# How Effective Are Illumina Methods for BGI-SEQ?

# Illumina - base by base



Genomic DNA → Cut DNA → Add Linkers → *In Situ* PCR → Sequencing → An image of hundreds of extended molecules

# BGI-SEQ (Nanoball) - probe by probe



Reading bases 1-5, e.g. position 5:

Probe ↓ NNNNCNNNN ACTGCTGACGTACTG
......... GCTAATCTGGGATAC TGACGACTGCATGACGC

Standard anchor ↓

Genomic sequence: ..5 4 3 2 1

DNB adaptor/anchor binding site

Common Probes
(5th base set shown):

5 4 3 2 1
NNNNANNNN
NNNNCNNNN
NNNNGNNNN
NNNNTNNNN

Reading bases 6-10, e.g. position 10:

Probe ↓ Degenerate anchor ↓ Standard anchor ↓

NNNNANNNN**NNNNNAC**TGCTGACGTAC
......... GCTAATCTGGGATAC TGACGACTGCATGACGC

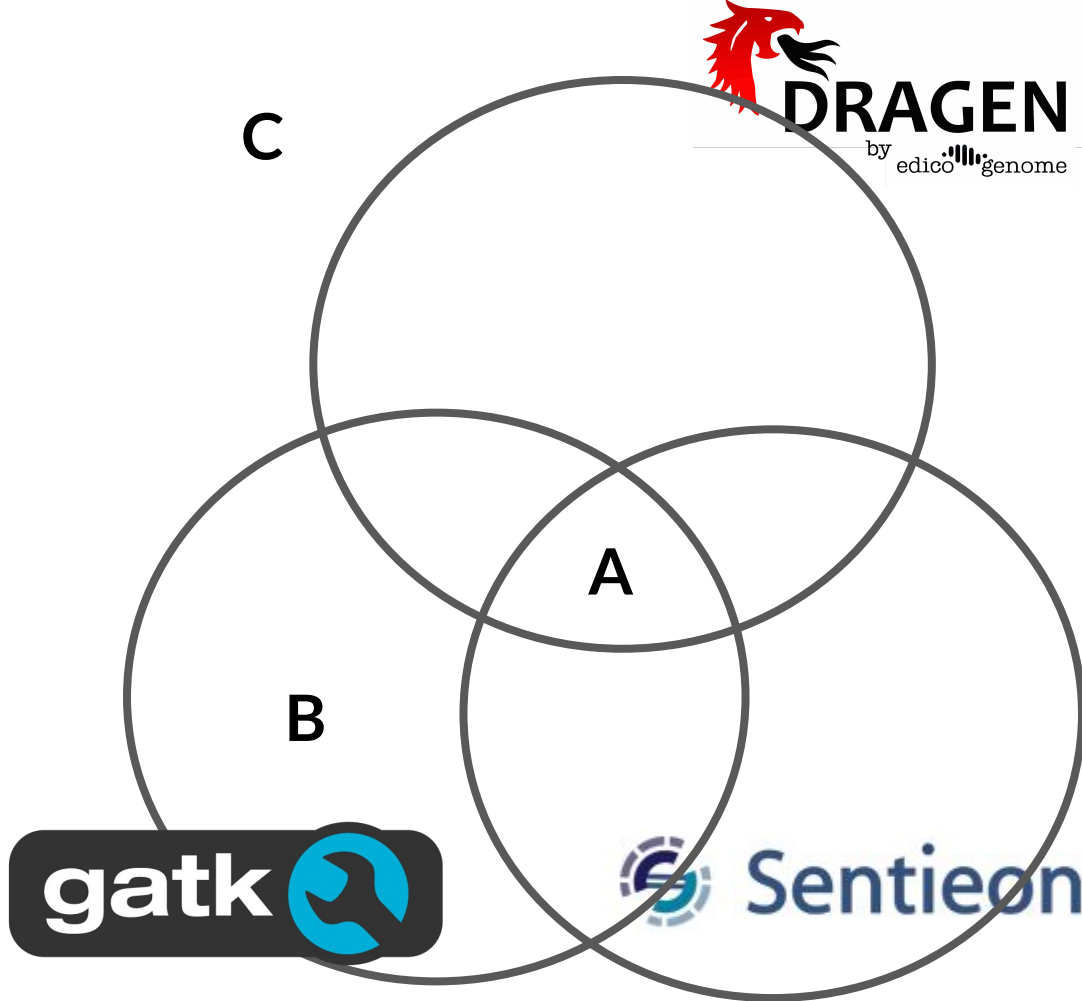Genomic sequence: ..10 9 8 7 6 5 4 3 2 1

DNB adaptor/anchor binding site

# Incredibly Important Resources





# Only 3 Genomes, Not Somatic, Not Bill
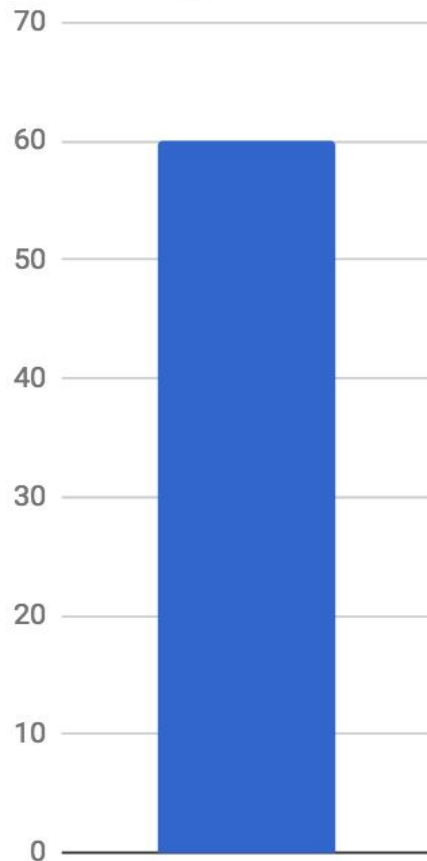
# New Benchmark Method - Drop Out of Mutual Agreement

