

vcfdist: accurately benchmarking phased variant calls

Tim Dunn

PhD Candidate

University of Michigan

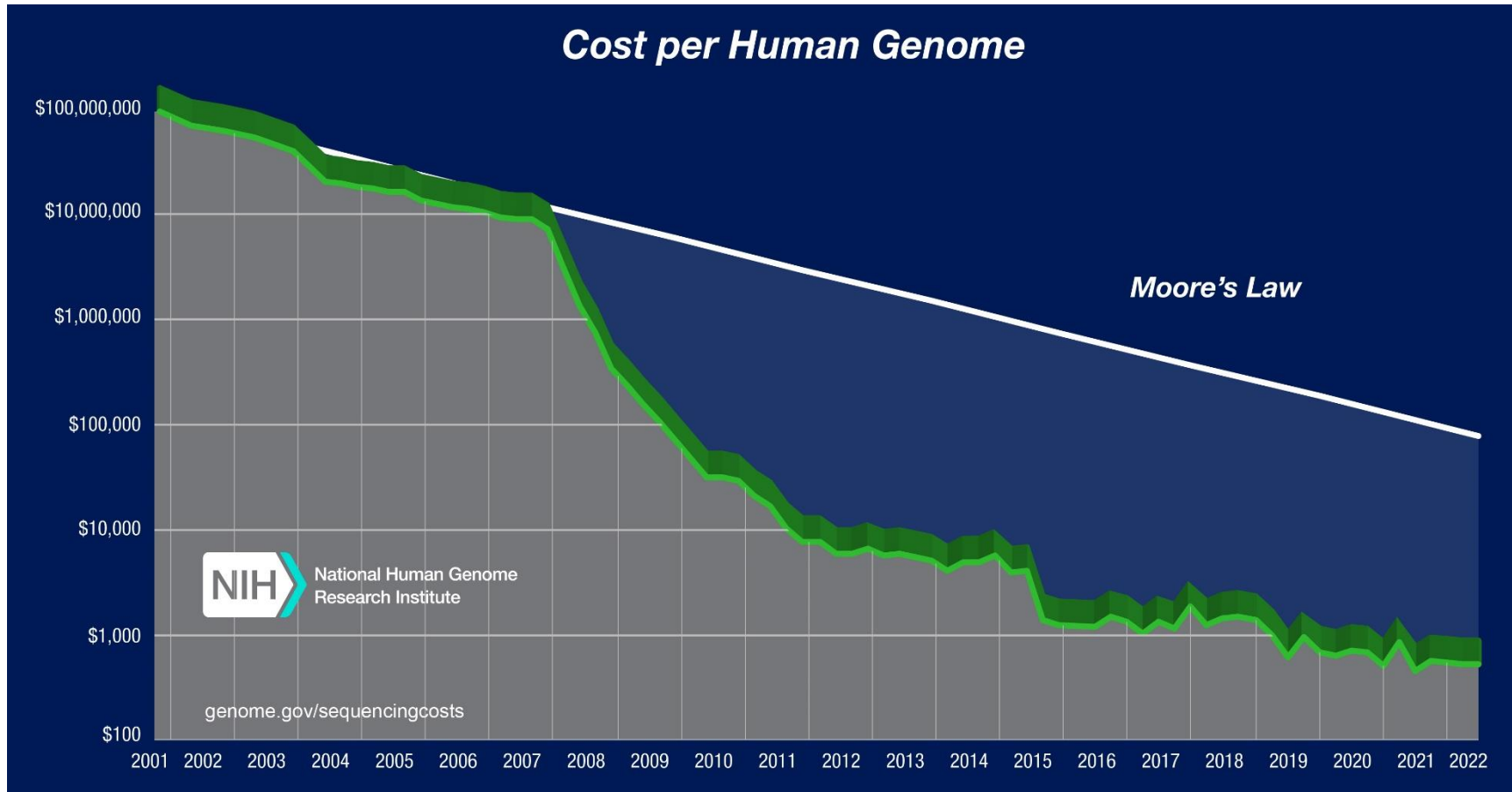
Outline

- 1. Background**
- 2. Problem #1**
- 3. Solution**
- 4. Results**
- 5. Problem #2**
- 6. Solution**
- 7. Results**

Outline

1. **Background:** whole genome sequencing evaluation
2. **Problem #1**
3. **Solution**
4. **Results**
5. **Problem #2**
6. **Solution**
7. **Results**

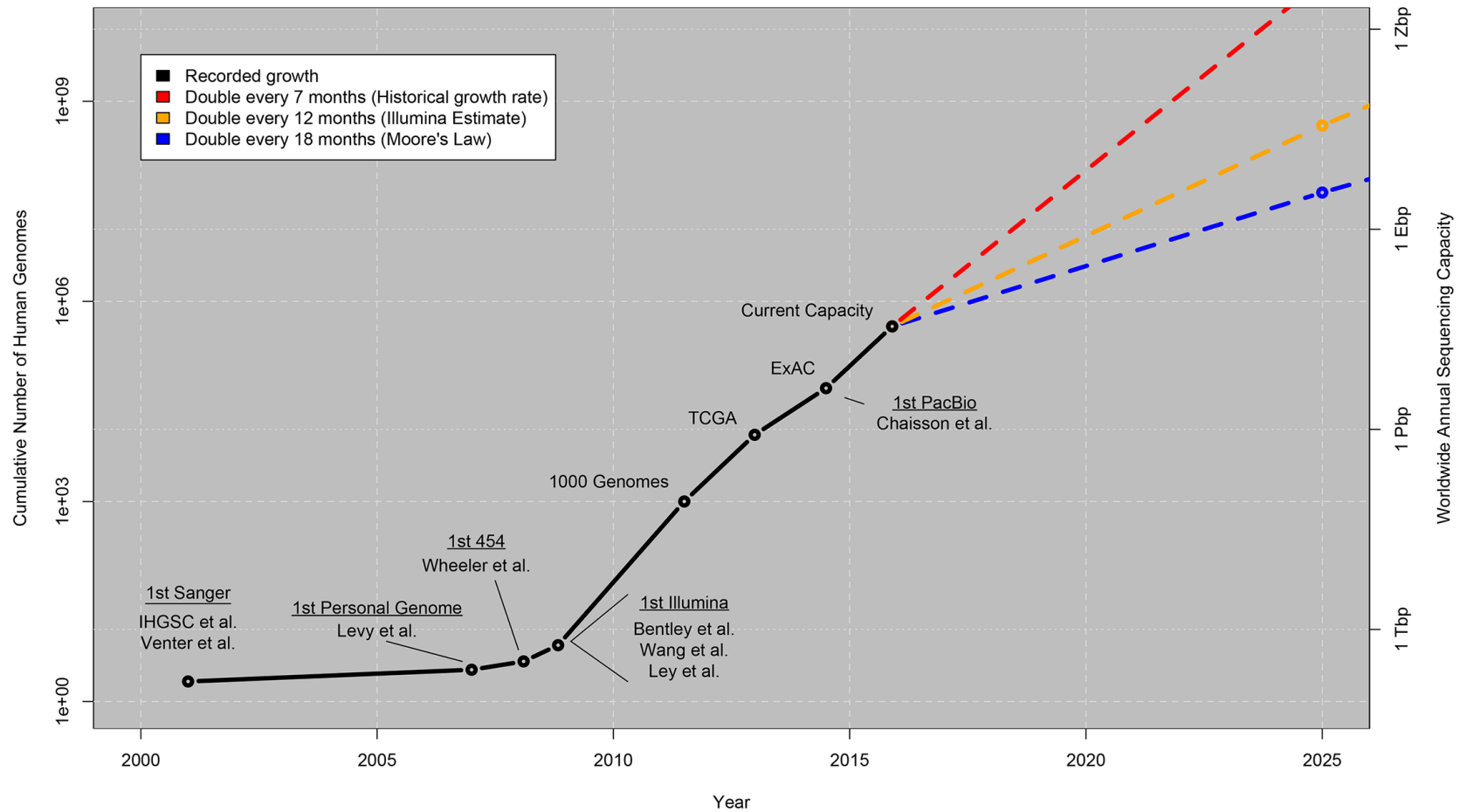
Sequencing: *cost is rapidly declining*



NHGRI. "DNA Sequencing Costs: Data". <https://www.genome.gov/about-genomics/fact-sheets/DNA-Sequencing-Costs-Data>, 2023.

Sequencing: *exponential growth in genomes*

Growth of DNA Sequencing



Stephens et al. "Big Data: Astronomical or Genomical?". PLOS Biology, 2015.

