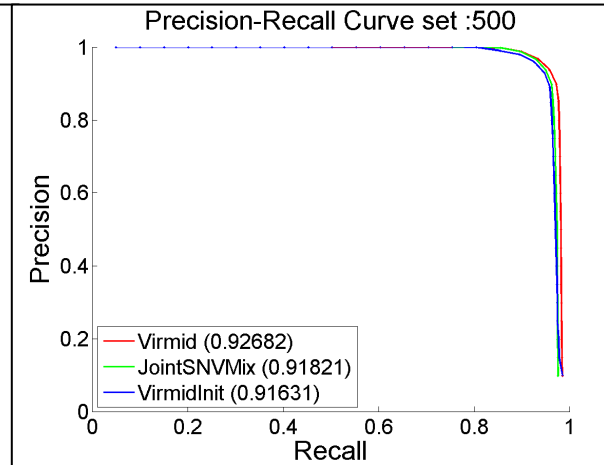


# PR (Precision-Recall) Curves

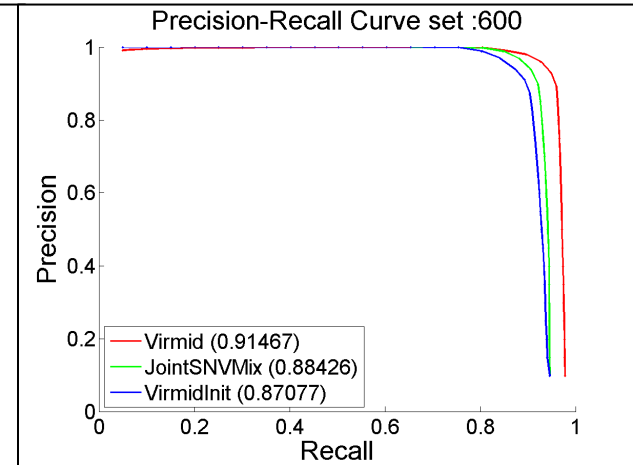
$\alpha=0.1\%$



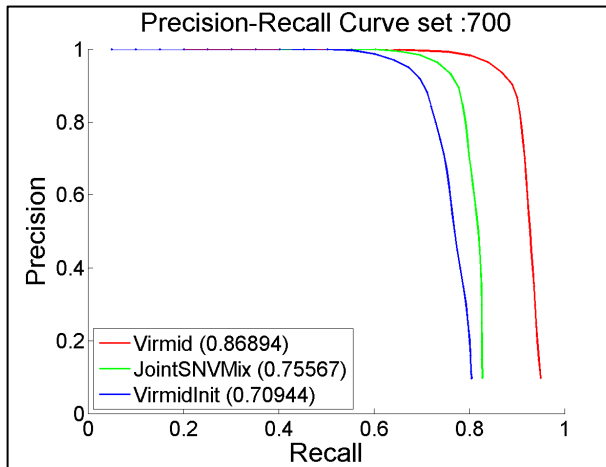
$\alpha=50\%$



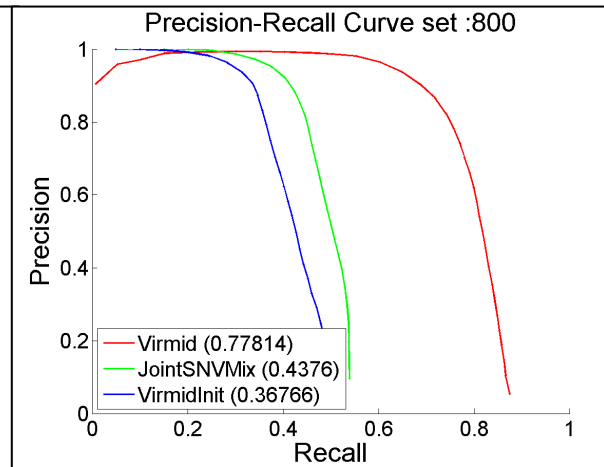
