

BIOINFORMATICS 101: GENOME ANALYSIS TOOLKIT (GATK) 4

W. Bailey Glen Jr.

02/26/2017

Broad Institute

- Independent research center partnered with Harvard and MIT
- Cambridge, MA
- Broad Genomics
 - Large Scale Sequencing Core
 - One 30X human whole genome every 12 minutes
 - Computation Model and Software Development



Software Tools

- GATK
- Integrated Genome Viewer
- Hail
- Tumor Portal

GATK

- Analyzing sequencing analysis for genetic variation
 - Primarily Focused on Small Variants
 - Copy number changes recently added
 - Primarily Focused on DNA
 - RNA best practices have been previously defined
 - Primarily Constitutional
 - Tumor/Normal developed
 - Tumor only workflow now described
- Not
 - RNA expression analysis
 - Linkage disequilibrium or association testing

GATK 4 – Refactoring and more

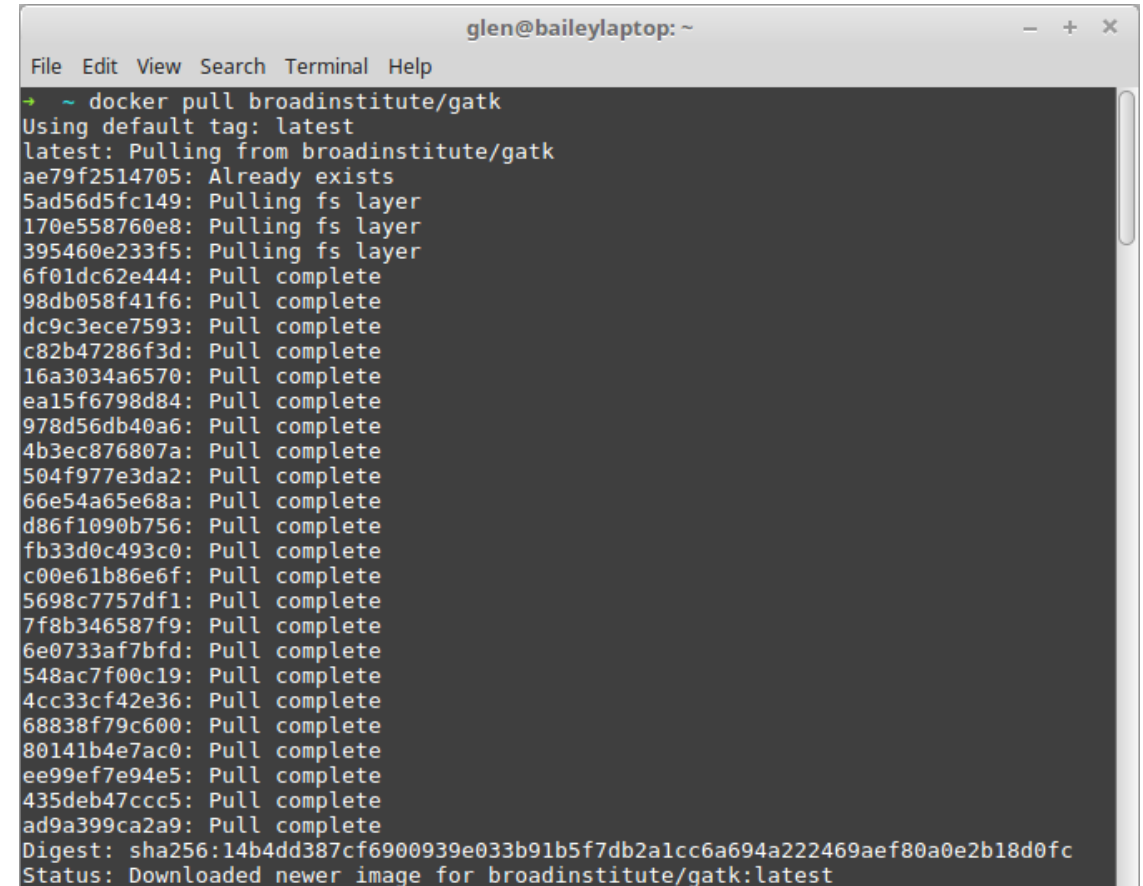
- In development for years
- Refactored GATK 3
 - In many cases, same commands producing equivalent results
- Partnered with Intel
- Large focus on computer science and software implementation
 - Performance
 - Deployment
 - Scalability
- Opensource

