# *n*PoRe: *n*-Polymer Realigner for improved pileup variant calling

Tools and Technology Seminar Series

September 15th, 2022

**Tim Dunn**

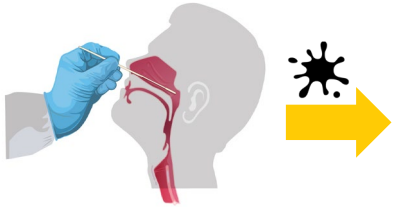UNIVERSITY OF MICHIGAN

# Overview

1. Background
   1. Whole Genome Sequencing
   2. Nanopore Sequencing
   3. Read Alignment
   4. Variant Calling

2. n-Polymer Realigner
   1. Motivation
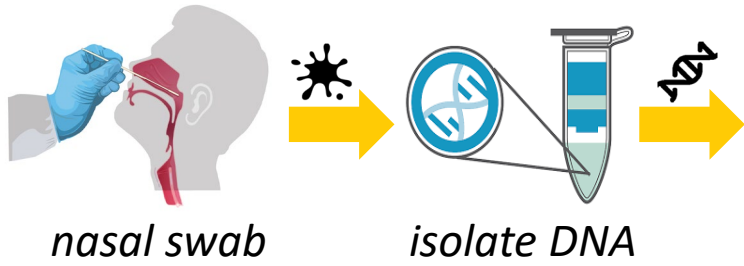   2. Algorithm
   3. Results
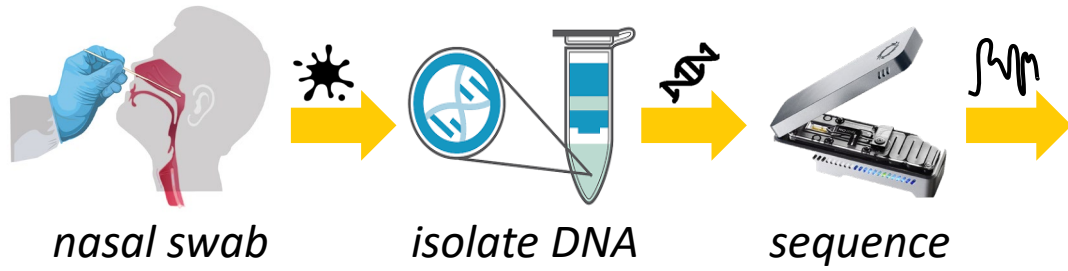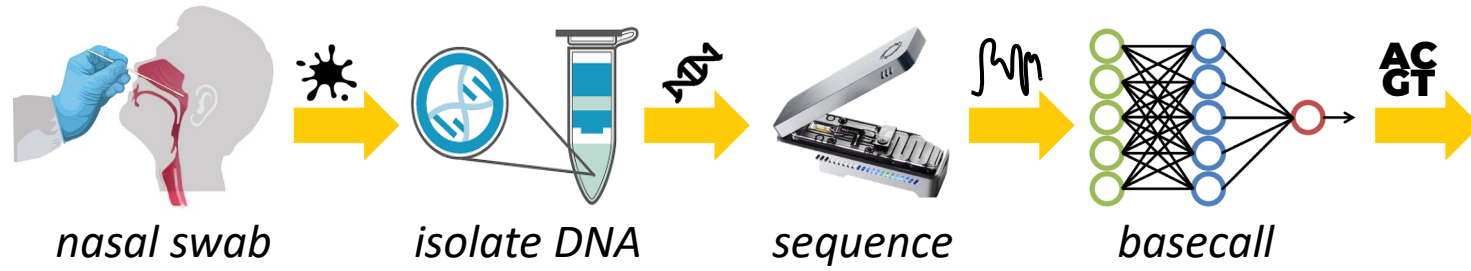
# Overview

# Whole Genome Sequencing: *Overview*

*nasal swab*

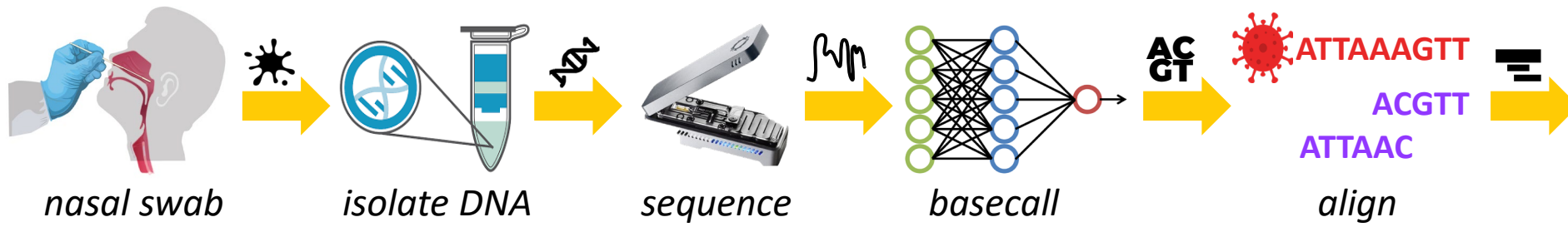*nasal swab*          *isolate DNA*

nasal swab      isolate DNA      sequence

nasal swab        isolate DNA        sequence        basecall

# Whole Genome Sequencing: *Overview*

nasal swab → isolate DNA → sequence → basecall → align

nasal swab     isolate DNA     sequence     basecall     align     variant call

# Whole Genome Sequencing: *Overview*



nasal swab → isolate DNA → sequence → basecall → align → variant call

**align:**
AC GT → ATTAAAGTT
ACGTT
ATTAAC

**variant call:** ATTAACGTT

nasal swab → isolate DNA → *sequence* → *basecall* → align → variant call

| **Illumina** | **Oxford Nanopore** | **Pacific Biosciences** |
|---|---|---|
| 100-1,000 bases | 1,000-1,000,000 bases | 1,000-50,000 bases |
| 99.9% accurate | 90-99.9% accurate | 99.9% accurate |
| Accuracy | Mapping | Accuracy in repetitive regions |
| Cost | Phasing | |

*nasal swab* → *isolate DNA* → *sequence* → *basecall* → *align* → *variant call*

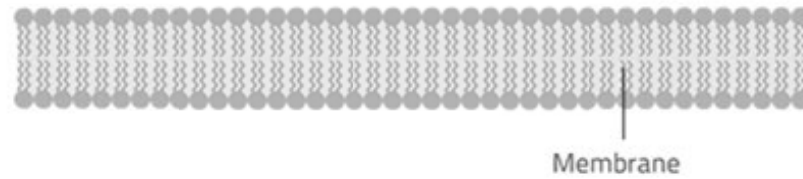| **Illumina** | **Oxford Nanopore** | **Pacific Biosciences** |
|---|---|---|
| 100-1,000 bases | 1,000-1,000,000 bases | 1,000-50,000 bases |
| 99.9% accurate | 90-99.9% accurate | 99.9% accurate |
| Accuracy<br>Cost | Mapping<br>Phasing | Accuracy in repetitive regions |

