

twitter



340 Million

Tweets per day

2.3 Billion

Queries per day

< 10 s

Indexing latency

50 ms

Avg. query response time

Earlybird - Realtime Search @twitter

Michael Busch

@michibusch

michael@twitter.com

buschmi@apache.org



Earlybird - Realtime Search @twitter

Agenda

- ▶ Introduction
- Search Architecture
- Inverted Index 101
- Memory Model & Concurrency
- What's next?

Introduction

Introduction

S U M M I Z E

Realtime Twitter Search

[Show Options](#)

Search

- Twitter acquired Summize in 2008
- 1st gen search engine based on MySQL

Introduction



- Next gen search engine based on Lucene
- Improves scalability and performance by orders or magnitude
- Open Source

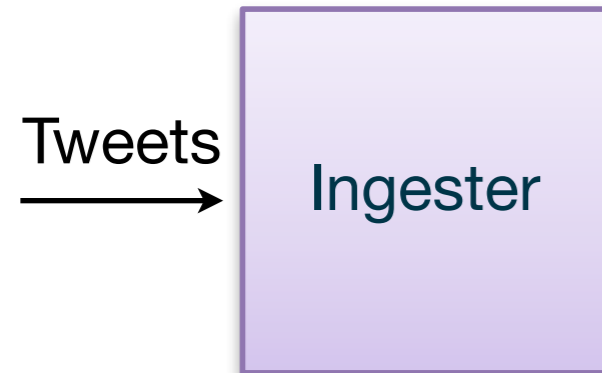
Realtime Search @twitter

Agenda

- Introduction
- ▶ **Search Architecture**
- Inverted Index 101
- Memory Model & Concurrency
- What's next?

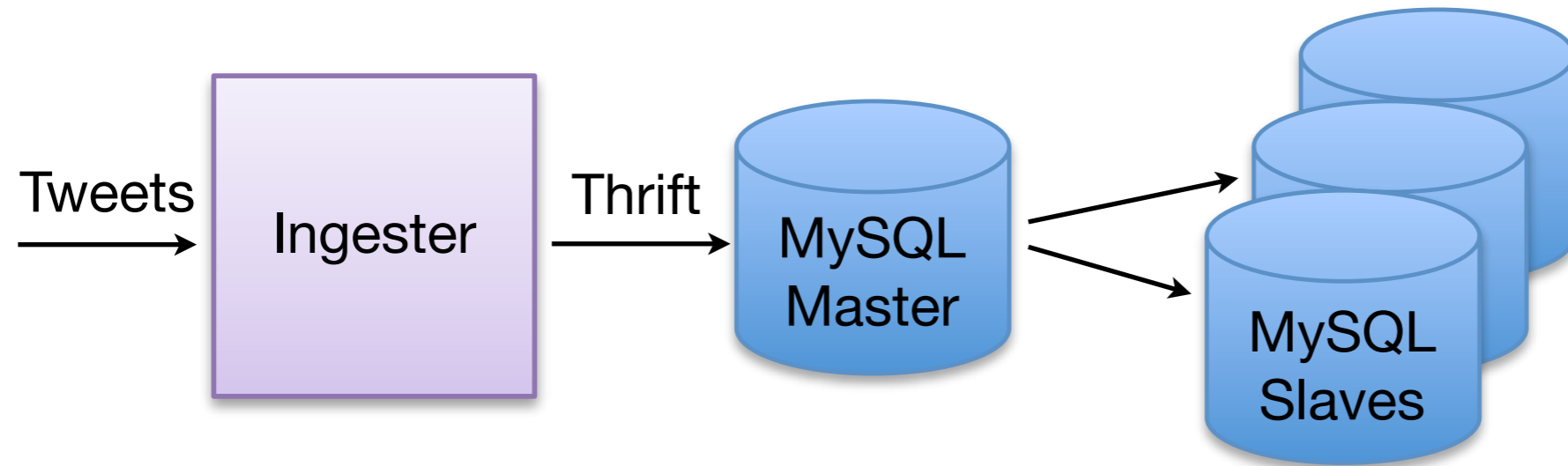
Search Architecture

Search Architecture



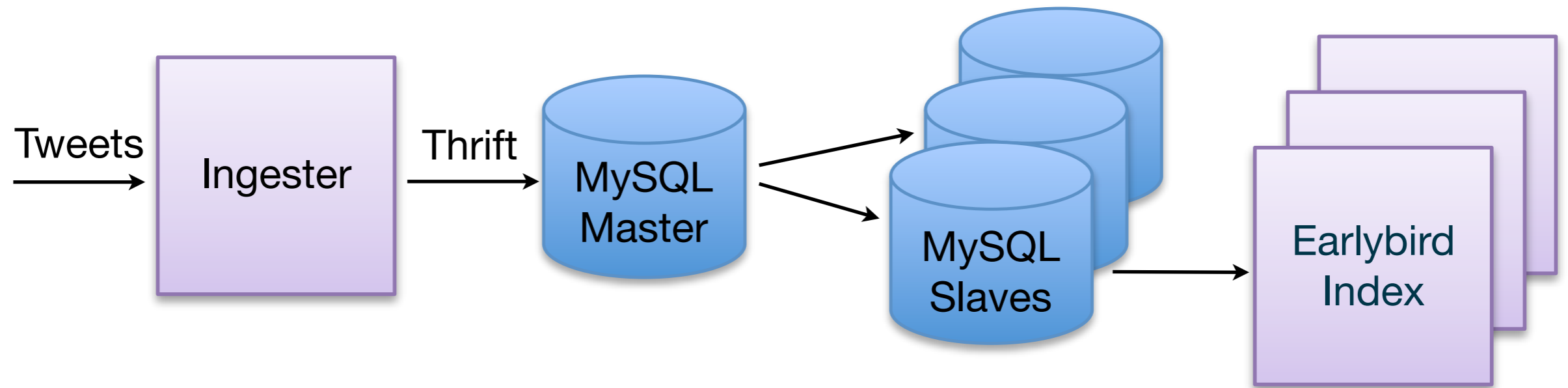
- Ingester pre-processes Tweets for search
- Geo-coding, URL expansion, tokenization, etc.

Search Architecture



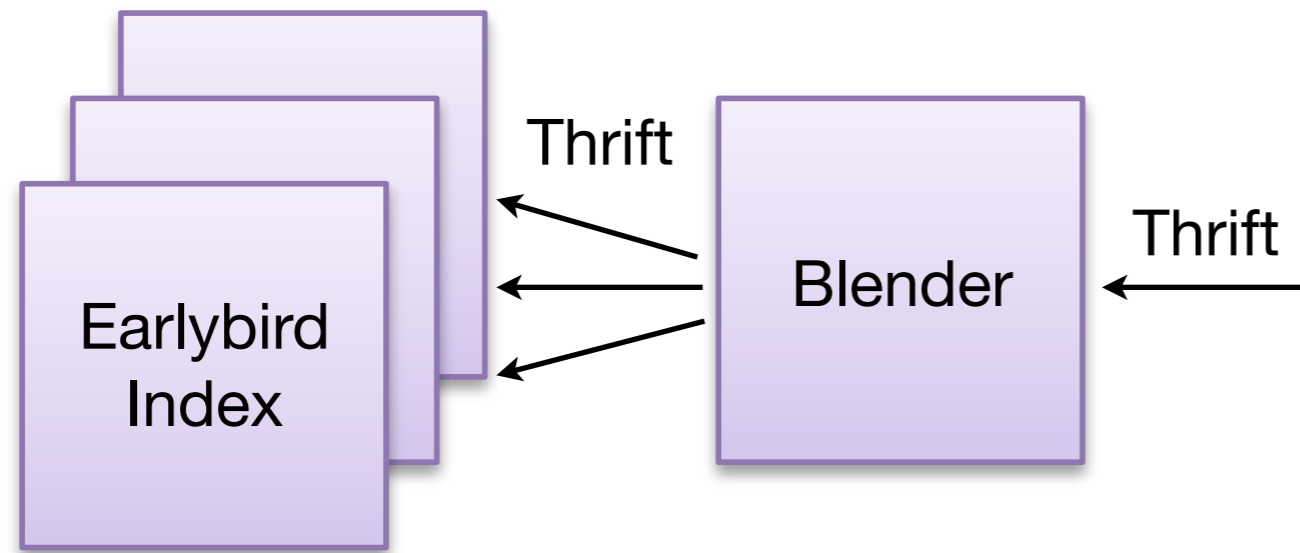
- Tweets are serialized to MySQL in Thrift format

Earlybird



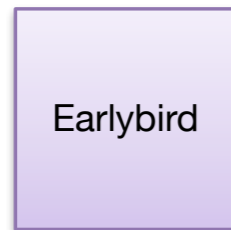
- Earlybird reads from MySQL slaves
- Builds an in-memory inverted index in real time

Blender

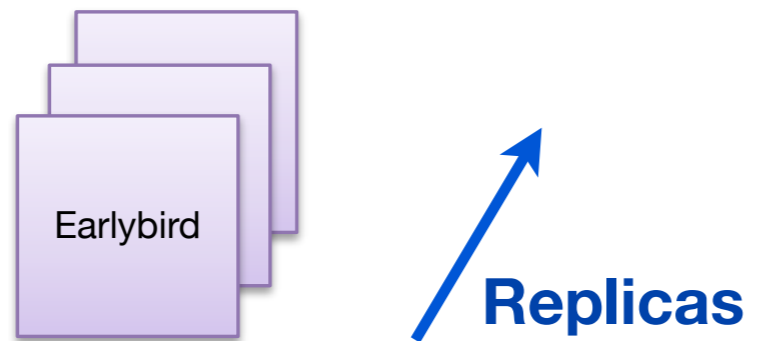


- Blender is our Thrift service aggregator
- Queries multiple Earlybirds, merges results

Cluster layout

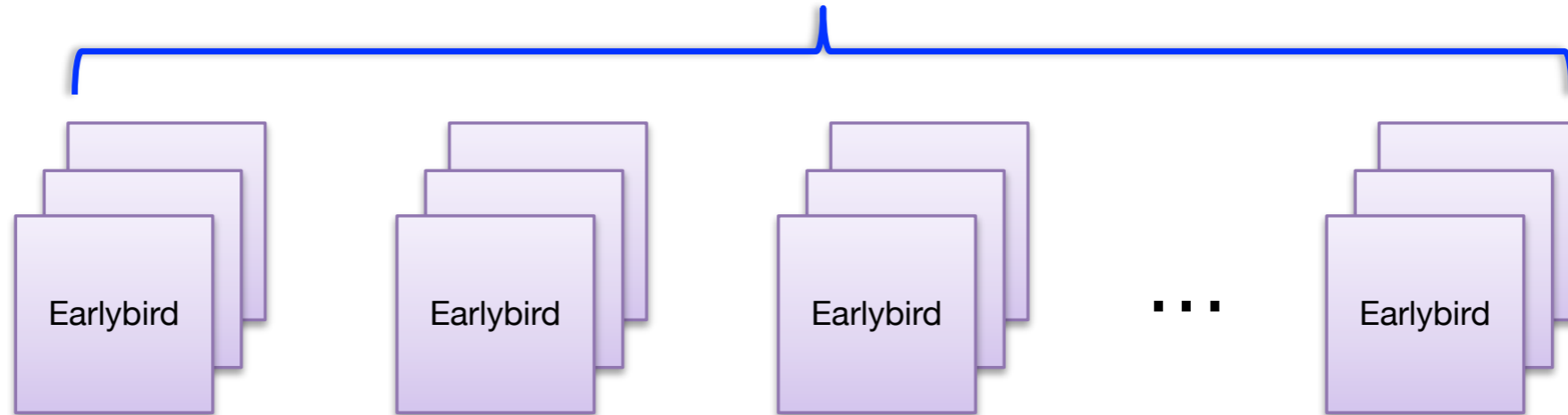


Cluster layout



Cluster layout

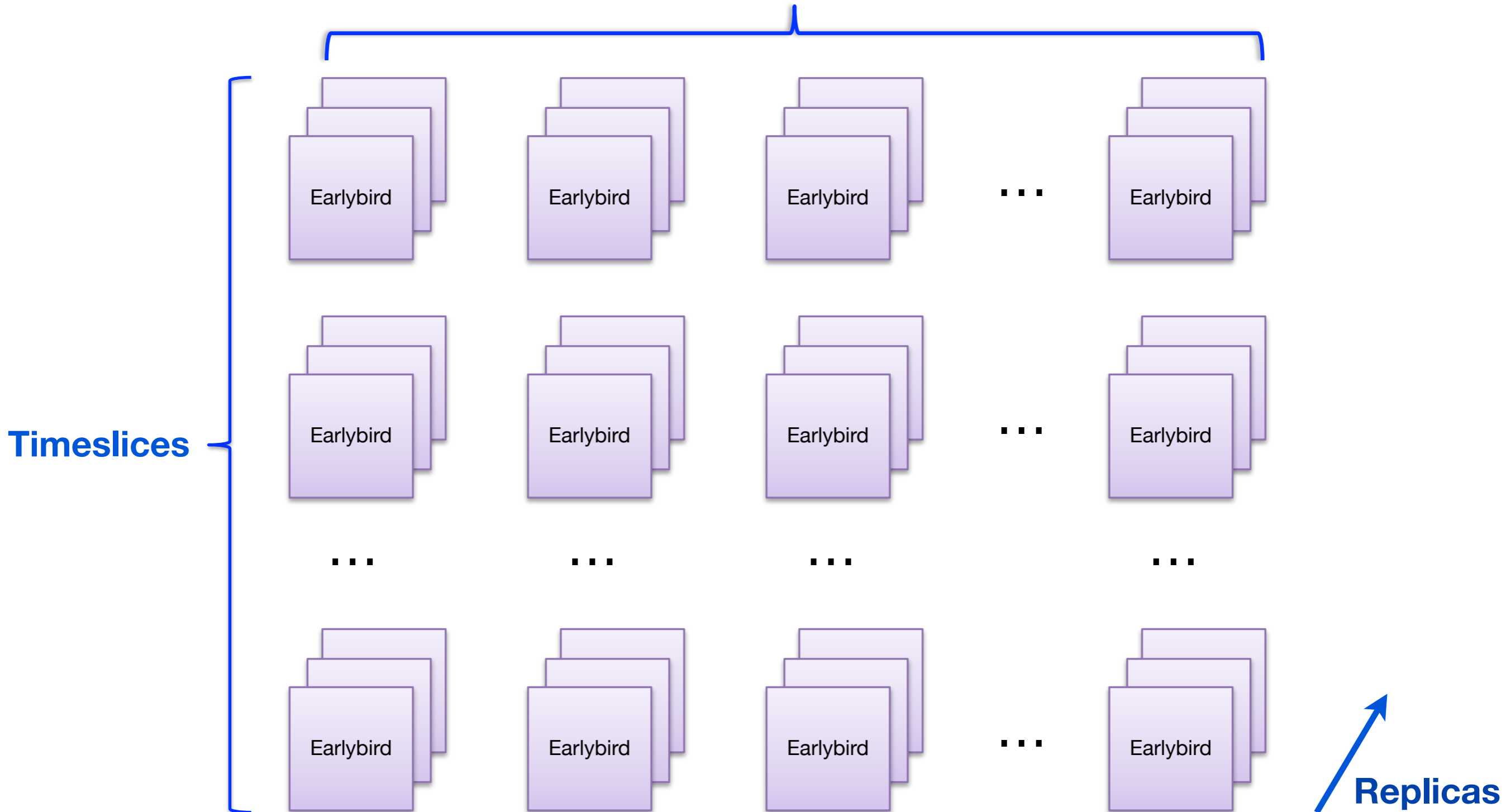
n hash partitions (docId % n)



Replicas

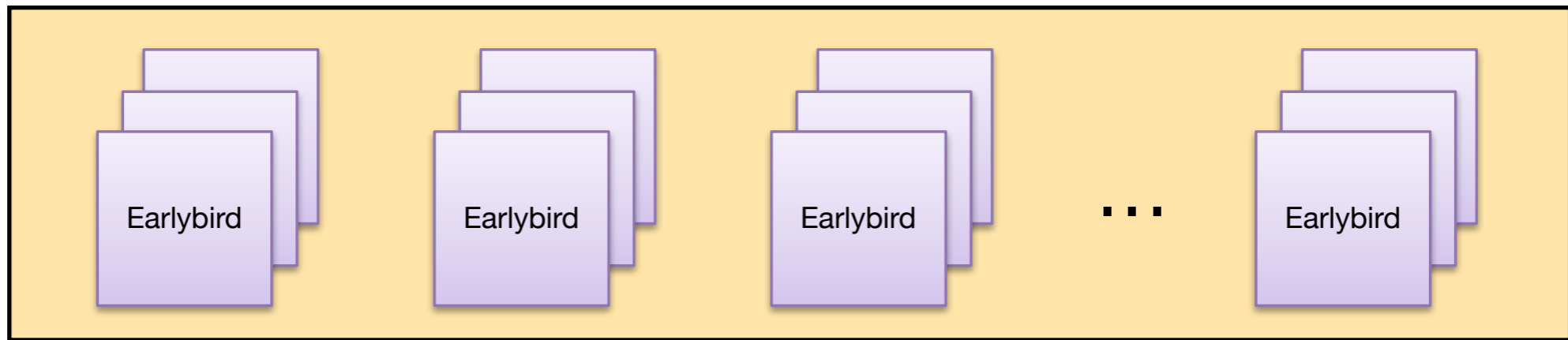
Cluster layout

n hash partitions ($\text{docId} \% n$)

