Shifted Hamming Distance (SHD): A Fast and Accurate SIMD-Friendly Filter

for Local Alignment in Read Mapping

Hongyi Xin¹, John Greth¹, John Emmons¹, Gennady Pekhimenko¹, Carl Kingsford¹, Can Alkan², Onur Mutlu¹

Shifted Hamming Distance (SHD):

quickly filter out incorrect mappings



¹ Departments of Computer Science and Electrical and Computer Engineering, Carnegie Mellon University, Pittsburgh, PA, USA

Key observations:

² Dept. of Computer Engineering, Bilkent University, Ankara, Turkey

Problem:

Carnegie Mellon

SAFARI

- NGS mappers can be divided into two categories: Suffix-array based and seed-and-extend based
- 1. Suffix-array based mappers (i.e. bwa, bowtie2, SOAP3) find the best mappings fast but lose high-error mappings
- 2. Seed-and-extend based mappers (i.e., mrfast, shrimp, RazerS3) finds all mappings but waste resources on rejecting incorrect mappings
- Our goal: Provide an effective filter to efficiently filter out incorrect mappings

Shifted Hamming Mask-set (SHM):

- Key idea: SHM identifies matching by incrementally shifting the read against the reference
- Mechanism: Use bit-wise XOR to find all matching bps. Then use bit-wise AND to merge them together (Fig 1)
- Mappings that contain more than e '1's in the final bit-vector must contain more than e errors
- Cons: SHD may let incorrect mappings pass through (Fig 2) because all '0's "survive" the AND operations

Speculative Removal of Short-matches (SRS):

 Key idea: SRS refines SHM by removing short stretches of matches (<3 bps) identified in the Hamming masks (Fig 3)

• Key idea: use simple bit-parallel and SIMD operations to

1. If two strings differ by $\leq e$ errors, then every non-

(e+1) identical sections (Pigeonhole Principle)

and speculative removal of short-matches (SRS)

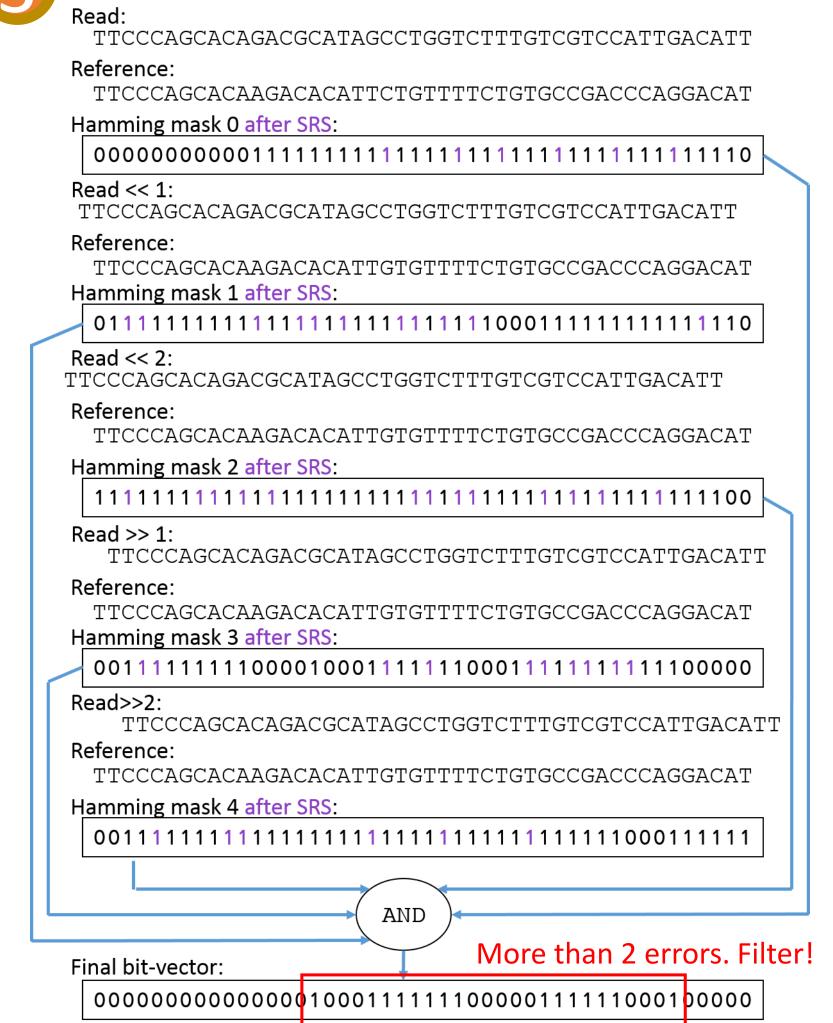
erroneous bp can be aligned in at most e shifts

