

# Benchmarking and Quality Control for genomic variant calling

2025

Johannes Köster

Bioinformatics and Computational Oncology

Institute for AI in Medicine (IKIM)

University of Duisburg-Essen

# Genomic variant calling and why it is challenging

# Seeking for genomic variants

DNA

⌋ ⌋ ⌋ ⌋ ⌋  
⌋ ⌋ ⌋ ⌋ ⌋ ⌋  
⌋ ⌋ ⌋ ⌋ ⌋  
⌋ ⌋ ⌋ ⌋ ⌋

sequencing



sequencing reads

AACCGATTAAACCGGAGTCCCTCGGTAGTTATTTACC  
AACCGGAGTCCCTCGGTAGTTATTTACCCTCTCCGC  
AGTCCCTCGGTAGTTATTTACCCTCTCCGCGTCCTTTC  
ATCCGGAGTCGCAACCGATTAAACCGGAGTCCCT  
GAGTCGCAACCGATTAAACCGGAGTCCCTCGGTAGTTAT

read alignment



aligned reads

AACCGATTAAACCGGAGTCCCGCGGTAGTTATTTACC  
AACCGGAGTCCCGCGGTAGTTATTGACCCTCTCCGC  
AGTCCCTCGGTAGTTATTTACCCTCTCCGCGTCCTTTC  
ATCCGGAGTCCAACCGATTAAACCGGAGTCCCT  
GAGTCGCAACCGATTAAACCGGAGTCCCTCGGTAGTTAT  
...GTAATCCGGAGTCGCAACCGATTAAACCGGAGTCCCGCGGTAGTTATTTACCCTCTCCGCGTCCTTTCTA...

variant calling

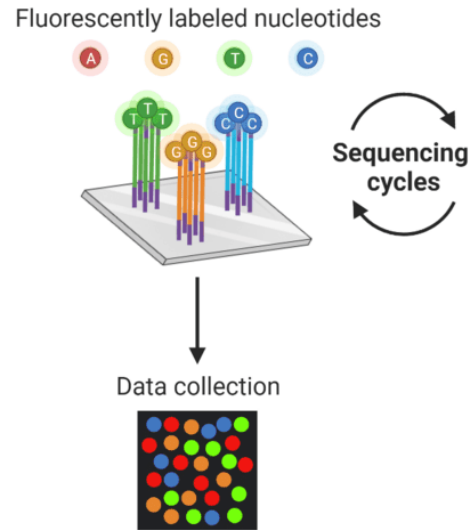


genomic variants

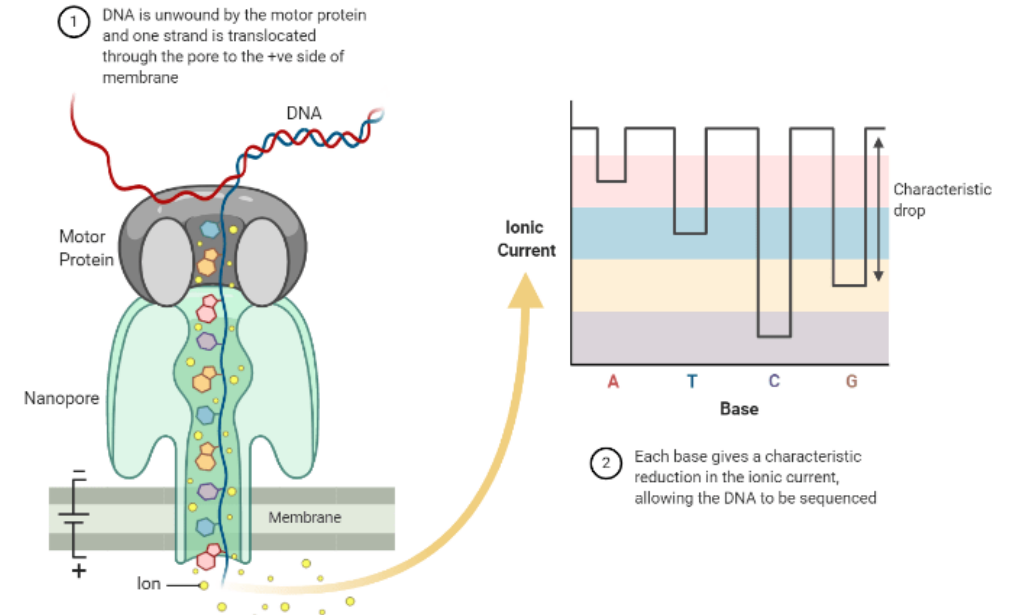
AACCGATTAAACCGGAGTCCCGCGGTAGTTATTTACC  
AACCGGAGTCCCGCGGTAGTTATTGACCCTCTCCGC  
AGTCCCTCGGTAGTTATTTACCCTCTCCGCGTCCTTTC  
ATCCGGAGTCCAACCGATTAAACCGGAGTCCCT  
GAGTCGCAACCGATTAAACCGGAGTCCCTCGGTAGTTAT  
...GTAATCCGGAGTCGCAACCGATTAAACCGGAGTCCCGCGGTAGTTATTTACCCTCTCCGCGTCCTTTCTA...<sup>3</sup>