# Whole Genome Sequencing: *Applications*

1. Viral/Bacterial Strain Typing
    1. Contact tracing
    2. Antibiotic resistance

2. Patient Diagnoses
    1. Cancer
    2. Genetic Diseases

3. Patient Treatment
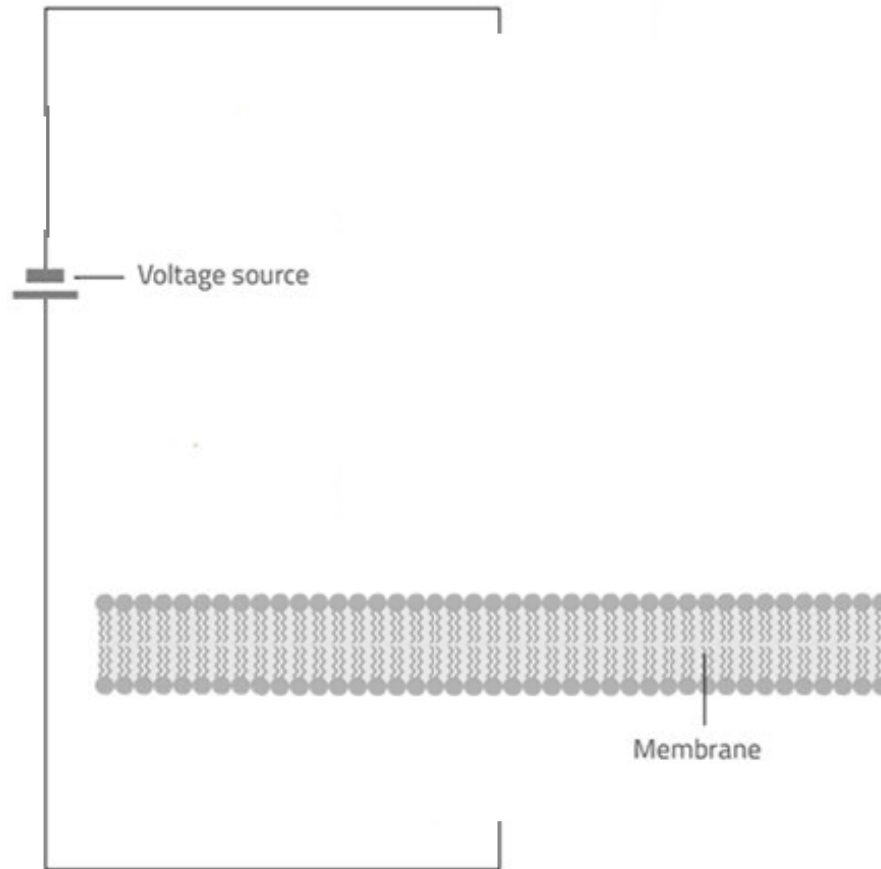    1. Predict drug response

# Overview

1. Background
   1. Whole Genome Sequencing
   2. **Nanopore Sequencing**
   3. Read Alignment
   4. Variant Calling
2. n-Polymer Realigner
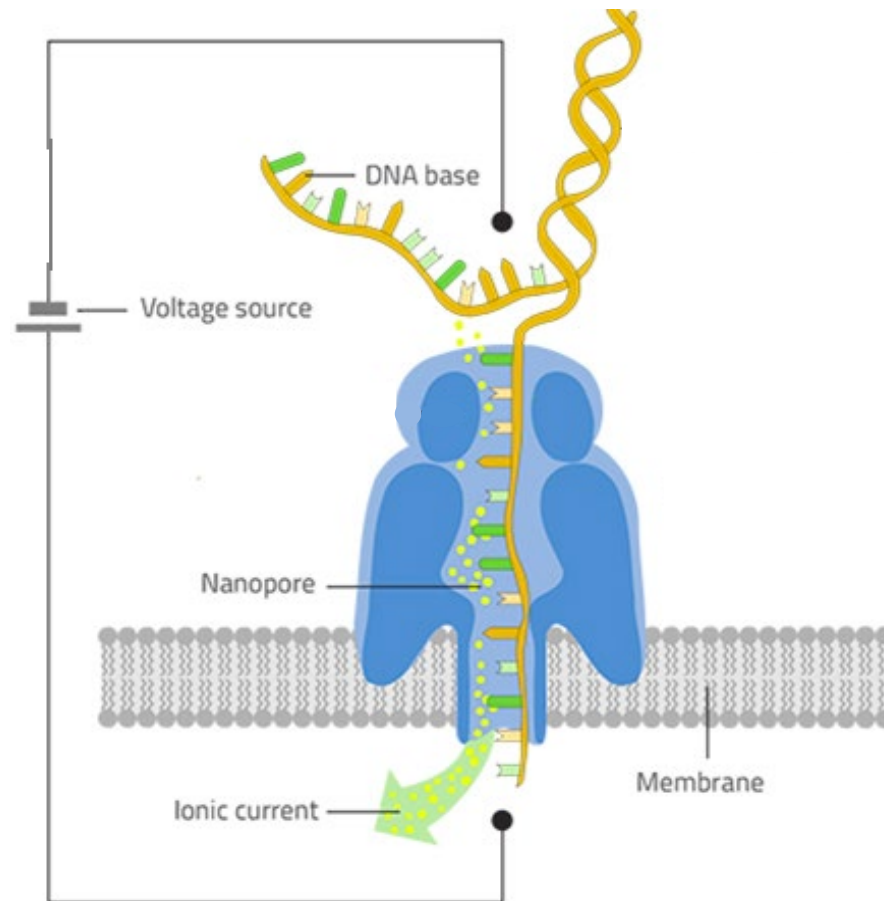   1. Motivation
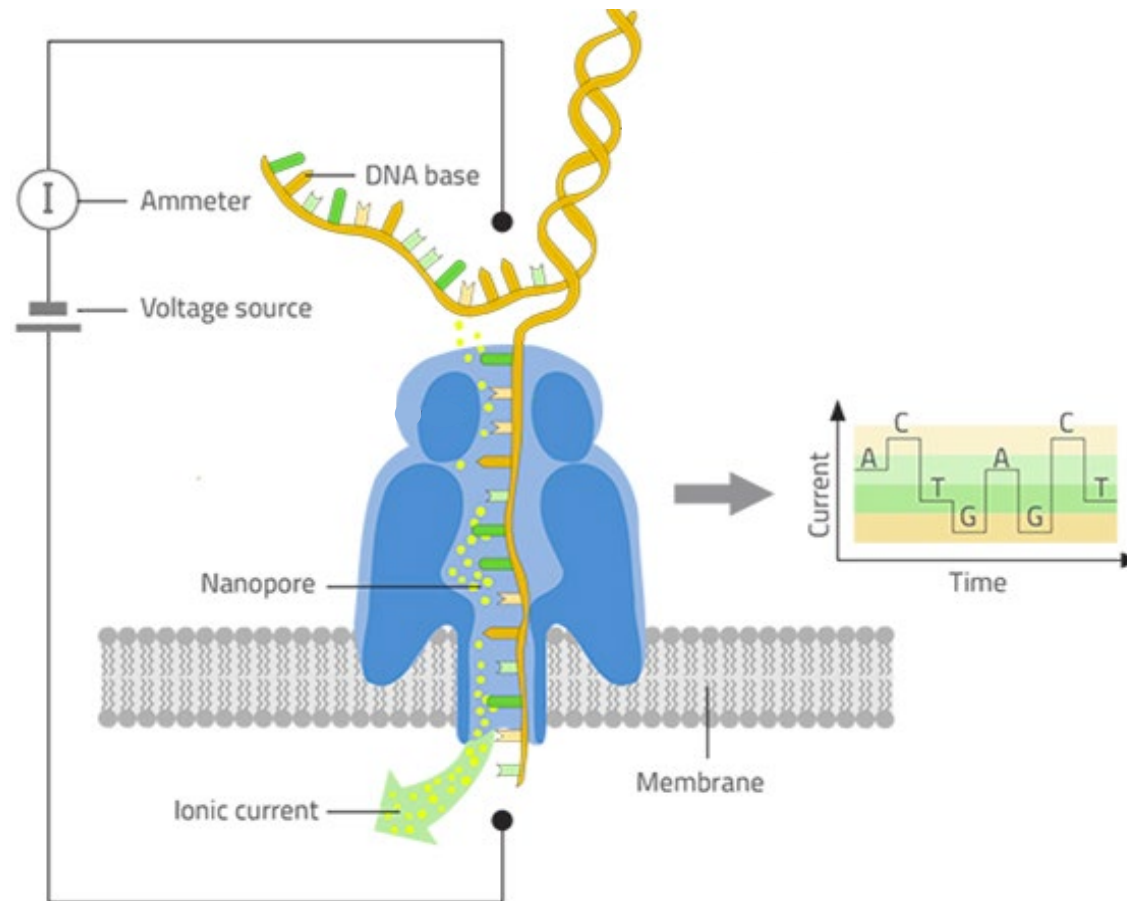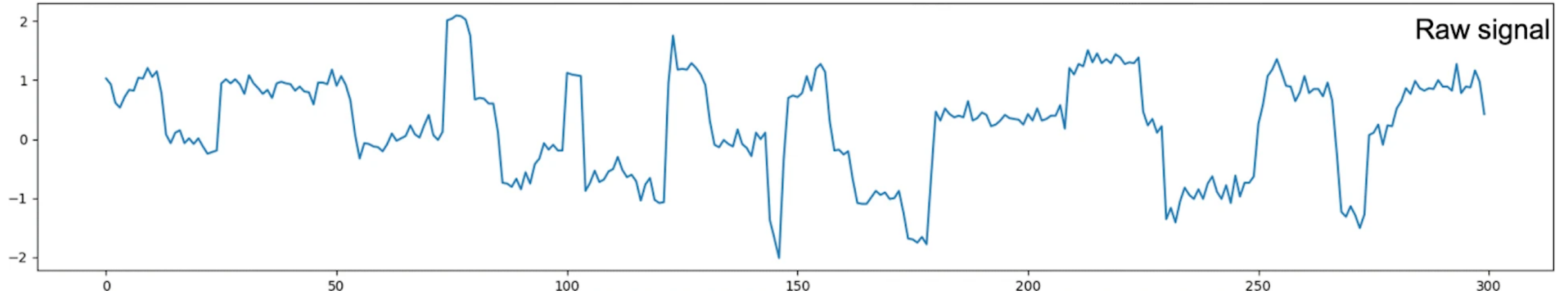   2. Algorithm
   3. Results

