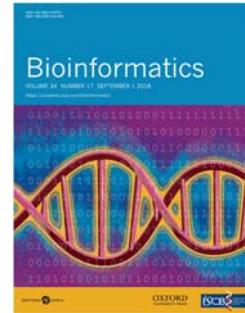




Binning directories: An efficient data structure to query k-mers in partitioned string sets

Temesgen H. Dadi, Enrico Siragusa, Vitor C. Piro, Andreas Andrusch, Enrico Seiler, Bernhard Y. Renard and Knut Reinert



Knut Reinert
Freie Universität Berlin

International Workshop on PDBs and Large Scale Search

Motivation



- The indexing step in read mapping:
 - Computationally very expensive
 - But amortizes in time due to the static nature of reference genome e.g., hg38 since 2013
- Read mapping in metagenomics
 - A large set of reference genomes
 - Frequently changing (Not static)
 - Undermines the purpose of indexing

Motivation



- Very large set of references



To build FM-Index

- ~100 GB RAM
- ~10-27 hours
- **Updating** takes the same amount of time because the index structures (CSA) are **static**

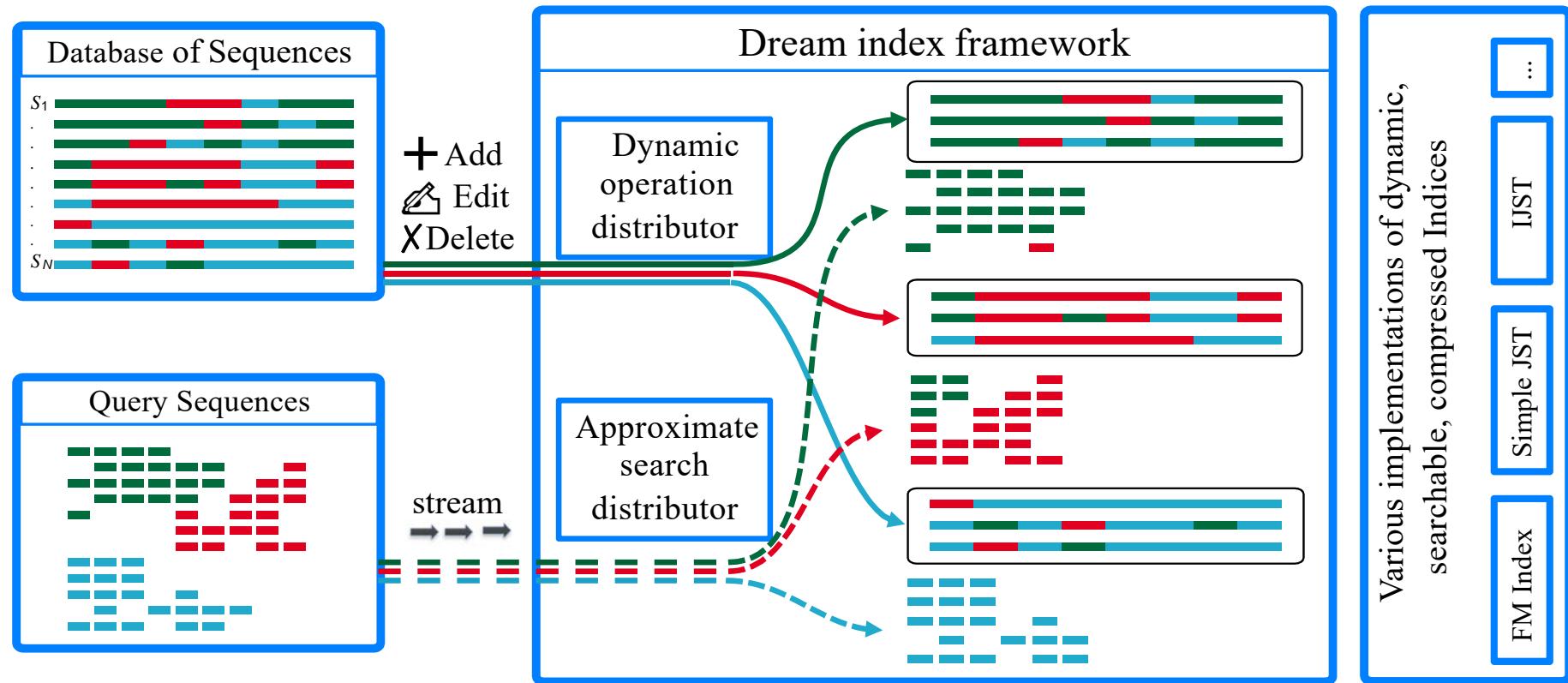
30 GB is “*not that big*”

- Partition the reference databases into smaller parts (bins)

DREAM Index framework



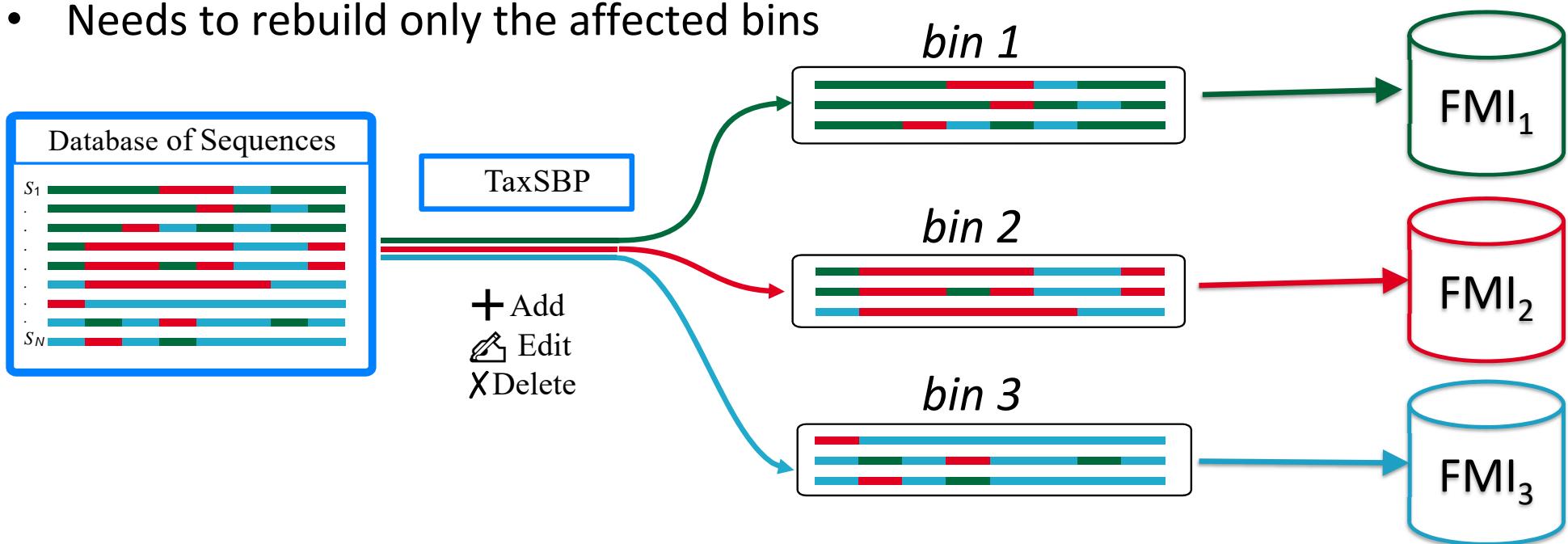
DREAM = Dynamic seaRchablE pArallel coMpressed index



