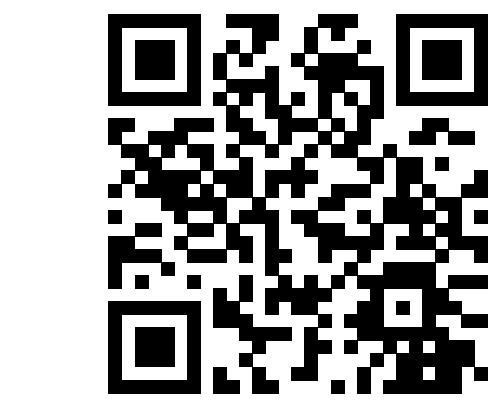


RawHash: Enabling Fast and Accurate Real-Time Analysis of Raw Nanopore Signals for Large Genomes

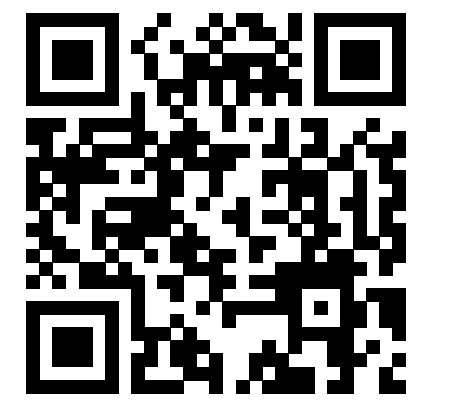
Can Firtina¹, Nika Mansouri Ghiasi¹, Joel Lindegger¹, Gagandeep Singh¹,
Meryem Banu Cavlak¹, Haiyu Mao¹ and Onur Mutlu¹

