Data Systems at Scale: Scaling Up by Scaling Down and Out

> Prashant Pandey VMware Research

Scalability challenge in a tweet



Mick W@tson @BioMickWatson-

Bioinformatics over the years:

1990s: doing a BLAST search 2000s: analysing 30 microarrays 2010s: nalysing 6Tb of NGS 2020s: creating a cloud the size of Netflix to reanalyse the whole of SRA for one figure

12:57 AM · 2/12/21 · Twitter Web App

117 Retweets 9 Quote Tweets 697 Likes









Michael Schatz @mike_schatz · 2h Replying to @BioMickWatson This is basically my life right now \mathbb{C}^{8} O_1 **1**↓ 企

 $\uparrow\uparrow$

Professor Bioinformatics and **Computational Biology** The University of Edinburgh

...

Associate Professor **Computational Biology** Johns Hopkins University

Sequence Read Archive (SRA) growth

SRA contains a lot of *diversity information*



Q: What if I find e.g., a new disease-related gene, and want to see if it appeared in other experiments?

Scalability is the bottleneck for data science

SRA contains a lot of *diversity information*



This renders what is otherwise an immensely valuable public resource *largely inert*

Scalability is a ubiquitous challenge

Cyber monitoring Internet of Things (IoT)

Financial tech

Social networks

AstroPhysics

Chemistry

Environmental science

• People generate 2.5 quintillion bytes of data each day. (IBM, 2016)

- More than 150 zettabytes (150 trillion gigabytes) of data will need analysis by 2025. (Forbes, 2019)
- 90 percent of the world's data was created between 2015 and 2016 alone. (IBM, 2016)
 https://learn.g2.com/big-data-statistics

2	4. 88% of data is ignored by companies.	https://leftronic.com/big-data-statistics/	
(orrester Research)	·	
A	A widely-quoted figure from a 2012 paper from Forrester Research says that, on average, companies analyze		
C	only 12% of the available data Reasons for this include a lack of analytics tools, repressive data silos, and the		
С	ifficulty in knowing which information is valuable	and which is worth leaving.	

Shrink it

Goal: make data smaller to fit in RAM

Techniques:

- Compact & succinct data structures
- Filters, e.g., Bloom, quotient, etc.



Goal: make data smaller to fit in RAM

Techniques:

- Compact & succinct data structures
- Filters, e.g.,
 Bloom,
 quotient, etc.

Organize it

Goal: organize data in a disk-friendly way

Techniques:

- B-tree
- B^{ε} -tree
- LSM-tree

Shrink it

Goal: make data smaller to fit in RAM

Techniques:

- Compact & succinct data structures
- Filters, e.g.,
 Bloom,
 quotient, etc.

Organize it

Goal: organize data in a disk-friendly way

Techniques:

- B-tree
- B^{ε} -tree
- LSM-tree

Distribute it

Goal: partition and distribute data on multiple nodes

Techniques:

- Distributed hash table
- Distributed key-value store

Research output Shrink Organize Distribute **Buffered CMS Quotient Filter Order Min Data structures** ESA '18, SIGMOD '17, Hash **Scalable MG** & Algorithms ISMB '19 SIGMOD '21 TODS '21 LERTs **B**etrFS file system SIGMOD '20 FAST '15, TOS 15, **Systems** FAST '16, TOS 16, Terrace SPAA '19, TOPC '21 SIGMOD '21 **LSM-Mantis** Squeakr, deBGR, Mantis, **VariantStore Bioinformatics '22 Rainbowfish**, MST-Mantis **Applications** Genome Biology '21 ISMB '17, WABI '17, **BIOINFORMATICS '17,** Computational **Distributed** *k*-mer **RECOMB '18, Cell Systems** biology counting '18, RECOMB '19, IPDPS '21 JCB '20

In this talk:



In this talk:





Filters

Time to change your filter

Dictionary data structure

A dictionary maintains a set S from universe U.





A dictionary supports membership queries on *S*.

Filter data structure

A filter is an *approximate* dictionary.



A filter supports *approximate* membership queries on *S*.



False-positive rate enables filters to be compact

space $\geq n \log(1/\epsilon)$

space
$$= \Omega(n \log |U|)$$





Filter

Dictionary

False-positive rate enables filters to be compact



	Optimal
Space (bits)	$pprox n \ \log(1/\epsilon) + \Omega(n)$
CPU cost	O(1)
Data locality	O(1) probes

Classic filter: The Bloom filter [Bloom '70]

Bloom filter: a bit array + *k* hash functions (here *k*=2)



Classic filter: The Bloom filter [Bloom '70]

Bloom filter: a bit array + *k* hash functions (here *k*=2)



Classic filter: The Bloom filter [Bloom '70]

Bloom filter: a bit array + *k* hash functions (here *k*=2)



Bloom filters are ubiquitous (> 4300 citations)



Bloom filters have suboptimal performance

	Bloom filter	Optimal
Space (bits)	$pprox 1.44 \; n \log(1/\epsilon)$	$lpha n \ \log(1/\epsilon) + \Omega(n)$
CPU cost	$\Omega(1/\epsilon)$	O(1)
Data locality	$\Omega(1/\epsilon)$ probes	O(1) probes

Applications often work around Bloom filter limitations

Limitations	Workarounds
No deletes	Rebuild
No resizes	Guess N, and rebuild if wrong
No filter merging or enumeration	???
No values associated with keys	Combine with another data structure

Bloom filter limitations increase system complexity, waste space, and slow down application performance

Quotienting is an alternative to Bloom filters [Knuth. Searching and Sorting Vol. 3, '97]

- Store fingerprints compactly in a hash table.
 - Take a fingerprint h(x) for each element x.