Bin with TaxSBP and then index



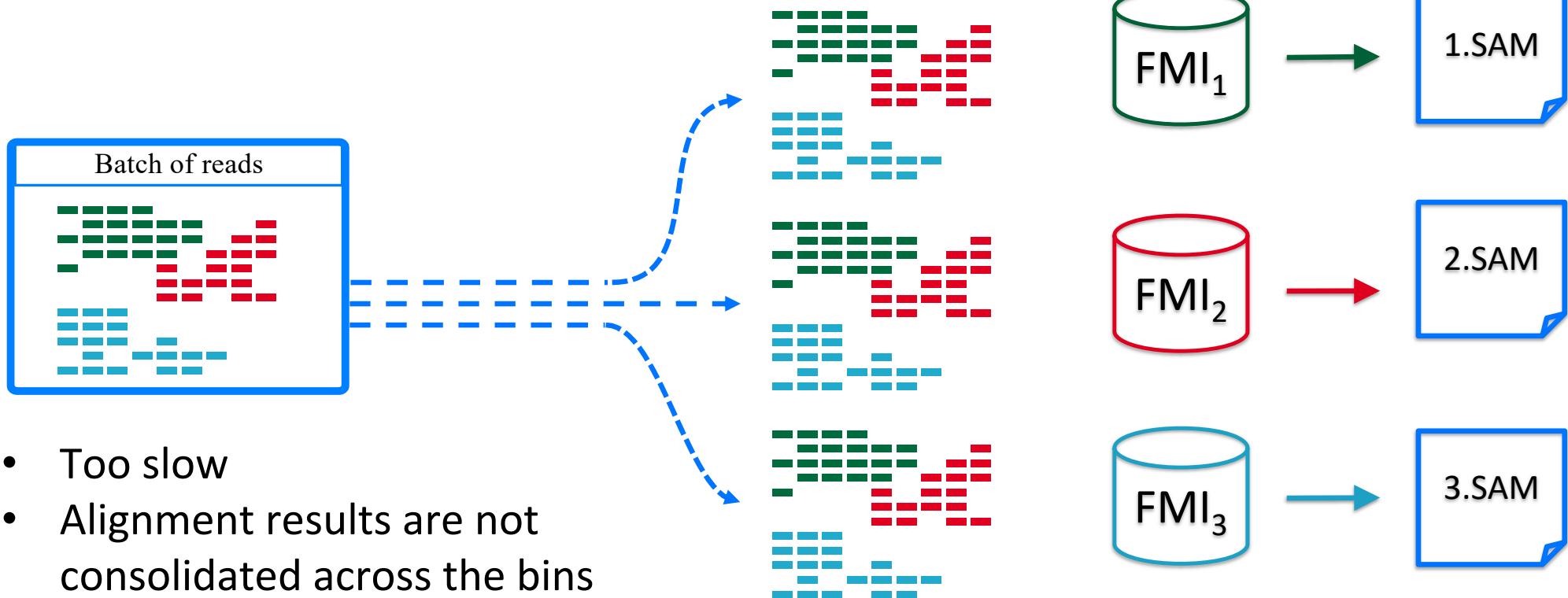
- Taxonomy based binning by V. Piro
- Hierarchically structured bin packing algorithm [*Codenotti et al.*]
- NCBI Taxonomy database
- Needs to rebuild only the affected bins



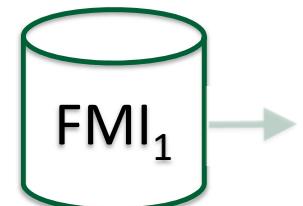
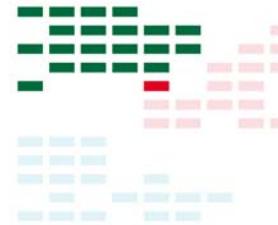
Trivial distributed mapping



Mapping all reads against each sub-index



Pre-filter reads before mapping

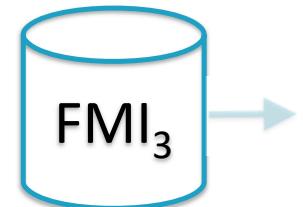
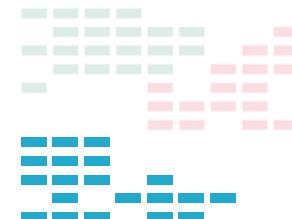


Lightweight

Problem:

We often need k-mers with $k > 20$

The data set contains often > 100 billion unique k-mers



**Returns binning
bitvector
indicating bin
membership**

Bloom Filter



1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25	26	27	28	29	30	31	32	33	34	35	36	37	38	39	40	...	n
0	0	0	0	1	1	0	1	1	0	0	0	0	0	1	1	0	0	0	1	0	0	0	0	0	0	1	1	1	0	0	0	0	0	1	1	1	0	...	0		

- A bit-vector of size n
- h different hash functions - $\{H_1, H_2, \dots, H_h\}$
- $H_i : K \rightarrow L_i$ where $0 \leq L_i < n$
- To add a ***k-mer*** we simply set h bit positions
- During lookup, we expect all of the h bit positions to be set.
- Collisions are allowed in the expense of false positive answers.

$$P_{fp} = \left(1 - \left(1 - \frac{1}{n} \right)^{h*m} \right)^h$$

Where:

m = the number of elements **Added**

Bloom Filter



Why multiple hash functions?

Size of bit-vector (n) = 1 Billion

m (Billion)	Number of hash functions			
	1	2	3	4
1.0000000	0.632120549	0.747645060	0.857951631	0.928725749
0.5000000	0.393469332	0.399576389	0.468861712	0.558973136
0.2500000	0.221199212	0.154818116	0.146891590	0.159661290
0.1250000	0.117503098	0.048929091	0.030579353	0.023968649
0.0625000	0.060586936	0.013806977	0.004997657	0.002394056
0.0312500	0.030766765	0.003670777	0.000716668	0.000190633
0.0156250	0.015503563	0.000946594	0.000096030	0.000013475
0.0078125	0.007782062	0.000240360	0.000012431	0.000000896
0.0039063	0.003898630	0.000060560	0.000001581	0.000000058

Probability of false positive answer (FP rate)

Smaller FP rate &
Add more elements.

