# Bioinformatics for NGS analysis
## From Reads to Variants

### Tim Yu, MD, PhD

*Division of Genetics, Boston Children's Hospital*
*Dept of Neurology, MGH*
*Harvard Medical School & the Broad Institute*

*2013 NSGC Annual Education Conference*

# Bioinformatics for NGS analysis
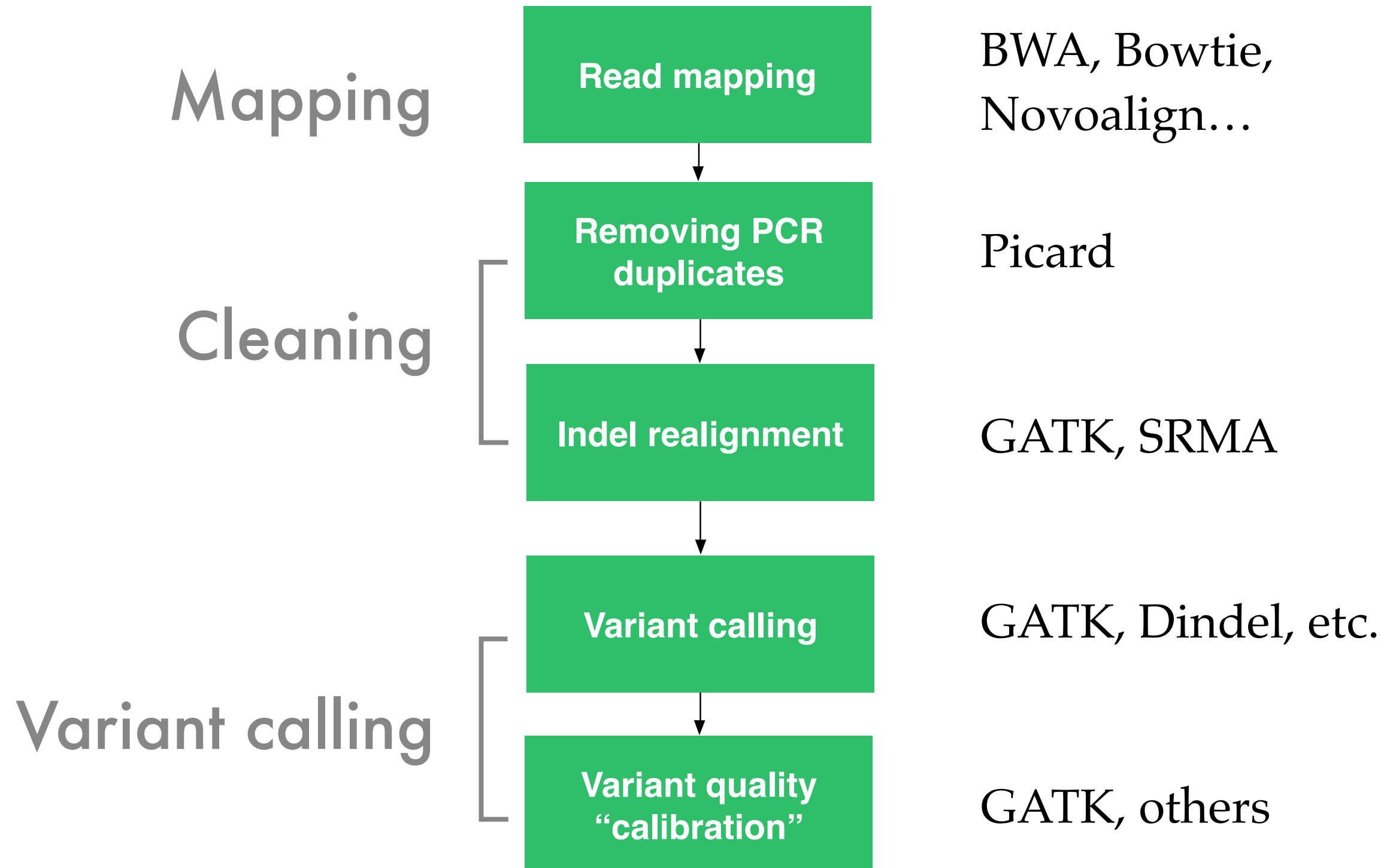## From Reads to Variants

**Tim Yu, MD, PhD**

Disclosure:
Co-founder & Principal consultant,
Claritas Genomics

- Goal: arm you with concepts and vocabulary to understand how NGS data is analyzed, and to ask critical questions

