## *n*PoRe: *n*-Polymer Realigner for improved pileup-based variant calling Tim Dunn, David Blaauw, Reetu Das, Satish Narayanasamy **T**University of Michigan, CSE

# **Background:**

Variant calling requires the identification of mutations present in sequenced reads relative to an expected reference genome. Small variants are classified into Single Nucleotide Polymorphisms, or SNPs, and Insertions/Deletions, or INDELs.



Although nanopore sequencing can identify SNPs with high accuracy, a current weakness is recalling small INDELs.

#### Nanopore R9.4.1 Variant Calling Accuracy<sup>[1]</sup>

## **Motivation:**

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

Nanopore sequencers frequently mis-call the length of repeated low-complexity DNA sequences. Existing aligners do not take this failure mode into account, leading to inconsistent alignments which complicate variant calling.



Most variant calling errors can be attributed to INDEL errors within n-polymer regions.



Most true INDEL mutations are variations in the number of copies of an n-polymer.

#### We define:

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

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



*n* = 2

- *n*, the length of the repeated unit
- *I*, the expected number of repeated copies
- *I'*, the measured number of repeated copies
- "n-polymer": repeated sequence where  $n \le 6$  and  $l \ge 3$

#### **Algorithm: Baseline: Affine Gap Smith-Waterman**



| 1. (  | Generate n-polymer reference annotations           | 2. Define conditions for valid n-polymer INDELs |   |
|-------|--|---|---|
| Refei | ence: A T A T A T A T T T T T T A A A G C G C G C  | $c_1 = l > 0$<br>$c_2 = l > 0$ and $idx == 0$   | start of repeat unit start of <i>n</i> -polymer |
| n=1   | <i>l</i> : 0 0 0 0 0 0 5 5 5 5 5 3 3 3 0 0 0 0 0 0 | $c_3 = R[j+1:j+1+n] =$                          |   |









Sequence is not considered



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

## bioRxiv Pre-Print

https://www.biorxiv.org/content/10.1101/2022.02.15.480561v1

## References

[1] Shafin, et. al. "Haplotype-aware variant calling with PEPPER Margin DeepVariant enables high accuracy in nanopore long-reads." Nature Methods, 2021. https://github.com/kishwarshafin/pepper [2] Zheng, et. al. "Symphonizing pileup and full-alignment for deep learning-based long-read variant calling." bioRxiv, 2022. [3] Li, Heng. "Minimap2: pairwise alignment for nucleotide sequences."

This project was supported by the Kahn Foundation, NSF Grant 2030454, and NSF Graduate Research Fellowship 1841052.