

# A unified haplotype-based method for accurate and comprehensive variant calling

Paper Presentation by Chuanyi Zhang

Daniel P Cooke<sup>1,2</sup>, David C Wedge<sup>2</sup>, and Gerton Lunter<sup>1</sup>

<sup>1</sup>Wellcome Centre for Human Genetics, University of Oxford, Oxford, UK

<sup>2</sup>Big Data Institute, University of Oxford, Oxford, UK

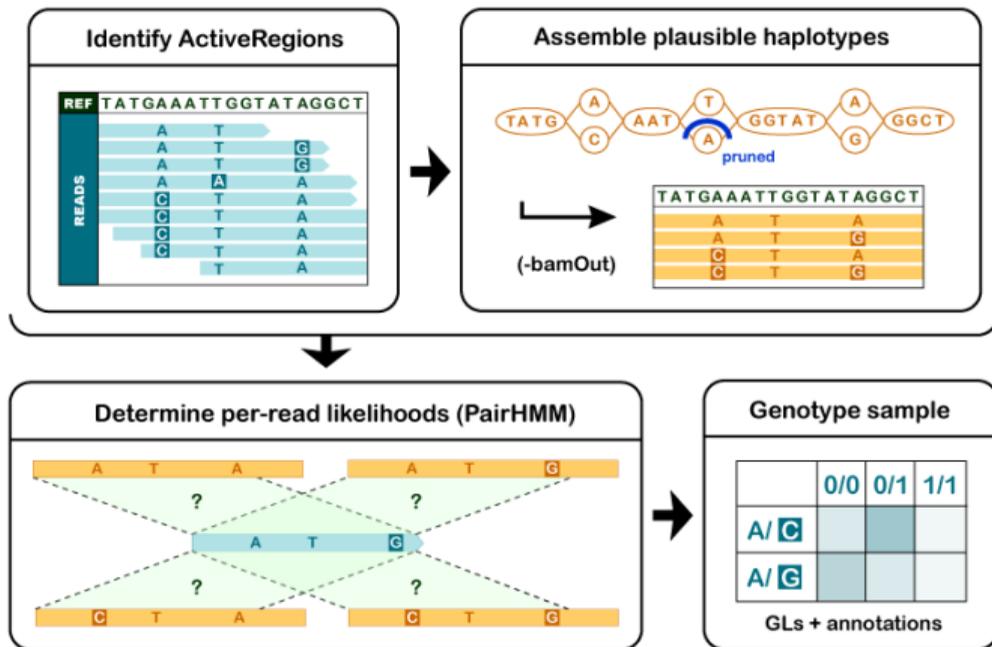
March 7, 2019



# Table of Contents

- 1 Motivation & Background
  - Variant Calling
  - Somatic Variant Calling
- 2 Method
  - Cancer Prior
  - Subclone genotype model
  - Cancer calling model
- 3 Result

# Variant Calling



- Ideal scenario: enough read depth
- (1) read processing,  
(2) mapping,  
(3) calling
- haplotype analysis:  
HaplotypeCaller in Genome  
Analysis Toolkit (GATK)