# Nanopore Sequencing



Membrane

# Nanopore Sequencing



Voltage source

Membrane

# Nanopore Sequencing
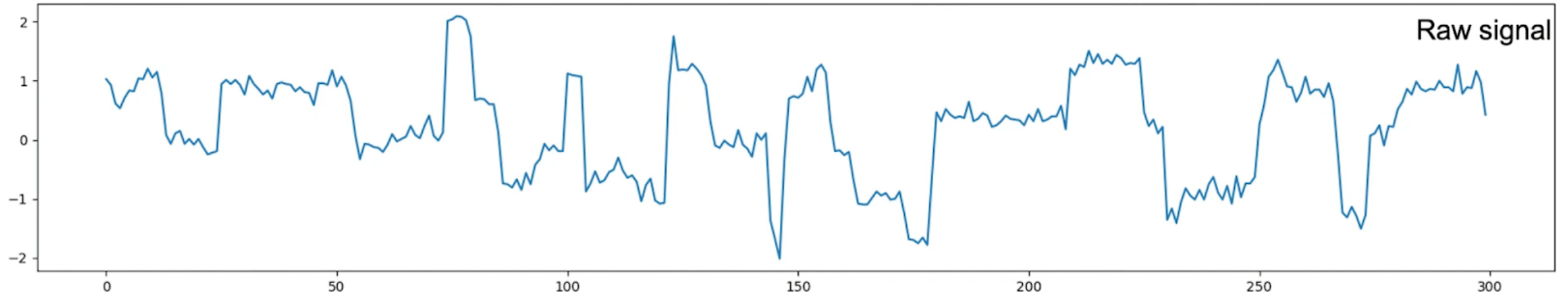
# Nanopore Sequencing

# Nanopore Sequencing

# Nanopore Sequencing

# Overview

1. Background
   1. Whole Genome Sequencing
   2. Nanopore Sequencing
   3. **Read Alignment**
   4. Variant Calling
2. n-Polymer Realigner
   1. Motivation
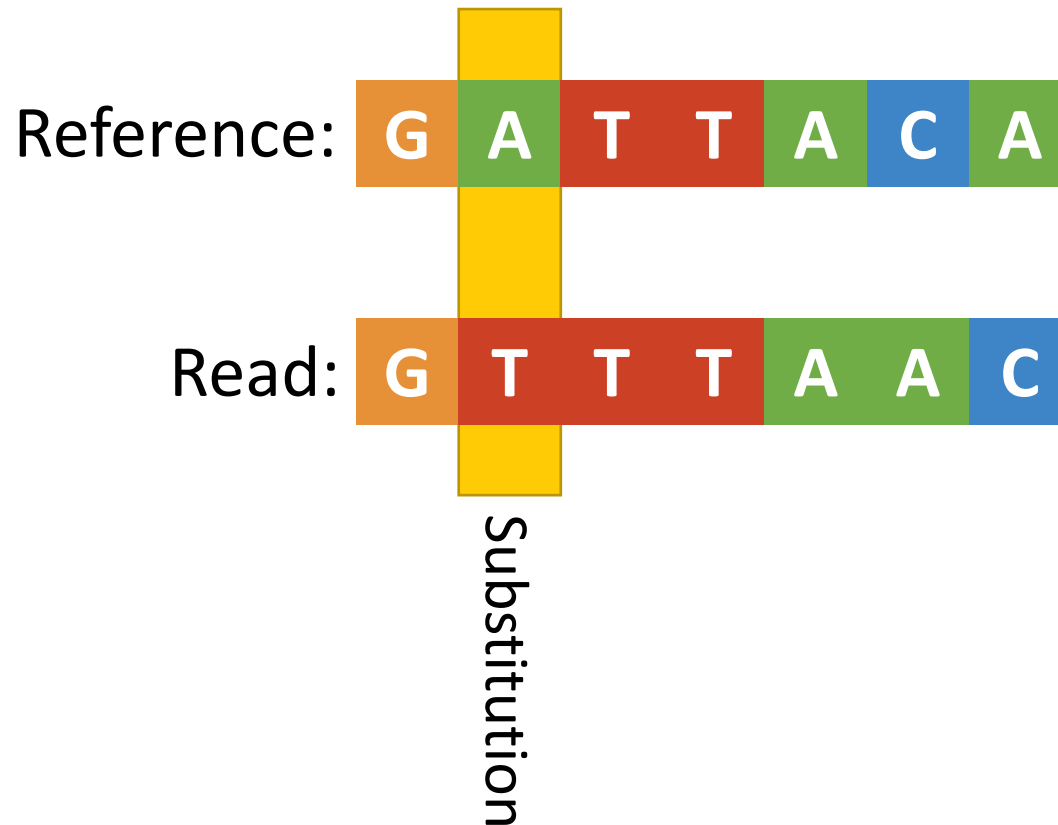   2. Algorithm
   3. Results

# Alignment: *Edit Distance*

Minimum number of edits required to transform one string to another

Reference: **G A T T A C A**

Read: **G T T T A A C**

Minimum number of edits required to transform one string to another

Reference: G A T T A C A
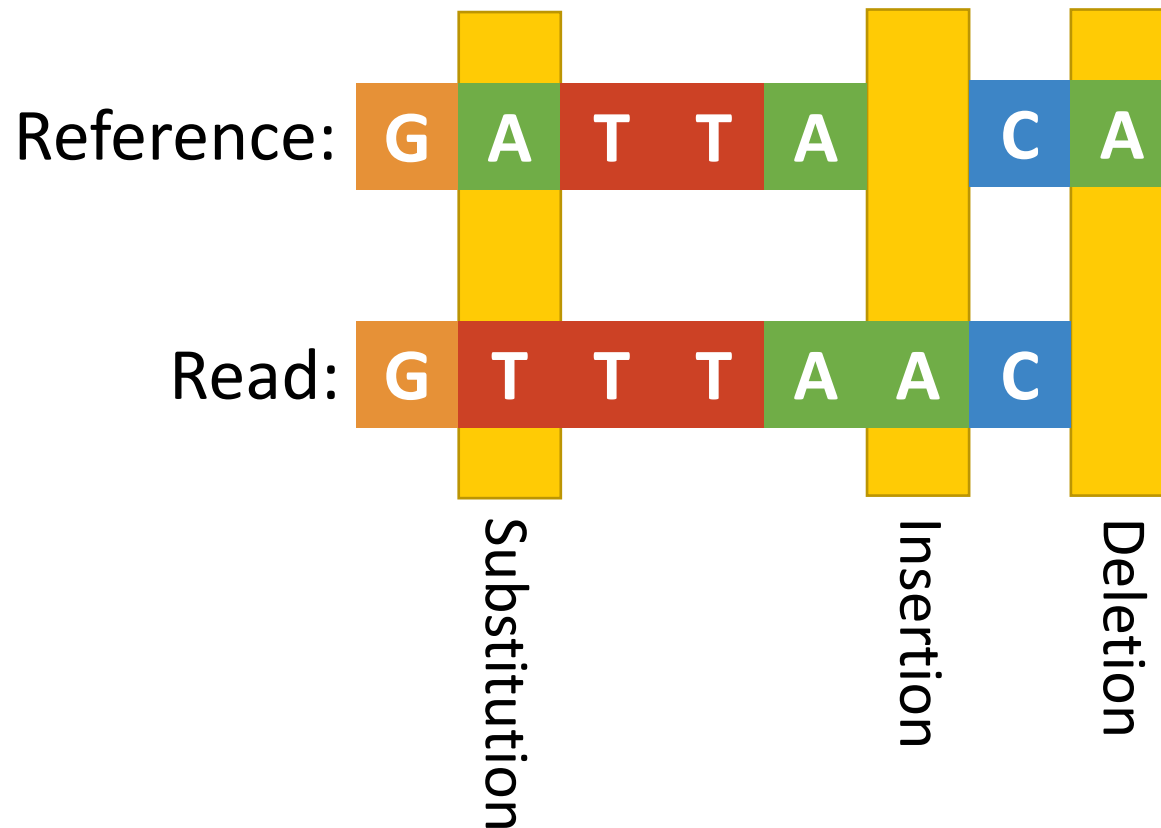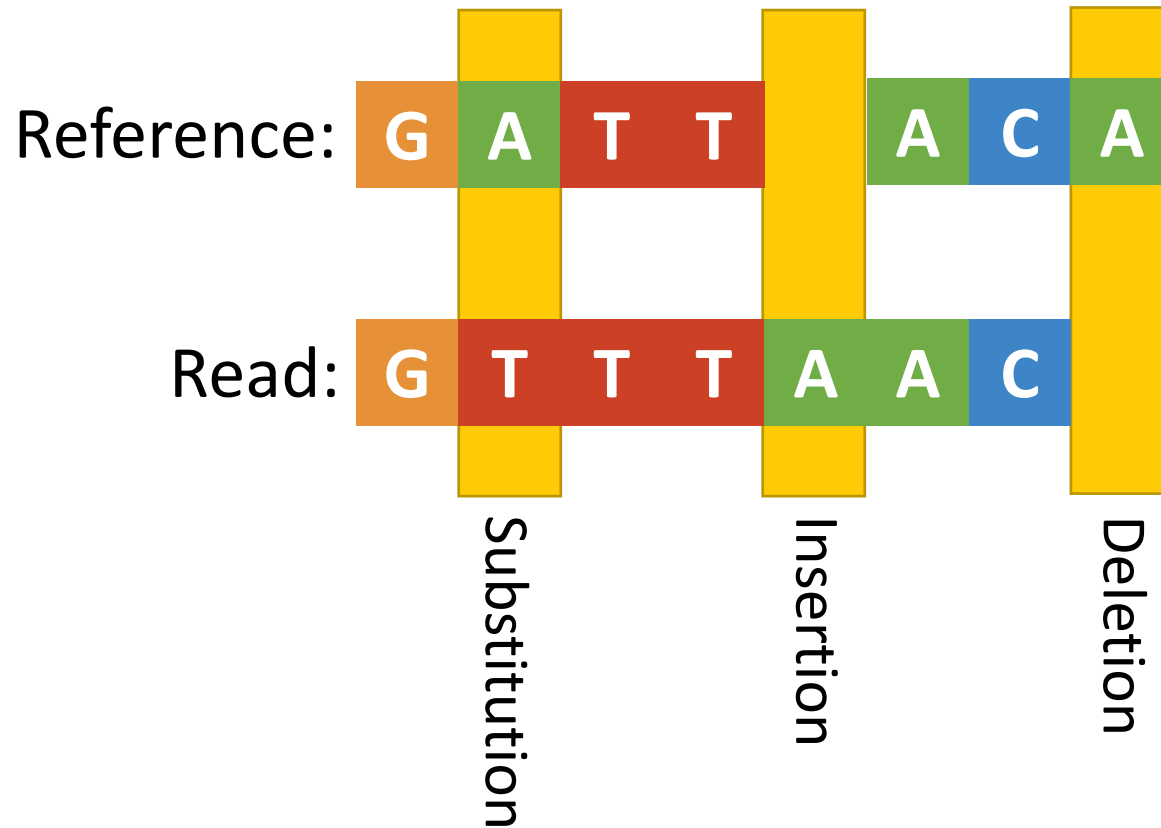
Read: G T T T A A C