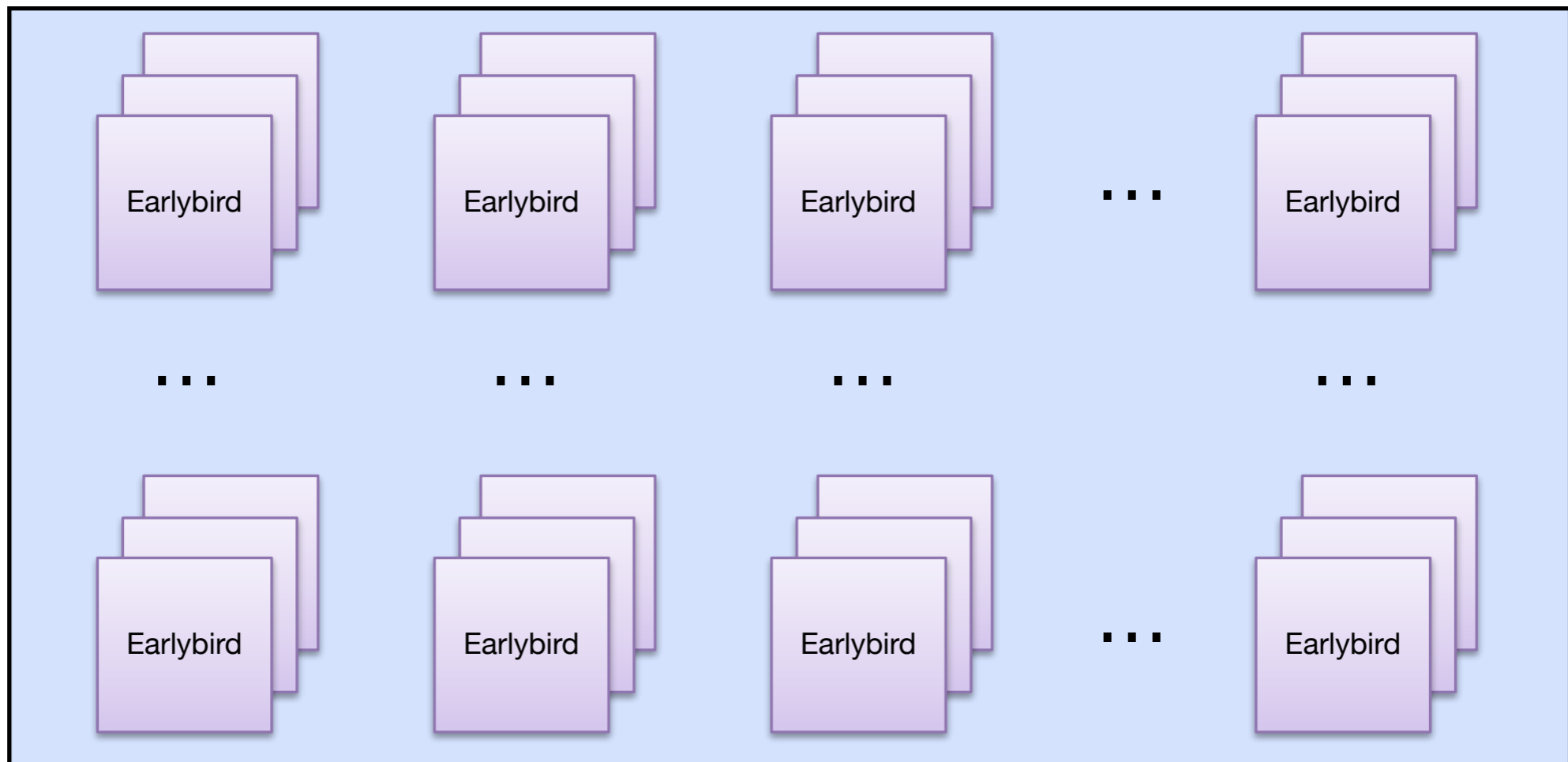


Cluster layout

**Writable
timeslice**



**Complete
timeslices**



Realtime Search @twitter

Agenda

- Introduction
- Search Architecture
- ▶ Inverted Index 101
- Memory Model & Concurrency
- What's next?

Inverted Index 101

Inverted Index 101

1	The old night keeper keeps the keep in the town
2	In the big old house in the big old gown.
3	The house in the town had the big old keep
4	Where the old night keeper never did sleep.
5	The night keeper keeps the keep in the night
6	And keeps in the dark and sleeps in the light.

Table with 6 documents

Example from:

*Justin Zobel , Alistair Moffat,
Inverted files for text search engines,
ACM Computing Surveys (CSUR)
v.38 n.2, p.6-es, 2006*

Inverted Index 101

1	The old night keeper keeps the keep in the town
2	In the big old house in the big old gown.
3	The house in the town had the big old keep
4	Where the old night keeper never did sleep.
5	The night keeper keeps the keep in the night
6	And keeps in the dark and sleeps in the light.

Table with 6 documents

term	freq	
and	1	<6>
big	2	<2> <3>
dark	1	<6>
did	1	<4>
gown	1	<2>
had	1	<3>
house	2	<2> <3>
in	5	<1> <2> <3> <5> <6>
keep	3	<1> <3> <5>
keeper	3	<1> <4> <5>
keeps	3	<1> <5> <6>
light	1	<6>
never	1	<4>
night	3	<1> <4> <5>
old	4	<1> <2> <3> <4>
sleep	1	<4>
sleeps	1	<6>
the	6	<1> <2> <3> <4> <5> <6>
town	2	<1> <3>
where	1	<4>

Dictionary and posting lists

Inverted Index 101

Query: keeper

1	The old night keeper keeps the keep in the town
2	In the big old house in the big old gown.
3	The house in the town had the big old keep
4	Where the old night keeper never did sleep.
5	The night keeper keeps the keep in the night
6	And keeps in the dark and sleeps in the light.

Table with 6 documents

term	freq	
and	1	<6>
big	2	<2> <3>
dark	1	<6>
did	1	<4>
gown	1	<2>
had	1	<3>
house	2	<2> <3>
in	5	<1> <2> <3> <5> <6>
keep	3	<1> <3> <5>
keeper	3	<1> <4> <5>
keeps	3	<1> <5> <6>
light	1	<6>
never	1	<4>
night	3	<1> <4> <5>
old	4	<1> <2> <3> <4>
sleep	1	<4>
sleeps	1	<6>
the	6	<1> <2> <3> <4> <5> <6>
town	2	<1> <3>
where	1	<4>

Dictionary and posting lists

Inverted Index 101

Query: keeper

1	The old night keeper keeps the keep in the town
2	In the big old house in the big old gown.
3	The house in the town had the big old keep
4	Where the old night keeper never did sleep.
5	The night keeper keeps the keep in the night
6	And keeps in the dark and sleeps in the light.

Table with 6 documents

term	freq	
and	1	<6>
big	2	<2> <3>
dark	1	<6>
did	1	<4>
gown	1	<2>
had	1	<3>
house	2	<2> <3>
in	5	<1> <2> <3> <5> <6>
keep	3	<1> <3> <5>
keeper	3	<1> <4> <5>
keeps	3	<1> <5> <6>
light	1	<6>
never	1	<4>
night	3	<1> <4> <5>
old	4	<1> <2> <3> <4>
sleep	1	<4>
sleeps	1	<6>
the	6	<1> <2> <3> <4> <5> <6>
town	2	<1> <3>
where	1	<4>

Dictionary and posting lists

Posting list encoding

Doc IDs to encode: 5, 15, 9000, 9002, 100000, 100090

Posting list encoding

Doc IDs to encode: 5, 15, 9000, 9002, 100000, 100090

Delta encoding:

5	10	8985	2	90998	90
---	----	------	---	-------	----

Posting list encoding

Doc IDs to encode: 5, 15, 9000, 9002, 100000, 100090

Delta encoding:

5	10	8985	2	90998	90
---	----	------	---	-------	----

Vint compression:

00000101

Values $0 \leq \text{delta} \leq 127$ need one byte

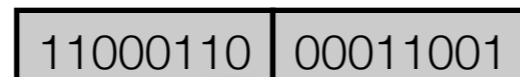
Posting list encoding

Doc IDs to encode: 5, 15, 9000, 9002, 100000, 100090

Delta encoding:



Vint compression:



Values $128 \leq \text{delta} \leq 16384$
need two bytes

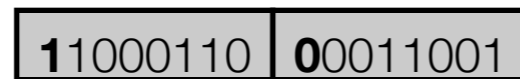
Posting list encoding

Doc IDs to encode: 5, 15, 9000, 9002, 100000, 100090

Delta encoding:



Vint compression:



First bit indicates whether next byte belongs to the same value

Posting list encoding

Doc IDs to encode: 5, 15, 9000, 9002, 100000, 100090

Delta encoding:

5	10	8985	2	90998	90
---	----	------	---	-------	----

VInt compression:

11000110	00011001
----------	----------

- Variable number of bytes - a VInt-encoded posting can not be written as a primitive Java type; therefore it can not be written atomically

Posting list encoding

Doc IDs to encode: 5, 15, 9000, 9002, 100000, 100090

Delta encoding:

5	10	8985	2	90998	90
---	----	------	---	-------	----



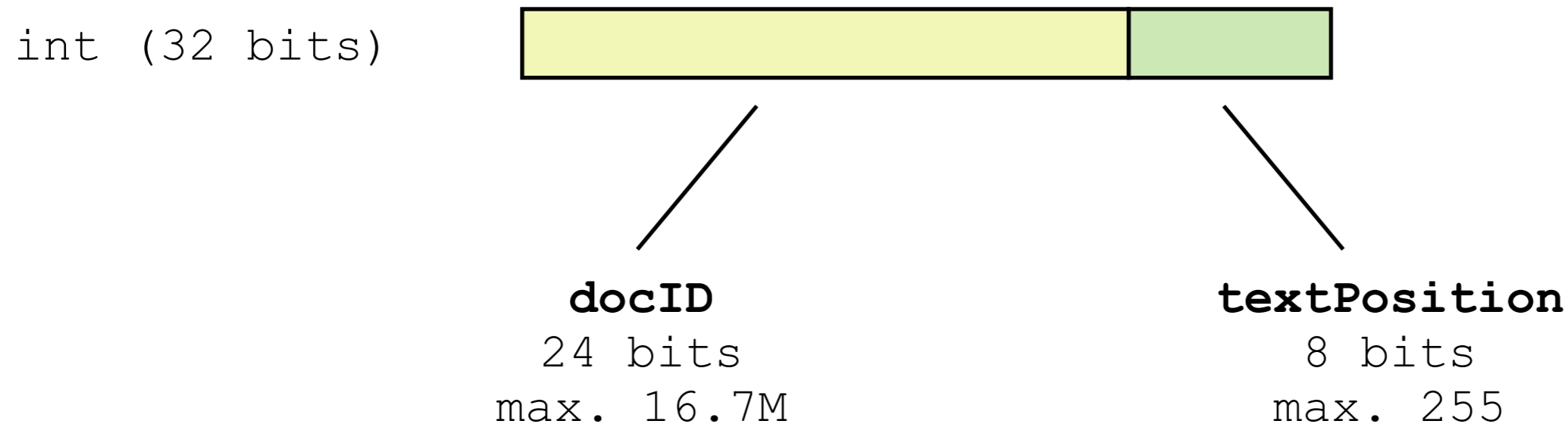
Read direction

- Each posting depends on previous one; decoding only possible in old-to-new direction
- With recency ranking (new-to-old) no early termination is possible

Posting list encoding

- By default Lucene uses a combination of delta encoding and VInt compression
- VInts are expensive to decode
- Problem 1: How to traverse posting lists backwards?
- Problem 2: How to write a posting atomically?

Posting list encoding in Earlybird

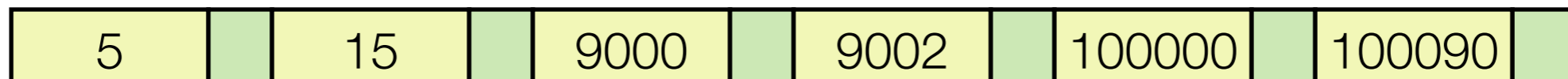


- Tweet text can only have 140 chars
- Decoding speed significantly improved compared to delta and VInt decoding (early experiments suggest 5x improvement compared to vanilla Lucene with FSDirectory)

Posting list encoding in Earlybird

Doc IDs to encode: 5, 15, 9000, 9002, 100000, 100090

Earlybird encoding:

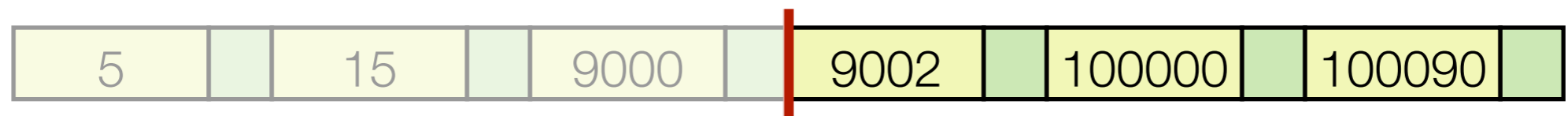


Read direction

Early query termination

Doc IDs to encode: 5, 15, 9000, 9002, 100000, 100090

Earlybird encoding:



Read direction

E.g. 3 result are requested: Here we can terminate after reading 3 postings

Posting list encoding - Summary

- ints can be written atomically in Java
- Backwards traversal easy on absolute docIDs (not deltas)
- Every posting is a possible entry point for a searcher
- Skipping can be done without additional data structures as binary search, even though there are better approaches which should be explored
- On tweet indexes we need about 30% more storage for docIDs compared to delta+Vints; compensated by compression of complete segments
- Max. segment size: $2^{24} = 16.7\text{M}$ tweets

Realtime Search @twitter

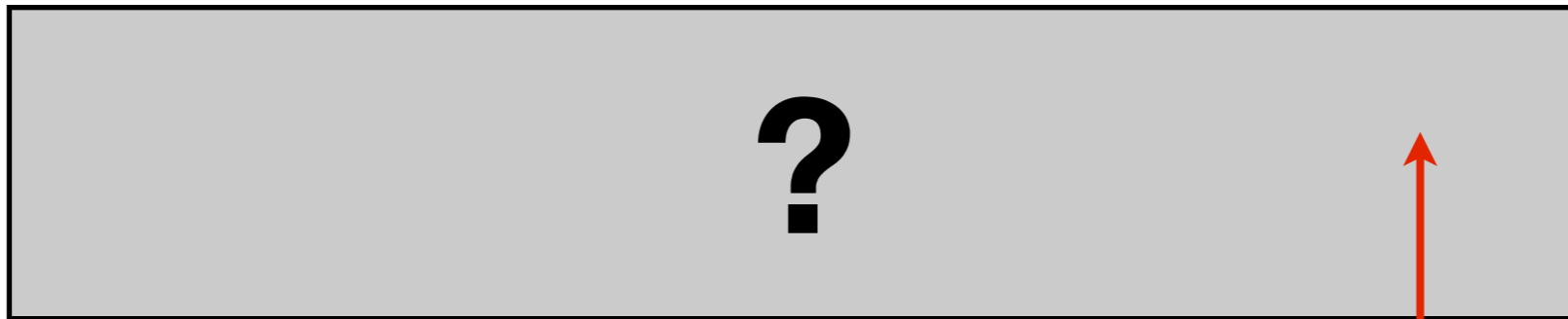
Agenda

- Introduction
- Search Architecture
- Inverted Index 101
- ▶ Memory Model & Concurrency
- What's next?

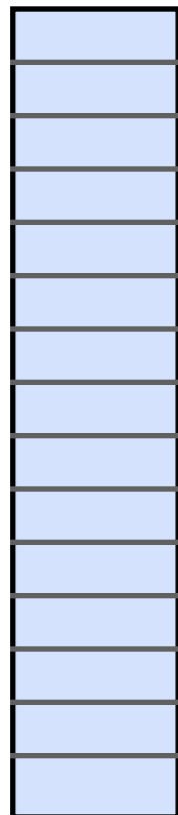
Memory Model & Concurrency

Inverted index components

Posting list storage

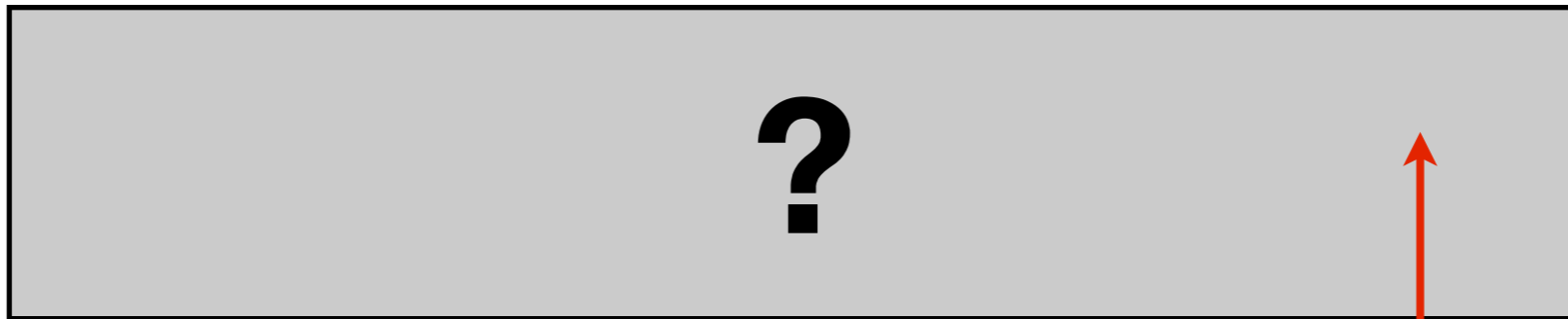


Dictionary



Inverted index components

Posting list storage



Dictionary



Inverted Index

1	The old night keeper keeps the keep in the town
2	In the big old house in the big old gown.
3	The house in the town had the big old keep
4	Where the old night keeper never did sleep.
5	The night keeper keeps the keep in the night
6	And keeps in the dark and sleeps in the light.