# A typical NGS processing pipeline

**Mapping**

Read mapping

BWA, Bowtie, Novoalign…

**Cleaning**

Removing PCR duplicates

Picard

Indel realignment

GATK, SRMA

**Variant calling**

Variant calling

GATK, Dindel, etc.

Variant quality "calibration"

GATK, others

NGS analysis is, in principle, a two step process

# 1. Millions/billions of reads are mapped *en masse* to a reference genome

# 1. Millions/billions of reads are mapped *en masse* to a reference genome

# 2. Variants are detected when enough reads disagree with reference

# Complicating factors

- **Mapping can be tricky**

- **Sequencing coverage is biased**

- **Not all variant calls are created equal**

- **Beyond SNPs and small indels**

# 1. Mapping can be tricky

- **Easy**: Perfect matches to unique genomic regions A, B, and C

Mapping and alignment algorithms

Enormous pile of short reads from NGS

Reference Genome

gene A    gene B    gene C

Reads mapped to reference

Image credit: Broad GSA Platform

# 1. Mapping can be tricky



**Harder**: Imperfect matches to unique genomic regions A, B, and C

Mapping and alignment algorithms

Enormous pile of short reads from NGS

Reference Genome

gene A    gene B    gene C

Reads mapped to reference

Image credit: Broad GSA Platform

# 1. Mapping can be tricky



- **Hardest:** Mapping to related genomic regions:
  - Gene families
  - Pseudogenes
  - Repeats / segmental dups
  - CNVs

Image credit: Broad GSA Platform

# Mapping confidence/mapability

- Mapping confidence is a prerequisite for good variant calls
- But mapability can vary quite a bit!



**Read length**

- 36 bp

- 75 bp

- 100 bp

# Mapping confidence/mapability

- Mapping confidence is a prerequisite for good variant calls
- But mapability can vary quite a bit!



**Read length**

- 36 bp
- 75 bp
- 100 bp

# Solutions to the mapability problem

- Longer reads
- Paired-end and mate-pair sequencing
- Better reference sequences (eg taking into account CN variable regions)

# 2. Sequencing coverage is biased

Read coverage on chr7 for a
a typical WES (whole exome sequencing) experiment



Refseq genes

Coverage

**Contributing factors:**
- intentional (eg, exome capture design)
- unintentional
  - PCR-related (eg, GC-rich regions)
  - mapability

# …& may result in gaps in coverage (insufficient breadth)



- **Example:** absent read coverage over *CFTR* exon 10

- **Consequence: variant dropout (false negatives)**

# ...or just inadequate coverage (insufficient depth)



- **Example:** low read coverage over *CFTR* exons 1 & 24

# ...or just inadequate coverage (insufficient depth)

- **Example:** low read coverage over *CFTR* exons 1 & 24



- SNP on 1/6 reads
- Is this a heterozygous variant?

- SNP on 5/6 reads
- Is this a heterozygous or homozygous variant?

- **Consequence:**
**Inaccurate genotyping in areas of low read depth**

# Excess coverage is sometimes a red flag



- Caused by:
- PCR duplicates
- CN expansions

- "Cleaning" alignments by **finding and removing PCR duplicates** evens out coverage, and reduces false positives

# Depth and breadth are usually a tradeoff

- **Given fixed $$$:  Depth or Breadth, choose one!**

<table>
<tr><td align="center"><b><u>Shallow & wide</u></b></td><td align="center"><b><u>Narrow & deep</u></b></td></tr>
<tr><td align="center">more variants</td><td align="center">fewer variants</td></tr>
<tr><td align="center">less accurate genotypes</td><td align="center">more accurate genotypes</td></tr>
<tr><td align="center">e.g., "exomes" at 50-150X</td><td align="center">e.g., "panels" at 500-1500X</td></tr>
</table>

- **Costs are gradually dropping so hopefully this tradeoff will become moot!**

# Solutions to the coverage bias problem

- [Optimize mapability (longer reads, paired end sequencing, etc.)]