# Somatic Variant Calling

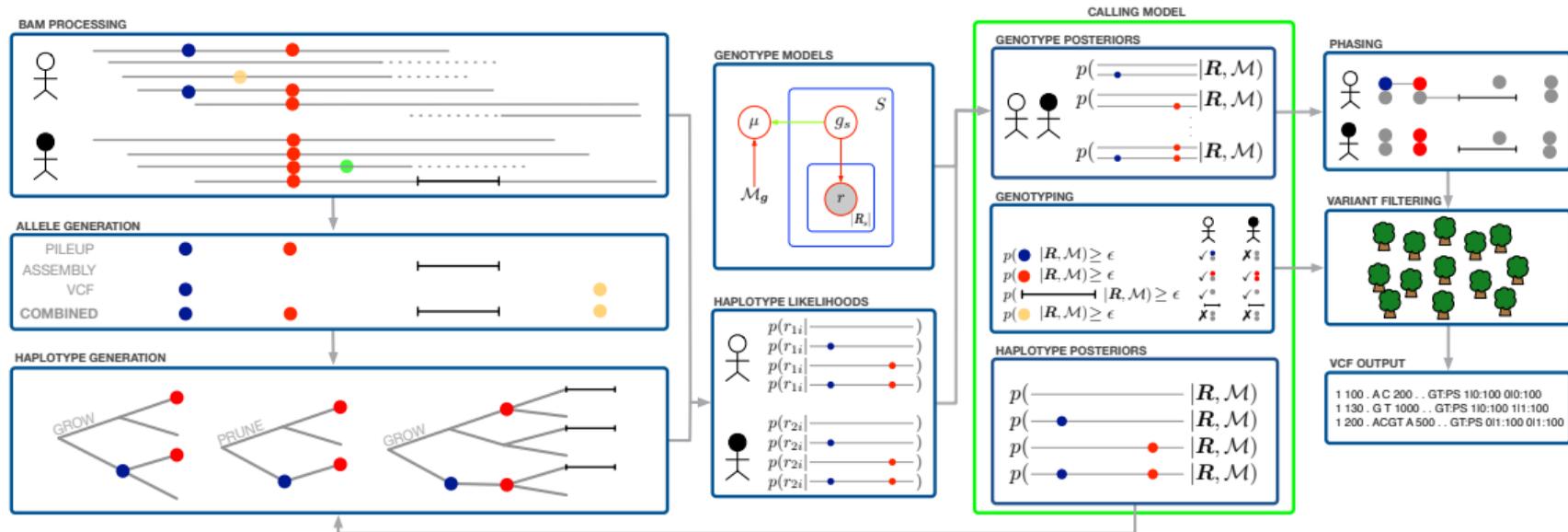
## Differences from germline calling

- Allele frequency assumption: purity, multiple subclones, CNA
- Low VAF vs. Artifacts
- Somatic vs. Germline: matched tumor-normal sample

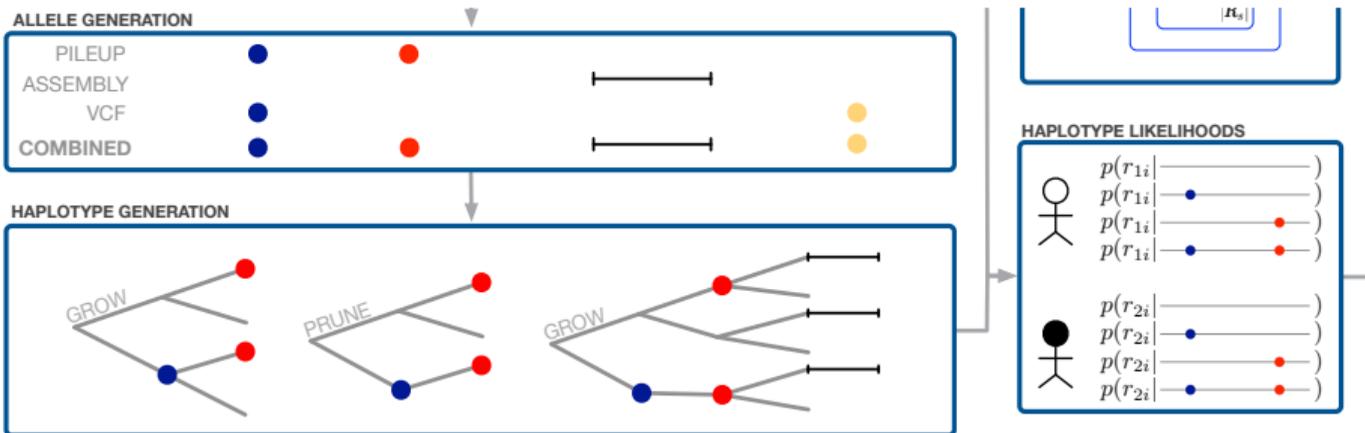
# Table of Contents

- 1 Motivation & Background
  - Variant Calling
  - Somatic Variant Calling
- 2 Method
  - Cancer Prior
  - Subclone genotype model
  - Cancer calling model
- 3 Result