$\alpha=60\%$



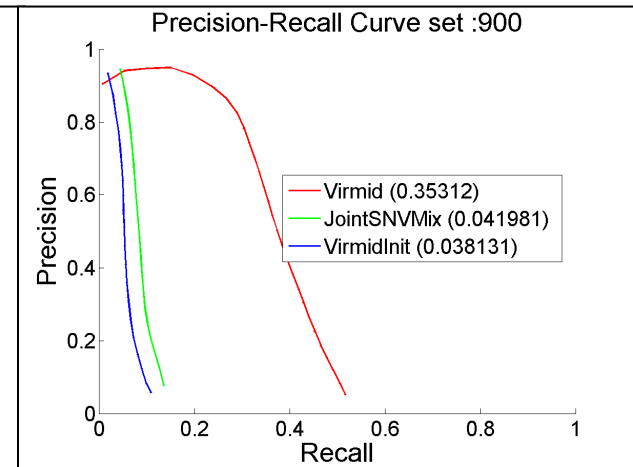
$\alpha=70\%$



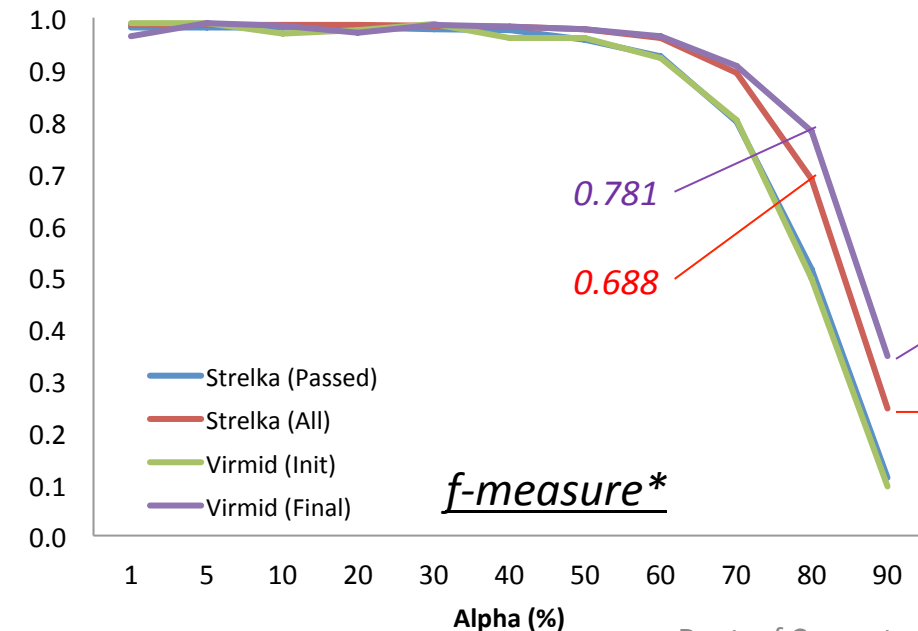
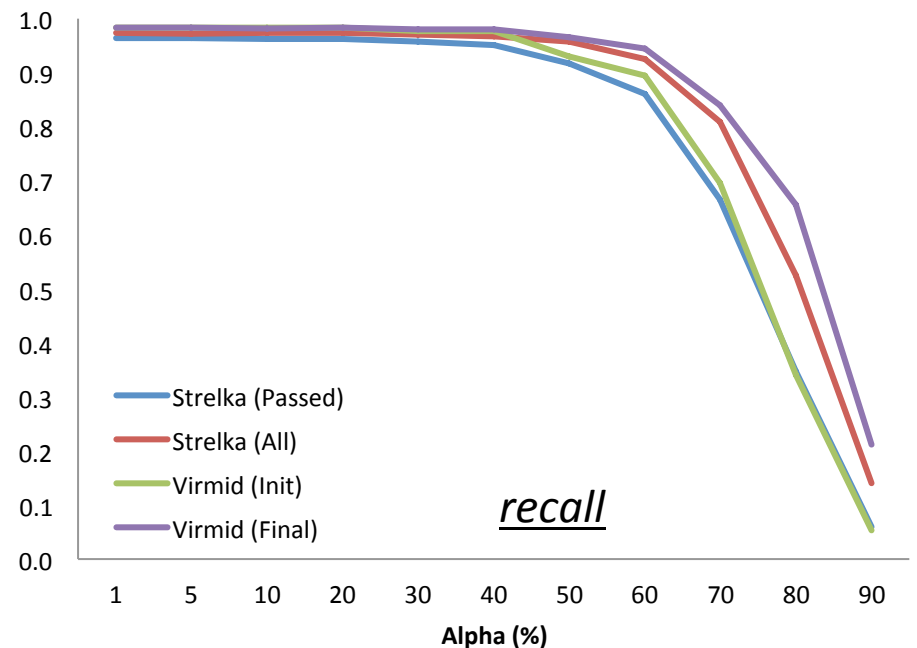
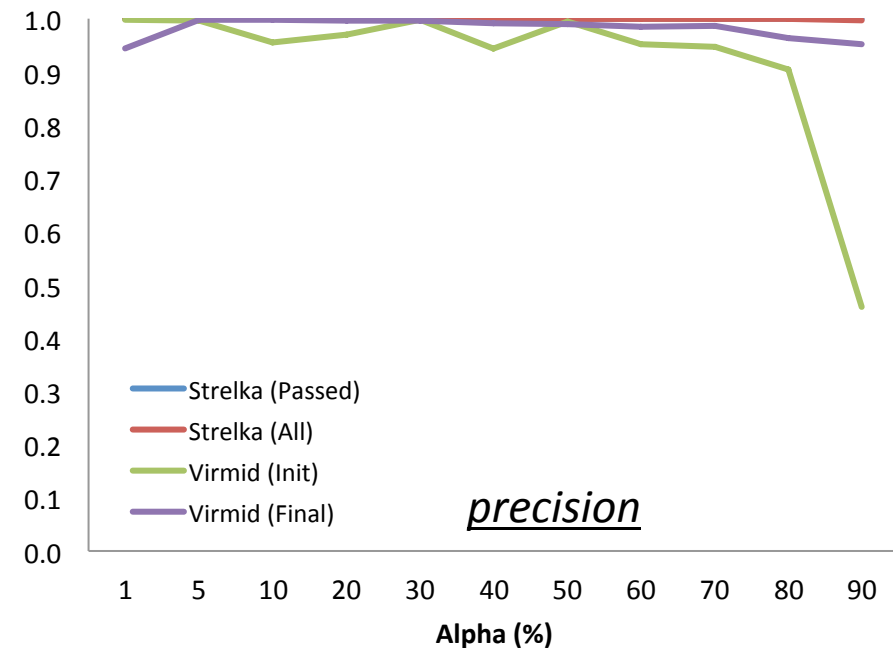
$\alpha=80\%$



