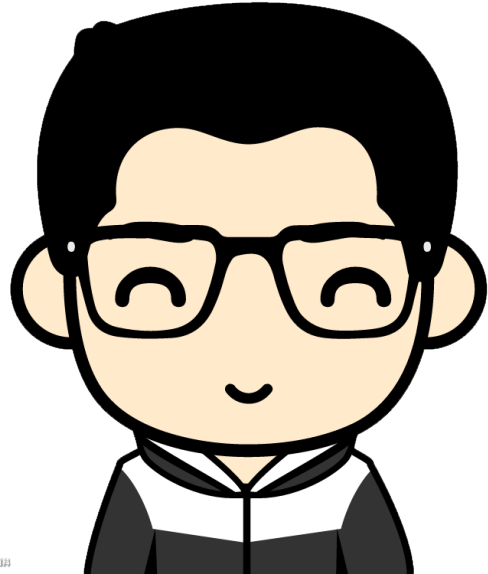


Why databases cry at night

Michael Yarichuk



Magic?



Nope. Not magic!
It's only an abstraction...

The Law of Leaky Abstractions

"All non-trivial abstractions, to some degree, are leaky."

- Joel Spolsky

The issue types

- Storage & relevant algorithms
- Indexing & Queries
- Network

1-a) Storage

It just works, no?

Here is a riddle... hint: storage!

RavenDB server-wide backup failed

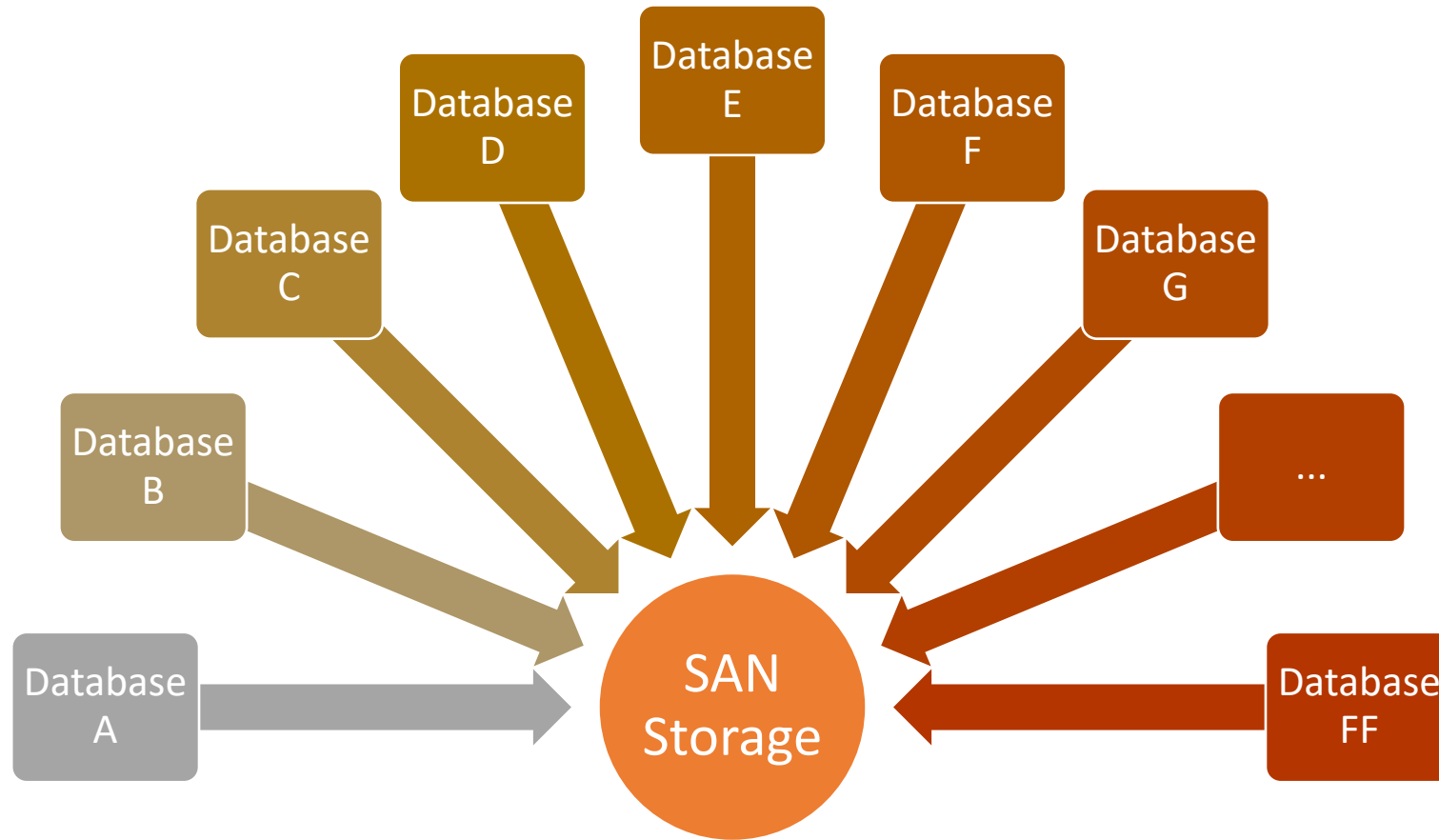
- The instance had multiple databases in single instance
- Plenty of memory and cores, resource usage is small
- Nothing else was running on the machine EXCEPT RavenDB
- Scheduled backup tasks fail soon after they started



The backup tasks *started* at the
same time!



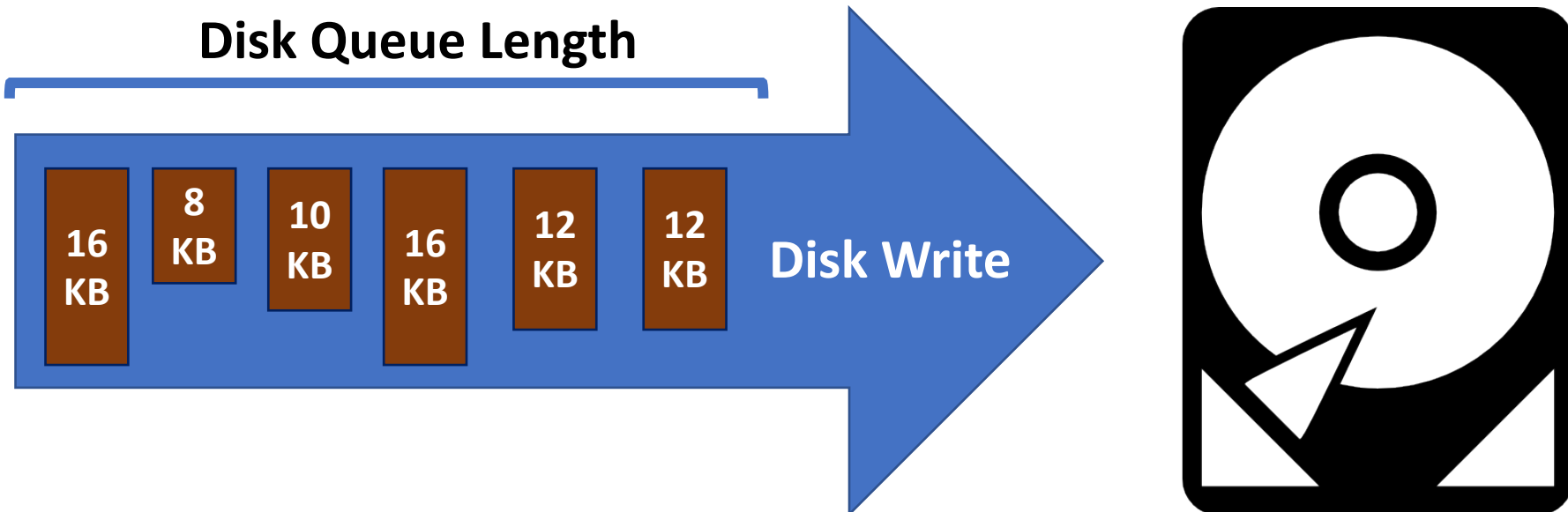
Also, there were gazillion of databases



Investigation: disk queue length

Bottleneck indicator

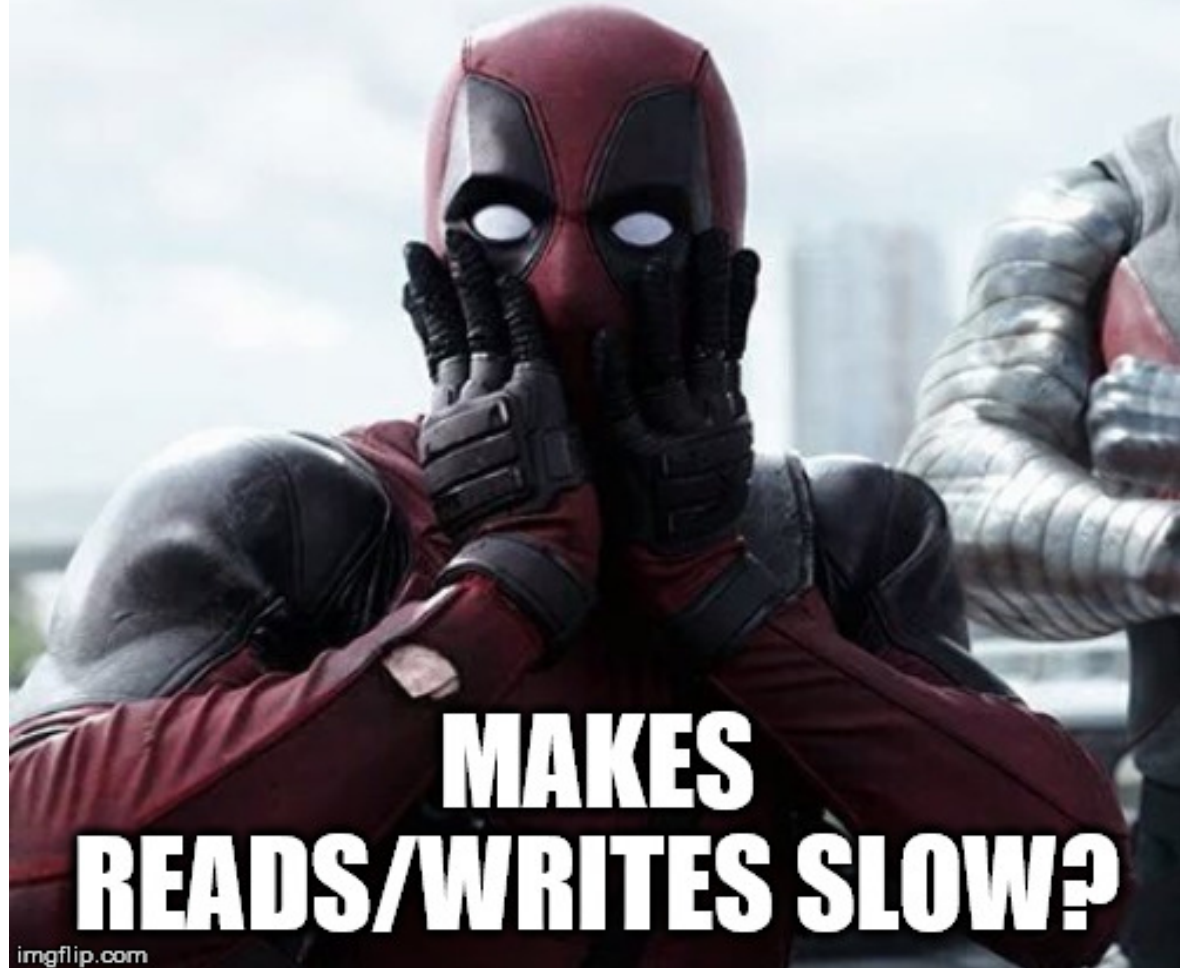
Active Time (%)	Available Space...	Total Space (MB)	Disk Queue Length
100.00	57,338	237,860	32.06
0.12	-	-	0.00



Disk Queue Length – what and how?