• Only source of false positives:

- Two distinct elements x and y, where h(x) = h(y)
- If x is stored and y isn't, query(y) gives a false positives

$$\Pr[x \text{ and } y \text{ collide}] = \frac{1}{2^p}$$









- b(x) = location in the hash table
 t(x) = tag stored in the hash table
- Collisions in the hash table?
- Linear probing
- Robin Hood hashing



(17] Resolving collisions in the QF [Fandey et al. Slowod

Implementation:

2 meta-bits per slot.

$$h(x) \longrightarrow h_{\theta}(x) \parallel h_{I}(x)$$



runends



(17] Resolving collisions in the QF [Fandey et al. Slowod

Implementation:

2 meta-bits per slot.

$$h(x) \longrightarrow h_{\theta}(x) \parallel h_{I}(x)$$



runends



Quotienting enables many features in the QF

- Good cache locality
- Efficient scaling out-of-RAM
- Deletions
- Enumerability/Mergeability
- Resizing
- Maintains count estimates or associate values
- Uses variable-sized encoding for counts [Counting quotient filter]
 - Asymptotically optimal space: $O(\sum |C(x)|)$



Quotient filters use less space than Bloom filters for all practical configurations

	Quotient filter	Bloom filter	Optimal
Space (bits)	$pprox n \ \log(1/\epsilon) + 2.125 n$	$pprox 1.44 \; n \log(1/\epsilon)$	$pprox n \ \log(1/\epsilon) + \Omega(n)$
CPU cost	O(1) expected	$\Omega(1/\epsilon)$	O(1)
Data locality	1 probe + scan	$\Omega(1/\epsilon)$ probes	O(1) probes

The quotient filter has theoretical advantages over the Bloom filter

Quotient filters use less space than Bloom filters for all practical configurations



Bloom filter: ~1.44 $log(1/\epsilon)$ bits/element. Quotient filter: ~2.125 + $log(1/\epsilon)$ bits/element.

Quotient filters perform better (or similar) to other non-counting filters



- Insert performance is similar to the state-of-the-art non-counting filters
- Query performance is significantly fast at low load-factors and slightly slower at higher load-factors

Summary of filters

- Bloom filters ^[Bloom '70]
- Quotient filters [Pagh et al. '05, Dillinger et al. '09, Bender et al. '12, Einziger et al. '15, Pandey et al. '17]
- Cuckoo/Morton filters [Fan et al. '14, Breslow & Jayasena '18]



- Others
 - Mostly based on perfect hashing and/or linear algebra
 - Mostly static
 - e.g., Xor filters ^[Graf & Lemire '20]

Current filter performance

Performance suffers due to high-overhead of collision resolution



Question: do you see a problem here??

Current filters have a problem..

Performance suffers due to high-overhead of collision resolution



Applications must choose between space and speed.

Current filters have a problem..

Performance suffers due to high-overhead of collision resolution



Update intensive applications maintain filters close to full.

Why quotient filters slow down

Quotient filters use Robin-Hood hashing (a variant of linear probing)

QFs use 2 bits/slot to keep track of runs.

To insert item *x*:

- 1. Find its run.
- 2. Shift other items down by 1 slot.

3. Store f(x).



As the QF fills, inserts have to do more shifting.









Note: $h_0(x)$ and $h_1(x)$ need to be dependent to support kicking.



As the CF fills, inserts have to do more kicking.

Note: $h_0(x)$ and $h_1(x)$ need to be dependent to support kicking.

Cuckoo filter performance [Fan et al. '14]

	Optimal	Cuckoo filter
Space (bits)	$lpha pprox n \ \log(1/\epsilon) + \Omega(n)$	$pprox n \ \log(1/\epsilon) + 3n$
CPU cost	O(1)	up to 500
Data locality	O(1) probes	random probes

 $s = \omega(\log \log n)$ slots/block (e.g., s=64)



Each block is a small quotient filter with false-positive rate $\varepsilon/2$ and capacity *s*.



 $s = \omega(\log \log n)$ slots/block (e.g., s=64)



Each block is a small quotient filter with false-positive rate $\varepsilon/2$ and capacity *s*.



 $s = \omega(\log \log n)$ slots/block (e.g., s=64)



To insert item *x*:

- 1. Compute $h_0(x)$ and $h_1(x)$.
- 2. Insert f(x) into emptier block.
- 3. Kick an item if needed.



Kick an item if needed.

2.



Kick an item if needed.

2.

No kicking $\Rightarrow h_0(x)$ and $h_1(x)$ can be independent for insert-only workload.







A vectorizable mini quotient filter

Each block has *b* logical buckets.

Fingerprints of each bucket are stored together.

We keep a bit vector of bucket boundaries.



A vectorizable mini quotient filter

Each block has *b* logical buckets.

Fingerprints of each bucket are stored together

Operations take constant time in a vector model of computation for vectors of size ω(log log n) ^[Bellloch '90]. Example, using AVX-512 instructions.



Vector quotient filter (VQF) performance

	Optimal	VQF
Space (bits)	$pprox n \; \log(1/\epsilon) + \Omega(n)$	$pprox n ~ \log(1/\epsilon) + 2.91 n$
CPU cost	O(1)	O(1)
Data locality	O(1) probes	2 probes

Evaluation: insertion



The vector quotient filter offers high performance at all load factors.

Evaluation: lookups



Evaluation: concurrency



Quotient filter's impact in computer science



Theoretically well-founded data structures can have a *big impact* on multiple subfields across *academia and industry*

Quotient filter's impact in computer science



Theoretically well-founded data structures can have a *big impact* on multiple subfields across *academia and industry*