$\alpha=90\%$



**precision** =  $\frac{\#right}{\#predicted} = \frac{\#TP}{\#TP + \#FP}$       **recall** =  $\frac{\#right}{\#true} = \frac{\#TP}{\#TP + \#FN}$



for each  $\alpha$ , **precision**, **recall**, **f-score\*** are calculated only from “final” outputs

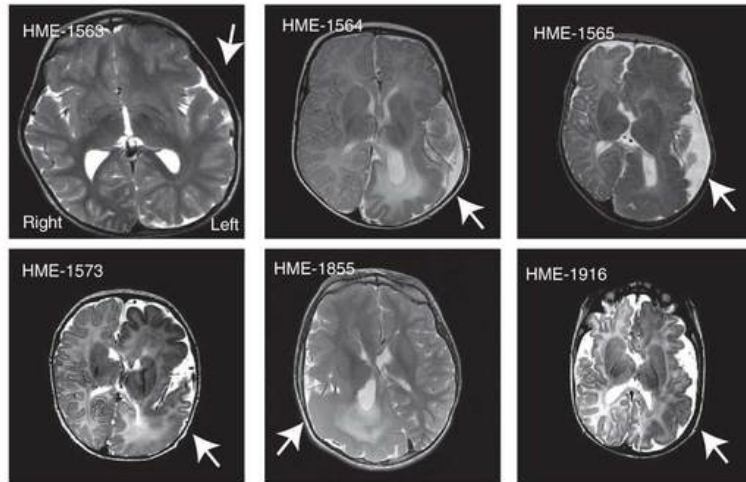
$$*f = \frac{2 \times \text{precision} \times \text{recall}}{\text{precision} + \text{recall}}$$

HME data

# APPLICATION

# HME-data

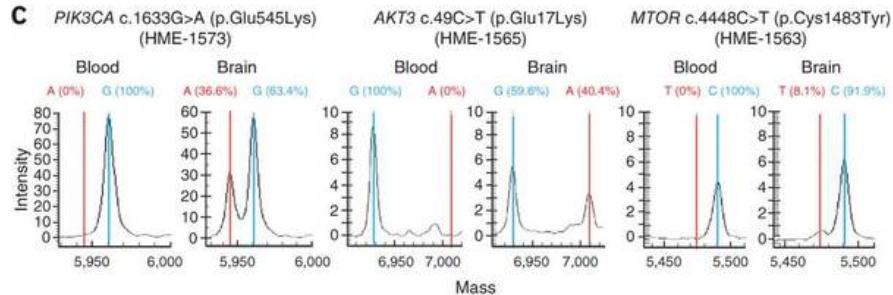
a



b

	Blood		Brain		
	Mut counts	Ref counts	Mut counts	Ref counts	Mut allele (%)
<i>PIK3CA</i> c.1633G>A (HME-1573)	0	121	9	47	16%
<i>AKT3</i> c.49C>T (HME-1565)	0	49	9	23	28%
<i>MTOR</i> c.4448C>T (HME-1563)	0	298	17	159	9.7%

c



Sample	Previous Mut.	Mut. site
HME-1563	mTOR	Cys1483Tyr
HME-1565	AKT3	Glu17Lys
HME-1573	PIK3CA	Glu545Lys
HME-1574	-	-
HME-1620	-	-

*SNPs called by JointSNVMix2*



# Variant calling with Virmid

Sample	Est. Alpha	Previous Calls	Virmid Calls	Misc.
HME-1563	67.1%	1547	1431	<i>mTOR in overlap</i>
HME-1565	66.6%	1328	1761	<i>AKT3 in overlap</i>
HME-1573	62.8%	1386	940	<i>PIK3CA in overlap</i>
HME-1574	66.4%	1440	3285	
HME-1620	65.8%	1335	4165	

The diagram illustrates the complex signaling pathways that regulate eIF4 and p70S6K. Key components and interactions include:

- Receptors and Ligands:** Growth Factors, Cytokines, Hormones & Neuropeptides; Integrin; Stress, Heat Shock, Anisomycin.
- Key Kinases and Phosphatases:**
  - PI3K Pathway:** PI3K (circled in yellow) → PIP3 → PDK1 → Akt/PKB (circled in yellow) → FRAP/mTOR (circled in yellow). Rapamycin inhibits FRAP/mTOR.
  - MAPK Pathway:** Ras → Raf → MEK → MAPK.
  - Other Kinases:** ERKs, PKC, PP2A, MNK1, MNK2, S6 (circled in red).
- Inhibitors:** Wortmannin, Ly294002 (inhibit PI3K); Rapamycin (inhibit FRAP/mTOR); PP2A (inhibit p70S6K).
- eIF4 Complexes:**
  - eIF4E/eIF4G/eIF4A/eIF4B:** The core eIF4 complex.
  - eIF4F Complex:** Includes eIF4E, eIF4G, and eIF4A.
  - 43S Complex:** Formed by the eIF4 complex and the 43S ribosomal subunit.
  - 48S Complex:** Formed by the 43S complex and mRNA.
- Translational Regulation:**
  - Translation Off:** Inhibited by eIF4E/eIF4BP and eIF4BP.
  - Translation On:** Promoted by the 48S complex.