- Dashboards
- Monitoring
- Windows Perf Monitor

SATURATED STORAGE



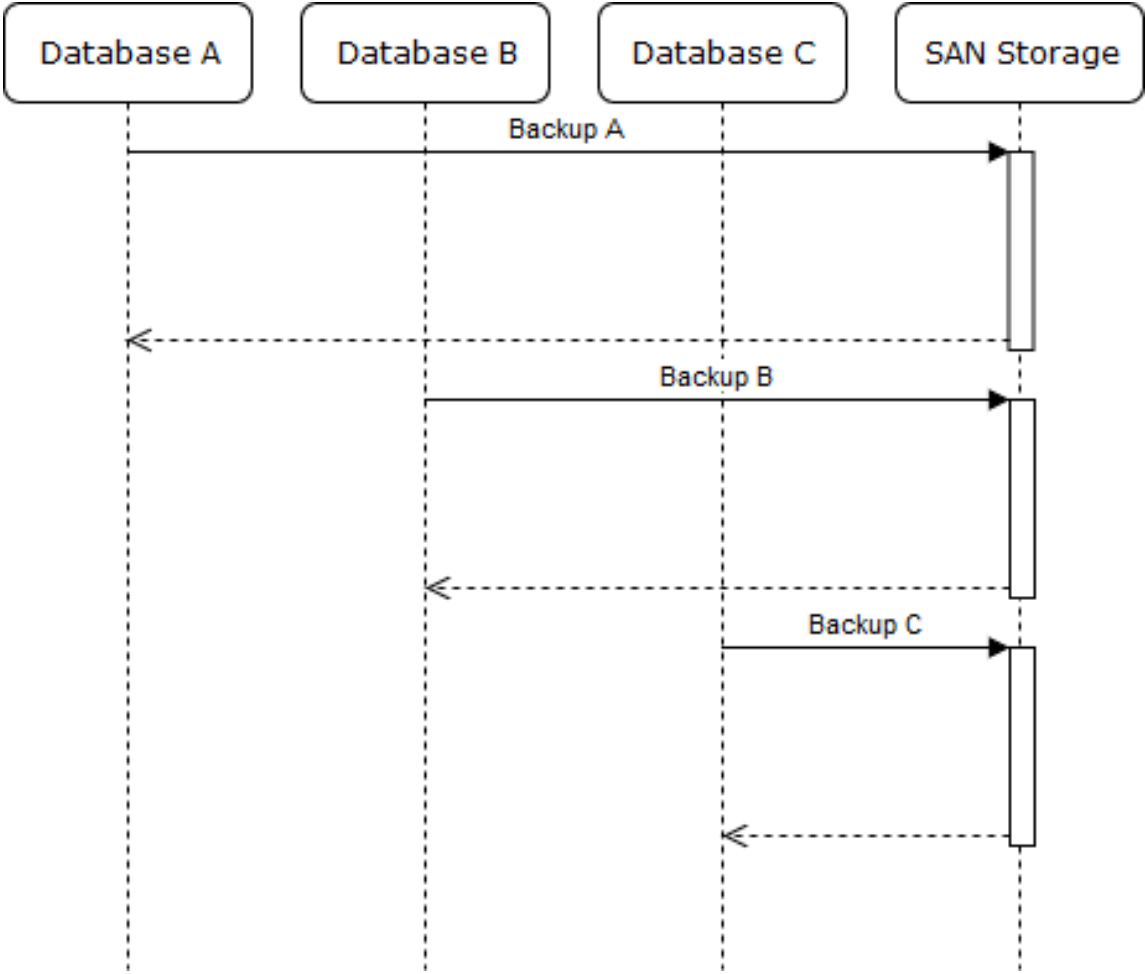
**MAKES
READS/WRITES SLOW?**

RavenDB's failing backups

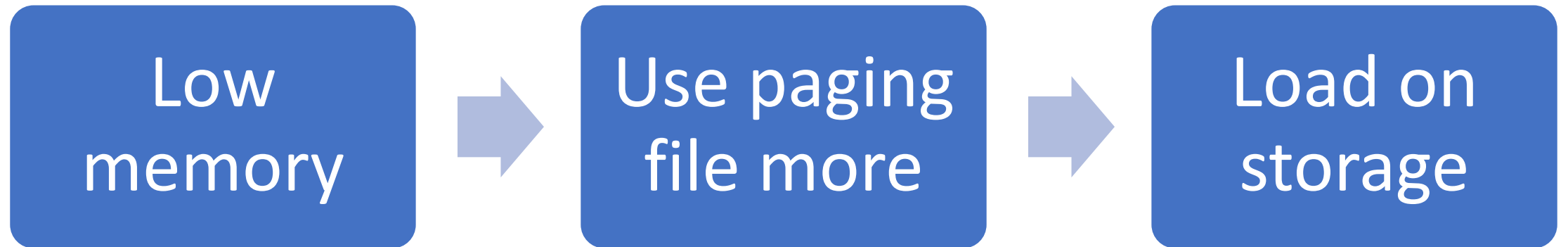
Approx. 200 databases doing backups at the same time WILL cause storage saturation!



The solution was rather simple



Also: low memory matters!



What can we do about storage issues?

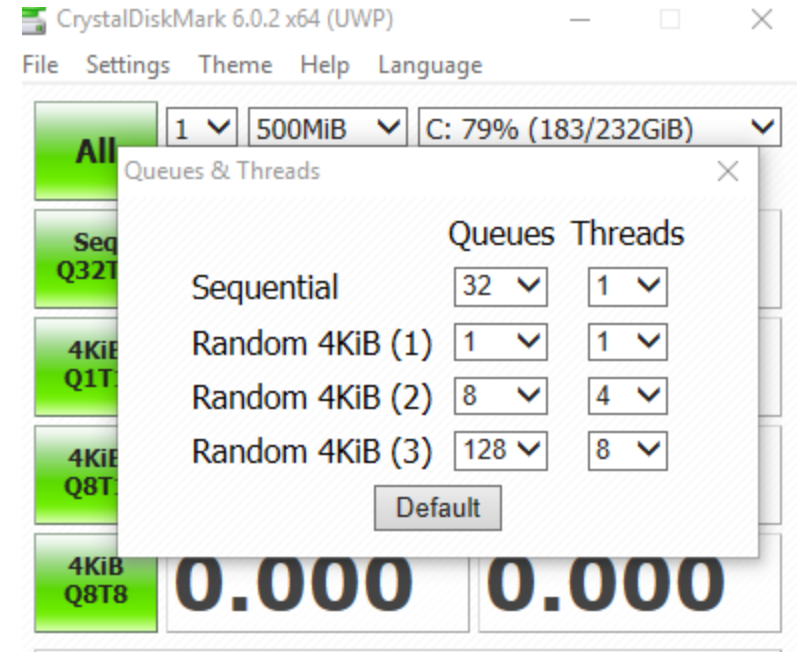
- Load test database-related code
 - **Write-through** throughput
 - Enough IOPS for expected production load (disk queue length is ≤ 2)
- Cloud \rightarrow provision IOPS (ensure disk performance)
- Load-test application to find limits of the system
- Make sure no low-memory situations happen

About storage benchmarks

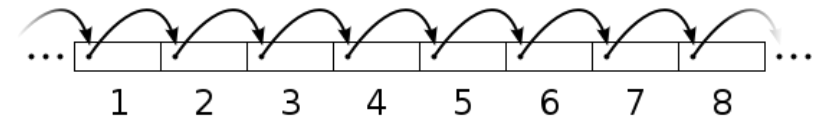
- Sysinternals Process Monitor
- CrystalDiskMark
- ATTO Disk Benchmark
- (Many) other tools

CrystalDiskMark

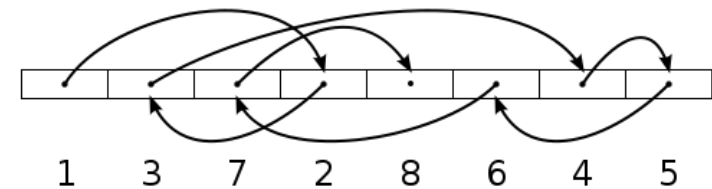
- Random/sequential I/O?
- Queues/Threads (queue depth/length)
- Size of each read/write



Sequential access



Random access



We will discuss some more
benchmarks later

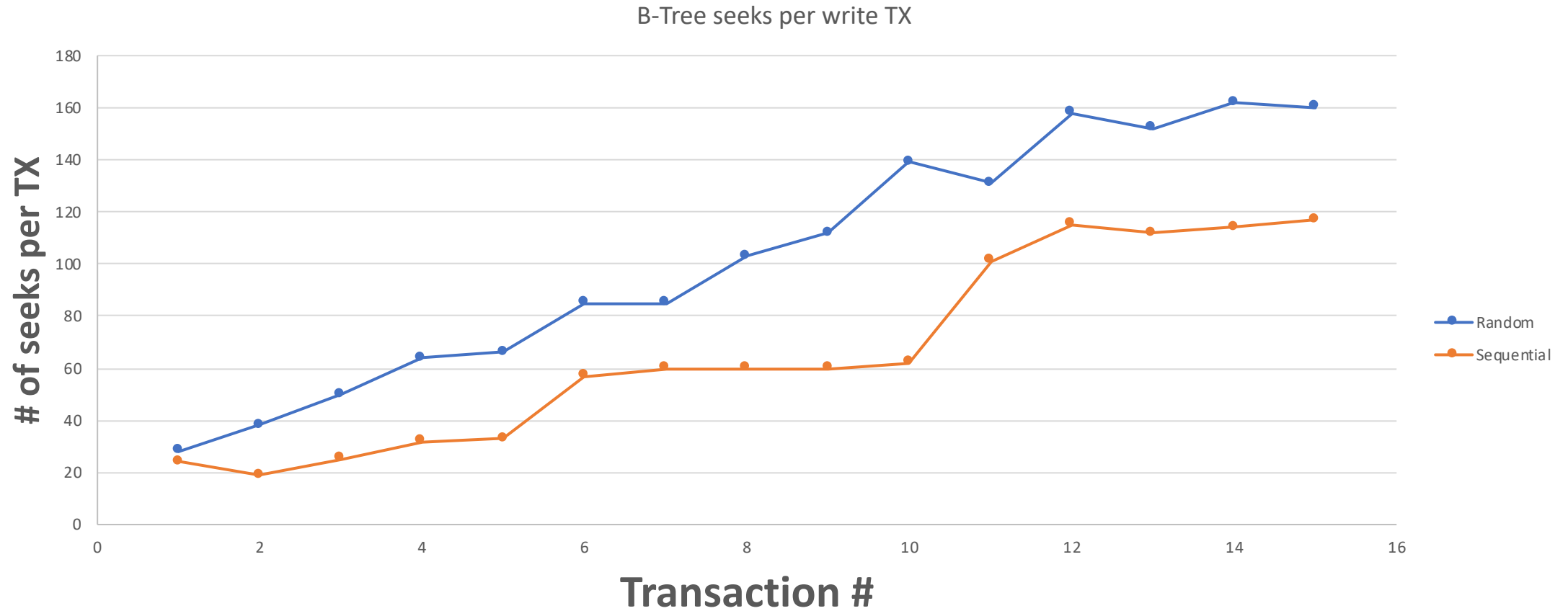
1-b) Storage

Effect of hardware of algorithm performance

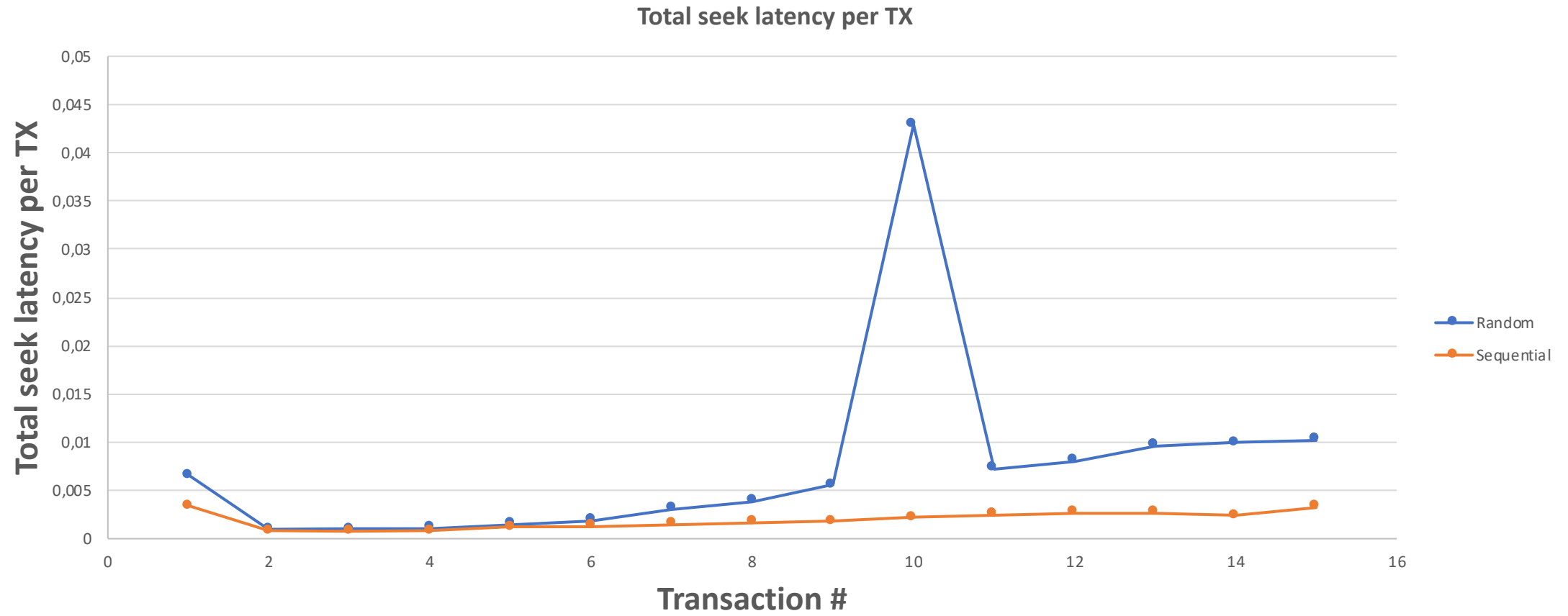
A tale of two primary keys

- One embedded transactional database engine (LMDB)
- 100 transactions, 100 key/value writes per transaction
- Two databases, keys and values have the same size
 - One uses sequential keys (using Win32's ***UuidCreateSequential()***)
 - One uses random keys (using Win32 ***UuidCreate()***)

A tale of two primary keys



A tale of two primary keys



Process Monitor

Types of OS operations to listen (file activity, network activity, etc)

