

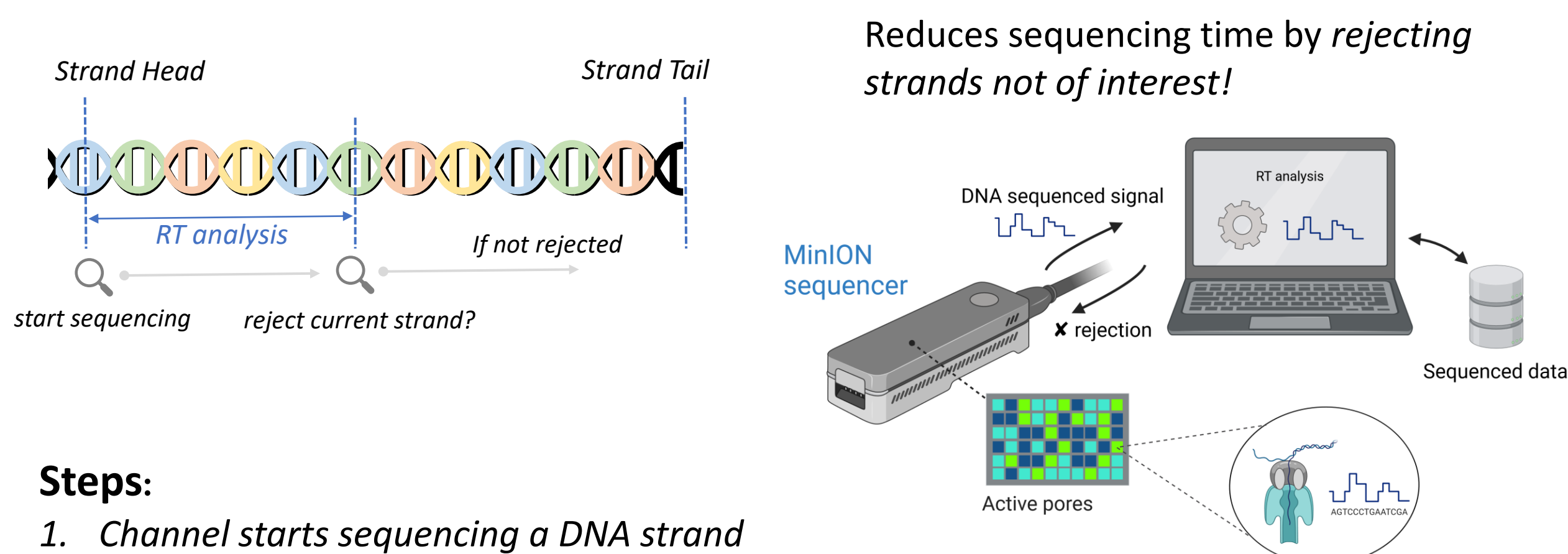
Efficient real-time selective genome sequencing on resource-constrained devices

Po Jui Shih¹, Hassaan Saadat², Sri Parameswaran³, and Hasindu Gamaarachchi^{1,4}

¹ School of Computer Science and Engineering, UNSW, Sydney ² School of Electrical Engineering and Telecommunications, UNSW, Sydney ³ School of Electrical and Information Engineering, University of Sydney, Sydney ⁴ Kinghorn Centre for Clinical Genomics, Garvan Institute of Medical Research, Sydney

Introduction

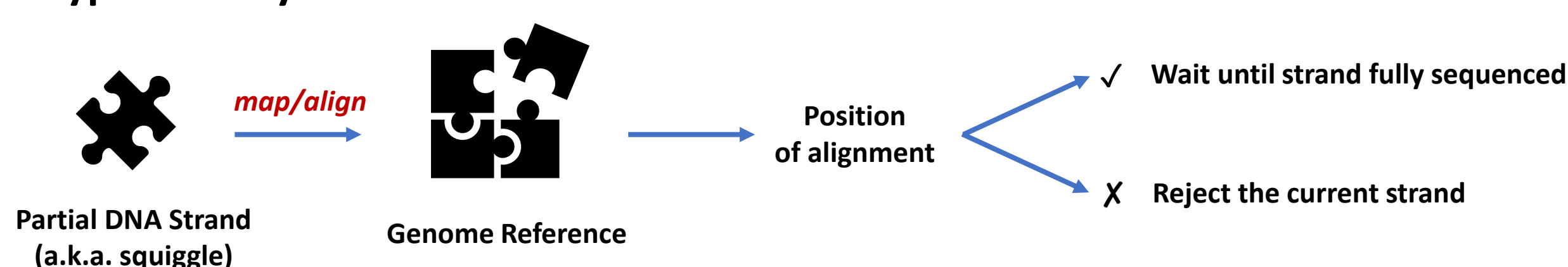
Nanopore sequencers provide **portable** long-read sequencing and the ability to **access, analyse, and filter reads in real-time** → **Read Until**