SAFARI

ETH zürich

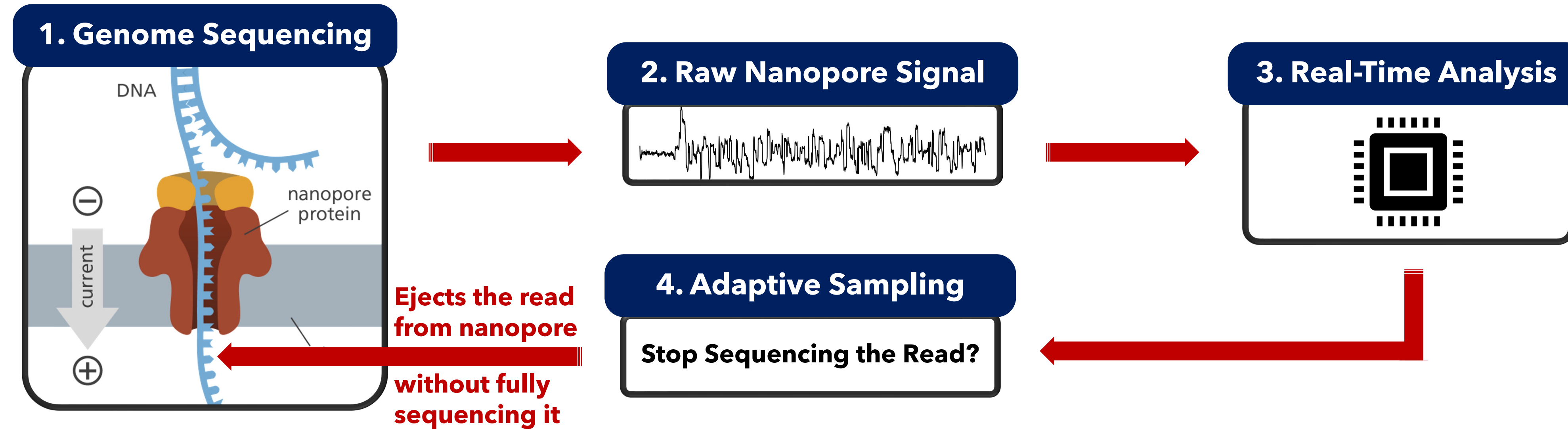


bioRxiv Preprint

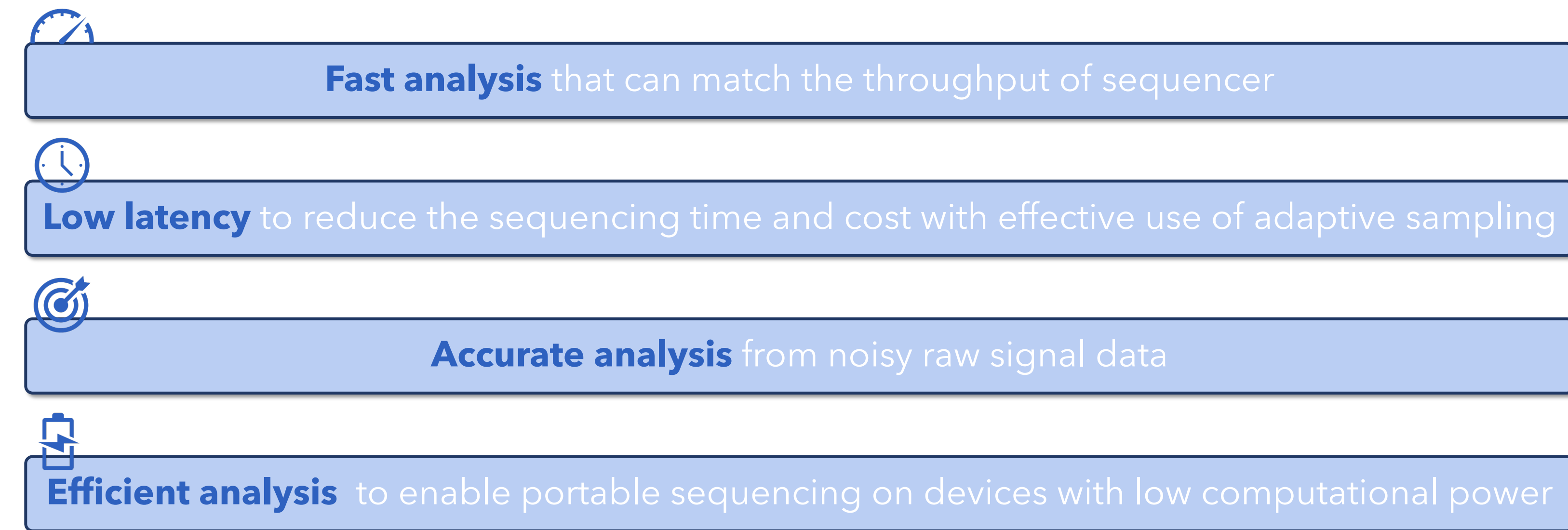


Source Code

1: Real-Time Genome Analysis with Adaptive Sampling



2: Challenges in Real-Time Genome Analysis



3: Problem

- Efficient tools (UNCALLED and Sigmap) **cannot provide** either
 - 1. Fast analysis or
 - 2. Accurate analysis for **large genomes**
- Accurate tools (e.g., ReadFish) **cannot provide**
 - 1. Efficient analysis

