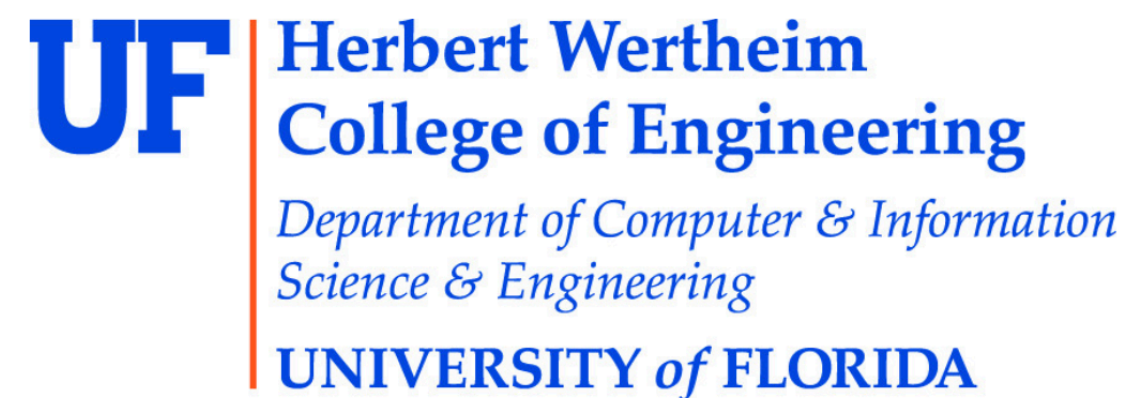


# Pan-genomic matching statistics for targeted nanopore sequencing

Omar Ahmed<sup>1</sup>, Massimiliano Rossi<sup>2</sup>, Sam Kovaka<sup>1</sup>, Michael C. Schatz<sup>1</sup>, Travis Gagie<sup>3</sup>, Christina Boucher<sup>2</sup>, and Ben Langmead<sup>1</sup>

RECOMB-Seq  
August 27-28, 2021



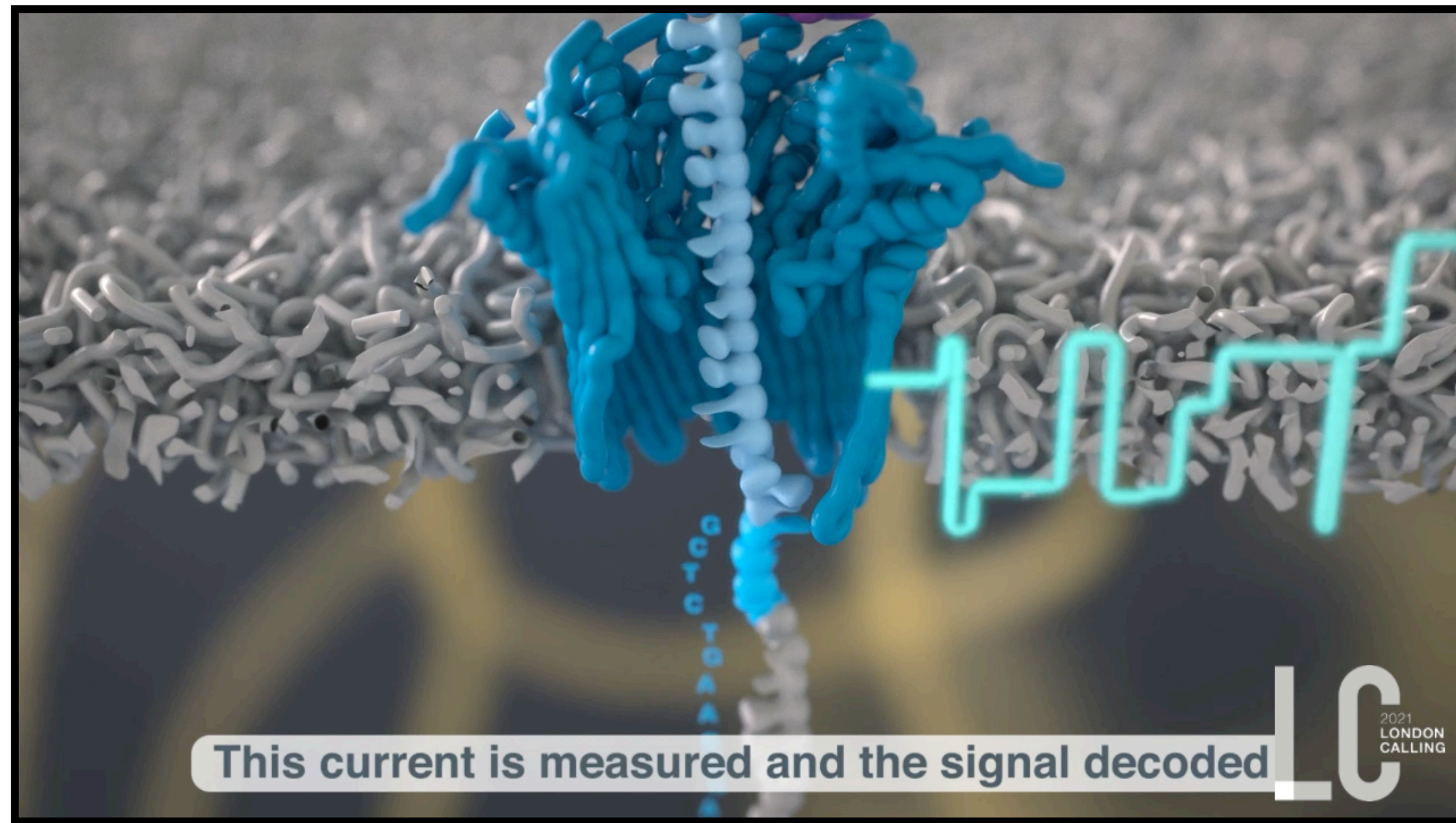
<sup>1</sup> Department of Computer Science, Johns Hopkins University, Baltimore, MD, USA

<sup>2</sup> Department of Computer and Information Science and Engineering, University of Florida, Gainesville, FL, USA