# Sequencing

## short reads:



## long reads:



## basecalling uncertainty:

posterior probability of incorrect base (base quality)

# Read alignment

```
AACCGATTAACCGGAGTCCCTCGGTAGTTATTTACC
AACCGGAGTCCCTCGGTAGTTATTTACCCTCTCCGC
AGTCCCTCGGTAGTTATTTACCCTCTCCGCGTCCTTTC
ATCCGGAGTCGCAACCGATTAACCGGAGTCCCT
GAGTCGCAACCGATTAACCGGAGTCCCTCGGTAGTTAT
```

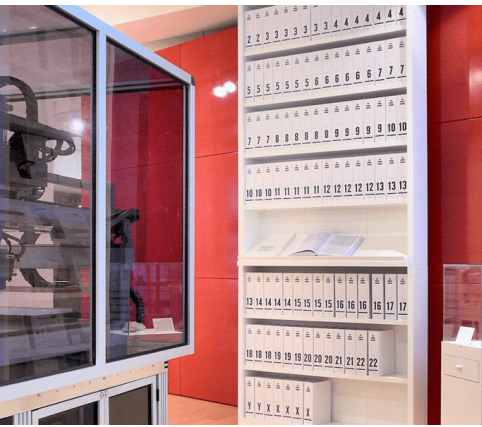
## for each read:

find best position of a short  
text in a very long text  
(alphabet: A,C,G,T)



## challenges:

- repetitive regions
- sequencing errors
- variants



```
AACCGATTAACCGGAGTCCCGCGGTAGTTATTTACC
AACCGGAGTCCCGCGGTAGTTATTGACCCTCTCCGC
AGTCCCTCGGTAGTTATTTACCCTCTCCGCGTCCTTTC
ATCCGGAGTCCCAACCGATTAACCGGAGTCCCTT
GAGTCGCAACCGATTAACCGGAGTCCCTCGGTAGTTAT
...GTAATCCGGAGTCGCAACCGATTAACCGGAGTCCCGCGGTAGTTATTTACCCTCTCCGCGTCCTTTCTA...
```

# Alignment uncertainty

## repetitive regions:

?

?

```
GAGTCGCAACCGATTAACCGGAGTCCCTCGGTAGTTAT      GAGTCGCAACCGATTAACCGGAGTCCCTCGGTAGTTAT
GTAATCCGGAGTCGCAACCGATTAACCGGAGTCCCGCGGTAGTTATTTACCCTCTCCGCGTCCTTTCTAGAGTCGCAACCGATTAACCGGAGTCCCGCGGTAGTTATGGCTGAT...
```

## sequencing errors:

?