Substitution

# Alignment: *Edit Distance*

Minimum number of edits required to transform one string to another

# Alignment: *Edit Distance*

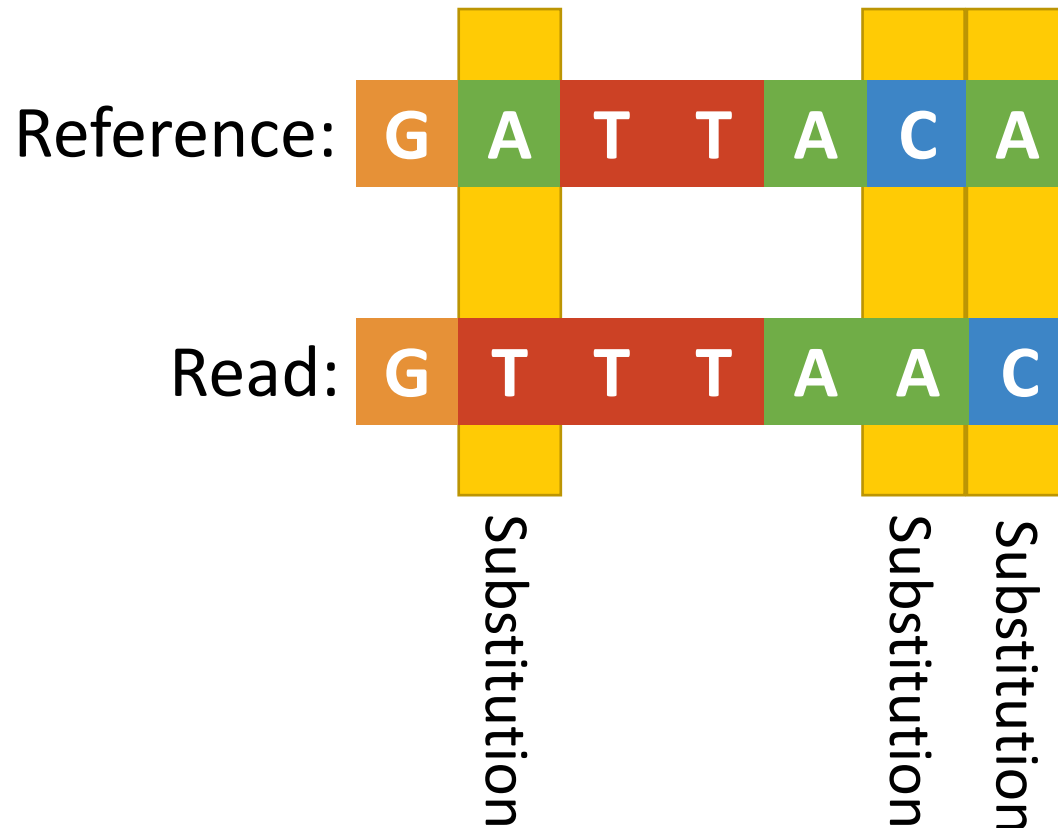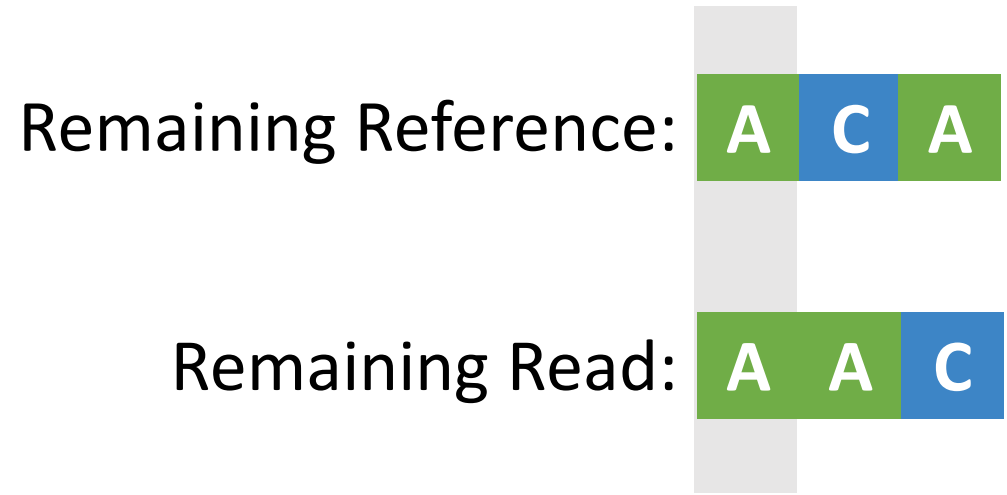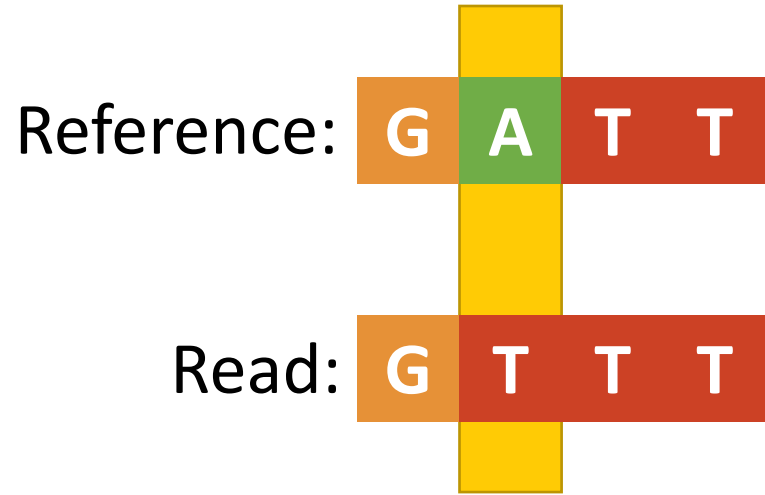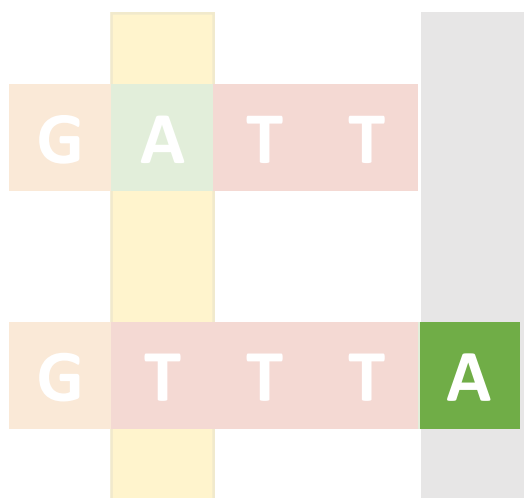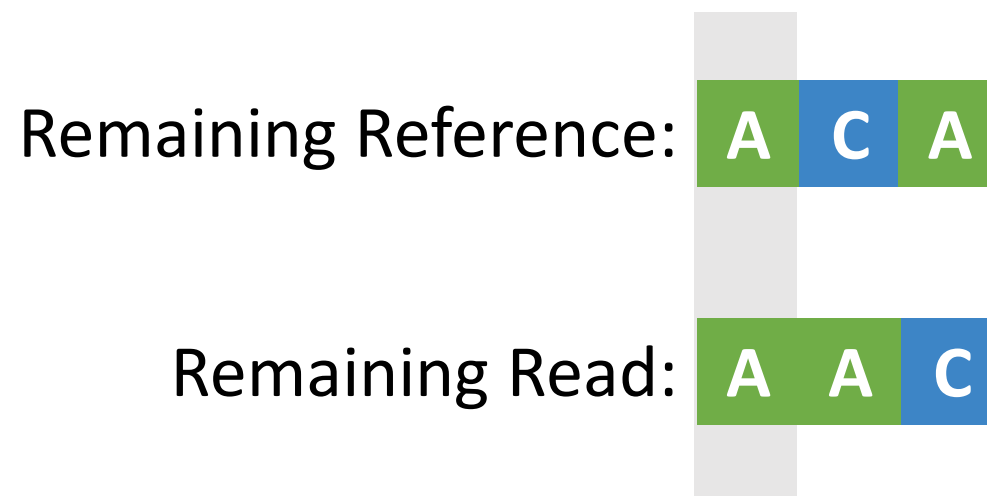Minimum number of edits required to transform one string to another