<sup>3</sup> Faculty of Computer Science, Dalhousie University, Halifax, NS, CAN

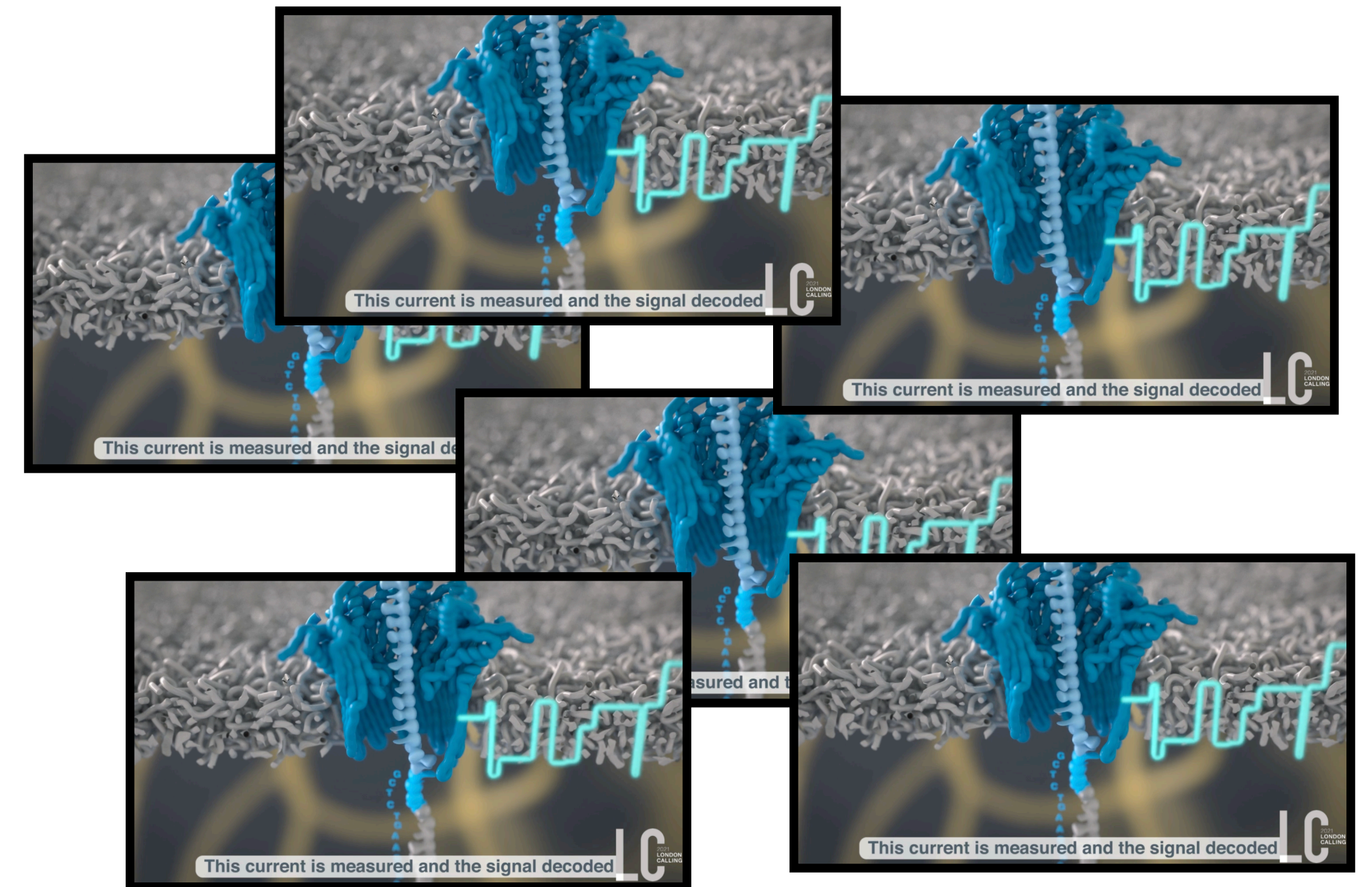


# Nanopore Sequencing



- ▶ Allows users to perform **targeted sequencing** using software

**Key Idea:** The software needs to be fast enough to keep up with incoming signal from numerous pores.



- ▶ MinION Flowcells have 512 pores<sup>1</sup>

<sup>1</sup><https://nanoporetech.com/products/comparison>

Video From: <https://nanoporetech.com/about-us/news/towards-real-time-targeting-enrichment-or-other-sampling-nanopore-sequencing-devices>

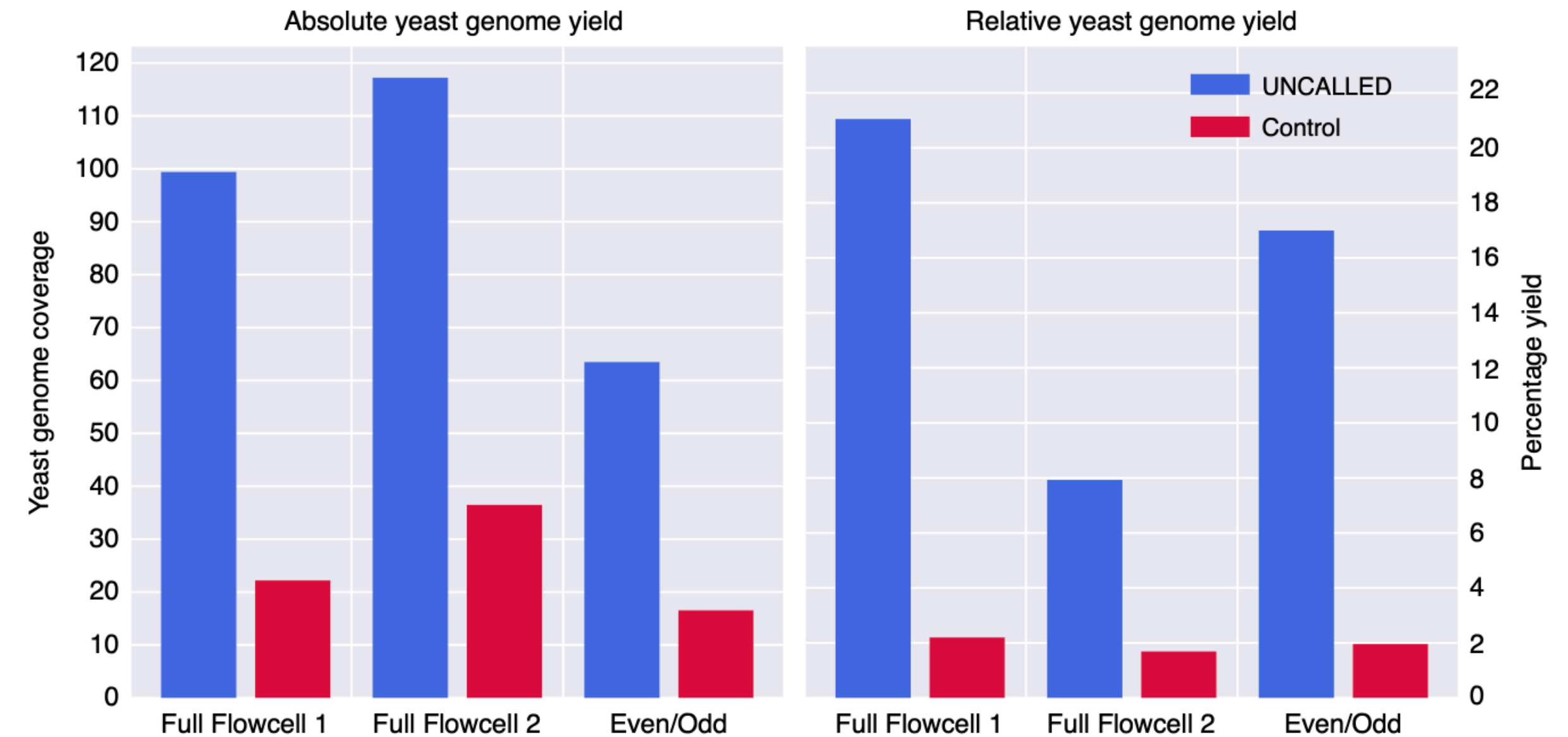


# Targeted Sequencing with Nanopore Sequencing

► Recent methods such as UNCALLED<sup>1</sup> and Readfish<sup>2</sup> have performed **targeted sequencing**

*Results from UNCALLED<sup>1</sup> paper*

- **Goal:** Sequence a mock community and target the **yeast** reads and eject the **microbial** reads
- Using targeted sequencing yielded **higher coverage** of the yeast genome, and **higher percent yield**



*“UNCALLED’s performance degrades as references become larger and more repetitive”<sup>1</sup>*



**Motivation:** A need for faster methods to classify reads against large, repetitive databases

**Why is it important?**

① Improves targeting accuracy

② Target genomes of unassembled strains

<sup>1</sup>Kovaka, S., Fan, Y., Ni, B. *et al.* Targeted nanopore sequencing by real-time mapping of raw electrical signal with UNCALLED. *Nat Biotechnol* **39**, 431–441 (2021).