Process Monitor - Sysinternals: www.sysinternals.com

File Edit Event Filter Tools Options Help

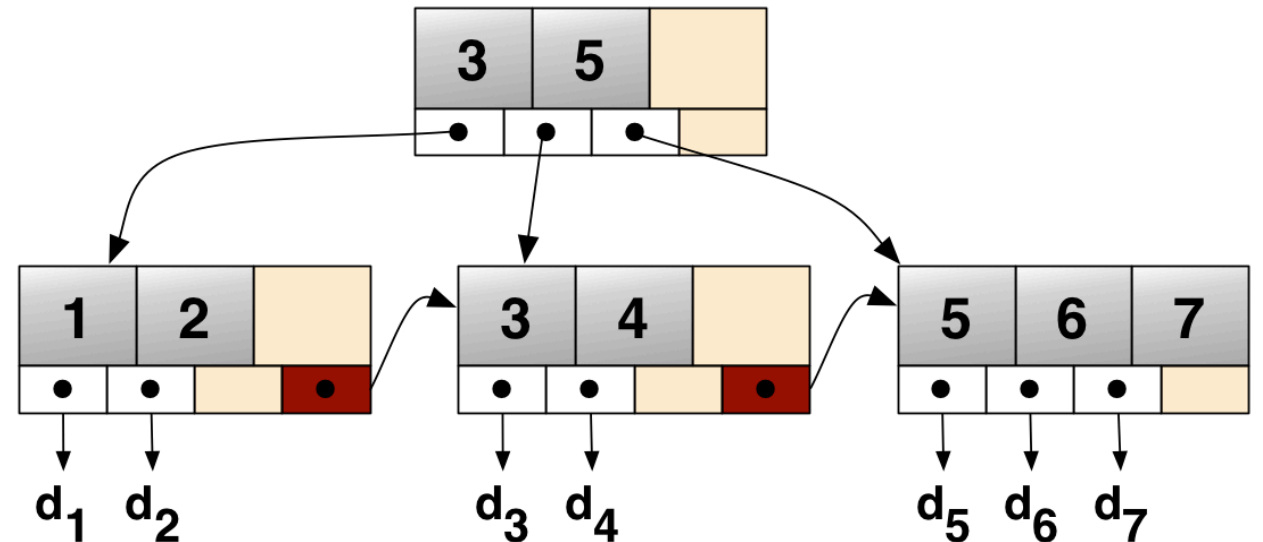
Types of operations

Time ...	Process Name	PID	Operation	Path	Result	Detail
12:10:...	SequentialVsR...	2216	QuerySecurityFile	C:\Windows\System32\msvcpl140.dll	SUCCESS	Information: Owner
12:10:...	SequentialVsR...	2216	CloseFile	C:\Windows\System32\msvcpl140.dll	SUCCESS	
12:10:...	SequentialVsR...	2216	QueryNameInformationFile	C:\Windows\System32\sechost.dll	SUCCESS	Name: \Windows\System32\sechost.dll
12:10:...	SequentialVsR...	2216	QueryNameInformationFile	C:\Windows\System32\sechost.dll	SUCCESS	Name: \Windows\System32\sechost.dll
12:10:...	SequentialVsR...	2216	ReadFile	C:\Windows\System32\msvcpl140.dll	SUCCESS	Offset: 405,504, Length: 4,096, I/O Flags:
12:10:...	SequentialVsR...	2216	ReadFile	C:\Windows\System32\msvcpl140.dll	SUCCESS	Offset: 405,504, Length: 16,384, I/O Flags:
12:10:...	SequentialVsR...	2216	QueryNameInformationFile	C:\Users\Michael.HRHINOS\source\re...	SUCCESS	Name: \Users\Michael.HRHINOS\source\
12:11:...	SequentialVsR...	2216	ReadFile	E:\\$Secure:\$SDS:\$DATA	SUCCESS	Offset: 32,768, Length: 4,096, I/O Flags: f
12:11:...	SequentialVsR...	2216	CreateFile	E:\data\mdb2\random\lock.mdb	SUCCESS	Desired Access: Generic Read/Write, Disp
12:11:...	SequentialVsR...	2216	Lock File	E:\data\mdb2\random\lock.mdb	SUCCESS	Exclusive: True, Offset: 0, Length: 1, Fail I
12:11:...	SequentialVsR...	2216	QueryStandardInformationFile	E:\data\mdb2\random\lock.mdb	SUCCESS	AllocationSize: 0, EndOfFile: 0, NumberOfL
12:11:...	SequentialVsR...	2216	SetEndOfFileInformationFile	E:\data\mdb2\random\lock.mdb	SUCCESS	EndOfFile: 8,192
12:11:...	SequentialVsR...	2216	SetAllocationInformationFile	E:\data\mdb2\random\lock.mdb	SUCCESS	AllocationSize: 8,192
12:11:...	SequentialVsR...	2216	CreateFile Mapping	E:\data\mdb2\random\lock.mdb	FILE LOCKED WI...	SyncType: SyncTypeCreateSection, Page
12:11:...	SequentialVsR...	2216	QueryStandardInformationFile	E:\data\mdb2\random\lock.mdb	SUCCESS	AllocationSize: 8,192, EndOfFile: 8,192, Ni
12:11:...	SequentialVsR...	2216	CreateFile Mapping	E:\data\mdb2\random\lock.mdb	SUCCESS	SyncType: SyncTypeOther
12:11:...	SequentialVsR...	2216	QueryInformationVolume	E:\data\mdb2\random\lock.mdb	BUFFER OVERFL...	VolumeCreationTime: 9/17/2018 7:25:21 I
12:11:...	SequentialVsR...	2216	QueryAllInformationFile	E:\data\mdb2\random\lock.mdb	BUFFER OVERFL...	CreationTime: 8/8/2019 12:11:02 PM, Las
12:11:...	SequentialVsR...	2216	ReadFile	E:\data\mdb2\random\lock.mdb	SUCCESS	Offset: 0, Length: 8,192, I/O Flags: Non-c
12:11:...	SequentialVsR...	2216	CreateFile	E:\data\mdb2\random\data.mdb	SUCCESS	Desired Access: Generic Read/Write, Disp
12:11:...	SequentialVsR...	2216	ReadFile	E:\data\mdb2\random\data.mdb	END OF FILE	Offset: 0, Length: 152, Priority: Normal
12:11:...	SequentialVsR...	2216	WriteFile	E:\data\mdb2\random\data.mdb	SUCCESS	Offset: 0, Length: 8,192, Priority: Normal

Why?

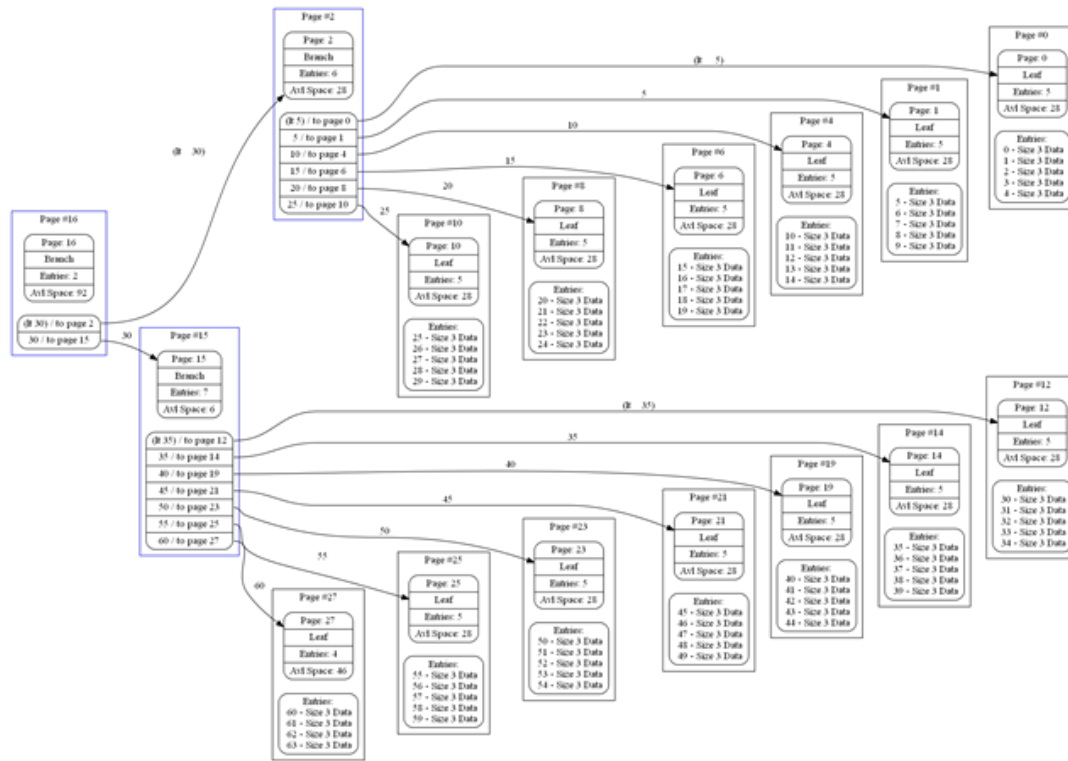
Storage algorithms (page-oriented storage engines)

- B-tree, B+ Tree
- Optimized for reads
- Optimized for sequential data

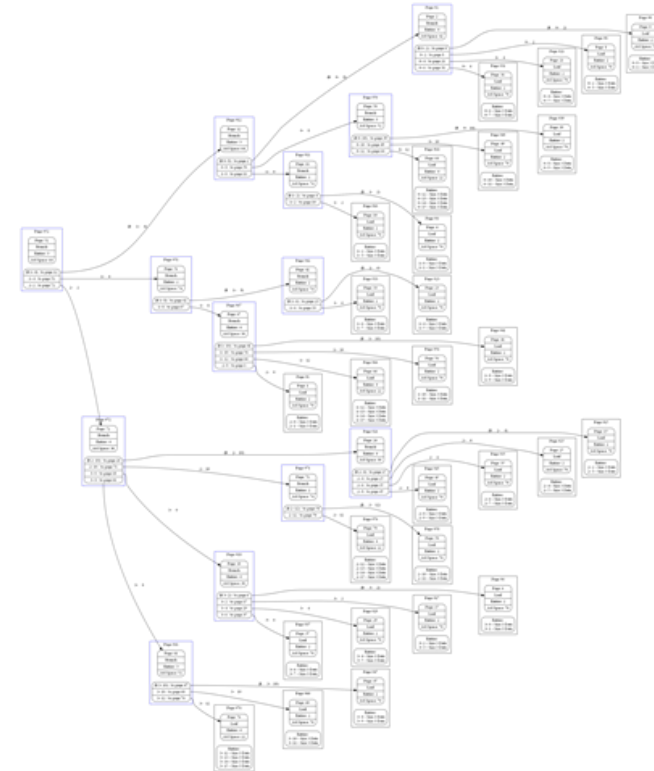


B+ tree keys

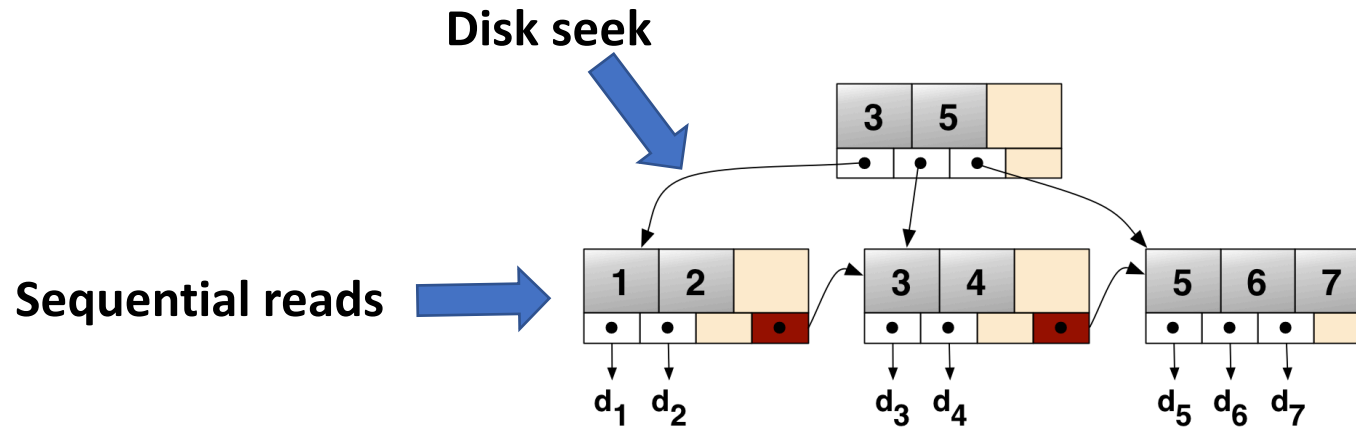
Sequential keys



Non-sequential keys



The cost of hops in the tree



SSD random read: $16,000\text{ns} \approx 16\mu\text{s}$



Disk seek: $3,000,000\text{ns} \approx 3\text{ms}$



Read 1,000,000 bytes sequentially from memory: $4,000\text{ns} \approx 4\mu\text{s}$



Read 1,000,000 bytes sequentially from disk: $947,000\text{ns} \approx 947\mu\text{s}$

Minimize performance impact of keys

- Sequential keys allow better performance
 - 1,2,3,4,5
 - Users/1, Users/2, Users/3
- B-Trees are used to store data AND indexes
 - Query performance!