Table with 6 documents

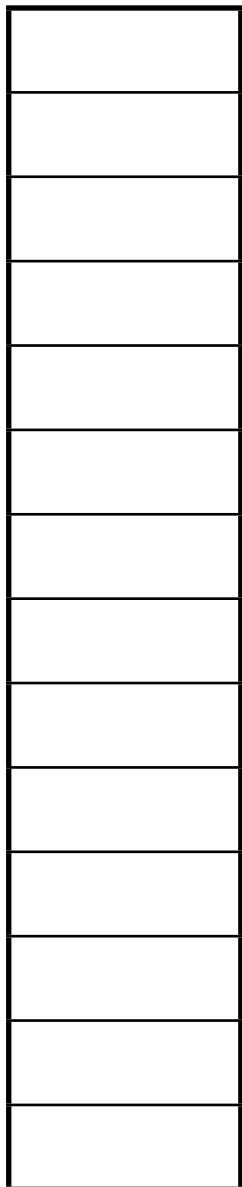
term	freq	
and	1	<6>
big	2	<2> <3>
dark	1	<6>
did	1	<4>
gown	1	<2>
had	1	<3>
house	2	<2> <3>
in	5	<1> <2> <3> <5> <6>
keep	3	<1> <3> <5>
keeper	3	<1> <4> <5>
keeps	3	<1> <5> <6>
light	1	<6>
never	1	<4>
night	3	<1> <4> <5>
old	4	<1> <2> <3> <4>
sleep	1	<4>
sleeps	1	<6>
the	6	<1> <2> <3> <4> <5> <6>
town	2	<1> <3>
where	1	<4>

Per term we store different kinds of metadata: text pointer, frequency, postings pointer, etc.

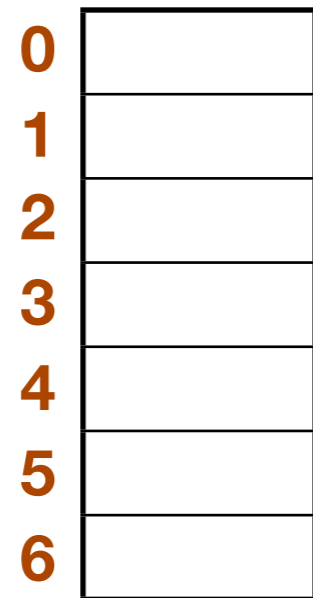
Dictionary and posting lists

Term dictionary

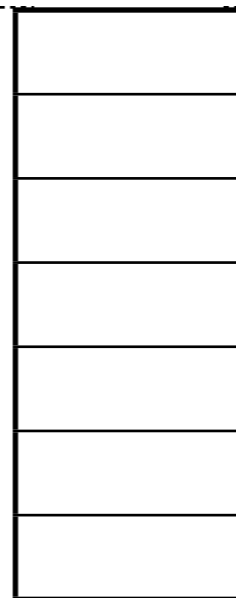
`termID`
`int[]`



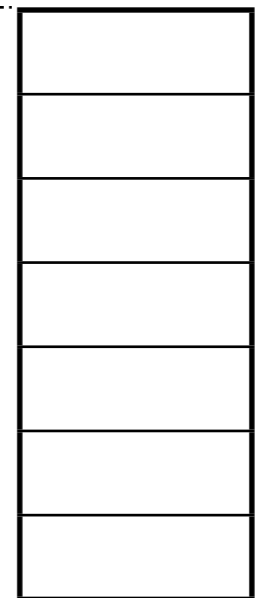
`textPointer;`
`int[]`



`postingsPointer;`
`int[]`

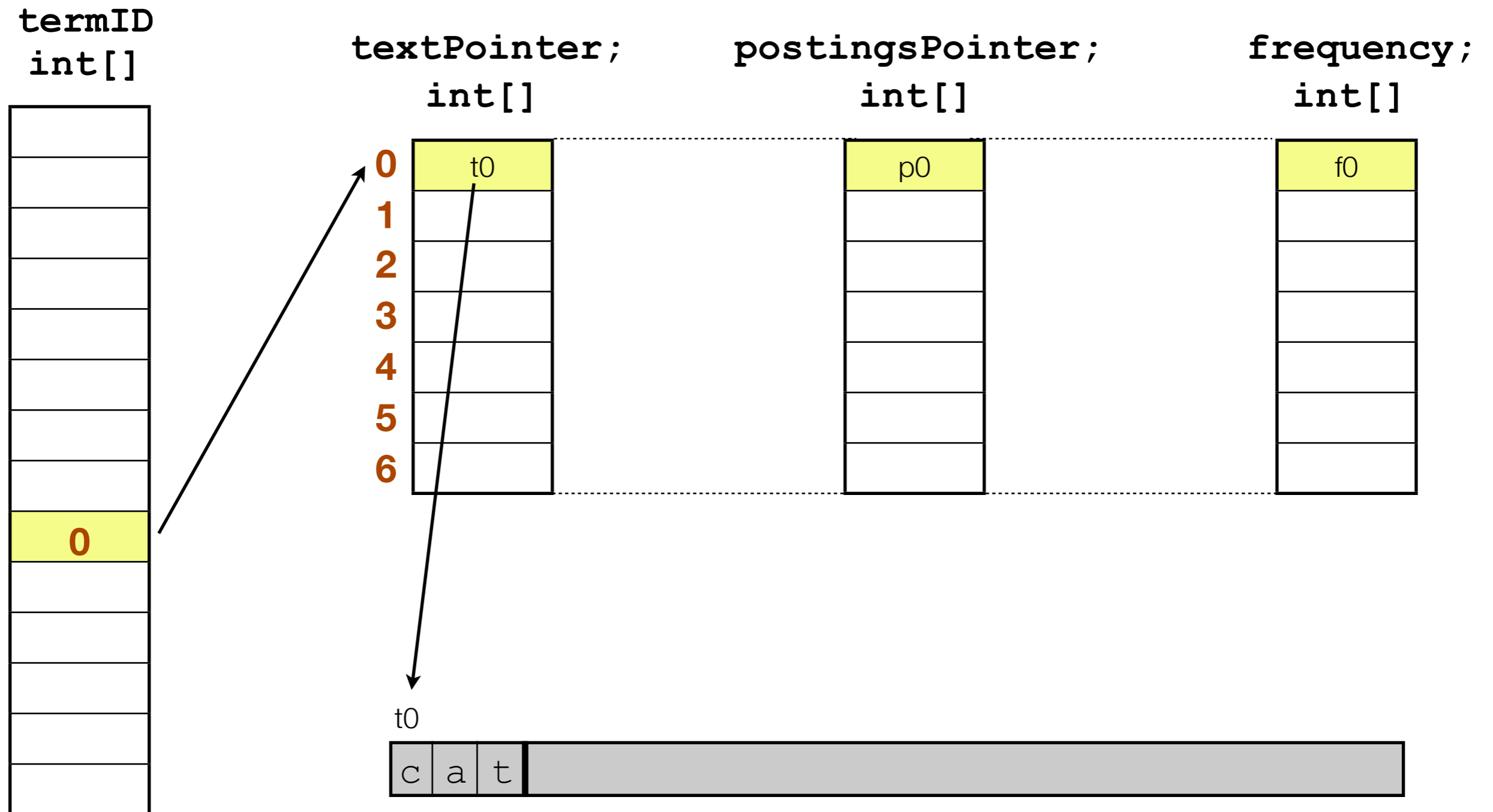


`frequency;`
`int[]`

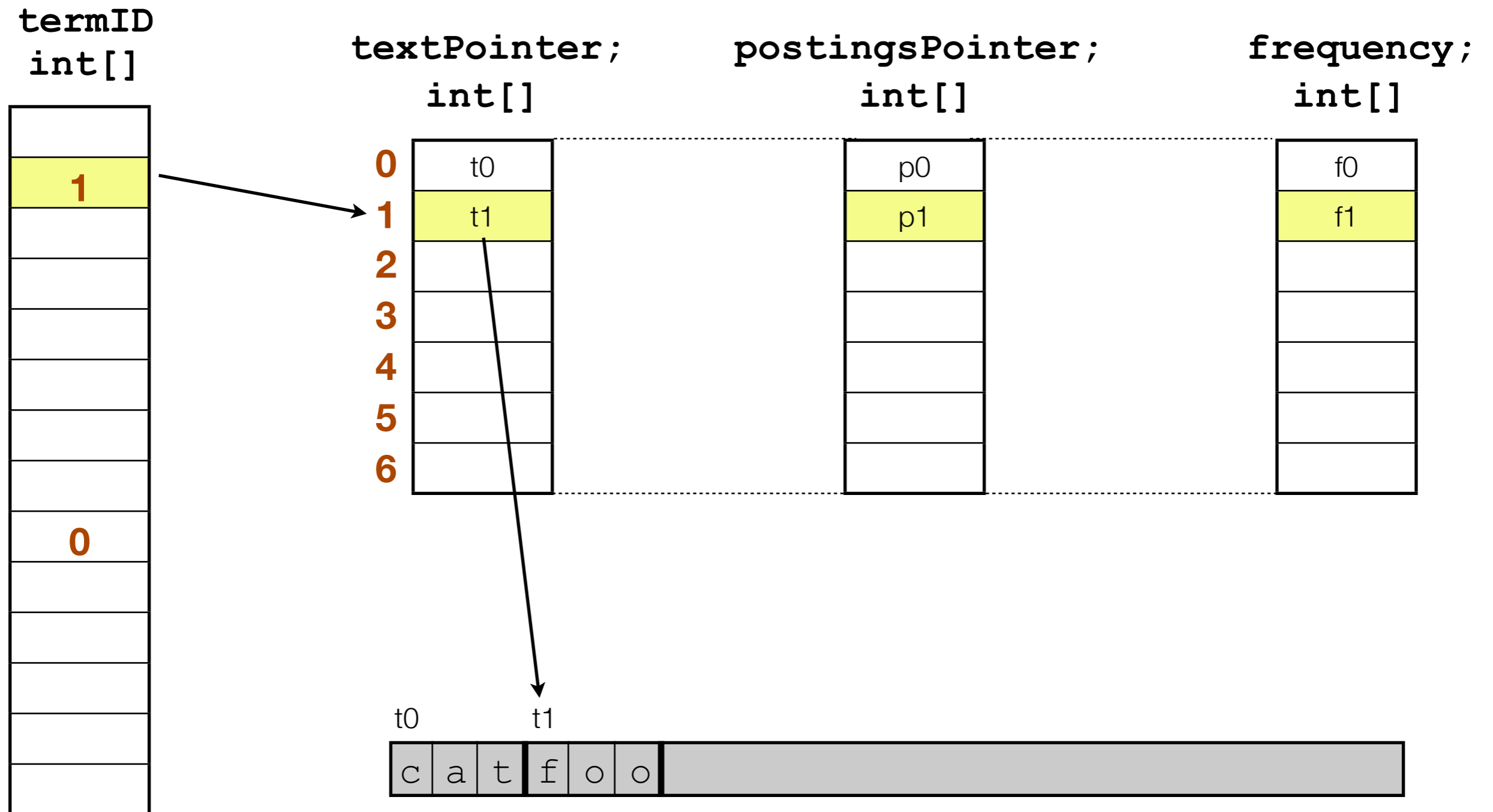


`term text pool`

Term dictionary



Term dictionary

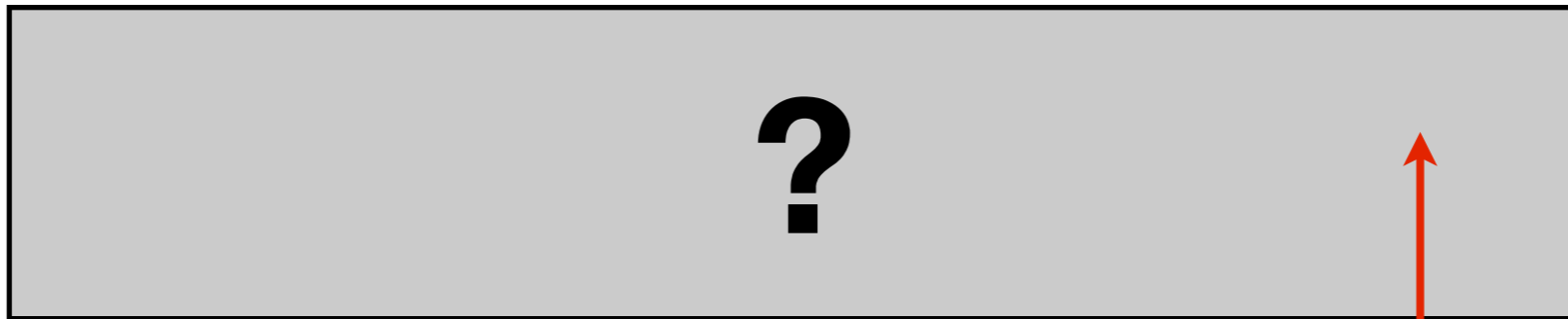


Term dictionary

- Number of objects \ll number of terms
- $O(1)$ lookups
- Easy to store more term metadata by adding additional parallel arrays