Minimum number of edits required to transform one string to another

# Alignment: *Edit Distance*

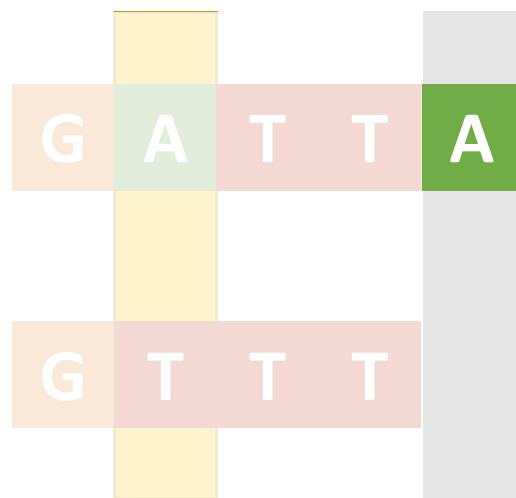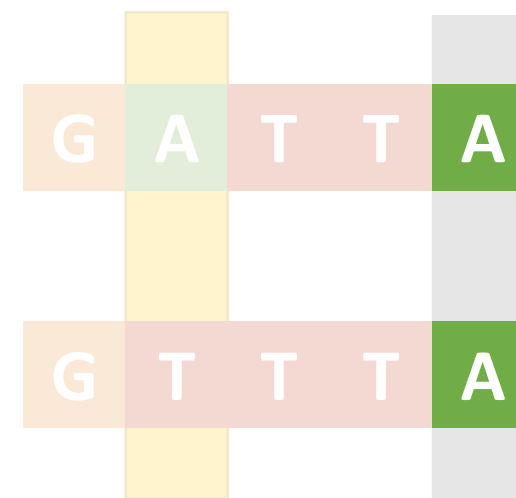Minimum number of edits required to transform one string to another

# Alignment: *Edit Distance*

Reference: G A T T

Read: G T T T

Remaining Reference: A C A

Remaining Read: A A C

# Alignment: *Edit Distance*

Reference: G A T T

Read: G T T T

Remaining Reference: A C A

Remaining Read: A A C

Insertion

Deletion

Match/Substitution

$$O(3^{|R|}) \text{ possible alignments!}$$

|R| = min( len(Read), len(Ref) )



Insertion



Deletion

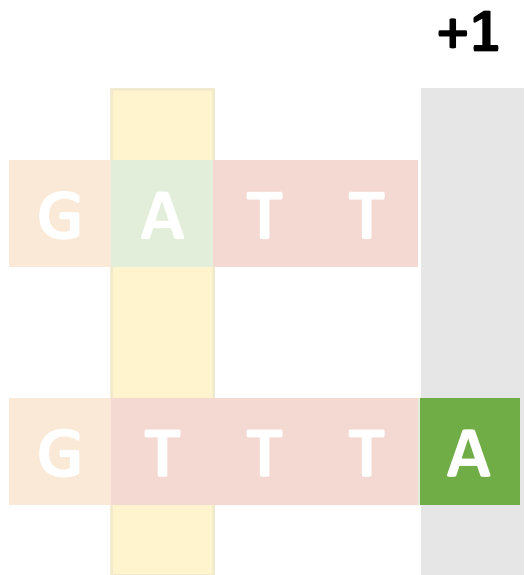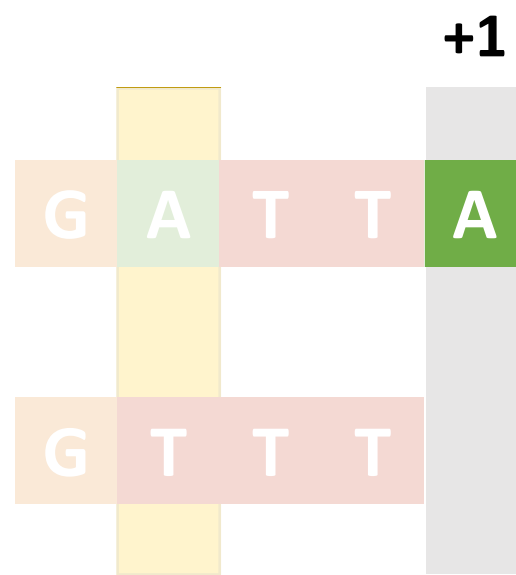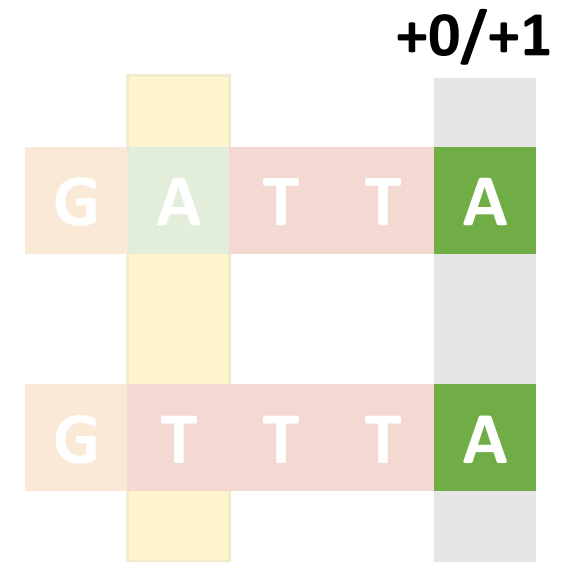

Match/Substitution

# Solution: Dynamic Programming

Edit distance is independent of prefix already aligned



Insertion

Deletion

Match/Substitution
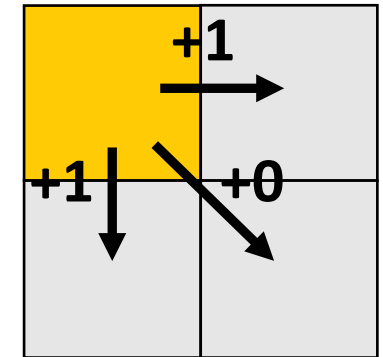
# Alignment: *Edit Distance*

# Alignment: *Edit Distance*

One larger insertion/deletion
is more likely than many small ones

# Alignment: *Affine Gap Scoring*

One larger insertion/deletion
is more likely than many small ones

# Alignment: *Affine Gap Scoring*



Normal

Inserting

Deleting

# Alignment: *Affine Gap Scoring*



Normal

Inserting

Deleting

\* 0 if match
  3 if mismatch

# Alignment: *Affine Gap Scoring*

Normal

Inserting

Deleting

# Alignment: *Affine Gap Scoring*

Normal

Inserting

Deleting

# Alignment: *Affine Gap Scoring*

One larger insertion/deletion
is more likely than many small ones

# Overview

1. Background
   1. Whole Genome Sequencing
   2. Nanopore Sequencing
   3. Read Alignment
   4. **Variant Calling**
2. n-Polymer Realigner
   1. Motivation
   2. Algorithm
   3. Results

## Pileup Calling

Identify Candidate Variants



## Full-Alignment Calling

Final Variant Calls

# Variant Calling: *Overview*

**Pileup Calling**

Identify Candidate Variants



**Full-Alignment Calling**

Final Variant Calls

# Variant Calling: *Pileup Calling*



Sliding window over a long sequence

# Overview

# Definition: *n-Polymers are homopolymers and STRs*

| T | T | T | T | T |

1-polymer          "homopolymer"

| A | T | A | T | A | T |

2-polymer          "simple tandem repeat
                   with repeat unit length 2"

| A | T | T | A | T | T |

3-polymer          ...

...