## Overview



# Haplotype generating



- build from candidate alleles
- haplotype tree
- prune, stages: (1) pre: haplotype likelihood, (2) post: haplotype posterior

# Genotype Prior Models

- 1 Uniform
- 2 Hardy-Weinberg-Equilibrium (HWE)
- 3 Coalescent-HWE
- 4 Trio
- 5 ★ Cancer

For ploidy  $m$ , genotypes:  $g = (h_1, \dots, h_m)$ ;

for  $n$  populations inside a tumor, joint genotypes:  $\mathbf{g} = (g_1, \dots, g_n)$ .

# Cancer

$$\mathcal{G}_{cancer} = (\mathcal{G}_{germ}, \mathcal{G}_{som})$$

$p(\mathcal{G}_{cancer} | \mathcal{M}_{cancer}) = p(\mathcal{G}_{germ} | \mathcal{M}_{germ})p(\mathcal{G}_{som} | \mathcal{G}_{germ}, \mathcal{M}_{som})$ ,  $\mathcal{M}_{germ}$  can be Coalescent-HWE prior model,  
if there's only 1 somatic haplotype

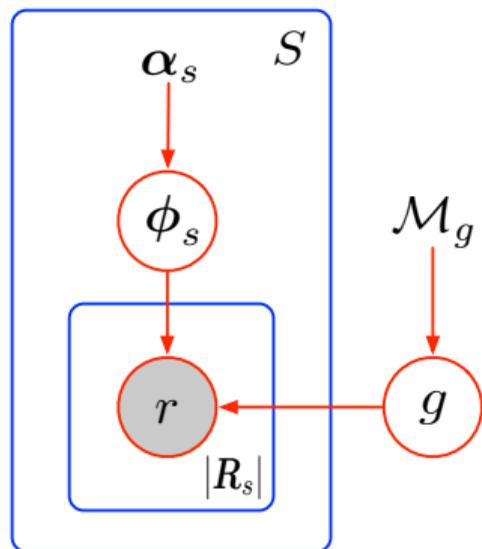
$$p(\mathcal{G}_{som} | \mathcal{G}_{germ}, \mathcal{M}_{som}) = \frac{1}{|\mathcal{G}_{germ}|} \sum_{i=1}^{|\mathcal{G}_{germ}|} p(\mathcal{G}_{som} | \mathcal{G}_{germ,i}, \mathcal{M}_{som})$$

if multi somatic haplotypes: (assume all haplotypes originate from germline, independently)

$$p(\mathcal{G}_{som} | \mathcal{G}_{germ}, \mathcal{M}_{som}) = \prod_{j=1}^{|\mathcal{G}_{som}|} p(\mathcal{G}_{som,j} | \mathcal{G}_{germ}) = \prod_{j=1}^{|\mathcal{G}_{som}|} \frac{1}{|\mathcal{G}_{germ,j}|} \sum_{i=1}^{|\mathcal{G}_{germ,j}|} p(\mathcal{G}_{som} | \mathcal{G}_{germ,j,i}, \mathcal{M}_{som})$$

