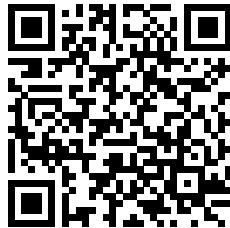




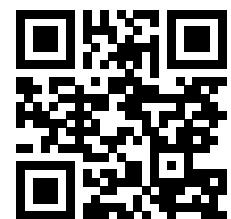
# BLEND

A Fast, Memory-efficient and Accurate Mechanism  
to Find Fuzzy Seed Matches in Genome Analysis

**Can Firtina**, Jisung Park, Mohammed Alser, Jeremie S. Kim, Damla Senol Cali,  
Taha Shahroodi, Nika Mansouri Ghiasi, Gagandeep Singh,  
Konstantinos Kanellopoulos, Can Alkan, Onur Mutlu



[Paper \(NARGAB\)](#)



[Source Code](#)

**SAFARI**

**ETH** zürich

**Carnegie Mellon**

**TU**Delft

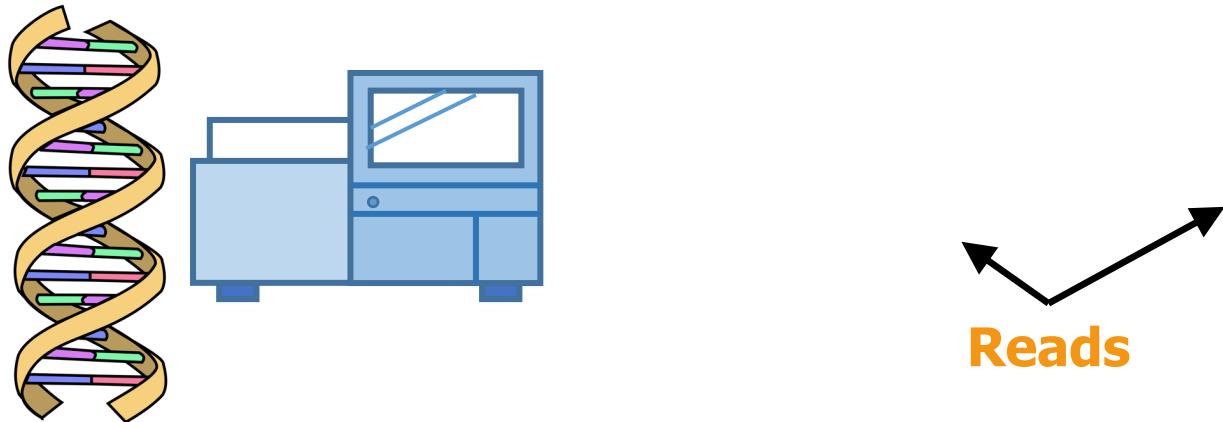
Delft University of Technology



**Bilkent University**

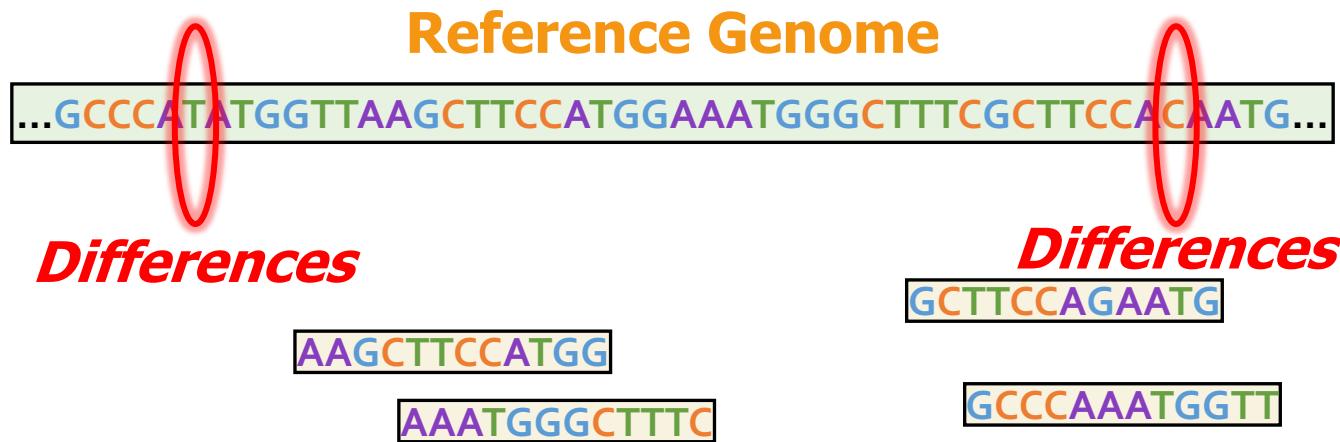
# Genome Analysis

- High-throughput sequencing machines extract smaller fragments of the original DNA sequence, known as **reads**
  - **Challenge:** Perform genome analysis from the small random pieces of genomes



# Identifying Sequence Similarities

1. Mapping reads to a **reference genome**
  - Finding the potential **matching locations** in the genome
  - Identifying **differences** between a read and a reference genome

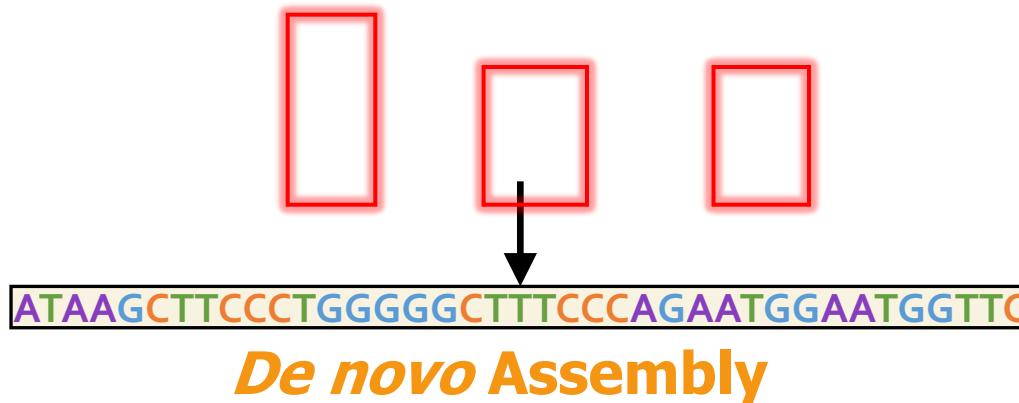


# Identifying Sequence Similarities

## 2. Overlapping reads **to each other**

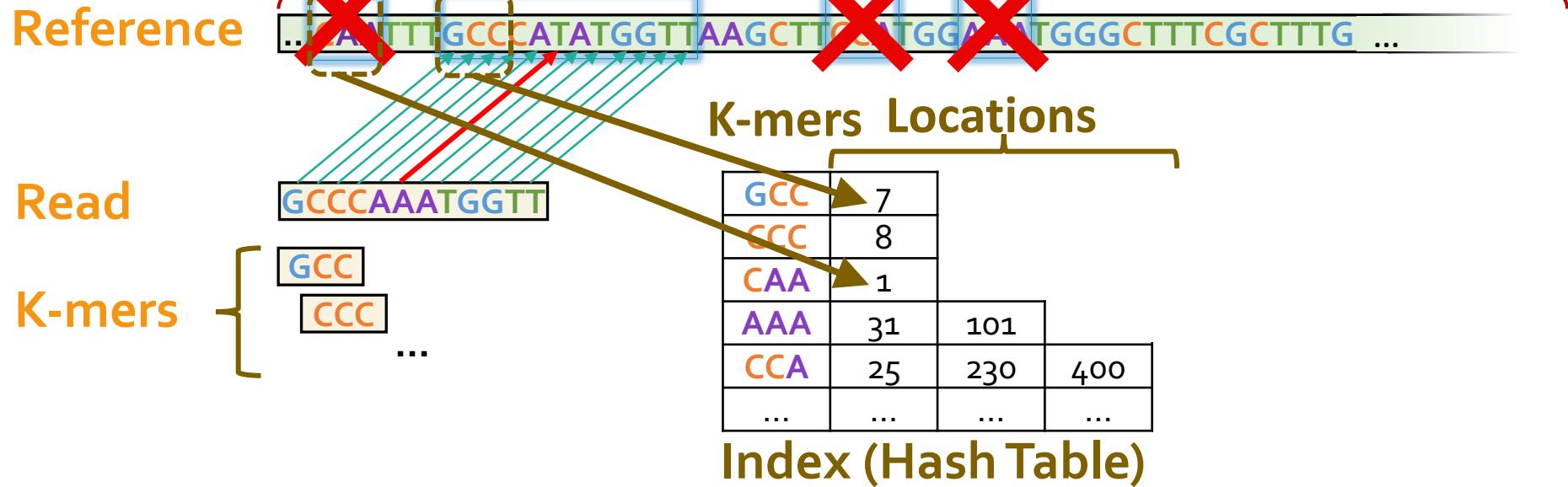
- Constructing assemblies

AAGCTTCCATGG ATAAGCTTCCCT CTGGGGGGCTTTC  
AGAACGTTACTT TTTCCCAGAACATG ATGGAATGGTTC



- Similarities must be identified **efficiently**
  - To facilitate a **practical search** among many sequence pairs