# Definition: *Copy number*

# Motivation: *INDEL Variant Calling Accuracy*



**Nanopore R9.4.1 SOTA Accuracy**

**SNPs:** 99.7% precision, 99.7% recall

**INDELs:** 92.8% precision, **76.0% recall**

# Motivation: *Inaccurate calling in repetitive regions*

# Motivation: *n-Polymer INDELs are important*

# Motivation: *Inconsistent alignments*

# Motivation: *Alignment doesn't reflect probability*

## Substitutions

|   | **A** | **C** | **G** | **T** |
|---|---|---|---|---|
| **A** | 0 | 3 | 3 | 3 |
| **C** | 3 | 0 | 3 | 3 |
| **G** | 3 | 3 | 0 | 3 |
| **T** | 3 | 3 | 3 | 0 |

## Insertions and Deletions
start: 5, extend: 2

# Overview

1. Background
   1. Whole Genome Sequencing
   2. Nanopore Sequencing
   3. Read Alignment
   4. Variant Calling

2. n-Polymer Realigner
   1. Motivation
   2. **Algorithm**
   3. Results

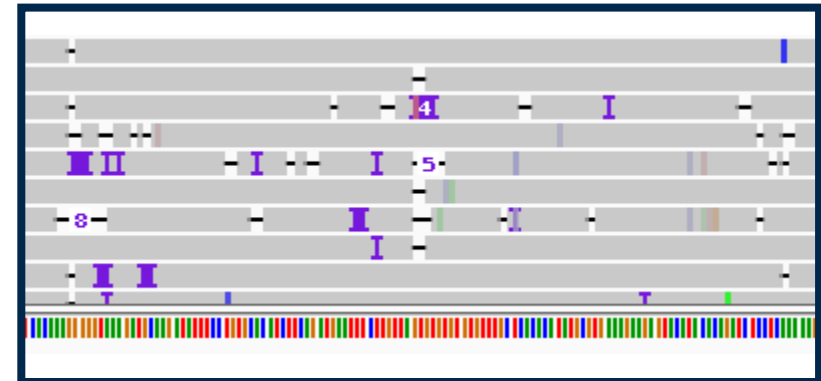# Algorithm: Use actual probabilities

a)

| | | Basecall | | |
|---|---|---|---|---|
| | A | C | G | T |
| A | 2210 | 4.75 | 17.99 | 5.12 |
| C | 4.88 | 2058 | 4.02 | 17.43 |
| G | 17.31 | 4.03 | 2064 | 4.91 |
| T | 5.17 | 18.23 | 4.86 | 2215 |

(Reference, vertical axis)

b)

| | | Basecall | | |
|---|---|---|---|---|
| | A | C | G | T |
| A | 0.01 | 6.16 | 4.82 | 6.08 |
| C | 6.06 | 0.01 | 6.25 | 4.78 |
| G | 4.79 | 6.25 | 0.01 | 6.05 |
| T | 6.07 | 4.81 | 6.13 | 0.01 |

(Reference, vertical axis)

**Figure 5: a) substitution confusion matrix $C_P$, count in millions, and b) resulting penalty matrix $P$.**

| a) | Basecall | | | |
|---|---|---|---|---|
| Reference | A | C | G | T |
| A | 2210 | 4.75 | 17.99 | 5.12 |
| C | 4.88 | 2058 | 4.02 | 17.43 |
| G | 17.31 | 4.03 | 2064 | 4.91 |
| T | 5.17 | 18.23 | 4.86 | 2215 |

| b) | Basecall | | | |
|---|---|---|---|---|
| Reference | A | C | G | T |
| A | 0.01 | 6.16 | 4.82 | 6.08 |
| C | 6.06 | 0.01 | 6.25 | 4.78 |
| G | 4.79 | 6.25 | 0.01 | 6.05 |
| T | 6.07 | 4.81 | 6.13 | 0.01 |

**Figure 5: a) substitution confusion matrix $C_P$, count in millions, and b) resulting penalty matrix $P$.**

$$P[i,j] \approx -\log \mathbb{P}(x[i] \rightarrow x[j]) \approx -\log \frac{C_P[i,j] + \epsilon}{\text{sum}(C_P[i,:]) + \epsilon}$$

# Algorithm: Context-dependent penalties

## Original

### Substitutions

|   | A | C | G | T |
|---|---|---|---|---|
| **A** | 0 | 3 | 3 | 3 |
| **C** | 3 | 0 | 3 | 3 |
| **G** | 3 | 3 | 0 | 3 |
| **T** | 3 | 3 | 3 | 0 |

### Insertions and Deletions
start: 5, extend: 2

## nPoRe

### Substitutions

|   | A | C | G | T |
|---|---|---|---|---|
| **A** | 0 | 6 | 5 | 6 |
| **C** | 5 | 0 | 6 | 4 |
| **G** | 4 | 6 | 0 | 6 |
| **T** | 6 | 5 | 6 | 0 |

### Insertions and Deletions
start: 7, extend: 2

### n-Polymers
lookup table

## Original

### Substitutions

|   | A | C | G | T |
|---|---|---|---|---|
| A | 0 | 3 | 3 | 3 |
| C | 3 | 0 | 3 | 3 |
| G | 3 | 3 | 0 | 3 |
| T | 3 | 3 | 3 | 0 |

### Insertions and Deletions
start: 5, extend: 2

## nPoRe

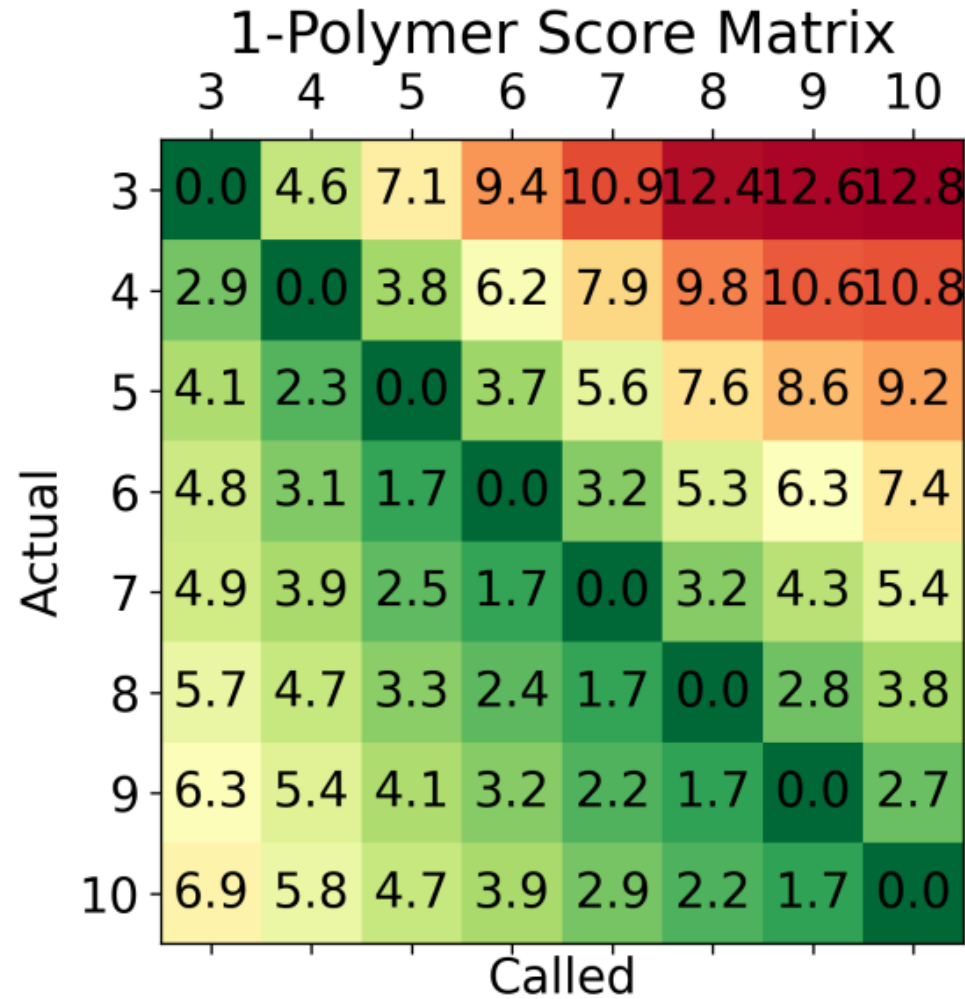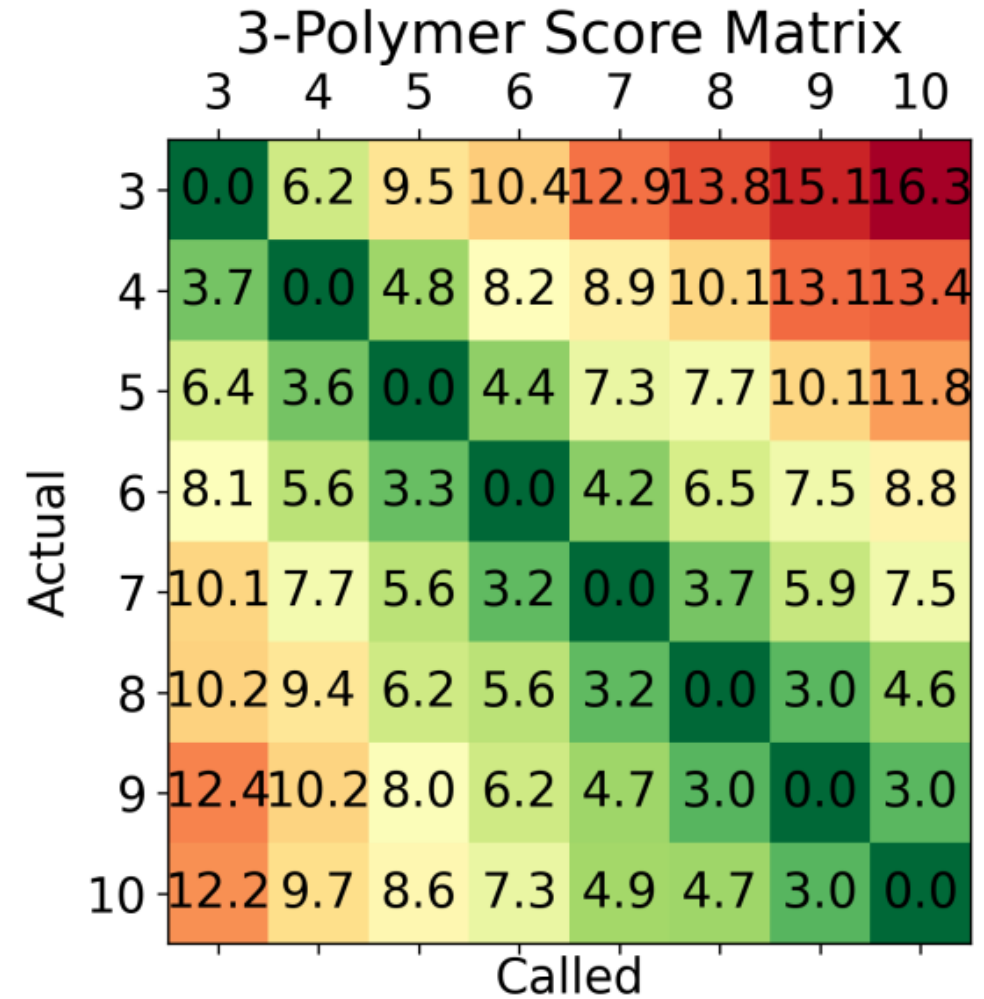### Substitutions