Component Software

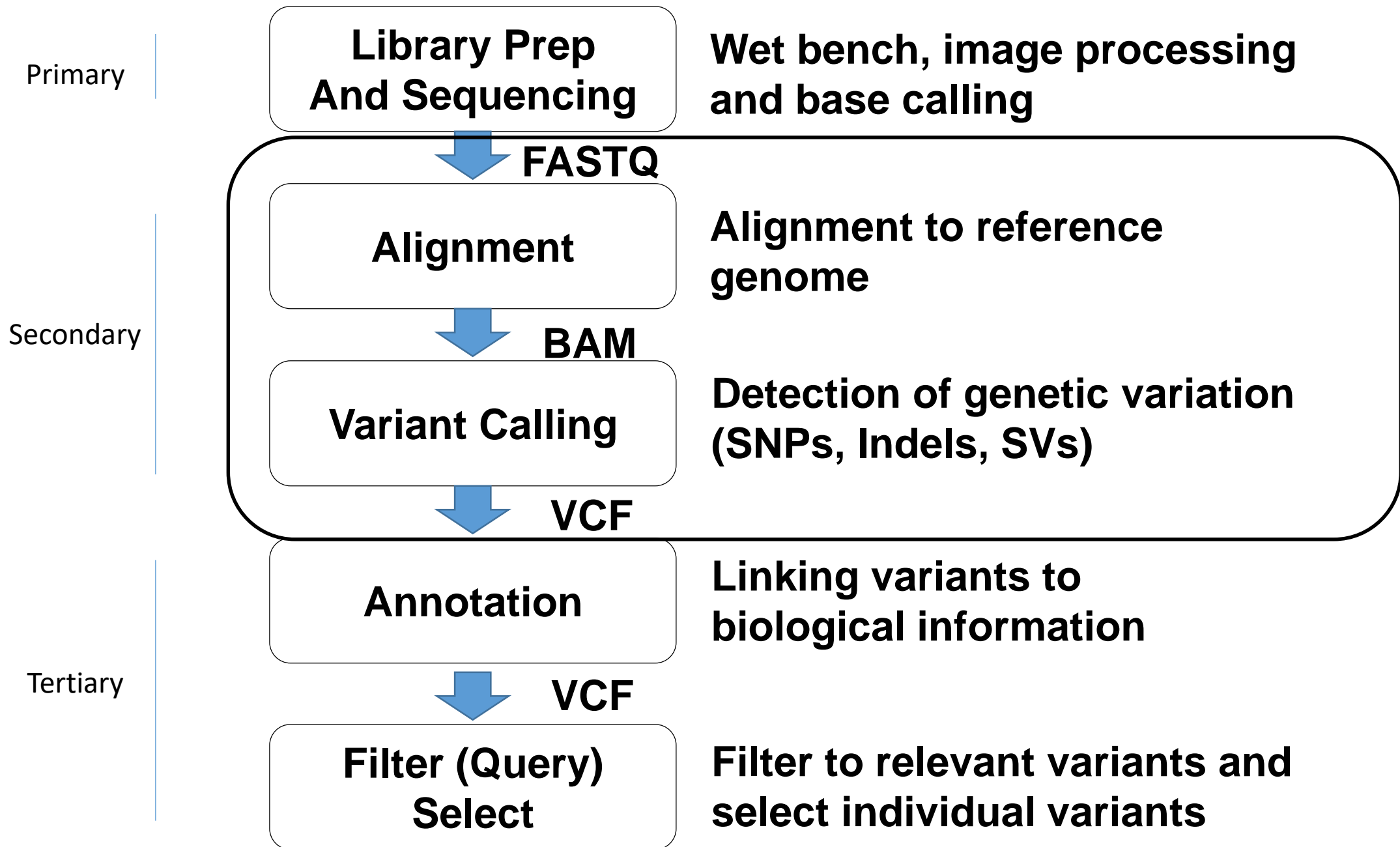
- GATK
 - Java
- Dependencies
 - JVM
 - BWA
 - Picard tools
- Available as Docker or JAR

Docker

- Container Virtualization
 - Runs on linux, mac, windows
 - Creates an isolated controlled virtual computer on your computer
 - Comes with packages installed and version controlled
 - Dockers are "easily" deployed
- Emerging standard for sharing bioinformatic tools
- One new tool to learn, improve using many tools