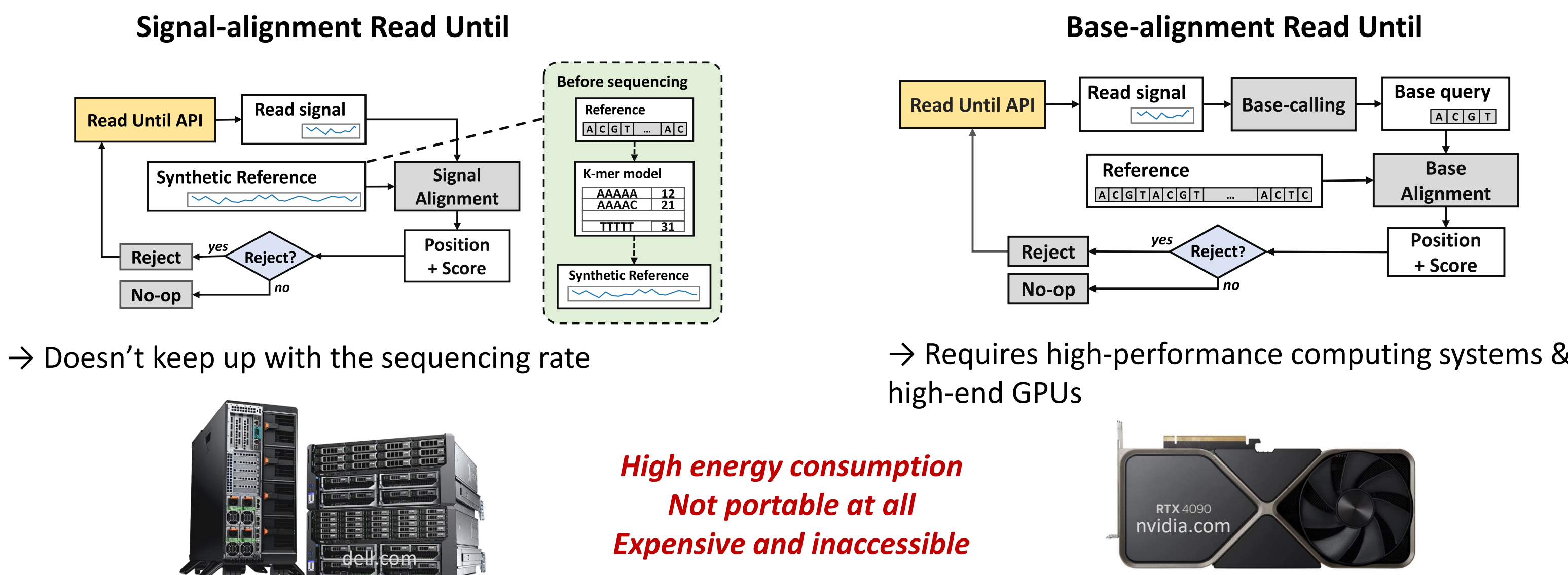
Steps:

1. Channel starts sequencing a DNA strand
2. Host machine performs **real-time analysis** as we sequence
3. If analysis results is to "skip", reject the strand
4. If analysis results is to "continue", let the pore finish sequencing

Typical analysis flow



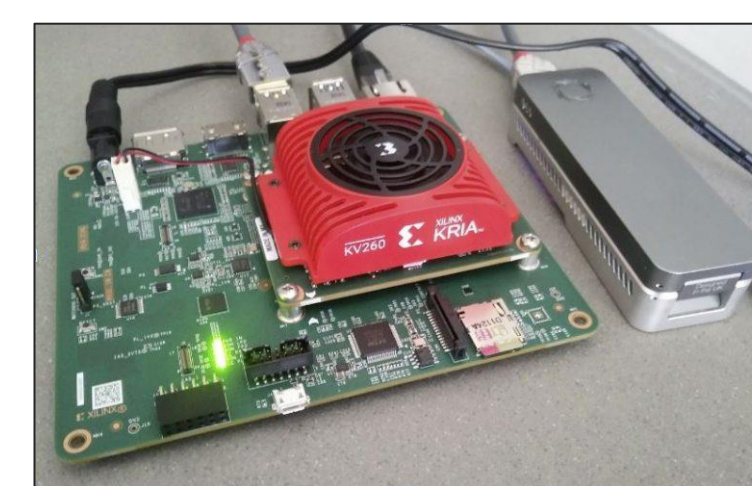
Problem Existing analysis pipelines are compute-intensive → *costly compute hardware requirements!*



Our goal Real-time analysis on an SoC for Read Until on Nanopore Sequencers!