|   | A | C | G | T |
|---|---|---|---|---|
| A | 0 | 6 | 5 | 6 |
| C | 5 | 0 | 6 | 4 |
| G | 4 | 6 | 0 | 6 |
| T | 6 | 5 | 6 | 0 |

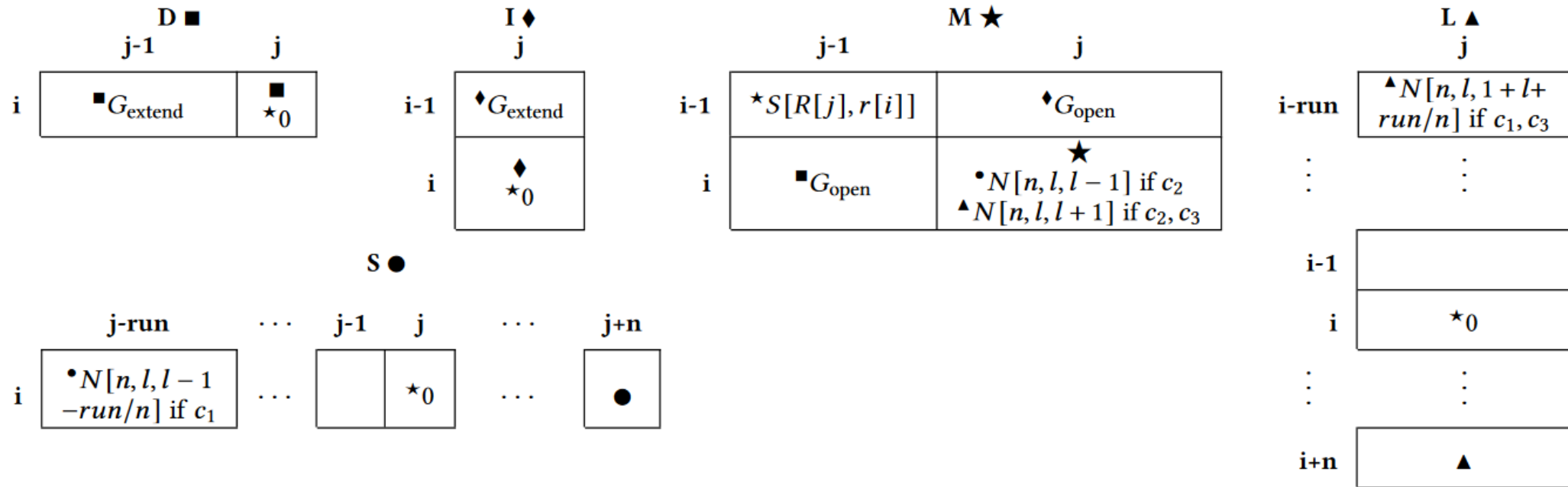### Insertions and Deletions
start: 7, extend: 2

### n-Polymers
lookup table

$$N[n, i, j] \approx -\log \mathbb{P}(n, i, j) \approx -\log \frac{C_N[n, i, j] + \epsilon}{\text{sum}(C_N[n, i, :]) + \epsilon}$$
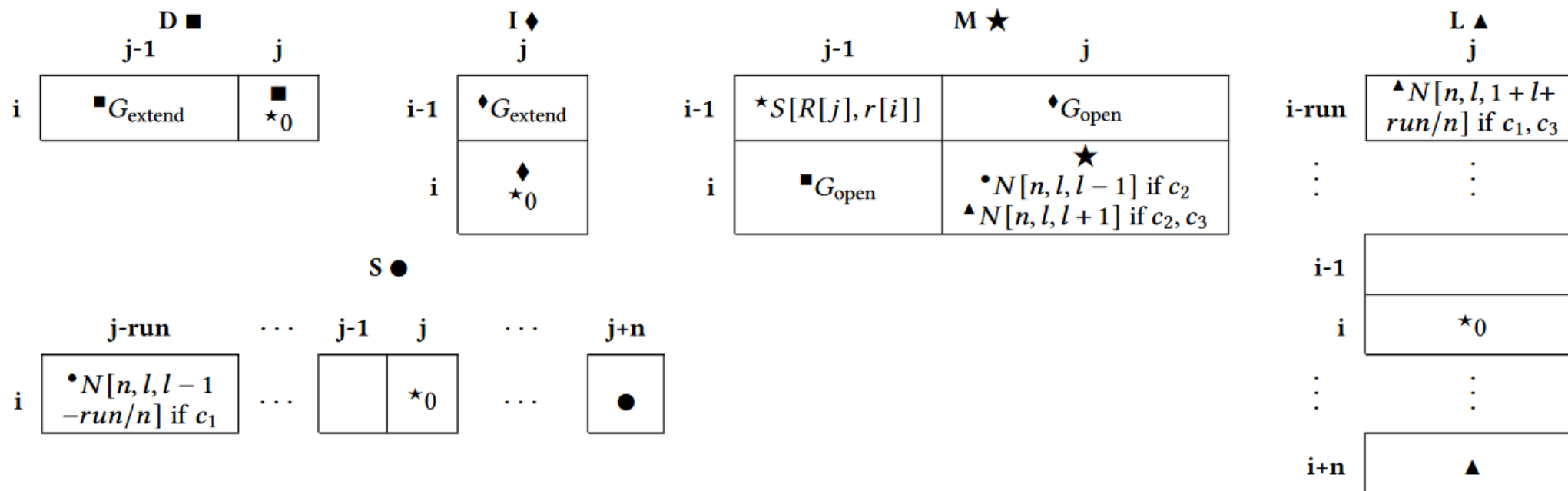
# Algorithm: *Lookup-based gap penalties*

# Algorithm: *Requires five state tables*

**D ■**

|  | j-1 | j |
|---|---|---|
| i | $^{■}G_{\text{extend}}$ | $^{★}0$ |

**I ♦**

|  | j |
|---|---|
| i-1 | $^{♦}G_{\text{extend}}$ |
| i | $^{♦}_{★}0$ |

**M ★**

|  | j-1 | j |
|---|---|---|
| i-1 | $^{★}S[R[j], r[i]]$ | $^{♦}G_{\text{open}}$ |
| i | $^{■}G_{\text{open}}$ | $^{★}$ $^{●}N[n, l, l-1]$ if $c_2$ $^{▲}N[n, l, l+1]$ if $c_2, c_3$ |