- Optimize library prep

  - Minimize PCR, or use PCR-free library prep methods

  - Normalize baits

- Informatically, find and remove PCR duplicates

# 3. Not all variant calls are created equal

- **We do quite well with SNPs (i.e., single base substitutions)**

- Calls are reliable: >99% concordance with chip-based SNP genotyping or other "truth sets"

# 3. Not all variant calls are created equal

- **But indels (i.e., small insertions or deletions) are significantly harder**

- It is computationally hard to map a 100bp read to the genome if you allow for gaps

- Sensitivity estimates vary hugely (50-90%), & 2-10X more false positives (compared to SNPs)

# Example: calling around homopolymers

- Small insertions/deletions (especially near the ends) can trick mappers into misaligning with mismatches

**10bp "T" homopolymer run**

ref: TGACTCGTAACCAGGCTTTTTTTTTTTGCGGGCCGAA

# Example: calling around homopolymers

- Small insertions/deletions (especially near the ends) can trick mappers into misaligning with mismatches

**10bp "T" homopolymer run**

```
ref: TGACTCGTAACCAGGCTTTTTTTTTTTGCGGGCCGAA
reads:    TCGTAACGAGGCTTTTTTTTTTTGCGGGC
                 AGGCTTTTTTTTTTTGCGGGCCGAA
          GACTCGTAACGAGGCTTTTTTTTTTTGC
                CGAGGCTTTTTTTTTTTGCGGGCCG
          TGACTCGTAACGAGGCTTTTTTTTTTTG
```

**many single-bp mismatches?**

# Example: calling around homopolymers

- Small insertions/deletions (especially near the ends) can trick mappers into misaligning with mismatches

**10bp "T" homopolymer run**

```
ref: TGACTCGTAACCAGGCTTTTTTTTTTTGCGGGCCGAA
reads:    TCGTAACGAGGCTTTTTTTTTT^GCGGGC
                 AGGCTTTTTTTTTT^GCGGGCCGAA
          GACTCGTAACGAGGCTTTTTTTTTT^GC
                 CGAGGCTTTTTTTTTT^GCGGGCCGAA
          TGACTCGTAACGAGGCTTTTTTTTTT^G
```

**Local realignment reveals a hidden 1bp delT**

# Red flags that a variant may be suspicious

- In fact, raw indel calls are infested with false positives
- Statistics can be calculated that predict problematic variants:
  - Low read depth
  - Strand bias
  - Low mapping quality
  - Clusters of nearby variants
  - Nearby homopolymer run/other repeats

# Variant Quality Scores

- **"Variant quality score"**: These statistics can be combined to derive a score that expresses the confidence in a particular call

# Solutions for calling difficult variants

- Increase coverage
- **Main advice:  Be aware that variant calling is imperfect**
  - SNPs pretty good
  - indels less so

- Investigational approaches:
  - Joint calling in large batches
  - Building custom references for specific difficult-to-catch variants

- Trust, but **verify!**

# 4. Beyond SNPs and small indels

- Algorithms for other variant classes are coming, but still largely investigational:

  - CNVs* and structural variants

  - Larger insertions (>20bp) or deletions (>50bp)

  - Repeat expansions/contractions

  - Transposable elements

# Take home points

# Take home points

Mapping

Cleaning

Variant calling

**Read mapping**

**Removing PCR duplicates**

**Indel realignment**

**Variant calling**

**Variant quality "calibration"**

- A proper analytic pipeline mitigates many of the complications of NGS analysis

# Take home points

- Ask not just about mean coverage, but coverage <u>breadth and depth</u> ("95% coverage at 30X")