Inverted index components

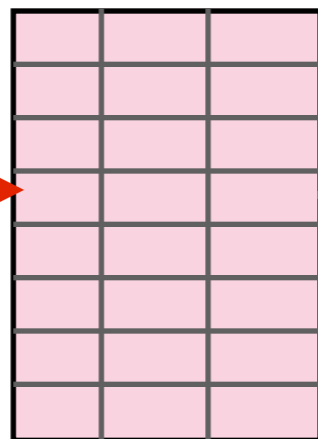
Posting list storage



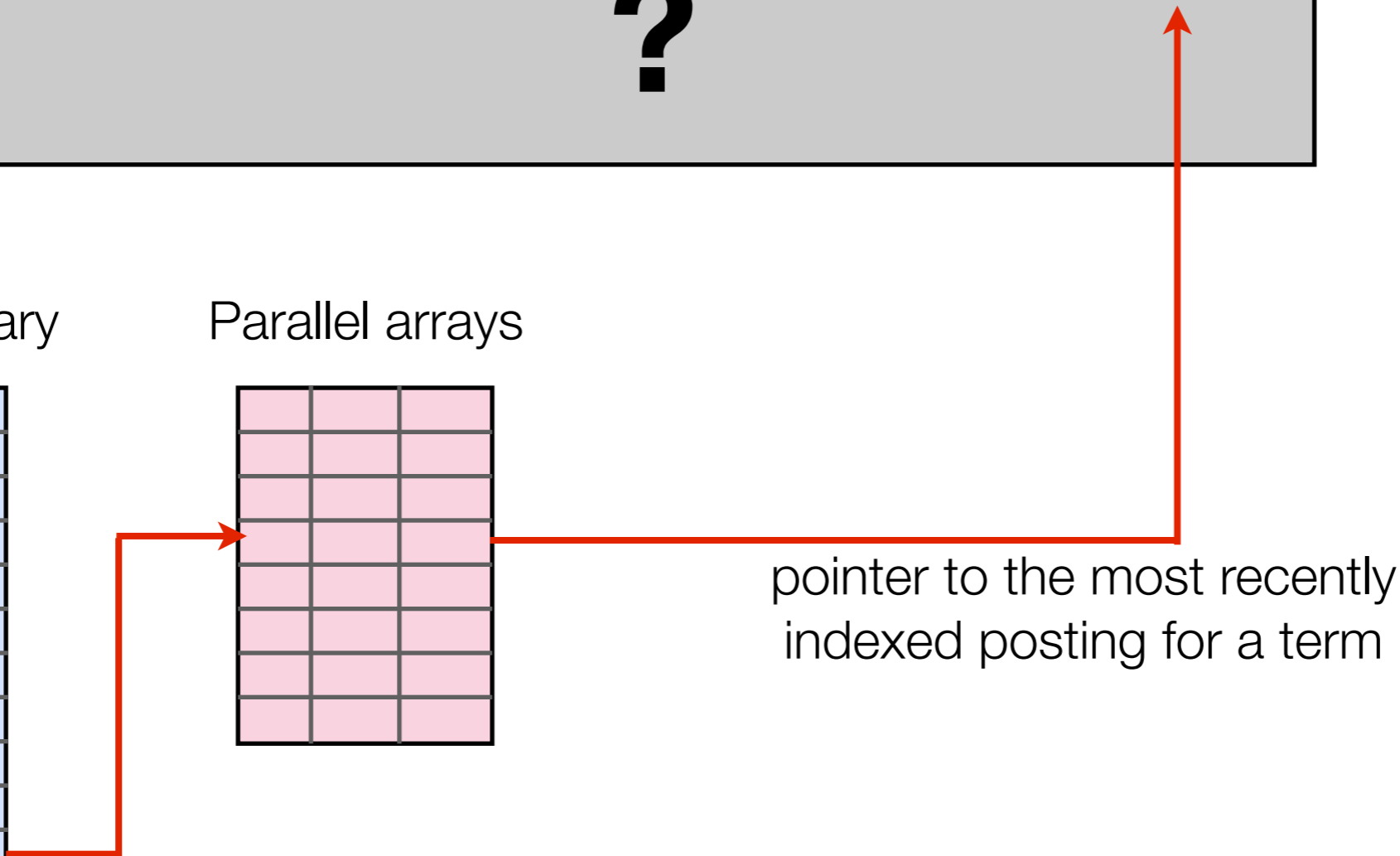
Dictionary



Parallel arrays



pointer to the most recently indexed posting for a term

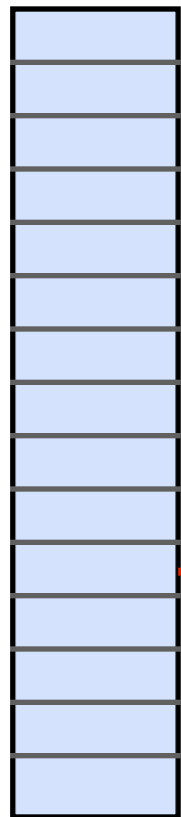


Inverted index components

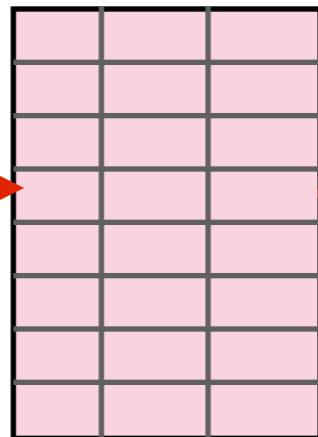
Posting list storage



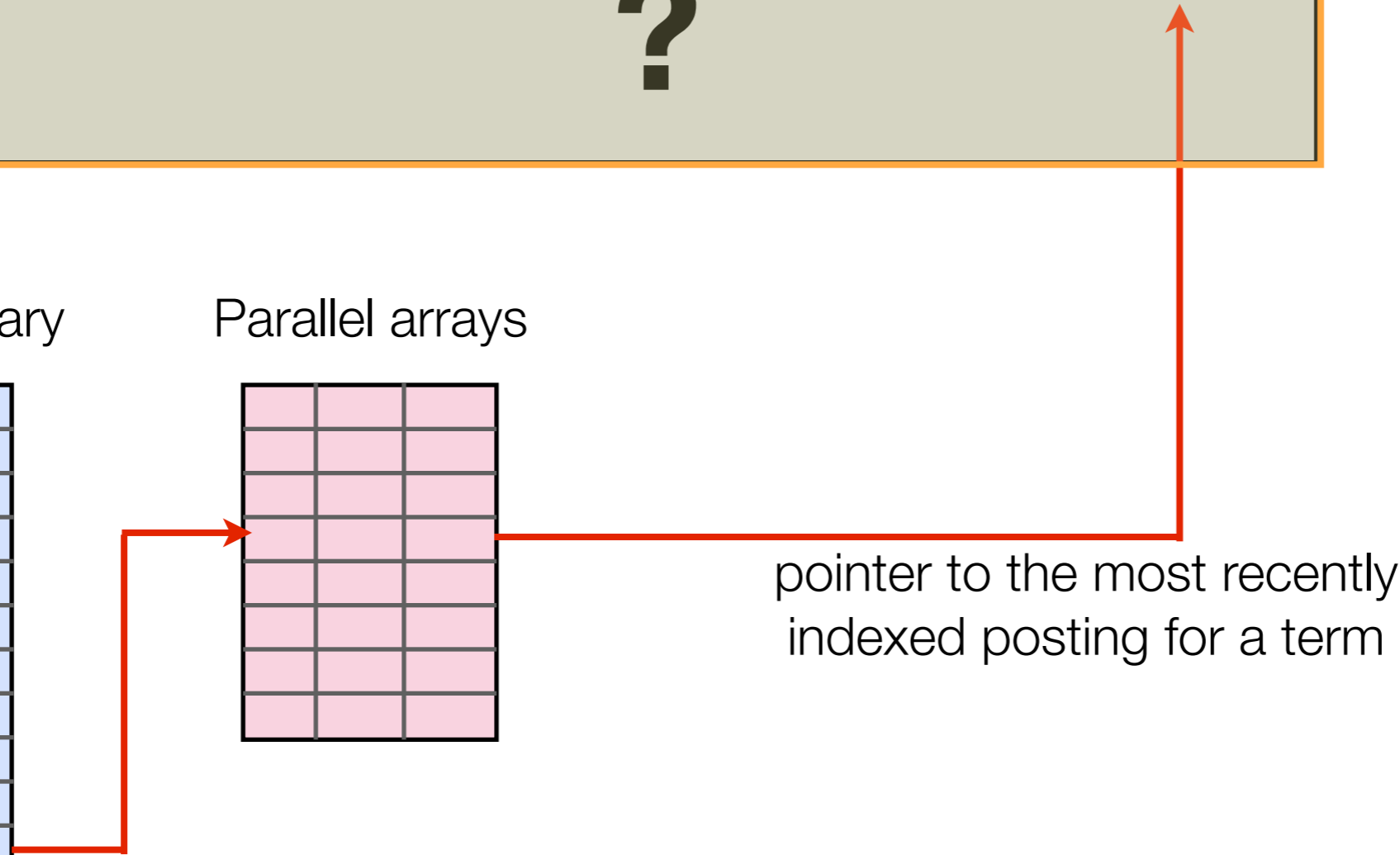
Dictionary



Parallel arrays



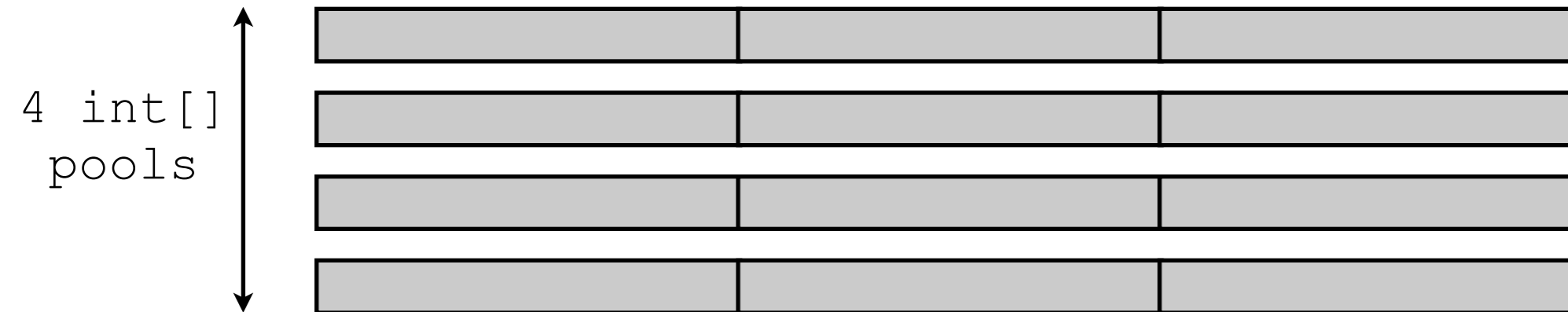
pointer to the most recently indexed posting for a term



Posting lists storage - Objectives

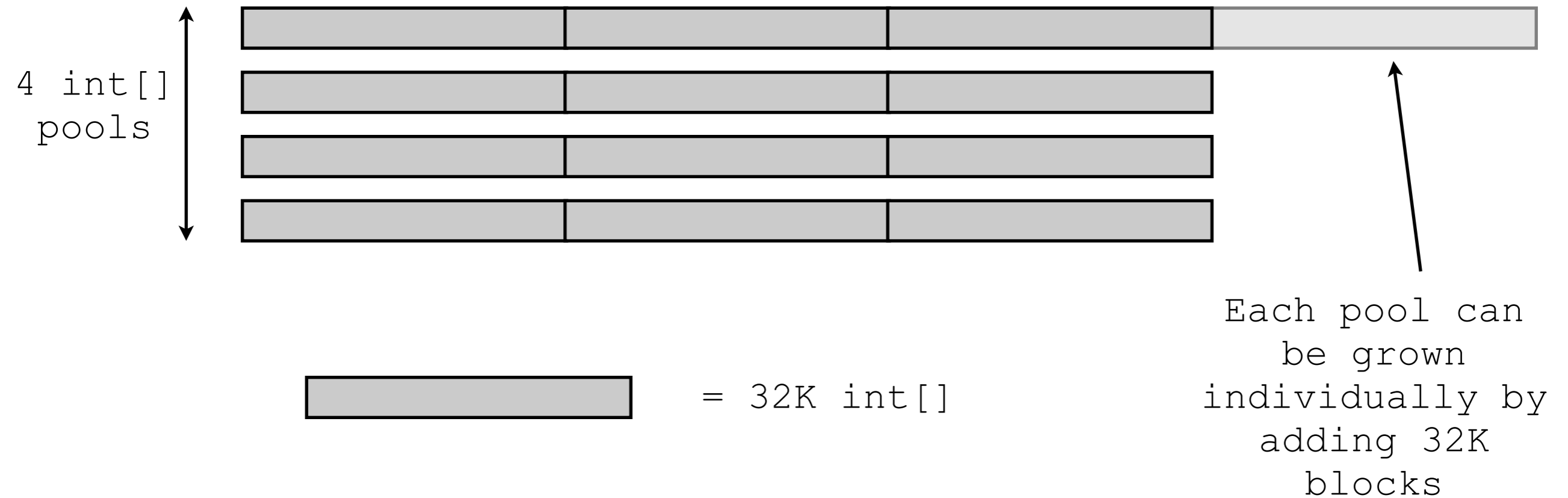
- Store many single-linked lists of different lengths space-efficiently
- The number of java objects should be independent of the number of lists or number of items in the lists
- Every item should be a possible entry point into the lists for iterators, i.e. items should not be dependent on other items (e.g. no delta encoding)
- Append and read possible by multiple threads in a lock-free fashion (single append thread, multiple reader threads)
- Traversal in backwards order

Memory management

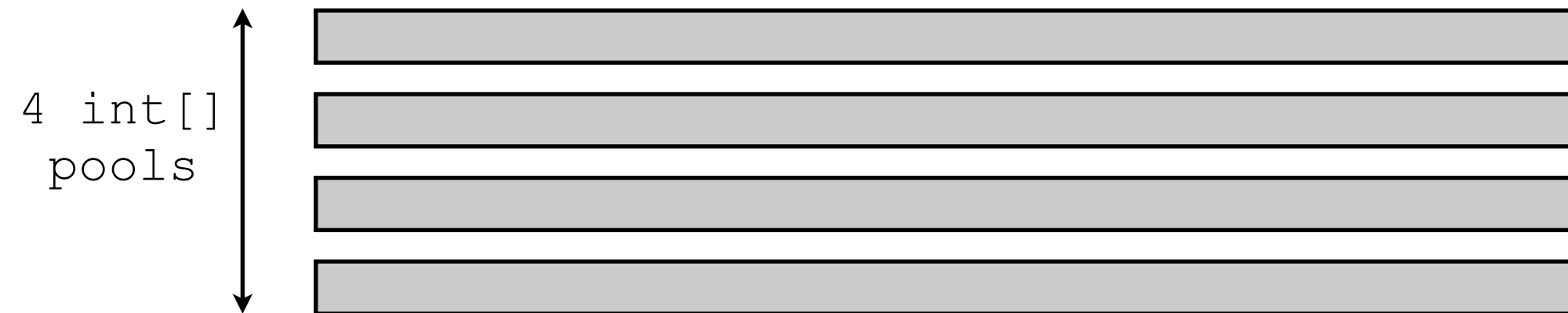


 = 32K int[]

Memory management



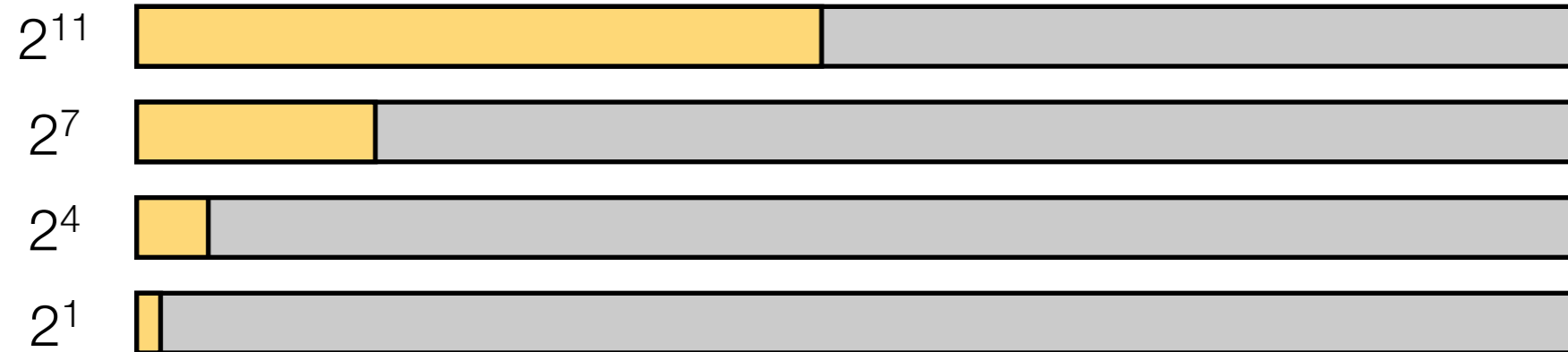
Memory management



- For simplicity we can forget about the blocks for now and think of the pools as continuous, unbounded `int[]` arrays
- Small total number of Java objects (each 32K block is one object)

Memory management

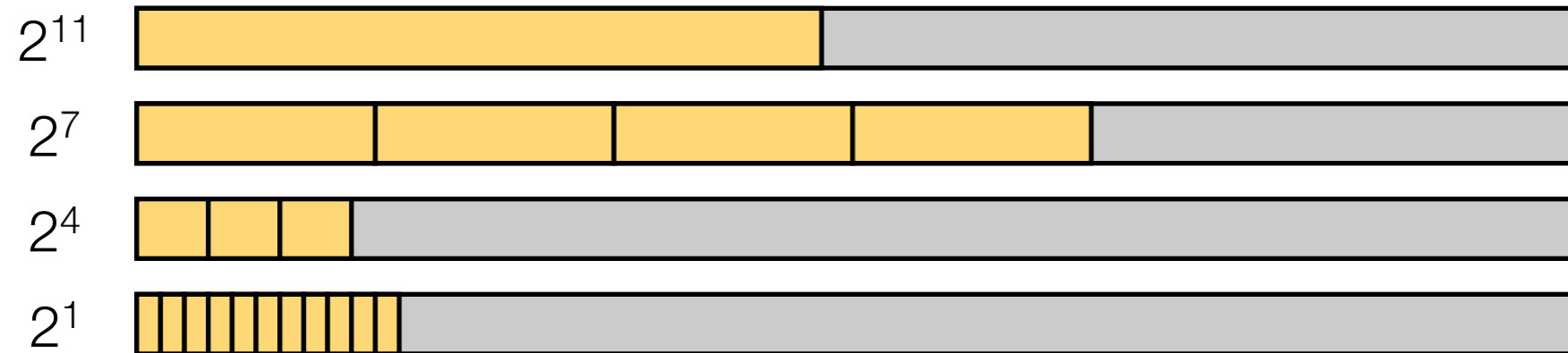
slice size



- Slices can be allocated in each pool
- Each pool has a different, but fixed slice size