# Practical Similarity Identification

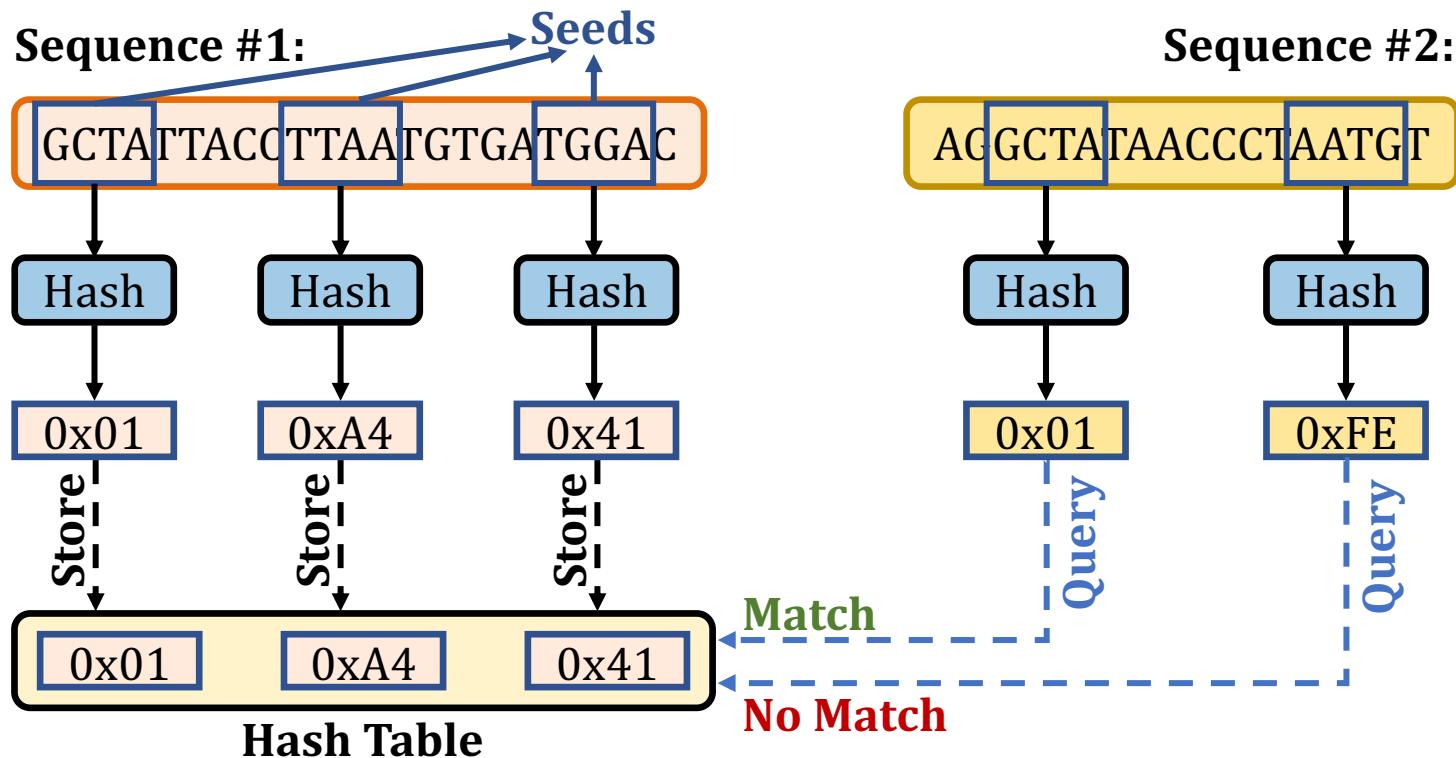


**Seeding** Determine potential matching regions (seeds) in the reference genome

**Seed Filtering (e.g., Chaining)** Prune some seeds in the reference genome

**Alignment** Determine the exact differences between the read and the reference genome

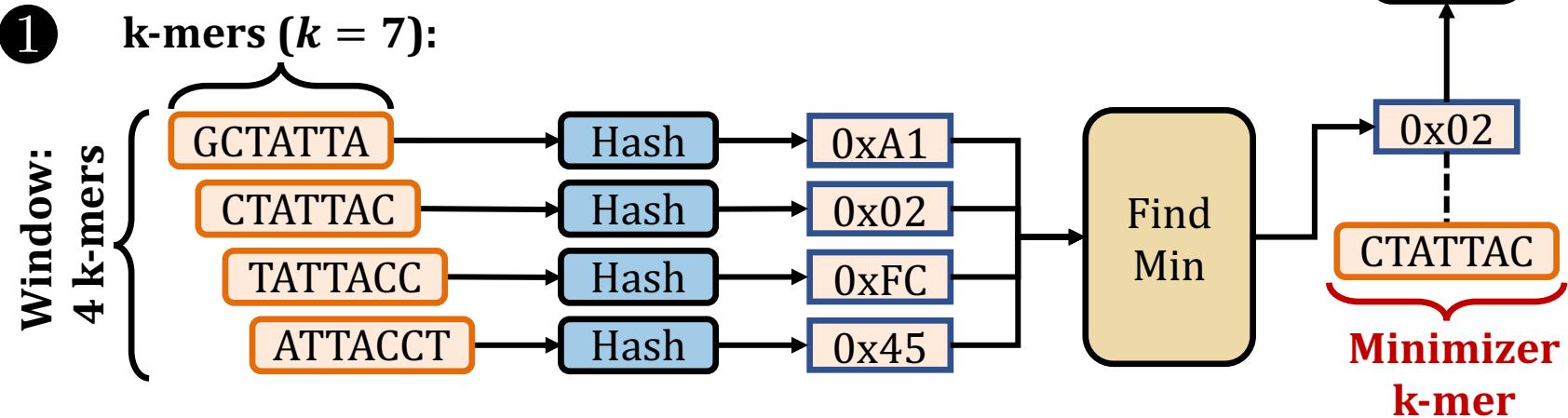
# Finding Seed Matches in Hash Tables



# Seed Matching Techniques

## 1. Sampling the overlapping k-mers

- **Minimizers**
- **Window length (w)**: accuracy & performance trade-off

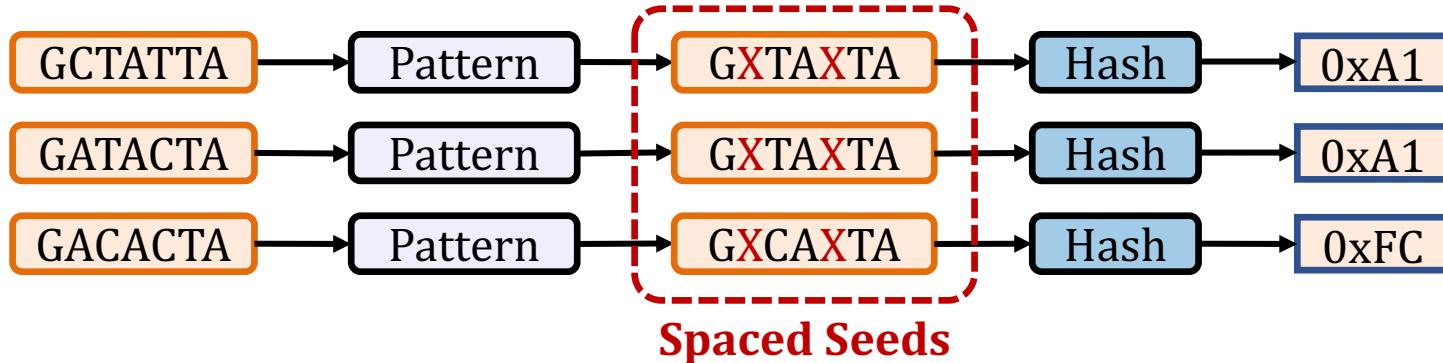


# Seed Matching Techniques

## 2. Allowing mismatches at **certain positions**

- **Spaced seeds**
- **Choice of pattern** is critical for the effectiveness of spaced seeds

2

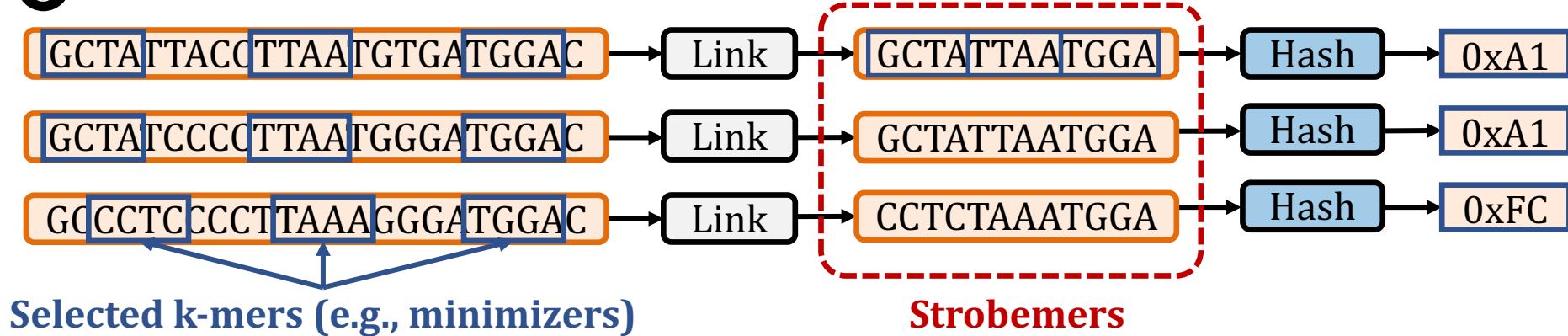


# Seed Matching Techniques

## 3. Allowing insertions and deletions

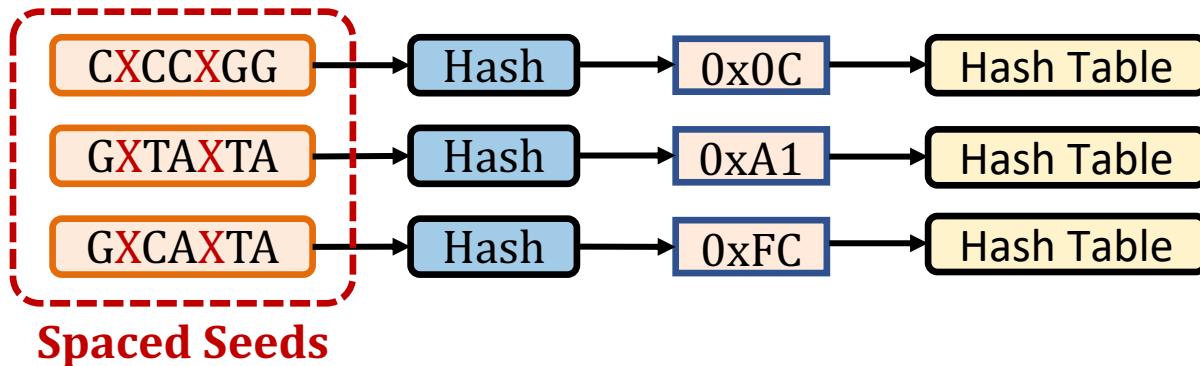
- **Linked k-mers** (e.g., strobemers)

③



# Hashing the Seeds

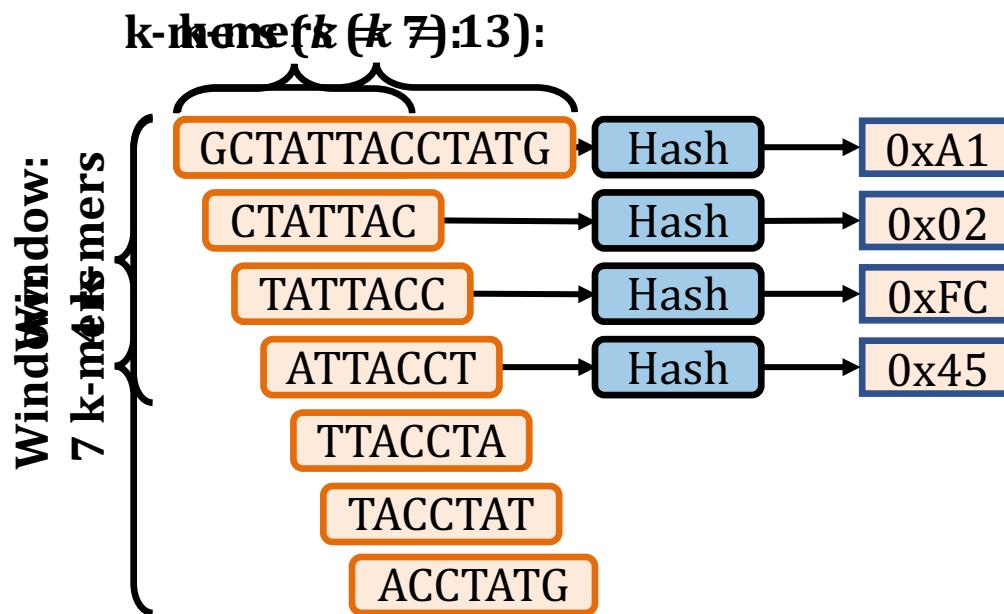
- Seeds are hashed
  - To find seed matches with a **direct lookup of hash values**



- **Low collision hash functions**
  - Seeds must **match exactly** to generate the **same hash value**
  - **Advantage:** **Dissimilar seeds** are **unlikely** to match
  - **Limitation:** **Highly similar seeds (fuzzy seeds)** are **unlikely** to match

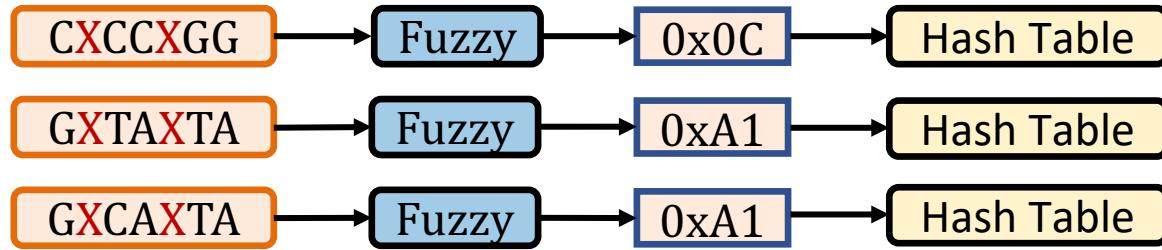
# Challenges - Exact-matching seeds

- **Limitations** when adjusting many **seeding parameters**
  - K-mer size
  - Window length (sampling ratio)
  - Determines the content of the hash table
- Trade-off between **performance, memory, and accuracy**



# Opportunities - Fuzzy matches of seeds

- A mechanism for **finding fuzzy seed matches** can enable
  - Assigning the **same** hash values to **highly similar seeds**
  - **Different** hash values for **dissimilar seeds**
  - **High performance** (e.g., no distance or similarity calculation) and
  - **Space-efficient** (no multiple hash functions for a single sketch) seed matching



- Finding **useful and novel seed matches** that cannot be identified when finding only exact-matching seeds
- **Rethinking the seeding parameters** to achieve better trade-off between
  - Performance, memory, and accuracy

# Outline

Background

Goal and Key Ideas

BLEND

Evaluation

Conclusions

# Our Goal

Enable finding **fuzzy matches of seeds**  
as well as **exact-matching seeds**  
with a **single lookup** of hash values of seeds



# BLEND

Uses a hashing mechanism, **SimHash**, that can generate **the same hash values** for **similar sets**

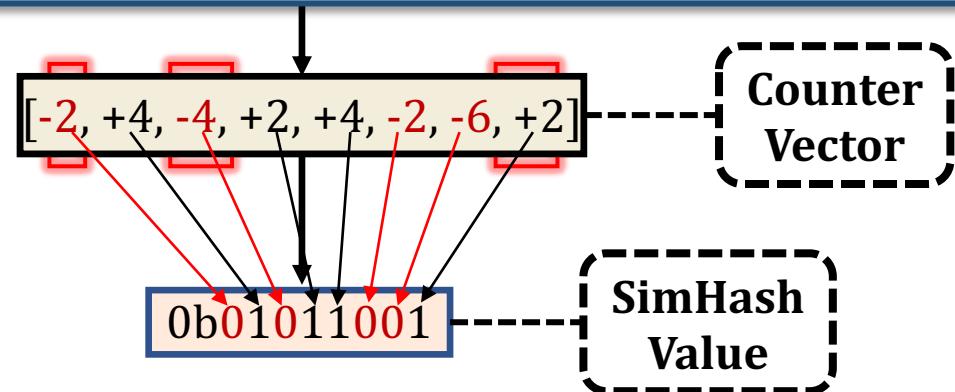
Provides the mechanisms for accurately and efficiently **converting seed sequences into set of items**