<sup>2</sup>Payne, A., Holmes, N., Clarke, T. *et al.* Readfish enables targeted nanopore sequencing of gigabase-sized genomes. *Nat Biotechnol* **39**, 442–450 (2021).

# Method Overview

## SPUMONI - Streaming Pseudo MONI<sup>1</sup>

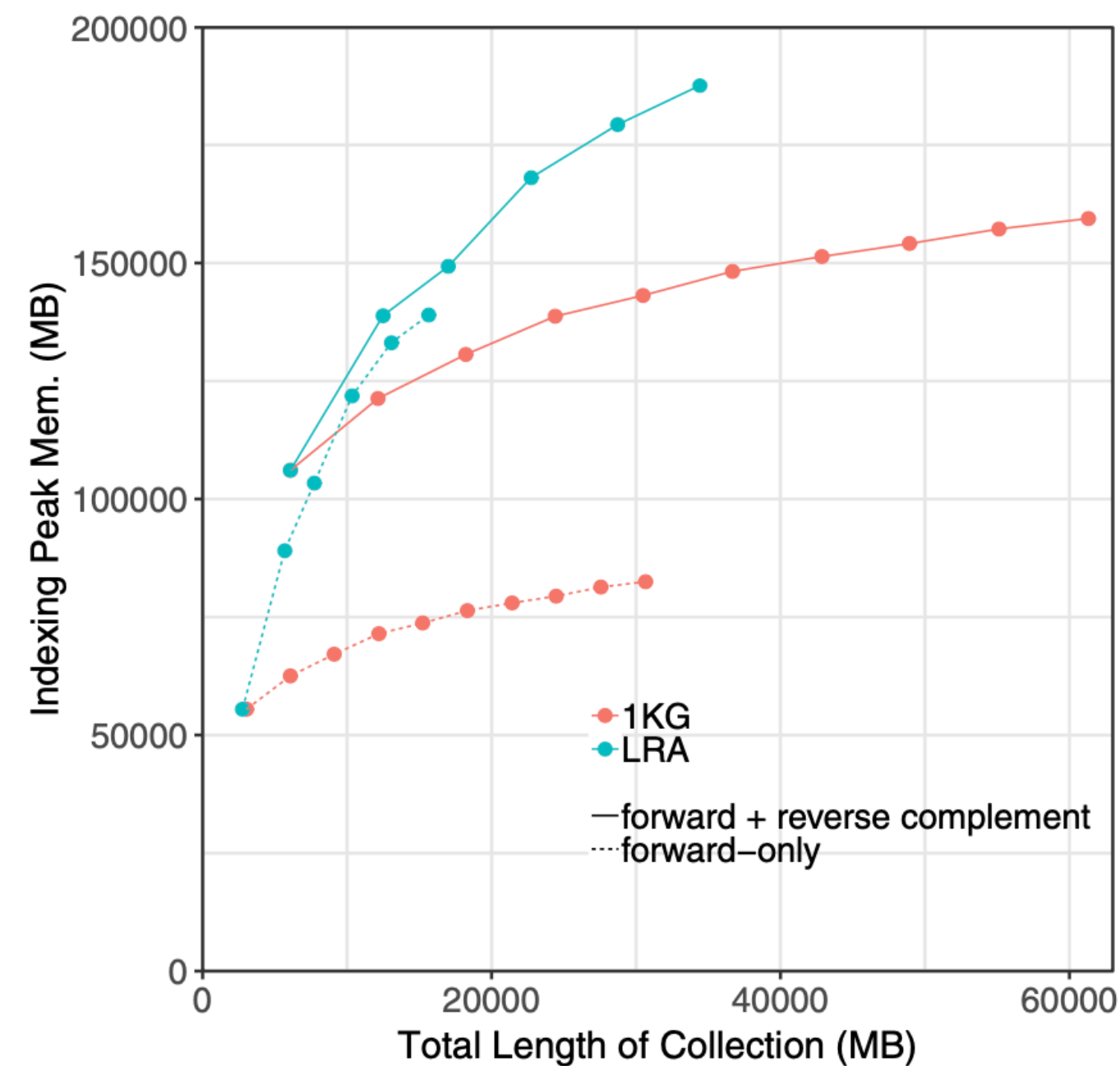
- ▶ Makes rapid targeting decisions based on input database
  - **Key Intuition:** A read's MS/PMLs with respect to a reference can reveal if there appears to be “good” approximate match to reference
- ▶ Uses the r-index<sup>2</sup> to enable efficient indexing of large, repetitive collections
  - Number of runs in BWT,  $r$ , typically grows sub-linearly w.r.t to length of input sequence,  $n$ .
- ▶ Extends MONI<sup>1</sup> in two key areas
  - Adds a “null index” and hypothesis testing framework for finding “significant” matches
  - Replaces MONI's “batch” matching statistic (MS) algorithm with a faster, streaming algorithm
    - ▶ Calculates new quantity called pseudo-matching lengths (PMLs)

<sup>1</sup>Rossi, M., Oliva, M., Langmead, B., Gagie, T., & Boucher, C. (2021). MONI: A pangenomics index for finding MEMs. *Proc. RECOMB*.

<sup>2</sup>Mun, T., Kuhnle, A., Boucher, C., Gagie, T., Langmead, B., and Manzini, G. (2020). Matching reads to many genomes with the r-index. *J. Comput. Biol.* 27, 514–518

# Strengths of the $r$ -index

- ▶ Mäkinen and Navarro's<sup>1</sup>  $O(r)$  rank data-structure over the BWT combined with Gagie et al.'s<sup>2</sup>  $O(r)$  suffix array sampling make up the  $r$ -index
  - This combination allows efficient queries to **count** the number of occurrences and **locate** those occurrences
- ▶ Peak memory of  $r$ -index construction shows a sub-linear trend showing the strength of the  $r$ -index for indexing large, repetitive collections (Mun et al.)<sup>3</sup>



<sup>1</sup> V. Mäkinen and G. Navarro. (2007). Rank and select revisited and extended. Theoretical Computer Science, 387. pp. 332–347.

<sup>2</sup> T. Gagie, G. Navarro, and N. Prezza. (2020). Fully functional suffix trees and optimal text searching in bwt-runs bounded space, J. ACM, 67 (2020), pp. 2:1–2:54.