DOMA-Assessment Germline SNP False Positives

Legend: ● GATK4 ● Sentieon ● DeepVariant ● DRAGEN ● Strelka2 ● Freebayes

X-axis: Fold-Coverage of Sample (BGI-SEQ 500)
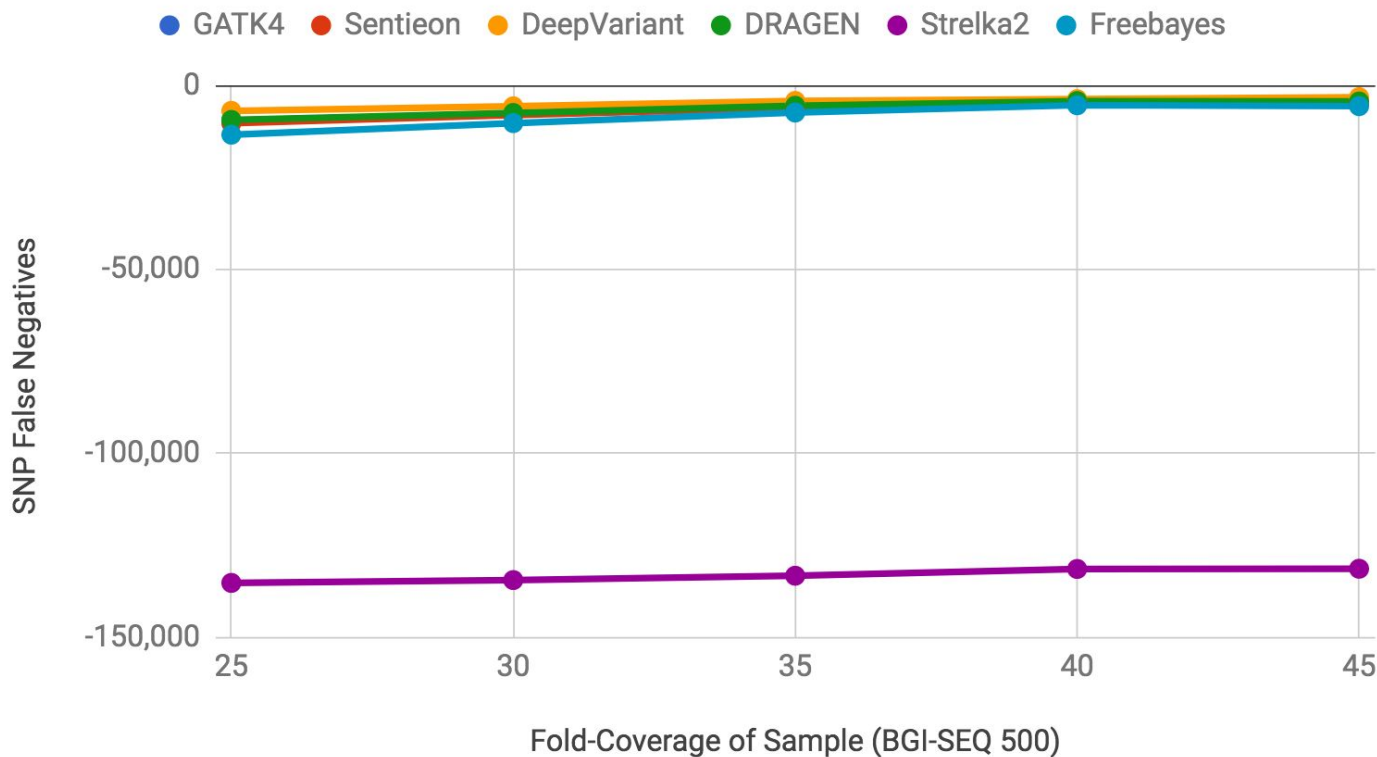Y-axis: SNP False Positives