Adding and appending to a list

slice size



Adding and appending to a list

slice size

2^{11}



2^7



2^4



2^1

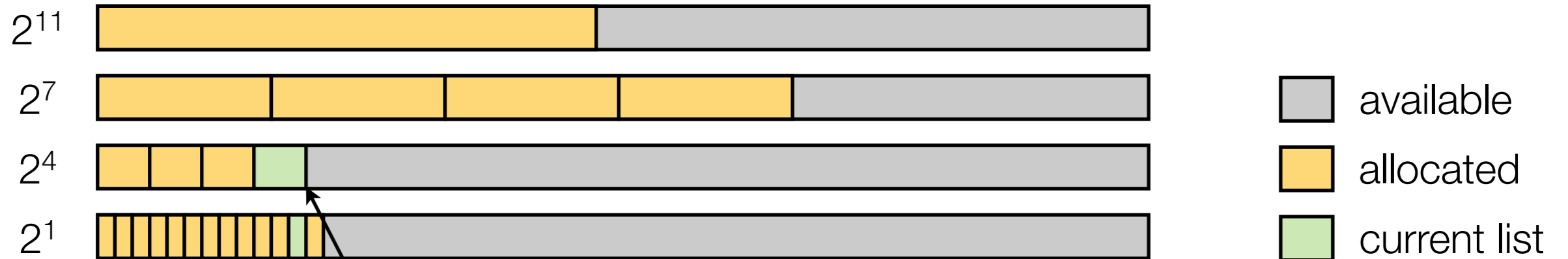


Store first two
postings in this slice



Adding and appending to a list

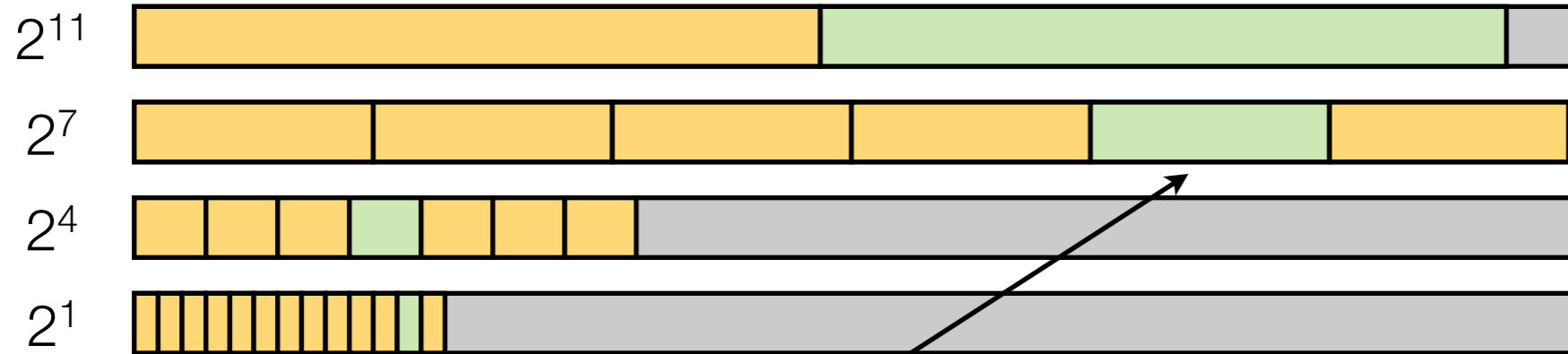
slice size



When first slice is full, allocate another one in second pool

Adding and appending to a list

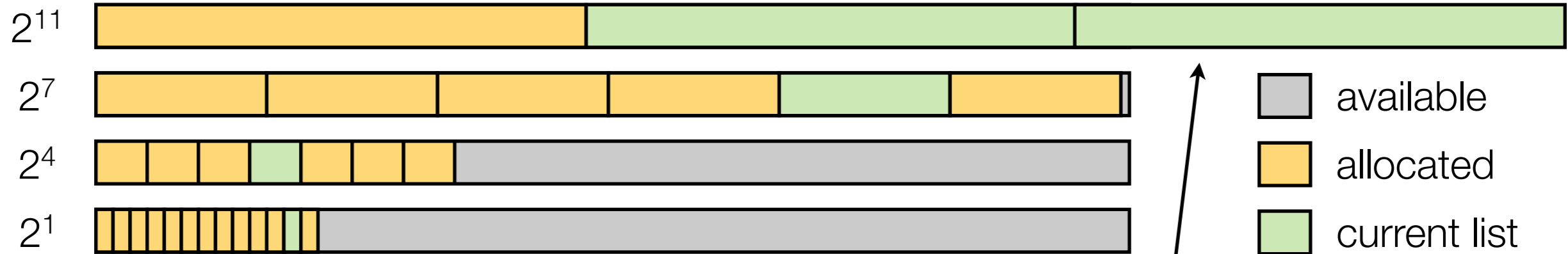
slice size



Allocate a slice on each level as list grows

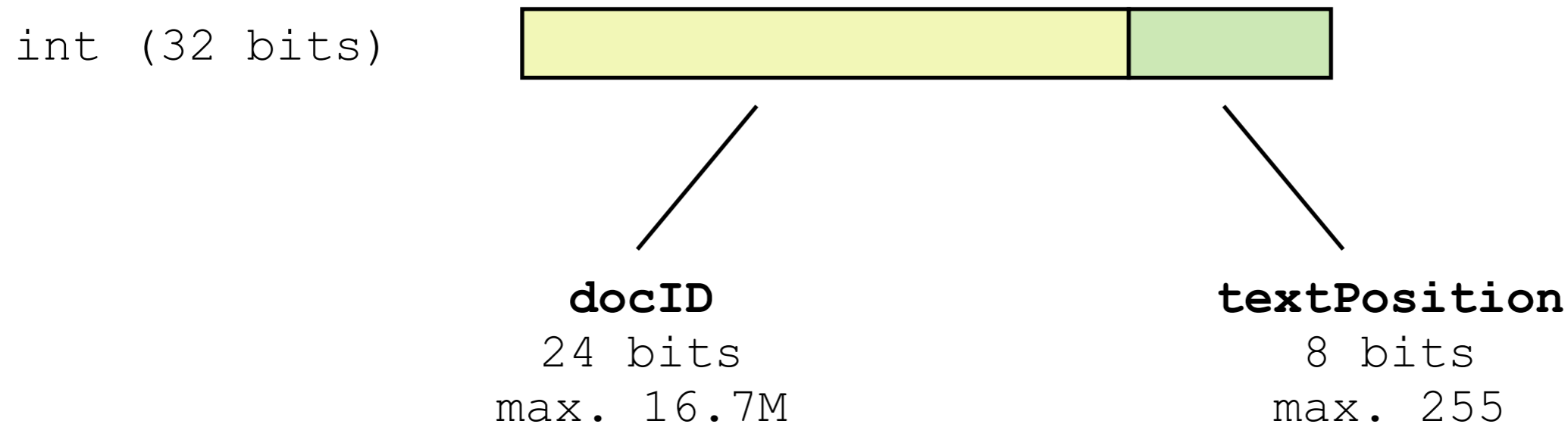
Adding and appending to a list

slice size



On upper most level one list can own multiple slices

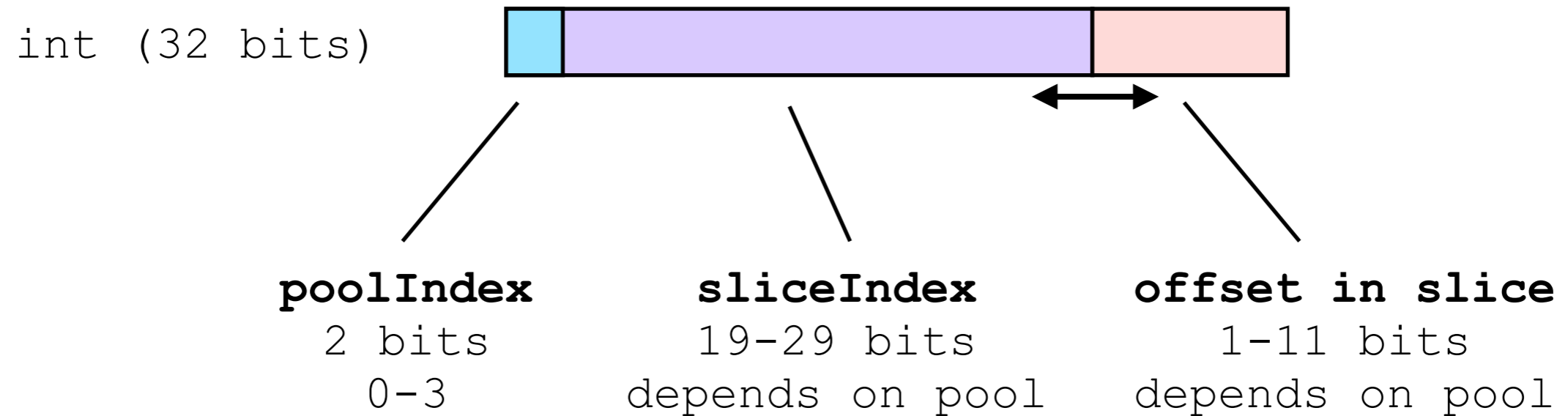
Posting list format



- Tweet text can only have 140 chars
- Decoding speed significantly improved compared to delta and VInt decoding (early experiments suggest 5x improvement compared to vanilla Lucene with FSDirectory)

Addressing items

- Use 32 bit (int) pointers to address any item in any list unambiguously:



- Nice symmetry: Postings and address pointers both fit into a 32 bit int

Linking the slices

slice size

2^{11}



2^7



2^4



2^1



-  available
-  allocated
-  current list

Linking the slices

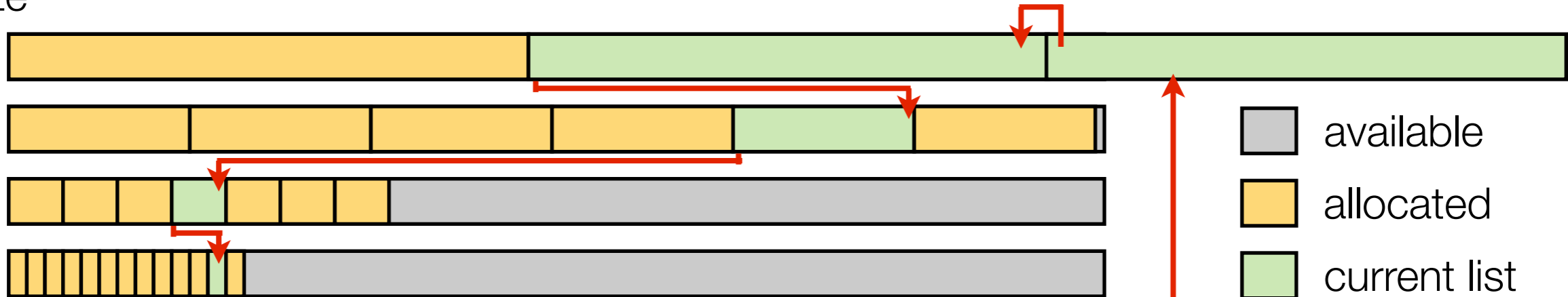
slice size

2^{11}

2^7

2^4

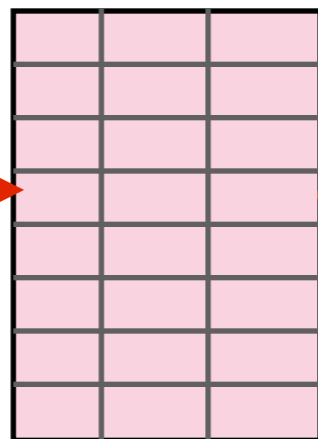
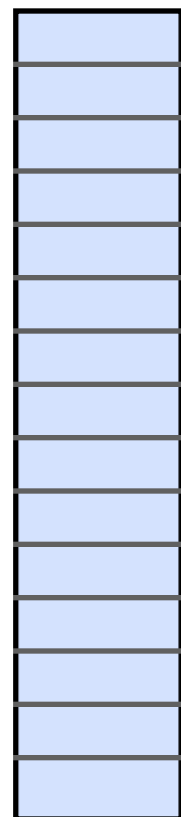
2^1



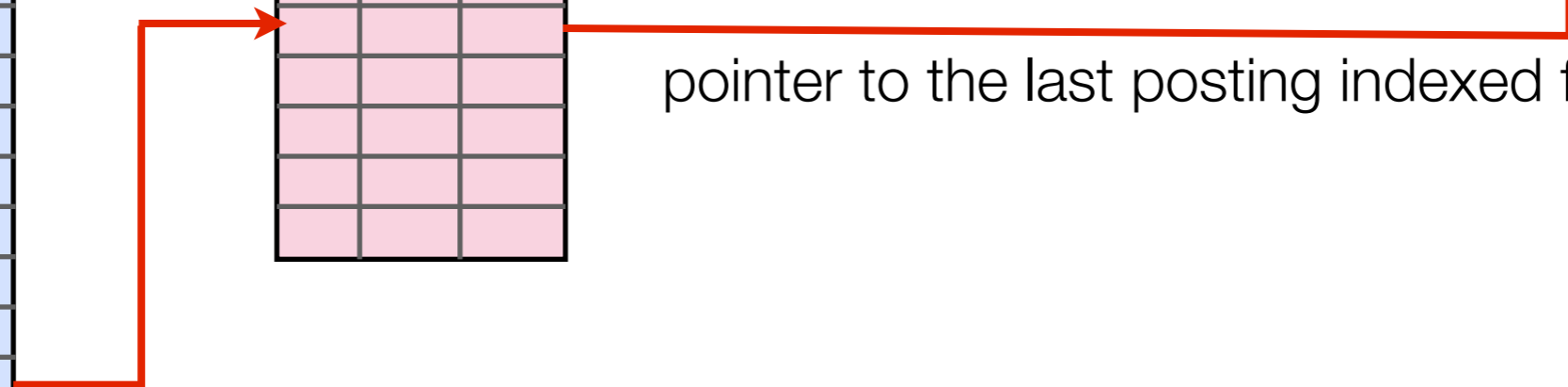
- available
- allocated
- current list

Dictionary

Parallel arrays



pointer to the last posting indexed for a term



Concurrency - Definitions

- **Pessimistic locking**

- A thread holds an exclusive lock on a resource, while an action is performed [mutual exclusion]
- Usually used when conflicts are expected to be likely

- **Optimistic locking**

- Operations are tried to be performed atomically without holding a lock; conflicts can be detected; retry logic is often used in case of conflicts
- Usually used when conflicts are expected to be the exception

Concurrency - Definitions

- **Non-blocking algorithm**

Ensures, that threads competing for shared resources do not have their execution indefinitely postponed by mutual exclusion.

- **Lock-free algorithm**

A non-blocking algorithm is lock-free if there is guaranteed system-wide progress.

- **Wait-free algorithm**

A non-blocking algorithm is wait-free, if there is guaranteed per-thread progress.

Concurrency

- Having a single writer thread simplifies our problem: no locks have to be used to protect data structures from corruption (only one thread modifies data)
- But: we have to make sure that all readers **always** see a consistent state of **all** data structures -> this is much harder than it sounds!
- In Java, it is not guaranteed that one thread will see changes that another thread makes in program execution order, unless the same memory barrier is crossed by both threads -> **safe publication**
- Safe publication can be achieved in different, subtle ways. Read the great book “Java concurrency in practice” by Brian Goetz for more information!

Java Memory Model

- **Program order rule**

Each action in a thread *happens-before* every action in that thread that comes later in the program order.

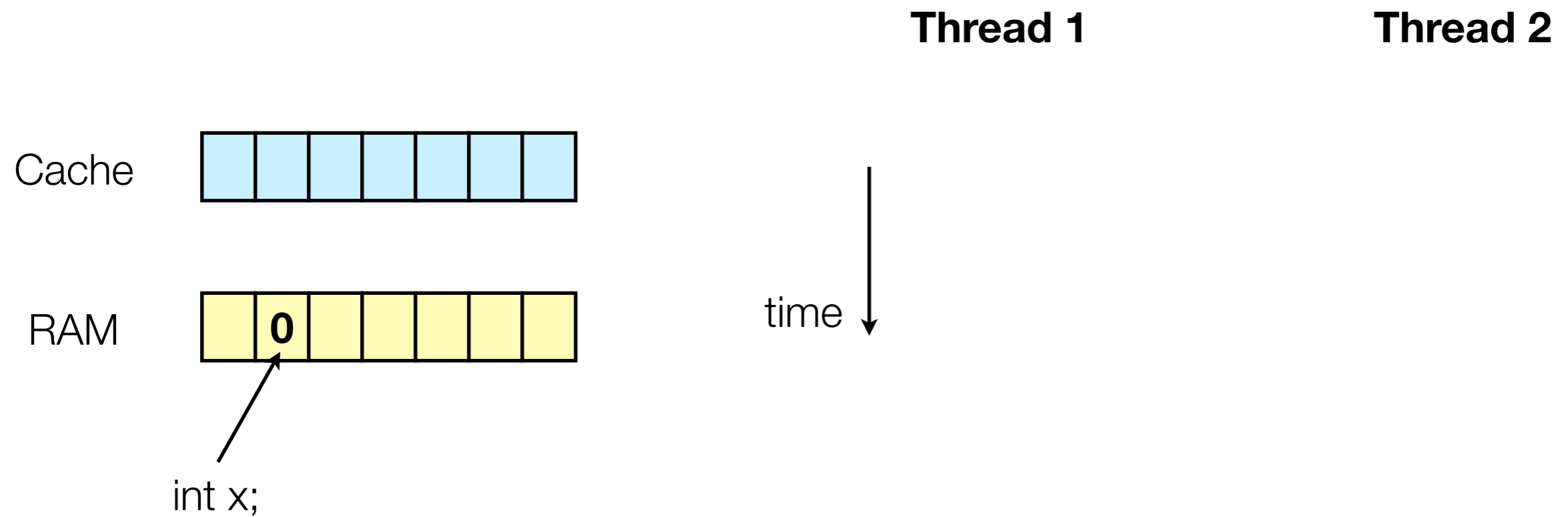
- **Volatile variable rule**

A write to a volatile field *happens-before* every subsequent read of that same field.

- **Transitivity**

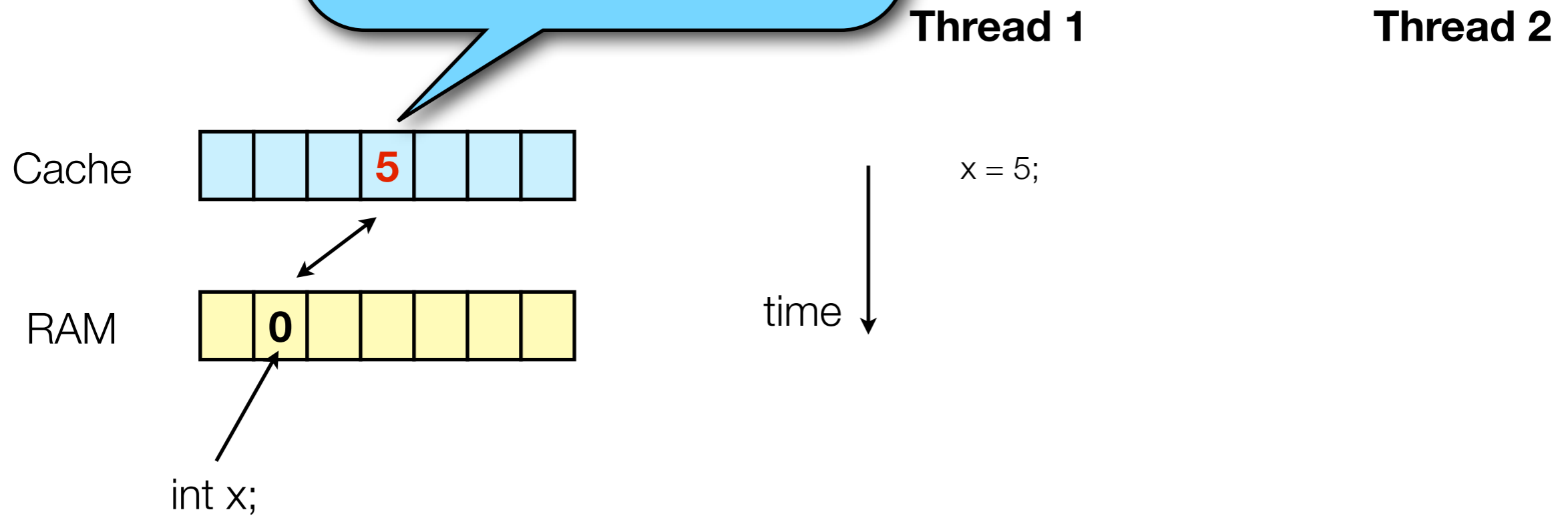
If A *happens-before* B, and B *happens-before* C, then A *happens-before* C.