- Ask for a list of <u>coverage dropouts</u>. *There is no such thing as a "whole" genome!*

- Weigh the pros and cons of maximizing breadth (exome) vs. depth (panels)

- SNPs are generally high quality, but it is still important to weigh <u>variant quality</u> and other red flags.  Especially for indels, trust but <u>verify</u>

- Recognize that CNV, SV, larger indels, repeat expansion/contractions, and mobile elements are out of the scope of most clinical NGS pipelines
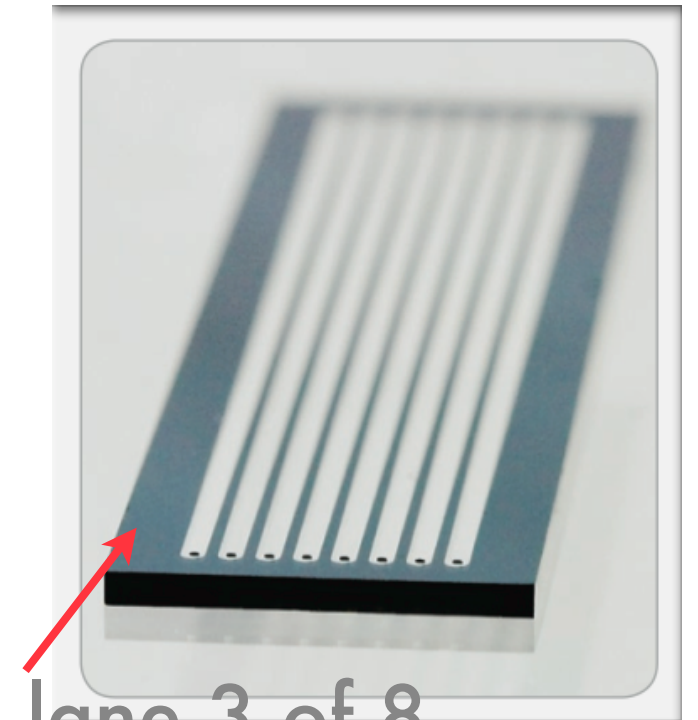
# Questions

# Sequence reads

# Anatomy of a read



lane 3 of 8

identifier

```
@HWI-EAS214_1:3:3:1375:979
CCCAACCAACCCNNCACATCCCAAACAACCCCAACC
+HWI-EAS214_1:3:3:1375:979
25 25 25 25 25 25 25 25 25 25 25 13 -2 -2 25 11 14 25 -2 25 25 4
```

# Anatomy of a read


lane 3 of 8

machine | lane | tile | X:Y

identifier

```
@HWI-EAS214_1:3:3:1375:979
CCCAACCAACCCNNCACATCCCAAACAACCCCAACC
+HWI-EAS214_1:3:3:1375:979
25 25 25 25 25 25 25 25 25 25 25 13 -2 -2 25 11 14 25 -2 25 25 4
```

# Anatomy of a read



lane 3 of 8

sequence →

```
@HWI-EAS214_1:3:3:1375:979
CCCAACCAACCCNNCACATCCCAAACAACCCCAACC
+HWI-EAS214_1:3:3:1375:979
25 25 25 25 25 25 25 25 25 25 25 13 -2 -2 25 11 14 25 -2 25 25 4
```
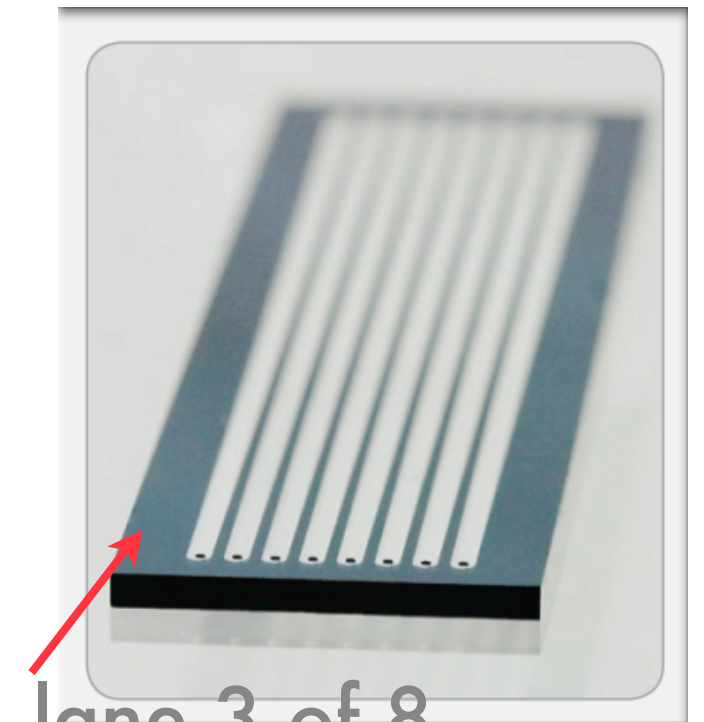
# Anatomy of a read