Applications

- Genome wide association studies
- Pharmacogenomics
- Clinical diagnostics
- Benchmarking new methods and tech

Applications: *genome comparison required*

- Genome wide association studies
- Pharmacogenomics
- Clinical diagnostics
- Benchmarking new methods and tech

Comparison: *genomes are mostly identical*

Reference:

ACCGTTGAAGGACGGCCATTTTTT AACTGAGCATCCATCTAAAAGCCTTTTAGCGGCGCCCTCTATAGAT

Query #1:

ACCGTTGAAGGACGGCCA TTTTTAAACTGAGCATCCATCTAAAAGCCTTTT



Variant Call Format: *difference-based*

Reference:

ACCGTTGAAGGACGGCCATTTTTT AACTGAGCATCCATCTAAAAGCCTTTTAGCGGCGCCCTCTATAGAT

Query #1:

ACCCTTGAAGGACGGCCA TTTTAAACTGAGCATCCATCTAAAAGCCTTTT



POSITION	REFERENCE	ALTERNATE
4	G	C
18	AT	A
25	T	TA
53	TAGCGGCGCCC...	T

Applications: *benchmarking*

- Genome wide association studies
- Pharmacogenomics
- Clinical diagnostics
- **Benchmarking new methods and tech**

Applications: *benchmarking*

- Genome wide association studies
- Pharmacogenomics
- Clinical diagnostics
- **Benchmarking new methods and tech**



Benchmarking: *a simple example*

Technology #1

Reference: ACCGTTGAAG

Query #1: ACAGTAGAAG

CHROM	POS	REF	ALT
chr14	3	C	A
chr14	6	T	A

Technology #2

Reference: ACCGTTGAAG

Query #2: ACCGTAGAGG

CHROM	POS	REF	ALT
chr14	6	T	A
chr14	9	A	G

Benchmarking: *a simple example*

Technology #1

Reference: ACCGTTGAAG

Query #1: ACAGTAGAAG

CHROM	POS	REF	ALT
chr14	3	C	A
chr14	6	T	A

Technology #2

Reference: ACCGTTGAAG

Query #2: ACCGTAGAGG

CHROM	POS	REF	ALT
chr14	6	T	A
chr14	9	A	G

Ground Truth

Reference: ACCGTTGAAG

Query: ACCGTAGAGG

CHROM	POS	REF	ALT
chr14	6	T	A
chr14	9	A	G

Benchmarking: *a simple example*

Technology #1

Reference: ACCGTTGAAG

Query #1: ACAGTAGAAG

CHROM	POS	REF	ALT
 chr14	3	C	A
 chr14	6	T	A

Technology #2

Reference: ACCGTTGAAG

Query #2: ACCGTAGAGG

CHROM	POS	REF	ALT
 chr14	6	T	A
 chr14	9	A	G

Ground Truth

Reference: ACCGTTGAAG

Query: ACCGTAGAGG

CHROM	POS	REF	ALT
chr14	6	T	A
chr14	9	A	G

Benchmarking: *stratification by variant type*

Reference:

ACCGTTGAAGGACGGCCATTTTTT AACTGAGCATCCATCTAAAAGCCTTTTAGCGGCGCCCTCTATAGAT

Query #1:

ACCCTTGAAGGACGGCCA TTTTTAAACTGAGCATCCATCTAAAAGCCTTTT

SNP

*single nucleotide
polymorphism*

substitution

INDEL

insertion/deletion

<50 basepairs

SV

structural variant

50+ basepairs

Benchmarking: *small variants only*

Reference:

ACCGTTGAAGGACGGCCATTTTTT AACTGAGCATCCATCTAAAAGCCTTTTAGCGGCGCCCTCTATAGAT

Query #1:

ACCCTTGAAGGACGGCCA TTTTTAAACTGAGCATCCATCTAAAAGCCTTTT

SNP

*single nucleotide
polymorphism*

substitution

INDEL

insertion/deletion

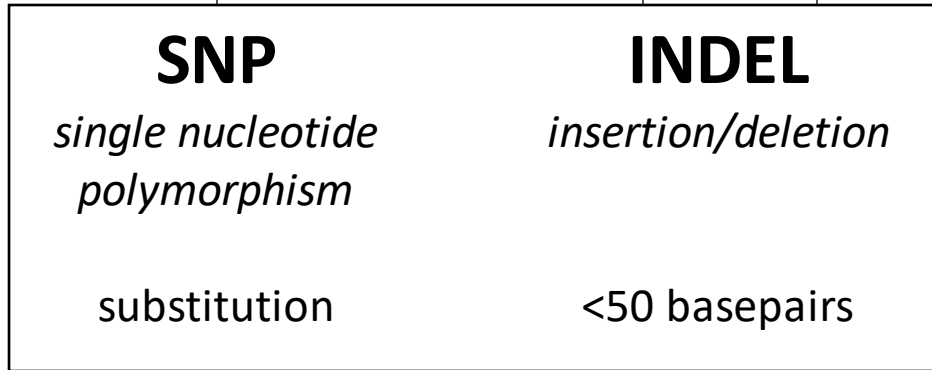
<50 basepairs



Benchmarking: *precision-recall curves*

Reference:
Query #1:

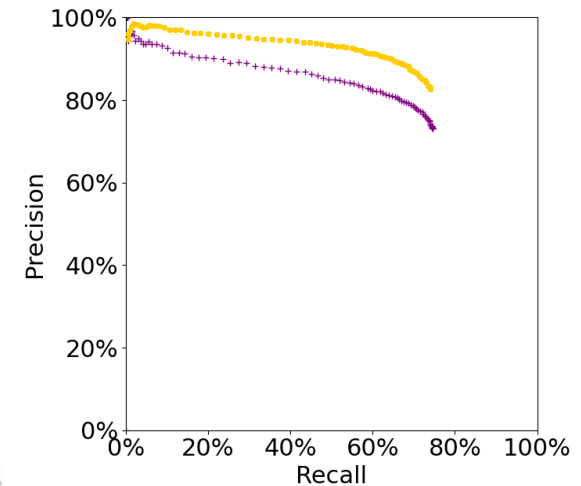
ACCGTTGAAGGACGGCCATTTTTT AACTGAGCATCCATCTAAAAGCCTTTTAGCGGCGCCCTCTATAGAT
ACCGTTGAAGGACGGCCA TTTTTAAACTGAGCATCCATCTAAAAGCCTTTT



$$Precision = \frac{TP}{TP + FP}$$

$$Recall = \frac{TP}{TP + FN}$$

$$F1 = 2 \cdot \frac{precision \cdot recall}{precision + recall}$$



Outline

1. **Background:** whole genome sequencing evaluation
2. **Problem #1:** complex variants
3. **Solution**
4. **Results**
5. **Problem #2**
6. **Solution**
7. **Results**

Problem: *variant representation matters*

Reference:

ACCGTTGAAGGACGGCCATTTTTT AACTGAGCATCCATCTAAAAGCCTTTTAGCGGCGCCCTCTATAGAT

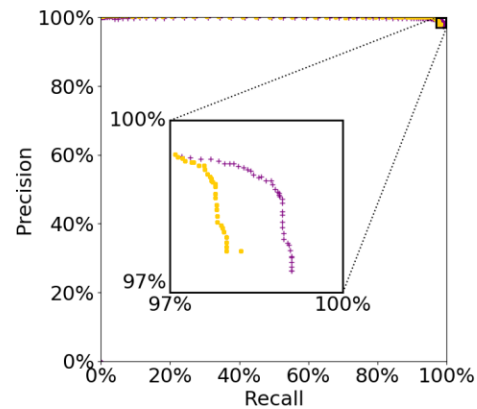
Query #1:

ACCGTTGAAGGACGGCCA TTTTTAAACTGAGCATCCATCTAAAAGCCTTTT

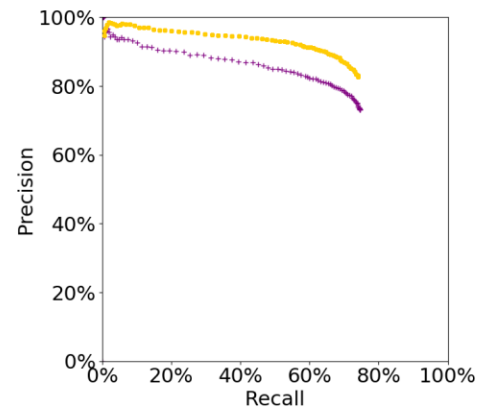
Query #2:

ACCGTTGAAGGACGGCCATTTTTA AACTGAGCATCCATCTAAAAGCCTTTT

SNP



INDEL



Prior work: **vcfeval**

$$Precision = \frac{TP}{TP + FP}$$

$$Recall = \frac{TP}{TP + FN}$$

$$F1 = 2 \cdot \frac{precision \cdot recall}{precision + recall}$$

Challenge: *comparing complex variants*

Ground Truth

Representation #1

Reference: AAGG AAATC

Truth: ATCGAAAATC

CHROM	POS	REF	ALT
chr14	2	A	T
chr14	3	G	C
chr14	4	G	GA

Representation #2

Reference: AAGG AAATC

Truth: A TCGAAAATC

CHROM	POS	REF	ALT
chr14	1	AAGG	A
chr14	1	A	ATCGA

Challenge: *comparing complex variants*

Ground Truth

Representation #1

Reference:	AAGG AAATC
Truth:	ATCGAAAATC

CHROM	POS	REF	ALT
chr14	2	A	T
chr14	3	G	C
chr14	4	G	GA

Representation #2

Reference:	AAGG	AAATC
Truth:	A	TCGAAAATC

CHROM	POS	REF	ALT
chr14	1	AAGG	A
chr14	1	A	ATCGA

Challenge: *comparing complex variants*

Technology #1

Reference: AAGGAAATC
 Query #1: ATCGAAATC

CHROM	POS	REF	ALT
chr14	2	A	T
chr14	3	G	C

Technology #2

Reference: AAGGAAATC
 Query #2: A AAATC

CHROM	POS	REF	ALT
chr14	1	AAGG	A

Ground Truth

Representation #1

Reference: AAGG AAATC
 Truth: ATCGAAATC

CHROM	POS	REF	ALT
chr14	2	A	T
chr14	3	G	C
chr14	4	G	GA

Representation #2

Reference: AAGG AAATC
 Truth: A TCGAAATC

CHROM	POS	REF	ALT
chr14	1	AAGG	A
chr14	1	A	ATCGA

Challenge: *comparing complex variants*

Technology #1

Reference: **AAGGAAATC**
 Query #1: **ATCGAAATC**

Technology #2

Reference: **AAGGAAATC**
 Query #2: **A AAATC**

Ground Truth

Representation #1

Reference: **AAGG AAATC**
 Truth: **ATCGAAATC**

CHROM	POS	REF	ALT
✓ chr14	2	A	T
✓ chr14	3	G	C

CHROM	POS	REF	ALT
✗ chr14	1	AAGG	A

CHROM	POS	REF	ALT
chr14	2	A	T
chr14	3	G	C
chr14	4	G	GA

SNP Precision: 100%
 SNP Recall: 100%
 INDEL Precision: NA
 INDEL Recall: 0%

SNP Precision: NA
 SNP Recall: 0%
 INDEL Precision: 0%
 INDEL Recall: 0%

Representation #2

Reference: **AAGG AAATC**
 Truth: **A TCGAAATC**

CHROM	POS	REF	ALT
chr14	1	AAGG	A
chr14	1	A	ATCGA

Challenge: *comparing complex variants*

Technology #1

Reference: **AAGGAAATC**
 Query #1: **ATCGAAATC**

CHROM	POS	REF	ALT
chr14	2	A	T
chr14	3	G	C

SNP Precision: 100%
 SNP Recall: 100%
 INDEL Precision: NA
 INDEL Recall: 0%

SNP Precision: 0%
 SNP Recall: NA
 INDEL Precision: NA
 INDEL Recall: 0%

Technology #2

Reference: **AAGGAAATC**
 Query #2: **A AAATC**

CHROM	POS	REF	ALT
✓ chr14	1	AAGG	A

SNP Precision: NA
 SNP Recall: 0%
 INDEL Precision: 0%
 INDEL Recall: 0%

SNP Precision: NA
 SNP Recall: NA
 INDEL Precision: 100%
 INDEL Recall: 50%

Ground Truth

Representation #1

Reference: **AAGG AAATC**
 Truth: **ATCGAAATC**

CHROM	POS	REF	ALT
chr14	2	A	T
chr14	3	G	C
chr14	4	G	GA

Representation #2

Reference: **AAGG AAATC**
 Truth: **A TCGAAATC**

CHROM	POS	REF	ALT
chr14	1	AAGG	A
chr14	1	A	ATCGA

Outline

1. **Background:** whole-genome sequencing evaluation
2. **Problem #1:** complex variants
3. **Solution:** variant normalization
4. **Results**
5. **Problem #2**
6. **Solution**
7. **Results**

Complex variants: *standard VCF normalization*

Reference AAGGAAATC

Query ATCGAAAATC

Heng Li. "Toward better understanding of artifacts in variant calling from high-coverage samples." *Bioinformatics*, 2014.
Tan et al. "Unified representation of genetic variants." *Bioinformatics*, 2015.

Complex variants: *standard VCF normalization*

Reference AAGGAAATC

Query ATCGAAAATC

Alignment

```
AAGGAAA-TC
. . . . .
ATCGAAAATC
```

VCF

POS	REF	ALT
2	AG	TC
6	AATC	AAATC

Original

Heng Li. "Toward better understanding of artifacts in variant calling from high-coverage samples." *Bioinformatics*, 2014.
Tan et al. "Unified representation of genetic variants." *Bioinformatics*, 2015.

Complex variants: *standard VCF normalization*

Reference AAGGAAATC

Query ATCGAAAATC

Alignment

AAGGAAA-TC
· · · · ·
ATCGAAAATC

AAGGAAA-TC
· · · · ·
ATCGAAAATC

VCF

POS	REF	ALT
2	AG	TC
6	AATC	AAATC

POS	REF	ALT
2	A	T
3	G	C
6	AATC	AAATC

Original

Decomposed

Heng Li. "Toward better understanding of artifacts in variant calling from high-coverage samples." *Bioinformatics*, 2014.

Tan et al. "Unified representation of genetic variants." *Bioinformatics*, 2015.

Complex variants: *standard VCF normalization*

Reference AAGGAAATC

Query ATCGAAAATC

Alignment

```
AAGGAAA-TC
. . . . .
ATCGAAAATC
```

```
AAGGAAA-TC
. . . . .
ATCGAAAATC
```

```
AAGGAAA-TC
. . . . .
ATCGAAAATC
```

VCF

POS	REF	ALT
2	AG	TC
6	AATC	AAATC

POS	REF	ALT
2	A	T
3	G	C
6	AATC	AAATC

POS	REF	ALT
2	A	T
3	G	C
7	A	AA

Original

Decomposed

Trimmed

Heng Li. "Toward better understanding of artifacts in variant calling from high-coverage samples." *Bioinformatics*, 2014.

Tan et al. "Unified representation of genetic variants." *Bioinformatics*, 2015.

Complex variants: *standard VCF normalization*

Reference AAGGAAATC

Query ATCGAAAATC

Alignment

```
AAGGAAA-TC
. . . . .
ATCGAAAATC
```

```
AAGGAAA-TC
. . . . .
ATCGAAAATC
```

```
AAGGAAA-TC
. . . . .
ATCGAAAATC
```

```
AAGG-AAATC
. . . . .
ATCGAAAATC
```

VCF

POS	REF	ALT
2	AG	TC
6	AATC	AAATC

POS	REF	ALT
2	A	T
3	G	C
6	AATC	AAATC

POS	REF	ALT
2	A	T
3	G	C
7	A	AA

POS	REF	ALT
2	A	T
3	G	C
4	G	GA

Original

Decomposed

Trimmed

Left shifted

Heng Li. "Toward better understanding of artifacts in variant calling from high-coverage samples." *Bioinformatics*, 2014.

Tan et al. "Unified representation of genetic variants." *Bioinformatics*, 2015.

Complex variants: *standard VCF normalization*

Reference AAGGAAATC

Query ATCGAAAATC

Alignment

```
AAGGAAA-TC
. . . . .
ATCGAAAATC
```

```
AAGGAAA-TC
. . . . .
ATCGAAAATC
```

```
AAGGAAA-TC
. . . . .
ATCGAAAATC
```

```
AAGG-AAATC
. . . . .
ATCGAAAATC
```

```
AAGG----AAATC
. . . . .
A---TCGAAAATC
```

VCF

POS	REF	ALT
2	AG	TC
6	AATC	AAATC

POS	REF	ALT
2	A	T
3	G	C
6	AATC	AAATC

POS	REF	ALT
2	A	T
3	G	C
7	A	AA

POS	REF	ALT
2	A	T
3	G	C
4	G	GA

POS	REF	ALT
1	AAGG	A
1	A	ATCGA

Original

Decomposed

Trimmed

Left shifted

Alternate

Heng Li. "Toward better understanding of artifacts in variant calling from high-coverage samples." *Bioinformatics*, 2014.

Tan et al. "Unified representation of genetic variants." *Bioinformatics*, 2015.

Complex variants: *standard VCF normalization*

Reference AAGGAAATC

Query ATCGAAAATC

Alignment

```
AAGGAAA-TC
. . . . .
ATCGAAAATC
```

```
AAGGAAA-TC
. . . . .
ATCGAAAATC
```

```
AAGGAAA-TC
. . . . .
ATCGAAAATC
```

```
AAGG-AAATC
. . . . .
ATCGAAAATC
```

```
AAGG----AAATC
. . . . .
A---TCGAAAATC
```

VCF

POS	REF	ALT
2	AG	TC
6	AATC	AAATC

POS	REF	ALT
2	A	T
3	G	C
6	AATC	AAATC

POS	REF	ALT
2	A	T
3	G	C
7	A	AA

POS	REF	ALT
2	A	T
3	G	C
4	G	GA

POS	REF	ALT
1	AAGG	A
1	A	ATCGA

Original

Decomposed

Trimmed

Left shifted

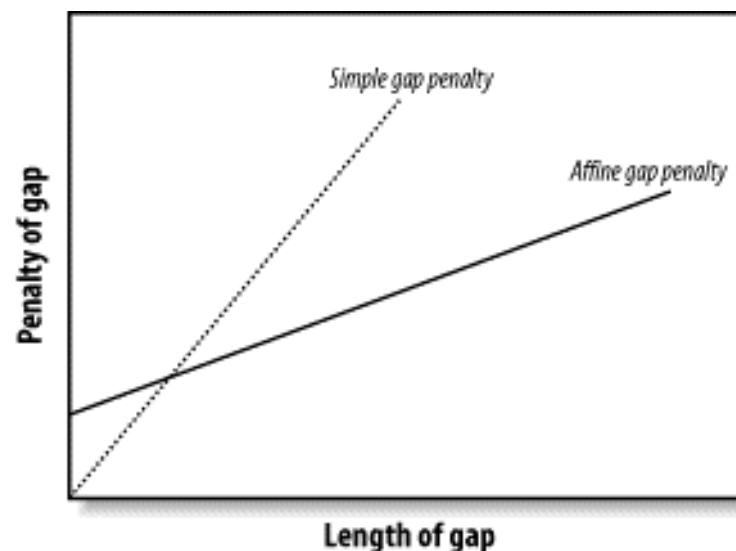
Alternate

Heng Li. "Toward better understanding of artifacts in variant calling from high-coverage samples." *Bioinformatics*, 2014.