?

```
GAGTCGCAACCGATTAACCGGAGTCCCTCGGTAGTTAT      GAGTCGCAACCGATTAACCGGAGTCCCTCGGTAGTTAT
GTAATCCGGAGTCGCAACCGATTAACCGGAGTCCCGCGGTAGTTATTTACCCTCTCCGCGTCCTTTCTAGAGTCGCAACCGATTAACCGGAGTCCCGCGGTAGTTATGGCTGAT...
```

## variants:

?

?

```
GAGTCGCAAC-----AACCGGAGTCCCGCGGTAGTTAT      GAGTCGCAACAACCGGAGTCCCGCGGTAGTTAT
GTAATCCGGAGTCGCAACCGATTAACCGGAGTCCCGCGGTAGTTATTTACCCTCTCCGCGTCCTTTCTAGAGTCGCAACAACCGGAGTCCCGCGGTAGTTATGGCTGAT...
```

# Global benchmarking of genomic variant calling

# Benchmarking genomic variant calling

## **Given:**

- a plethora of methods to choose from, each with a plethora of parameters
- continuous innovations and method improvements

## **Thus:**

- take (real) datasets with known truth
- fight overfitting by using many benchmark datasets
- continuous benchmarking!!



# Exemplary benchmarks

# Genome in a Bottle



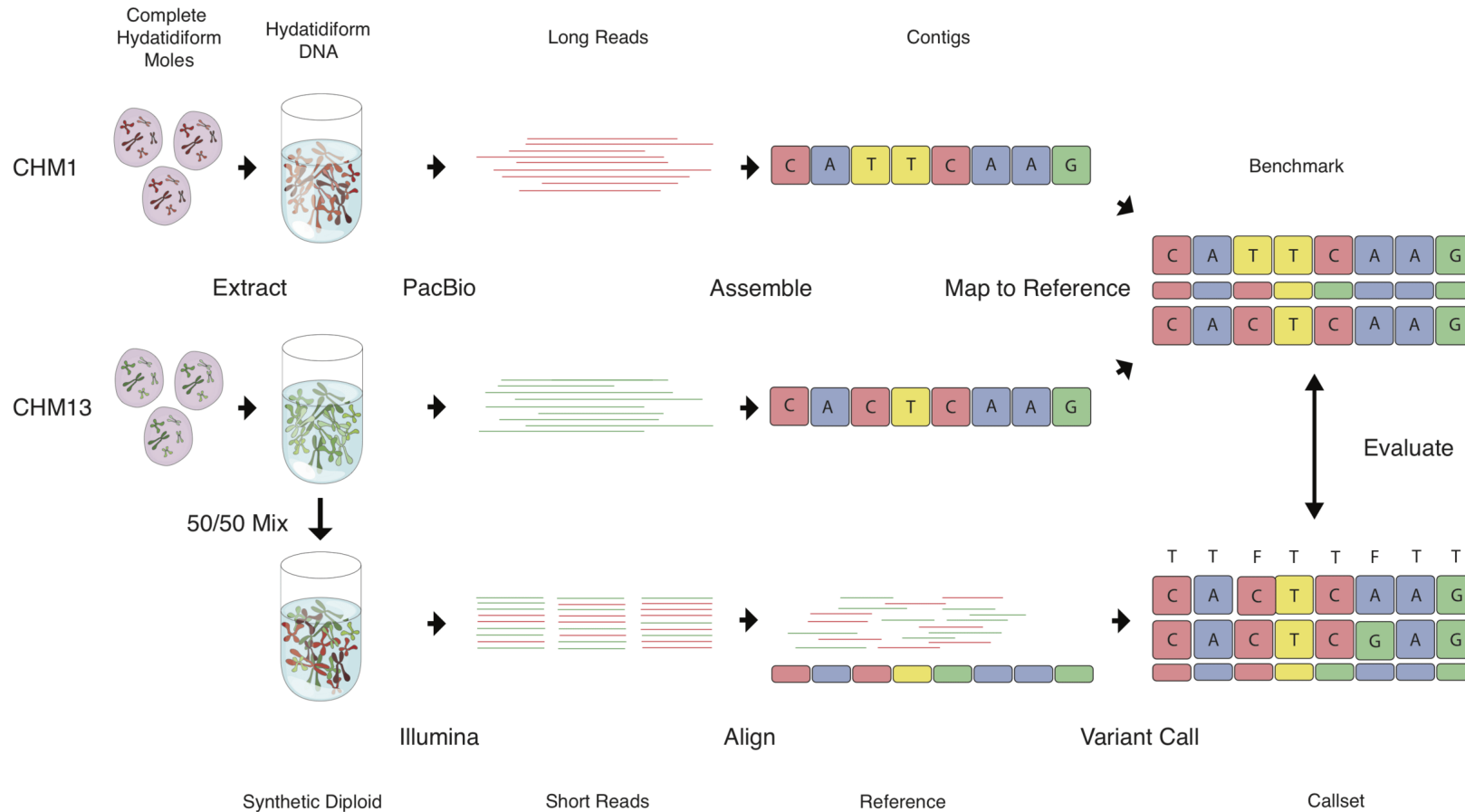
- NIST-lead consortium for generating benchmark datasets for genomic variant calling
- NA12878/HG001 (well-studied sample from the HapMap project)
- three further sets of samples from family pedigrees
- Snakemake-based pipeline for consensus HQ variant calls from multiple technologies and callers