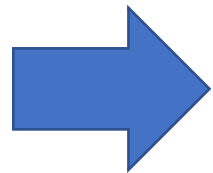
1-c) Storage


Data structure performance


The tale of an occasionally slow database

- Sometimes, Cassandra database was fast and sometimes not
- This happened non-deterministically

A page in
Cassandra
documentation!




Search 

[Search other guides](#)  [Search tips](#)

[Linux related problems](#) / [Reads are getting slower while writes are still fast](#)

Reads are getting slower while writes are still fast

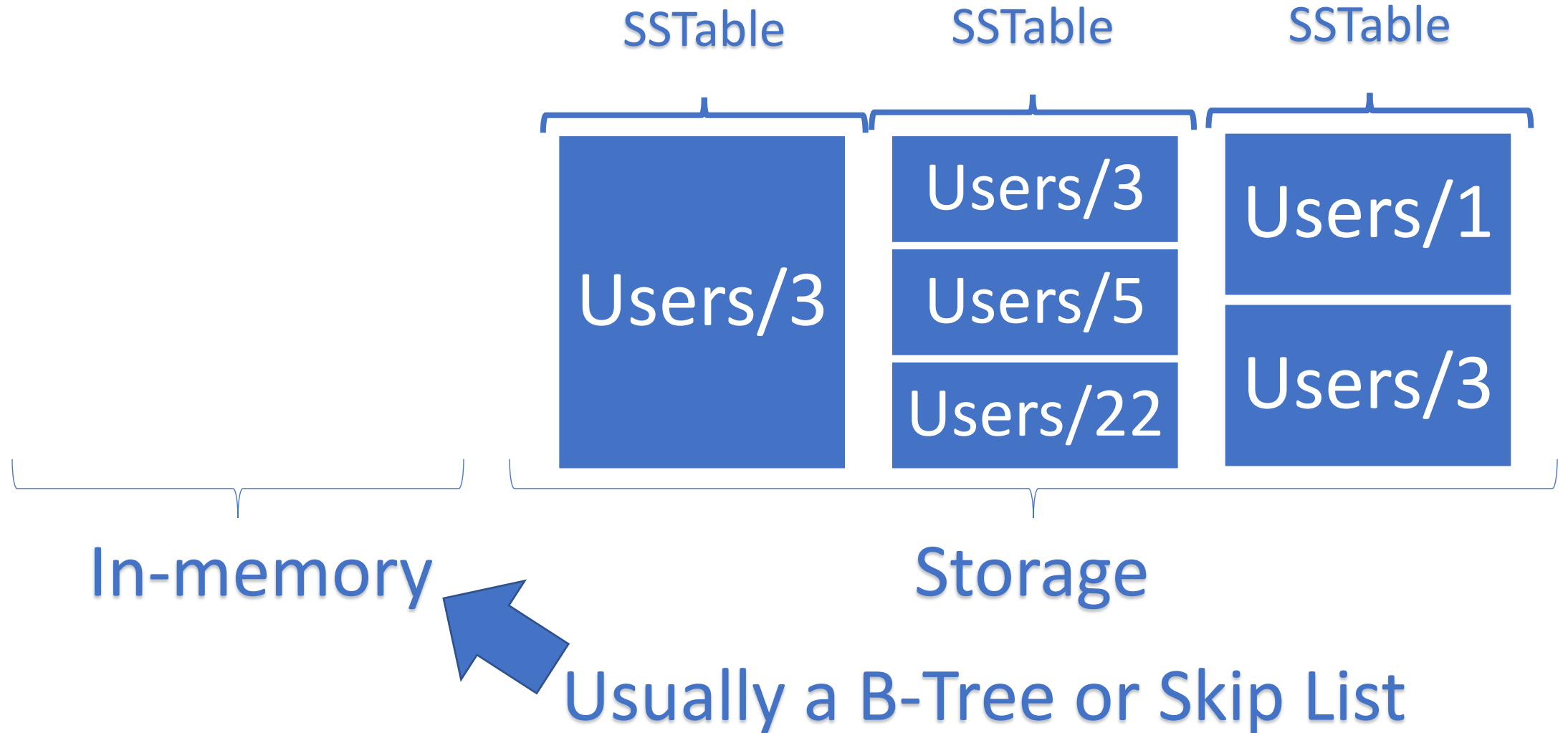
 The [DataStax Enterprise Help Center](#) also provides troubleshooting information.

Too many SSTables can cause slow reads. Take the following steps to determine and correct slow reads:

Determine the total number of SSTables for each table.
Check this number with [nodetool tablestats](#).

Get the number of SSTables consulted for each read.
Check this number with [nodetool tablehistograms](#). A median value over 2 or 3 is likely causing problems.

Storage algorithms (log-structured storage engines)



Storage algorithms (log-structured storage engines)



Insert *users/44*

Storage algorithms (log-structured storage engines)



Update *users/3*

Storage algorithms (log-structured storage engines)



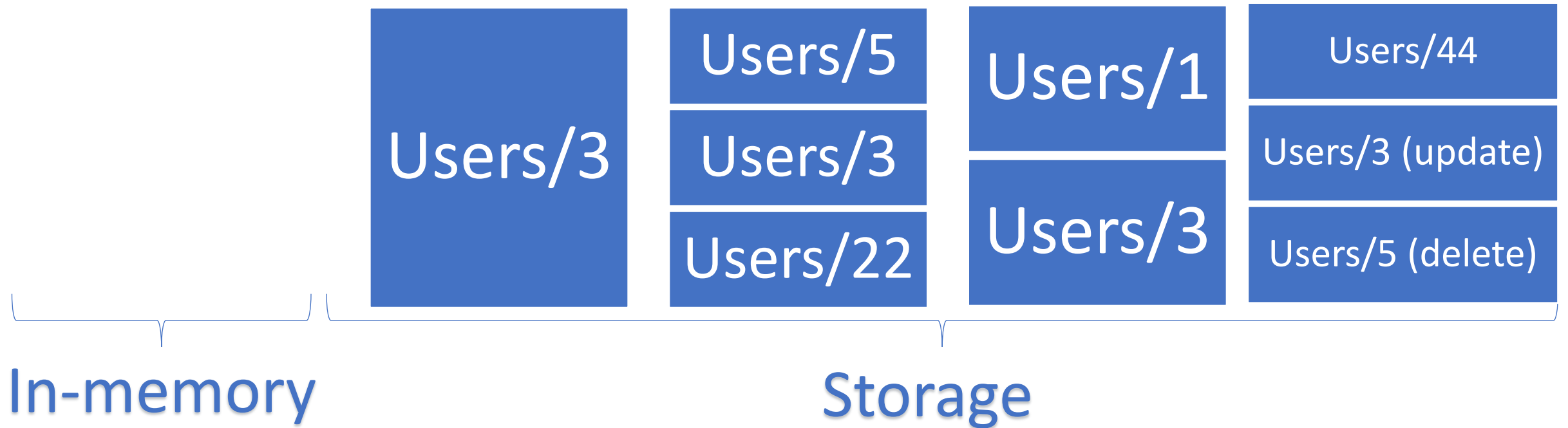
Delete *users/5*

Storage algorithms (log-structured storage engines)



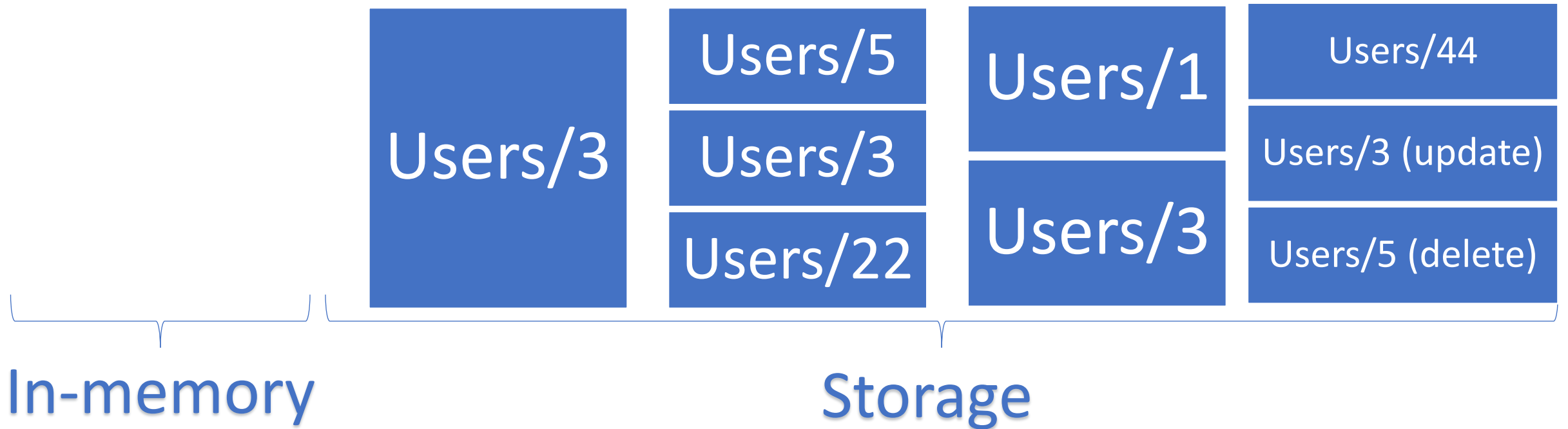
Flush!

Storage algorithms (log-structured storage engines)



(appended at the end)

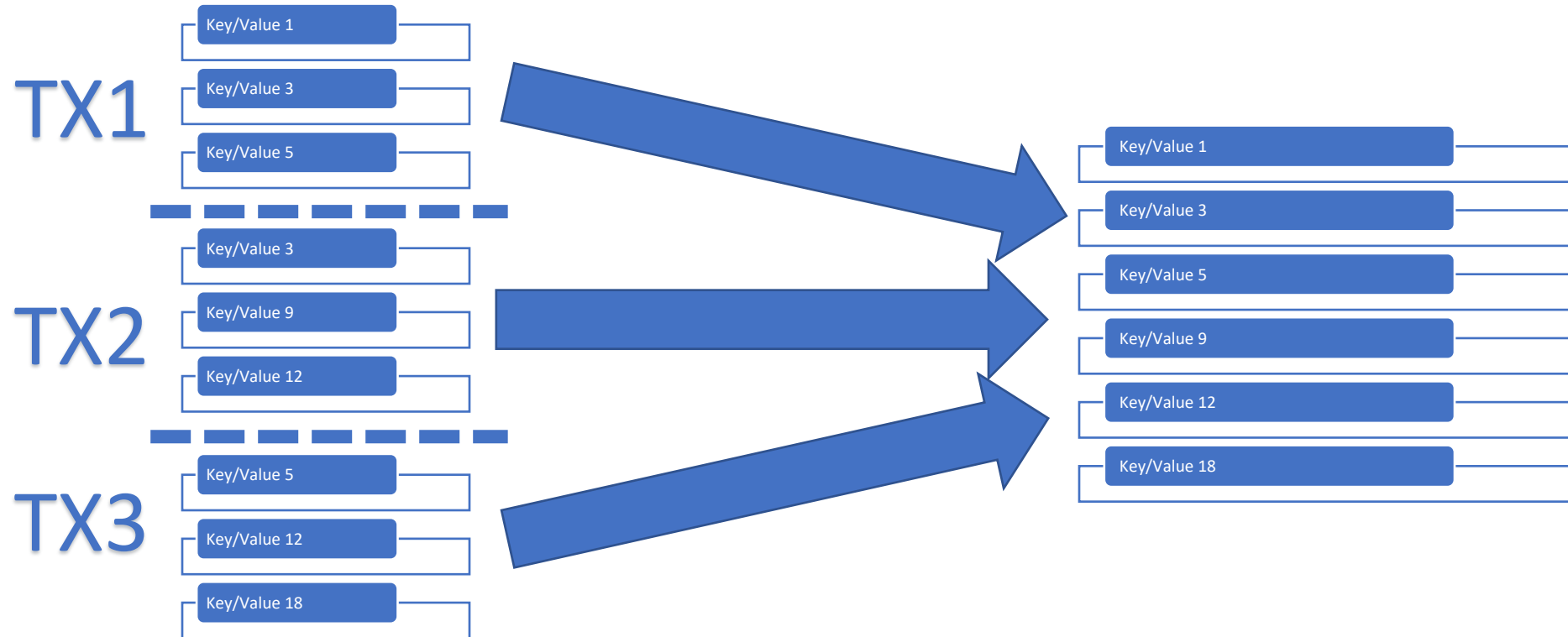
Storage algorithms (log-structured storage engines)