General principle:

$$h \leq \frac{n}{m}$$

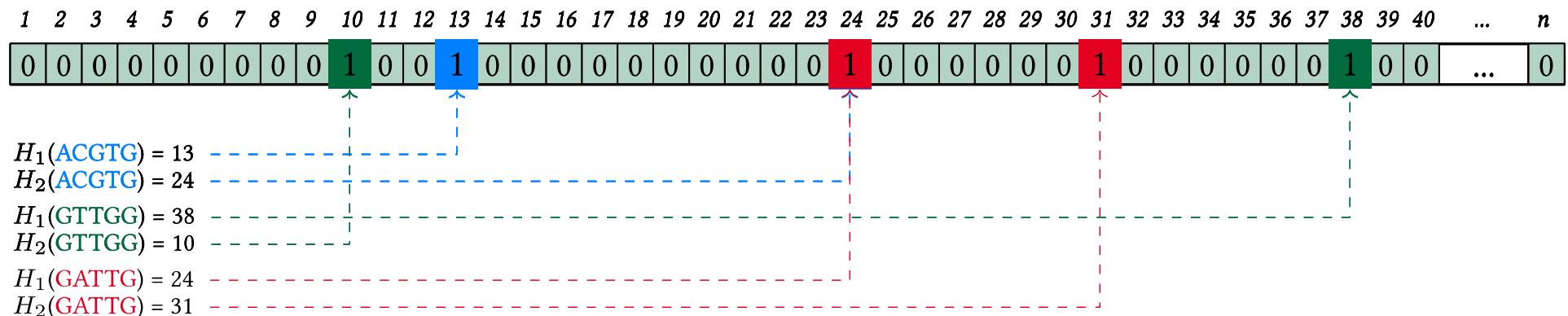
Bloom Filter



E.g. Add k-mers

ACGTG GTTGG GATTG

2 hash functions: $\{H_1, H_2\}$



Bloom Filter for each bins?



	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25	26	27	28	29	30	31	32	33	34	35	36	37	38	39	40	...	n
BF_b	0	0	0	0	1	1	0	1	1	0	0	0	0	1	1	0	0	0	1	0	0	0	0	0	0	1	1	1	0	0	0	0	0	1	1	1	0	...	0			

...

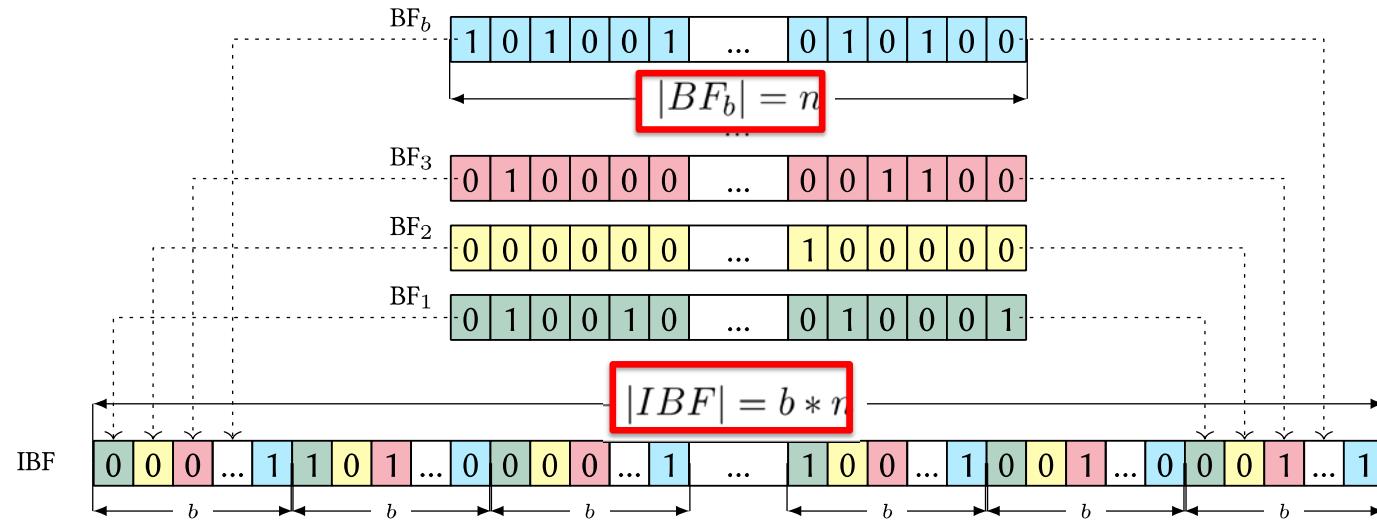
BF_3	0	0	0	0	1	1	0	1	1	0	0	0	0	1	1	0	0	0	1	0	0	0	0	0	0	1	1	1	0	...	0
--------	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	-----	---

BF_2	0	0	0	0	1	1	0	1	1	0	0	0	0	1	1	0	0	0	1	0	0	0	0	0	0	1	1	1	0	...	0
--------	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	-----	---

BF_1	0	0	0	0	1	1	0	1	1	0	0	0	0	1	1	0	0	0	1	0	0	0	0	0	0	1	1	1	0	...	0
--------	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	-----	---

- We need to check every k-mer of a read against each Bloom filter and construct binning vector

Interleaved Bloom Filter (IBF)



To add a k-mer (\mathbf{K}) from bin j

We set position

$$\mathbf{b} \times H_i(\mathbf{K}) + j$$

To check in which bins a k-mer is present we use the bit blocks

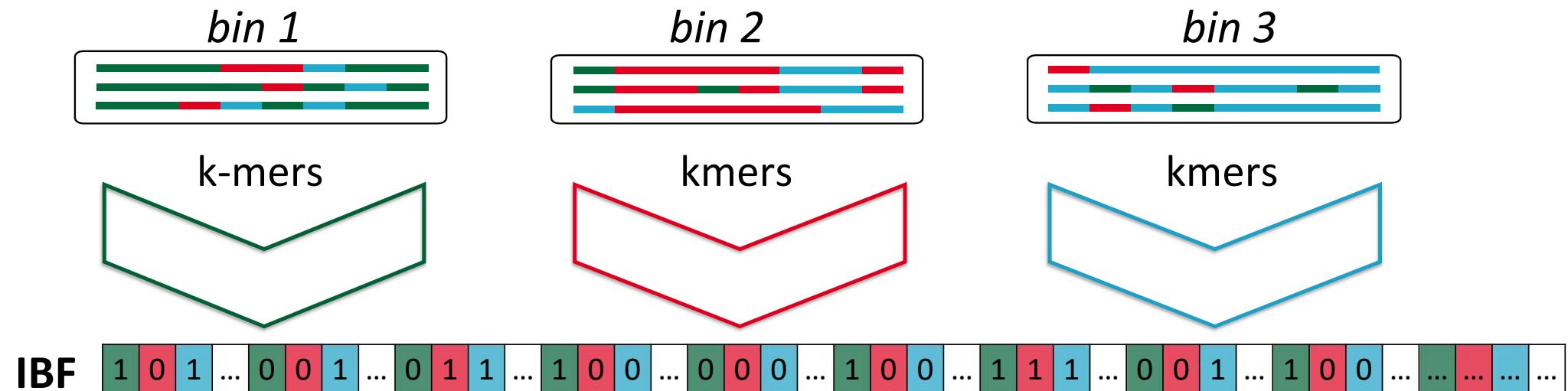
$$[H_i(\text{kmer}) \times \mathbf{b}, H_i(\text{kmer}) \times (\mathbf{b} + 1))$$

Interleaved Bloom Filter (IBF)