4: Goal

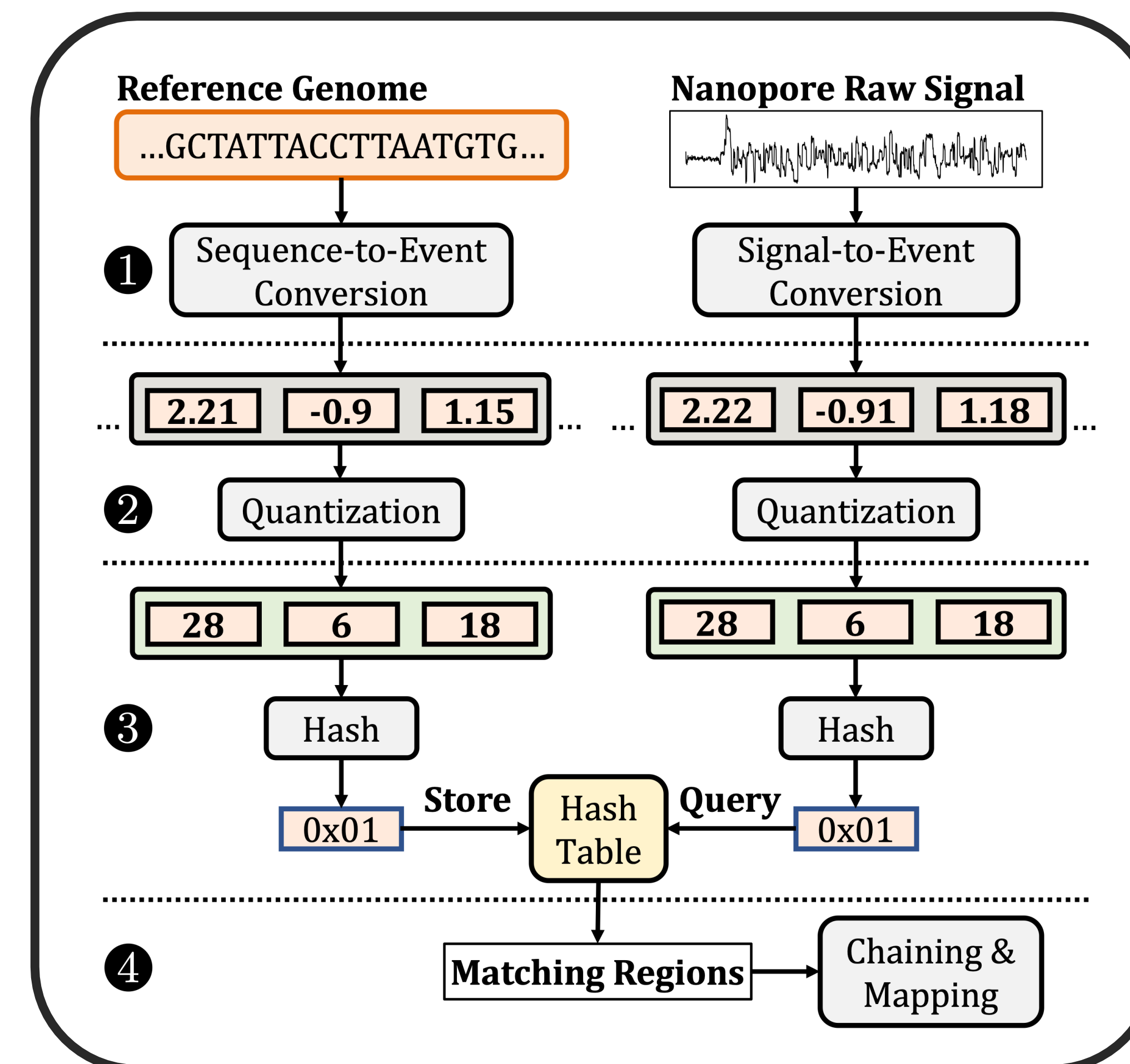
- Fast analysis that can scale to large genomes
- Low latency to make quick decisions on adaptive sampling
- Accurate analysis for large genomes
- Efficient analysis that can be used with portable devices