<sup>3</sup> Mun, T., Kuhnle, A., Boucher, C., Gagie, T., Langmead, B., and Manzini, G. (2020). Matching reads to many genomes with the  $r$ -index. J. Comput. Biol. 27, 514–518

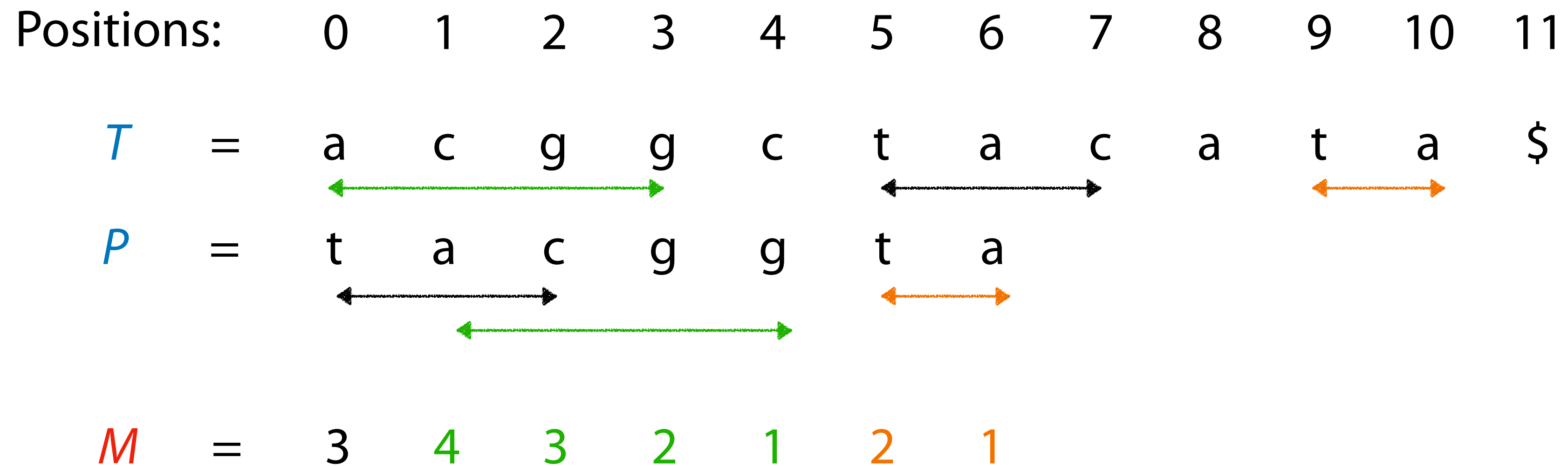
# Matching Statistics

▶ The matching statistics of  $P$  w.r.t to  $T$  is an array  $M$  of length  $m$  where  $M[i]$  is the length of the longest prefix of the pattern  $P[i..m]$  that occurs in text  $T$

○ Let  $T$  be a text of length  $n$ , and  $P$  be a pattern of length  $m$

▶ Think of matching statistics like they are half MEMs (Maximally Exact Matches)

○ Example:



# Matching Statistics - Comparing Approaches

► MONI<sup>1</sup> - Two Pass Algorithm → Calculates Matching Statistics

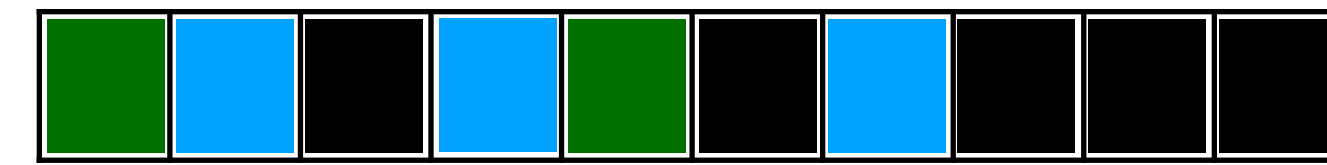
Case 1  $p[i] = BWT[j]$

Case 2a  $p[i] \neq BWT[j]$

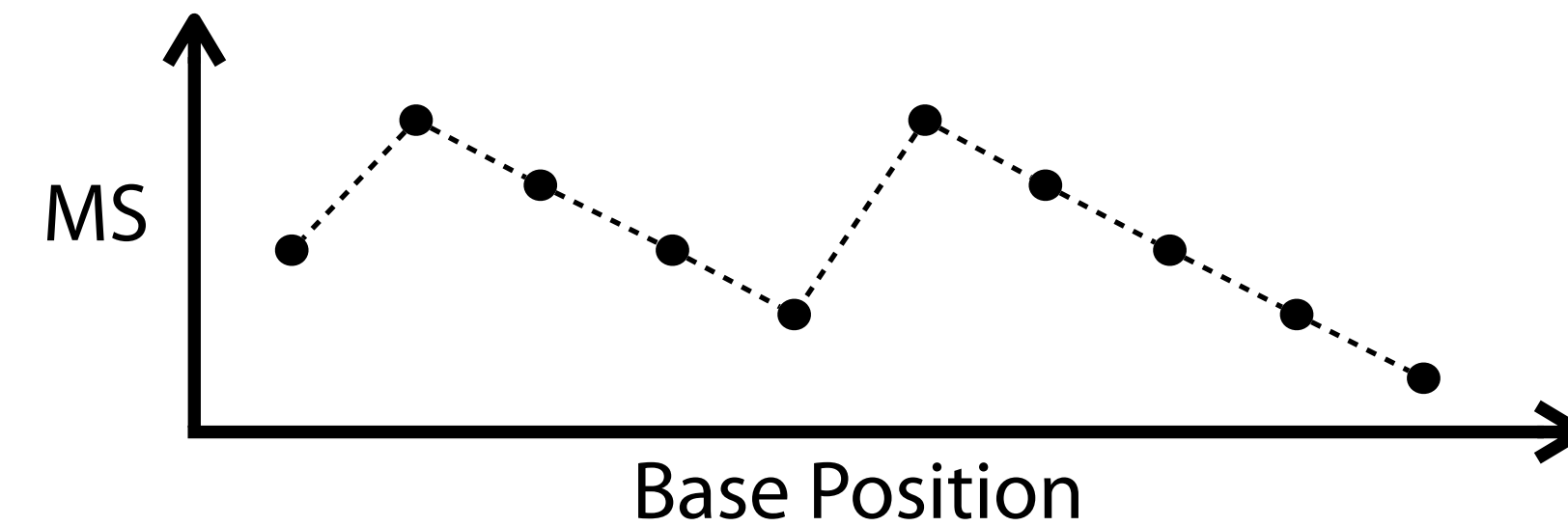
*but can still extend current match by 1*