- Low energy
- Ultra portable and accessible
- Scalable & adaptive to future sequencing advancements
- Fully working from end-to-end
- HIGH THROUGHPUT!

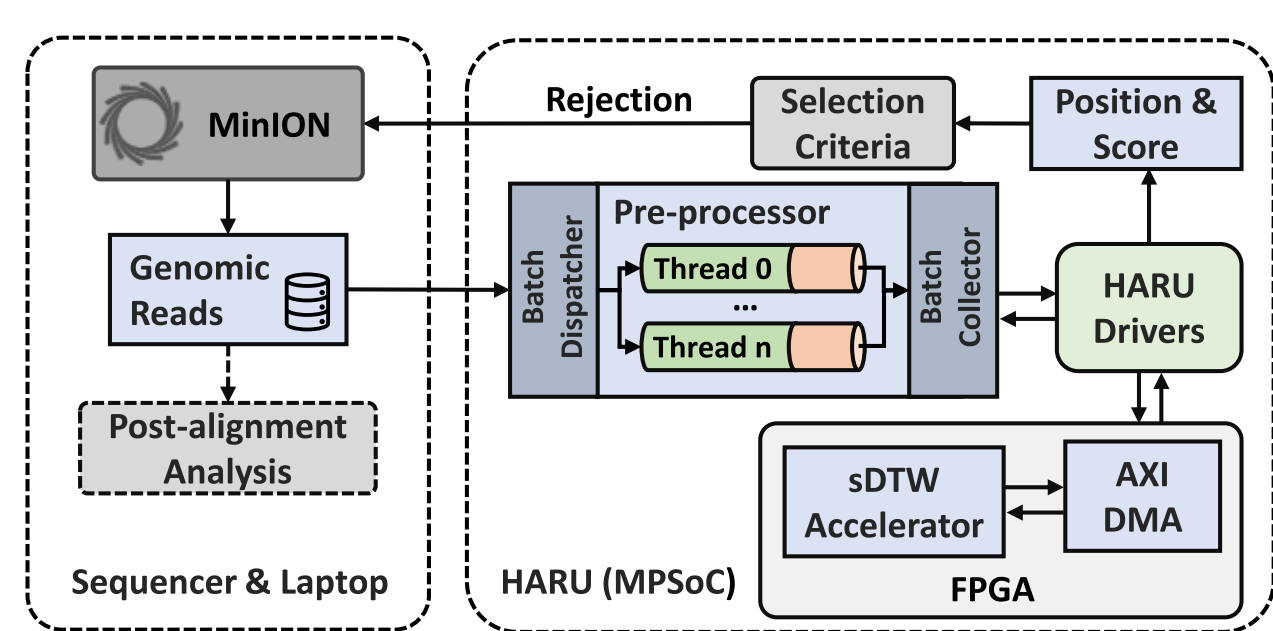
Hardware Accelerated Read Until



System Design and Methodologies

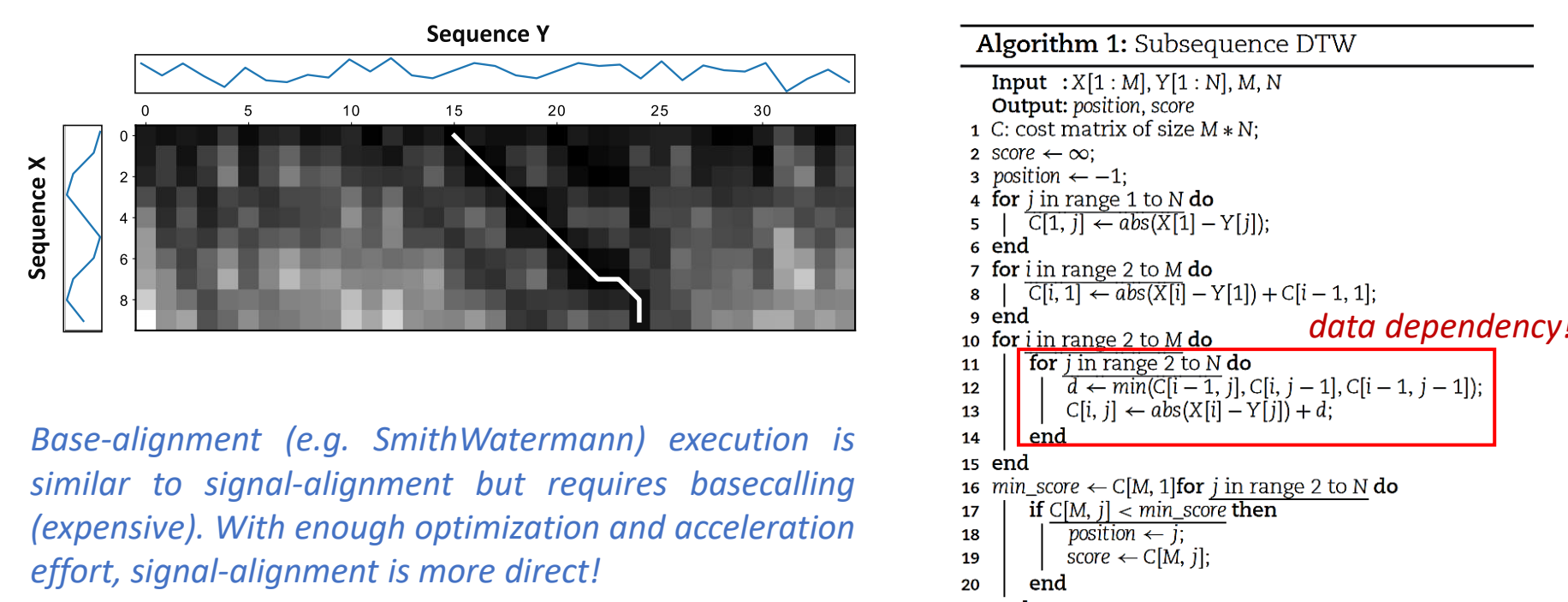
HARU Overview

HARU is a **signal-alignment** hardware-software co-design pipeline for Read Until. It features an efficient **subsequence dynamic time warping (sDTW) hardware