# Graphical model & joint posterior

We want to know joint posterior distribution



$$p(g, \pi \mid \mathcal{R}, \alpha, \mathcal{M}_g) = \frac{p(\pi, g, \alpha, \mathcal{M}_g, \mathcal{R})}{p(\alpha, \mathcal{M}_g, \mathcal{R})}$$

$$= \frac{p(\mathcal{R} \mid \pi, g)p(\pi \mid \alpha)p(g \mid \mathcal{M}_g)}{p(\alpha, \mathcal{M}_g, \mathcal{R})}$$

$$\propto p(g \mid \mathcal{M}_g) \prod_{s=1}^S p(\mathcal{R}_s \mid \pi_s, g)p(\pi_s \mid \alpha_s)$$

$$= p(g \mid \mathcal{M}_g) \prod_{s=1}^S \int p(\mathcal{R}_s \mid \phi_s, g)p(\phi_s \mid \alpha_s) d\phi_s$$

$$= p(g \mid \mathcal{M}_g) \prod_{s=1}^S \int \prod_{r \in \mathcal{R}_s} \sum_{i=1}^{|\mathcal{G}|} \phi_{si} p(r \mid h_i) p(\phi_s \mid \alpha_s) d\phi_s$$

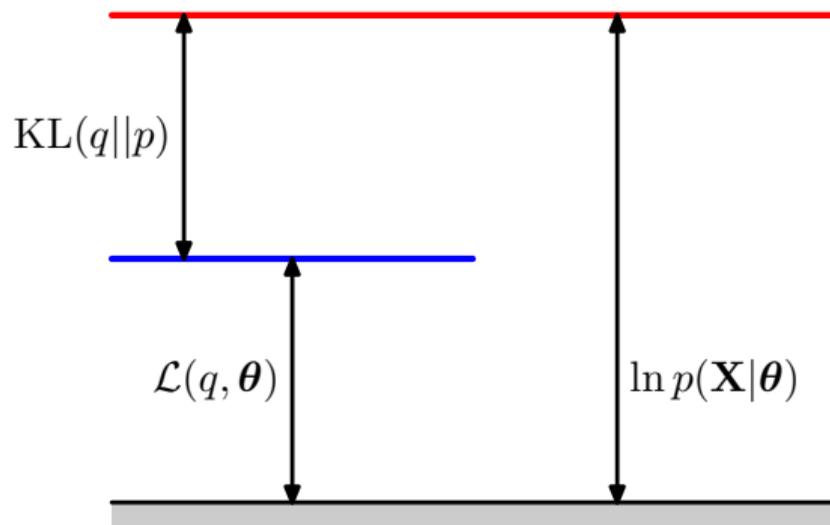
# Problem of computing

- This posterior is intractable Since  $\phi_s \sim \text{Dir}(\alpha_s)$  So the integration over  $\phi_s$  is intractable.  $\phi$  is latent variables.
- Using **Variational Bayes** (VB)

Approximate  $p^*(x) \triangleq p(x | \mathcal{D})$  (intractable posterior) with  $q(x)$ . Maximize  $L(q) \triangleq -D_{KL}(q || \tilde{p})$  (not  $D_{KL}(p^* || p)$ ), where  $\tilde{p} = p(x, \mathcal{D}) = p^*(x)p(\mathcal{D})$

$$\begin{aligned}
 L(q) &= -\mathbb{E}_q \left[ \log \frac{q}{\tilde{p}} \right] = - \int q(x) \log \frac{q(x)}{p^*(x)p(\mathcal{D})} d\mu(x) \\
 &= - \int q(x) \log \frac{q(x)}{p^*(x)} - q(x) \log p(\mathcal{D}) d\mu(x) \\
 &= -D_{KL}(q || p^*) + \log p(\mathcal{D}) \\
 &\leq \log p(\mathcal{D})
 \end{aligned}$$

## ELBO



$L(q)$  is *evidence* lower bound (ELBO). Maximizer is  $q = p^*$ .