Case 2b  $p[i] \neq BWT[j]$

*but can possibly retain some of current match*



3 5 4 3 2 5 4 3 2 1



► SPUMONI - Online Algorithm → Calculates Pseudo-Matching Lengths

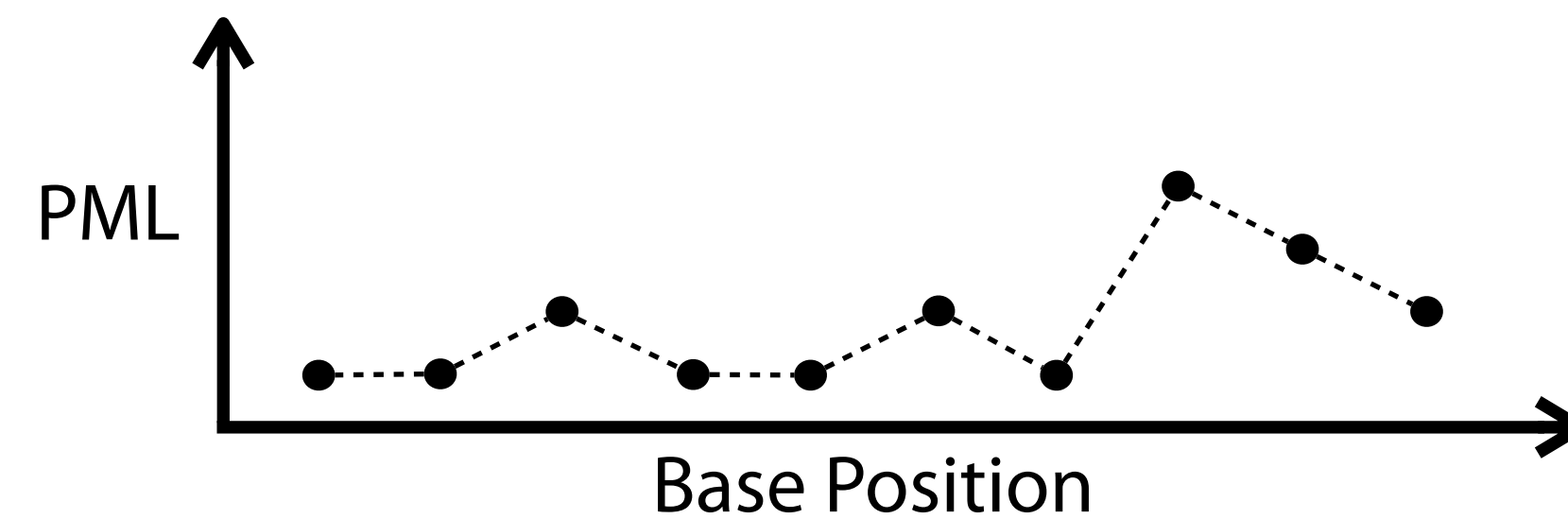
Case 1  $p[i] = BWT[j]$

Case 2  $p[i] \neq BWT[j]$

*so we set the length down to 0*



0 0 1 0 0 1 0 3 2 1



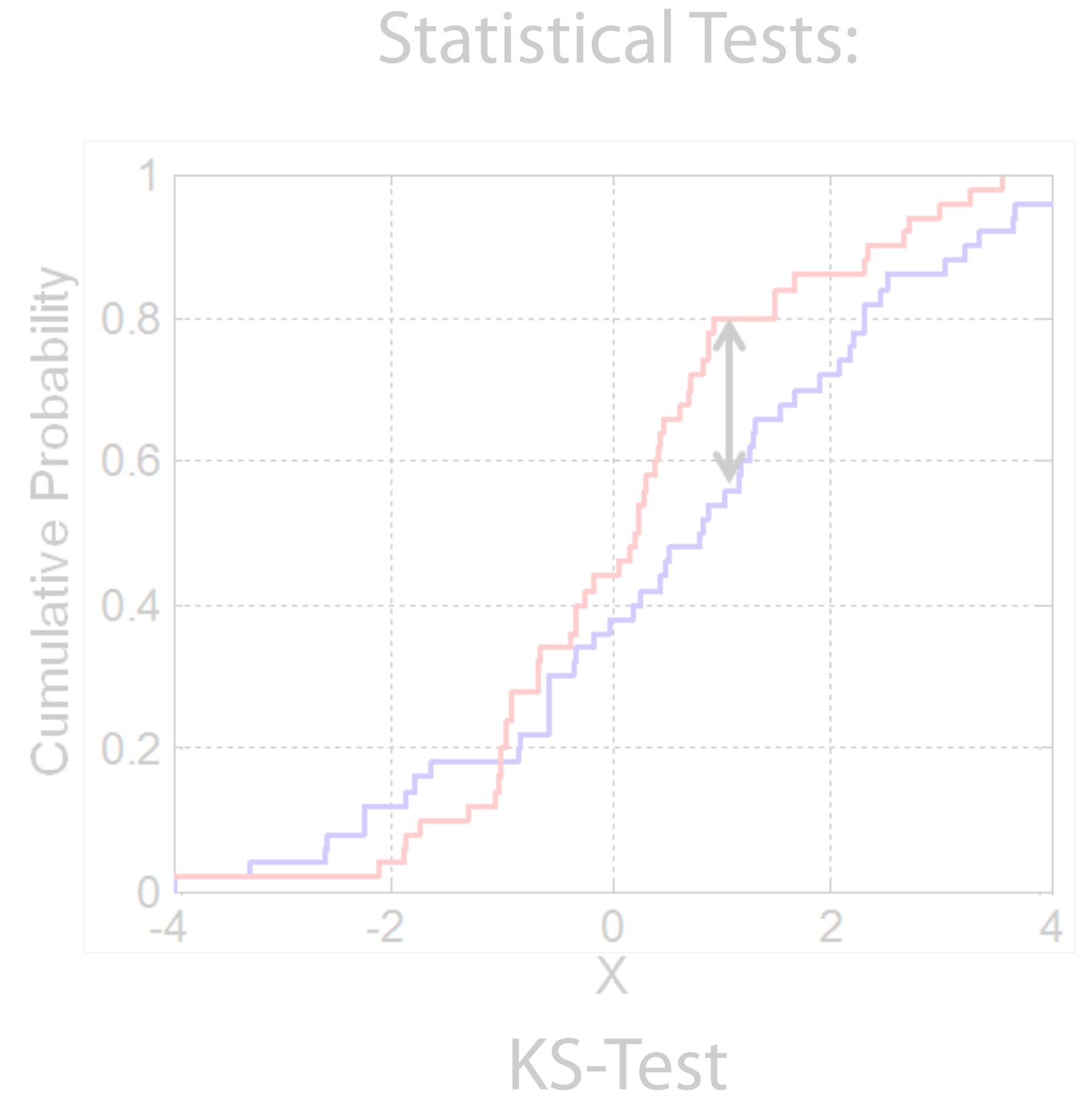
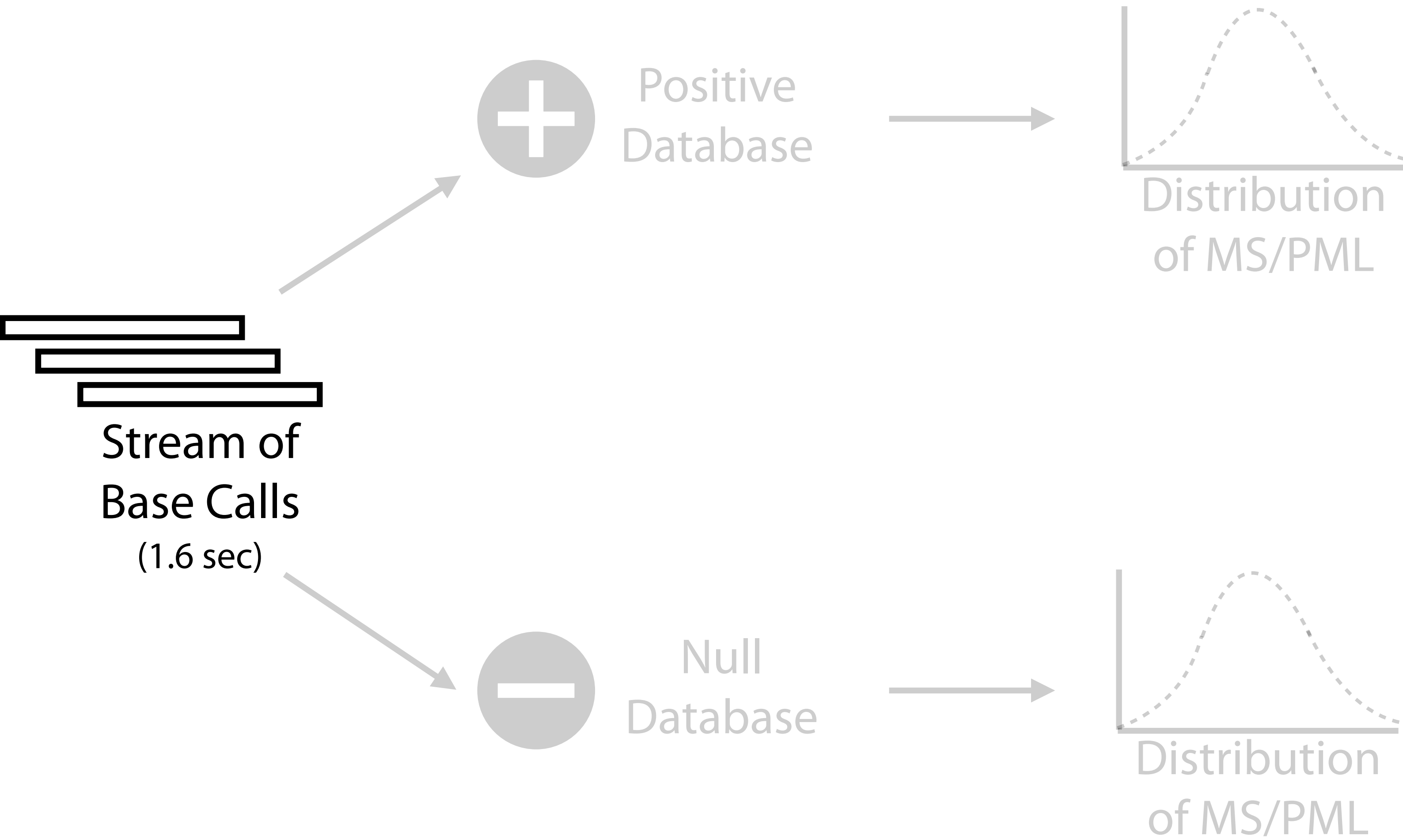
► Simplification leads to smaller indexes, so SPUMONI uses ~3X less memory and is ~3X faster