2. If two strings differ by ≤e errors, then they share at most

- Key observations:
- 1. Identical sections tend to be long (≥ 3 bps)
- 2. Short stretches of matches (streaks of '0' <3 bps) are likely ! to be random matches of bps (generate spurious '0's)
- Mechanism: Amend short streak of '0's into '1's while count errors conservatively in the final bit-vector (Fig 4)

SHM fails to identify an incorrect mapping due to random matches (e=2) TTCCCAGCACAGACGCATAGCCTGGTCTTTGTCGTCCATTGACATT TTCCCAGCACAAGACACATTCTGTTTTCTGTGCCGACCCAGGACAT Hamming mask 0: 0000000000111111111011101101110111011101110 TTCCCAGCACAGACGCATAGCCTGGTCTTTGTCGTCCATTGACATT TTCCCAGCACAAGACACATTGTGTTTTCTGTGCCGACCCAGGACAT Hamming mask 1: Read << 2: TTCCCAGCACAGACGCATAGCCTGGTCTTTGTCGTCCATTGACATT TTCCCAGCACAAGACACATTGTGTTTTCTGTGCCGACCCAGGACAT Refined by **SRS** Hamming mask 2: Read >> 1: TTCCCAGCACAGACGCATAGCCTGGTCTTTGTCGTCCATTGACATT Reference: TTCCCAGCACAAGACACATTGTGTTTTCTGTGCCGACCCAGGACAT Hamming mask 3: Read>>2: TTCCCAGCACAGACGCATAGCCTGGTCTTTGTCGTCCATTGACATT TTCCCAGCACAAGACACATTGTGTTTTCTGTGCCGACCCAGGACAT Hamming mask 4: AND No errors? Pass. Oops!

SRS removes short random matches from the Hamming masks (e=2)



Results and Conclusion:

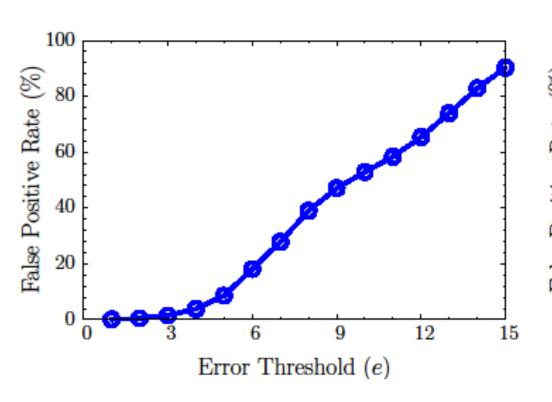
Final bit-vector:

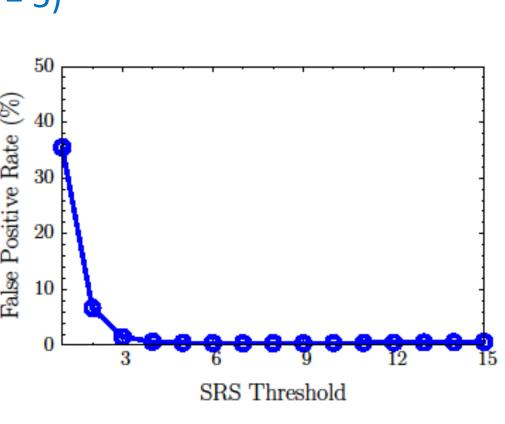
• SHM and SRS are implemented using bit-parallel and SIMD operations (with Intel SSE, details in upcoming paper)