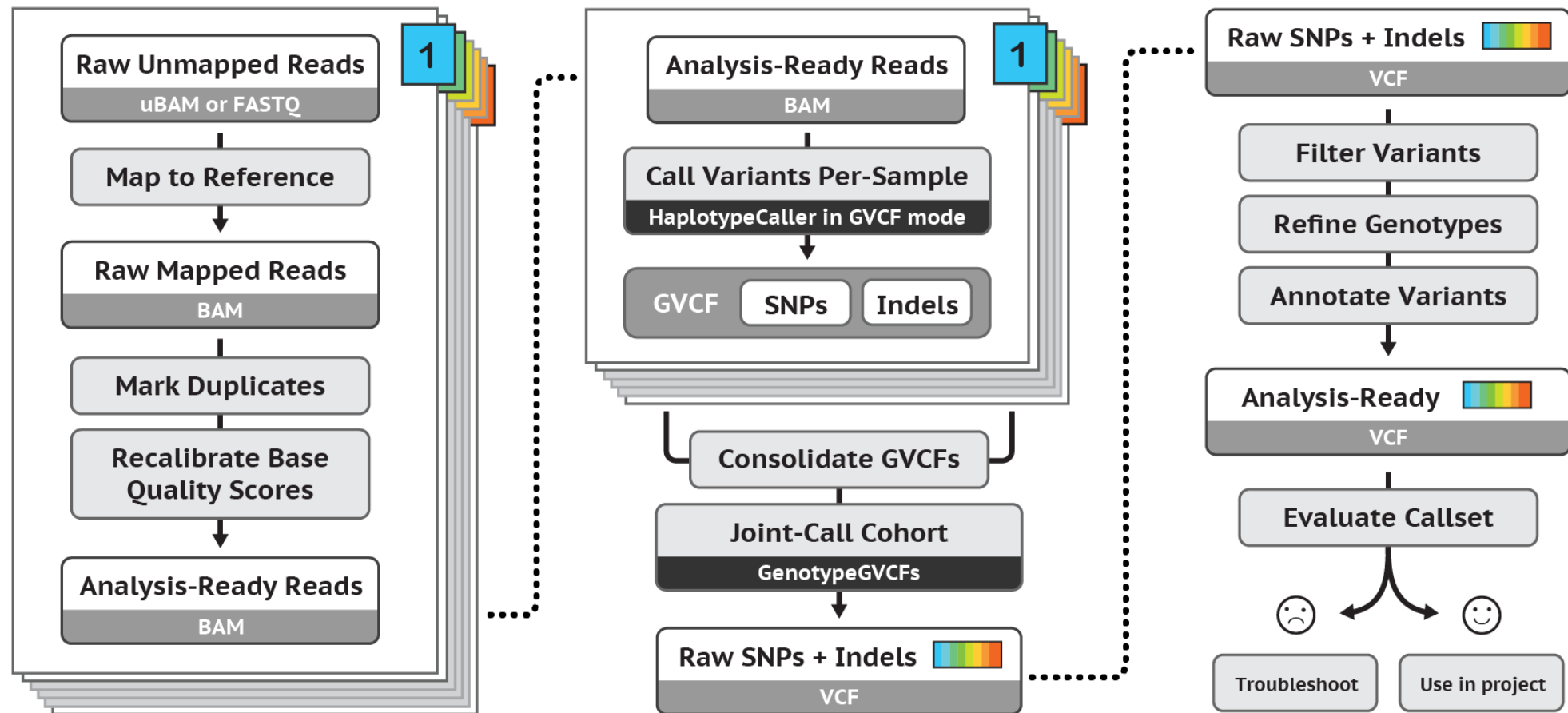
```
glen@baileylaptop: ~  
File Edit View Search Terminal Help  
→ ~ docker pull broadinstitute/gatk  
Using default tag: latest  
latest: Pulling from broadinstitute/gatk  
ae79f2514705: Already exists  
5ad56d5fc149: Pulling fs layer  
170e558760e8: Pulling fs layer  
395460e233f5: Pulling fs layer  
6f01dc62e444: Pull complete  
98db058f41f6: Pull complete  
dc9c3ece7593: Pull complete  
c82b47286f3d: Pull complete  
16a3034a6570: Pull complete  
ea15f6798d84: Pull complete  
978d56db40a6: Pull complete  
4b3ec876807a: Pull complete  
504f977e3da2: Pull complete  
66e54a65e68a: Pull complete  
d86f1090b756: Pull complete  
fb33d0c493c0: Pull complete  
c00e61b86e6f: Pull complete  
5698c7757df1: Pull complete  
7f8b346587f9: Pull complete  
6e0733af7bfd: Pull complete  
548ac7f00c19: Pull complete  
4cc33cf42e36: Pull complete  
68838f79c600: Pull complete  
80141b4e7ac0: Pull complete  
ee99ef7e94e5: Pull complete  
435deb47ccc5: Pull complete  
ad9a399ca2a9: Pull complete  
Digest: sha256:14b4dd387cf6900939e033b91b5f7db2a1cc6a694a222469aef80a0e2b18d0fc  
Status: Downloaded newer image for broadinstitute/gatk:latest
```