How will each quantity perform when it comes to classification?

<sup>1</sup>Rossi, M., Oliva, M., Langmead, B., Gagie, T., & Boucher, C. (2021). MONI: A pangenomics index for finding MEMs. *Proc. RECOMB.*



# SPUMONI Approach

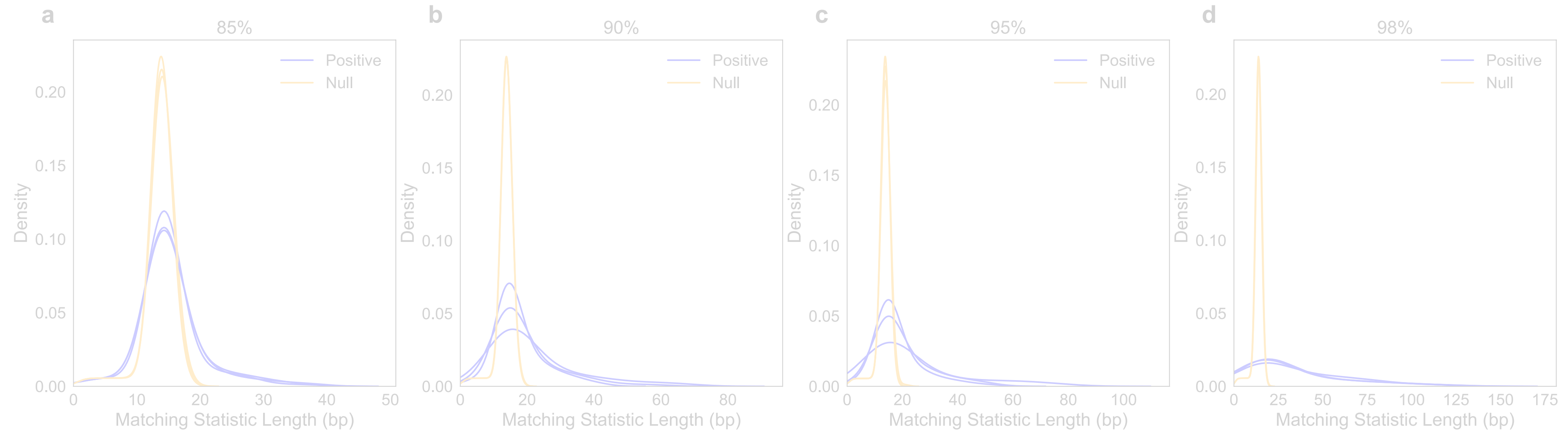




# Results - Varying Mean Read Accuracy

► **Question:** How does sequencing error affect SPUMONI's classification?

*Distributions of MS for E. coli reads against a Zymo Mock Community Index at Varying Read Accuracy*



► **Answer:** SPUMONI's accuracy  $\geq 99.5\%$  for reads with 90% accuracy or more

○ At 85% read accuracy, SPUMONI's accuracy was 95.43 and minimap2's accuracy was 99.16

# Results - Real Mock Community Experiment

- ▶ **Question:** Can using a pan-genome reference allow us to target a particular strain that is not present in the reference? and how does it compare on time and memory?
- Using **real** mock community reads<sup>1</sup> where we want to “target” the **yeast** reads, and eject all the **microbial** species

**Table 3. Comparing SPUMONI and minimap2 across various metrics on Real ZymoMC Reads**

Accuracy, throughput and index size on real mock community reads

Reference:	One genome ref			Pan-genome ref			
Reference size:	56 MB	56 MB	28 MB	31 GB	31 GB	16 GB	
Approach:	SPUMONI-ms	SPUMONI	minimap2	SPUMONI-ms	SPUMONI	minimap2	
Accuracy	81.64	86.72	87.82	94.62	96.02	97.52	
Precision	100.00	100.00	100.00	100.00	100.00	99.96	
Recall	81.39	86.54	87.66	94.55	95.97	97.53	
Specificity	100.00	100.00	100.00	100.00	100.00	96.97	
F1-score	89.74	92.79	93.42	97.20	97.94	98.73	
Peak RSS (GB)	0.63	0.08	0.17	6.24	1.90	8.07	<b>4.2X</b> smaller RSS
Index size (GB) <sup>a</sup>	0.68	0.09	0.10	6.20	1.90	31.00	<b>16.3X</b> smaller index
Throughput (bp/s)	252,974	901,609	851,869	64,384	185,618	15,570	<b>11.9X</b> faster

<sup>a</sup>The index sizes for SPUMONI-ms and SPUMONI are only for the positive index because the null index can be used offline and removed.