accelerator** running in the FPGA of an AMD Zynq MPSoC (an SoC with an ARM processor and FPGA).

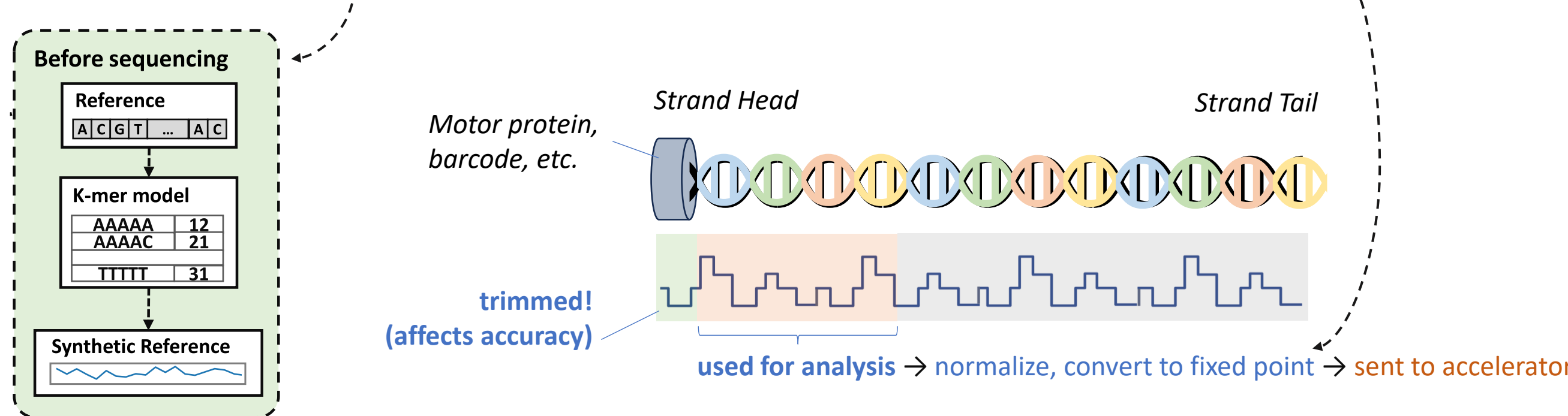
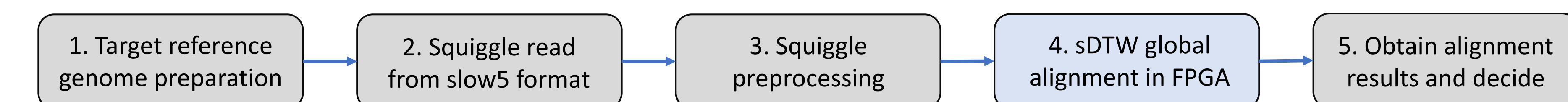


Subsequence DTW

Finds the optimal (warped) alignment of a subsequence to reference → sequences are represented in raw signal values (as opposed to bases) → **computationally intensive!** Naively quadratic for both space and time complexity → takes up 98% of processing time for signal-alignment Read Until!