GATK Best Practices

- <https://software.broadinstitute.org/gatk/best-practices/>
- Small Nucleotide Polymorphisms
 - Germline SNPs + Indels
 - Somatic SNVs + Indels
 - RNAseq SNPs + Indels
- Copy Number Variations
 - Germline CNVs
 - Somatic CNVs

Example Genome Analysis Toolkit Workflow



Major File Types

- Unaligned Reads

- FASTQ or uBAM

```
@HWI-ST1001:97:D0E5LABXX:7:1101:1180:2041 1:Y:0:  
AATCNNTAGTTGGTGGCGTAAGGTCGCAAAGTAAGAGCTTCTCGGGCTGCGTCAGGATAGGTCGTATTTGCTCATTTCCTCCCTTCAGTCCACTTCGATTT  
+  
(079##2(3@@9)@(@#####
```

- Aligned Reads

- BAM

```
HWI-ST1001:97:D0E5LABXX:1:1102:11138:107030  
163  
chrM  
1  
41  
28S70M3S  
=  
1  
70  
CGTTCCCCTTAAATAAGACATCACGATGGATCACAGGTCTATCACCTATTAACCACTCACGGGAGCTCTCCATGCATTTGGTATTTTCGTCTGGGGGCTA  
;39<====>=<<<<<<>=<==<=<6=====<=<=>=====<=====>=>><=<=<=<6=>=>=>=>??=??><?><<?>>8<=?>>=;9:8
```