<https://www.nist.gov/programs-projects/genome-bottle>

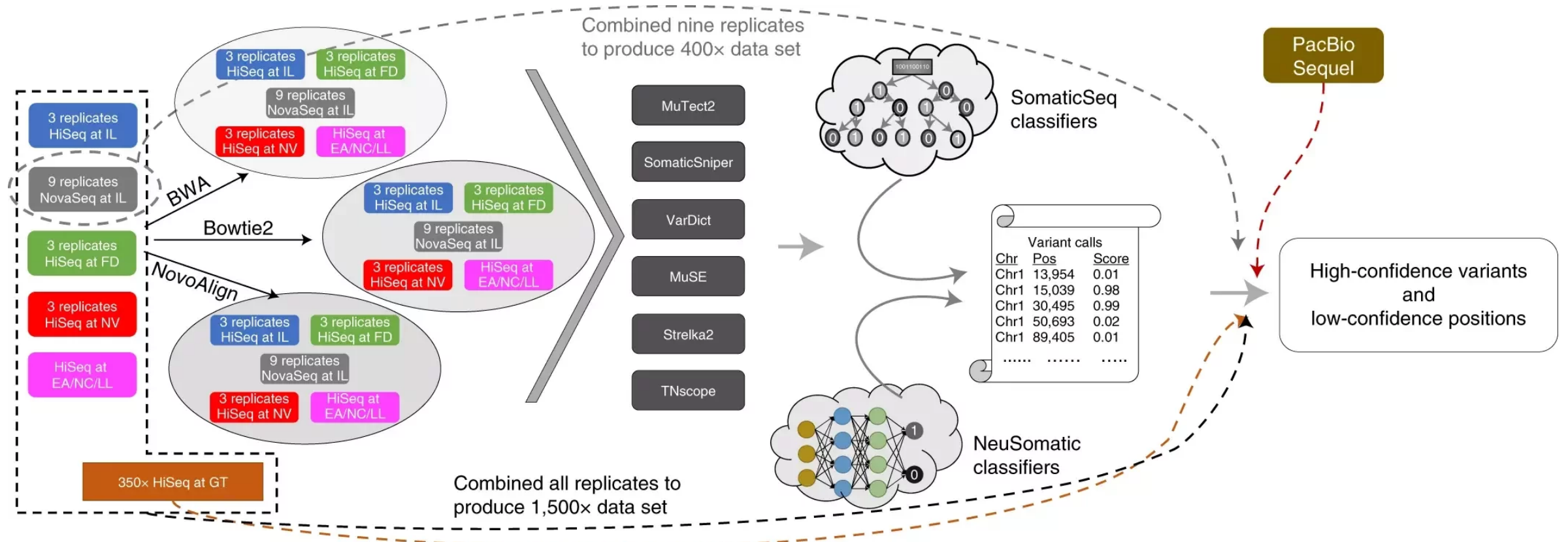
<https://github.com/usnistgov/defrabb>

# CHM-eval

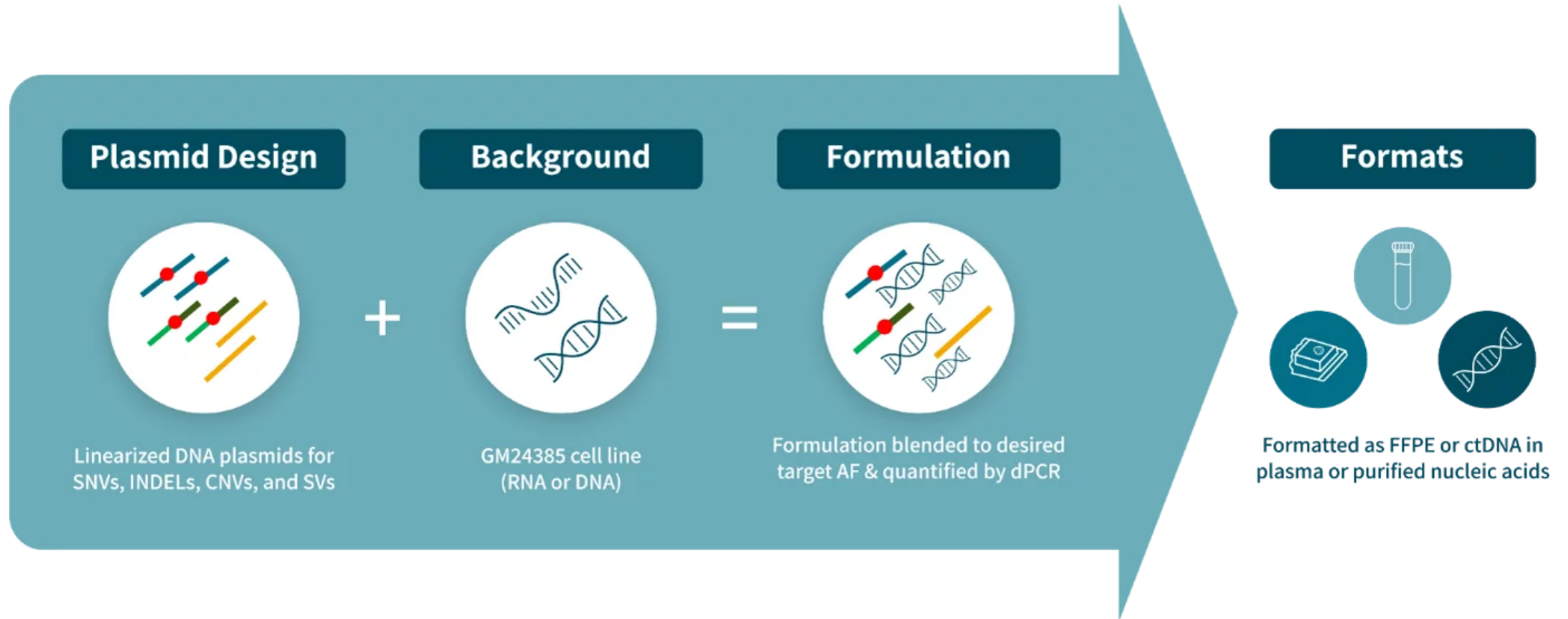
Synthetic, pseudo-diploid benchmark  
sample as mixture of two haploid cell lines  
(CHM13, CHM1)