Concurrency

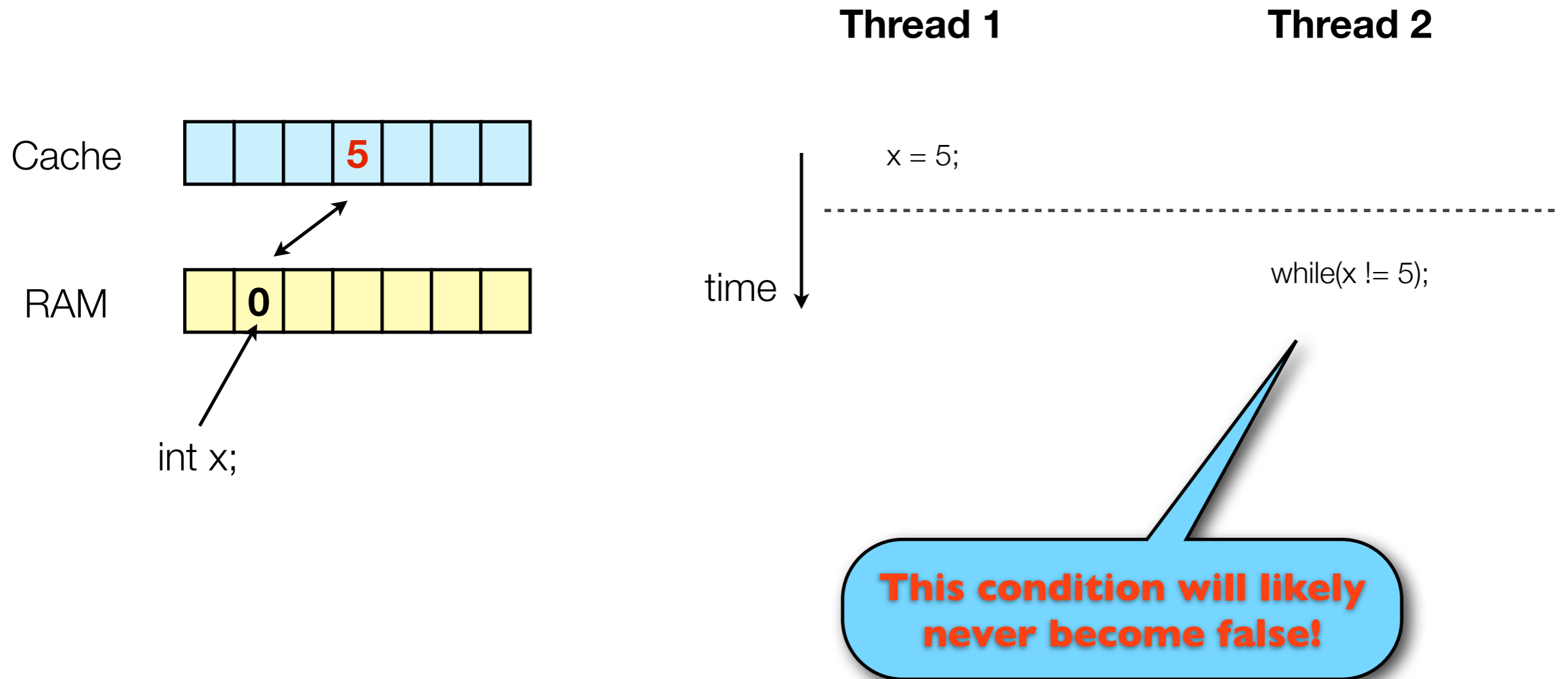


Concurrency

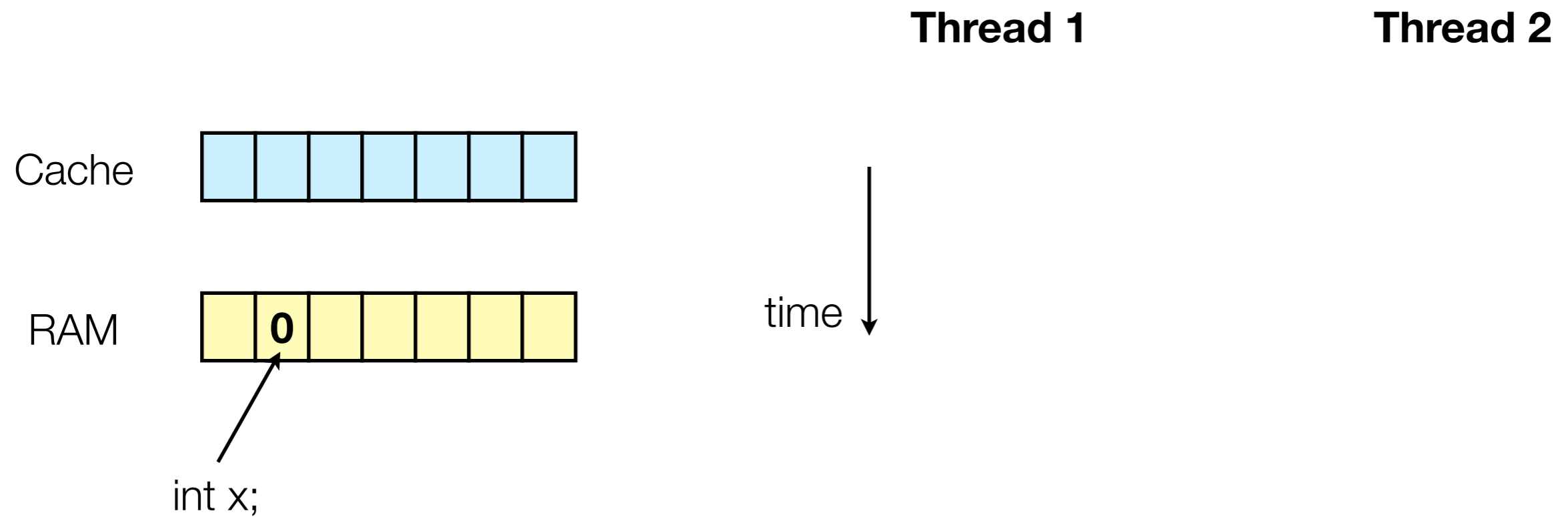
Thread A writes $x=5$ to cache



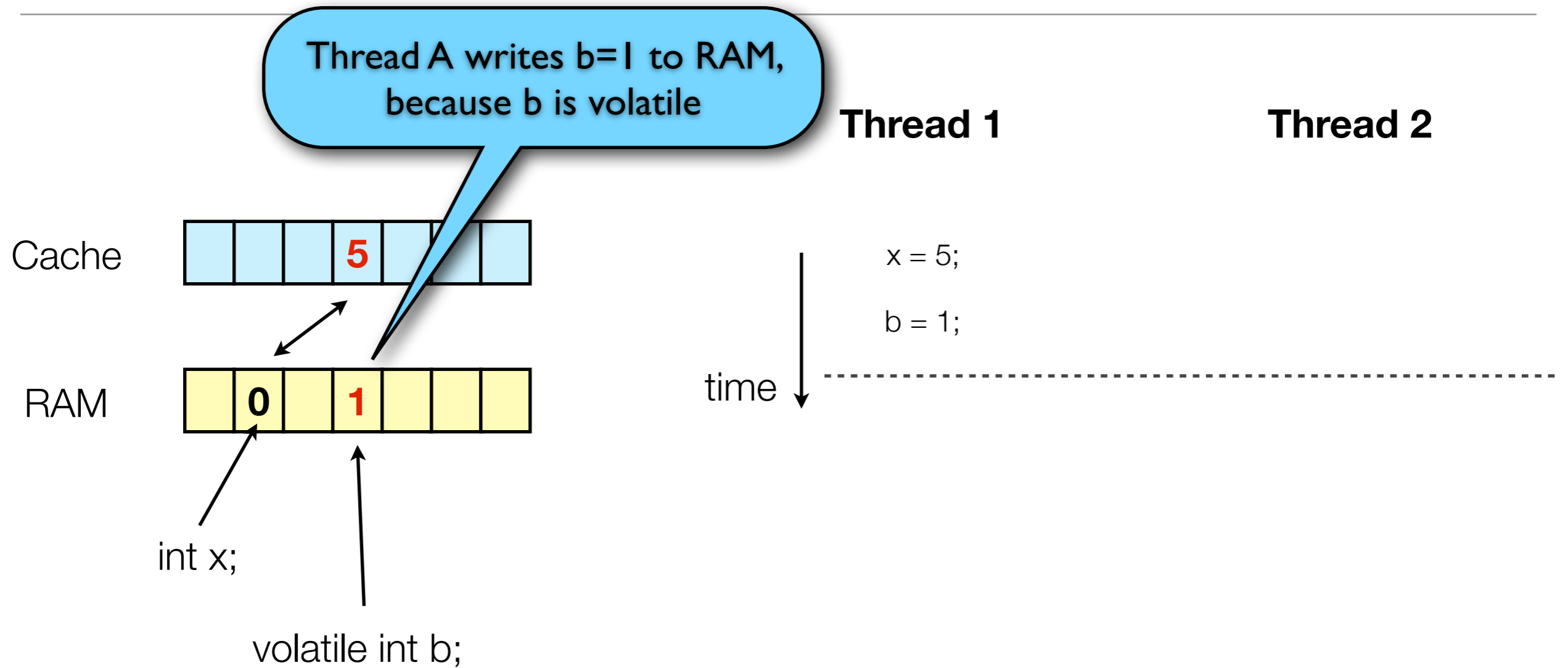
Concurrency



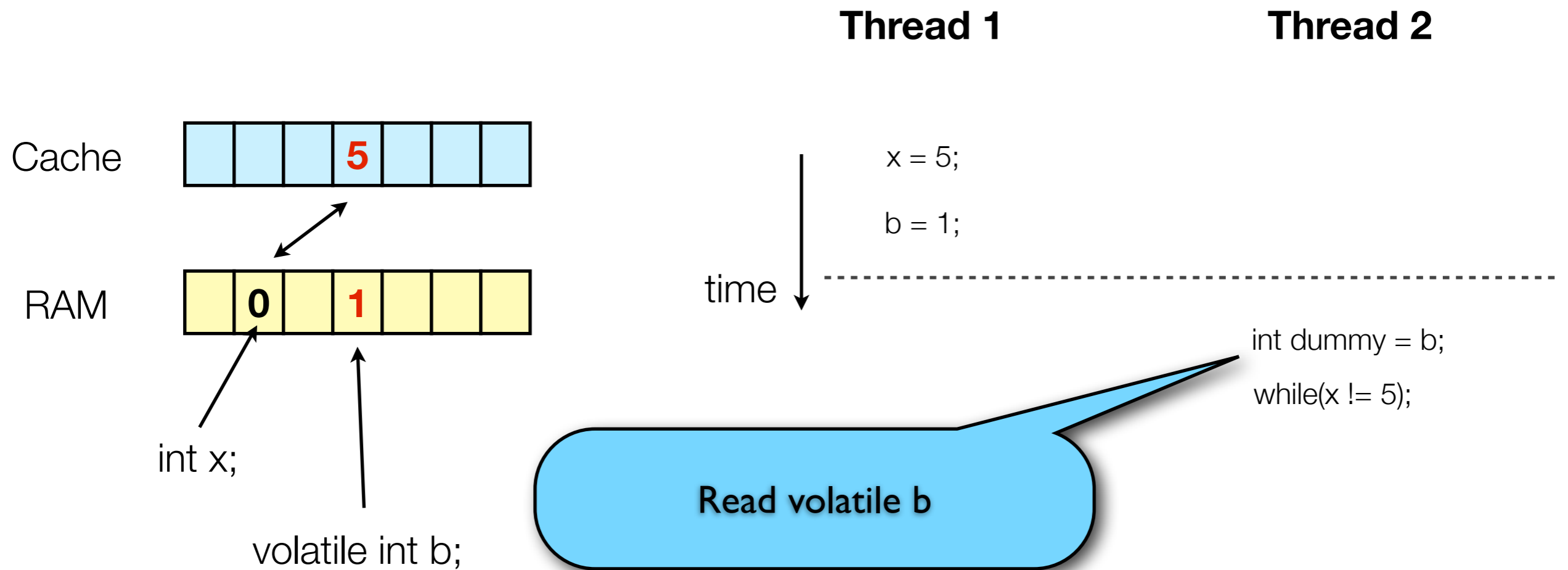
Concurrency



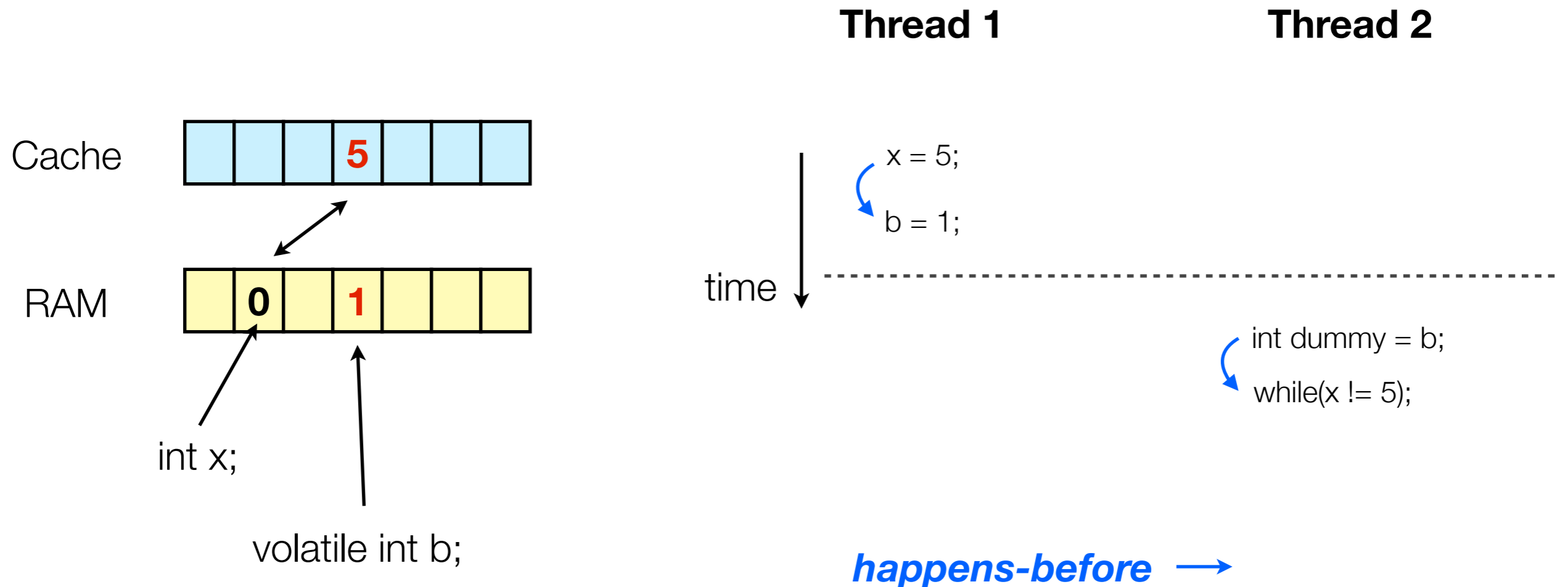
Concurrency



Concurrency

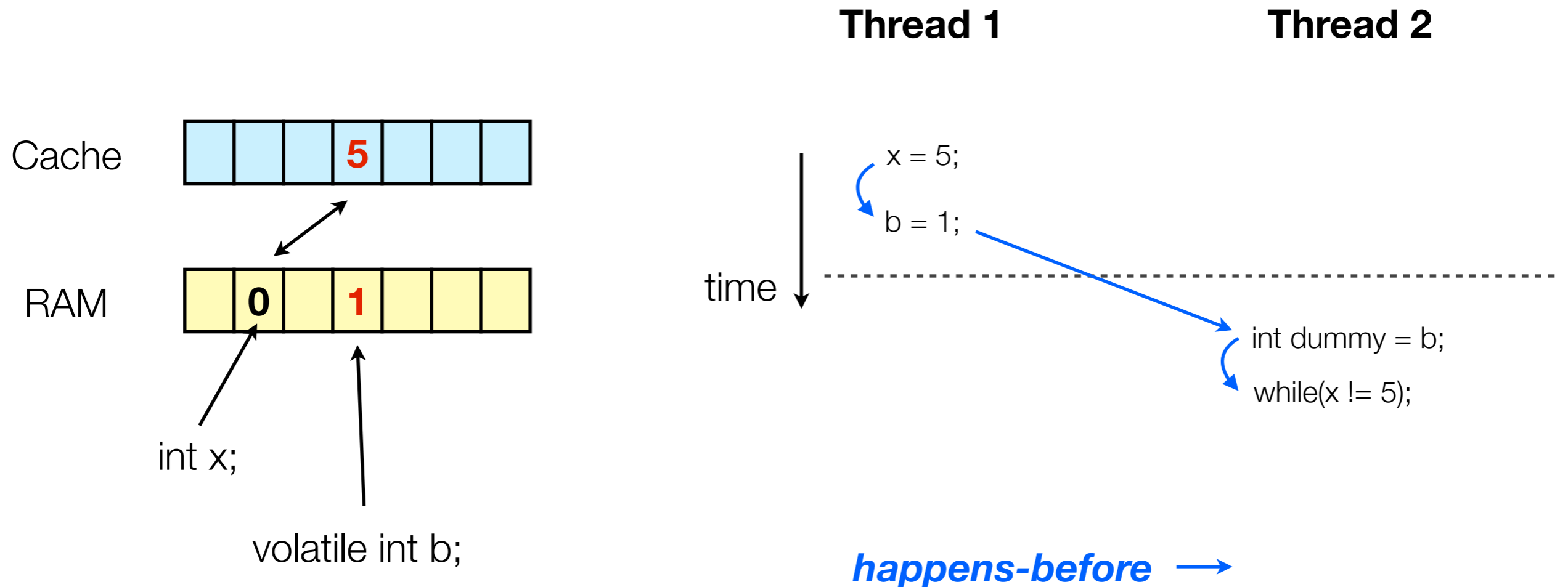


Concurrency



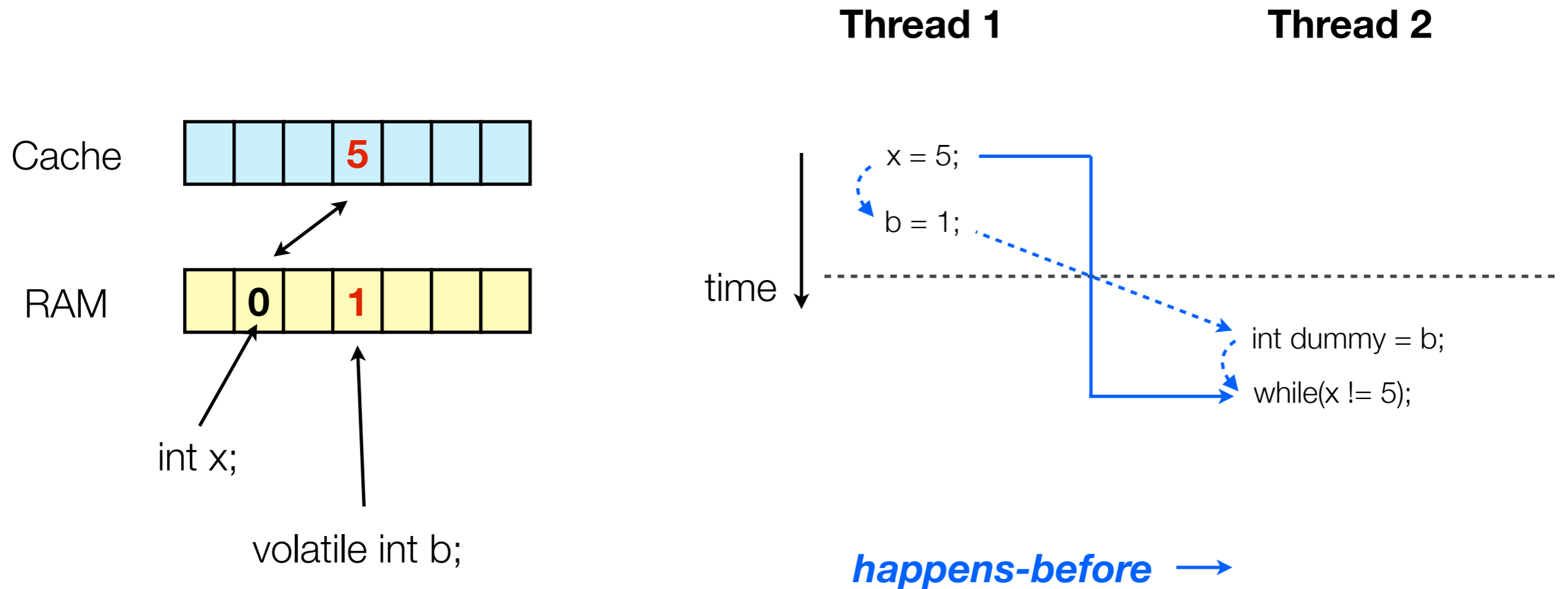
- **Program order rule:** Each action in a thread *happens-before* every action in that thread that comes later in the program order.

Concurrency



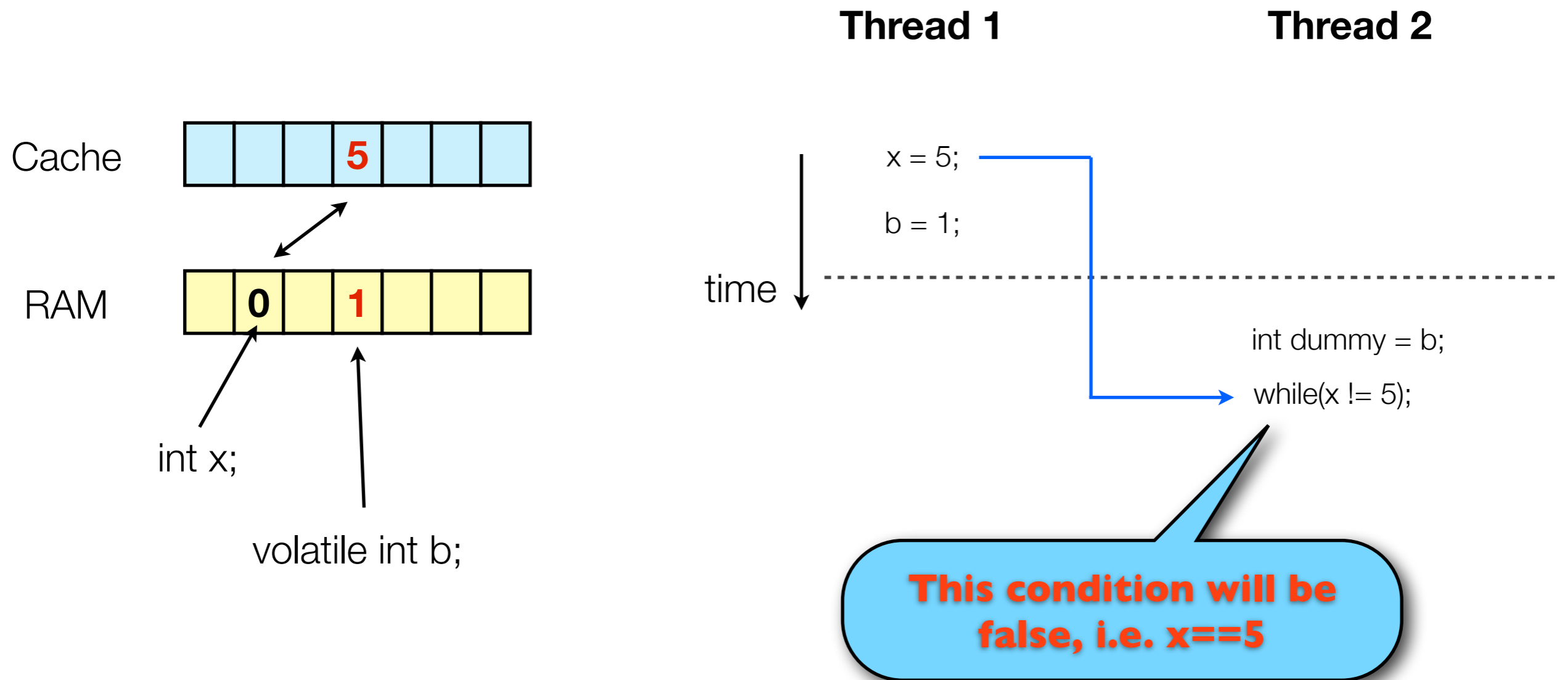
- **Volatile variable rule:** A write to a volatile field *happens-before* every subsequent read of that same field.

Concurrency



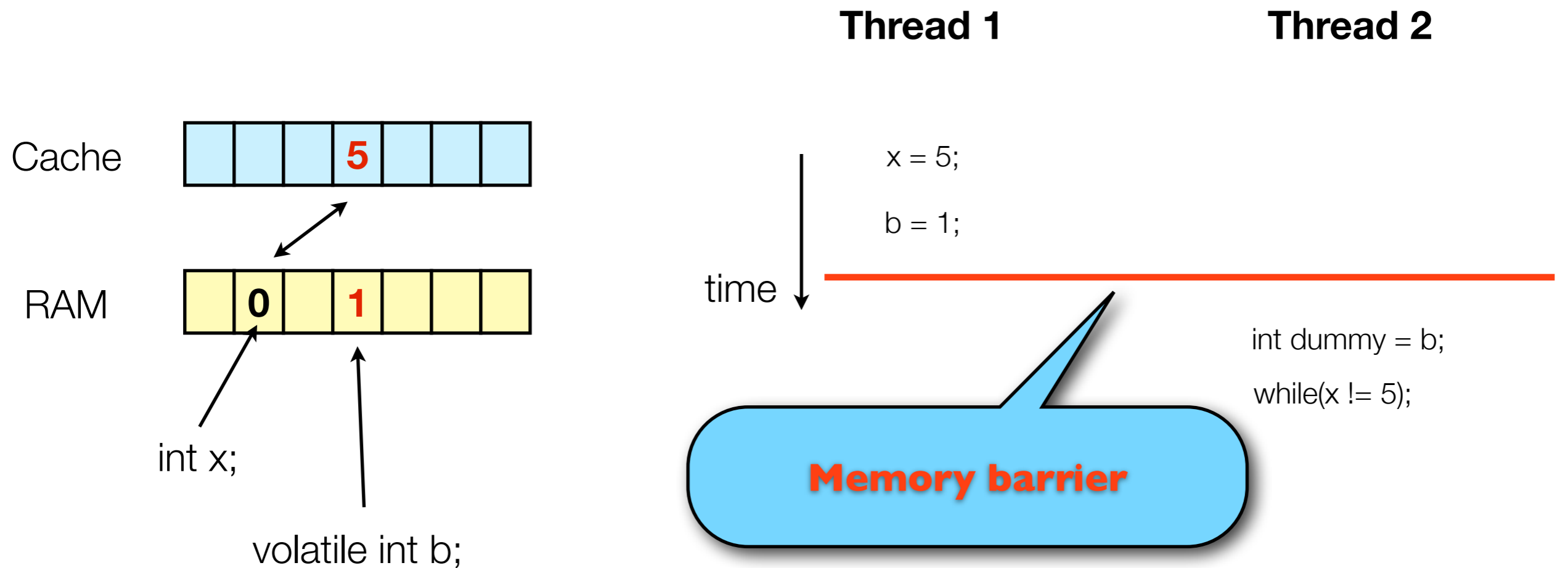