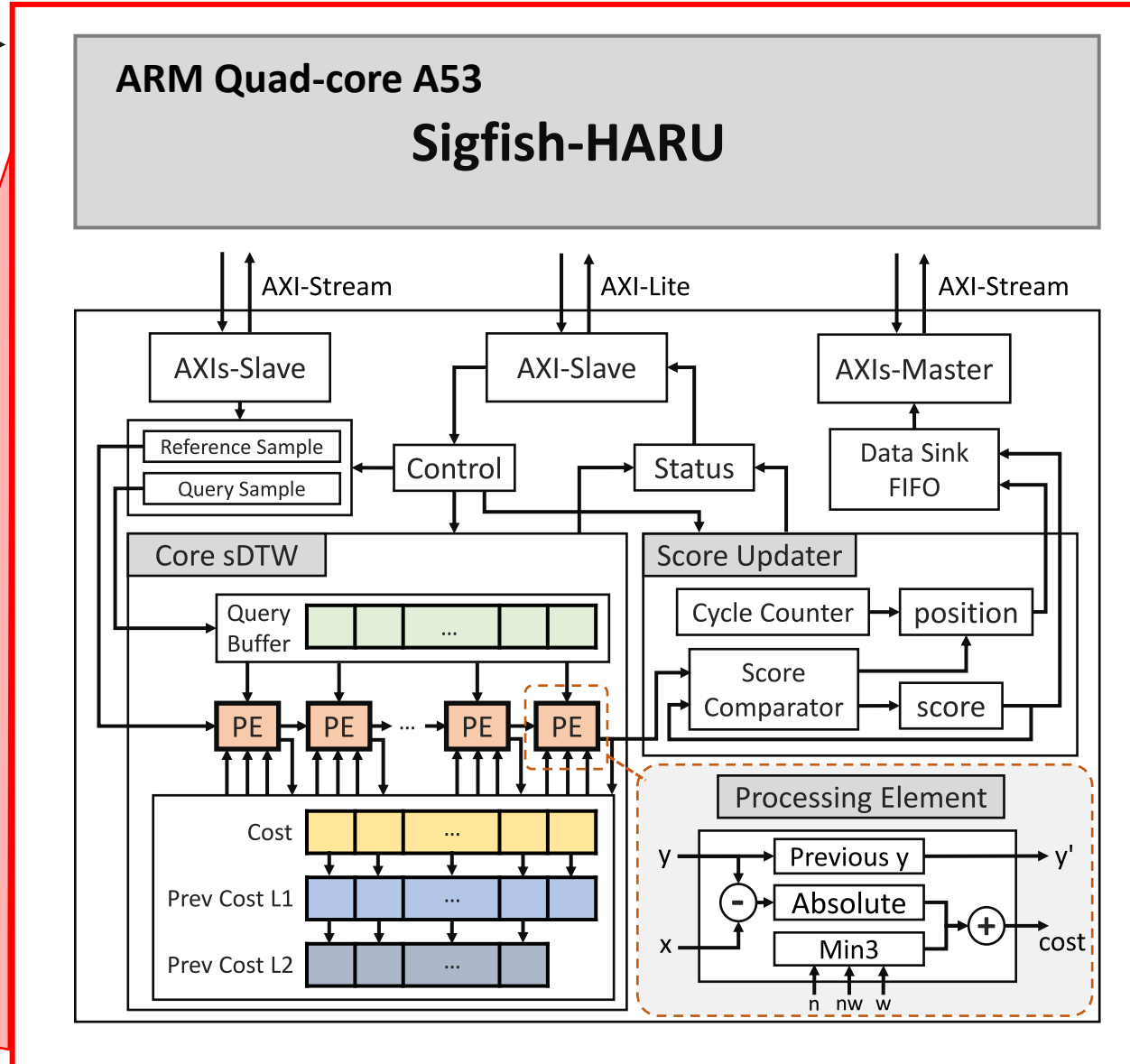
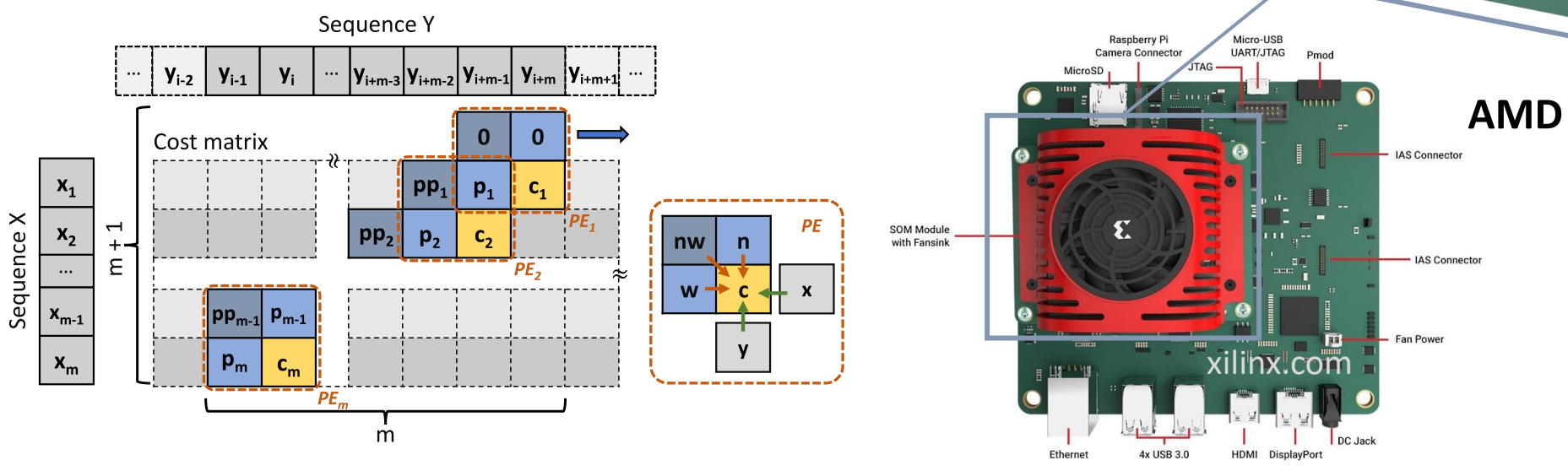