# SEQC2



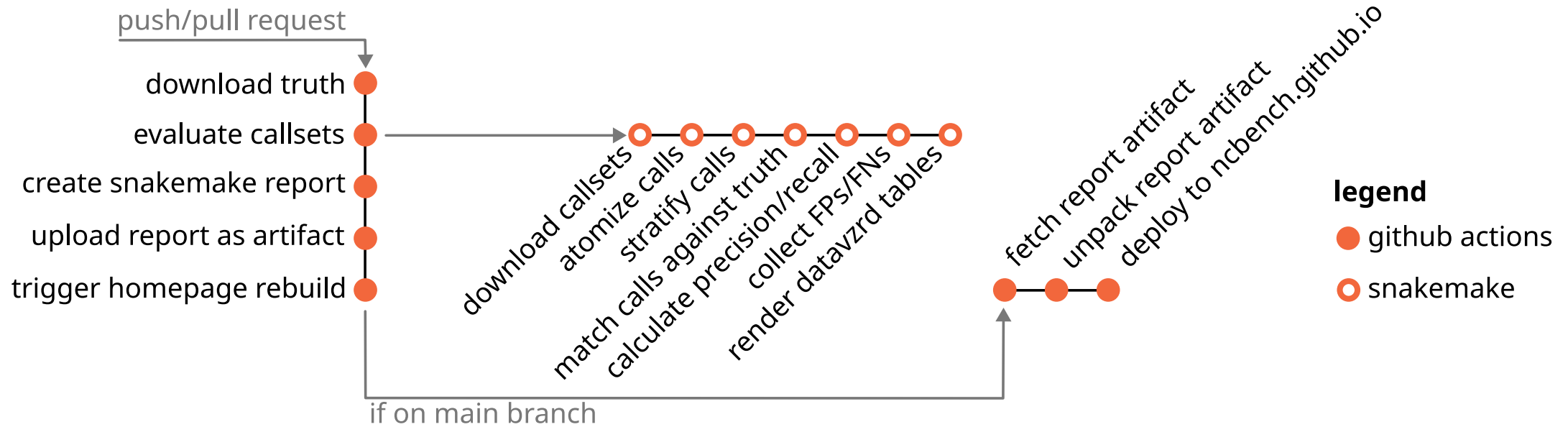
# Seracare "somatic cancer" reference material



# NCBench



# NCbench in a nutshell



F1000Research

F1000Research 2024, 12:1125 Last updated: 13 NOV 2024



RESEARCH ARTICLE

**REVISED** NCBench: providing an open, reproducible, transparent, adaptable, and continuous benchmark approach for DNA-sequencing-based variant calling

[version 2; peer review: 2 approved]

Friederike Hanssen<sup>1</sup>, Gisela Gabernet<sup>1</sup>, Famke Bäuerle<sup>1-4</sup>, Bianca Stöcker<sup>5</sup>, Felix Wiegand<sup>5</sup>, Nicholas H. Smith<sup>6</sup>, Christian Mertes<sup>6-8</sup>, Avirup Guha Neogi<sup>9</sup>, Leon Brandhoff<sup>9,10</sup>, Anna Ossowski<sup>9</sup>, Janine Altmueller<sup>9,11,12</sup>, Kerstin Becker<sup>9</sup>, Andreas Petzold<sup>13</sup>, Marc Sturm<sup>14</sup>, Tyll Stöcker<sup>15</sup>, Sugirthan Sivalingam<sup>16</sup>, Fabian Brand<sup>17</sup>, Axel Schmidt<sup>18</sup>, Andreas Bunes<sup>19</sup>, Alexander J. Probst<sup>20</sup>, Susanne Motamenv<sup>9,10</sup>, Johannes Köster<sup>5,21</sup>

## Datasets

- GIAB
- CHM-eval
- SEQC2
- Seracare (upcoming)