- **Transitivity:** If A *happens-before* B, and B *happens-before* C, then A *happens-before* C.

Concurrency



- **Note:** `x` itself doesn't have to be volatile. There can be many variables like `x`, but we need only a single volatile field.

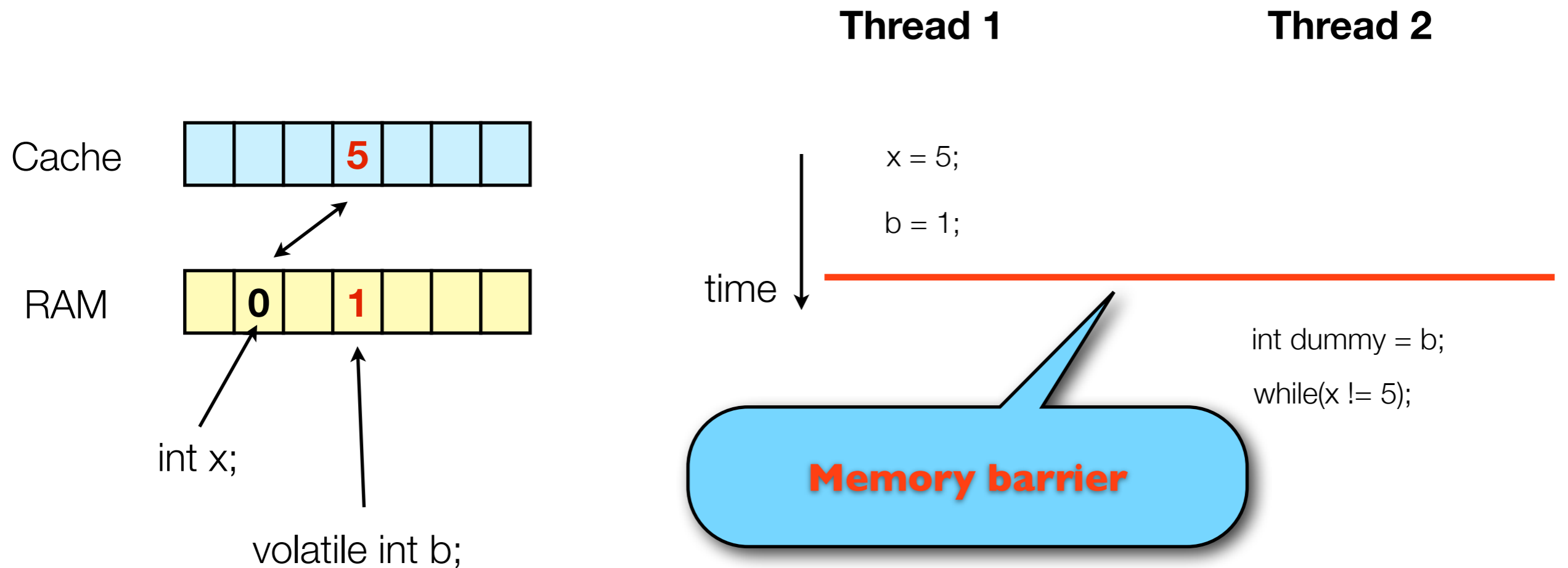
Concurrency



- **Note:** `x` itself doesn't have to be volatile. There can be many variables like `x`, but we need only a single volatile field.

Demo

Concurrency



- **Note:** `x` itself doesn't have to be volatile. There can be many variables like `x`, but we need only a single volatile field.

Concurrency

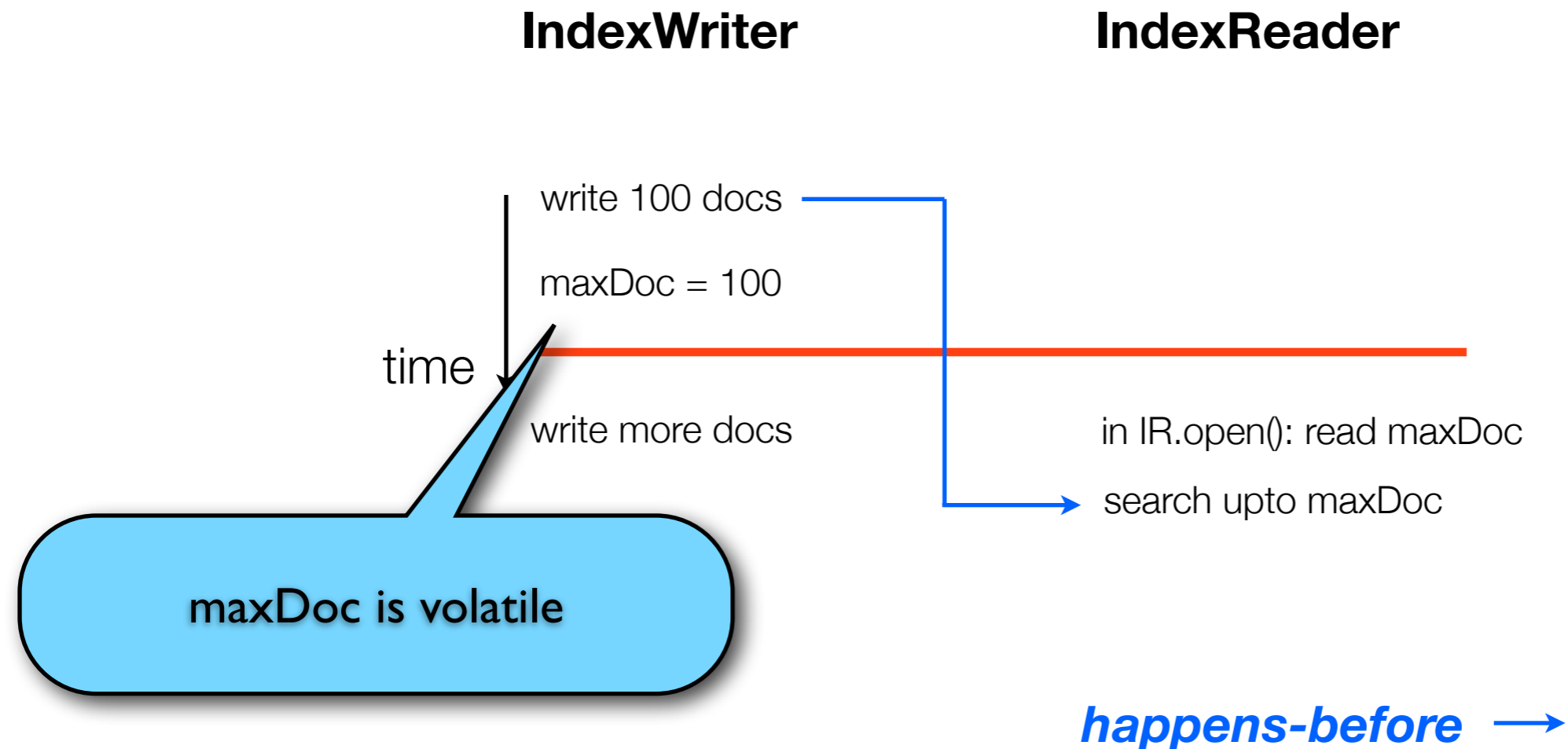
IndexWriter

IndexReader



maxDoc is volatile

Concurrency



- **Only maxDoc is volatile.** All other fields that IW writes to and IR reads from don't need to be!