A light weight filtering data structure

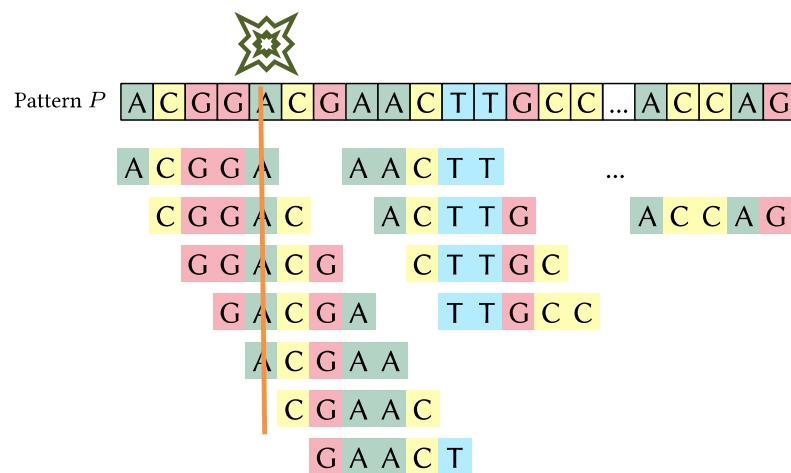
- Based on bloom filters
- Contains the kmer information of all bins
- Together with q-gram lemma → very fast pre-filter
- **37 times faster than the trivial method (1024 bins)**



K-mer Counting Lemma



For a given k and number of errors e , there are $k_p = |P| - k + 1$ many k-mers in pattern P and an approximate occurrence of P in T has to share at least $t = (k_p - k \cdot e)$ k-mers



Total number of k-mers

An error can destroy at most k k-mers
 e errors can destroy $k \cdot e$ k-mers

Limitation when e is bigger
 $t = |P| - (e+1) \cdot k + 1$

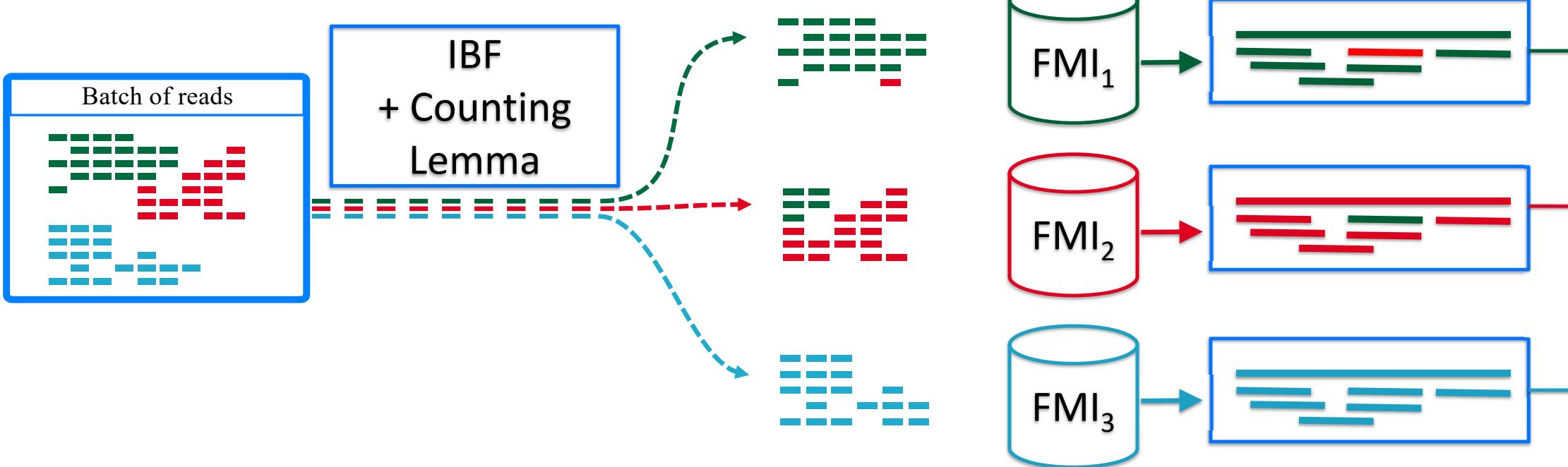
The Yara read mapper

Freie Universität Berlin



- State of the art read mapper of the SeqAn library
- Both **all mapping and strata mapping**
- Yara is a **fully sensitive** read mapper
- Fast (2X Bowtie2, 3X BWA-MEM)

DREAM-Yara (*Mapping*)



Benchmark



Infrastructure

- Intel(R) Xeon(R) CPU E5-2650 v3 2.30GHz processor
- 130GB of memory
- 8 threads whenever the tool allows

Dataset (2017-09-26)

- 15,250 sequences
- 2,991 species
- 31.34 GB

Update (2017-12-19)