5: Key Contributions

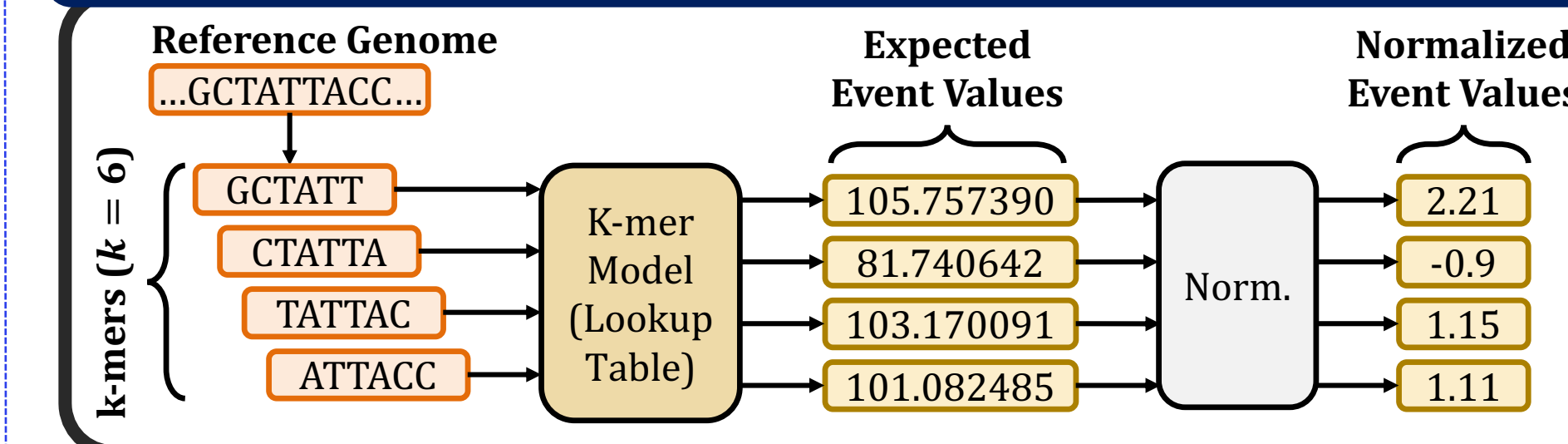
First mechanism that can **efficiently and accurately map raw signals to large reference genomes**

Proposes a novel mechanism, **Sequence Until**, that can dynamically **decide if further sequencing of reads is unnecessary to stop the entire sequencing run**