## Key technologies:

- Snakemake
- Datavzrd
- Github Actions
- Github Pages

<https://snakemake.github.io>

<https://datavzrd.github.io>

# Representing and comparing variants

1      10167      .      CCTAACCTAACCTA   CCTTAACCTTAACCTA

CCTAACCC TAACCTA  
CCTTAACCTTAACCTA

CC TAACCCTAACCTA  
CCTTAACCTTAACCTA



1	10167	CCTAACCCTAA	C
1	10170	A	T
1	10172	C	A
1	10174	C	CT

1	10167	CCTAACCCTAA	C
1	10168	C	CT
1	10174	C	T



- apply bcftools norm with same version to truth and callset
- use rtg vcfeval for determining TP, FP, FN

<https://samtools.github.io/bcftools>

<https://github.com/RealTimeGenomics/rtg-tools>



# Calling precision and recall

## **recall:**

should go down with weaker evidence

## **precision:**

should be independent of evidence strength

## **Hence:**

stratification by read depth



- determine read depth with mosdepth
- stratify truth and callsets with bedtools

<https://github.com/brentp/mosdepth>

<https://bedtools.readthedocs.io>

# Genotyping precision and recall

For genotyping, precision may decrease with lower read depth (because it is MLE), for calling not.

► Genotyping and calling should be analyzed separately.