## VB cont'

Bayes:

$$p(x | \mathcal{D}) = \frac{p(\mathcal{D} | x)p(x)}{p(\mathcal{D})}, \left( \text{post} = \frac{\text{likelihood} \cdot \text{prior}}{\text{evidence}} \right)$$

And by assuming this factorization

$$q(\mathbf{g}, \mathbf{Z}, \phi) = q(\mathbf{g}) \prod_{s=1}^S q(\mathbf{Z}_s) q(\phi_s)$$

where we introduce the latent binary matrix  $\mathbf{Z}_s$ ,  $q(Z_{snk})$  are so-called *responsibilities* of assuming haplotype  $k$  for read  $n$  in sample  $s$ . By this factorization (mean field) we can optimize on them alternately. Moreover, if we assume these priors are Dirichlet, then prior and posterior are *conjugated*. Categorical (likelihood) and Dirichlet are conjugate distributions.

# Calling Model

Assume 3 possible cases:

- 1 No somatic mutations, clean germline, the individual model with any germline prior (merge)  $\mathcal{M}_{ind}$
- 2 Copy number changes, but no somatic, the *subclone model* with germline prior (e.g. Coalescent-HWE)  $\mathcal{M}_{ind}$
- 3 Somatic occurs, possible CNA, the *subclone model* with cancer genotype prior.

# Calling Model

Germline genotype posterior

$$\begin{aligned}
 p(g | \mathcal{R}) &= \sum_x p(g, \mathcal{M}_x | \mathcal{R}) \\
 &= \sum_x p(g, | \mathcal{M}_x, \mathcal{R}) p(\mathcal{M}_x | \mathcal{R}) \\
 &= p(g, | \mathcal{M}_{ind}) p(\mathcal{M}_{ind} | \mathcal{R}) \\
 &\quad + p(g, | \mathcal{M}_{CNV}) p(\mathcal{M}_{CNV} | \mathcal{R}) \\
 &\quad + p(g, | \mathcal{M}_{somatic}) p(\mathcal{M}_{somatic} | \mathcal{R})
 \end{aligned}$$

where  $p(\mathcal{M}_x | \mathcal{R}) = p(\mathcal{M}_x) p(\mathcal{R} | \mathcal{M}_x)$ , and  $p(\mathcal{R} | \mathcal{M}_x)$  is the “evidence”; and  $p(g | \mathcal{M}_{somatic}) = \sum_{\tilde{g}: g \in \tilde{g}} p(\tilde{g} | \mathcal{M}_{somatic})$ ,  $\tilde{g} = (g_{germ}, g_{som})$ , from *cancer prior* Germline allele posterior  $p(a | \mathcal{R}) = \sum_{g: a \in g} p(g | \mathcal{R})$ .

## Credible somatic mass

$$p_{\text{somatic}}(a | \mathcal{R}) \leftarrow \sum_{\tilde{g}} p(\tilde{g} | \mathcal{M}_{\text{somatic}}, \text{credible})$$

There are  $K$  somatic haplotypes, then the *credible somatic frequencies* satisfy

$$p(\phi_{sk} > \tau | \mathcal{M}_{\text{somatic}}) = \int_{\tau}^1 \text{Beta}\left(\theta; \alpha_{P+1}, \sum_{i=1}^P \alpha_i\right) d\theta$$

where  $\phi_{sk} \sim \text{Beta}(\alpha_k, \sum \alpha - \alpha_k)$  since  $\phi_s \sim \text{Dir}(\alpha_s)$  i.e.  $p(\phi_s) = \frac{1}{B(\alpha)} \prod_{k=1}^K \phi_{sk}^{\alpha_k - 1}$ .

The *credible somatic mass* is

$$\lambda_s = 1 - \prod_k 1 - p(\phi_{sk} > \tau | \mathcal{M}_{\text{somatic}})$$

means the probability mass of  $\exists$  at least 1 credible in  $K$  somatic haplotypes.

# Calling allele

Then  $\lambda = 1 - \prod_s \lambda_s$ .  $(\overline{\exists_1 \wedge \dots \wedge \exists_S} = \#_1 \vee \dots \vee \#_S)$

$$p_{somatic}(a | \mathcal{R}) = \lambda \left( 1 - \prod_s \sum_a \mathbb{1}_{\{a \notin \check{g}.germ \wedge a \in \check{g}.som\}} p(\check{g} | \mathcal{R}, \mathcal{M}_{somatic}) \right)$$

Might be a typo?