**L ▲**

|  | j |
|---|---|
| i-run | $^{▲}N[n, l, 1+l+$ $run/n]$ if $c_1, c_3$ |
| ⋮ | ⋮ |
| i-1 |  |
| i | $^{★}0$ |
| ⋮ | ⋮ |
| i+n | ▲ |

**S ●**

|  | j-run | ⋯ | j-1 | j | ⋯ | j+n |
|---|---|---|---|---|---|---|
| i | $^{●}N[n, l, l-1$ $-run/n]$ if $c_1$ | ⋯ |  | $^{★}0$ | ⋯ | ● |

$$c_1 = \quad l > 0 \qquad\qquad\qquad \text{start of repeat unit}$$
$$c_2 = \quad l > 0 \text{ and } idx == 0 \qquad \text{start of } n\text{-polymer}$$
$$c_3 = \quad R[j+1 : j+1+n] ==$$
$$r[i+1 : i+1+n] \qquad \text{next } n \text{ bases of } r \text{ match } R$$

# Algorithm: *Properties*



Penalty depends on *l*, *l'*

Penalty

*l* = 3  A A A
*l* = 5  A A A A A
*l* = 8  A A A A A A A A

Penalty depends on *n*, only CNVs allowed

Penalty

*n* = 1  A A A A A A A A
*n* = 2  A T A T A T A T
*n* = 3  A A T A A T A A T

# Algorithm: *Properties*



Sequence is not considered



Overlaps are allowed

**Figure 7: Follow banding matrix transformation** $A \rightarrow B$

# Overview

1. Background
   1. Whole Genome Sequencing
   2. Nanopore Sequencing
   3. Read Alignment
   4. Variant Calling
2. n-Polymer Realigner
   1. Motivation
   2. Algorithm
   3. **Results**

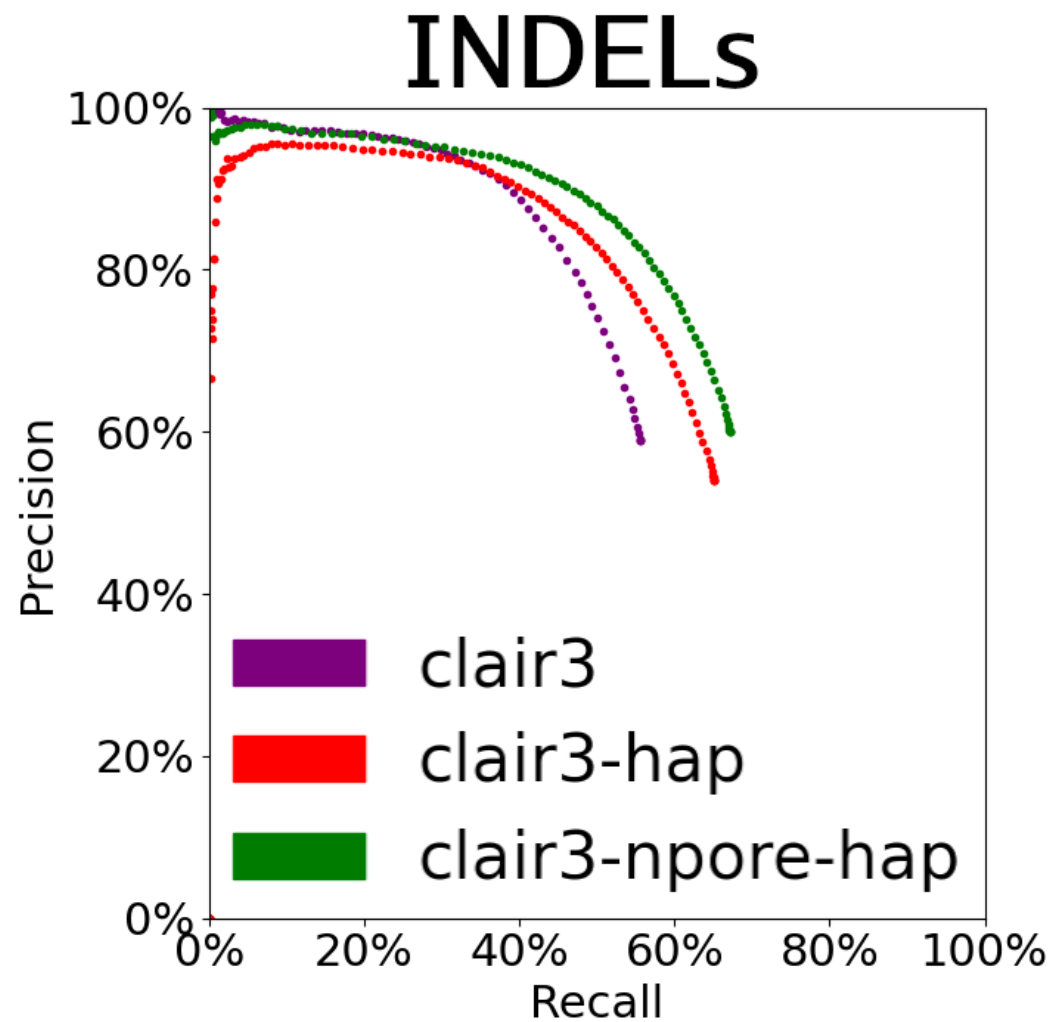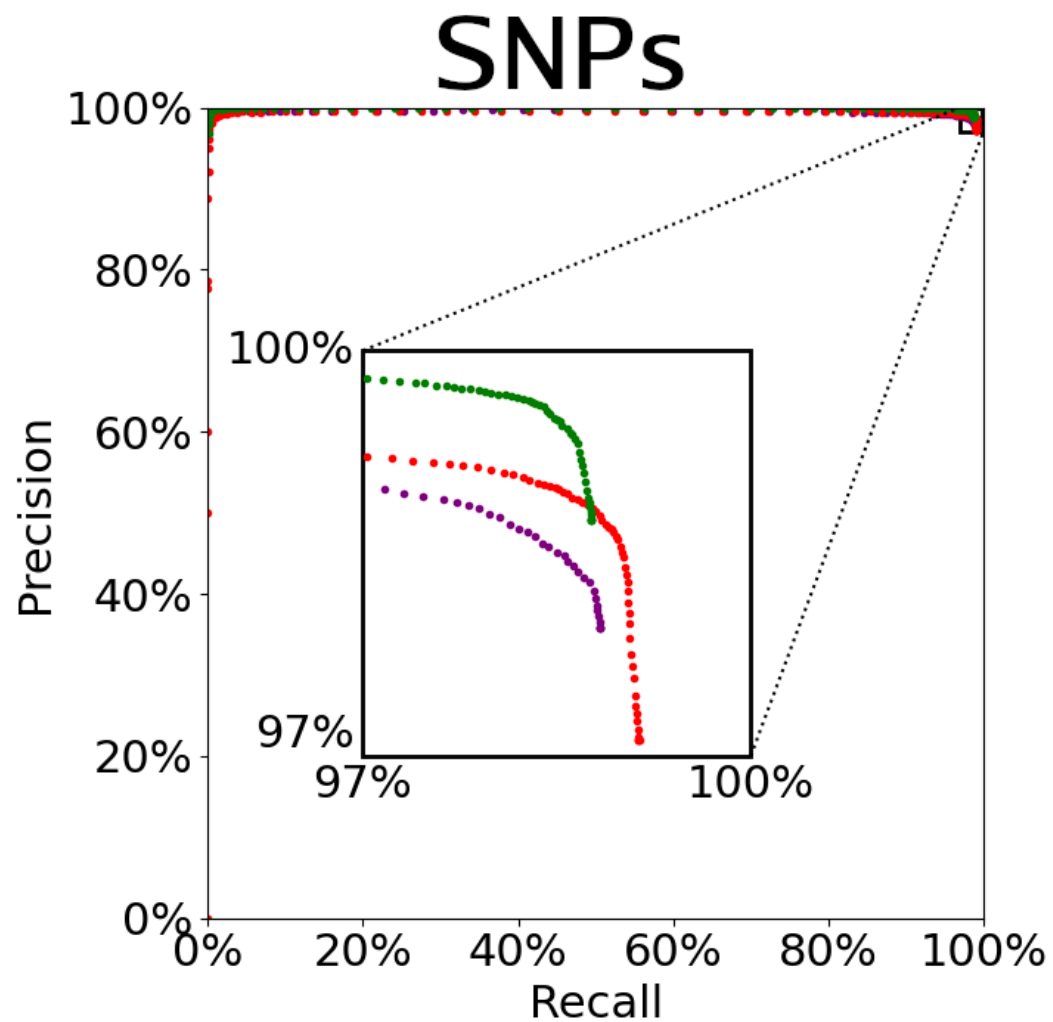**Figure 9: Read concordance: Gini purity histograms for a) pileup columns and b) insertions**

$$GP = \sum_{i=1}^{n} P(i)^2$$

# Results: *GitHub code*

*https://github.com/TimD1/nPoRe*

# Funding

# Questions?

Thanks for listening!

*If you'd like to talk about research, please reach out!*

Email: timdunn@umich.edu

Twitter: @T1MD1