Example Conclusion : GATK4 and Sentieon are "the same" except Sentieon doesn't downsample and is reproducible

This amount of separation is very weird in Illumina benchmarks. Looks like coverage is more variable in BGI-SEQ
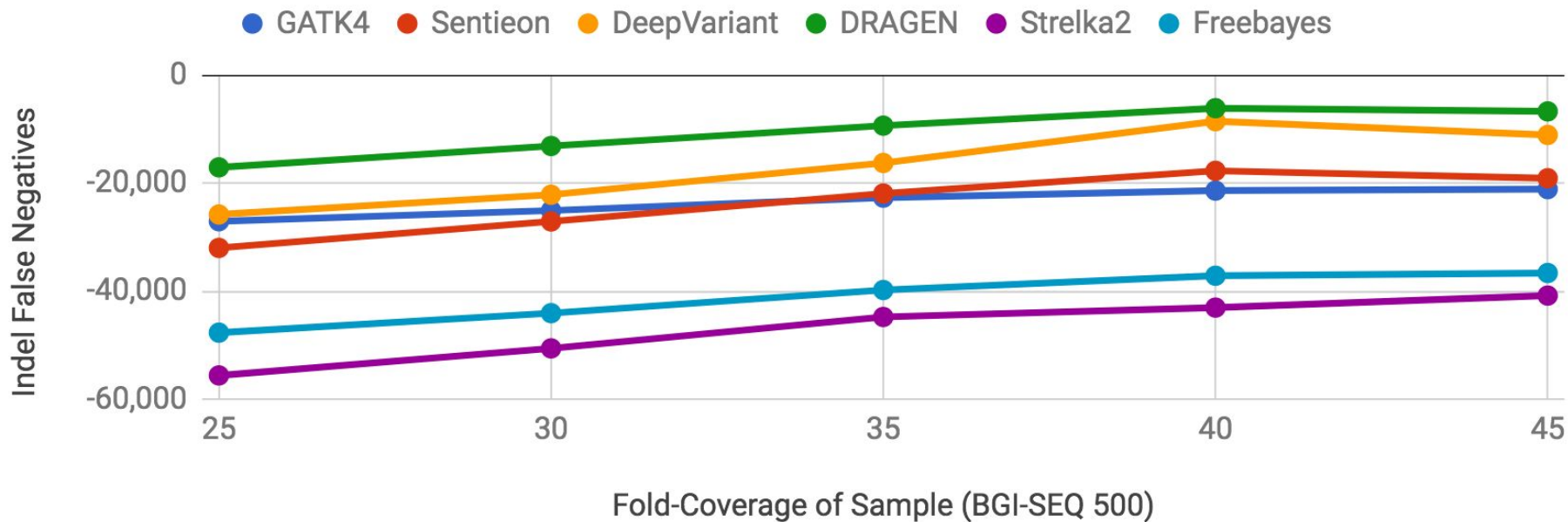
DOMA-Assessment Germline SNP False Negatives