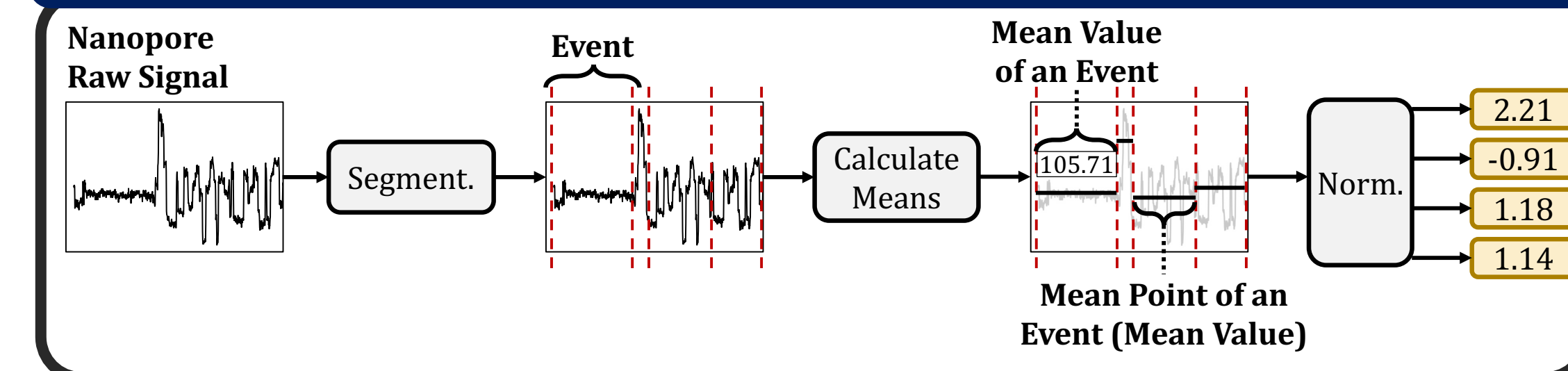
6: RawHash



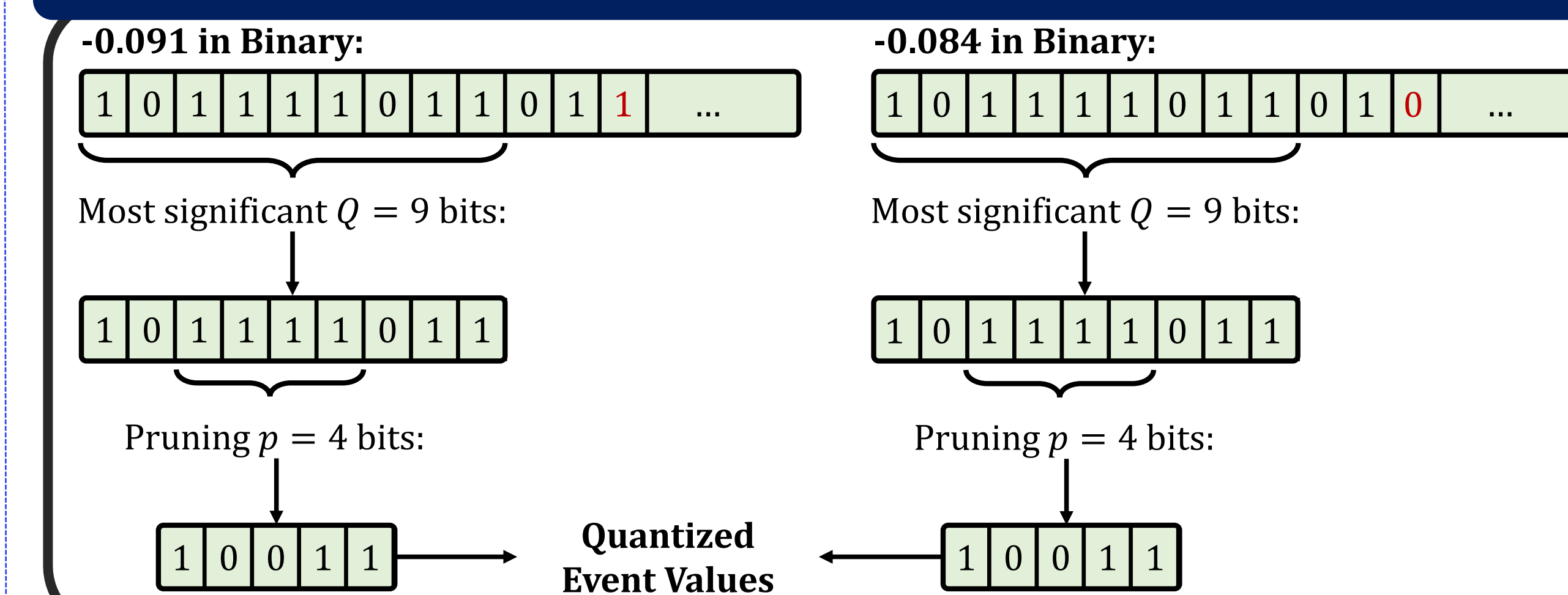
7.1: Sequence-to-Event Conversion



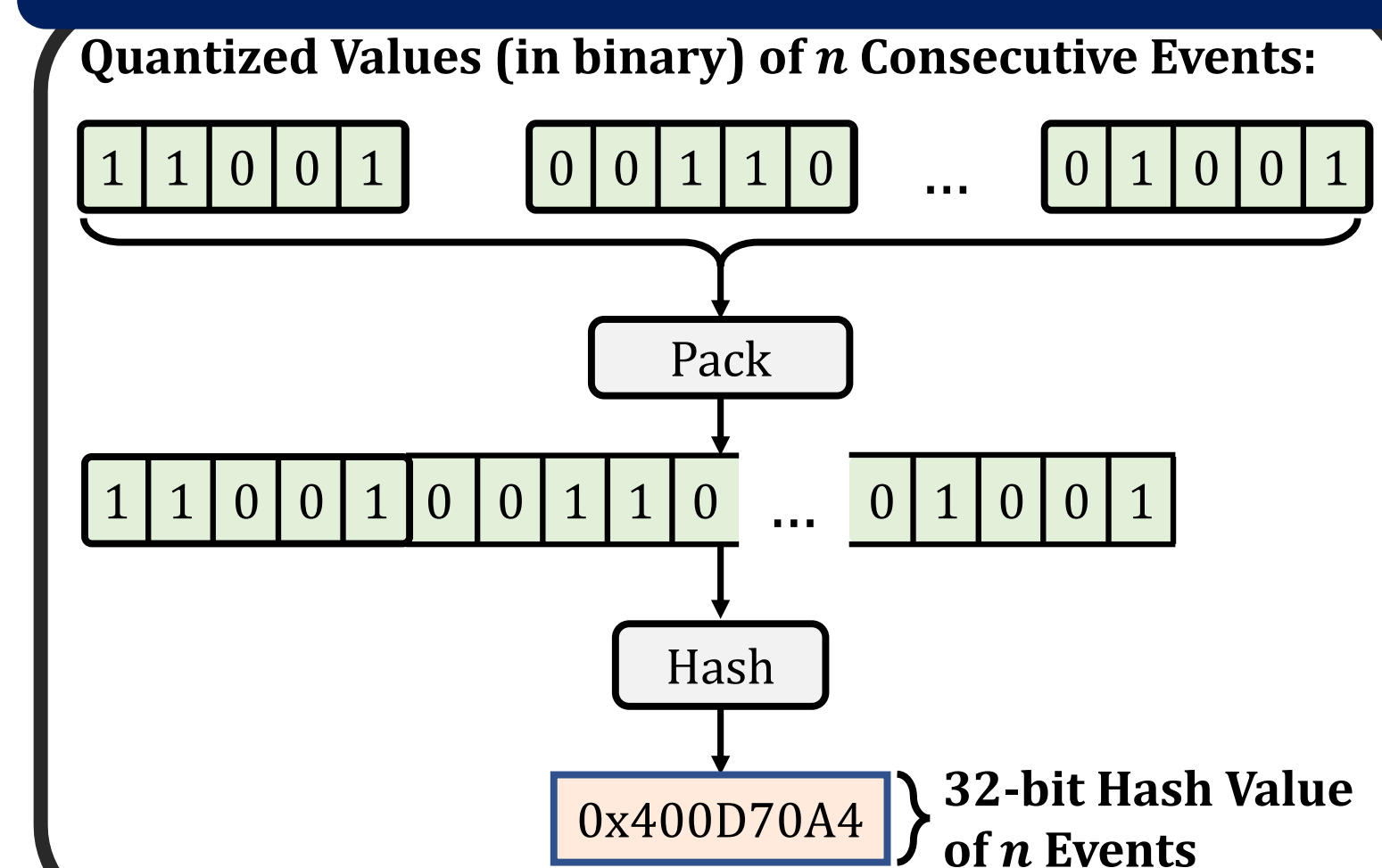
7.2: Signal-to-Event Conversion



8: Quantizing the Event Values

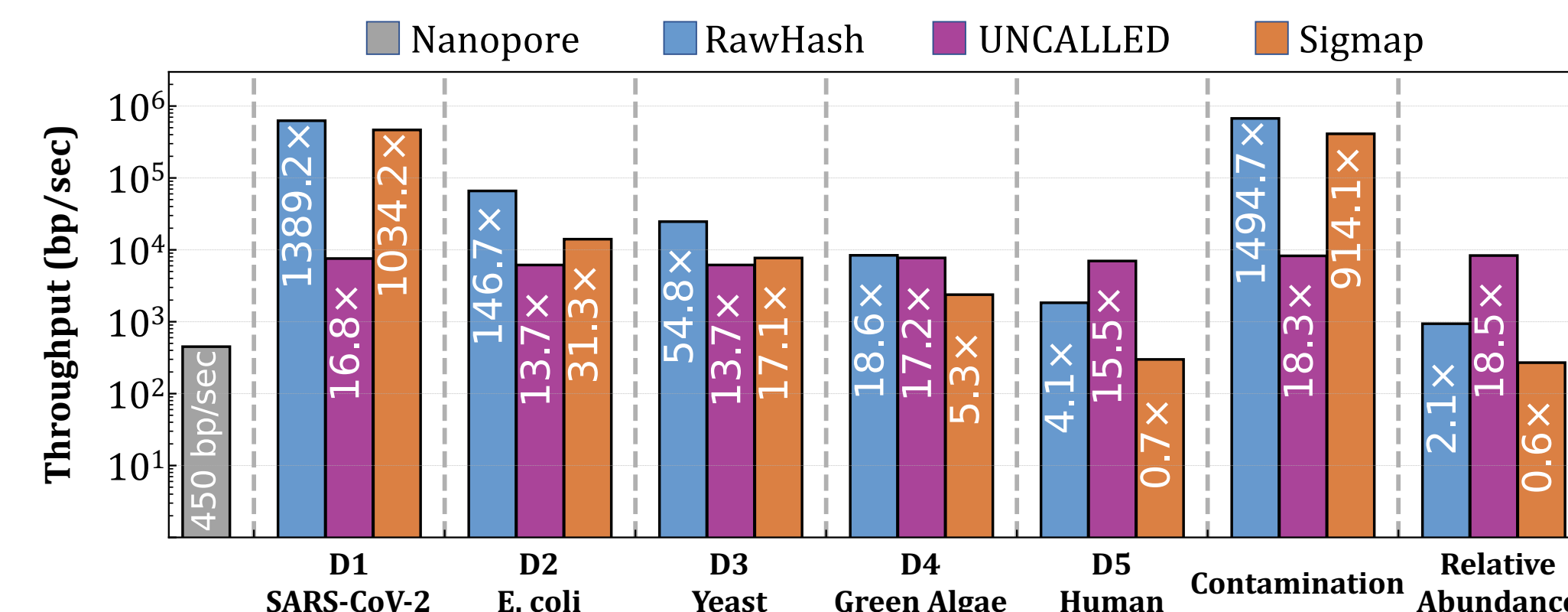


9: Hashing for Efficient Search



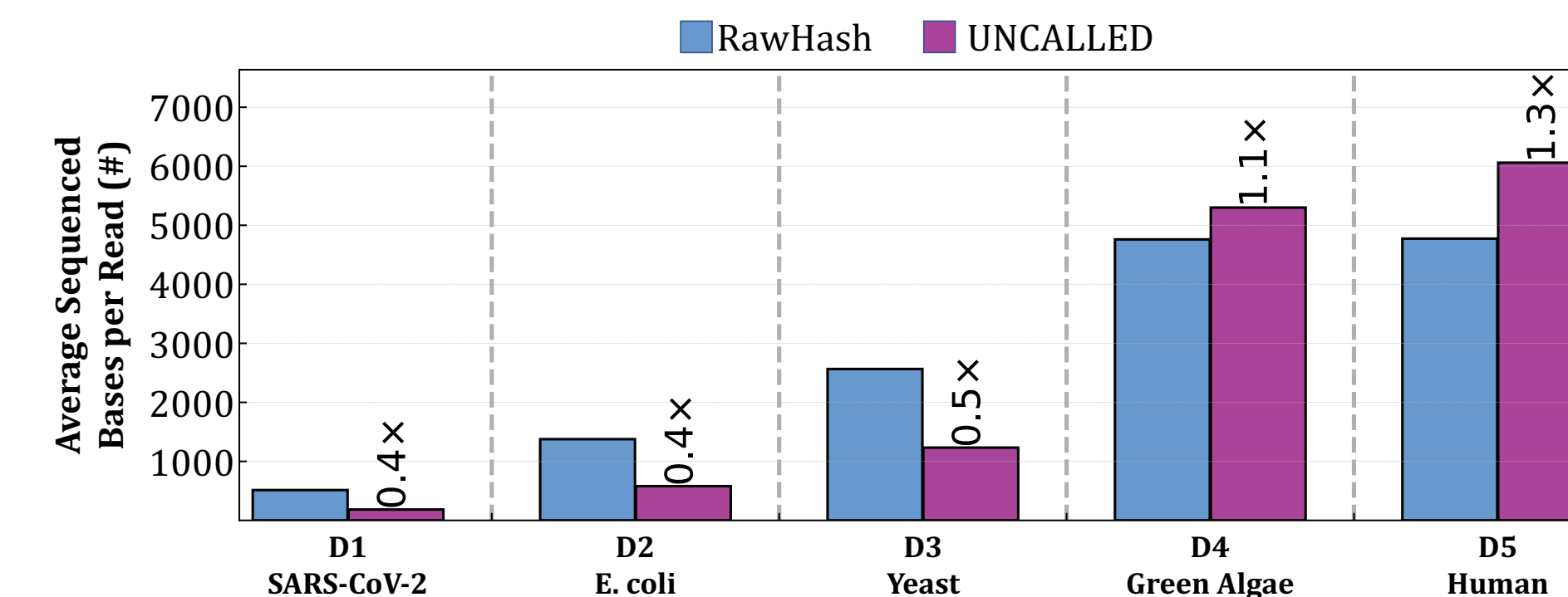
10: Evaluation Methodology

- Datasets from very small (viral) to large genomes (human and metagenomics)
- Compared with UNCALLED and Sigmap
- Use cases
 - 1. Read mapping
 - 2. Relative abundance estimation
 - 3. Contamination analysis