# The SimHash Technique – Example

- Goal: Generate the **same hash value** for **similar set of items**
  - **Example input:** A sentence (a set of items)
  - **Items:** Words in a sentence (hash values of items)
- Count the net difference between 0s and 1s at each position

**Challenge:** Efficiently and accurately  
convert seeds to set of items  
to use with SimHash

Example	0b11110000
sentence	0b11110000
to	0b01011101
generate	0b10011100
a	0b11000001
SimHash	0b01001001
value	0b00101011



# Outline

Background

Goal and Key Ideas

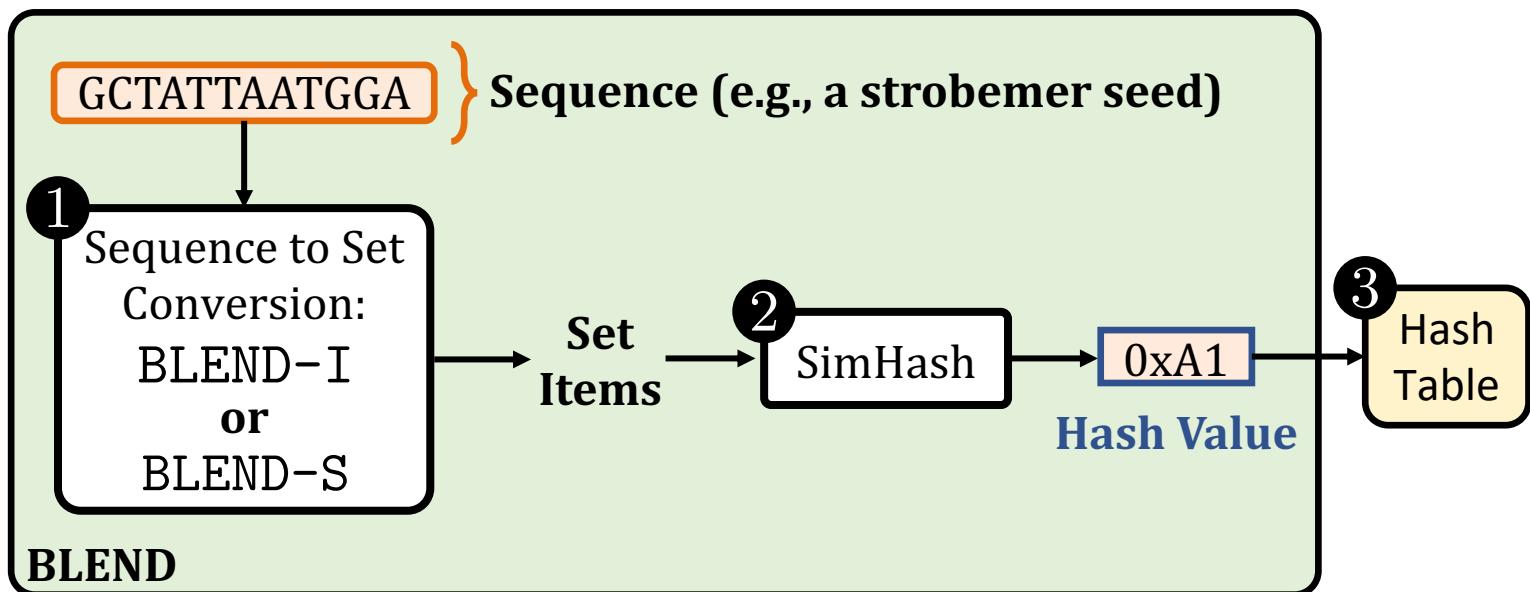
BLEND

Evaluation

Conclusions

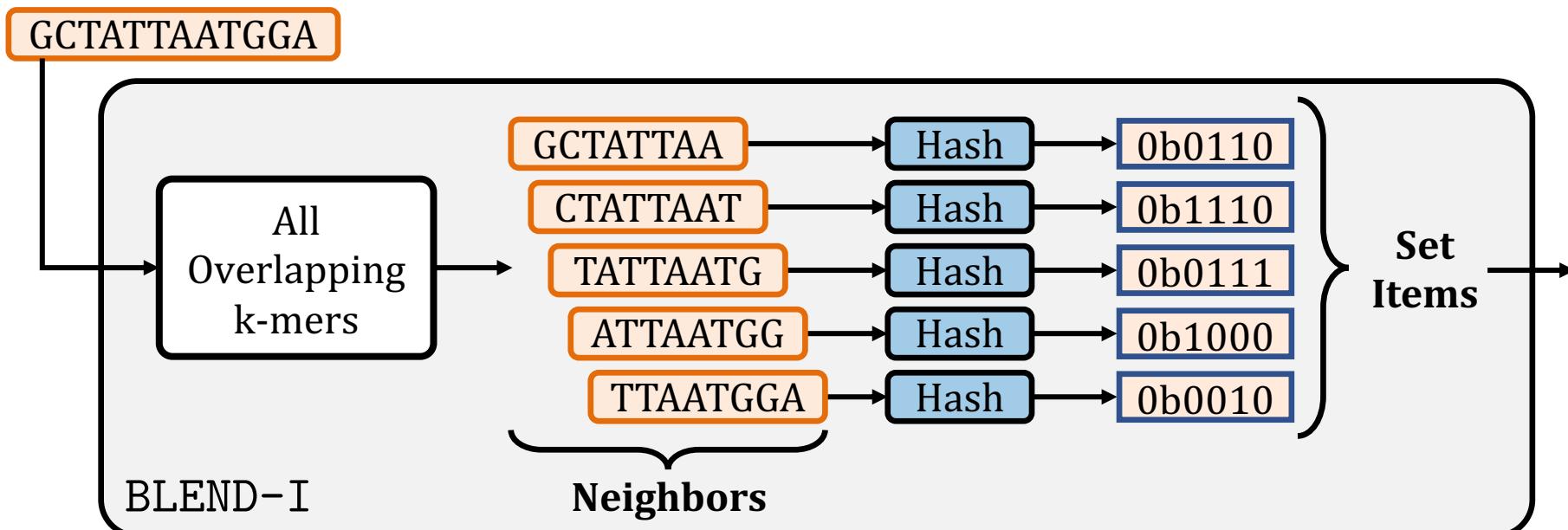
# BLEND Overview

- **Goal:** Efficiently find fuzzy seed matches with a single lookup
  - **Input:** A fixed-length sequence (seed sequence)
- 1. Efficiently and accurately **convert the seed** to its **set of items**
  - Two conversion mechanisms: **BLEND-I** and **BLEND-S**
- 2. Generate the SimHash value of the seed
- 3. Efficiently identify fuzzy seed matches by matching the SimHash values using a hash table



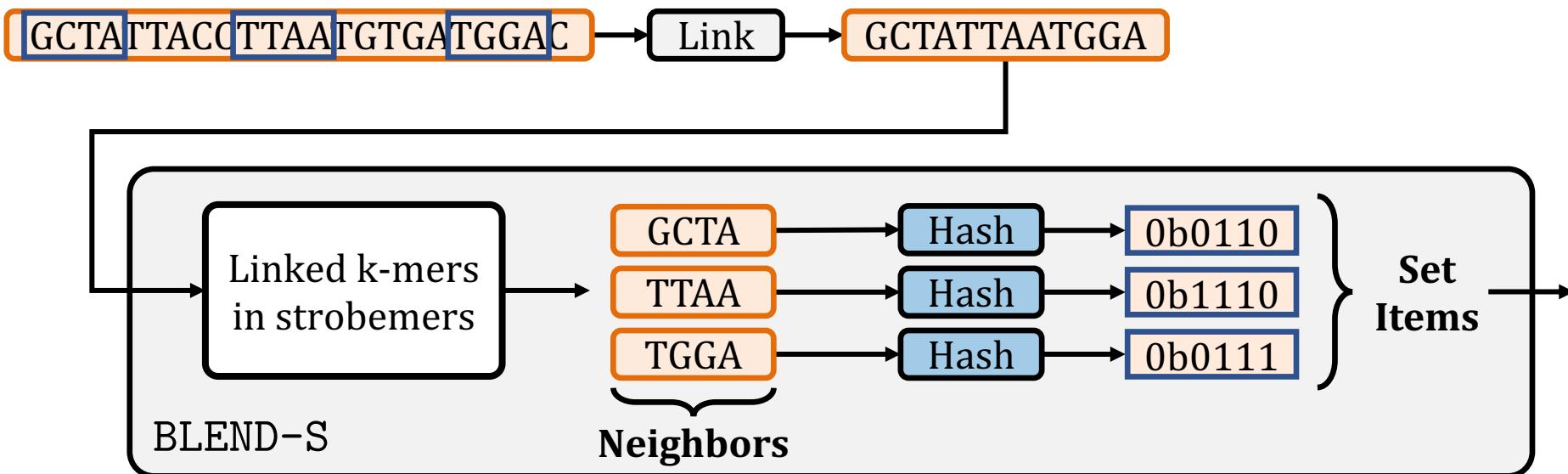
# Sequence-to-Set Conversion (BLEND-I)

- **Goal:** Convert seed sequences into set of items
    - **Input:** A fixed-length sequence (seed sequence)
1. Extract **all overlapping k-mers** of the seed (**neighbors**)
  2. Generate the **hash values of neighbors** using any hash function
  3. **Set items:** Hash values of neighbors



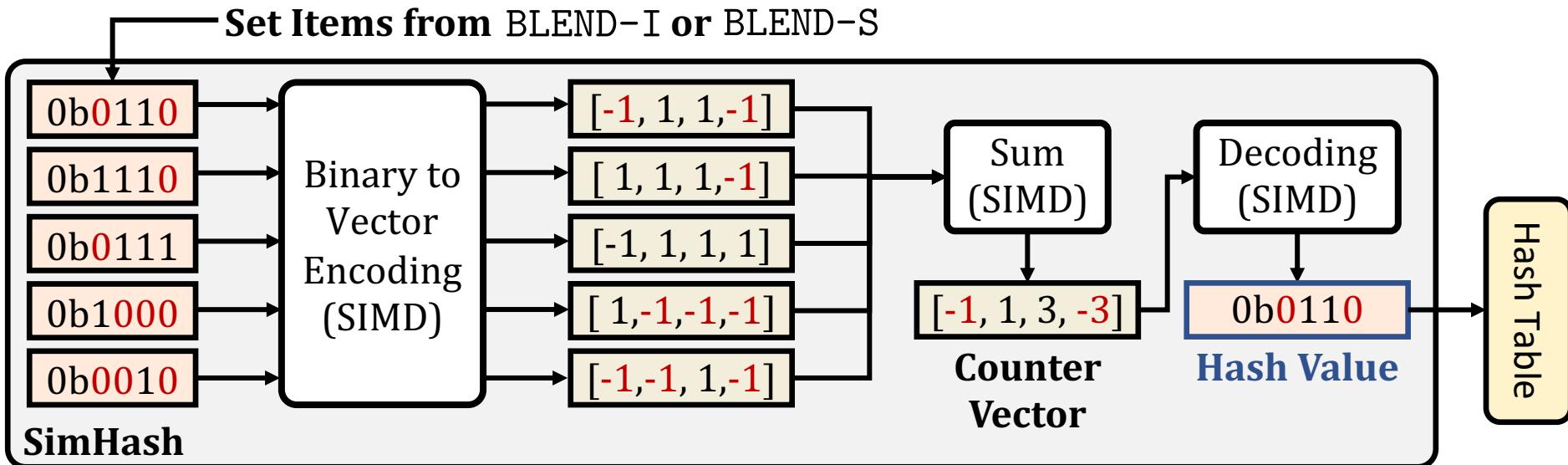
# Sequence-to-Set Conversion (BLEND-S)

- **Goal:** Convert seed sequences into set of items
    - **Input:** A fixed-length sequence (seed sequence)
1. Extract **all linked k-mers** of the seed (**neighbors**)
  2. Generate the **hash values of neighbors** using any hash function
  3. **Set items:** Hash values of neighbors