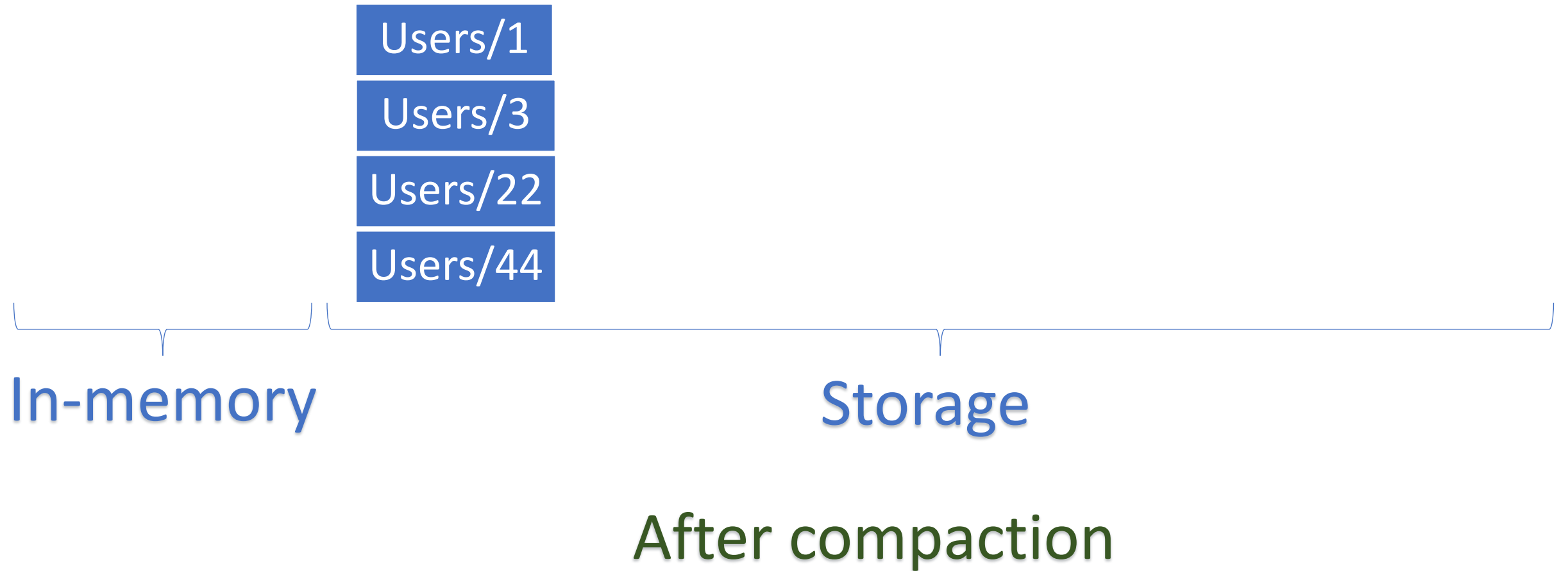
Reading requires searching ALL SSTables!

Storage algorithms – LSM Tree compaction

This is roughly $O(n \cdot \log(n))$ operation!



Storage algorithms – LSM Tree compaction



Performance

- Compact

- Time-

Optimization for different usage patterns!

```
ALTER TABLE users  
WITH compaction =  
{'class' : 'SizeTieredCompactionStrategy', 'min_threshold' : 6 }
```

Different databases have specific options to optimize performance

For example

- MS-SQL index optimization
- RavenDB custom indexes
- MongoDB Aggregation Pipelines

1-d) Storage

Transaction implementation & performance

ACID guarantees!

All operations either succeed or not

- Atomicity

Data is consistent BEFORE and AFTER transaction

- Consistency

Multiple transactions do not interfere each other

- Isolation

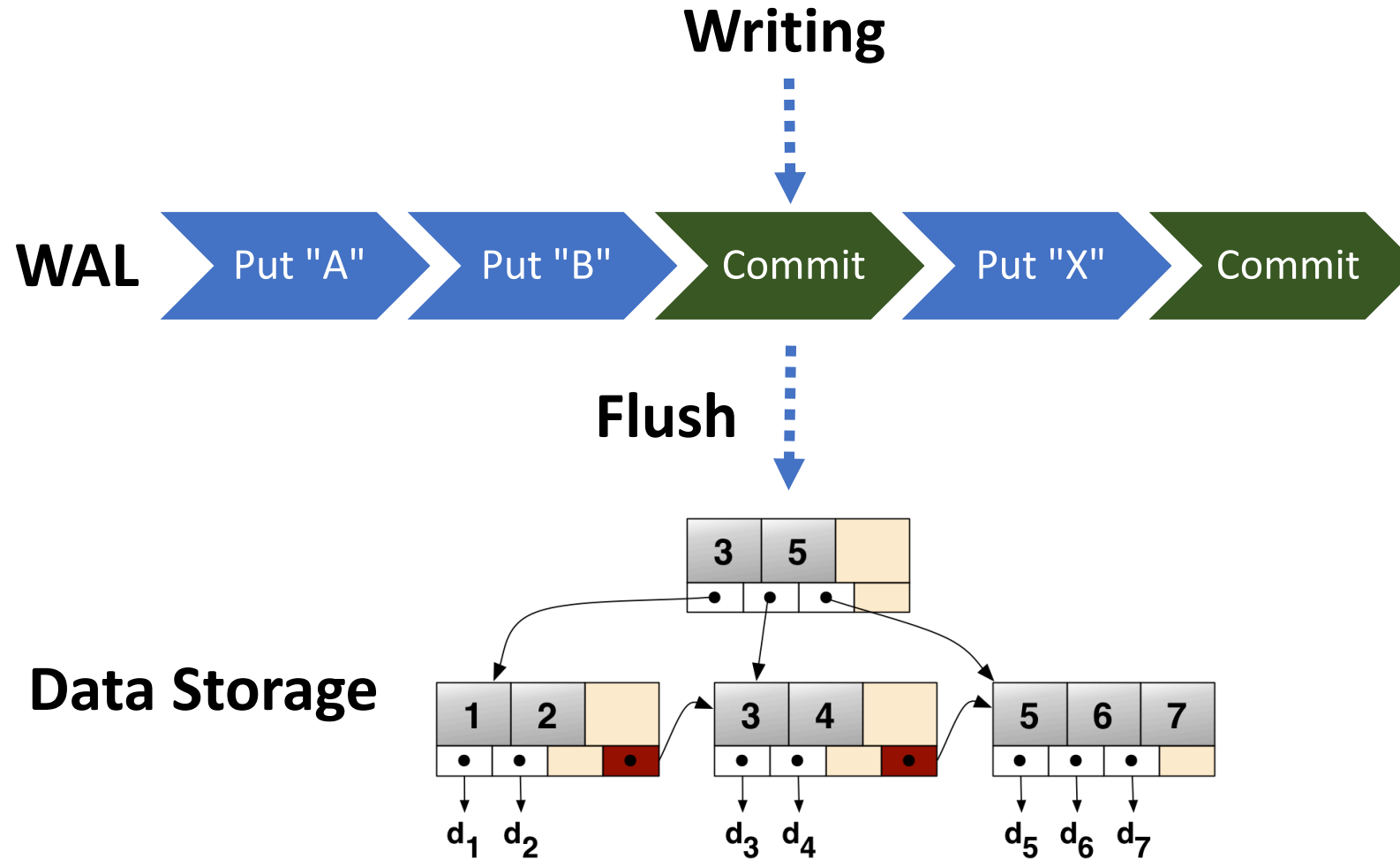
Even if system failure happens, transaction is recorded

- Durability

ACID guarantees!

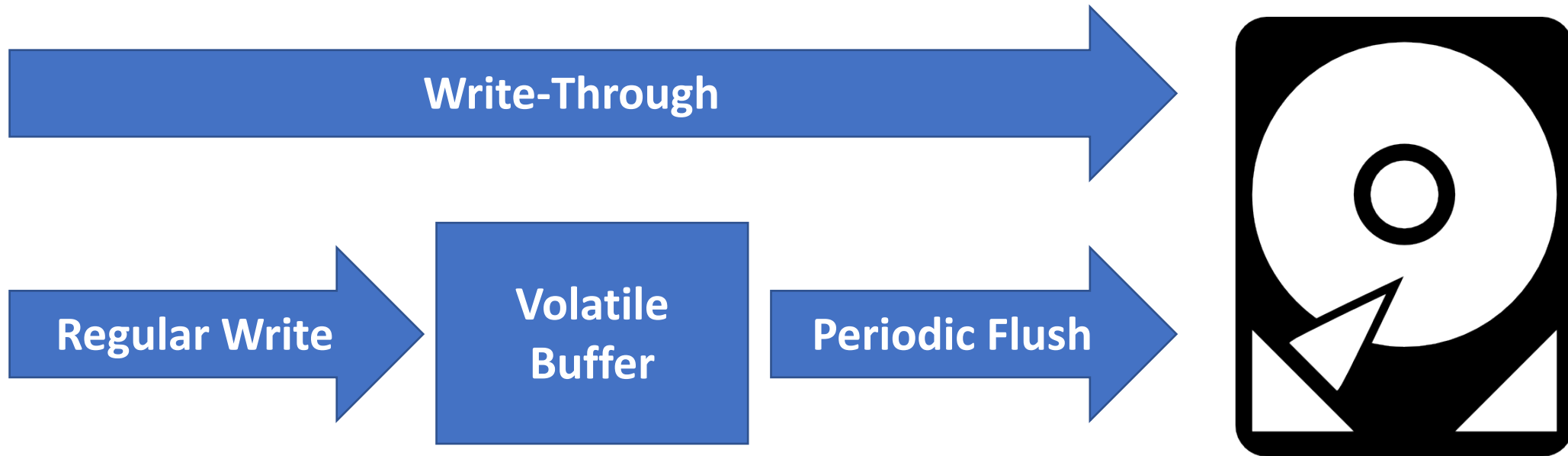
- Note: not all DBs support it
 - All of RDBMS
 - Some NoSQLs – RavenDB, LevelDB, LMDB

Write-ahead log (WAL) - Atomicity, Durability



Write-ahead log (WAL)

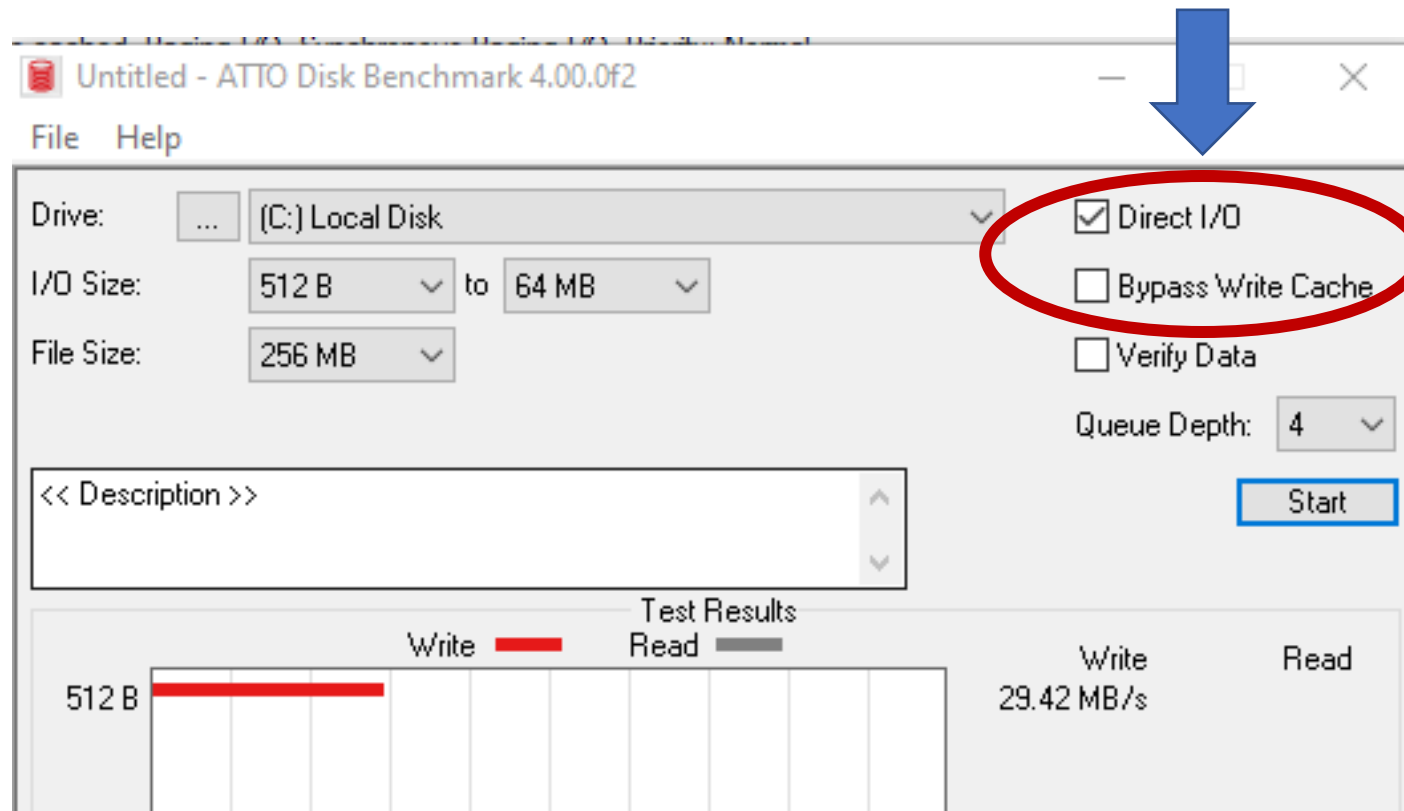
- "Write-Through" writes – no caching (otherwise no durability!)
- Lots of small writes (overhead of each write)