lane 3 of 8

**identifier (again!)**

```
@HWI-EAS214_1:3:3:1375:979
CCCAACCAACCCNNCACATCCCAAACAACCCCAACC
+HWI-EAS214_1:3:3:1375:979
25 25 25 25 25 25 25 25 25 25 25 13 -2 -2 25 11 14 25 -2 25 25 4
```

# Anatomy of a read


lane 3 of 8

base qualities
(higher=better)

```
@HWI-EAS214_1:3:3:1375:979
CCCAACCAACCCNNCACATCCCAAACAACCCCAACC
+HWI-EAS214_1:3:3:1375:979
25  25  25  25  25  25  25  25  25  25  25  13  -2  -2  25  11  14  25  -2  25  25  4
```
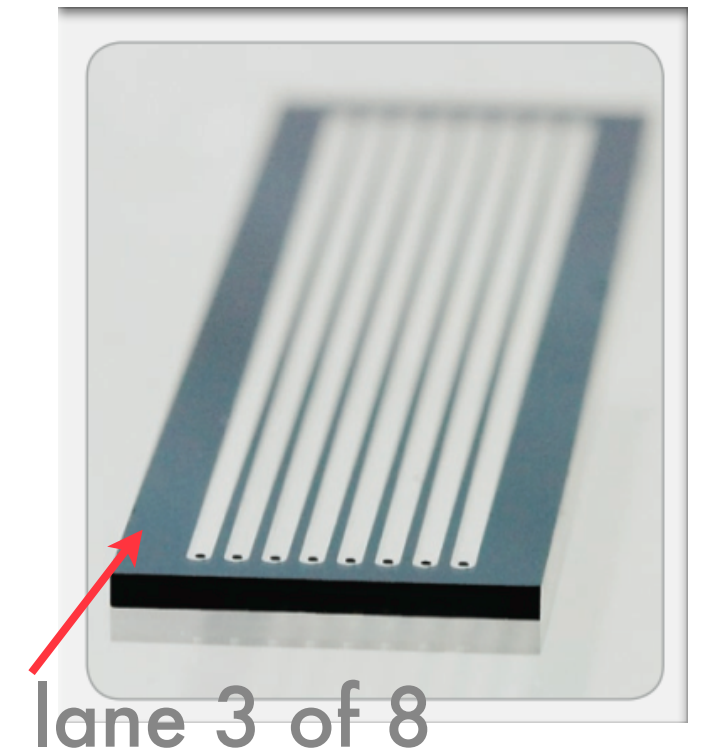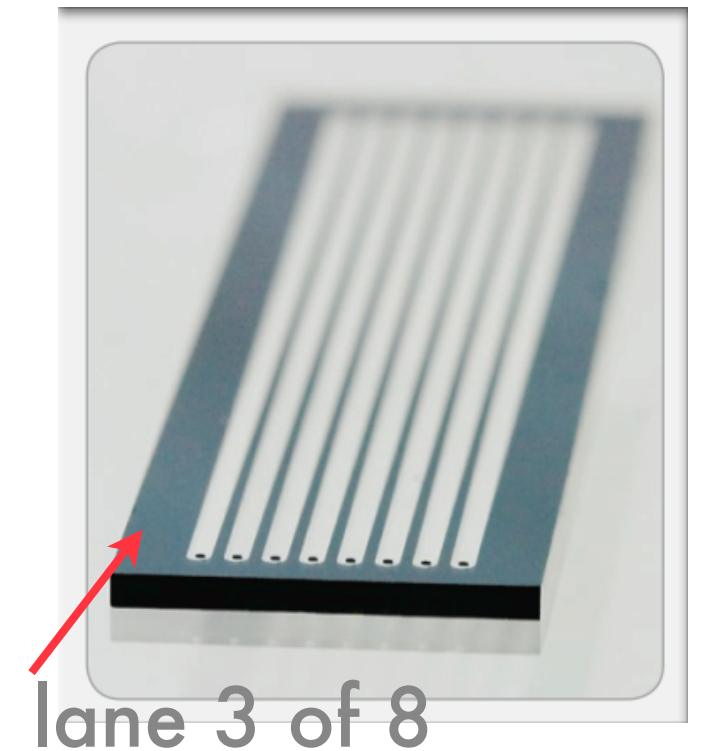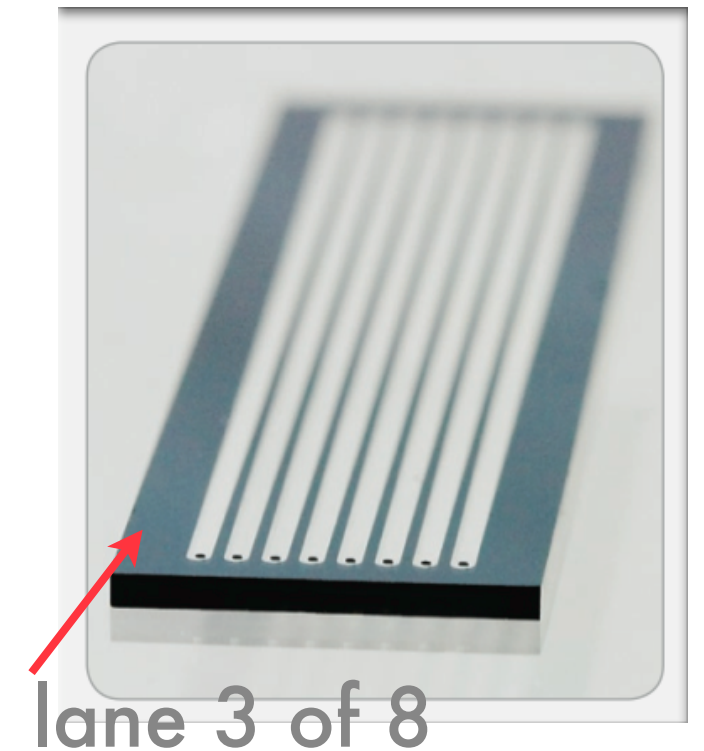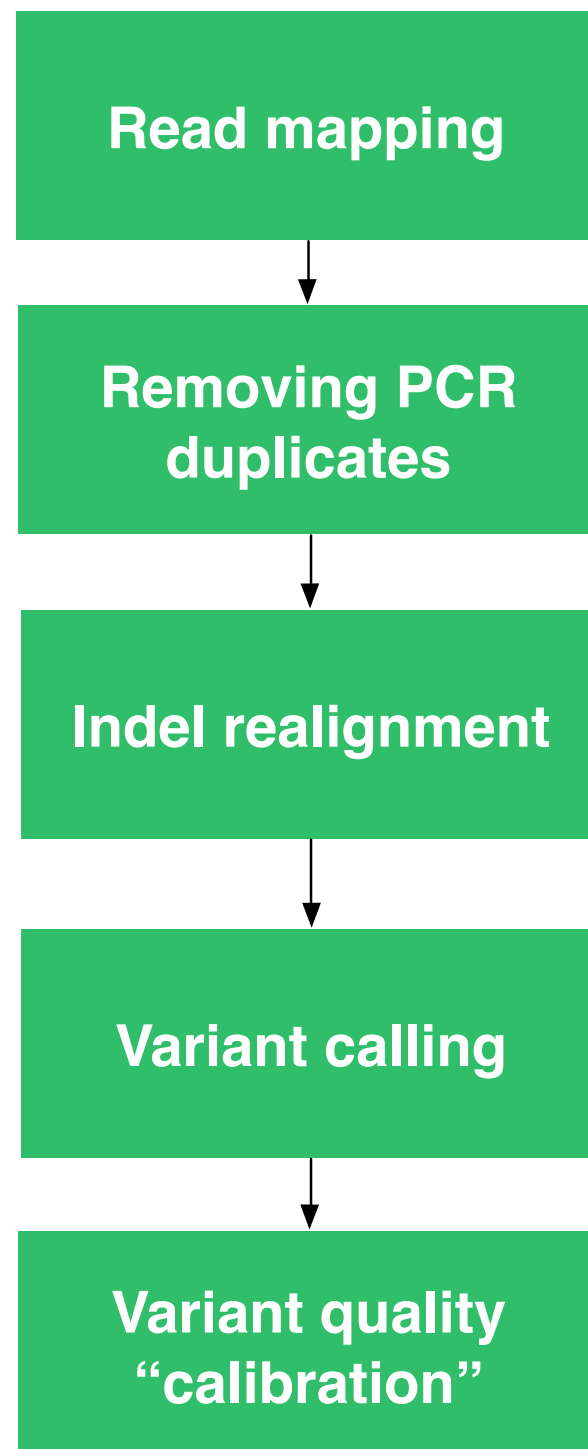
# Questions to ask