Sample & Assay Technologies

# CONCLUSIONS

# Conclusions

- Within individual contamination seriously affects somatic variant calling
- Virmid accurately infers the proportion of non-disease sample in a potentially mixed disease sample
- Virmid increases accuracy (precision and recall) by considering the within individual contamination
- By applying Virmid on disease samples with heterogeneity issues, we can identify more somatic variants to correlate with phenotypes

# Acknowledgements

## University of California, San Diego

Dept. of Computer Science and Engineering

**Vineet Bafna, Ph.D.**

**Kunal Bhutani**

Dept. of Electrical and Computer Engineering

**Kyowon Jeong**

Institute for Genomic Medicine, Rady Children's Hospital

**Joseph Gleeson, M.D, Ph.D.**

**Jeong Ho Lee, M.D, Ph.D. (KAIST)**

Dept. of Bioengineering

**Hojung Nam, Ph.D.**

## Stony Brook University

Dept. of Computer Science

**Hayan Lee**

**Funding:**

P01 HD070494-01 and RO1 HG004962



**Thank you**

# **SUPPLEMENTARY SLIDES**

# Estimating calls

A  
A  
A  
A

A


T

A

sequencing error = e

- True genotype = “AA” and no sequencing error
  - $P_{1-e} g=AA P(g=AA)$
- True genotype = “AT”
  - Read was generated from ‘A’ allele and no sequencing error
  - $1/2 * P_{1-e} g=AT P(g=AT)$
  - Read was generated from ‘T’ allele and sequencing error and ‘A’ was generated by chance
  - $1/2 * 1/4 * P_{e} g=AT P(g=AT)$
- True genotype = “TT” and sequencing error
  - $P_{e} g=TT P(g=TT)$

# Estimating calls (cont'd)

A  
A  
A  
  
A

sequencing error = e

- True genotype = “AA” and sequencing error
  - $P_{eg=AA} P(g=AA)$
- True genotype = “AT”
  - Read was generated from ‘A’ allele and sequencing error and ‘T’ was generated by chance
    - $1/2 * 1/4 * P_{eg=AT} P(g=AT)$
  - Read was generated from ‘T’ allele and no sequencing error
    - $1/2 * P_{1-e} P(g=AT)$
- True genotype = “TT” and no sequencing error
  - $P_{eg=TT} P(g=TT)$

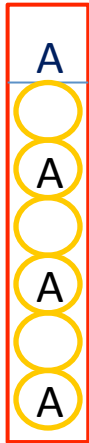
A

T

A



# Estimating calls (Cont'd)



A

T

A

$$P(AA|D) = PDAA P(AA) / PDAA P(AA) + PDAB P(AB) + PDBB P(BB)$$

$$= PDAA P(AA) / \sum_{gi \in G} PDgi P(gi)$$

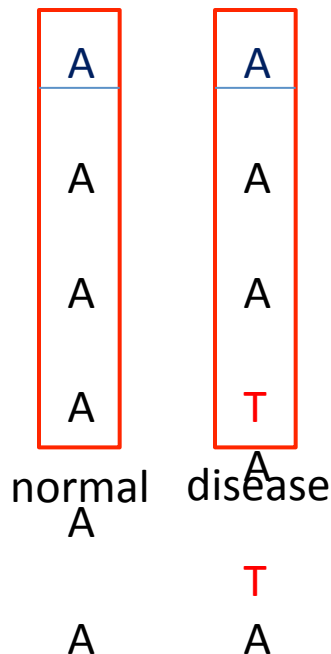
$$P(AB|D) = \dots$$

$$P(BB|D) = \dots$$

	G=AA	G=AB	G=BB
Probability	56%	40%	4%

# Calling somatic mutations

- Call variants for disease sample
  - And filter out if there's a same variant in normal sample (GATK-UnifiedGenotyper)
- Calculate a joint probability for totally 9 (3 x 3) genotypes (jointSNVMix, Strelka)