# Table of Contents

- 1 Motivation & Background
  - Variant Calling
  - Somatic Variant Calling
- 2 Method
  - Cancer Prior
  - Subclone genotype model
  - Cancer calling model
- 3 Result

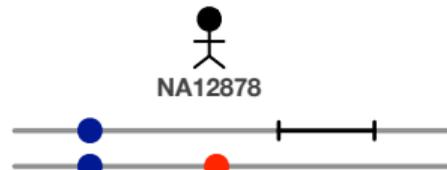
# Synthetic Tumors

Evaluation of somatic mutation calling is challenging

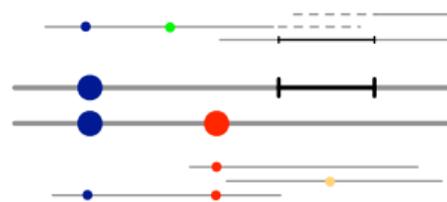
- Real tumor with manual inspected mutations
- Mix reads from unrelated individuals
- ★ Spike mutations directly into raw sequencing reads from healthy tissue

# Synthetic Tumors

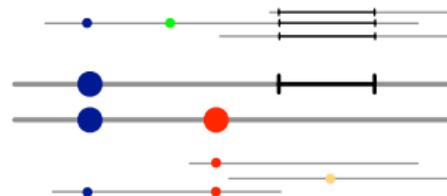
## 1. Select sample with known germline



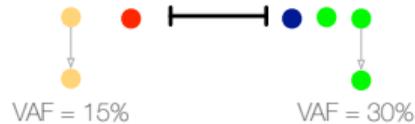
## 2. Assign reads to germline haplotypes



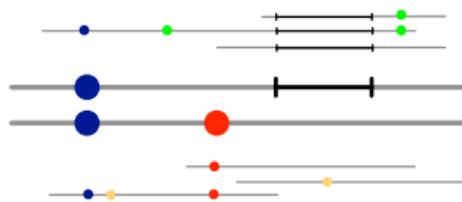
## 3. Realign reads to germline haplotypes



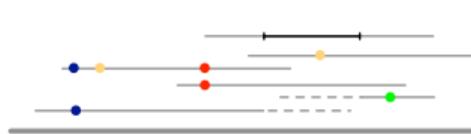
## 4. Sample PCAWG tumour-specific calls