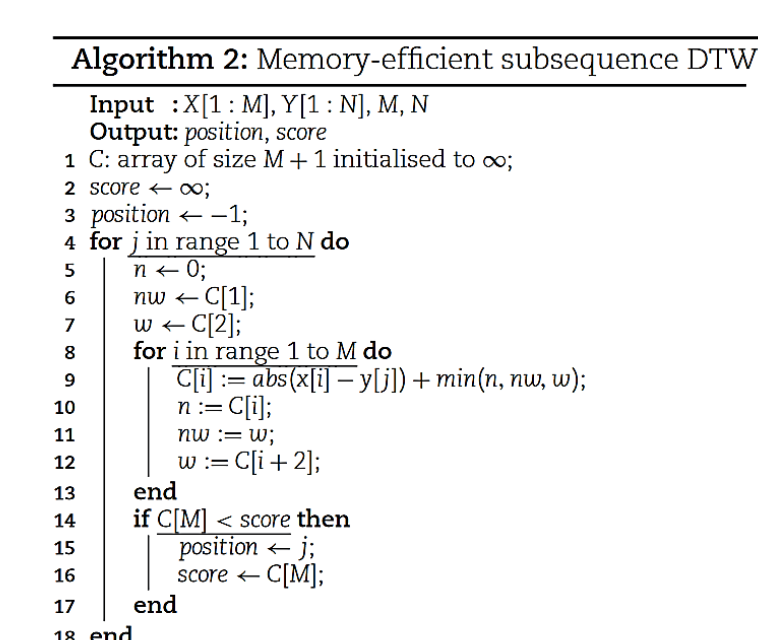