	$G_T=AA$	$G_T=AB$	$G_T=BB$
$G_N=AA$	10%	80%	2%
$G_N=AB$	0.5%	5%	1%
$G_N=BB$	0.1%	0.4%	1%

# Within individual contamination



normal



genotype = AA

A
A
A
A
A
A
A
A
A
A



disease



AA AB BB

A	A	A
A	A	B
A	A	B
A	A	B
A	A	B
A	A	B
A	B	B
A	B	B
A	B	B
A	B	B



normal



genotype = AA

A
A
A
A
A
A
A
A
A
A



disease

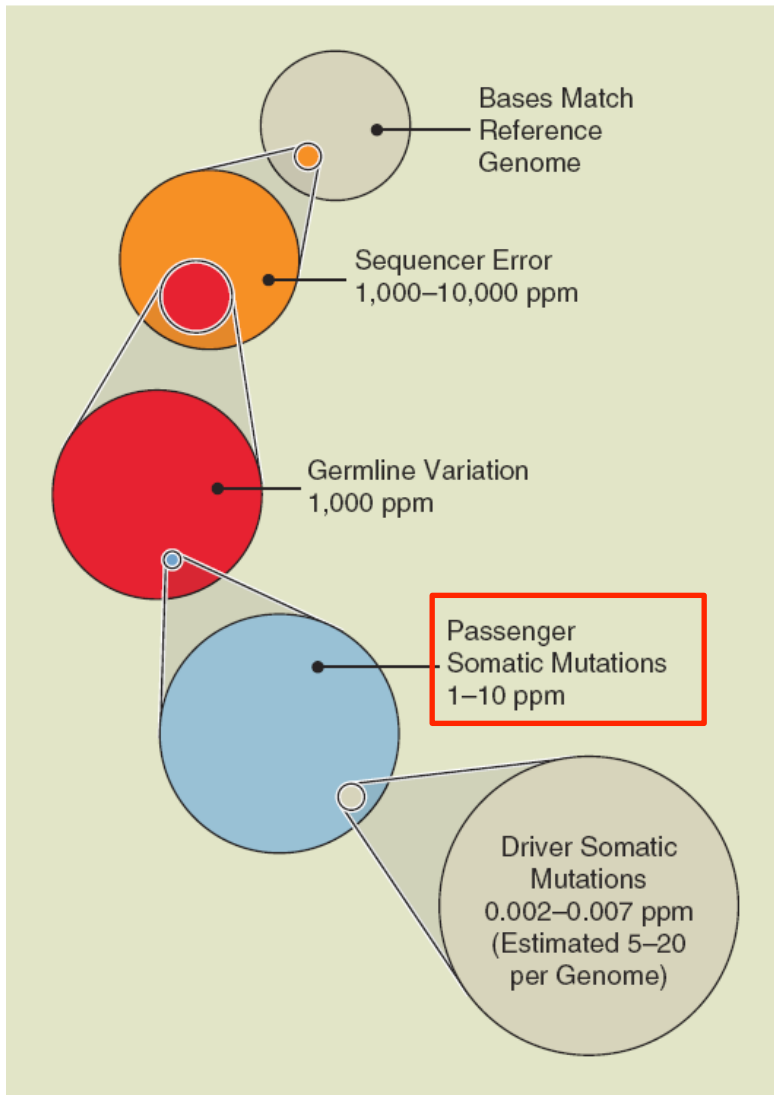


AA AB BB

A	A	A
A	A	A
A	A	A
A	A	A
A	A	A
A	A	B
A	A	B
A	B	B
A	B	B
A	B	B

read from  
normal

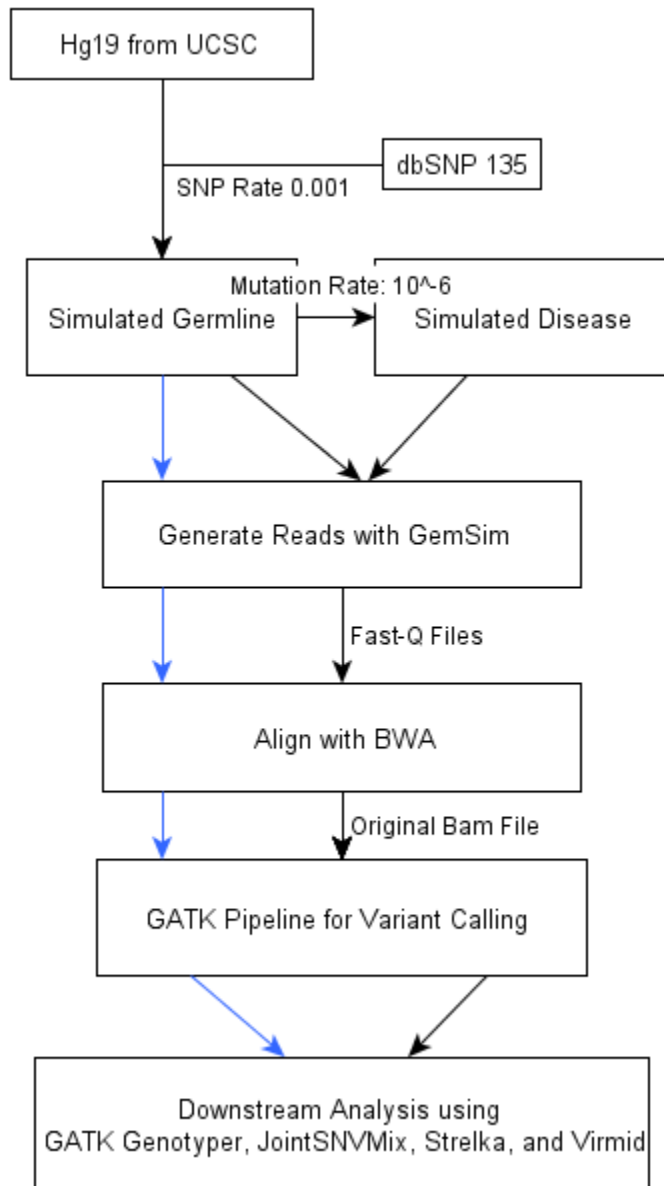
# Somatic mutations are rare and difficult to find



[Signal Processing Magazine, IEEE](#) **Developing Algorithms to Discover Novel Cancer Genes: A look at the challenges and approaches**

# Competitor Tools

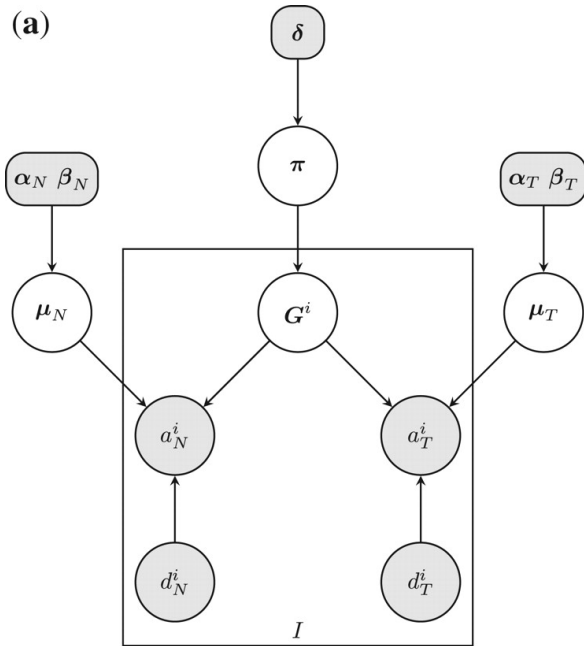
- Other NextGen Tools
  - VarScan
  - Somatic Sniper
- SNP-Array Tools
  - Absolute
  - OncoSNP