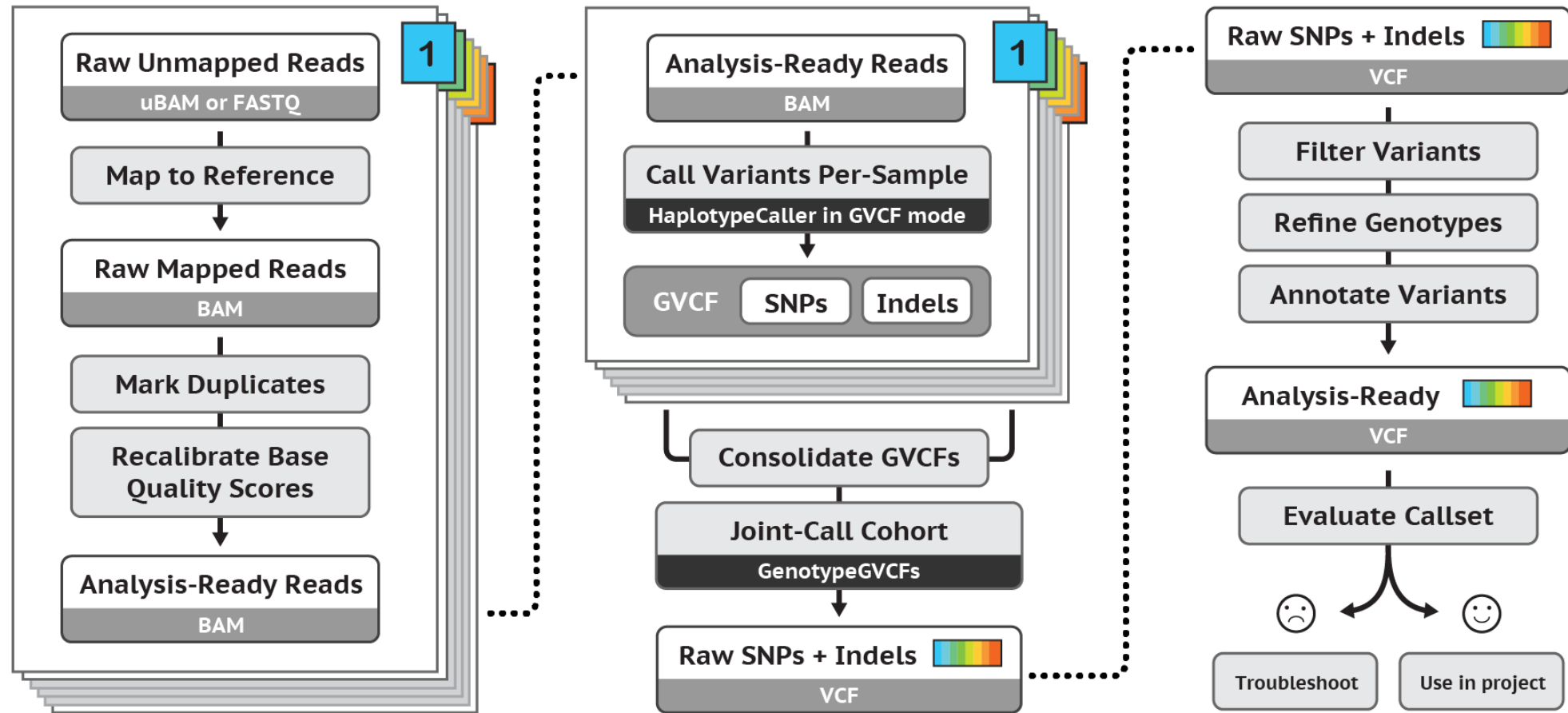
- Variant Call Files

- GVCf (intermediate)
- VCF

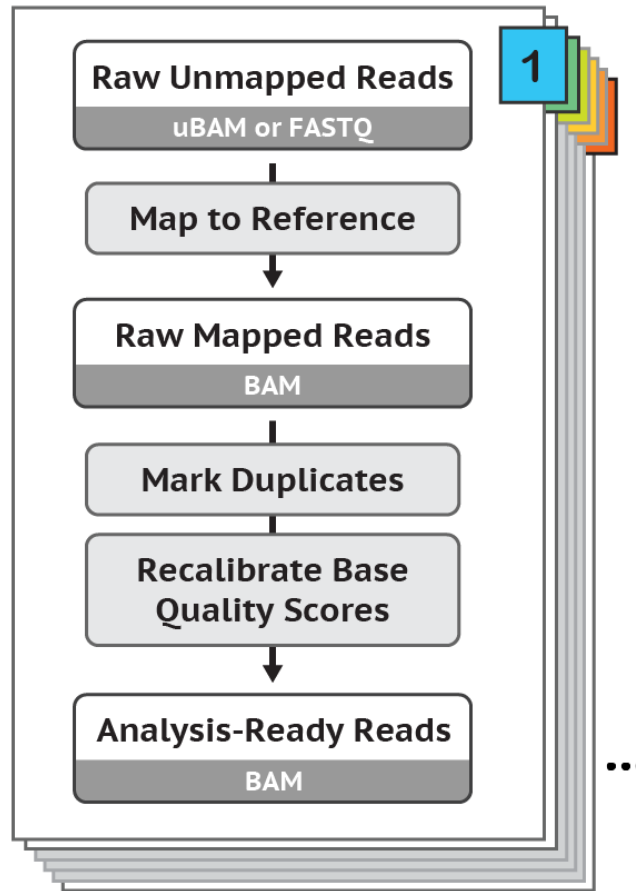
Alignment Software

- Aligns reads to Reference Genome
 - Database Search Alignment
 - BLAST
 - Short Read Sequence Alignment
 - BWA-MEM
 - BOWTIE
- Reads can also be aligned to themselves if a reference alignment is missing (*de novo* assembly)
 - Long reads