HARU Co-Design Execution Breakdown



HARU's sDTW accelerator

- Optimizations**
- Column cells computed in parallel (via pipelining)
- fixed-point data & tuned scaling factor (reduce computation cost)
- data reusing, no backtracing (keeping only necessary data of the DP memoization)
- *Becomes a chain of Processing Elements (PEs) sliding through the reference (linear)*



AMD Kria AI Starter Kit (249 USD)

Scan to access paper and codebase

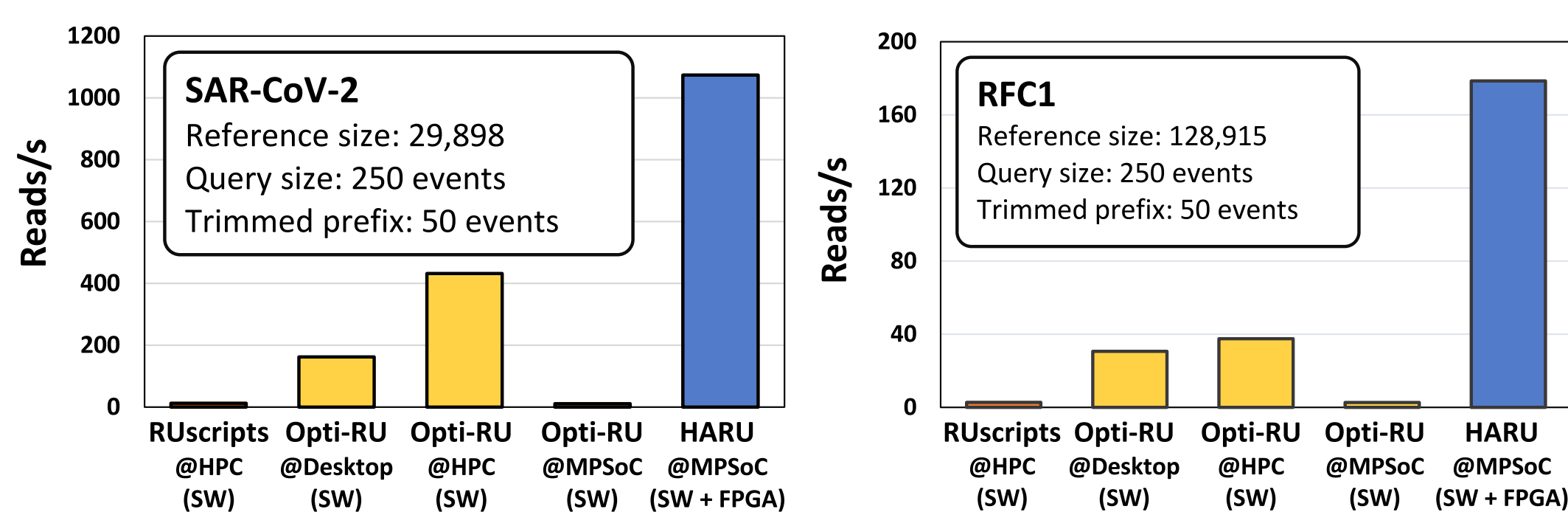


<https://doi.org/10.1093/gigascience/giad046>

HARU Sigfish-HARU

Results and Outcomes

Processing capability (HARU vs sDTW implementations)



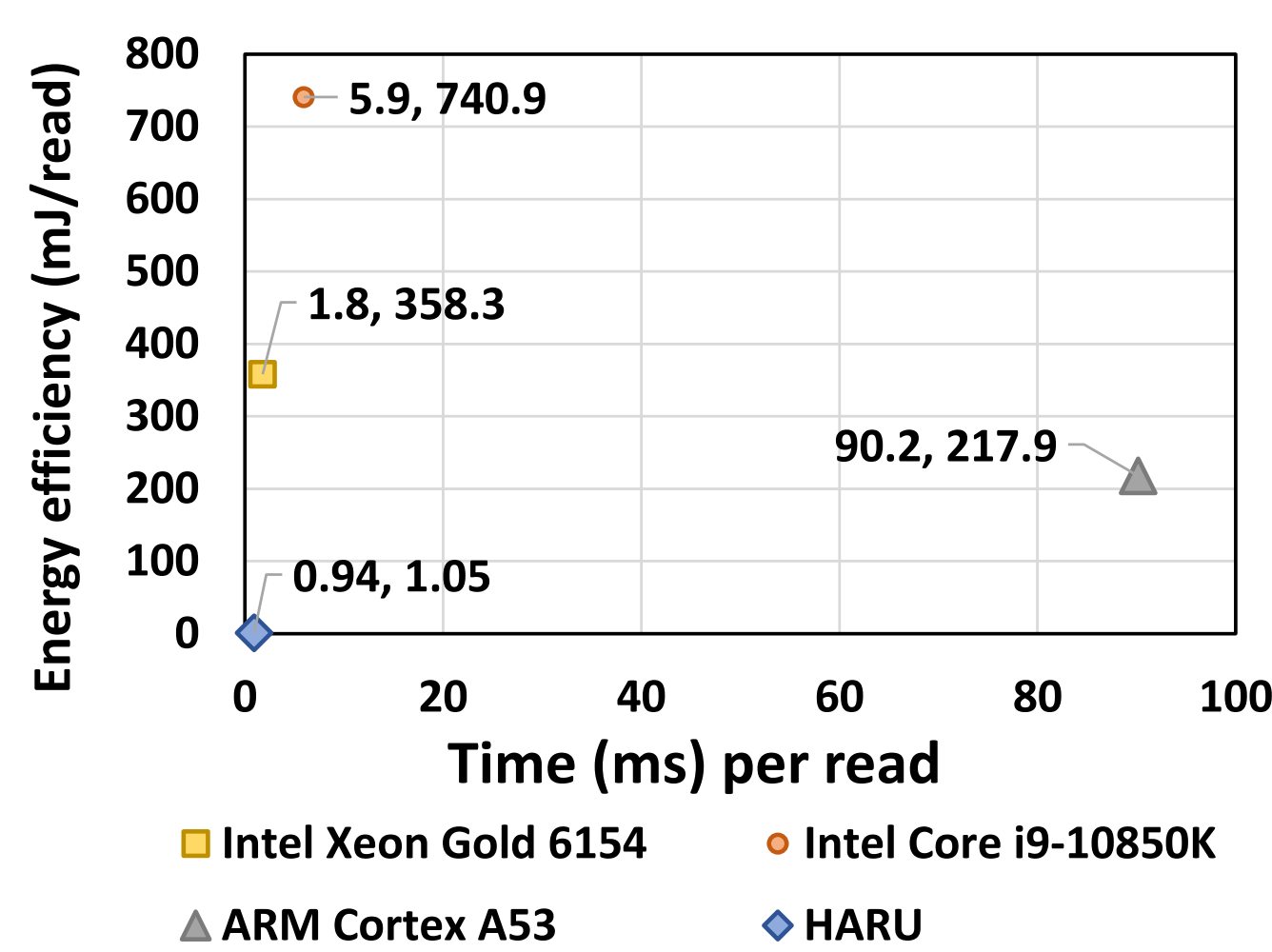
Desktop - Intel Core i9-10850K (10 cores, 32GB RAM)
HPC - Intel Xeon Gold 6154 (36 cores, 377GB RAM)
MPSoC - ARM Cortex A53 (4 cores, 4GB RAM)