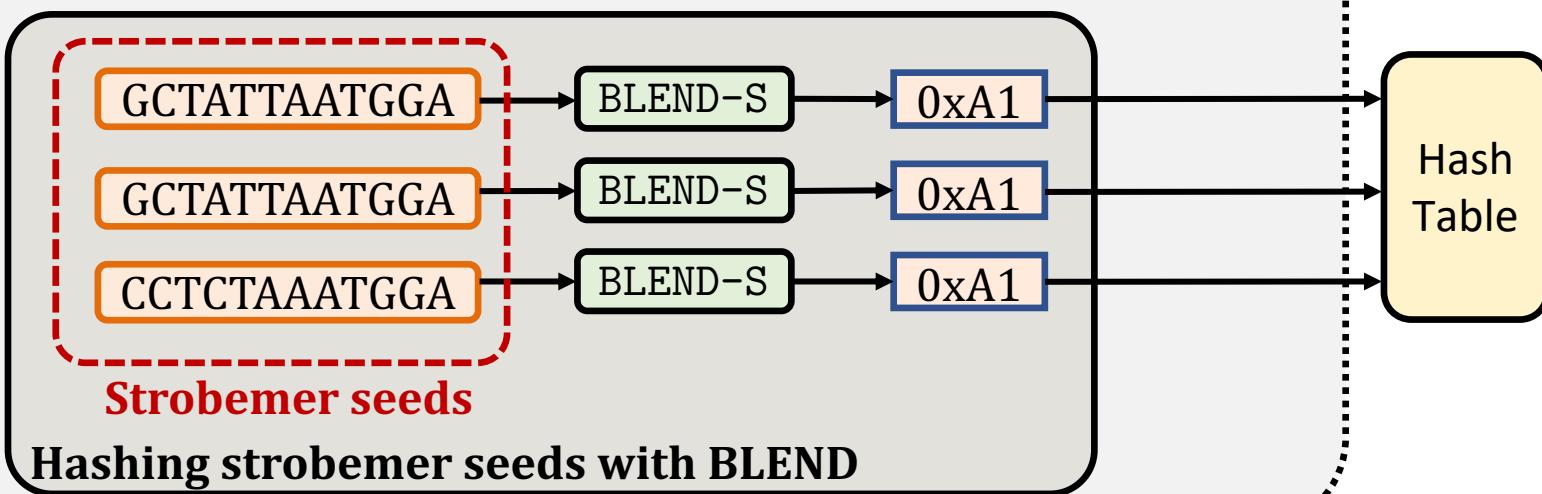
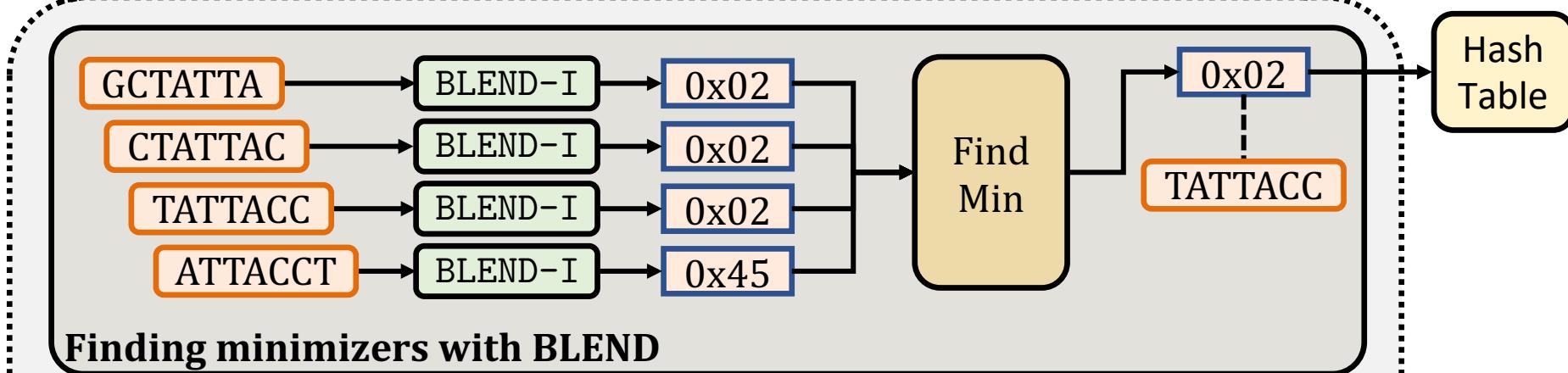


# Generating the SimHash values

- **Goal:** Generate the SimHash value of a seed
    - **Input:** Set items from BLEND-I or BLEND-S
1. Encode hash values using vectors of **-1s** and **+1s**
  2. Bitwise sum in SimHash: **Vector summation**
  3. Decode the counter vector into a **SimHash value for the seed**



# Integrating BLEND for Seeding



# Outline

Background

Goal and Key Ideas

BLEND

Evaluation

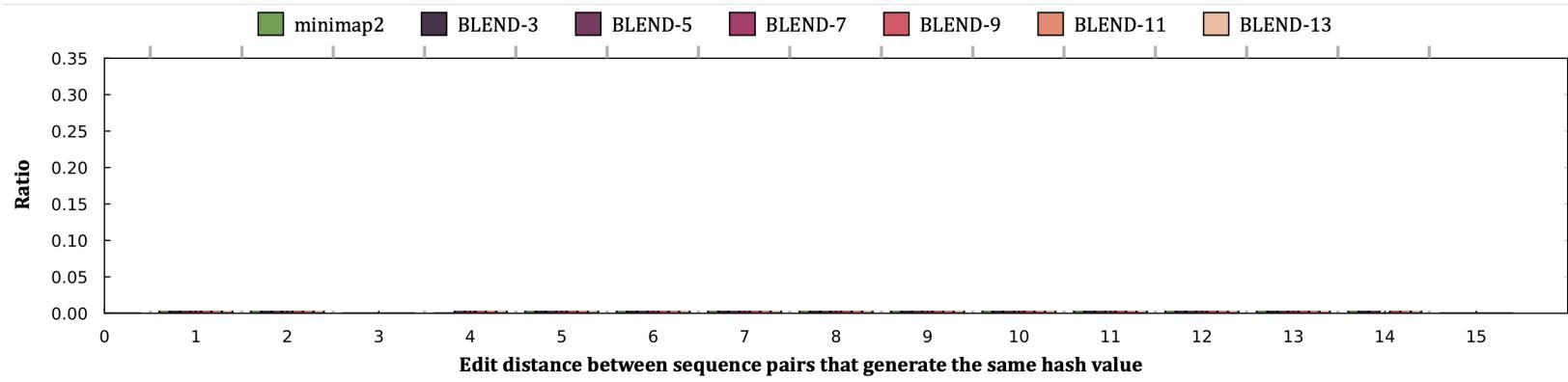
Conclusions

# Evaluation Methodology

- **Integrated into minimap2** to perform end-to-end mapping
- Real and simulated datasets from
  - PacBio (HiFi and CLR), ONT, and Illumina reads
  - Human CHM13 and HG002, Fruit fly, Yeast, and E. coli genomes
- Use case 1: **Read overlapping** (all-vs-all overlapping)
  - Evaluated the **accuracy, completeness, and contiguity** of *de novo* assemblies generated from overlaps
- Use case 2: **Read mapping** to a reference genome
  - Mapping accuracy from simulated reads
  - Structural variant calling

# Empirical Analysis on Fuzzy Seed Matching

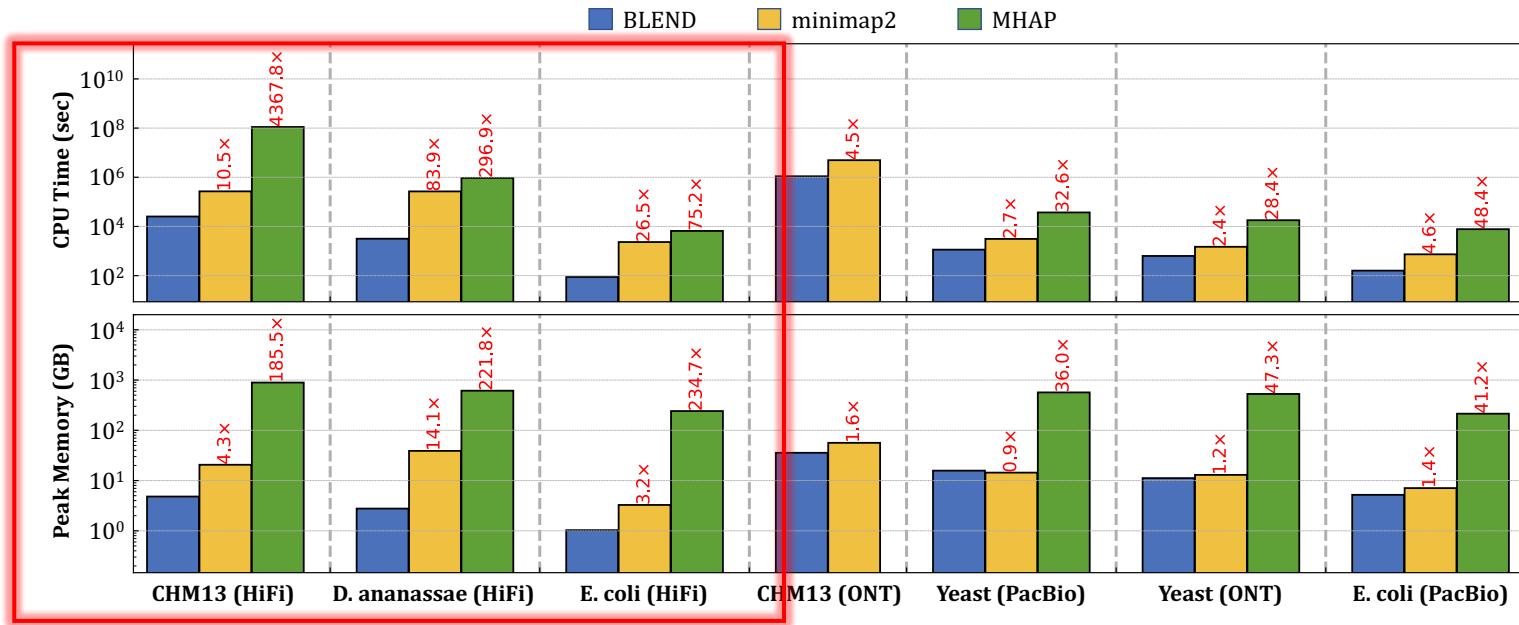
- Non-identical minimizers with the same hash value (collision)
  - Edit distance between minimizer k-mers with the same hash value
  - Ratio of collisions with a certain edit distance using minimap2 and BLEND
  - **Goal:** Increase the collision rate for highly similar seeds



BLEND **increases the collision rate** for highly similar seeds

while keeping a **collision rate similar to minimap2** for dissimilar seeds

# Read Overlapping – Performance & Memory



For HiFi: Average **speedup** of 40.3x (minimap2)

Reducing the **memory** footprint by 7.2x

Improving critical parameters without hurting the accuracy:

**Window length** (200) and **seed length** (31-mers)

# Read Overlapping – Assembly Evaluation

Dataset	Tool	Average Identity (%)	Genome Fraction (%)	k-mer Compl. (%)
<i>CHM13</i> (HiFi)	BLEND	<b>99.8526</b>	<b>98.4847</b>	<b>90.15</b>
	minimap2	99.7421	97.1493	83.05
	MHAP	N/A	N/A	N/A
	Reference	100	100	100
<i>D. ananassae</i> (HiFi)	BLEND	<b>99.7856</b>	<b>97.2308</b>	<b>86.43</b>
	minimap2	99.7044	96.3190	72.33
	MHAP	99.5551	0.7276	0.21
	Reference	100	100	100
<i>E. coli</i> (HiFi)	BLEND	<b>99.8320</b>	<b>99.8801</b>	<b>87.91</b>
	minimap2	99.7064	99.8748	79.27
	MHAP	N/A	N/A	N/A
	Reference	100	100	100

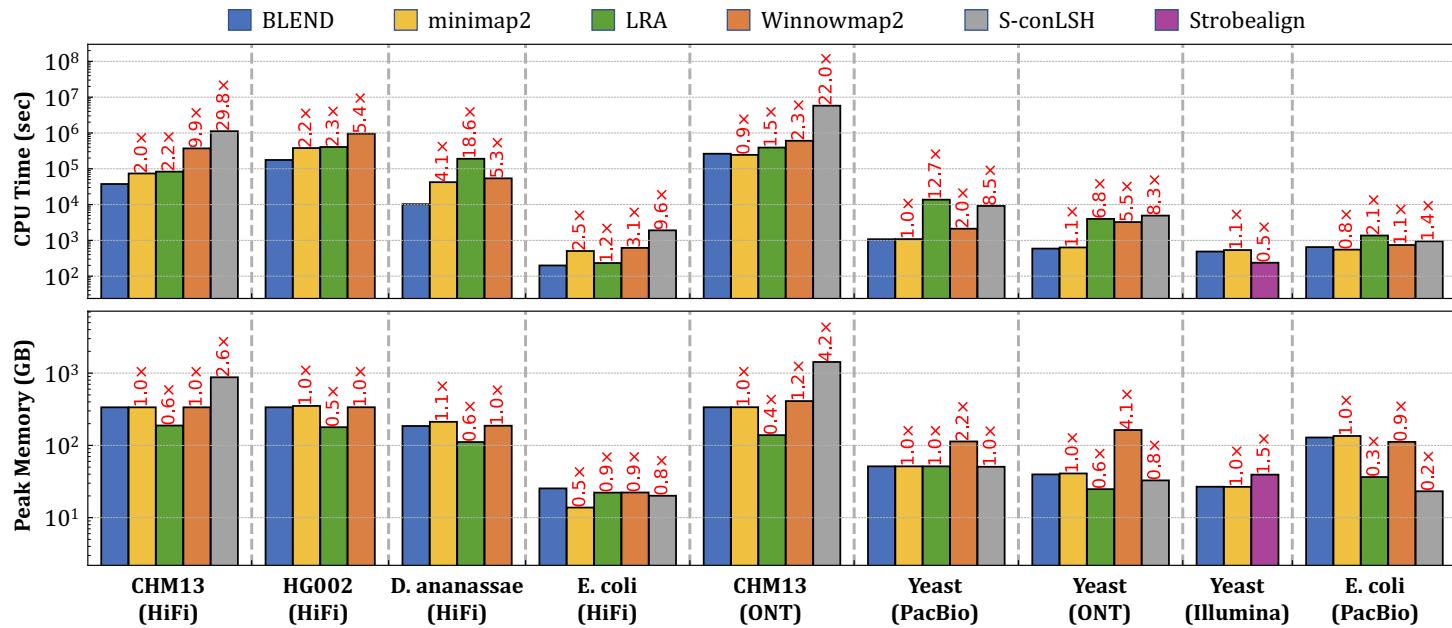
Largest Contig (Mb)	NGA50 (kb)	NG50 (kb)
22.8397	5442.25	5442.31
<b>47.1387</b>	<b>7133.43</b>	<b>7134.31</b>
N/A	N/A	N/A
248.387	154 260	154 260
<b>6.23256</b>	<b>792.407</b>	<b>798.913</b>
4.43396	273.398	278.775
0.028586	N/A	N/A
30.6728	26 427.4	26 427.4
<b>3.41699</b>	<b>3416.99</b>	<b>3416.99</b>
3.08849	3087.05	3087.05
N/A	N/A	N/A
4.94446	4944.46	4944.46