Wait-free

- Not a single exclusive lock
- Writer thread can always make progress
- Optimistic locking (retry-logic) in a few places for searcher thread
- Retry logic very simple and guaranteed to always make progress

Realtime Search @twitter

Agenda

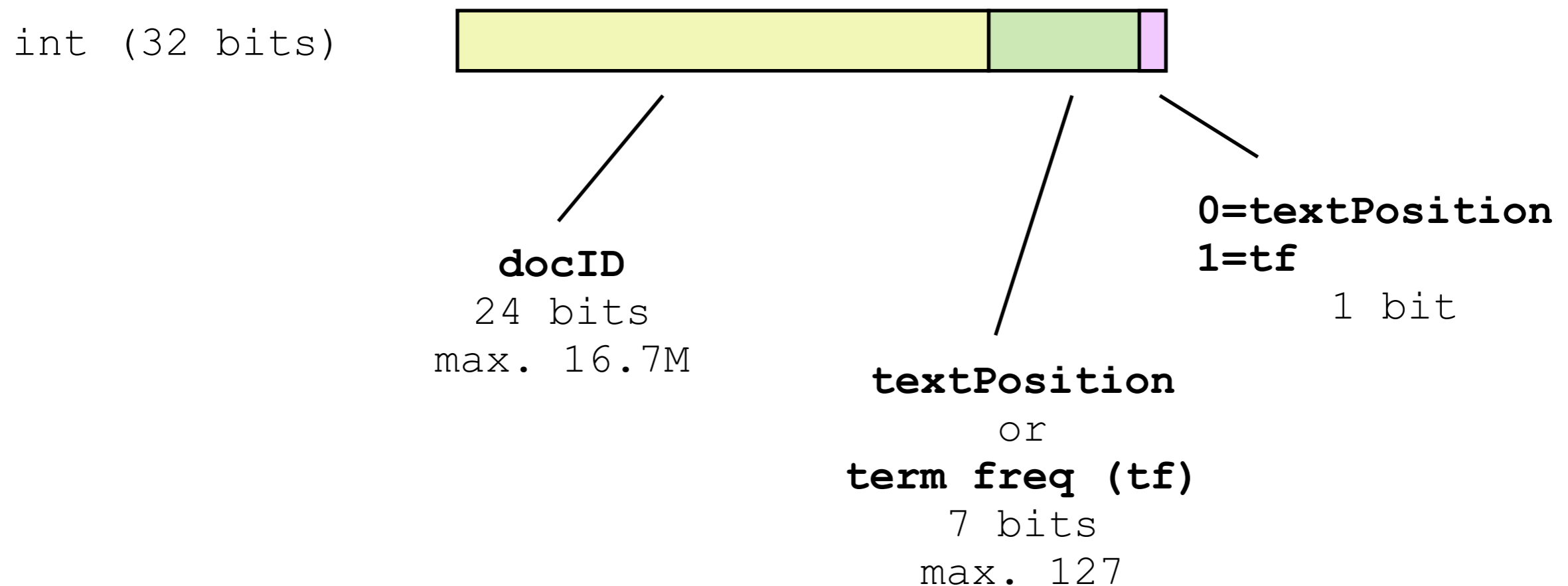
- Introduction
- Search Architecture
- Inverted Index 101
- Memory Model & Concurrency
- ▶ What's next?

What's next?

What's next?

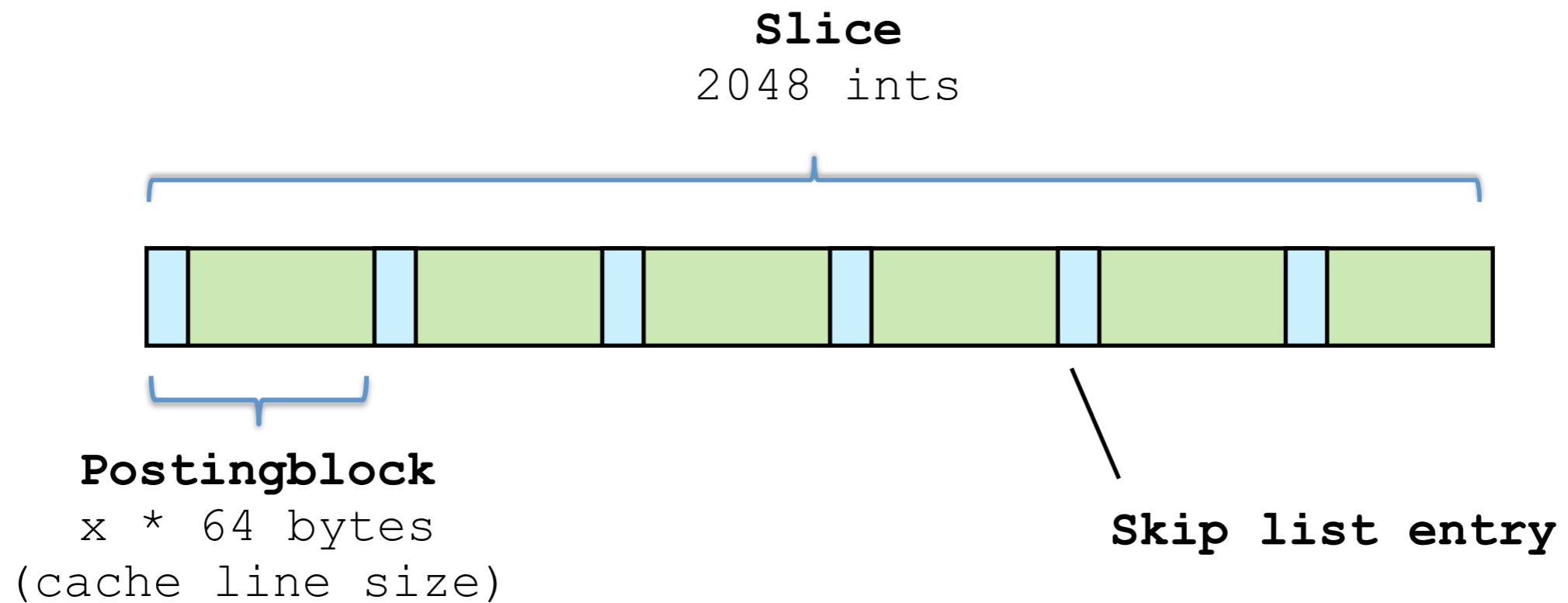
- Posting list format that supports
 - Positions > 255
 - Payloads
 - Point-in-time document frequencies (df)
- Code complete
- Performance tests soon!

DocID Posting



- textPosition is only stored inline, if (tf == 1 && textPosition <=127 && hasPayload == false)
- change from old format: for tf > 1 we don't repeat the docID anymore

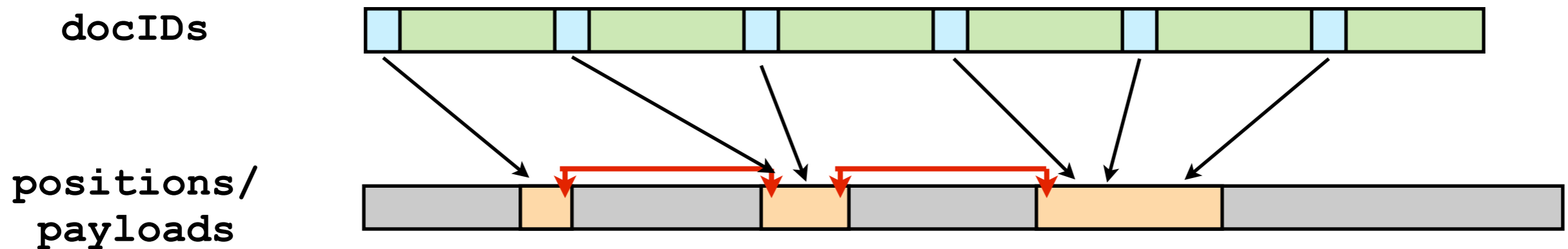
Embedded Skip Lists



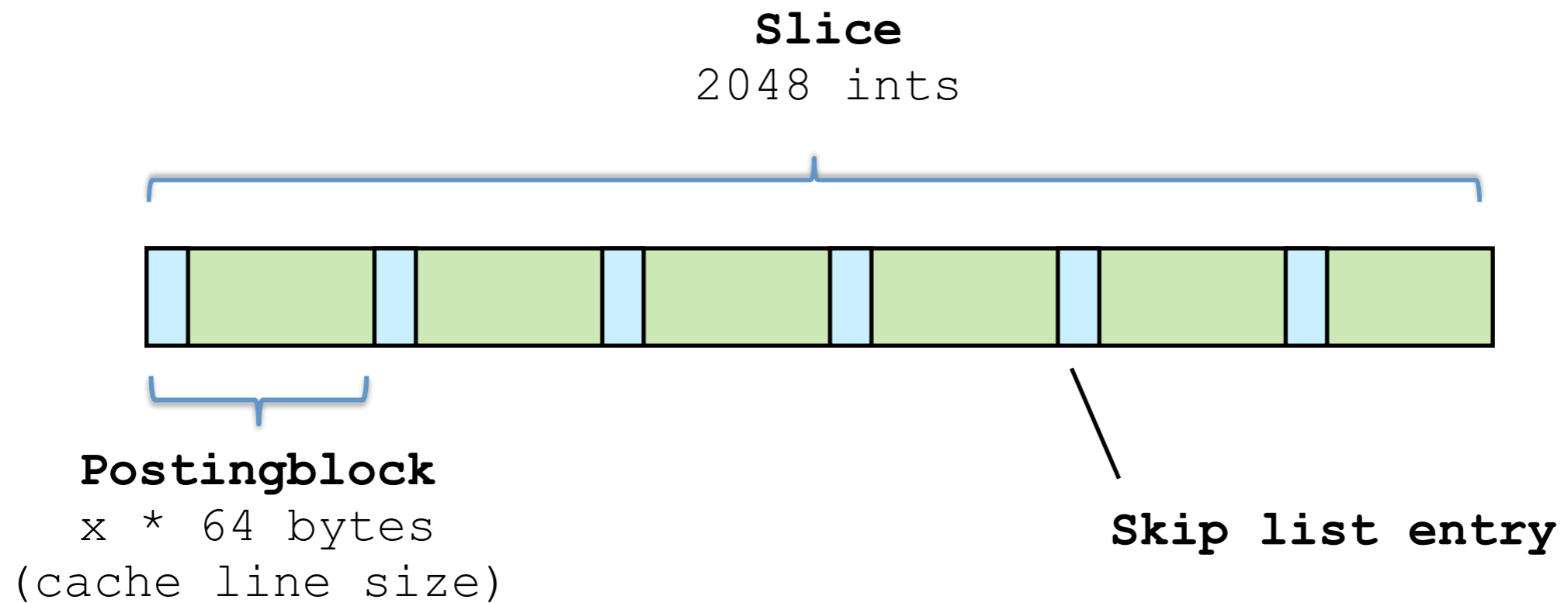
- x to be determined in performance tests

Storing Positions and Payloads

- Use Lucene's ByteBlockPool: It builds up a linked list of **byte slices** of increasing lengths
- Unlike Lucene's ByteBlockPool we use double-linked lists, i.e. the slices will have pointers to the previous **and** next slices

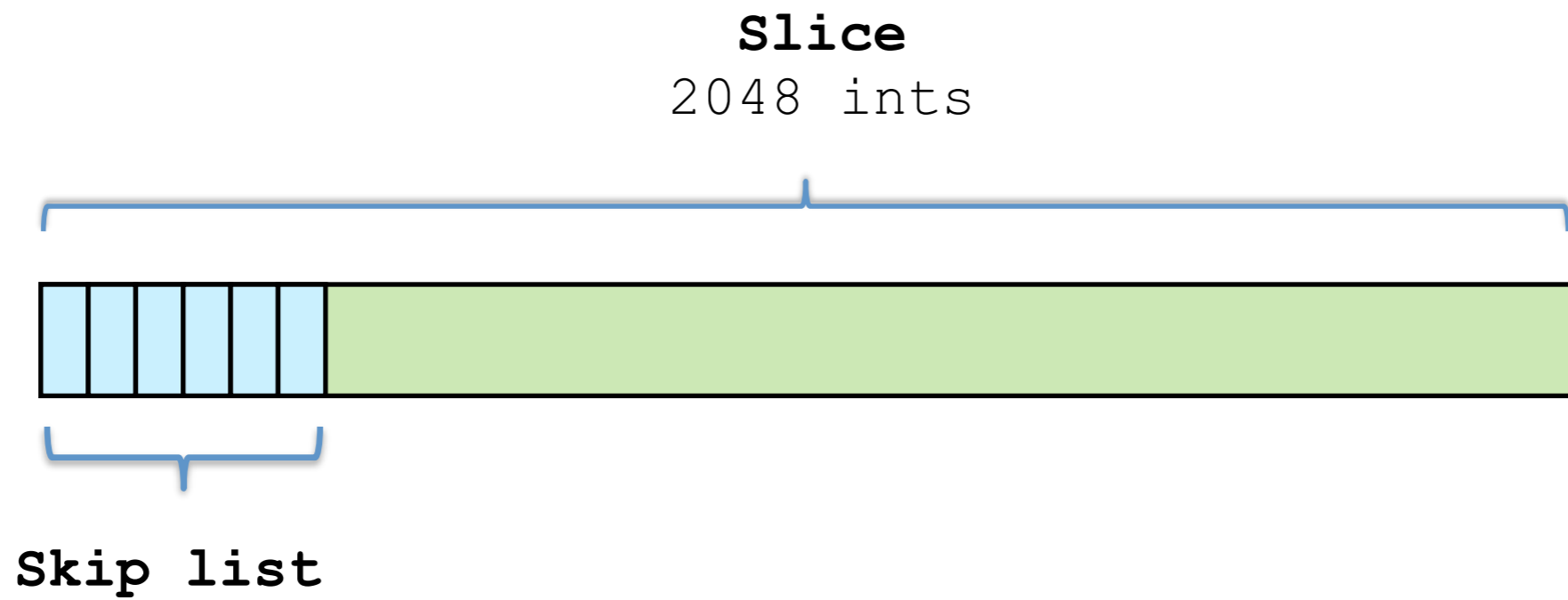


Embedded Skip Lists

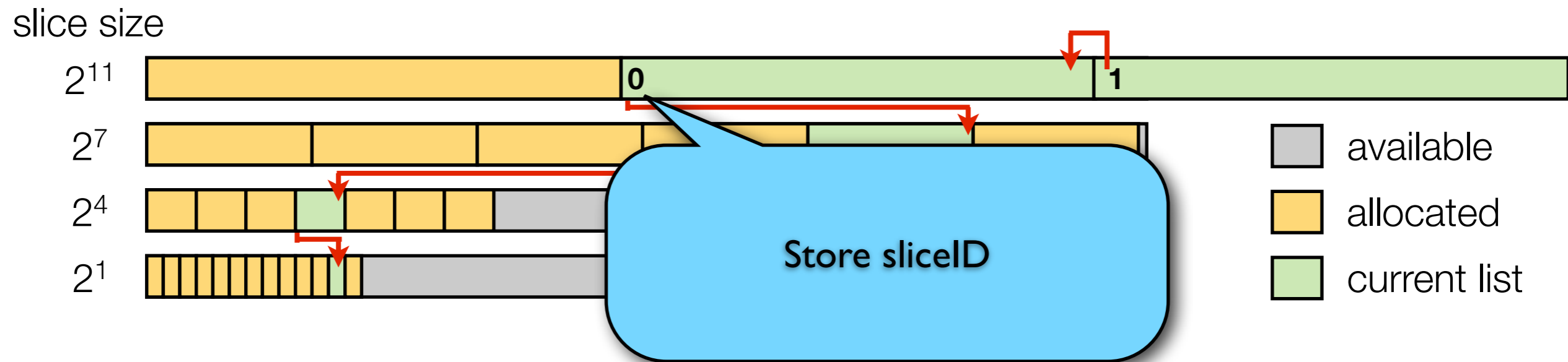


- x to be determined in performance tests

Embedded Skip Lists



Document Frequencies



- SliceIDs only stored for slices on highest level

E.g. $DF = 2^1 + 2^4 + 2^7 + (\text{sliceID} * 2^{11}) + \text{offsetWithinSlice}$

Summary

- Efficient for small documents due to position inlining
- Position/payload encoding size comparable to vanilla Lucene for bigger documents
- Concurrency model unchanged
- A reader thread will never try to access positions/payloads that have not been safely published yet
- Document frequencies can be looked up in constant time (even worst case)

twitter 

Questions?