Constitutional DNA Best Practices



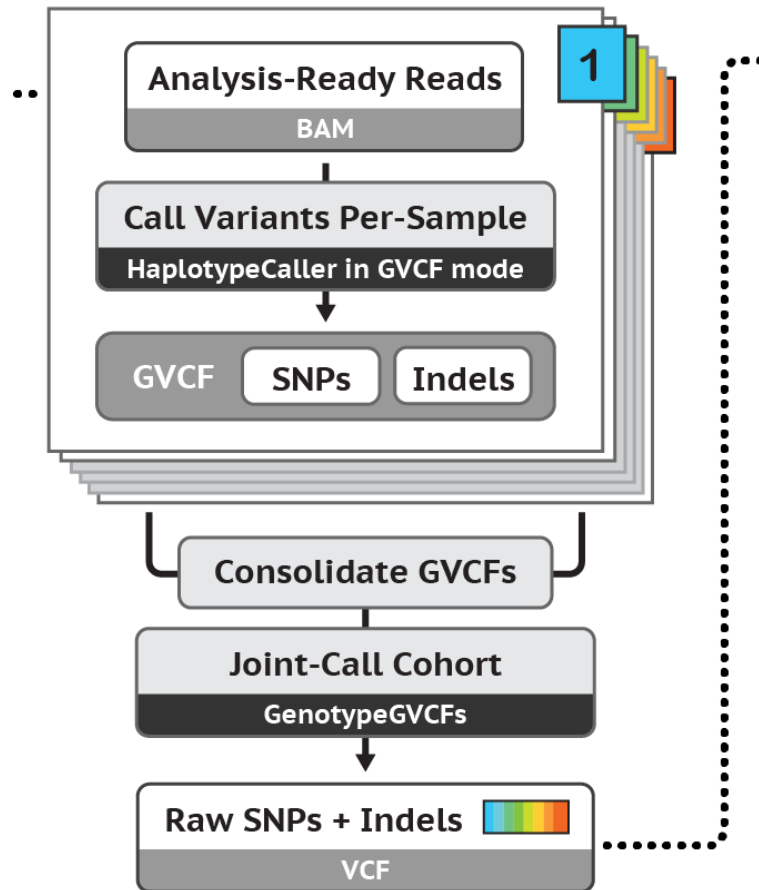
Preprocessing (Alignment)



Starting with unaligned reads

- Align to Reference Gene
 - BWA-MEM for DNA
 - STAR/BOWTIE for RNA
- Remove Duplicate Reads
 - Unless using an amplicon, most duplicate reads are sequencing errors
- Recalibrate Base Quality Scores
 - Statistical approach to improve base calls after completing a run

Alignment



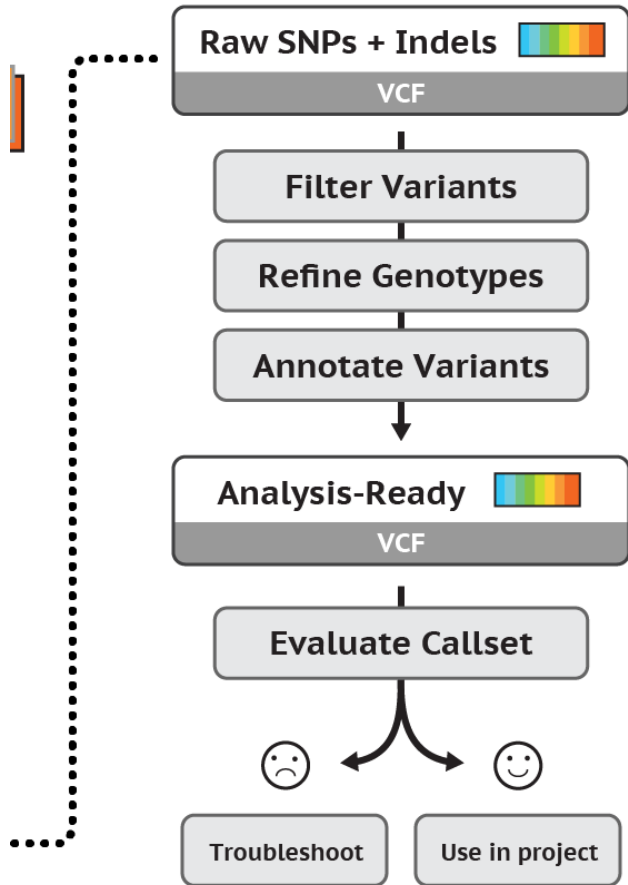
Starting with preprocessed bams

- Generate gVCF for each sample
- Join call gVCFs

OR

- Joint call bams to make VCF

Variant Calling



Starting with VCF

- Filter by Quality
- Refine Variants
 - Phasing (Trio, local within sample)
- Annotate Variants
 - Provide genomic context
- Perform Tertiary