For HiFi: More accurate and complete assemblies

Better contiguity in most cases

Fuzzy seed matches can lead to finding **novel and useful overlaps**

# Read Mapping – Performance & Memory



Average **speedup** of 1.7x (minimap2) and similar peak memory usage

The computational cost of sequence alignment **slightly hinders**  
the benefits of fuzzy seed matching

# Read Mapping – SV Calling

- Structural variant (SV) calling using read mappings from each tool
  - Sniffles2 to call SVs from HG002 long read mappings
  - Truvari to compare the resulting SVs with the benchmarking SV set (Tier 1 set from GIAB)

Tool	HG002 SVs (high-confidence tier 1 SV set)					
	TP (#)	FP (#)	FN (#)	Precision	Recall	$F_1$
BLEND	<b>9229</b>	855	<b>412</b>	0.9152	<b>0.9573</b>	<b>0.9358</b>
minimap2	9222	915	419	0.9097	0.9565	0.9326
LRA	9155	<b>830</b>	486	<b>0.9169</b>	0.9496	0.9329
Winnowmap2	9170	1029	471	0.8991	0.9511	0.9244

Best overall accuracy in downstream analysis

# Outline

Background

Goal and Key Ideas

BLEND

Evaluation

Conclusions

# BLEND Summary

## Problem

Finding exact-matching seeds introduce limitations in further improving the performance and accuracy of genome analysis

## Goal

Enable finding the fuzzy seed matches as well as the exact-matching seeds accurately and efficiently

## BLEND

- Provides effective mechanisms for converting seeds into set of items to use with the SimHash technique

## Key Results

- **Significant speedups** and **lower memory footprint** especially when using **HiFi reads**
- **Improves the accuracy** of important applications in genome analysis

# BLEND

- Can Firtina, Jisung Park, Mohammed Alser, Jeremie S. Kim, Damla Senol Cali, Taha Shahroodi, Nika Mansouri Ghiasi, Gagandeep Singh, Konstantinos Kanellopoulos, Can Alkan, and Onur Mutlu,

**"BLEND: A Fast, Memory-Efficient, and Accurate Mechanism to Find Fuzzy Seed Matches in Genome Analysis"**

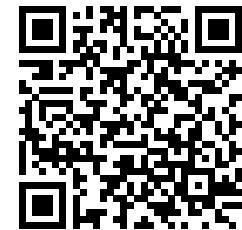
***NAR Genomics and Bioinformatics***, March 2023.

[[Online link at NAR Genomics and Bioinformatics Journal](#)]

[[arXiv preprint](#)]

[[biorXiv preprint](#)]

[[BLEND Source Code](#)]



Paper (NARGAB)



Volume 5, Issue 1  
March 2023

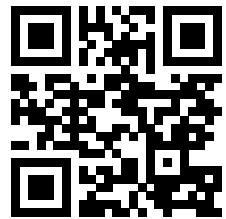
#### JOURNAL ARTICLE

## BLEND: a fast, memory-efficient and accurate mechanism to find fuzzy seed matches in genome analysis

Can Firtina , Jisung Park, Mohammed Alser, Jeremie S Kim, Damla Senol Cali, Taha Shahroodi, Nika Mansouri Ghiasi, Gagandeep Singh, Konstantinos Kanellopoulos, Can Alkan, Onur Mutlu 

*NAR Genomics and Bioinformatics*, Volume 5, Issue 1, March 2023, lqad004,

# BLEND Source Code



[Source Code](#)

Screenshot of the BLEND GitHub repository page:

**Code** | **Issues** | **Pull requests** | **Actions** | **Projects** | **Wiki** | **Security** | **Insights** | **Settings**

master | 1 branch | 1 tag | Go to file | Add file | <> Code

**About**

BLEND is a mechanism that can efficiently find fuzzy seed matches between sequences to significantly improve the performance and accuracy while reducing the memory space usage of two important applications: 1) finding overlapping reads and 2) read mapping. Described by Firtina et al. (published in NARGAB)  
<https://doi.org/10.1093/nargab/lqad004>

**bioinformatics** | **genome-analysis**  
**genome-assembly** | **blend**  
**read-mapping** | **de-novo-assembly**  
**minimizers** | **strobemers** | **seed-matching**  
**fuzzy-seeds** | **read-overlapping**  
**spaced-seeds**

**README.md**

**BLEND: A Fast, Memory-Efficient, and Accurate Mechanism to Find Fuzzy Seed Matches in Genome Analysis**

Readme | View license | Code of conduct | 24 stars | 12 watching | 2 forks | Report repository

<https://github.com/CMU-SAFARI/BLEND>



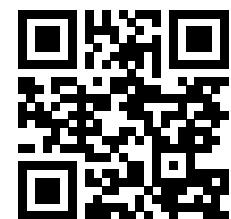
# BLEND

A Fast, Memory-efficient and Accurate Mechanism  
to Find Fuzzy Seed Matches in Genome Analysis

**Can Firtina**, Jisung Park, Mohammed Alser, Jeremie S. Kim, Damla Senol Cali,  
Taha Shahroodi, Nika Mansouri Ghiasi, Gagandeep Singh,  
Konstantinos Kanellopoulos, Can Alkan, Onur Mutlu



[Paper \(NARGAB\)](#)



[Source Code](#)