- Evaluating
 - 1. Throughput (bp/sec)
 - 2. Overall Runtime (sec)
 - 3. Memory usage (GB)
 - 4. Number of sequenced bases before ejecting reads (bases)
 - 5. Accuracy (baseline: minimap2 mappings)
 - 6. Sequence Until benefits



- Efficient Analysis: RawHash, UNCALLED, and Sigmap do not require powerful computational resources (e.g., GPUs)
- Fast Analysis: Both RawHash and UNCALLED can match the throughput of nanopore
- Sigmap falls behind the throughput of nanopores for larger genomes

11: Results



- Low Latency: RawHash utilizes adaptive sampling faster for large genomes than UNCALLED
- UNCALLED more bases due as it cannot make accurate decisions for large genomes using fewer bases

Dataset		UNCALLED	Sigmap	RawHash
Read Mapping				
D1 SARS-CoV-2	Precision	0.9547	0.9929	0.9868
	Recall	0.9910	0.5540	0.8735
	F ₁	0.9725	0.7112	0.9267
D2 E. coli	Precision	0.9816	0.9842	0.9573
	Recall	0.9647	0.9504	0.9009
	F ₁	0.9731	0.9670	0.9282
D3 Yeast	Precision	0.9459	0.9856	0.9862
	Recall	0.9366	0.9123	0.8412
	F ₁	0.9412	0.9475	0.9079
D4 Green Algae	Precision	0.8836	0.9741	0.9691
	Recall	0.7778	0.8987	0.7015
	F ₁	0.8273	0.9349	0.8139
D5 Human HG001	Precision	0.4867	0.4287	0.8959
	Recall	0.2379	0.2641	0.4054
	F ₁	0.3196	0.3268	0.5582
Relative Abundance Estimation				
D1-D5	Precision	0.7683	0.7928	0.9484
	Recall	0.1273	0.2739	0.3076
	F ₁	0.2184	0.4072	0.4645
Contamination Analysis				
D1, D5	Precision	0.9378	0.7856	0.8733
	Recall	0.9910	0.5540	0.8735
	F ₁	0.9637	0.6498	0.8734

- Accurate Analysis: RawHash provides the best accuracy for large genomes
- Sequence Until dynamically stops the entire sequencing after sequencing **only 7% of the entire sample**
- Sequence Until least to almost as accurate relative abundance estimation as using the entire (100%) sample.