<https://ncbench.github.io>

# Local benchmarking

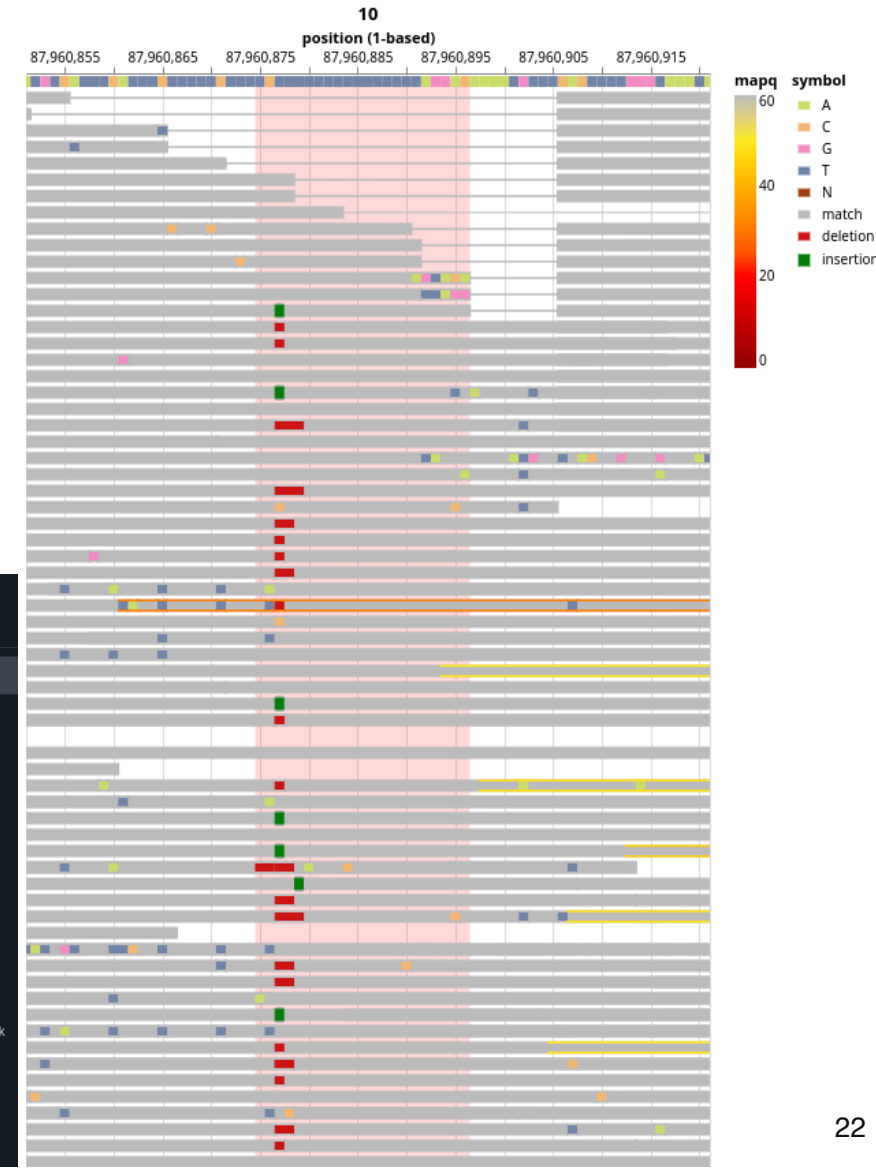
# Issues with global benchmarks

- Truth can be biased towards used tools
- Limited set of confidence regions
- At individual loci, performance may differ a lot from a global measure
- Expensive to compute ▶ rarely repeated

# Good complement: local benchmarks

- individual variants *or non-variants*
- truth is known by prior knowledge, orthogonal validation, or manual inspection
- fast to test
- can be embedded in continuous integration platforms like github actions

```
Summary
Jobs
  Formatting
  Testing
Run details
  Usage
  Workflow file
Testing
succeeded on Mar 5 in 5m 19s
  Run cargo test
    test test_giab_20_exact_mode ... ok
    test test_giab_27_exact_mode ... ok
    test test_giab_22_exact_mode ... ok
    test test_giab_29_exact_mode ... ok
    test test_giab_31_exact_mode ... ok
    test test_giab_33_exact_mode ... ok
    test test_giab_30_exact_mode ... ok
    test test_giab_32_exact_mode ... ok
    test test_haplotype_absent_exact_mode ... ok
    test test_haplotype_present_exact_mode ... ok
    test test_haplotype_singleton_exact_mode ... ok
    test test_giab_35_exact_mode ... ok
    test test_imprecise_fusion_absent_exact_mode ... ok
    test test_imprecise_fusion_exact_mode ... ok
    test test_giab_25_exact_mode ... ok
    test test_hiv_vaf_higher_than_expected_exact_mode ... ok
    test test_l2fc_exact_mode ... ok
    test test_long_pattern_exact_mode ... ok