**SAFARI**

**ETH** zürich

**Carnegie Mellon**

**TU**Delft

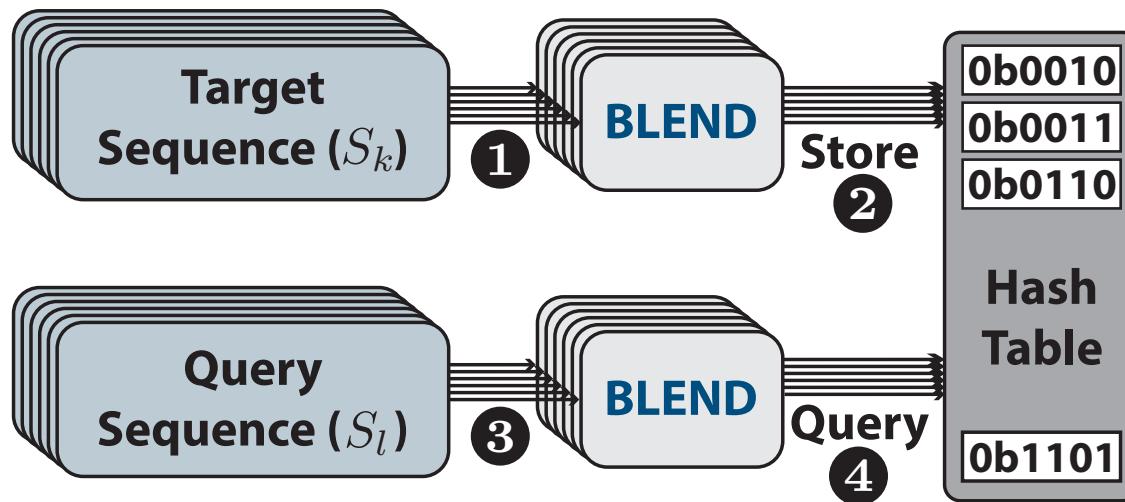
Delft University of Technology



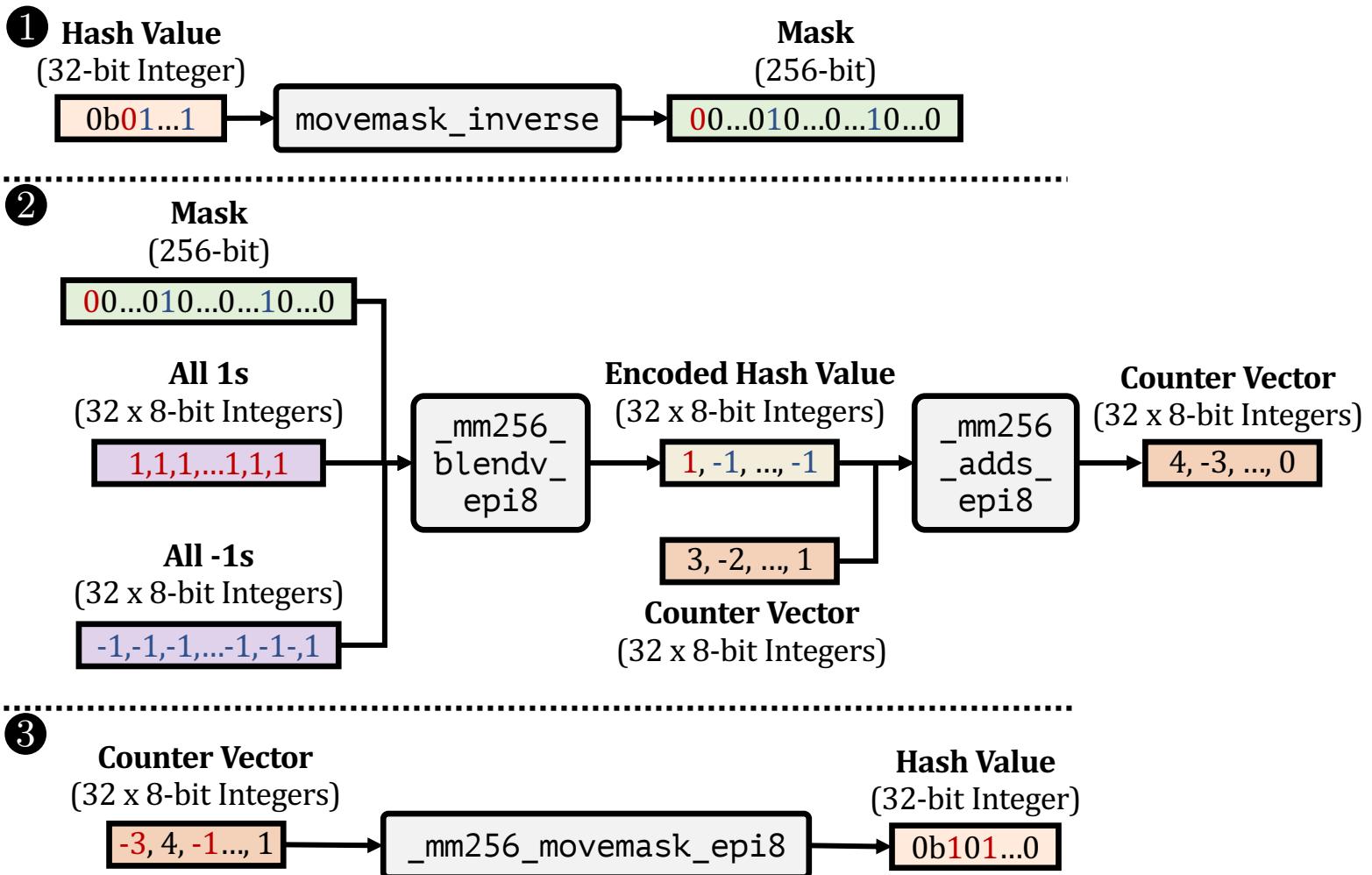
**Bilkent University**

# Backup Slides

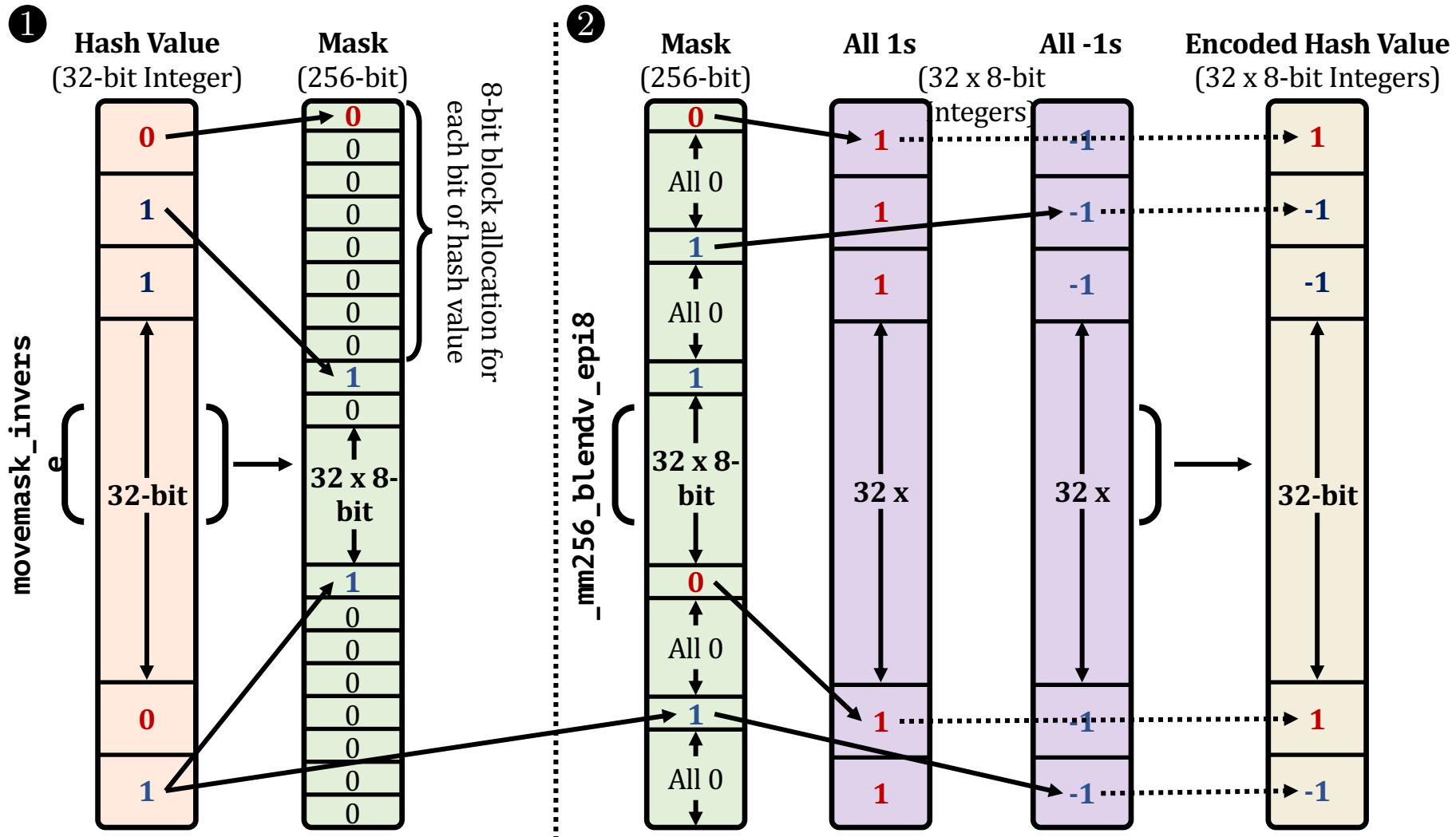
# Integrating BLEND in Read Mappers



# Calculating SimHash using SIMD



# Encoding Hash Values into Vectors using SIMD



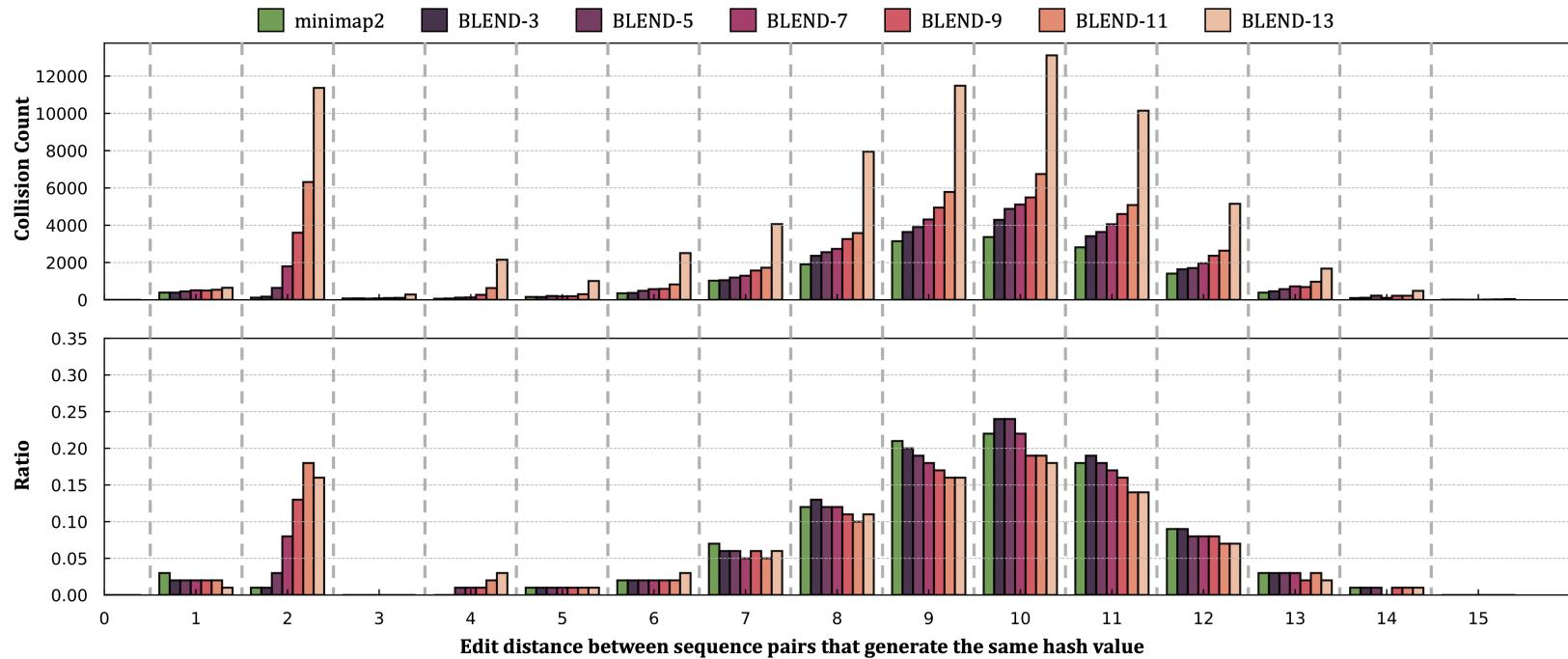
# Dataset

Organism	Library	Reads (#)	Seq. Depth	SRA Accession	Reference Genome
<i>Human CHM13</i>	PacBio HiFi	3 167 477	16	SRR11292122-3	T2T-CHM13 (v1.1)
	ONT*	10 380 693	30	Simulated R9.5	T2T-CHM13 (v2.0)
<i>Human HG002</i>	PacBio HiFi	11 714 594	52	SRR10382244-9	GRCh37
<i>D. ananassae</i>	PacBio HiFi	1 195 370	50	SRR11442117	(90)
<i>Yeast</i>	PacBio CLR*	270 849	200	Simulated P6-C4	GCA_000146045.2
	ONT*	135 296	100	Simulated R9.5	GCA_000146045.2
	Illumina MiSeq	3 318 467	80	ERR1938683	GCA_000146045.2
<i>E. coli</i>	PacBio HiFi	38 703	100	SRR11434954	(90)
	PacBio CLR	76 279	112	SRR1509640	GCA_000732965.1

\* We use PBSIM2 to generate the simulated PacBio and ONT reads.  
We show the simulated chemistry under the SRA Accession column.

# Empirical Analysis on Fuzzy Seed Matching

- Minimizer collisions and the edit distance between collisions



BLEND **increases the collision rate** for **highly similar seeds**  
while keeping a **collision rate similar to minimap2**  
for dissimilar seeds

# Finding Fuzzy Seed Matching between Minimizers

Tool	Number of Minimizers	Number of Collisions	Collision/Minimizer Ratio	Avg. Edit Distance Between Minimizers With Collision
minimap2	903,043	15,306	0.016949	9.327061
BLEND-3	1,014,173	18,224	0.017969	9.393437
BLEND-5	1,090,468	20,659	0.018945	9.213660
BLEND-7	1,140,254	23,591	0.020689	8.874698
BLEND-9	1,173,198	28,411	0.024217	8.495301
BLEND-11	1,186,687	35,500	0.029915	8.067549
BLEND-13	1,197,966	72,078	0.060167	8.075918

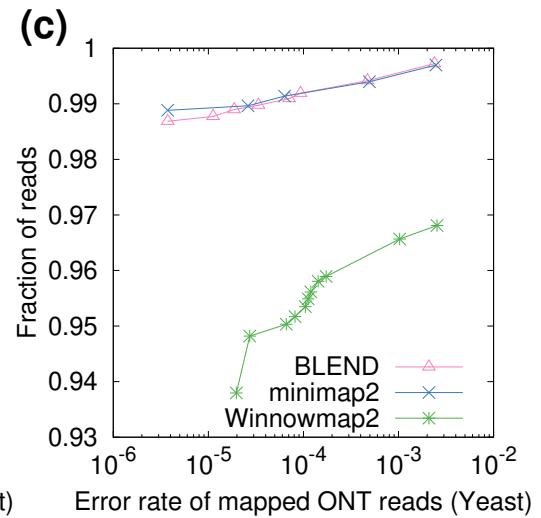
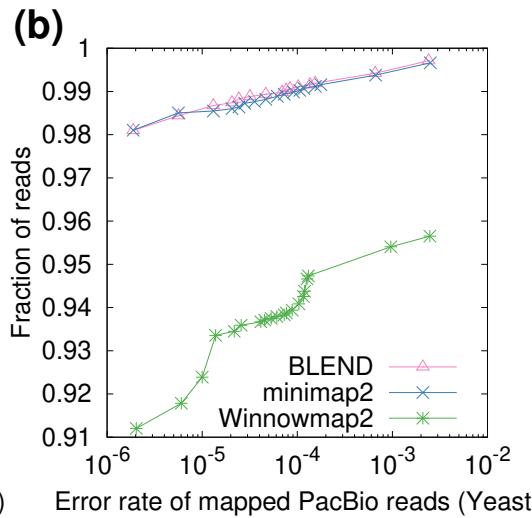
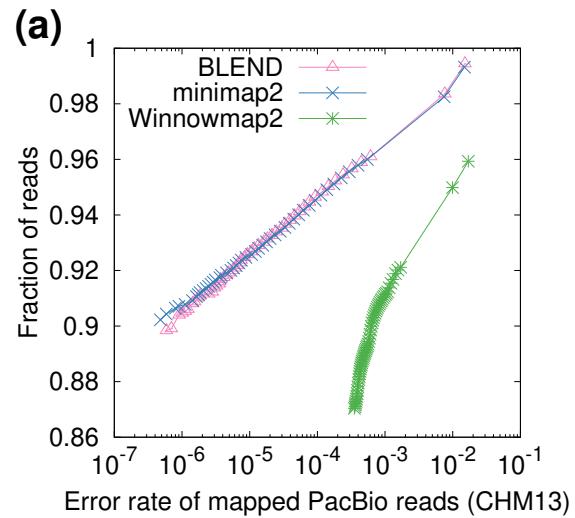
# Finding Fuzzy Seed Matching between Similar Sequences

Tool	Number of Sequences	Number of Sequences with Collision	Collision/Sequence Ratio	Avg. Edit Distance Between K-mers With Collision
minimap2	4,130	0	0	N/A
BLEND-3	4,130	0	0	N/A
BLEND-5	4,130	11	0.00263663	1.45455
BLEND-7	4,130	50	0.0119847	1.5
BLEND-9	4,130	77	0.0184564	2.01299
BLEND-11	4,130	273	0.0654362	2.80952
BLEND-13	4,130	329	0.0788591	2.20669