- 155 new sequences of E. coli
- 42/1024 bins affected
- 230 MB

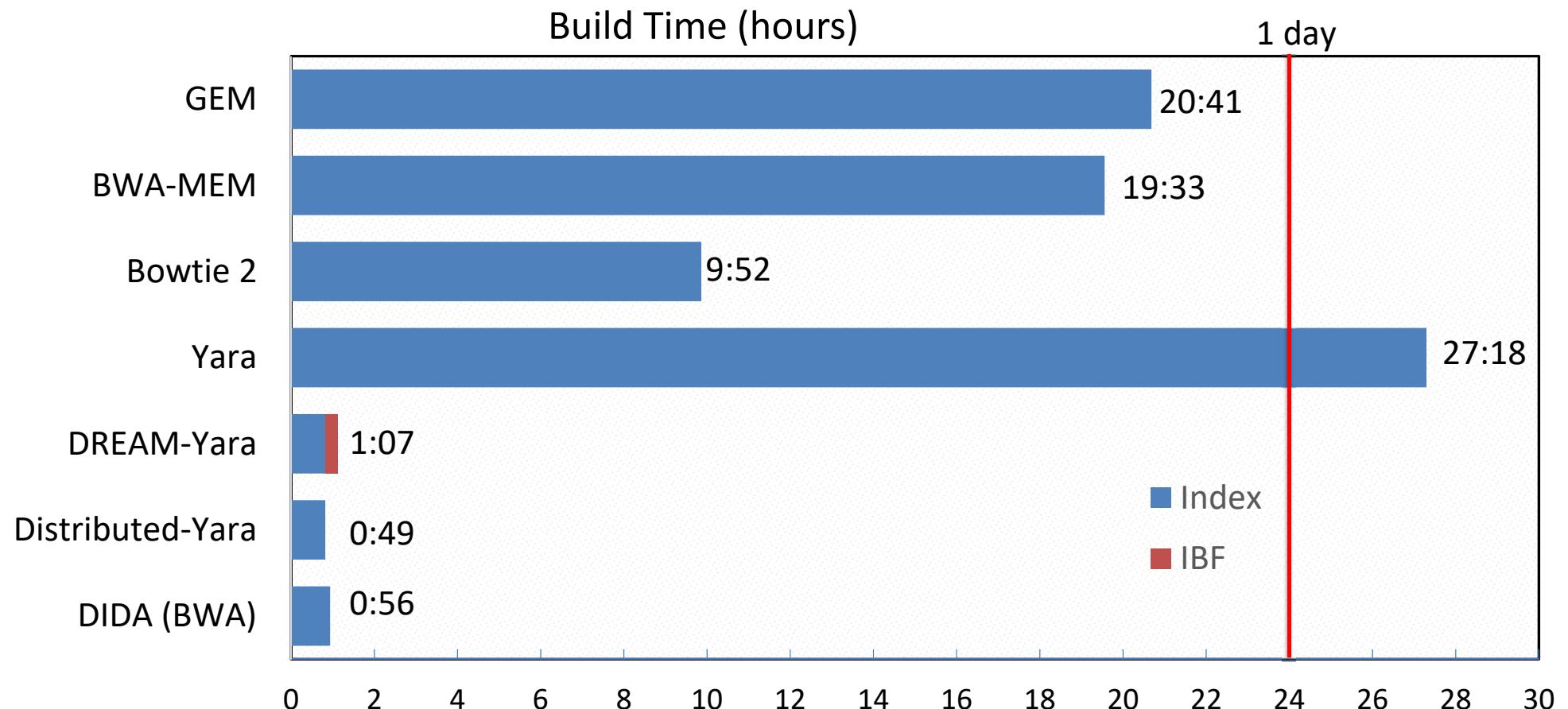
BWA-MEM

Bowtie2

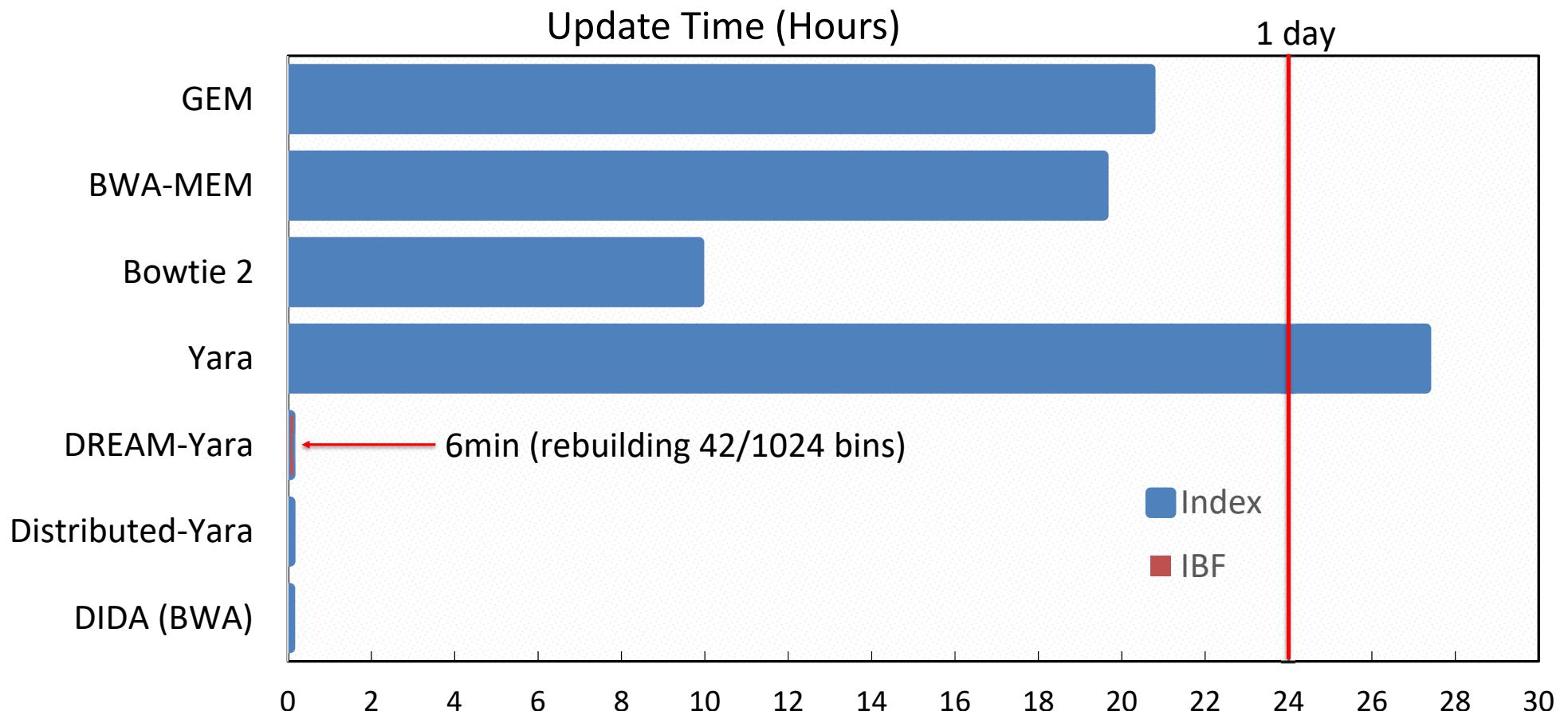
GEM

DIDA (BWA)