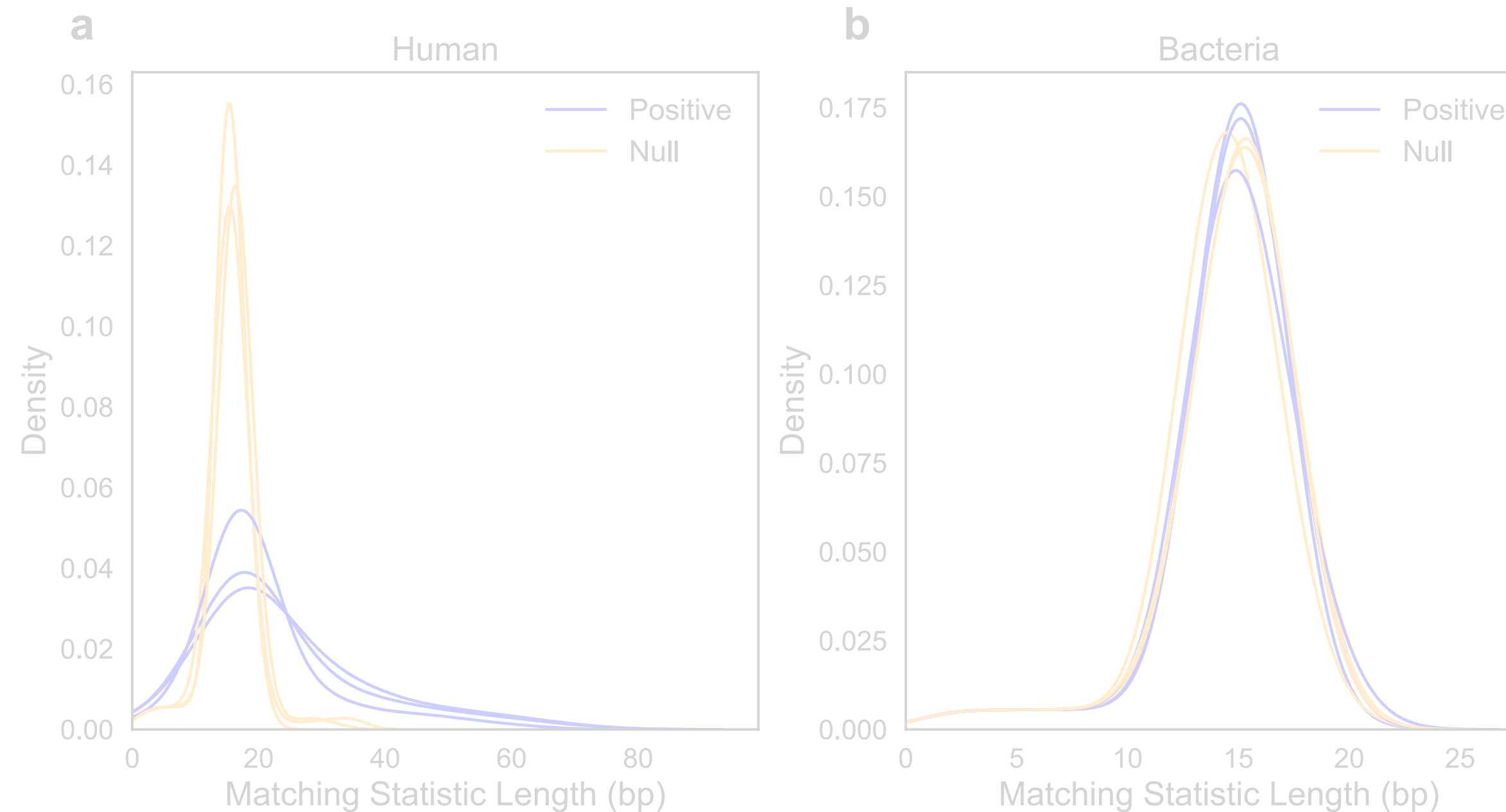
- ▶ **Answer:** ① Yes, using a pan-genome reference, allowed us to target the ZymoMC strains
- ② Faster and uses less memory than minimap2 with similar classification metric

<sup>1</sup>Kovaka, S., Fan, Y., Ni, B. *et al.* Targeted nanopore sequencing by real-time mapping of raw electrical signal with UNCALLED. *Nat Biotechnol* **39**, 431–441 (2021).

# Results - Human Microbiome Experiment

► **Question:** Does SPUMONI's extend to other scenarios<sup>1</sup>, in particular using human genomes?

*Distributions of MS for Human Microbiome Reads Against a Human Index*



► **Answer:** SPUMONI's accuracy is 99.44, while minimap2's accuracy is 99.84 while SPUMONI is 2X faster.

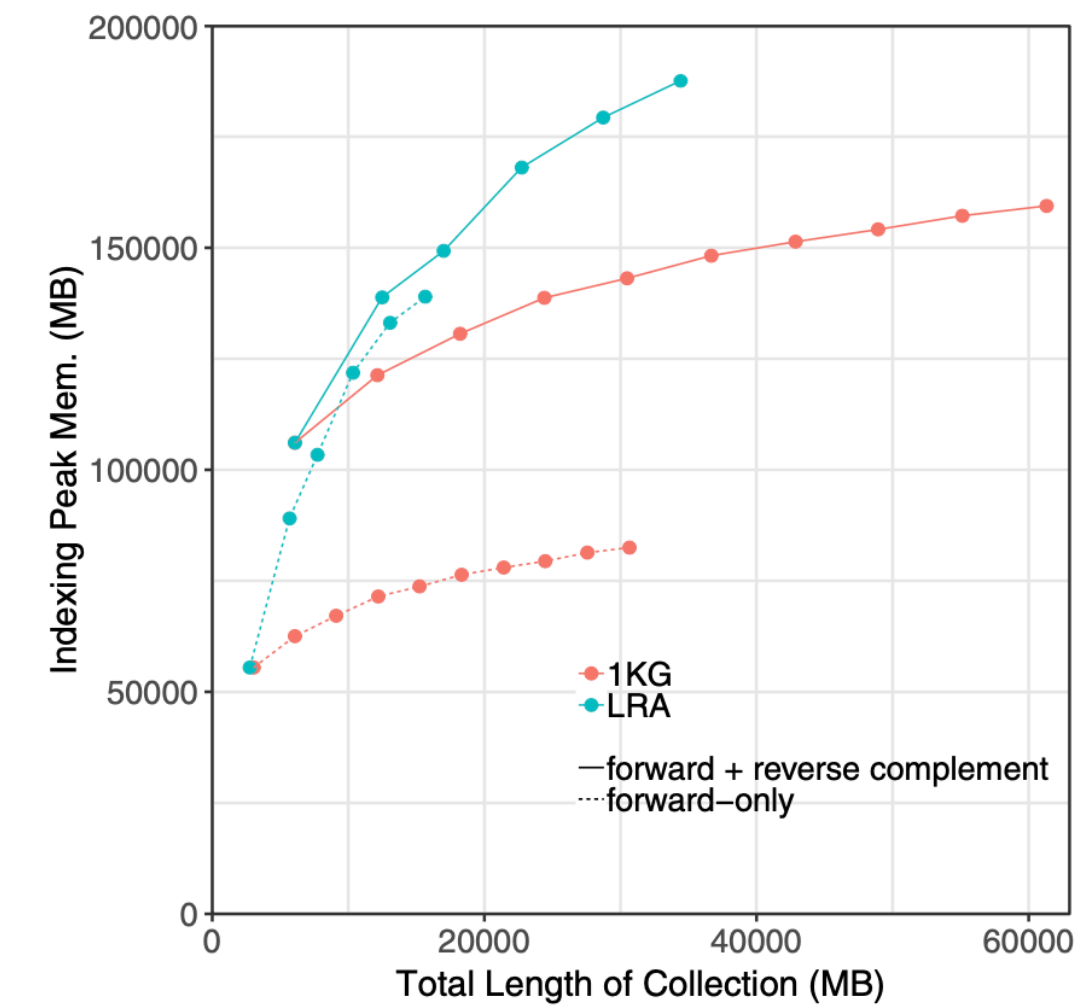
<sup>1</sup>Moss, E.L., Maghini, D.G., and Bhatt, A.S. (2020). Complete, closed bacterial genomes from microbiomes using nanopore sequencing. Nat. Biotech. 38, 701–707.