## 5. Spike PCAWG mutations onto reads



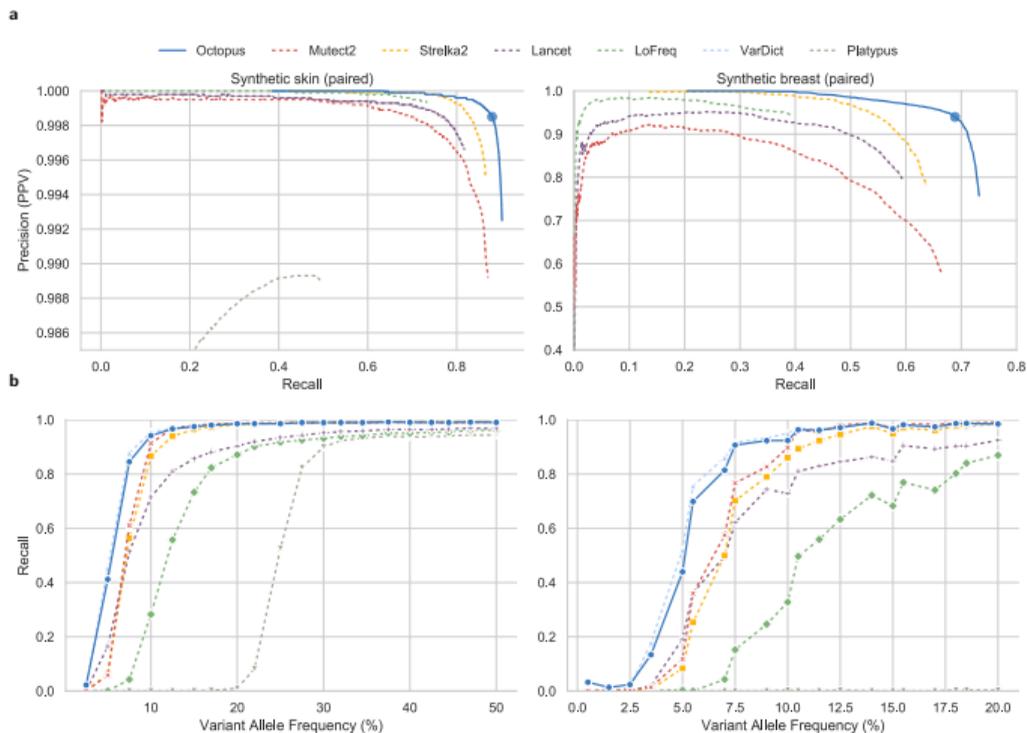
## 6. Remap spiked reads



- Reads from NA12878  $\sim 300\times \rightarrow \{30\times, 35\times, 60\times, 65\times\}$
- Assign and realign to make sure spiked mutations fall on same haplotype, and position consistent
- Sample mutations from pan-cancer analysis of whole genomes (PCAWG) uniformly

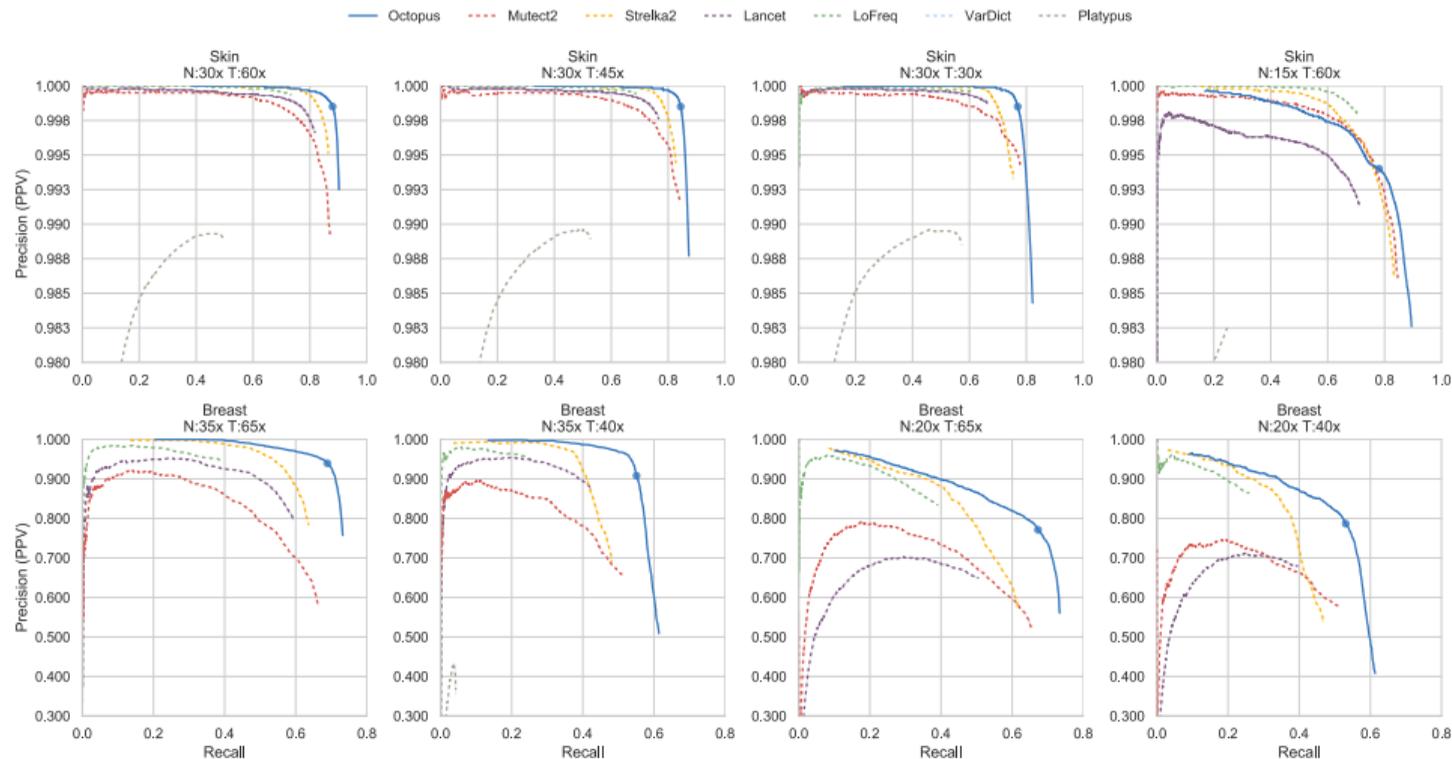
Produce raw unmapped reads (FASTQ) files.