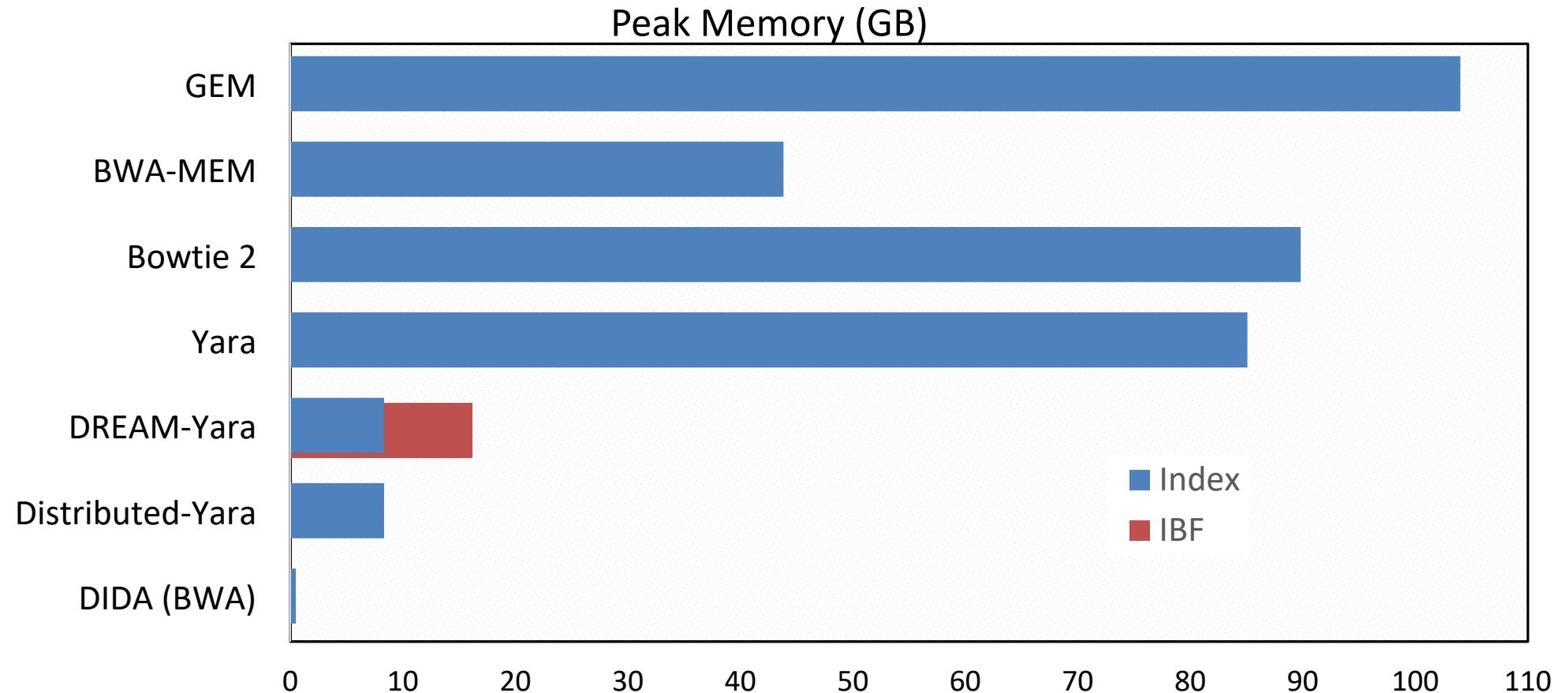
Indexing – build time



Indexing – update time



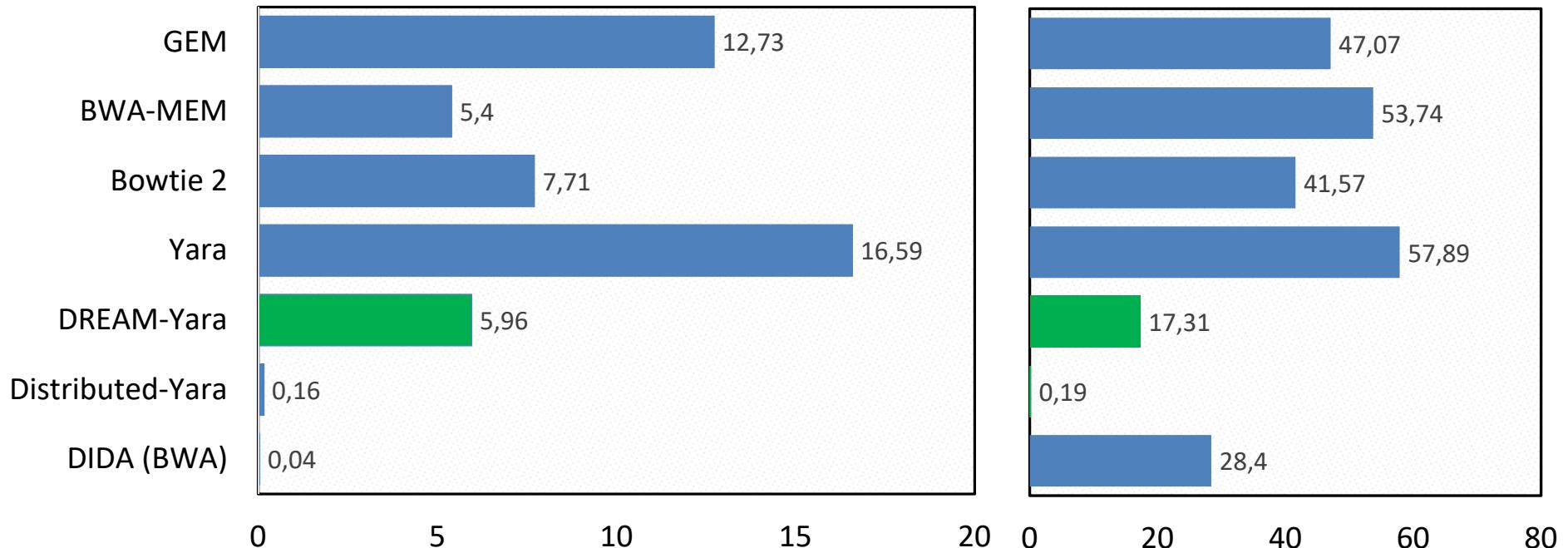
Indexing – memory footprint



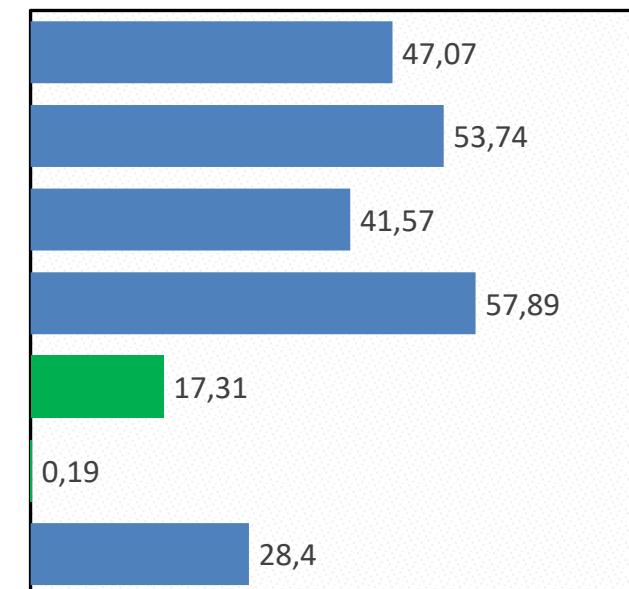
Mapping speed and memory



Throughput (Gbp/h)



Memory (GB)



Rabema benchmarks (Sensitivity)



Read Set: SRR6504858 (human gut metagenome sample)

	Co-opt. locations			
	[% Normalized]	[% Absolute]		
Distributed-Yara	100.0 100.0 100.0 100.0 100.0 100.0 100.0	100.0 100.0 100.0 100.0 100.0 100.0 100.0	100.0 99.8 95.4 74.2 67.6 79.3	e less
DREAM-Yara	100.0 100.0 100.0 100.0 100.0 100.0 100.0	100.0 100.0 100.0 100.0 100.0 100.0 100.0	100.0 99.8 95.4 74.2 67.6 79.3	
Yara	100.0 100.0 100.0 100.0 100.0 100.0 100.0	100.0 100.0 100.0 100.0 100.0 100.0 100.0	100.0 99.8 95.4 74.2 67.6 79.3	
GEM	99.8 99.5 99.7 94.8	96.4 74.2 67.6 79.3		Absolute: All locations equally matter
Bowtie 2	99.8 99.7 99.5 97.3	81.6 67.6 63.4 81.3		
BWA-MEM	99.3 98.2 93.8 90.1	85.2 81.5 72.7 84.5		
DIDA-BWA	99.5 98.2 93.8 -	85.2 83.6 92.1 93.6 81.5 72.7 84.5		

Summary



- DREAM framework for distributed read mapping which can support **fast updates** of the sub-indices.
- A fast and light weight **pre-filter using IBF** and q-gram lemma ($q=19$)
- DREAM-Yara, implementation of an **in-memory distributed version** of Yara.
 - **Competitive** mapping speed
(faster than BWA and slightly slower the Bowtie2)
 - The pre-filter introduced a $\sim 37X$ speedup.
 - It can conduct a typical batch index update in **6 minutes**
(Bowtie2 ≈ 10 hr; BWA & GEM ≈ 20 hr; Yara ≈ 27 hr)

Thank you!