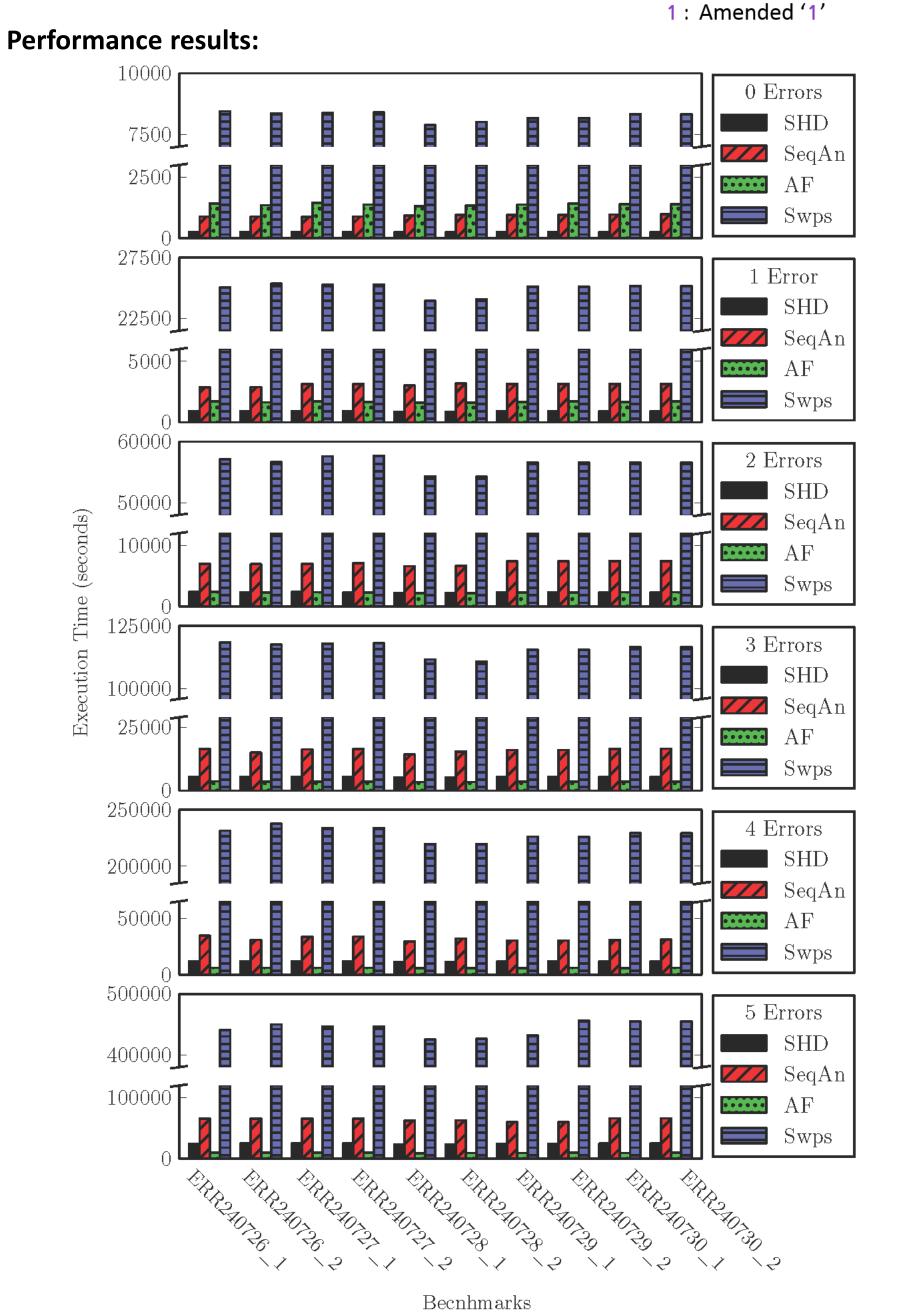
0: Spurious '0'