**Read mapping**

**Removing PCR duplicates**

**Indel realignment**

**Variant calling**

**Variant quality "calibration"**

- Was sufficient breadth & depth of coverage achieved?
  - "85-95% coverage at >30X"

- What regions were missed?

# Questions to ask

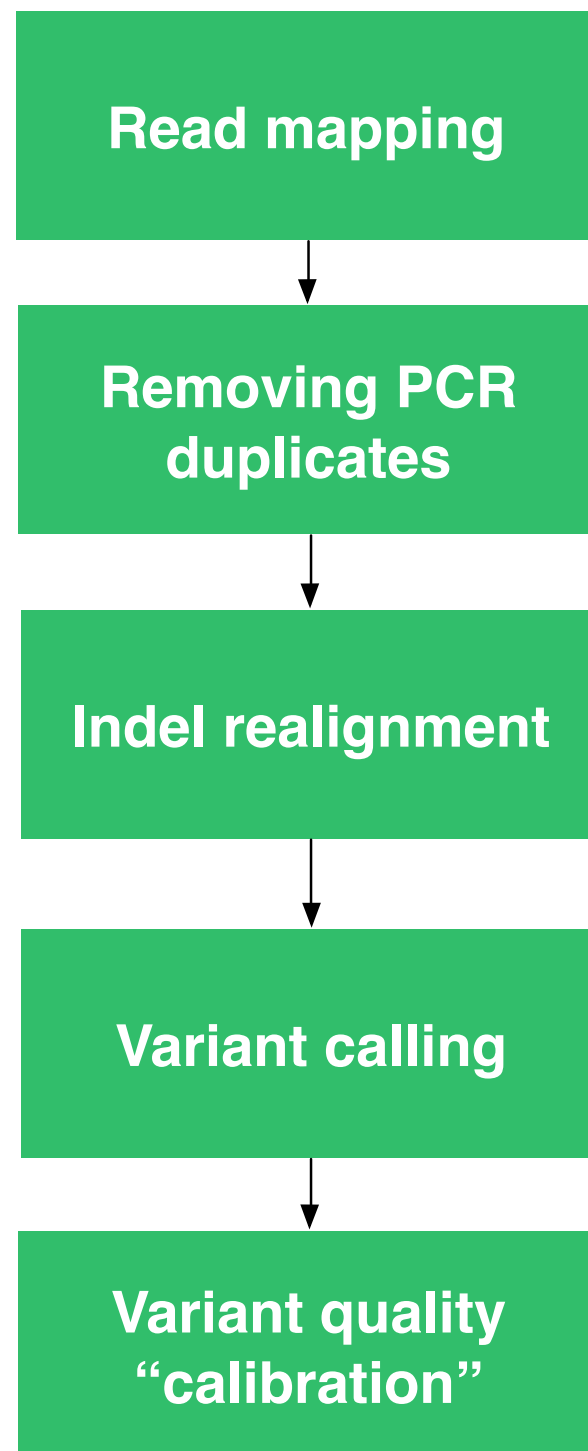Read mapping

Removing PCR duplicates

Indel realignment