Co-Authors

Enrico Siragusa,
Vitor C. Piro,
Andreas Andrusch,
Enrico Seiler,
Bernhard Y. Renard and
Knut Reinert

Freie Universität Berlin



**Algorithmic
Bioinformatics Group**

ROBERT KOCH INSTITUT



Outline



- **Introduction**
 - Motivation
 - The DREAM index framework
- **Methods**
 - Binning sequences
 - Interleaved Bloom Filter
 - The Yara read mapper
 - DREAM-Yara
- **Evaluation**
 - Index creation and updating benchmark
 - Rabema benchmarks (Sensitivity)
 - Read mapping speed

The Yara read mapper



- State of the art read mapper of the SeqAn Library
- Yara is an **exact** read mapper (**All mapper**)
- Fast (2X Bowtie2, 3X BWA-MEM)
- Implements strata based mapping

E.g. Mapping a read (*100 bp long*) up to 5 errors (*edit distance*)

Best map has 2 errors

0 strata => 2e

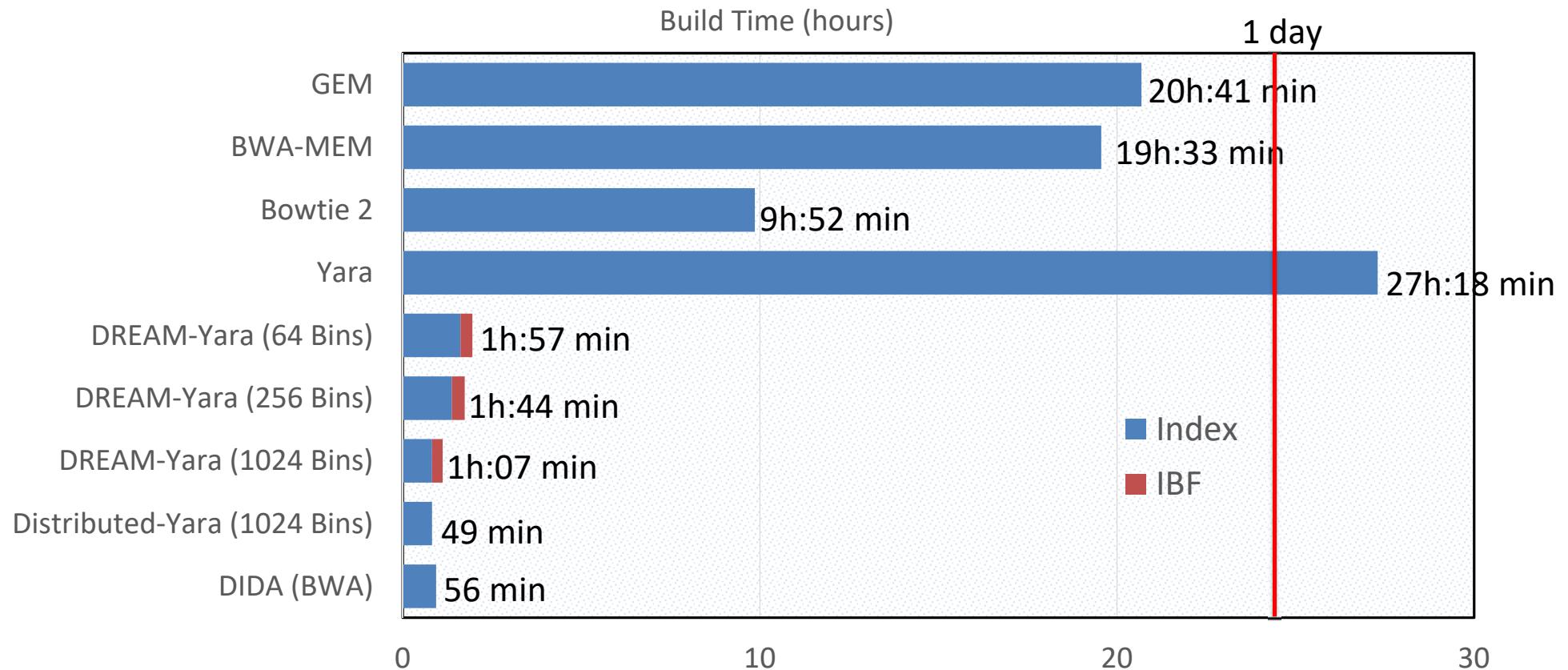
1 strata => 2e, 3e