    test test_low_cov_vaf_exact_mode ... ok
    test test_long_pattern_fast_mode ... ok
    test test_meth_candidates1 ... ok
    test test_mendelian_prior_exact_mode ... ok
```



# Automatic testcase generation in Varlociraptor

Automatic test case generation:

```
varlociraptor call variants --testcase-prefix testcase --testcase-locus CHROM:POS generic \  
--scenario scenario.yaml --obs tumor=tumor.bcf normal=normal.bcf
```

**145**

public testcases  
(simulated + real  
benchmarks)

**66**

private testcases  
(clinical)

Köster et al. Genome Biology 2020

Lähnemann et al. Nature Communications 2021<sup>281</sup>

# Conclusion

- Genomic variant calling is still a hard problem
- Global benchmarks help to understand performance of method and parameter combinations
- As methods evolve fast, it is crucial to apply them continuously
- Combination of Github/Snakemake/Datavzrd makes that easy
- Local benchmarks complement global ones by being faster and specific for problematic regions

## Acknowledgements

Susanne Motameny, Gisela Gabernet, Friederike Hansen, Famke Bäuerle, Bianca Stöcker, Hamdiye Uzuner, Adrian Prinz, Felix Mölder, Felix Wiegand, Nicholas H. Smith, Christian Mertes, Avirup Guha Neogi, Leon Brandhoff, Anna Ossowski, Janine Altmueller, Kerstin Becker, Andreas Petzold, Marc Sturm, Tyll Stöcker, Sugirthan Sivalingam, Fabian Brand, Axel Schmidt, Andreas Bunniss, Alexander J. Probst, Narci Kübra, Anna Bennet-

Pages