# Workflow

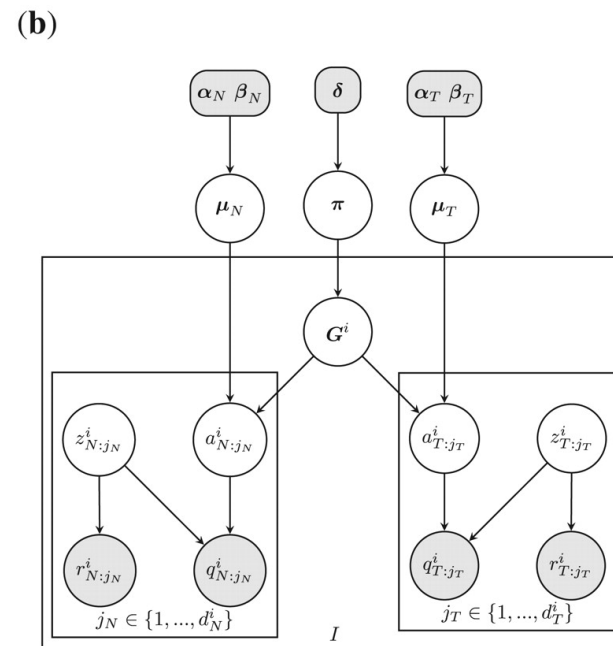
- GATK Unified Genotyper:
- JointSNVMix:
  - Default Settings with the `joint_snv_mix_two` option
- Strelka
  - Default settings with the `bwa` option

# JointSNVMix1



$$\begin{aligned}\pi &\sim \text{Dirichlet}(\pi|\delta) \\ \mathbf{G}^i|\pi &\sim \text{Multinomial}(\mathbf{G}^i|\pi) \\ a_N^i|G_{(g_N, g_T)}^i = 1, \mu_N, d_N^i &\sim \text{Binomial}(a_N^i|d_N^i, \mu_{N:g_N}) \\ \mu_{x:g}|\alpha_{x:g}, \beta_{x:g} &\sim \text{Beta}(\mu_{x:g}|\alpha_{x:g}, \beta_{x:g})\end{aligned}$$

# JointSNVMix2



$$\begin{aligned}\pi &\sim \text{Dirichlet}(\pi|\delta) \\ \mathbf{G}^i|\pi &\sim \text{Multinomial}(\mathbf{G}^i|\pi) \\ a_{N:j_N}^i|G_{(g_N, g_T)}^i = 1, \mu_N &\sim \text{Bernoulli}(a_{N:j_N}^i|\mu_{N:g_N}) \\ z_{N:j_N}^i &\sim \text{Bernoulli}(z_{N:j_N}^i|0.5) \\ q_{N:j_N}^i|a_{N:j_N}^i, z_{N:j_N}^i &\sim f(q_{N:j_N}^i|a_{N:j_N}^i, z_{N:j_N}^i) \\ r_{N:j_N}^i|z_{N:j_N}^i &\sim g(r_{N:j_N}^i|z_{N:j_N}^i) \\ \mu_{x:g}|\alpha_{x:g}, \beta_{x:g} &\sim \text{Beta}(\mu_{x:g}|\alpha_{x:g}, \beta_{x:g})\end{aligned}$$