Tan et al. "Unified representation of genetic variants." *Bioinformatics*, 2015.

Choosing representations: *best-alignment normalization*

m = match
 x = mis-match
 o = gap opening
 e = gap extension



Choosing representations: *best-alignment normalization*

m = match
 x = mis-match
 o = gap opening
 e = gap extension

Option #1

```
AAGG-AAATC
. . . . .
ATCGAAAATC
```

POS	REF	ALT
2	A	T
3	G	C
4	G	GA

$$x + x + (o+e)$$

Option #2

```
AAGG----AAATC
. . . . .
A---TCGAAAATC
```

POS	REF	ALT
1	AAGG	A
1	A	ATCGA

$$(o+3e) + (o+4e)$$

Choosing representations: *best-alignment normalization*

$m = 0$ = match
 $x = 5$ = mis-match
 $o = 6$ = gap opening
 $e = 2$ = gap extension

Option #1

```
AAGG-AAATC
. . . . .
ATCGAAAATC
```

POS	REF	ALT
2	A	T
3	G	C
4	G	GA

$x + x + (o+e)$
18

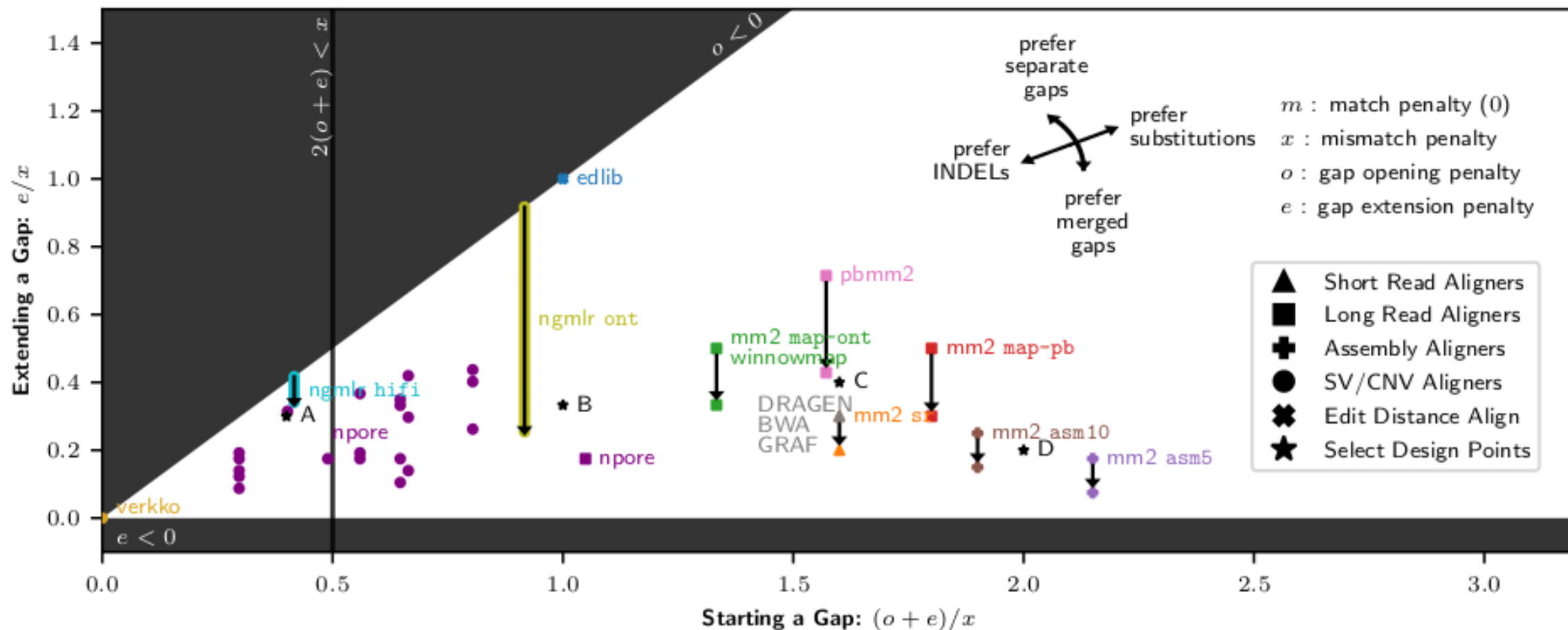
Option #2

```
AAGG----AAATC
. . . . .
A---TCGAAAATC
```

POS	REF	ALT
1	AAGG	A
1	A	ATCGA

$(o+3e) + (o+4e)$
26

Alignment-based normalization design space



Tim Dunn, Satish Narayanasamy. "vcfdist: accurately benchmarking phased small variant calls in human genomes". Nature Communications, 2023.

Outline

1. **Background:** whole-genome sequencing evaluation
2. **Problem #1:** complex variants
3. **Solution:** variant normalization
4. **Results:** stable evaluations
5. **Problem #2**
6. **Solution**
7. **Results**

Example: *tandem repeat benchmark representation*

Dataset	SNPs	INDELS
Original Representation	917,255	431,545
Normalized At Point C	502,076	461,258

Example: *tandem repeat benchmark representation*

Original VCF: *GIAB Tandem Repeats*

```
chr20 278985   A C
chr20 278986   C G
chr20 278990   G C
chr20 278993   C A
chr20 278994   G GGGAGGGAGGGCGGGACGGAGGGA
GGGAGGGAGGGACGGAGGGCGGGACGGCGGGAGGGCGGGAC
GGAGGGACGGAGGGAGGGCGGGACGGAGGGCGGGAGGGCGG
GACGGAGGGAGGGAGGGAGGGAGGGCGGGACGGAGGGAGGG
AGGGCGGGACGGAGGGAGGGAGGGAGGGACGGAGGGCGGGGA
CGGCGGGAGGGCGGGACGGAGGGACGGAGGGAGGGCGGGAC
GGAGGGCGGGAGGGCGGGACGGAGGGAGGGAGGGCGGGACG
GAGGGACGGAGGGAGGGAGGGCGGGACGGAGGGAGGGAGGG
CGGGACGGAGGGACGGAGGGAGGGAGGGCGGGACGGAGGGA
GGGAGGGAGGGACGGAGGGCGGGACGGCGGGAGGGCGGGAC
GGAGGGACGGAGGGAGGGCGGGACGGAGGGCGGGAGGGAGG
GAGGGCGGGACGGAGGGAGGGAGGGAGGGAGGGCGGGACGG
AGGGAGGGAGGGAGGGAGGGACGGAGGGACGGAGGGAGGGA
GGGAGGGAGGGACGGAGGGCGGGACGGAGGGAGGGAGGGCG
GAGGGAGGGAGGGCGGGACGGAGGGAGGGAGGGAGGGACGG
AGGGCGGGACGGAGGGAGGGAGGGC
```

chr20 278998 C G
chr20 279001 C A
chr20 279022 C G
chr20 279029 A C
chr20 279033 C A
chr20 279038 C T
chr20 279045 C A
chr20 279069 A C

12 SNPs
1 INS (622bp)

Normalized VCF: *vcfdist design point C*

```
chr20 278984   G GCGGGACGGAGGGAGGGAGGGCG
GGACGGAGGGAGGGAGGGAGGGACGGAGGGCGGGA
CGGCGGGAGGGCGGGACGGAGGGACGGAGGGAGGG
CGGGACGGAGGGCGGGAGGGCGGGACGGAGGGAGG
GAGGGAGGGAGGGCGGGACGGAGGGAGGGAGGGCG
GGACGGAGGGAGGGAGGGAGGG
chr20 279069   A AGGGCGGGACGGAGGGACGGAGG
GAGGGAGGGCGGGACGGAGGGAGGGAGGGCGGGAC
GGAGGGACGGAGGGAGGGAGGGCGGGACGGAGGGGA
GGGAGGGAGGGACGGAGGGCGGGACGGCGGGAGGG
CGGGACGGAGGGACGGAGGGAGGGCGGGACGGAGG
GCGGGAGGGAGGGAGGGCGGGACGGAGGGAGGGAG
GGAGGGAGGGCGGGACGGAGGGAGGGAGGGAGGGA
GGGACGGAGGGACGGAGGGAGGGAGGGAGGGAGGG
ACGGAGGGCGGGACGGAGGGAGGGAGGGCGGAGGG
AGGGAGGGCGGGACGGAGGGAGGGAGGGAGGGACG
GAGGGCGGGACGGAGGGAGGGAGGGCGGAGGGAGG
GAGGGCGGGACGGAGGGAGGGAGGGCGGGAGGGAT
GGAGGGAGGGAGGGCGGGACGGAGGGAGGGC
```

2 INS (438bp, 184bp)

Other Contributions

- Efficient variant clustering
- Allow inexact variant matches
- Use phasing information for evaluation
- Distance based evaluation metrics