Relevant storage benchmarks

ATTO Disk Benchmark

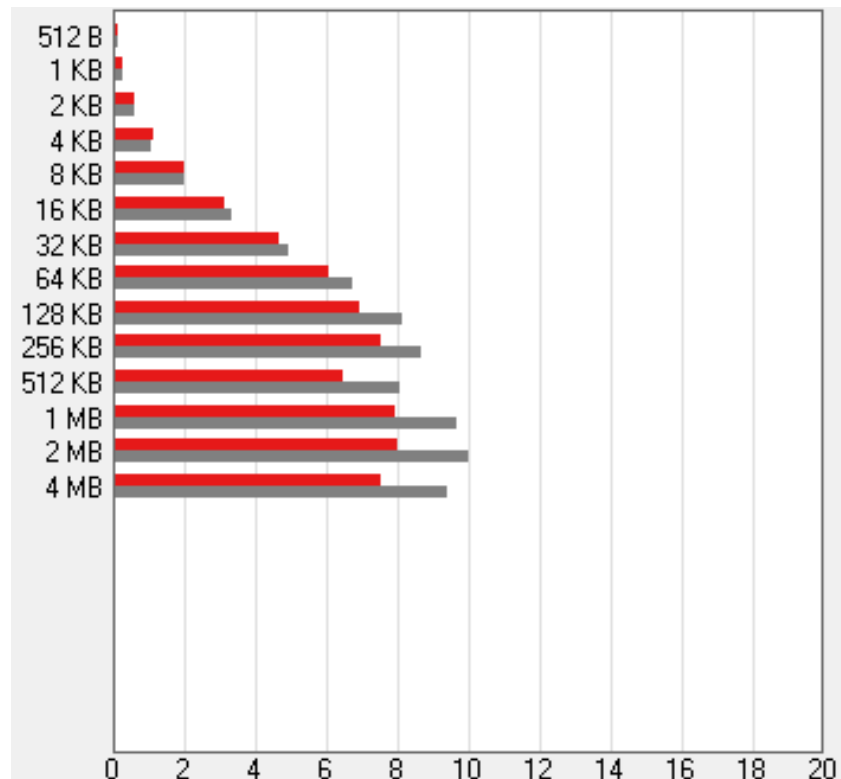
Benchmarking modes



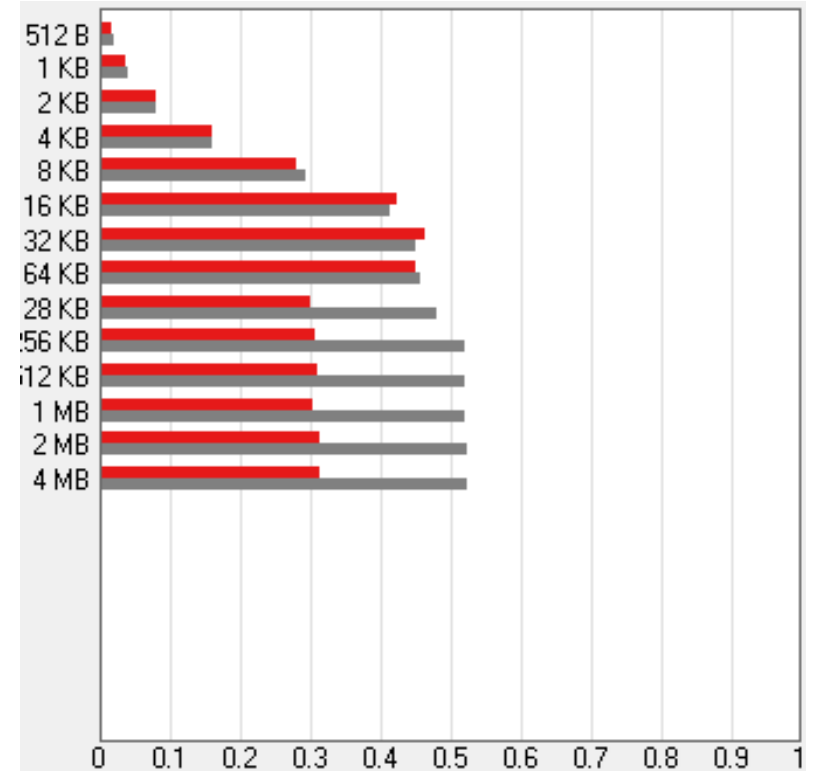
Buffered vs. Write-through

NVMe SSD (Samsung 860 EVO m.2)

With buffers and caching (GBs/sec)



No caching, Write-Through (MBs/sec)



Storage effect on transactions

- Slow storage throughput = bottleneck on transactions
- Write-through performance = transaction throughput

2-a) Indexing & Queries

Query complexity

Are those queries different?



```
SELECT * FROM Orders
WHERE ShipToCity IN ("Paris", "Lyon") AND
      OrderedAt >= '2010-05-01'
```



```
db.Orders.find({
  "ShipTo.City": {
    "$in": ["Paris", "Lyon"]
  },
  "OrderedAt": {
    "$gte": "2010-05-1"
  }
});
```



```
from Orders
where ShipTo.City in ("Paris", "Lyon") and
      OrderedAt >= '2010-05-01'
```

Let's start from something... simple.

```
{
  "Company": "companies/73-A",
  "Employee": "employees/7-A",
  "Freight": 18.44,
  "Lines": [
    {
      "Discount": 0.05,
      "PricePerUnit": 17.45,
      "Product": "products/16-A",
      "ProductName": "Pavlova",
      "Quantity": 14
    }
  ],
  "OrderedAt": "1998-05-06T00:00:00.0000000",
  "RequireAt": "1998-06-03T00:00:00.0000000",
  "ShipTo": {
    "City": "Kobenhavn",
    "Country": "Denmark",
    "Line1": "Vinbæltet 34",
    "Line2": null,
    "Location": null,
    "PostalCode": "1734",
    "Region": null
  }
},
```


First, we define an index

We create an index that covers **city** and **country** fields of *ShipTo*

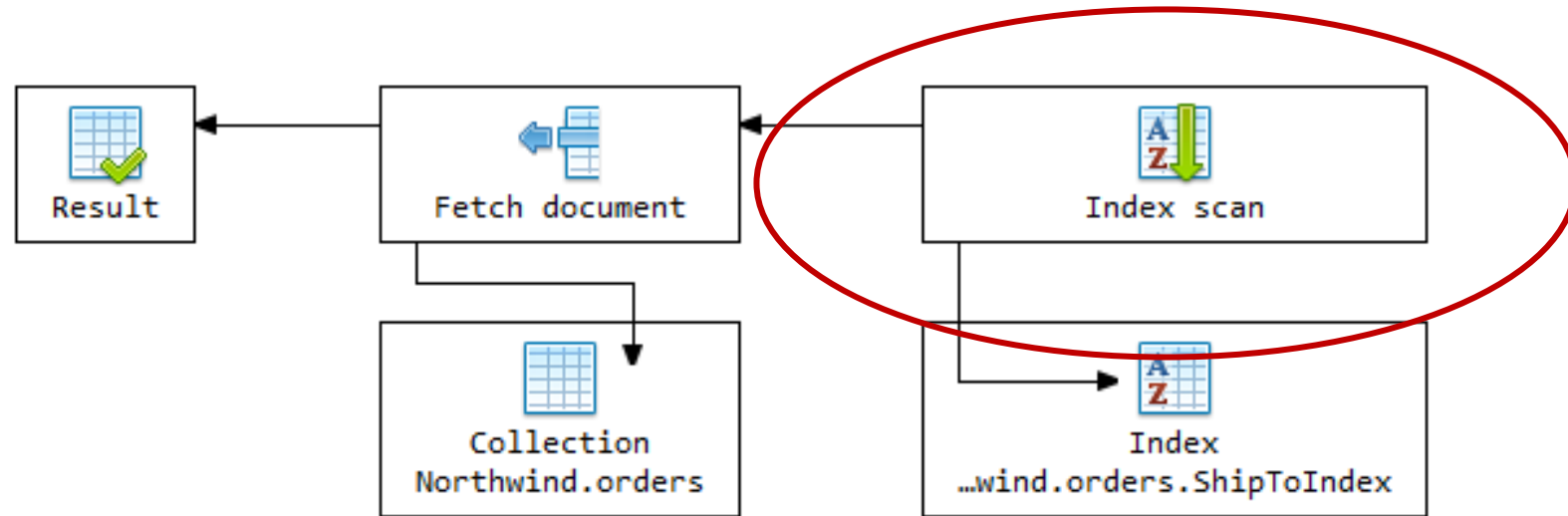
```
> db.orders.getIndexes()
[
  {
    "v" : 2,
    "key" : {
      "_id" : 1
    },
    "name" : "_id_",
    "ns" : "Northwind.orders"
  },
  {
    "v" : 2,
    "key" : {
      "ShipTo.City" : 1,
      "ShipTo.Country" : 1
    },
    "name" : "ShipToIndex",
    "ns" : "Northwind.orders"
  }
]
```

Then we do some queries

Fetching orders that were shipped to Paris or Lyon

```
db.orders.find({
  "ShipTo.City": {
    "$in": ["Paris", "Lyon"]
  }
});
```

So far so good...

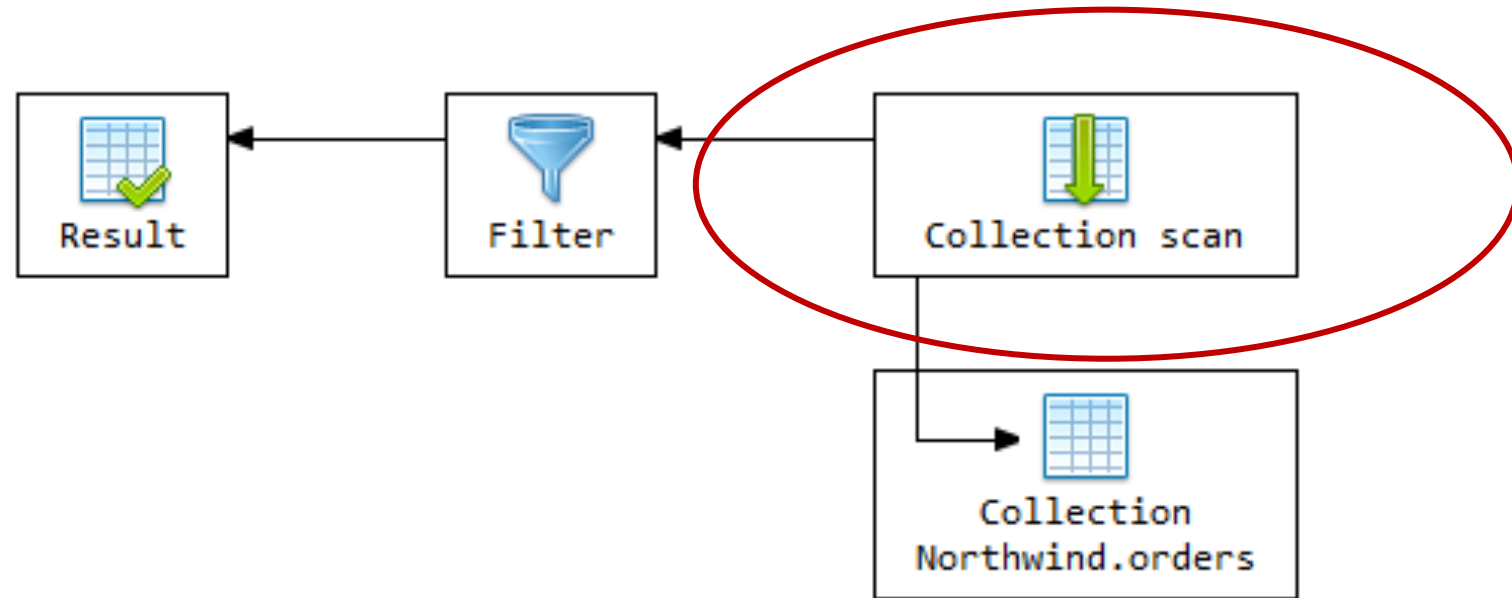


And another query

Fetching orders that were shipped to all of France

```
db.orders.find({  
    "ShipTo.Country": "France"  
});
```

It's a gotcha!



Index Scan vs. Collection Scan

Collection Scan

- $O(n)$ scan

Index Scan

- $O(\log(n))$ seek

Why?

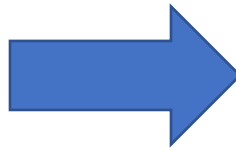
In some databases (like MongoDB)

- Indexed fields are concatenated into single index key
- Filtering only by prefix

The values are concatenated!