# Read Overlapping – Assembly Accuracy

Dataset	Tool	Mismatch									
		Average Identity (%)	Genome Fraction (%)	k-mer Compl. (%)	Aligned Length (Mb)	per 100 kb (#)	Average GC (%)	Assembly Length (Mb)	Largest Contig (Mb)	NGA50 (kb)	NG50 (kb)
<i>CHM13</i> (HiFi)	BLEND	<b>99.8526</b>	<b>98.4847</b>	<b>90.15</b>	3092.54	<b>22.02</b>	<b>40.78</b>	<b>3095.21</b>	22.8397	5442.25	5442.31
	minimap2	99.7421	97.1493	83.05	<b>3094.79</b>	55.96	40.71	3100.97	<b>47.1387</b>	<b>7133.43</b>	<b>7134.31</b>
	MHAP	N/A	N/A	N/A	N/A	N/A	N/A	N/A	N/A	N/A	N/A
	Reference	100	100	100	3054.83	0.00	40.85	3054.83	248.387	154 260	154 260
<i>D. ananassae</i> (HiFi)	BLEND	<b>99.7856</b>	<b>97.2308</b>	<b>86.43</b>	240.391	<b>143.13</b>	<b>41.75</b>	<b>247.153</b>	<b>6.23256</b>	<b>792.407</b>	<b>798.913</b>
	minimap2	99.7044	96.3190	72.33	<b>289.453</b>	191.53	41.68	298.28	4.43396	273.398	278.775
	MHAP	99.5551	0.7276	0.21	2.29	239.76	42.07	2.34951	0.028586	N/A	N/A
	Reference	100	100	100	213.805	0.00	41.81	213.818	30.6728	26 427.4	26 427.4
<i>E. coli</i> (HiFi)	BLEND	<b>99.8320</b>	<b>99.8801</b>	<b>87.91</b>	<b>5.12155</b>	<b>3.77</b>	<b>50.53</b>	5.12155	<b>3.41699</b>	<b>3416.99</b>	<b>3416.99</b>
	minimap2	99.7064	99.8748	79.27	5.09249	19.71	50.47	<b>5.09436</b>	3.08849	3087.05	3087.05
	MHAP	N/A	N/A	N/A	N/A	N/A	N/A	N/A	N/A	N/A	N/A
	Reference	100	100	100	5.04628	0.00	50.52	5.04628	4.94446	4944.46	4944.46
<i>CHM13</i> (ONT)	BLEND	N/A	N/A	<b>29.26</b>	<b>2891.28</b>	<b>4077.53</b>	<b>41.32</b>	2897.87	25.2071	5061.52	5178.59
	minimap2	N/A	N/A	28.32	2860.26	4660.73	41.36	<b>2908.55</b>	<b>66.7564</b>	<b>13 189.2</b>	<b>13 820.3</b>
	Reference	100	100	100	3117.29	0.00	40.75	3117.29	248.387	150 617	150 617
<i>Yeast</i> (PacBio)	BLEND	89.1677	<b>97.0854</b>	<b>33.81</b>	<b>12.3938</b>	2672.37	38.84	12.4176	1.54807	635.966	636.669
	minimap2	88.9002	96.9709	33.38	12.0128	2684.38	38.85	<b>12.3325</b>	<b>1.56078</b>	<b>810.046</b>	<b>828.212</b>
	MHAP	<b>89.2182</b>	88.5928	32.39	10.9039	<b>2552.05</b>	<b>38.81</b>	10.9896	1.02375	85.081	436.285
	Reference	100	100	100	12.1571	0.00	38.15	12.1571	1.53193	924.431	924.431
<i>Yeast</i> (ONT)	BLEND	<b>89.6889</b>	99.2974	<b>35.95</b>	<b>12.3222</b>	<b>2529.47</b>	<b>38.64</b>	<b>12.3225</b>	1.10582	793.046	793.046
	minimap2	88.9393	<b>99.6878</b>	34.84	12.304	2782.59	38.74	12.3725	<b>1.56005</b>	<b>796.718</b>	<b>941.588</b>
	MHAP	89.1970	89.2785	33.58	10.8302	2647.19	38.84	10.9201	1.44328	118.886	618.908
	Reference	100	100	100	12.1571	0.00	38.15	12.1571	1.53193	924.431	924.431
<i>E. coli</i> (PacBio)	BLEND	<b>88.5806</b>	<b>96.5238</b>	<b>32.32</b>	<b>5.90024</b>	<b>1857.56</b>	<b>49.81</b>	6.21598	2.40671	<b>769.981</b>	2060.4
	minimap2	88.1365	92.7603	30.74	5.37728	2005.72	49.66	<b>6.02707</b>	<b>3.77098</b>	367.442	<b>3770.98</b>
	MHAP	88.4883	90.5533	31.32	5.75159	1999.48	49.69	6.26216	1.04286	110.535	456.01
	Reference	100	100	100	5.6394	0.00	50.43	5.6394	5.54732	5547.32	5547.32

# Read Mapping – Mapping Accuracy



# Read Mapping – Mapping Accuracy

Dataset	Overall error rate (%)		
	BLEND	minimap2	Winnowmap2
<i>CHM13</i> (ONT)	1.5168427	<b>1.4914009</b>	1.7001222
<i>Yeast</i> (PacBio)	<b>0.2403134</b>	0.2504307	0.2474206
<i>Yeast</i> (ONT)	<b>0.2386617</b>	0.2468770	0.2534777

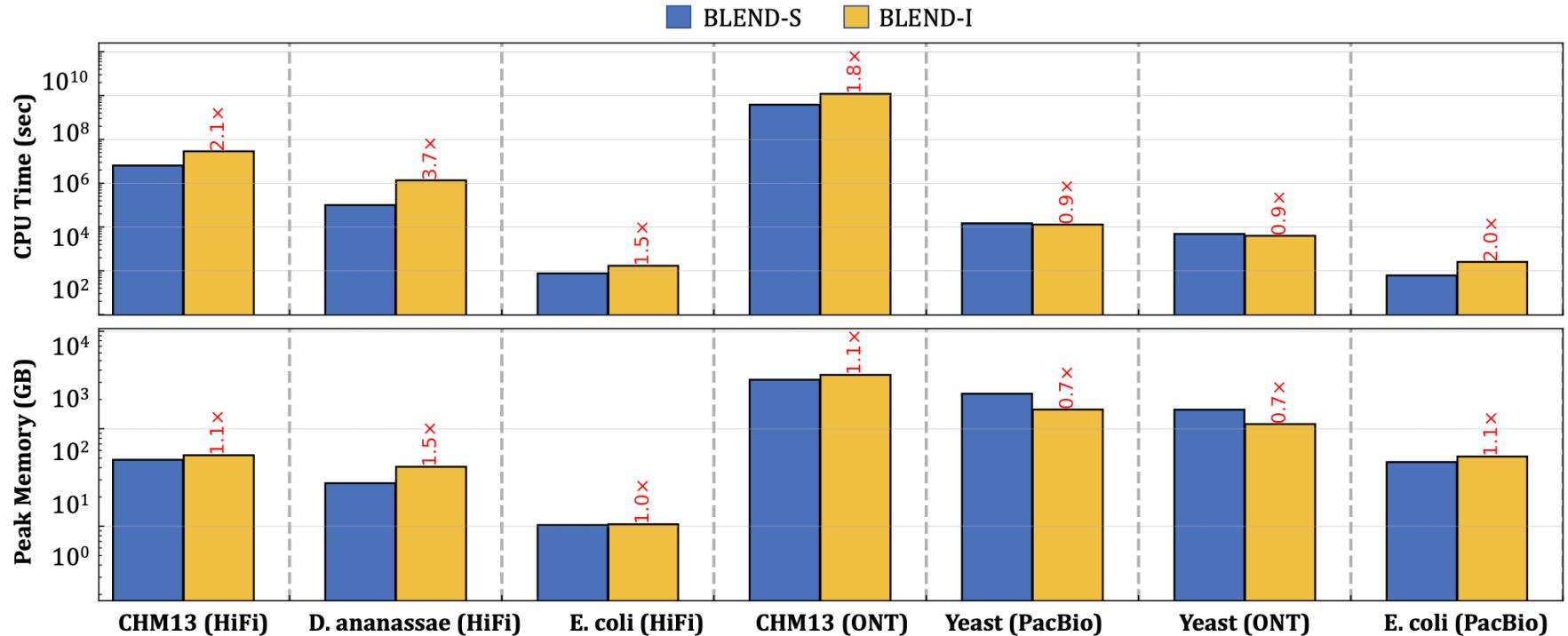
# Read Mapping Quality

Dataset	Tool	Average depth of Cov. (x)	Breadth of coverage (%)	Aligned reads (#)	Properly paired (%)
<i>CHM13</i> (HiFi)	BLEND	<b>16.58</b>	<b>99.991</b>	3 171 916	NA
	minimap2	<b>16.58</b>	<b>99.991</b>	<b>3 172 261</b>	NA
	LRA	16.37	99.064	3 137 631	NA
	Winnowmap2	<b>16.58</b>	99.990	3 171 313	NA
<i>HG002</i> (HiFi)	BLEND	51.25	92.245	11 424 762	NA
	minimap2	53.08	92.242	12 407 589	NA
	LRA	52.48	<b>92.275</b>	<b>13 015 195</b>	NA
	Winnowmap2	<b>53.81</b>	92.248	12 547 868	NA
<i>D. ananassae</i> (HiFi)	BLEND	57.37	99.662	1 223 388	NA
	minimap2	<b>57.57</b>	<b>99.665</b>	1 245 931	NA
	LRA	57.06	99.599	1 235 098	NA
	Winnowmap2	57.40	99.663	<b>1 249 575</b>	NA
<i>E. coli</i> (HiFi)	BLEND	<b>99.14</b>	99.897	39 048	NA
	minimap2	<b>99.14</b>	99.897	<b>39 065</b>	NA
	LRA	99.10	99.897	39 063	NA
	Winnowmap2	<b>99.14</b>	99.897	39 036	NA
<i>CHM13</i> (ONT)	BLEND	<b>29.34</b>	<b>99.999</b>	<b>10 322 767</b>	NA
	minimap2	29.33	<b>99.999</b>	10 310 182	NA
	LRA	28.84	99.948	9 999 432	NA
	Winnowmap2	28.98	99.936	9 958 402	NA
<i>Yeast</i> (PacBio)	BLEND	<b>195.87</b>	<b>99.980</b>	<b>270 064</b>	NA
	minimap2	195.86	<b>99.980</b>	269 935	NA
	LRA	194.65	99.967	267 399	NA
	Winnowmap2	192.35	99.977	259 073	NA
<i>Yeast</i> (ONT)	BLEND	<b>97.88</b>	<b>99.964</b>	<b>134 919</b>	NA
	minimap2	<b>97.88</b>	<b>99.964</b>	134 885	NA
	LRA	97.25	99.952	132 862	NA
	Winnowmap2	97.04	99.963	130 978	NA
<i>Yeast</i> (Illumina)	BLEND	<b>79.92</b>	<b>99.975</b>	6 493 730	95.88
	minimap2	79.91	99.974	6 492 994	95.89
	Strobealign	<b>79.92</b>	99.970	<b>6 498 380</b>	<b>97.59</b>
<i>E. coli</i> (PacBio)	BLEND	<b>97.51</b>	100	83 924	NA
	minimap2	97.29	100	<b>85 326</b>	NA
	LRA	93.61	100	80 802	NA
	Winnowmap2	89.78	100	69 884	NA

# BLEND-I vs. BLEND-S (Accuracy)

Dataset	Tool	Average Identity (%)	Genome Fraction (%)	K-mer Compl. (%)	Aligned Length (Mbp)	Mismatch per 100Kbp (#)	Average GC (%)	Assembly Length (Mbp)	Largest Contig (Mbp)	NGA50 (Kbp)	NG50 (Kbp)
<i>CHM13</i> (HiFi)	BLEND-I	99.7535	96.7203	83.65	3,054.49	48.49	<b>40.79</b>	<b>3,059.29</b>	<b>41.8342</b>	<b>8,507.53</b>	<b>8,508.92</b>
	BLEND-S	<b>99.8526</b>	<b>98.4847</b>	<b>90.15</b>	<b>3,092.54</b>	<b>22.02</b>	40.78	3,095.21	22.8397	5,442.25	5,442.31
	Reference	100	100	100	3,054.83	0.00	40.85	3,054.83	248.387	154,260	154,260
<i>D. ananassae</i> (HiFi)	BLEND-I	99.6890	97.2290	77.85	<b>270.218</b>	233.18	41.95	280.388	5.01099	356.745	356.745
	BLEND-S	<b>99.7856</b>	<b>97.2308</b>	<b>86.43</b>	240.391	<b>143.13</b>	<b>41.75</b>	<b>247.153</b>	<b>6.23256</b>	<b>792.407</b>	<b>798.913</b>
	Reference	100	100	100	213.805	0.00	41.81	213.818	30.6728	26,427.4	26,427.4
<i>E. coli</i> (HiFi)	BLEND-I	99.6902	<b>99.8824</b>	79.36	5.04157	17.92	<b>50.52</b>	<b>5.04263</b>	<b>4.94601</b>	<b>4,025.48</b>	<b>4,946.01</b>
	BLEND-S	<b>99.8320</b>	99.8801	<b>87.91</b>	<b>5.12155</b>	<b>3.77</b>	50.53	5.12155	3.41699	3,416.99	3,416.99
	Reference	100	100	100	5.04628	0.00	50.52	5.04628	4.94446	4,944.46	4,944.46
<i>CHM13</i> (ONT)	BLEND-I	N/A	N/A	<b>29.26</b>	<b>2,891.28</b>	4,077.53	<b>41.32</b>	<b>2,897.87</b>	<b>25.2071</b>	<b>5,061.52</b>	<b>5,178.59</b>
	BLEND-S	N/A	N/A	0	0.010546	<b>3,250.70</b>	51.30	0.010548	0.010548	0	0
	Reference	100	100	100	3,117.29	0.00	40.75	3,117.29	248.387	150,617	150,617
<i>Yeast</i> (PacBio)	BLEND-I	89.1677	<b>97.0854</b>	<b>33.81</b>	12.3938	<b>2,672.37</b>	38.84	<b>12.4176</b>	<b>1.54807</b>	<b>635.966</b>	<b>636.669</b>
	BLEND-S	<b>90.3347</b>	83.8814	33.17	<b>22.9473</b>	4,795.58	<b>38.71</b>	22.9523	0.265118	114.125	116.143
	Reference	100	100	100	12.1571	0.00	38.15	12.1571	1.53193	924.431	924.431
<i>Yeast</i> (ONT)	BLEND-I	89.6889	<b>99.2974</b>	<b>35.95</b>	<b>12.3222</b>	2,529.47	38.64	<b>12.3225</b>	<b>1.10582</b>	<b>793.046</b>	<b>793.046</b>
	BLEND-S	<b>91.0865</b>	7.9798	4.90	0.898565	<b>2,006.91</b>	<b>38.35</b>	0.899654	0.043321	0	0
	Reference	100	100	100	12.1571	0.00	38.15	12.1571	1.53193	924.431	924.431
<i>E. coli</i> (PacBio)	BLEND-I	88.5806	<b>96.5238</b>	<b>32.32</b>	<b>5.90024</b>	1,857.56	<b>49.81</b>	<b>6.21598</b>	<b>2.40671</b>	<b>769.981</b>	<b>2,060.4</b>
	BLEND-S	<b>90.3551</b>	36.6230	17.07	2.10137	<b>1,299.50</b>	48.91	2.10704	0.095505	0	0
	Reference	100	100	100	5.6394	0.00	50.43	5.6394	5.54732	5,547.32	5,547.32