Example: *inexact complex variant matches*

Query:

CHROM	POS	REF	ALT	CALL	CREDIT
chr1	976722	C	CAGGAACCGCCTCCCACTCCCCCA	CAACCCCGGGGAACCGCCTCCCACTC	
			CCCCCGCAACCCCGGGGAACCGCCTCCCACTCCCCCGCAACCCC	INS TP	0.979
chr1	976745	G	A	SNP TP	0.979

Truth:

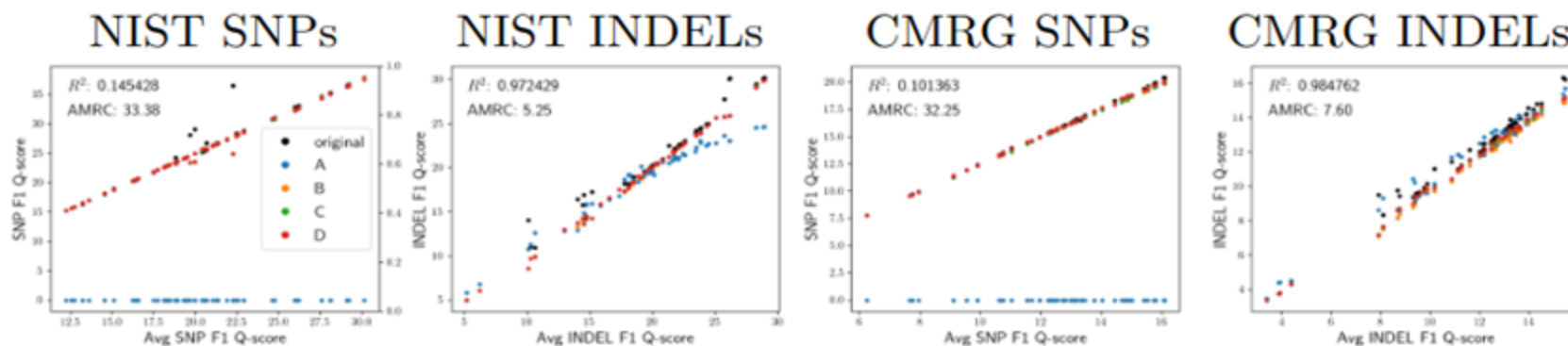
CHROM	POS	REF	ALT	CALL	CREDIT
chr1	976715	A	CAACCCAGGAACCGCCTCCCACTCCCCCA	INS TP	0.979
chr1	976747	C	CAACCCCGGGGAACCGCCTCCCACTCCCCCG	INS TP	0.979
chr1	976777	G	A	SNP TP	0.979
chr1	976811	C	CAACCCCGGGGAACCGCCTCCCACTCCCCCG	INS TP	0.979
chr1	976840	C	G	SNP TP	0.979
chr1	976841	G	A	SNP TP	0.979

Dataset: *PrecisionFDA Truth Challenge V2*

- 64 whole genome sequencing submissions
- Illumina, PacBio, ONT, and Multi-tech

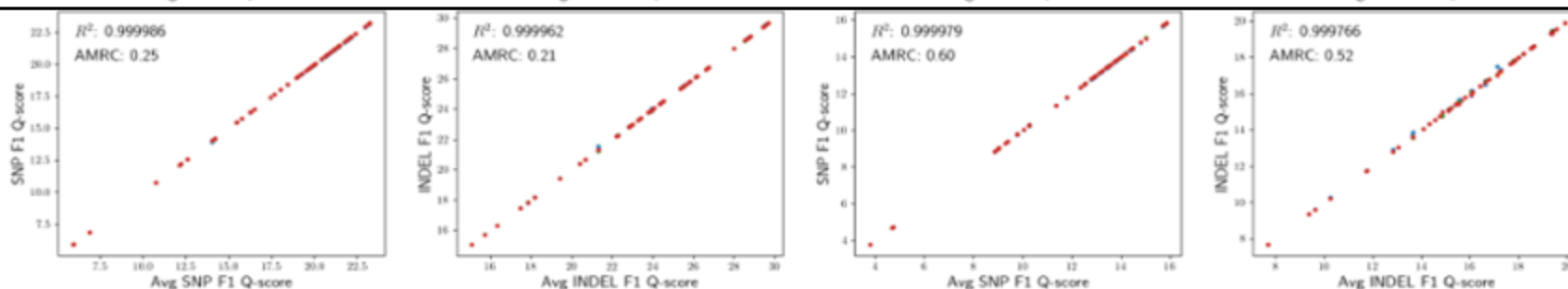
Results: *stable performance across representations*

vcfeval
precision/recall
metrics



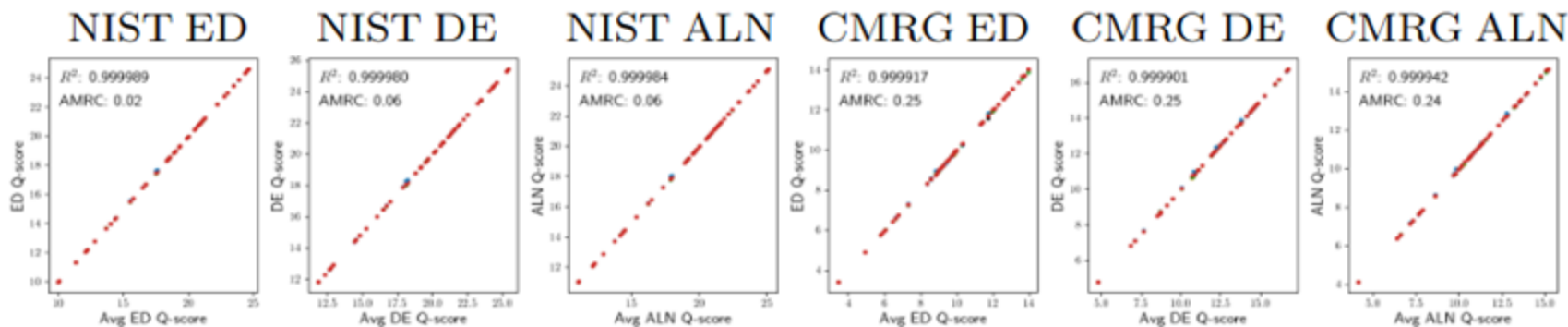
Prior Work

vcfdist
precision/recall
metrics



This Work

vcfdist
distance
metrics



Outline

1. **Background:** whole-genome sequencing evaluation
2. **Problem #1:** complex variants
3. **Solution:** variant normalization
4. **Results:** stable evaluations
5. **Problem #2:** separate evaluation of small and structural variants
6. **Solution**
7. **Results**

Problem: *separate evaluations for small and SVs*

Reference:

ACCGTTGAAGGACGGCCATTTTTT AACTGAGCATCCATCTAAAAGCCTTTTAGCGGCGCCCTCTATAGAT

Query #1:

ACCCTTGAAGGACGGCCA TTTTTAAACTGAGCATCCATCTAAAAGCCTTTT

SNP

*single nucleotide
polymorphism*

substitution

INDEL

insertion/deletion

<50 basepairs

SV

structural variant

50+ basepairs

Example: *variant matches across size categories*

Query:

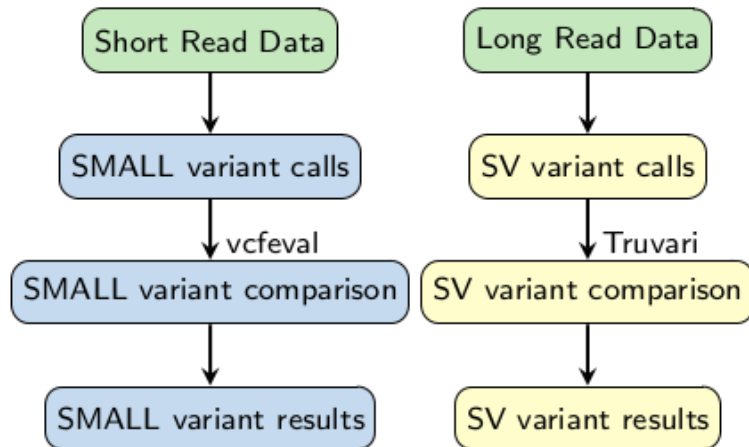
CHROM	POS	REF	ALT	CALL	CREDIT
chr1	976722	C	CAGGAACCGCCTCCCACTCCCCCA	CAACCCCGGGGAACCGCCTCCCACTC	
			CCCCCGCAACCCCGGGGAACCGCCTCCCACTCCCCCGCAACCCC	INS TP	0.979
chr1	976745	G	A	SNP TP	0.979

Truth:

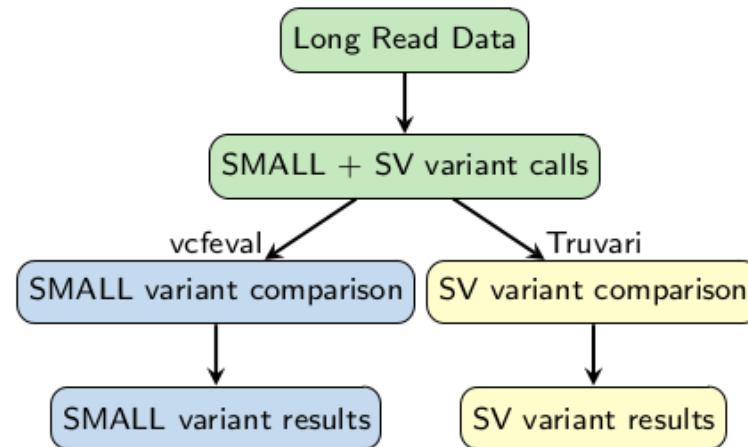
CHROM	POS	REF	ALT	CALL	CREDIT
chr1	976715	A	CAACCCAGGAACCGCCTCCCACTCCCCCA	INS TP	0.979
chr1	976747	C	CAACCCCGGGGAACCGCCTCCCACTCCCCCG	INS TP	0.979
chr1	976777	G	A	SNP TP	0.979
chr1	976811	C	CAACCCCGGGGAACCGCCTCCCACTCCCCCG	INS TP	0.979
chr1	976840	C	G	SNP TP	0.979
chr1	976841	G	A	SNP TP	0.979

Why: *short-read mapping, inaccurate long reads*

Historically



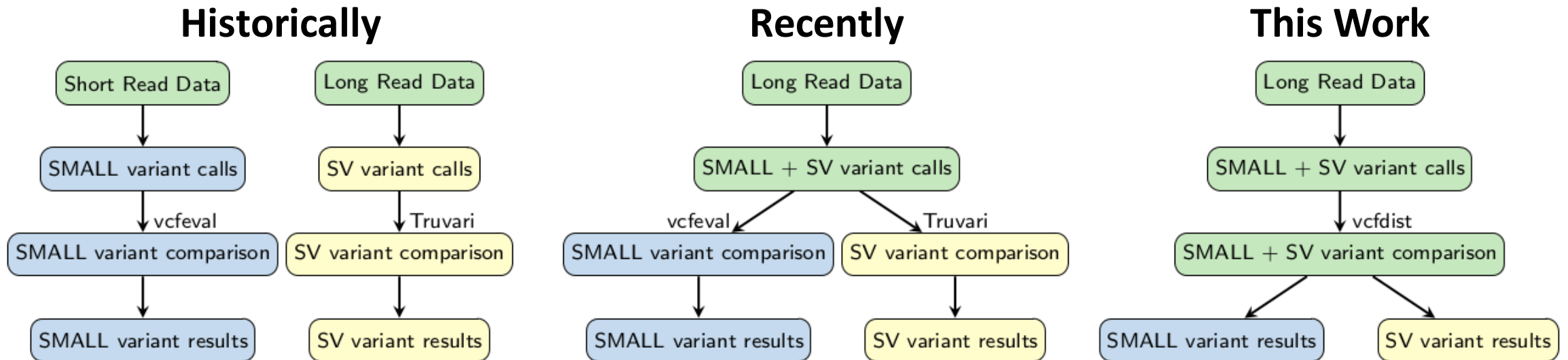
Recently



Outline

1. **Background:** whole-genome sequencing evaluation
2. **Problem #1:** complex variants
3. **Solution:** variant normalization
4. **Results:** stable evaluations
5. **Problem #2:** separate evaluation of small and structural variants
6. **Solution:** joint evaluation of small and structural variants
7. **Results**

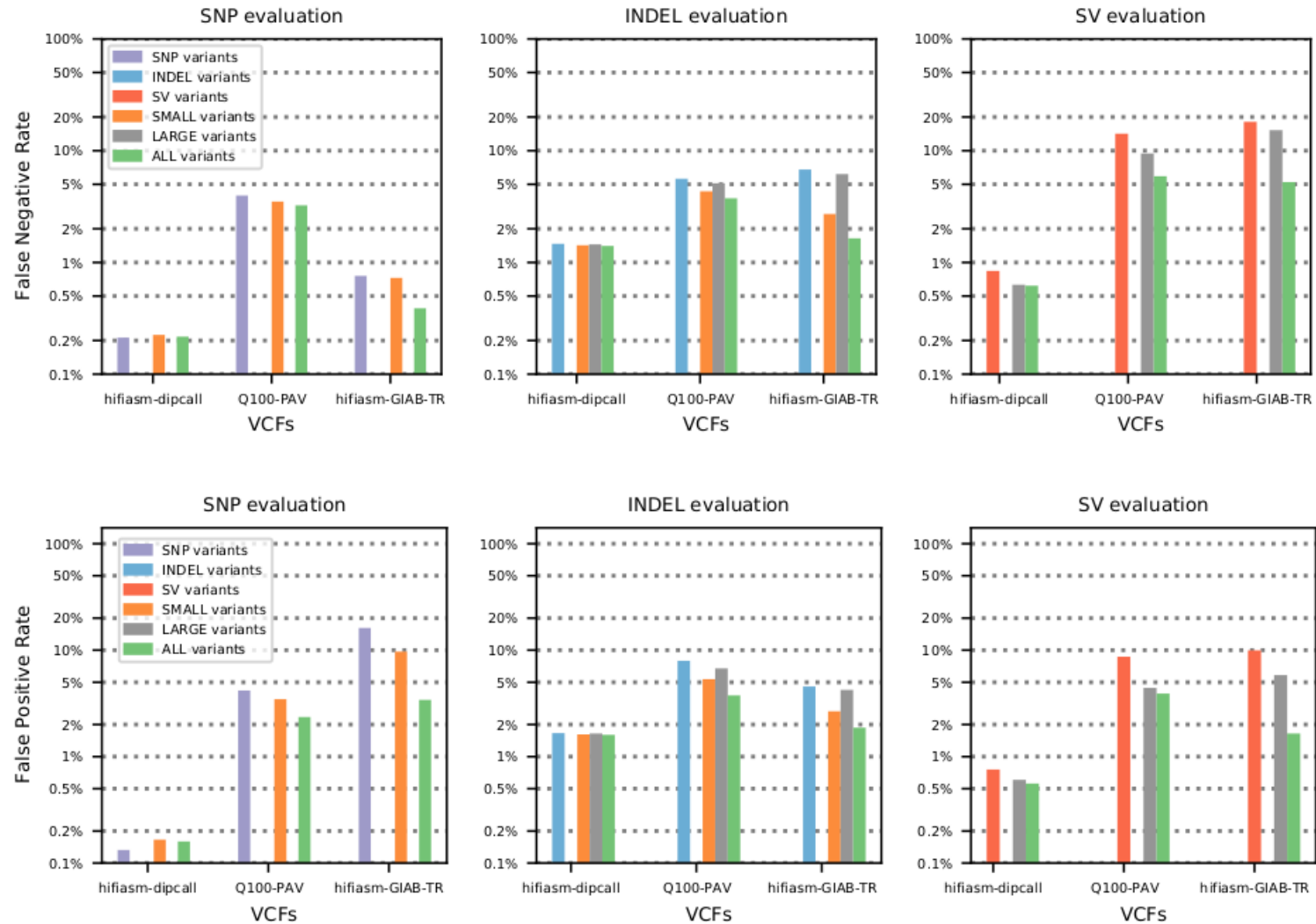
Solution: *joint small and SV evaluation*