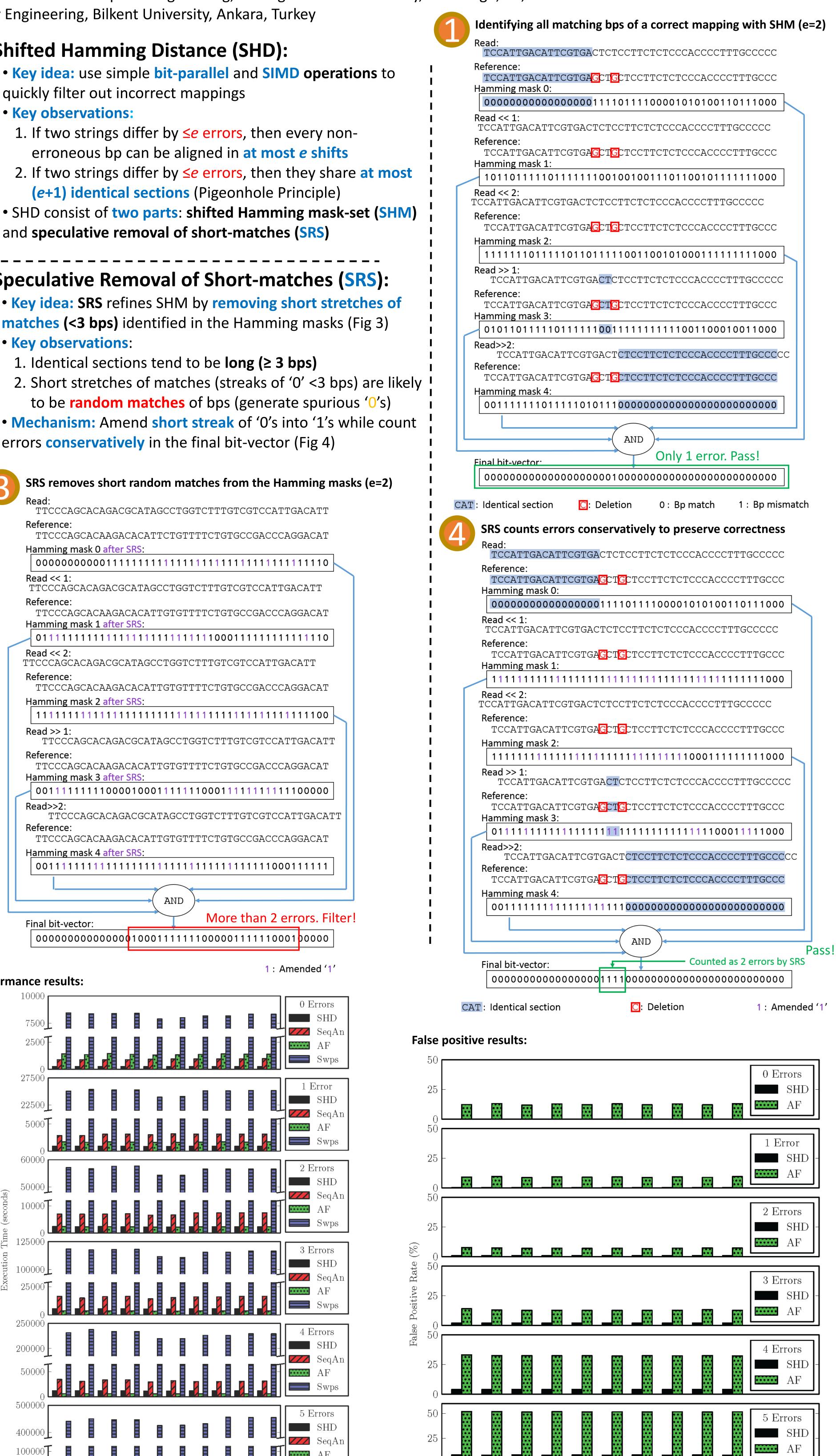
• The threshold of SRS is platform dependent (3 bps at maximum on Intel platforms)

- We compare SHD against:
 - SeqAn: Gene Myers' bit-vector algorithm
 - Swps: A SIMD implementation of Smith-Waterman algorithm
 - AF: A k-mer locality based filter, FastHASH
- We used mrFast to retrieve all potential mappings (readreference pairs) from ten real read sets from 1000 **Genomes Project**
- The false positive rate of SHD increases with larger error thresholds. SHD is effective with up to 5% error rate
- **Key Conclusion:** SHD is **3x faster** than the best previous implementation of edit-distance calculation, while having a false positive rate of only 7% (e = 5)









ERR240730

ERR240730

ERR240739

ERREAD 20

ERR240738

Benchmarks