2 strata => 2e, 3e, 4e

3 strata => 2e, 3e, 4e, 5e

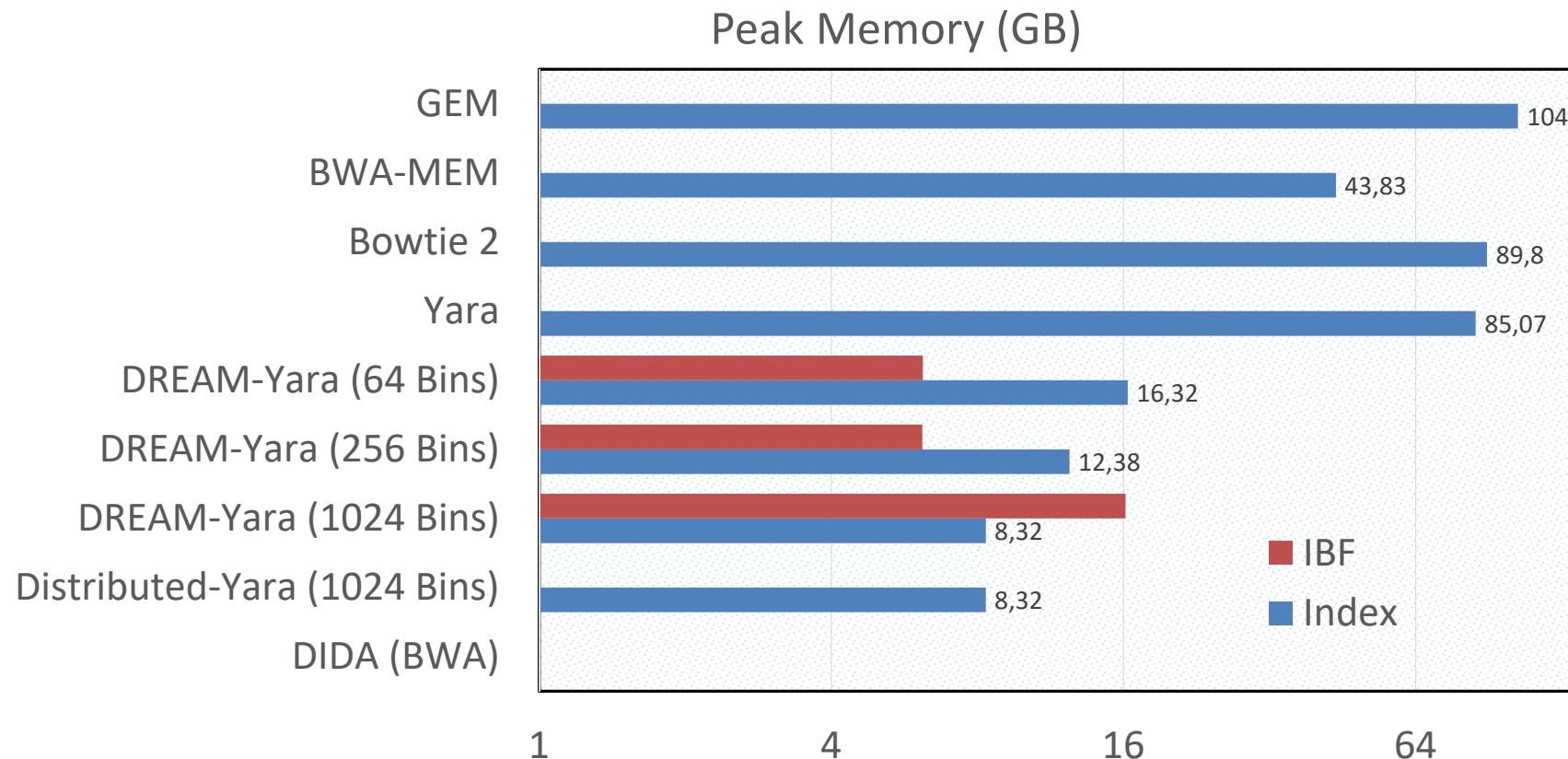


Benchmark (Index creation)



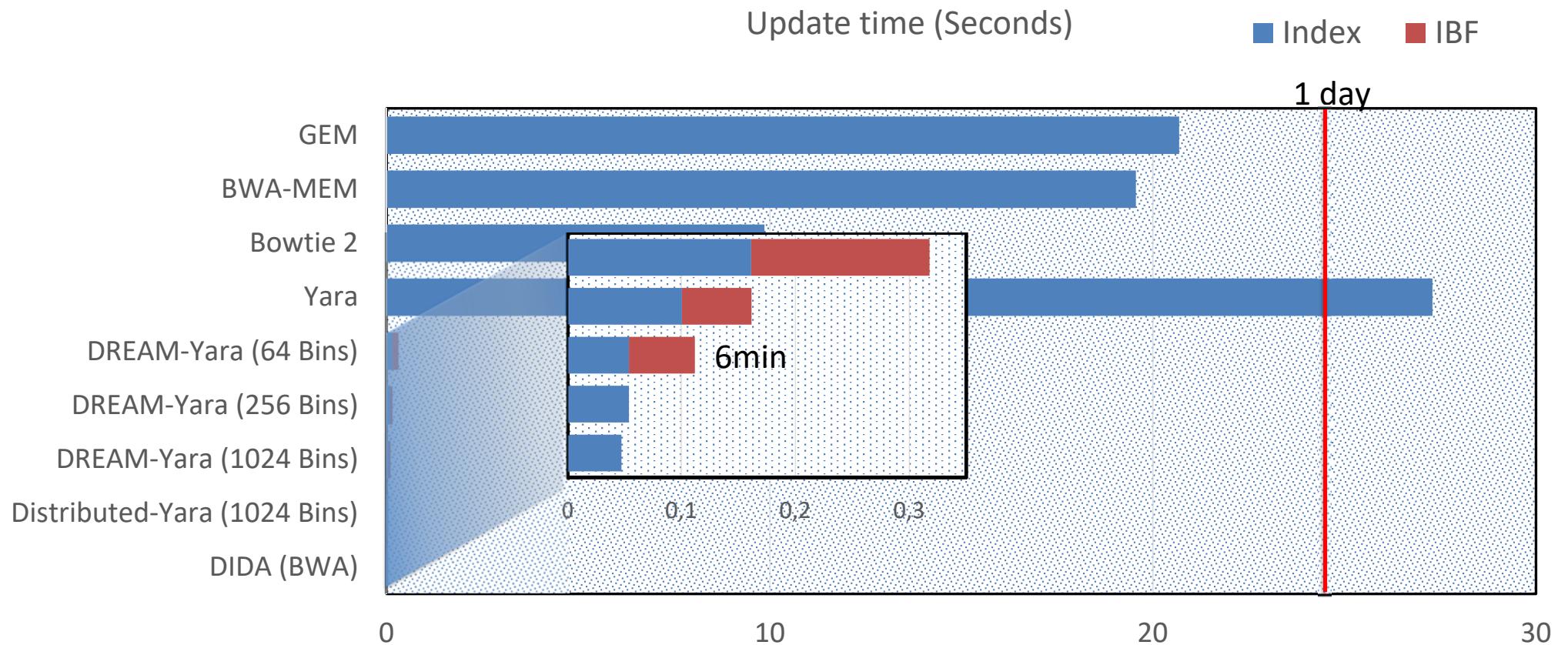


Benchmark (Indexing - memory)





Benchmark (Index updating)



IBF + Counting Lemma



Interleaved Bloom Filter IBF

0	0	0	0	0	0	...	1	1	0	1	0	1	...	0	0	0	0	0	0	...	1	0	0	0	...
---	---	---	---	---	---	-----	---	---	---	---	---	---	-----	---	---	---	---	---	---	-----	---	---	---	---	-----

Pattern p

A C G G A C G A ... A C C A G

A C G G A

$IBF(k_1)$ 0 0 0 0 1 ... 1

C G G A C

$IBF(k_2)$ 1 0 1 0 1 ... 1

G G A C G

$IBF(k_3)$ 1 1 1 1 0 ... 0

...

A C C A G

$IBF(k_n)$ 1 0 0 0 1 ... 0

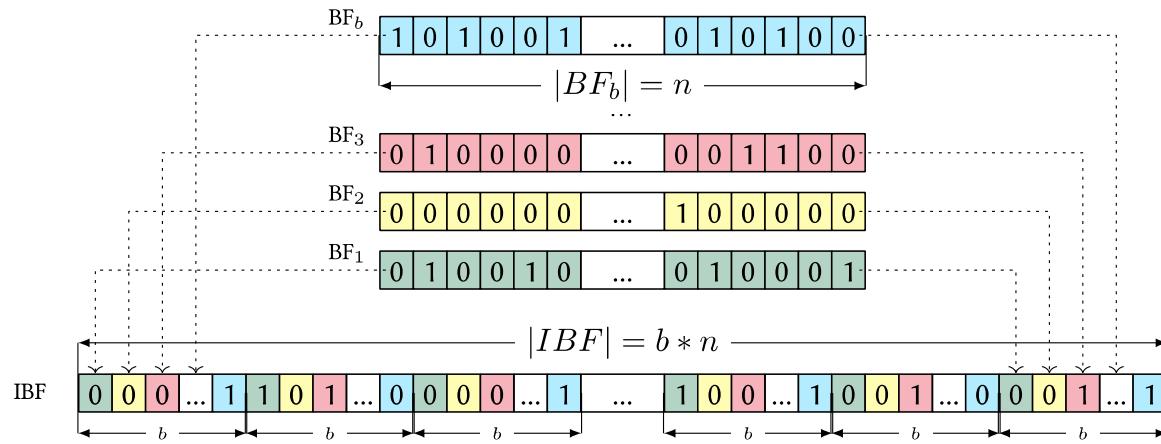
$Count(P)$ 5 2 4 0 3 ... 3

potential bins for pattern p (threshold = 4)

✓ ✓

Approximate
search
distributor

IBF Caveats



Works best when

- Bins are equal in text size