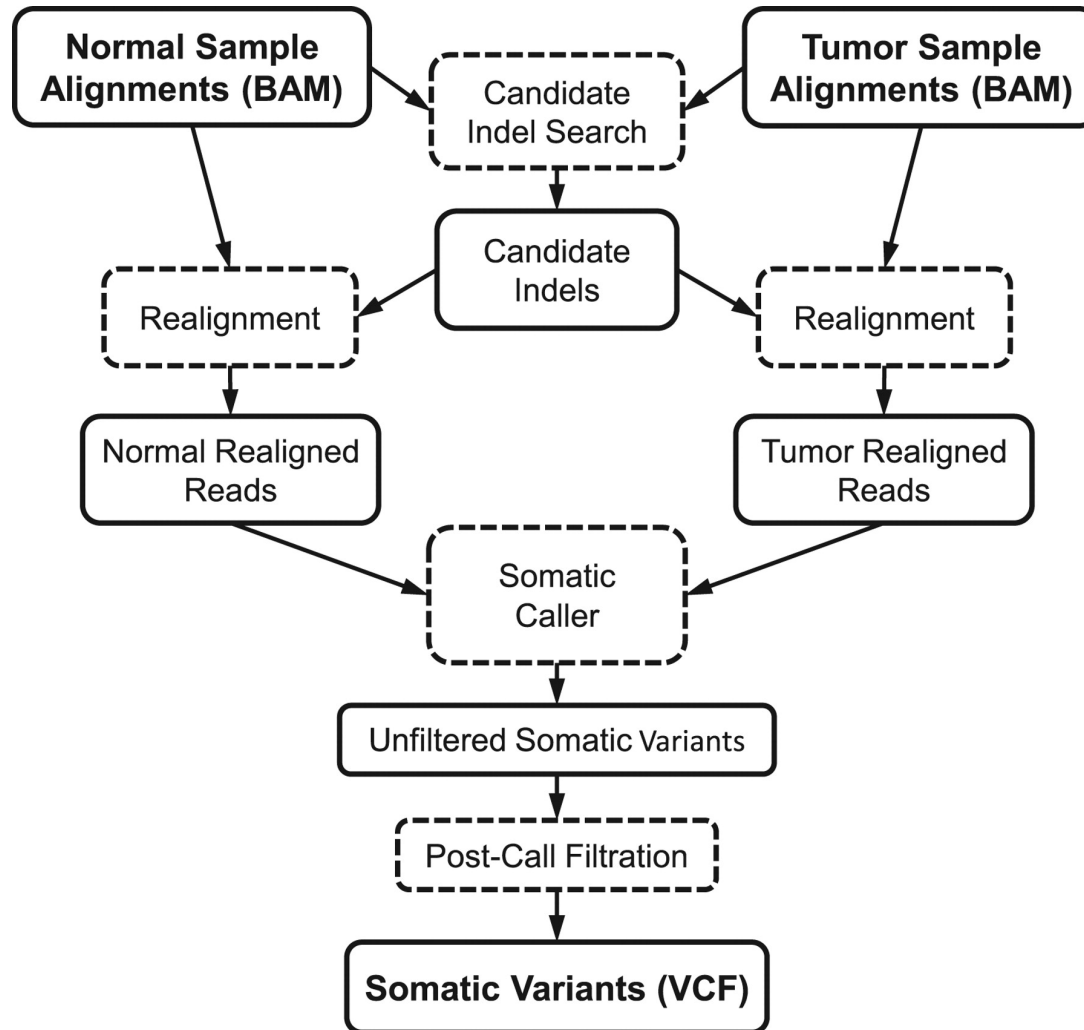
Roth A et al. Bioinformatics 2012;28:907-913

# Strelka

- Developed at Illumina
- Maximizes posterior probability of the joint tumor and normal allele frequencies
  - Other tools only look at genotypes in terms of matching reference or not.
- Is able to call somatic variants as well as detect insertions and deletions
- Does not implicitly calculate contamination or tumor impurity
- High filtering: less false positives

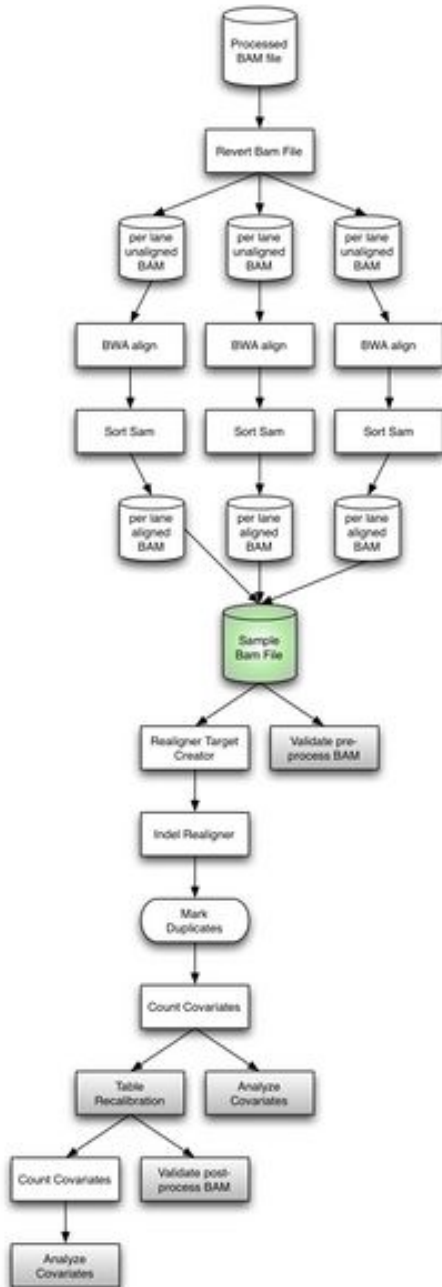


# Strelka



Saunders C T et al. *Bioinformatics* 2012;28:1811-1817

# GATK Best Practices



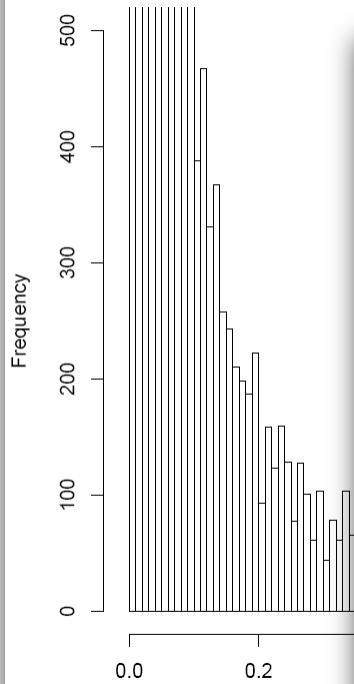
[http://www.broadinstitute.org/gsa/wiki/index.php/Data\\_Processing\\_Pipeline](http://www.broadinstitute.org/gsa/wiki/index.php/Data_Processing_Pipeline)