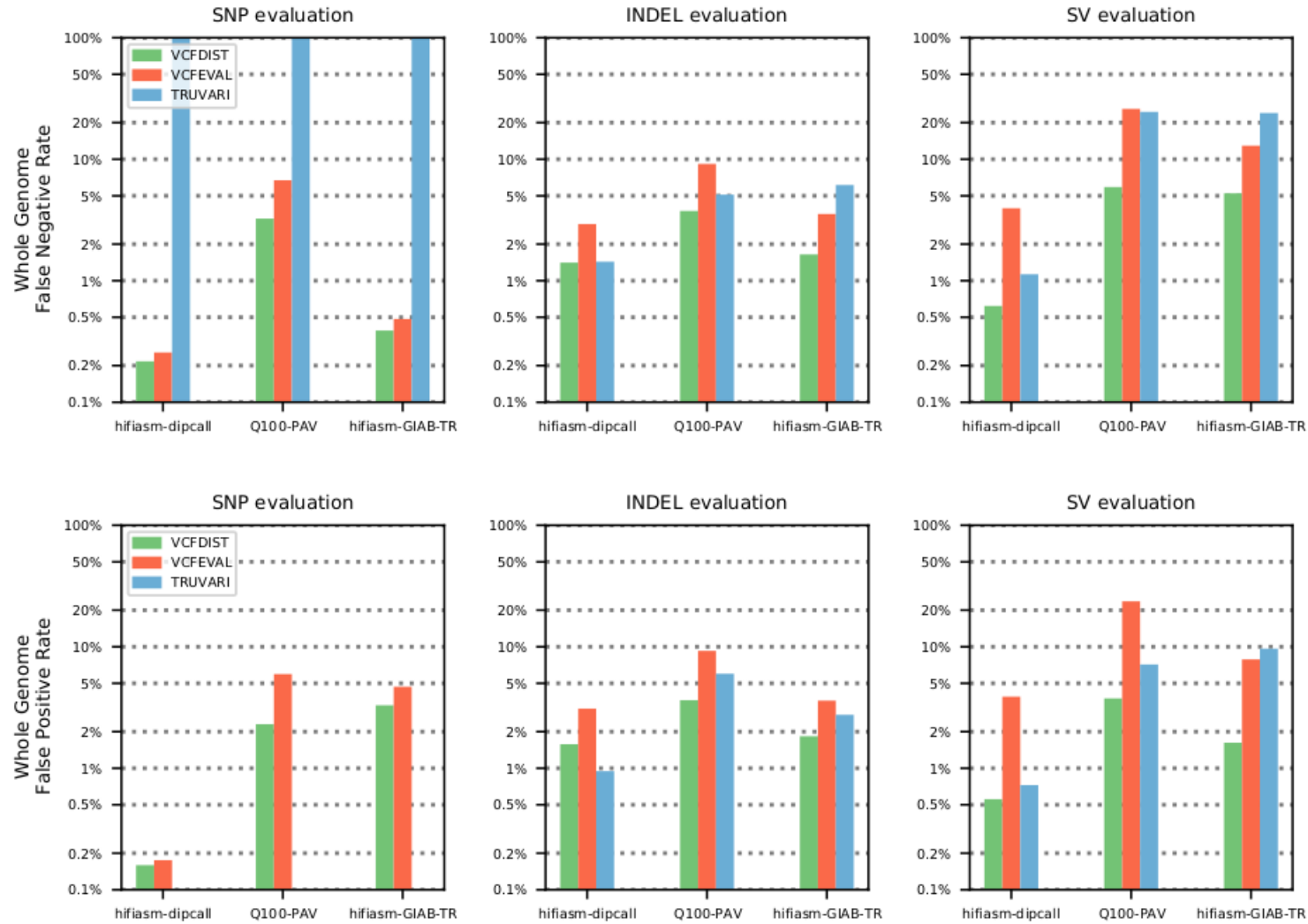
Outline

1. **Background:** whole-genome sequencing evaluation
2. **Problem #1:** complex variants
3. **Solution:** variant normalization
4. **Results:** stable evaluations
5. **Problem #2:** separate evaluation of small and structural variants
6. **Solution:** joint evaluation of small and structural variants
7. **Results:** accurate evaluations

Results: *joint small and SV evaluation*



Results: *comparison to prior work*



Results: *better phasing evaluations*

Dataset	Tool	Switches	Flips
hifiasm-dipcall	WhatsHap	610	396
	vcfdist	494	390
Q100-PAV	WhatsHap	324	433
	vcfdist	6	52
hifiasm-GIAB-TR	WhatsHap	1074	1004
	vcfdist	494	396

Results: *better phasing evaluations*

CONTIG	POS	REF	ALT	FORMAT	TRUTH	QUERY
chr1	32,653,646	T	G	GT:BD:BC	0 1:TP:1.0	0 1:TP:1.0
chr1	32,653,657	TTTG	T	GT:BD:BC	0 1:TP:1.0
chr1	32,653,658	TTG	T	GT:BD:BC	0 1:TP:1.0
chr1	32,653,658	TTG	T	GT:BD:BC	1 0:TP:1.0
chr1	32,653,659	TG	T	GT:BD:BC	1 0:TP:1.0
chr1	32,653,665	TG	T	GT:BD:BC	1 1:TP:1.0
chr1	32,653,666	G	T	GT:BD:BC	1 1:TP:1.0

	Position	37	46	47	58	59	60	61	66	67
	Reference	GTTTTTTTT	T	TTTTTTTTTTTT	T	T	G	TTTTTT	G	TTTT
Haplotype 1	Truth	GTTTTTTTT	T	TTTTTTTTTTTT	T			TTTTTT	T	TTTT
	Query	GTTTTTTTT	T	TTTTTTTTTTTT	T	T		TTTTTT		TTTT
Haplotype 2	Truth	GTTTTTTTT	G	TTTTTTTTTTTT				TTTTTT	T	TTTT
	Query	GTTTTTTTT	G	TTTTTTTTTTTT	T			TTTTTT		TTTT

Conclusion



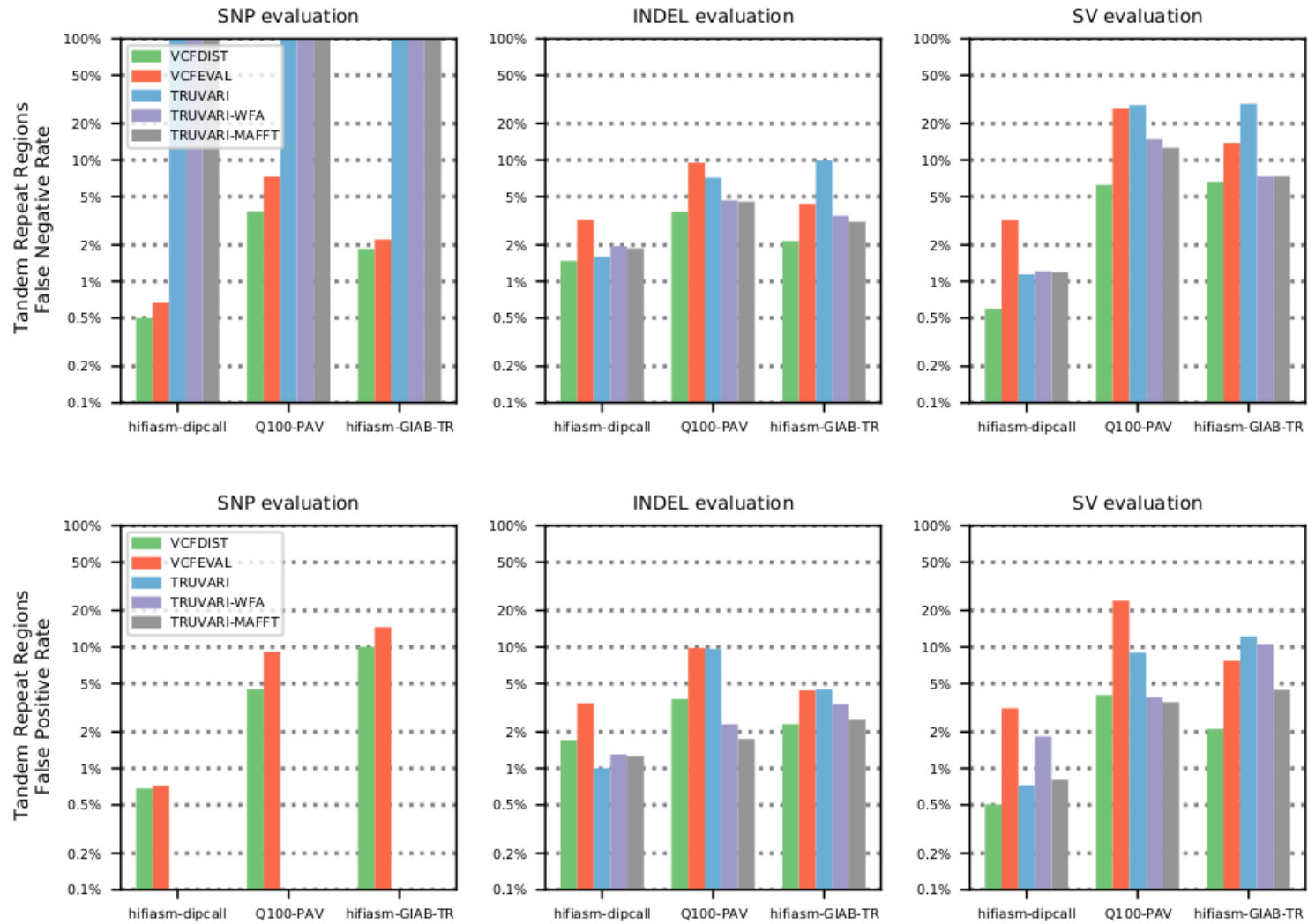
`github.com/TimD1/vcfdist`



This project was supported by the National Science Foundation Graduate Research Fellowship under Grant 1841052. Any opinion, findings, and conclusions or recommendations expressed in this material are those of the authors and do not necessarily reflect the views of the National Science Foundation