# BLEND-I vs. BLEND-S (Performance and Memory)



# Read Overlapping Parameters

Tool	Dataset	Parameters
BLEND	<i>CHM13 (HiFi)</i>	-x ava-hifi -t 32
BLEND	<i>D. ananassae (HiFi)</i>	-x ava-hifi -t 32
BLEND	<i>E. coli (HiFi)</i>	-x ava-hifi -t 32
BLEND	<i>CHM13 (ONT)</i>	-x ava-ont -t 32
BLEND	<i>Yeast (PacBio)</i>	-x ava-pb -t 32
BLEND	<i>Yeast (ONT)</i>	-x ava-ont -t 32
BLEND	<i>E. coli (PacBio)</i>	-x ava-pb -t 32
minimap2	<i>CHM13 (HiFi)</i>	-x ava-pb -Hk21 -w14 -t 32
minimap2	<i>D. ananassae (HiFi)</i>	-x ava-pb -Hk21 -w14 -t 32
minimap2	<i>E. coli (HiFi)</i>	-x ava-pb -Hk21 -w14 -t 32
minimap2	<i>CHM13 (ONT)</i>	-x ava-ont -t 32
minimap2	<i>Yeast (PacBio)</i>	-x ava-pb -t 32
minimap2	<i>Yeast (ONT)</i>	-x ava-ont -t 32
minimap2	<i>E. coli (PacBio)</i>	-x ava-pb -t 32
minimap2-Eq	<i>CHM13 (ONT)</i>	-x ava-ont -k19 -w10 -t 32
minimap2-Eq	<i>Yeast (PacBio)</i>	-x ava-pb -k23 -w10 -t 32
minimap2-Eq	<i>Yeast (ONT)</i>	-x ava-ont -k19 -w10 -t 32
minimap2-Eq	<i>E. coli (PacBio)</i>	-x ava-pb -k23 -w10 -t 32
MHAP	<i>CHM13 (HiFi)</i>	--store-full-id --ordered-kmer-size 18 --num-hashes 128 --num-min-matches 5 --ordered-sketch-size 1000 --threshold 0.95 --num-threads 32
MHAP	<i>D. ananassae (HiFi)</i>	--store-full-id --ordered-kmer-size 18 --num-hashes 128 --num-min-matches 5 --ordered-sketch-size 1000 --threshold 0.95 --num-threads 32
MHAP	<i>E. coli (HiFi)</i>	--store-full-id --ordered-kmer-size 18 --num-hashes 128 --num-min-matches 5 --ordered-sketch-size 1000 --threshold 0.95 --num-threads 32
MHAP	<i>Yeast (PacBio)</i>	--store-full-id --num-threads 32
MHAP	<i>Yeast (ONT)</i>	--store-full-id --num-threads 32
MHAP	<i>E. coli (PacBio)</i>	--store-full-id --num-threads 32

# BLEND Parameter Definitions

Parameter	Definition
<code>-strobemers</code>	Use the BLEND-S mechanism when generating the list of k-mers of a seed
<code>-immediate</code>	Use the BLEND-I mechanism when generating the list of k-mers of a seed
<code>-H</code>	Use homopolymer-compressed k-mers
<code>-w INT</code>	Window size used when finding minimizers.
<code>-k INT</code>	k-mer size used when generating the list of k-mers of a seed
<code>-neighbors INT</code>	Number of k-mers included in the list of seeds. Combination of both <code>-k</code> ( $k$ ) and <code>-neighbors</code> ( $n$ ) determines the seed length. Seed length in BLEND-S is calculated as: $k \times n$ Seed length in BLEND-I is calculated as: $k + (n - 1)$
<code>-fixed-bits INT</code>	Bit length of hash values that BLEND generates for each seed. Setting it to $2 \times k$ is the default behavior.
<code>-t INT</code>	Number of CPU threads to use.
<code>-x STR</code>	Preset for setting the default parameters given the use case (STR)
<code>-x map-ont</code>	Preset for mapping ONT reads. It uses the following parameters: <code>-immediate -w 10 -k 9 -neighbors 7 -fixed-bits 30</code>
<code>-x map-pb</code>	Preset for mapping erroneous PacBio reads. It uses the following parameters: <code>-immediate -H -w 10 -k 13 -neighbors 7 -fixed-bits 32</code>
<code>-x map-hifi</code>	Preset for mapping accurate long (HiFi) reads. It uses the following parameters: <code>-strobemers -w 50 -k 19 -neighbors 5 -fixed-bits 38</code>
<code>-x sr</code>	Preset for mapping short reads. It uses the following parameters: <code>-immediate -w 11 -k 21 -neighbors 5 -fixed-bits 32</code>
<code>-x ava-ont</code>	Preset for overlapping ONT reads. It uses the following parameters: <code>-immediate -w 10 -k 15 -neighbors 5 -fixed-bits 30</code>
<code>-x ava-pb</code>	Preset for overlapping erroneous PacBio reads. It uses the following parameters: <code>-immediate -H -w 10 -k 19 -neighbors 5 -fixed-bits 38</code>
<code>-x ava-hifi</code>	Preset for overlapping accurate long (HiFi) reads. It uses the following parameters: <code>-strobemers -w 200 -k 25 -neighbors 7 -fixed-bits 50</code>

# Read Mapping Parameters

Tool	Dataset	Parameters
BLEND	<i>CHM13 (HiFi)</i>	-ax map-hifi -t 32 --secondary=no
BLEND	<i>HG002 (HiFi)</i>	-ax map-hifi -t 32 --secondary=no
BLEND	<i>D. ananassae (HiFi)</i>	-ax map-hifi -t 32 --secondary=no
BLEND	<i>E. coli (HiFi)</i>	-ax map-hifi -t 32 --secondary=no
BLEND	<i>CHM13 (ONT)</i>	-ax map-ont -t 32 --secondary=no
BLEND	<i>Yeast (PacBio)</i>	-ax map-pb -t 32 --secondary=no
BLEND	<i>Yeast (ONT)</i>	-ax map-ont -t 32 --secondary=no
BLEND	<i>Yeast (Illumina)</i>	-ax sr -t 32
BLEND	<i>E. coli (PacBio)</i>	-ax map-pb -t 32 --secondary=no
minimap2	<i>CHM13 (HiFi)</i>	-ax map-hifi -t 32 --secondary=no
minimap2	<i>HG002 (HiFi)</i>	-ax map-hifi -t 32 --secondary=no
minimap2	<i>D. ananassae (HiFi)</i>	-ax map-hifi -t 32 --secondary=no
minimap2	<i>E. coli (HiFi)</i>	-ax map-hifi -t 32 --secondary=no
minimap2	<i>CHM13 (ONT)</i>	-ax map-ont -t 32 --secondary=no
minimap2	<i>Yeast (PacBio)</i>	-ax map-pb -t 32 --secondary=no
minimap2	<i>Yeast (ONT)</i>	-ax map-ont -t 32 --secondary=no
minimap2	<i>Yeast (Illumina)</i>	-ax sr -t 32
minimap2	<i>E. coli (PacBio)</i>	-ax map-pb -t 32 --secondary=no
Winnowmap2	<i>CHM13 (HiFi)</i>	meryl count k=15 meryl print greater-than distinct=0.9998
Winnowmap2	<i>HG002 (HiFi)</i>	-ax map-pb -t 32 meryl count k=15 meryl print greater-than distinct=0.9998
Winnowmap2	<i>D. ananassae (HiFi)</i>	-ax map-pb -t 32 meryl count k=15 meryl print greater-than distinct=0.9998
Winnowmap2	<i>E. coli (HiFi)</i>	-ax map-pb -t 32 meryl count k=15 meryl print greater-than distinct=0.9998
Winnowmap2	<i>CHM13 (ONT)</i>	-ax map-pb -t 32 meryl count k=15 meryl print greater-than distinct=0.9998
Winnowmap2	<i>Yeast (PacBio)</i>	-ax map-ont -t 32 meryl count k=15 meryl print greater-than distinct=0.9998
Winnowmap2	<i>Yeast (ONT)</i>	-ax map-pb-clr -t 32 meryl count k=15 meryl print greater-than distinct=0.9998
Winnowmap2	<i>E. coli (PacBio)</i>	-ax map-ont -t 32 meryl count k=15 meryl print greater-than distinct=0.9998 -ax map-pb-clr -t 32
LRA	<i>CHM13 (HiFi)</i>	align -CCS -t 32 -p s
LRA	<i>HG002 (HiFi)</i>	align -CCS -t 32 -p s
LRA	<i>D. ananassae (HiFi)</i>	align -CCS -t 32 -p s
LRA	<i>E. coli (HiFi)</i>	align -CCS -t 32 -p s
LRA	<i>CHM13 (ONT)</i>	align -ONT -t 32 -p s
LRA	<i>Yeast (PacBio)</i>	align -CLR -t 32 -p s
LRA	<i>Yeast (ONT)</i>	align -ONT -t 32 -p s
LRA	<i>E. coli (PacBio)</i>	align -CLR -t 32 -p s
S-conLSH	<i>CHM13 (HiFi)</i>	-threads 32 -align 1
S-conLSH	<i>E. coli (HiFi)</i>	-threads 32 -align 1
S-conLSH	<i>CHM13 (ONT)</i>	-threads 32 -align 1
S-conLSH	<i>Yeast (PacBio)</i>	-threads 32 -align 1
S-conLSH	<i>Yeast (ONT)</i>	-threads 32 -align 1
S-conLSH	<i>E. coli (PacBio)</i>	-threads 32 -align 1
Strobealign	<i>Yeast (Illumina)</i>	-t 32

# Versions of each Tool

Tool	Version	GitHub or Conda Link to the Version
BLEND	1.0	<a href="https://github.com/CMU-SAFARI/BLEND">https://github.com/CMU-SAFARI/BLEND</a>
minimap2	2.24	<a href="https://github.com/lh3/minimap2/releases/tag/v2.24">https://github.com/lh3/minimap2/releases/tag/v2.24</a>
MHAP	2.1.3	<a href="https://anaconda.org/bioconda/mhap/2.1.3/download/noarch/mhap-2.1.3-hdfd78af_1.tar.bz2">https://anaconda.org/bioconda/mhap/2.1.3/download/noarch/mhap-2.1.3-hdfd78af_1.tar.bz2</a>
LRA	1.3.2	<a href="https://anaconda.org/bioconda/lra/1.3.2/download/linux-64/lra-1.3.2-ha140323_0.tar.bz2">https://anaconda.org/bioconda/lra/1.3.2/download/linux-64/lra-1.3.2-ha140323_0.tar.bz2</a>
Winnowmap2	2.03	<a href="https://anaconda.org/bioconda/Winnowmap/2.03/download/linux-64/Winnowmap2-2.03-h2e03b76_0.tar.bz2">https://anaconda.org/bioconda/Winnowmap/2.03/download/linux-64/Winnowmap2-2.03-h2e03b76_0.tar.bz2</a>
S-conLSH	2.0	<a href="https://github.com/anganachakraborty/S-conLSH-2.0/tree/292fbe0405f10b3ab63fc3a86cba2807597b582e">https://github.com/anganachakraborty/S-conLSH-2.0/tree/292fbe0405f10b3ab63fc3a86cba2807597b582e</a>
Strobealign	0.7.1	<a href="https://anaconda.org/bioconda/strobealign/0.7.1/download/linux-64/strobealign-0.7.1-hd03093a_1.tar.bz2">https://anaconda.org/bioconda/strobealign/0.7.1/download/linux-64/strobealign-0.7.1-hd03093a_1.tar.bz2</a>

# BLEND Backup Slide