Field 1	Field 2
Lyon	France
Paris	France
Oslo	Norway

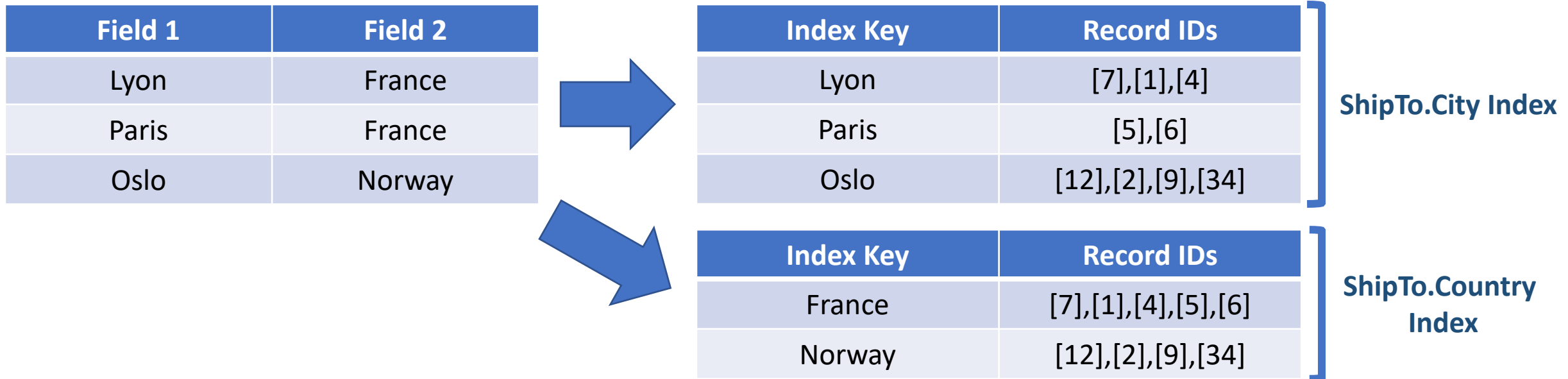


Index Key	Record IDs
LyonFrance	[7],[1],[4]
ParisFrance	[5],[6]
OsloNorway	[12],[2],[9],[34]

Why?

In some databases (RavenDB, any Lucene-based index)

- Indexed terms stored separate
- Filtering by one or both fields in any order (union/intersect as needed)



Collection/table scans are easily overlooked

Collection scan - development

- Small amount of data
- Extremely small query latency

Collection scan - production

- Large amount of data
- HUGE latency (quite often!)

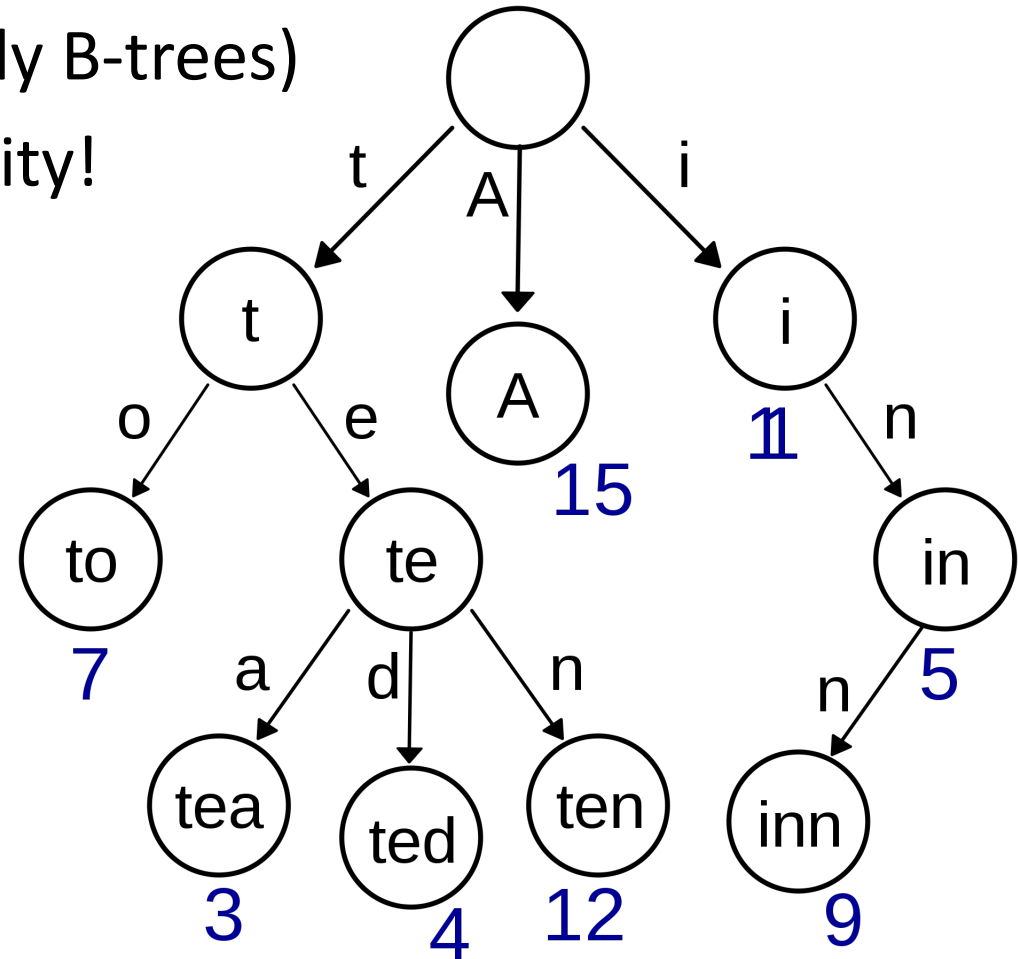
Latency: 50ms vs 50 hours

2-b) Indexing & Queries

More about indexing

Indexing

- Indexes are stored as trees (usually B-trees)
- Updates have non-trivial complexity!



Indexing

```
from Orders
where ShipTo.City in ("Paris", "Lyon") and
      OrderedAt >= '2010-05-01'
```



Search time complexity (WHERE clause):

$O(\log(N)) + O(\log(M)) + O(\text{Max}(K,P))$

Where:

- *N and M are amount of rows in indexes*
- *K and P are result sets of index searches*

And if we use RDBMS, things become even more interesting...

Join Algorithm	Complexity
Merge Join	$O(n \cdot \log(n) + m \cdot \log(m))$
Hash Join	$O(n + m)$
Index Join	$O(m \cdot \log(n))$

And if we have a non-trivial query...

```
SELECT
  e.employee_id AS "Employee #"
  , e.first_name || ' ' || e.last_name AS "Name"
  , e.email AS "Email"
  , e.phone_number AS "Phone"
  , TO_CHAR(e.hire_date, 'MM/DD/YYYY') AS "Hire Date"
  , TO_CHAR(e.salary, 'L99G999D99', 'NLS_NUMERIC_CHARACTERS = ''.,'' NLS_CURRENCY = '$''') AS "Salary"
  , e.commission_pct AS "Commission %"
  , 'works as ' || j.job_title || ' in ' || d.department_name || ' department (manager: ' || dm.first_name || ' and immediate manager: ' || m.first_name || ' )' || ' ' || e.first_name || ' ' || e.last_name AS "Current Job"
  , TO_CHAR(j.min_salary, 'L99G999D99', 'NLS_NUMERIC_CHARACTERS = ''.,'' NLS_CURRENCY = '$''') || ' - ' || TO_CHAR(j.max_salary, 'L99G999D99', 'NLS_NUMERIC_CHARACTERS = ''.,'' NLS_CURRENCY = '$''') AS "Current Salary"
  , l.street_address || ', ' || l.postal_code || ', ' || l.city || ', ' || l.state_province || ', ' || c.country_name || ' (' || r.region_name || ' )' AS "Location"
  , jh.job_id AS "History Job ID"
  , 'worked from ' || TO_CHAR(jh.start_date, 'MM/DD/YYYY') || ' to ' || TO_CHAR(jh.end_date, 'MM/DD/YYYY') || ' as ' || jj.job_title || ' in ' || dd.department_name || ' department' AS "History Job Title"
FROM employees e JOIN jobs j ON e.job_id = j.job_id
LEFT JOIN employees m ON e.manager_id = m.employee_id
LEFT JOIN departments d ON d.department_id = e.department_id
LEFT JOIN employees dm ON d.manager_id = dm.employee_id
LEFT JOIN locations l ON d.location_id = l.location_id
LEFT JOIN countries c ON l.country_id = c.country_id
LEFT JOIN regions r ON c.region_id = r.region_id
LEFT JOIN job_history jh ON e.employee_id = jh.employee_id
LEFT JOIN jobs jj ON jj.job_id = jh.job_id
LEFT JOIN departments dd ON dd.department_id = jh.department_id
ORDER BY e.employee_id;
```

Usually, there is also a WHERE clause, so...

Those are 10 JOIN statements!

...we have complexity between
 $O(\log(n))$ and $O(\textit{too much})!$

More often than not it is $O(\textit{too much})...$

What can (should!) we do?

- RDBMS

- Proper indexing (kinda obvious, but still 😊)
- Optimize (remove unnecessary JOINS – depends on business logic)
- Reduce query complexity
 - Replace 'row by row' cursors with set based queries
 - Reduce the amount of work queries do (for example, unnecessary sub-queries)
 - Remove ORDER BY where it makes sense (huge overhead)
 - Other optimizations are possible

- NoSQL

- Proper modeling
- Well planned indexing

2-c) Indexing & Queries

Indexes (sometimes) have complexity too!

Indexing complexity

```
db.runCommand( {
  mapReduce: "blogs",
  map: function(){
    for (let index = 0; index < this.authors.length; ++index) {
      let author = this.authors[ index ];
      emit( author.firstName + " " + author.lastName, 1 );
    }
  },
  reduce: function(author, counters){
    count = 0;

    for (let index = 0; index < counters.length; ++index) {
      count += counters[index];
    }

    return count;
  },
  out: { inline: 1 }
} )
```

Indexing complexity

Maps

Syntax ?

```
1 from order in docs.Orders
2 select new
3 {
4     order.Company,
5     Count = 1,
6     Total = order.Lines.Sum(l => (l.Quantity * l.PricePerUnit) * (1 - l.Discount))
7 }
```

Press Shift+F1

+ Add map

Reduce

Syntax ?

```
1 from result in results
2 group result by result.Company
3 into g
4 select new
5 {
6     Company = g.Key,
7     Count = g.Sum(x => x.Count),
8     Total = g.Sum(x => x.Total)
9 }
```

Press Shift+F1

Indexing complexity

- As you can see, indexing has its own complexity
- More often than not it can be optimized

3-a) Network

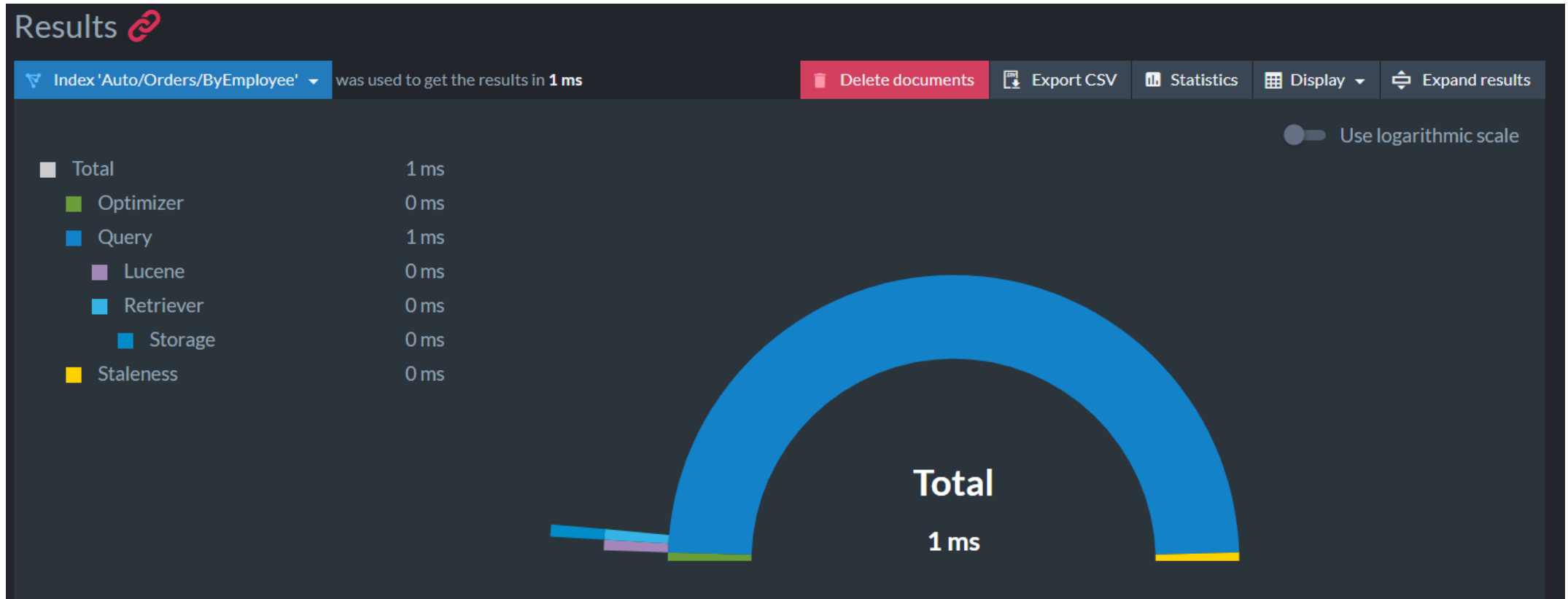
Distributed system fallacy: Bandwidth is infinite

*Here is a riddle: why a query with 100 results
takes several seconds to complete?*

Hint: the request spends < 10ms on the server

The investigation

1. Look at query latency on the server



The investigation

2. Look at Fiddler timings

Response Size

The screenshot displays the Fiddler Orchestra Beta interface. At the top, there are tabs for Statistics, Inspectors, AutoRe, Composer, and Fiddler Orchestra Beta. Below the tabs, the following statistics are shown:

- Request Count: 1
- Bytes Sent: 848 (headers: 848; body: 0)
- Bytes Received: 3,287,227 (headers: 267; body: 3,286,960)