Speedups (HARU vs):

- RUScripts [1] @HPC → 85.8 x
- Optimized RUScripts @Desktop → 6.6 x
- Optimized RUScripts @HPC → 2.5 x

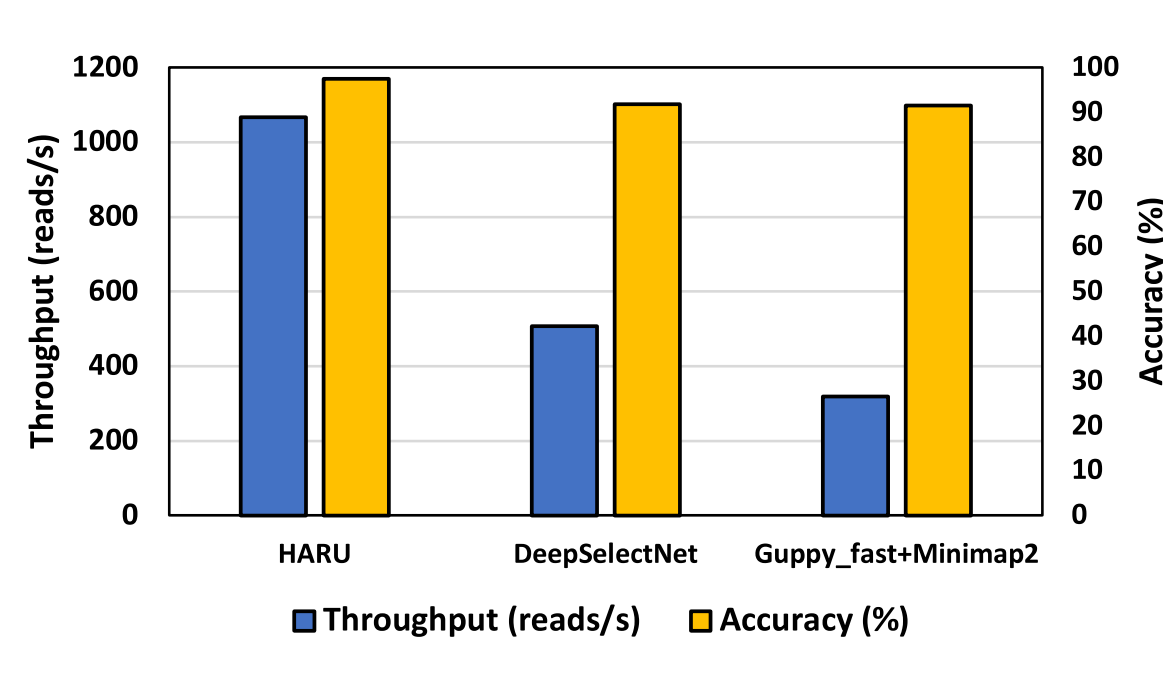
Smaller, cheaper, but higher throughput!

Energy-delay product (HARU vs sDTW implementations)



Significant improvement in energy efficiency!

Processing capability & Accuracy (HARU vs other SOTA implementations)



DeepSelectNet [2] running on HPC + Telsa V100 GPU
Guppy_fast + Minimap2 [3] running on Nvidia Jetson Xavier