# AUC (Area under curve)

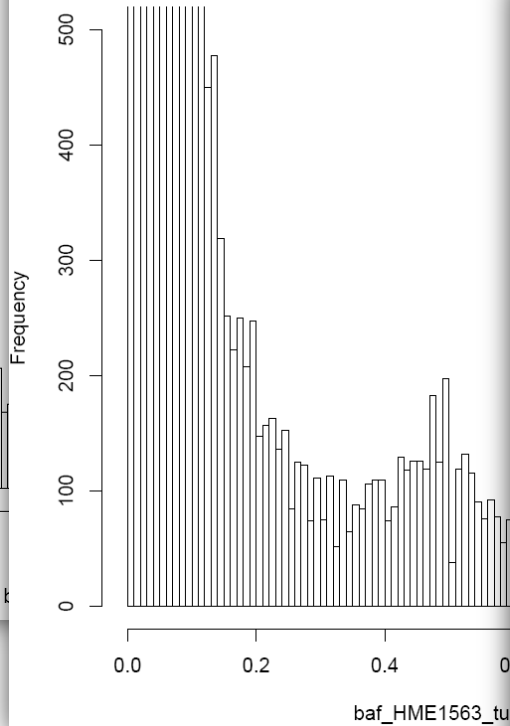
Alpha (%)	JointSNVMix	Virmid( <i>init</i> )	Virmid	AUC diff.
1	0.930	<b>0.935</b>	<b>0.935</b>	< 10 <sup>-3</sup>
5	0.932	<b>0.936</b>	0.935	
10	0.929	0.934	<b>0.935</b>	
20	0.930	<b>0.933</b>	<b>0.933</b>	
30	0.926	0.930	<b>0.931</b>	
40	0.926	<b>0.931</b>	0.930	<b>0.008</b>
50	0.918	0.915	<b>0.927</b>	
60	0.884	0.870	<b>0.910</b>	
70	0.756	0.702	<b>0.854</b>	
80	0.438	0.356	<b>0.741</b>	
90	N/A	0.040	<b>0.302</b>	N/A

# WIC in HME data

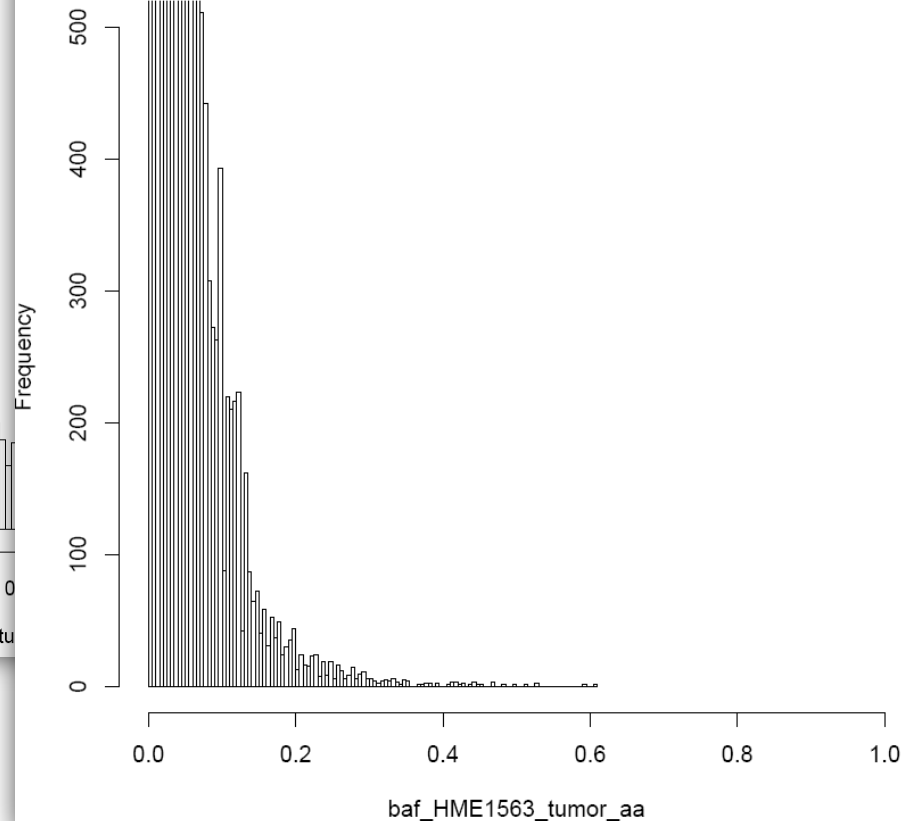
Histogram of baf\_HME1563\_normal



Histogram of baf\_HME1563\_tumor



Histogram of baf\_HME1563\_tumor\_aa



# Bias 1 - Loss of Reads



$x_{\downarrow a} = p(\text{a read that passes } g_{\downarrow 1} \text{ being unmapped})$

$= p(r_{\downarrow 1} \text{ has } d+1 \text{ or more variants in the remaining sites})$

$x_{\downarrow b} = p(\text{a read that passes } g_{\downarrow 2} \text{ being unmapped})$

$= p(r_{\downarrow 2} \text{ has } d+1 \text{ or more variants in the remaining sites})$

$$x_{\downarrow a} = 1 - \sum_{i=0}^{d-1} \binom{l-1}{i} p^i (1-p)^{l-1-i}$$

$$x_{\downarrow b} = 1 - \sum_{i=0}^{d-1} \binom{l-1}{i} p^i (1-p)^{l-1-i}$$

,where  $d$ =maximum edit distance,  $l$ =read length, and  $p$ =frequency of variation

# Bias 2 - Loss of variants