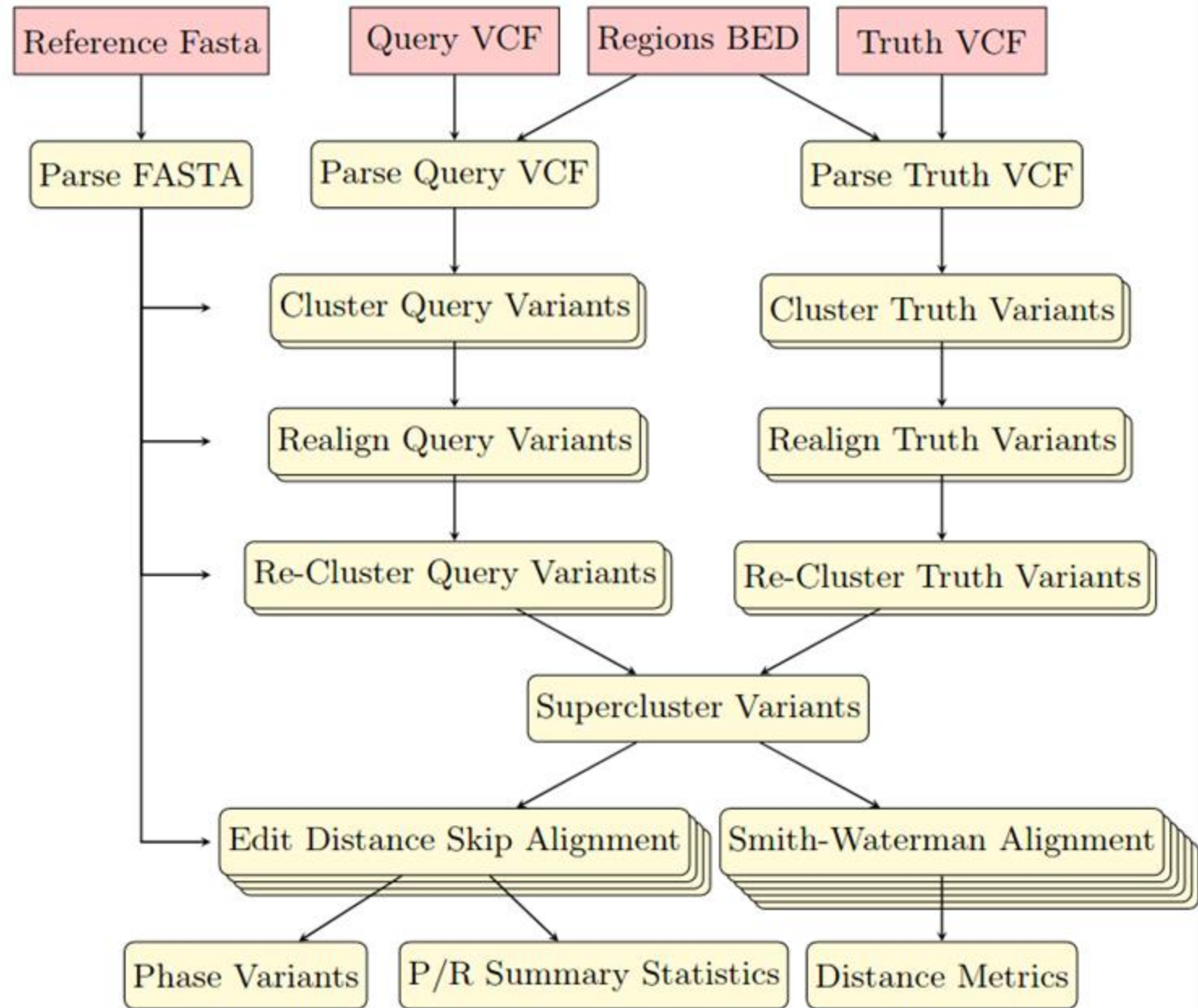
Supplementary Slides

Results: *comparison to prior work*

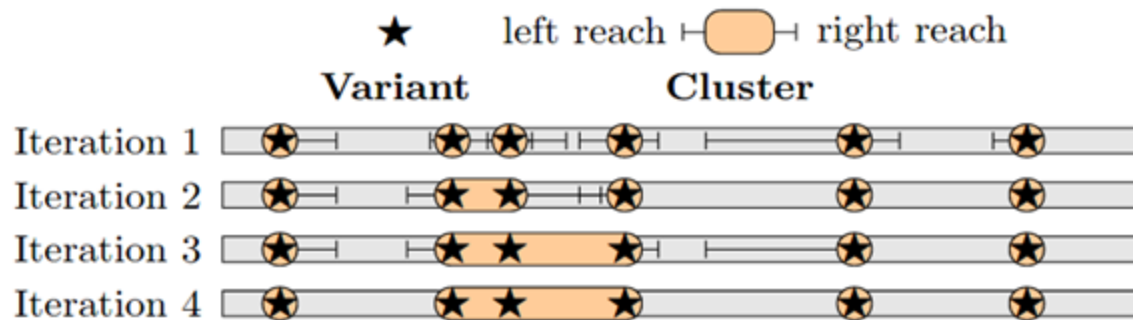


Supplementary Slides: *Implementation*

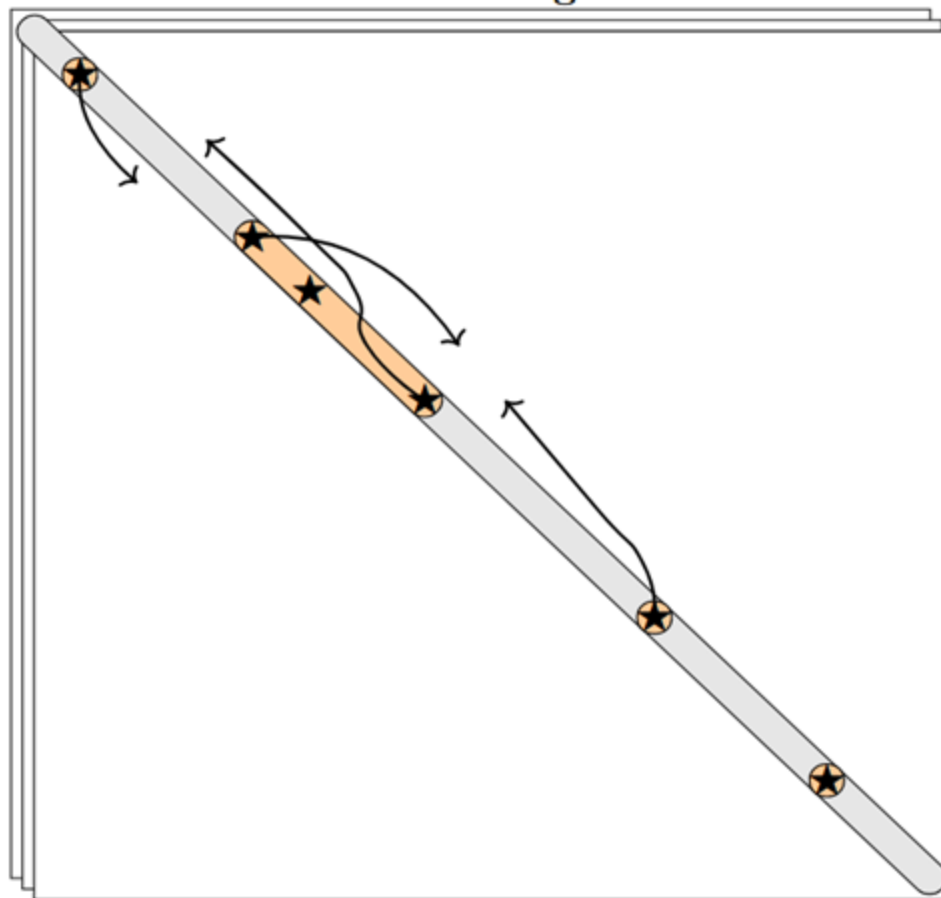
Overview



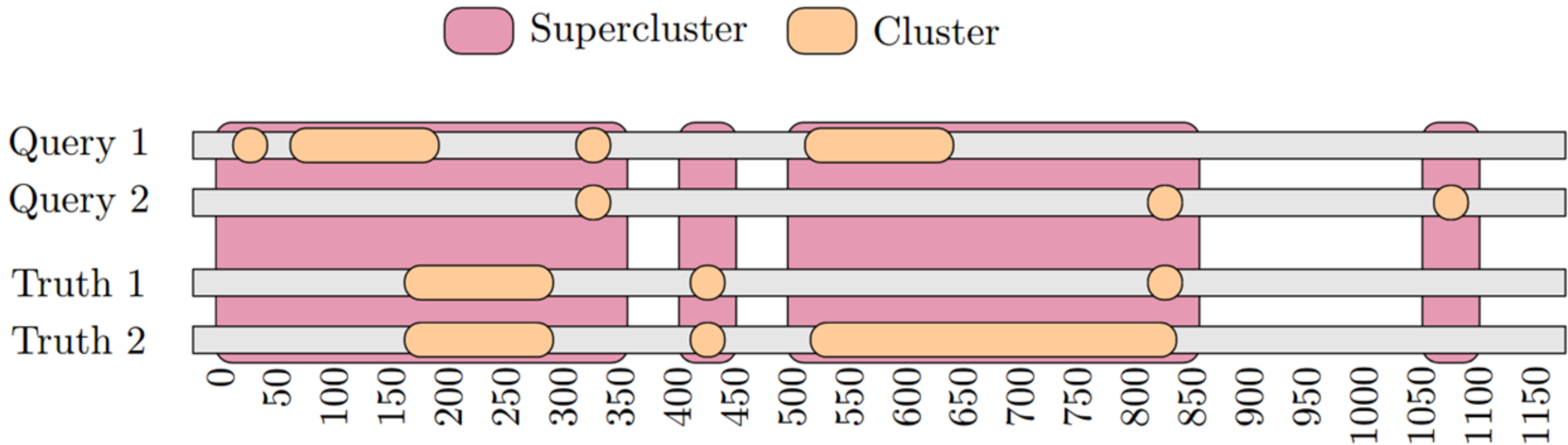
Clustering: *iterative* *bi-directional* *wavefront* *alignment*



Iteration 3 Alignment



Superclustering: *simple reference distance heuristic*



Precision/Recall:

edit distance, allows skipping FP query variants, backtracking

(a) Reference ATGCTCC

Query VCF				Truth VCF			
POS	REF	ALT	GT	POS	REF	ALT	GT
2	T	C	1 0	2	T	C	1 0
4	C	CG	1 1	4	C	CGG	1 1
6	C	A	1 0	7	C	G	1 1