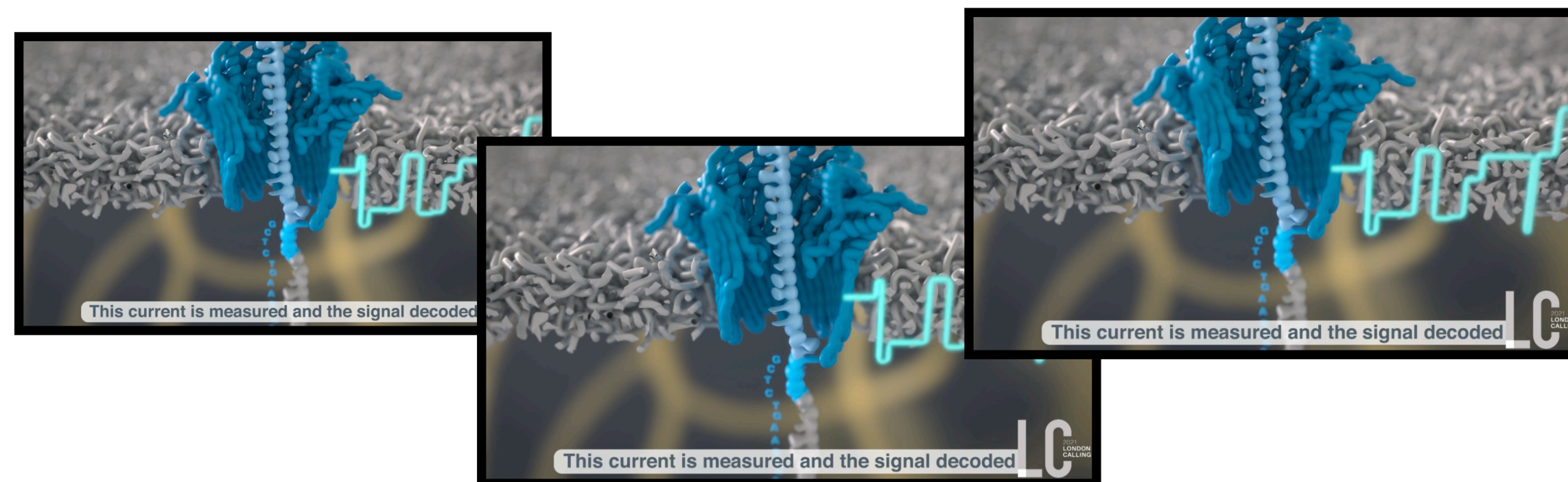
# Conclusion

- ▶ SPUMONI is a streaming algorithm for targeted nanopore sequencing that uses a read's MS or PMLs to classify it in real time

- Use of the r-index and MONI's thresholds allow SPUMONI to index **large, repetitive** references more efficiently than competing approaches

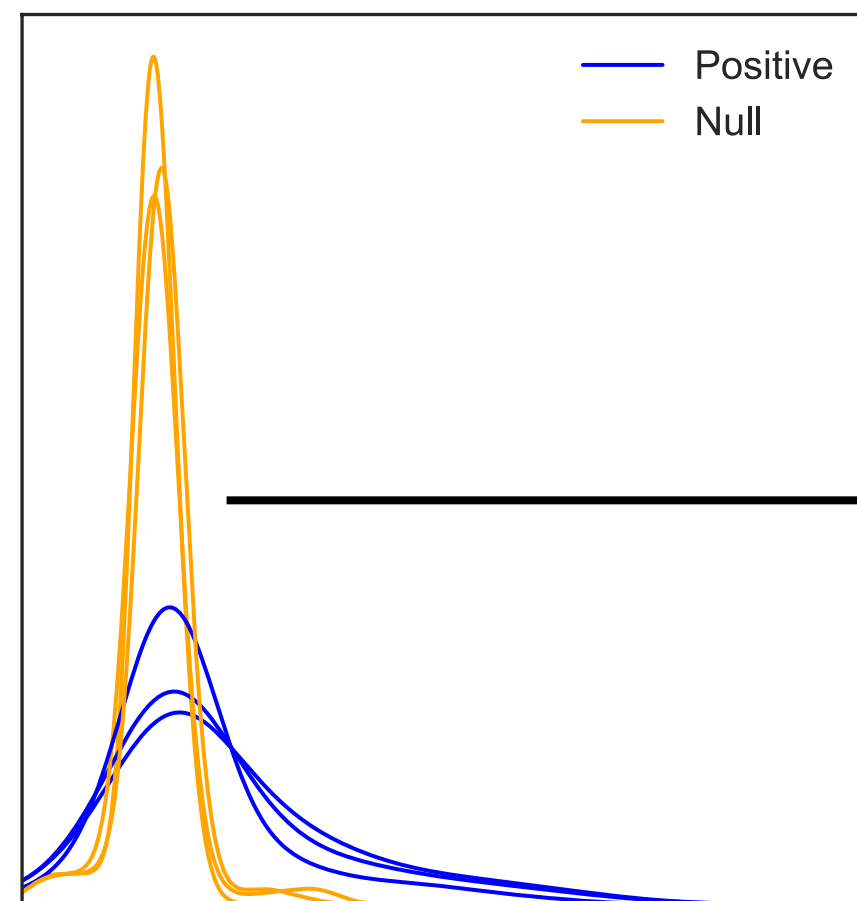


- ▶ SPUMONI's computational efficiency makes it well-suited for nanopore's portable sequencers like MinION or Flongle, and there numerous channels in single flow-cell



# Conclusion

- ▶ SPUMONI's approach opens the door to doing targeted sequencing for strains that have not been assembled previously, or cannot be cultured
- ▶ SPUMONI is not limited by the use of a k-mer length setting or other simple threshold



“Null” distribution allows the notion of significance to be a function of database sequences, and the error rate of the query read.

- ▶ Future work will consist of implementing the needed software for SPUMONI to interact with the Read Until API

# Thank you!

Contact: [oahmed6@jhu.edu](mailto:oahmed6@jhu.edu)

GitHub: <https://github.com/oma219/spumoni>

## Acknowledgements:

- ▶ Massimiliano Rossi, Sam Kovaka, Michael C. Schatz, Travis Gagie, Christina Boucher & Ben Langmead for help & assistance on the project
- ▶ Nae-Chyun Chen, Daniel Baker, Taher Mun, Kathleen Newcomer, Anna Liebhoff & Dominik Kempa from Langmead Lab
- ▶ Marco Oliva from Boucher Lab, and Jarno Niklas Alanko from Gagie Lab

## Funding:



***NSERC***  
***CRSNG***