We may have observed a negative interaction with Strelka2 germline heuristics and BGI-SEQ data

In retrospect, It may be better to run BGI-SEQ with --exome parameters even for WGS.

DOMA-Assessment Germline Indel False Negatives

GATK4 • Sentieon • DeepVariant • DRAGEN • Strelka2 • Freebayes

For Indels, it seems all methods have issue with Recall. This may be the dominant error mode. (y-axis is 3x SNP)

# Train DeepVariant for BGI-Seq data

- DeepVariant is trained on data from **Illumina** sequenceers
- How does it perform on **BGI-Seq** data?
- Can we make it better?

# Train DeepVariant for BGI-Seq data

- DeepVariant is trained on data from Illumina sequenceers
- How does it perform on BGI-Seq data?
- Can we make it better? → YES!!
- Indel F1

  94.28% → 98.10%
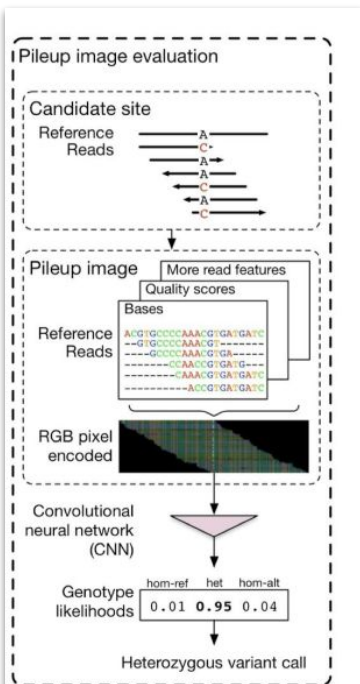- SNP F1:

  99.83% → 99.89%

# Train DeepVariant for BGI-Seq data

- DeepVariant is trained on data from Illumina sequenceers
- How does it perform on BGI-Seq data?
- Can we make it better? → YES!!
- Indel F1:

  94.28% → 98.10%
- SNP F1:

  99.83% → 99.89%
- Step-by-step (you can do it too!) http://bit.ly/train-deepvariant
- Please also check out: https://github.com/google/**nucleus**