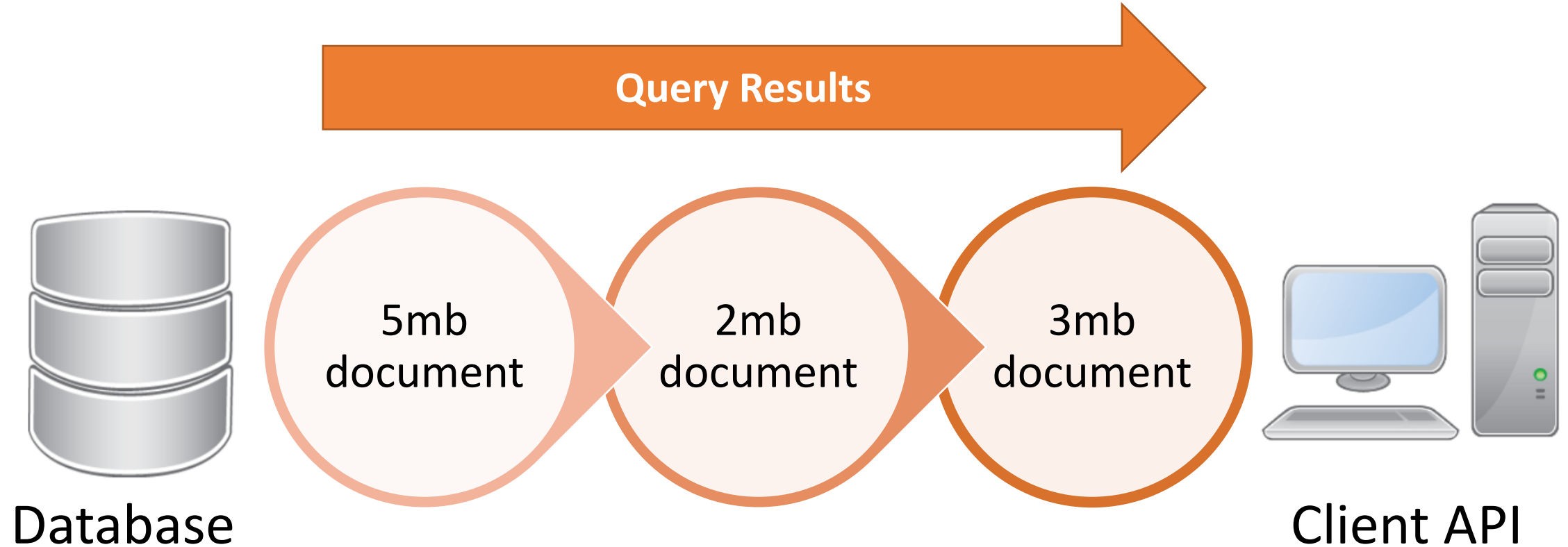
The 'Bytes Received' value is circled in red. Below this, the 'ACTUAL PERFORMANCE' section is shown with a list of events and their corresponding timestamps:

```
ACTUAL PERFORMANCE
-----
ClientConnected:      13:42:01.036
ClientBeginRequest:  13:42:38.211
GotRequestHeaders:   13:42:38.211
ClientDoneRequest:   13:42:38.211
Determine Gateway:   0ms
DNS Lookup:          0ms
TCP/IP Connect:      199ms
HTTPS Handshake:     0ms
ServerConnected:     13:42:38.411
FiddlerBeginRequest: 13:42:38.411
ServerGotRequest:    13:42:38.411
ServerBeginResponse: 13:42:38.644
GotResponseHeaders:  13:42:38.644
ServerDoneResponse:  13:42:41.273
ClientBeginResponse: 13:42:41.273
ClientDoneResponse:  13:42:41.273
```

At the bottom, the 'Overall Elapsed' time is shown as 0:00:03.061, which is also circled in red. Two blue arrows point from the text 'Response Size' and 'Latency' to the circled values.

Latency

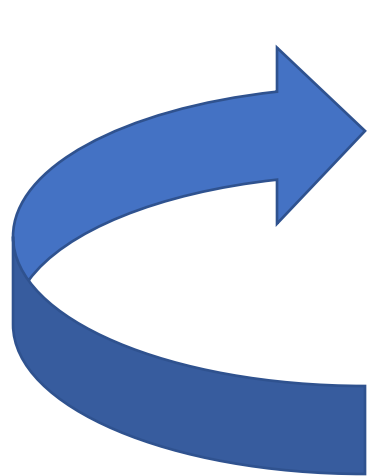
Network bandwidth is not infinite!



Solution: server-side projections (NoSQL)

RQL

```
from index 'Orders/ByShipment/Location'  
where spatial.within(ShipmentLocation,  
                    spatial.wkt('Circle(2.349014 48.864716 d=1000.0000)'))  
load Company as c, Employee as e  
select {  
  CompanyName: c.Name,  
  EmployeeName: e.FirstName + " " + e.LastName  
}
```



Server-side Projection

Solution: server-side projections (NoSQL)

MongoDB API

```
var findDocuments = function(db, callback) {  
  var collection = db.collection( 'restaurants' );  
  // Find some documents  
  collection.find({ 'cuisine' : 'Brazilian' }, { 'name' : 1, 'cuisine' : 1 })  
    .toArray(function(err, docs) {  
      console.log("Found the following records");  
      //do something with found records  
    });  
}
```

Server-side Projection

3-b) Network

Distributed system fallacy: Latency is zero

Also... database requests can be an interesting issue...

Support » Plugin: WooCommerce » Site being slowed by too many (5000+) database queries

Site being slowed by too many (5000+) database queries



dynit (@dynit)
3 months, 3 weeks ago

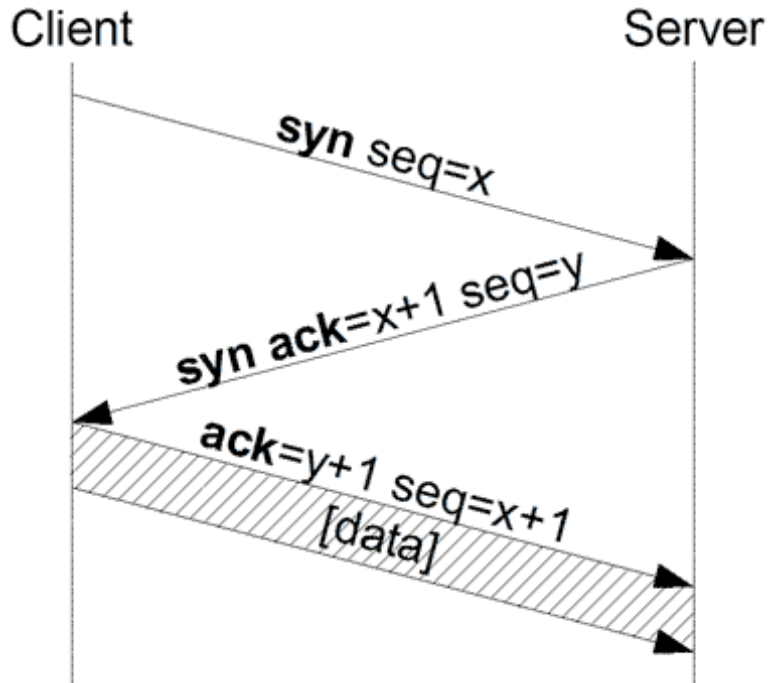
✓ Resolved

Hello,

One of our sites is really getting slowed with too many Woocommerce database queries. I checked via Query Monitor and just browsing through categories (on the default shop page with the post type archive for products) fires 5000+ database queries per page load. The site has 66 categories, about 600 products, 2200 variations, and 6 custom attributes per product. Single product pages are 'only' 150-300 queries each, no performance issue there.

Network overhead

TCP handshake



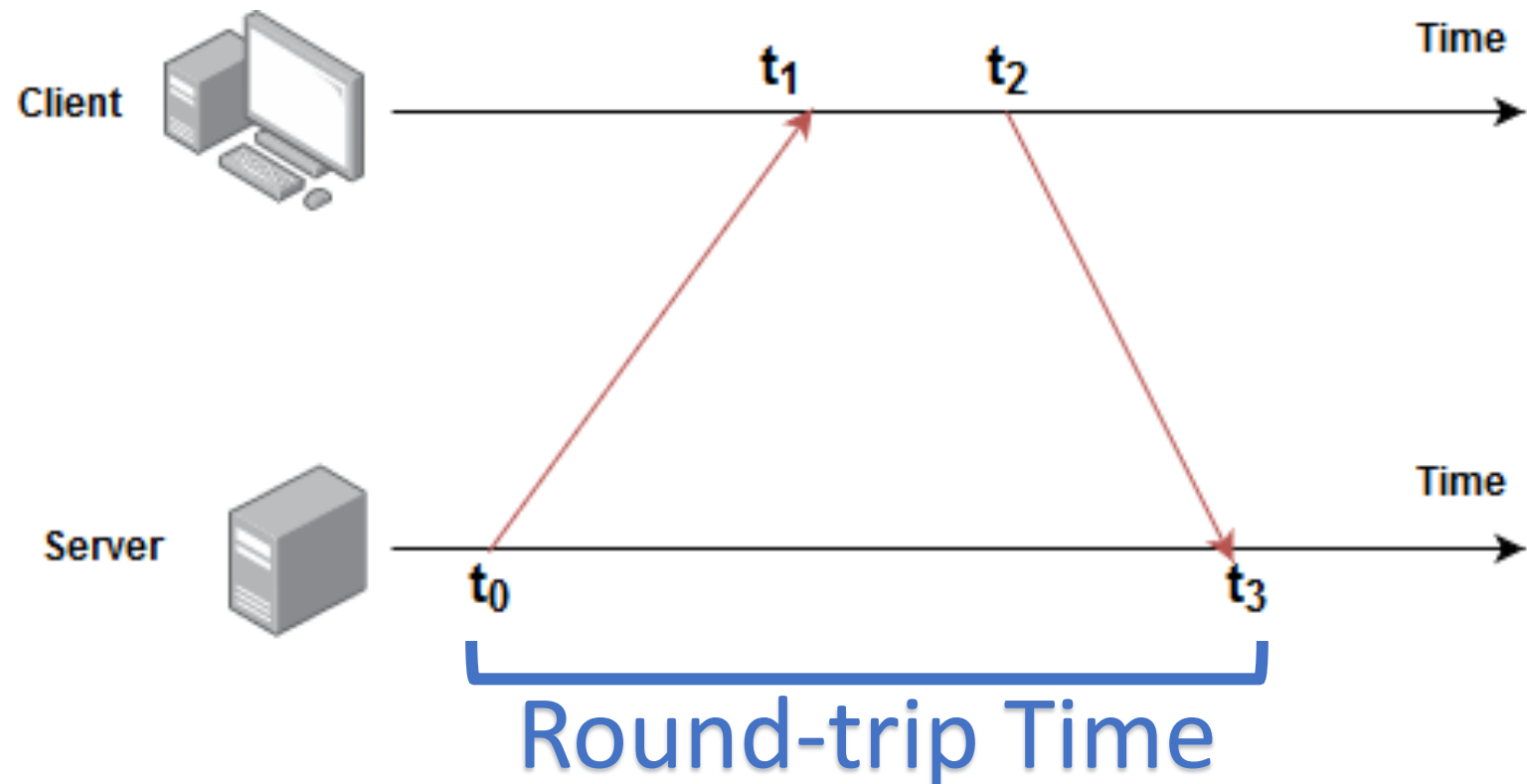
```
Route.Add("/hello",  
    (request, response, args) =>  
        response.AsText("world!"));  
  
var cts = new CancellationTokenSource();  
var _ = HttpServer.ListenAsync(8080,  
    cts.Token, Route.OnHttpRequestAsync);  
  
Console.ReadKey();  
  
cts.Cancel();
```

```
michael@michael-pc:~$ wrk -t16 -c1000 -d30s http://127.0.0.1:8080/hello  
Running 30s test @ http://127.0.0.1:8080/hello  
16 threads and 1000 connections  
Thread Stats Avg Stdev Max +/- Stdev  
Latency 21.62ms 93ms 61.91ms 93.59%  
Req/Sec 2.87 290.86 18.69k 93.98%  
1369876 requests in 30.09s, 178.98MB read  
Requests/sec: 45524.71  
Transfer/sec: 5.95MB
```

Network overhead

Round Trip Time (RTT)

- Physical distance (insignificant for LANs)
- Bandwidth
- Network hops



What can we do?

- Refactor to reduce number of requests (kinda obvious, but still...)
- NHibernate – Future Queries
- Entity Framework - QueryFuture
- RavenDB – Lazy Queries

```
using(var session = store.OpenSession())
{
    var lazyUser = session.Advanced.Lazily.Load<User>("users/michael");
    var lazyPosts = session.Query<Posts>().Take(30).Lazily();

    session.Advanced.Eagerly.ExecuteAllPendingLazyOperations();

    //do something with lazyUser and lazyPosts
}
```


May sound trivial, but...

*Do take a look at database traffic while **stress testing** and if possible in **production** too.*

- Fiddler
- Wireshark
- Profilers
- Any other tool to inspect traffic

To sum it up

- Databases are abstractions
- Abstractions are leaky and might be the cause of perf issues
- Such perf issues can be dealt with (if we know about the "leak"!)

Questions?

michael.yarichuk@hibernatingrhinos.com

@myarichuk

<https://github.com/ravendb/ravendb>

<https://github.com/myarichuk/PerfDemo-Sequential-vs-Random-Key>