# Somatic Mutations Calling Accuracy



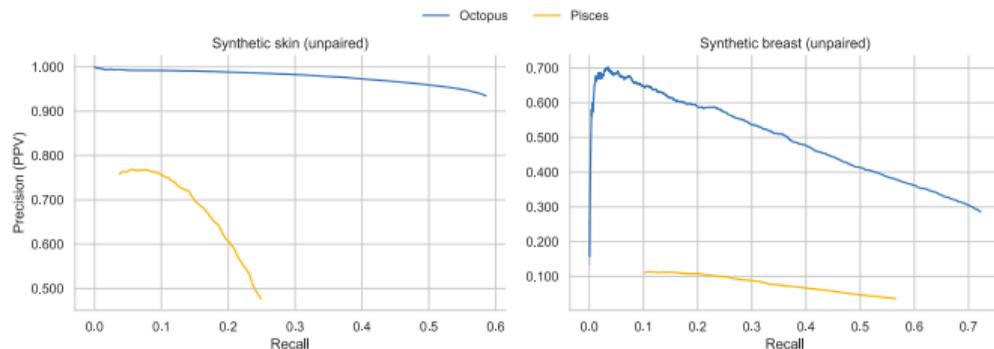
- 1 Precision-Recall curve: top-right is optimal
- 2 Recalls for each VAF, using PASS variants
- 3 Most differences in recall is due to low frequencies

## Low coverage

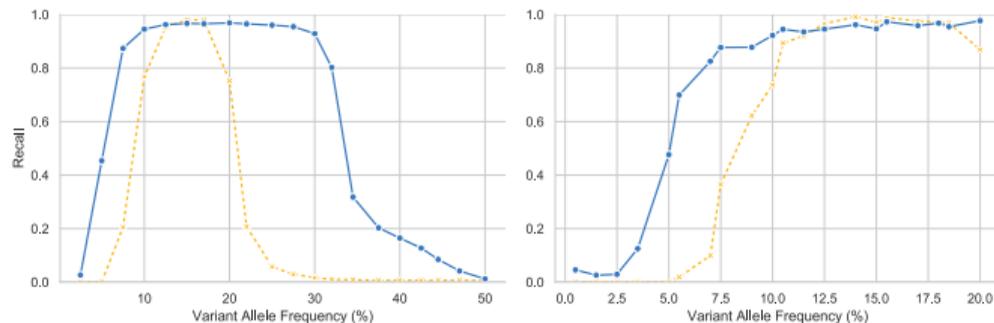


# Somatic Mutations Calling Accuracy without paired normal

a



b



- 1 Octopus is able to discover mutations even without paired normal sample