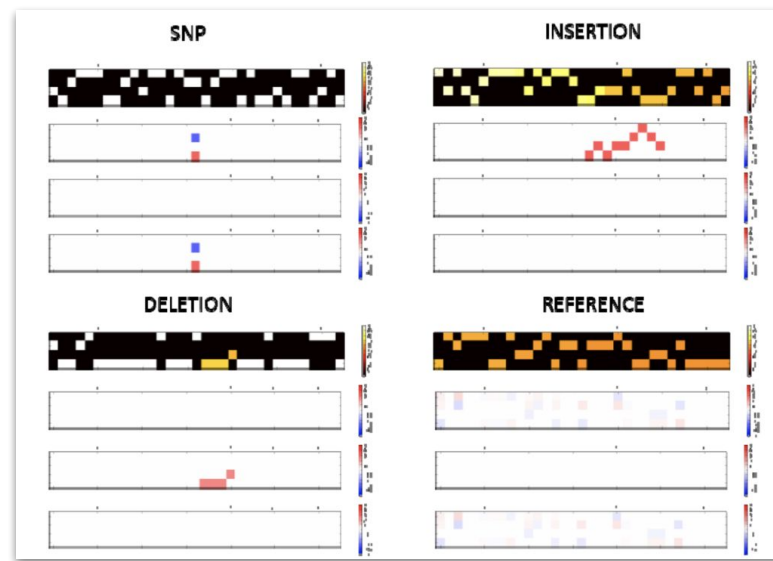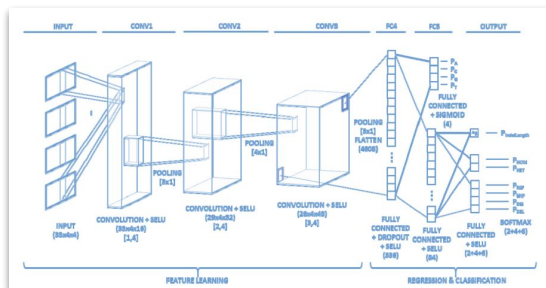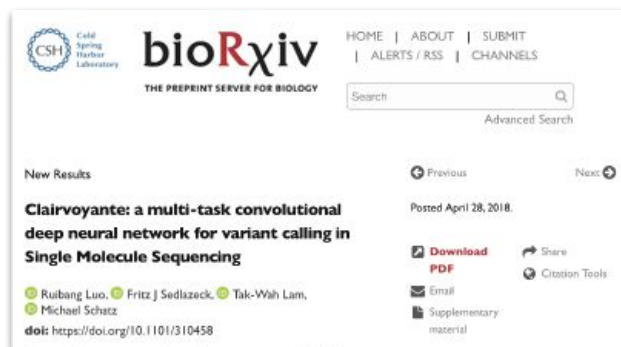  (the library that supports the genomics file formats for DeepVariant!)

# Germline Variant Calling with Various Deep Network Architecture

Different technologies could have different error profiles. Deep learning can help to adapt variant calling for different technologies. We have re-trained DeepVariant and Clairvoyante for BGI-Seq data.
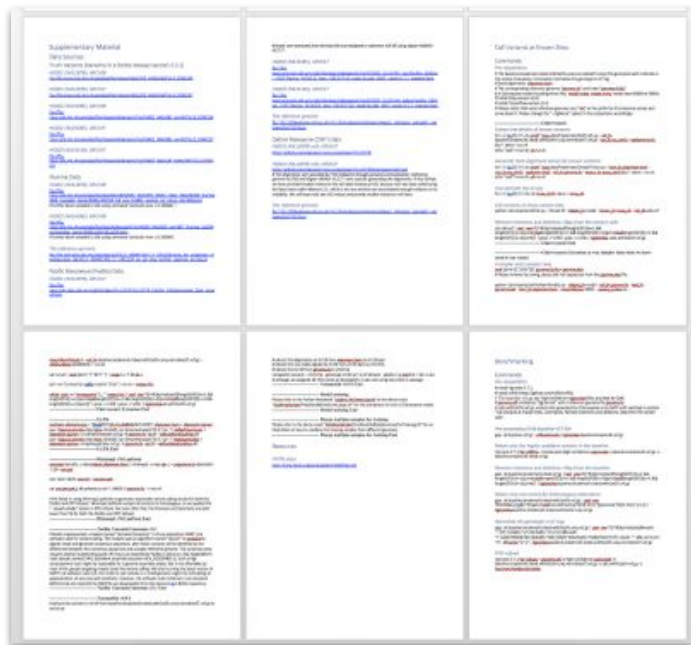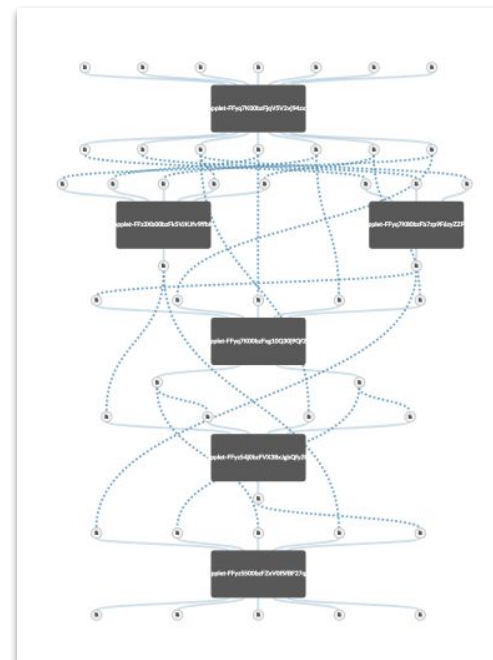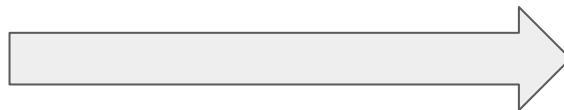
## DeepVariant