- Was appropriate cleaning performed?

Variant calling

Variant quality "calibration"

# Questions to ask

Read mapping

↓

Removing PCR duplicates

↓

Indel realignment

↓

Variant calling

↓

Variant quality "calibration"

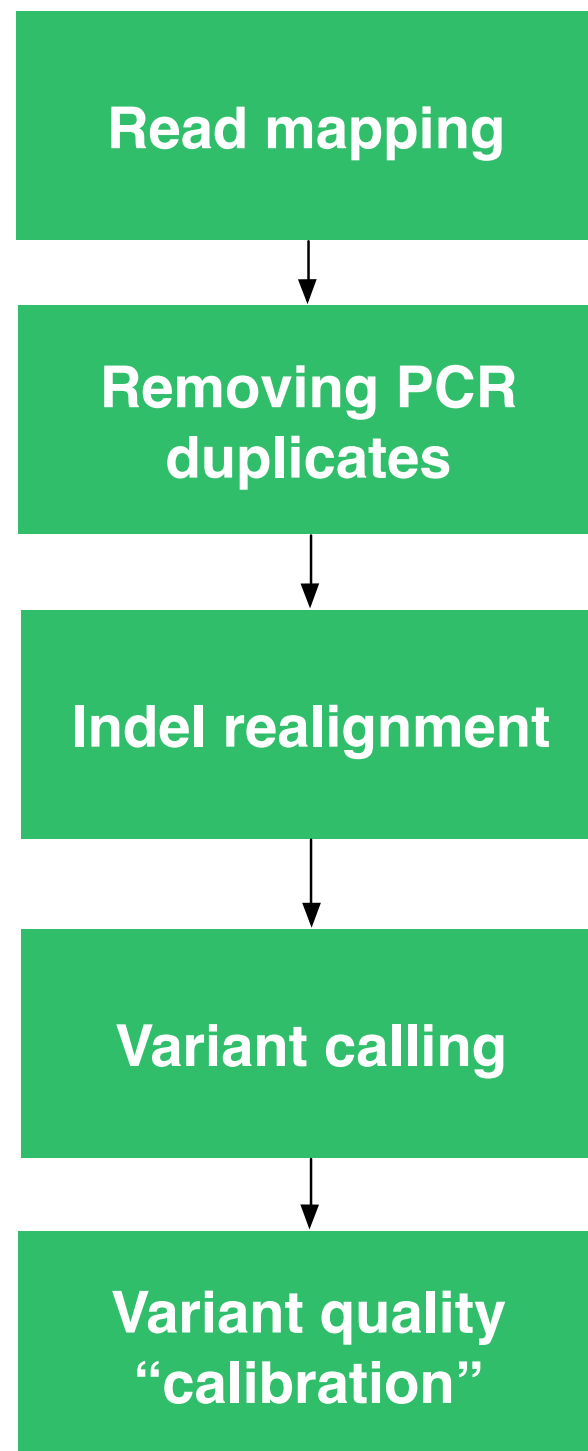- Are the #s of variants called reasonable?
- Especially indels
- Is the percentage of "known SNPs" reasonable (98% in dbSNP)?