## Clairvoyante

# Workflow to create training data for Clairvoyante

"Supplementary Material"



"Crystalize" the workflow/process documented in a "Supplementary Material" to a fully reproducible and re-usable deep-learning-in-cloud workflow.



Results: Two new set deep learning based germline variation calls at 45x and 35x generated
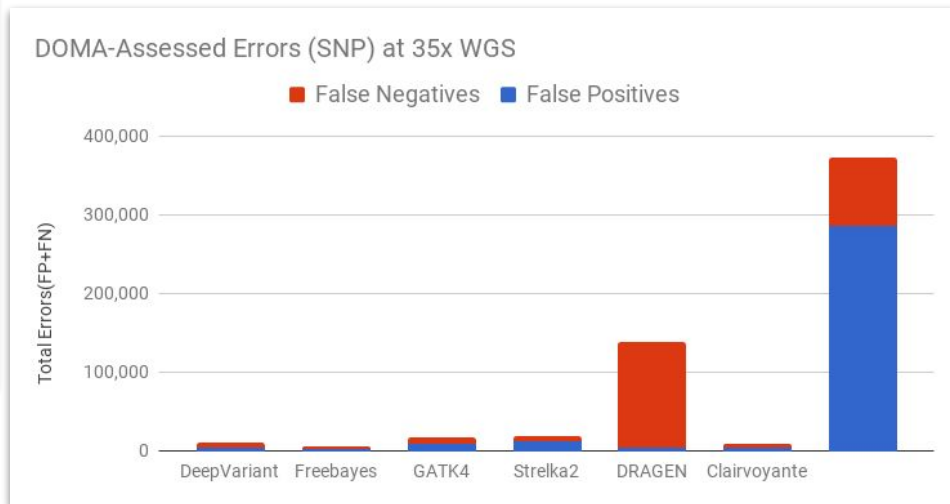
# Quick evaluations

Open source collaboration:

preprint -> reproduce results on a different technology with reproducible workflow in short time.

Making improve the deep learning model easier.

# Conclusion

- When using a new technology, you should always carefully understand your data and how methods perform in it.

- These differences represent gaps that current methods can exploit to improve

- The promise of deep learning based methods allows rapid improvement (over 1 weekend)

Andrew Carroll: https://www.linkedin.com/in/andrewcarrolldna/
Pi-Chuan Chang: https://www.linkedin.com/in/pichuan/
Jason Chin:  https://www.linkedin.com/in/chenshanchin/ Twitter: @infoecho

DNAnexus