

# Next Generation Cassandra Compaction, Going beyond LCS

Joey Lynch



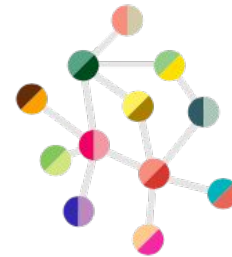
## Speaker

Joey Lynch



Senior Software Engineer  
Cloud Data Engineering at Netflix

Distributed system addict and data  
wrangler



# Outline

Main Compaction challenges

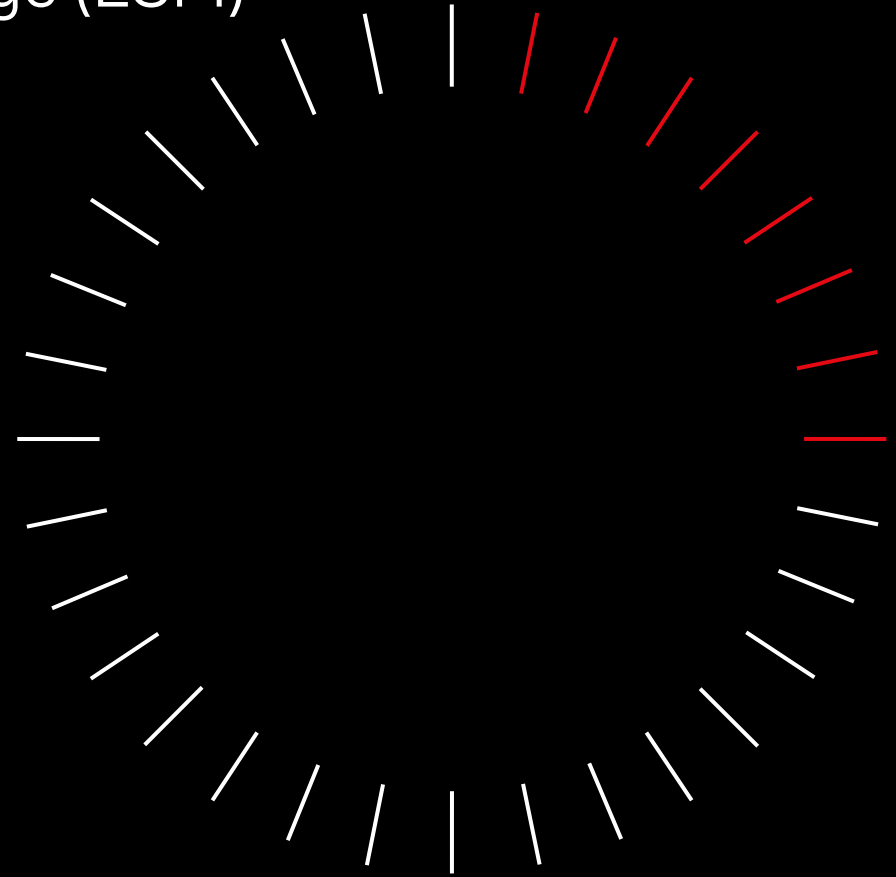
General Purpose Compaction Strategies:

- Size Tiered
- Leveled

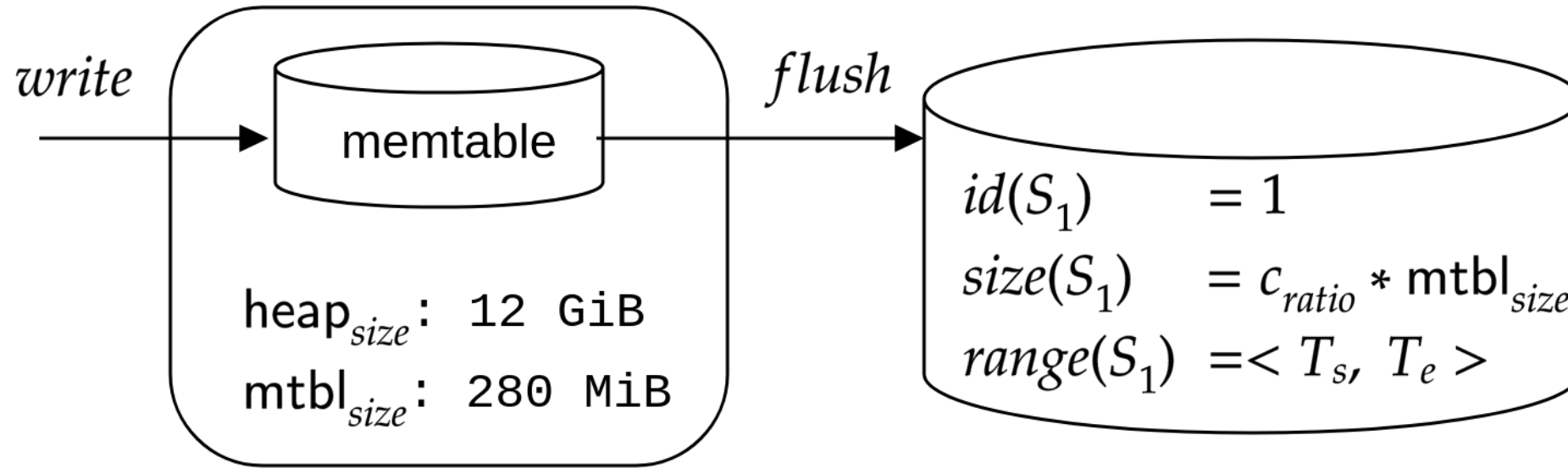
Next Generation Compaction Strategy

# Compaction

Putting the merge  
in log structured  
merge (LSM)



# Path to Compaction: Flush



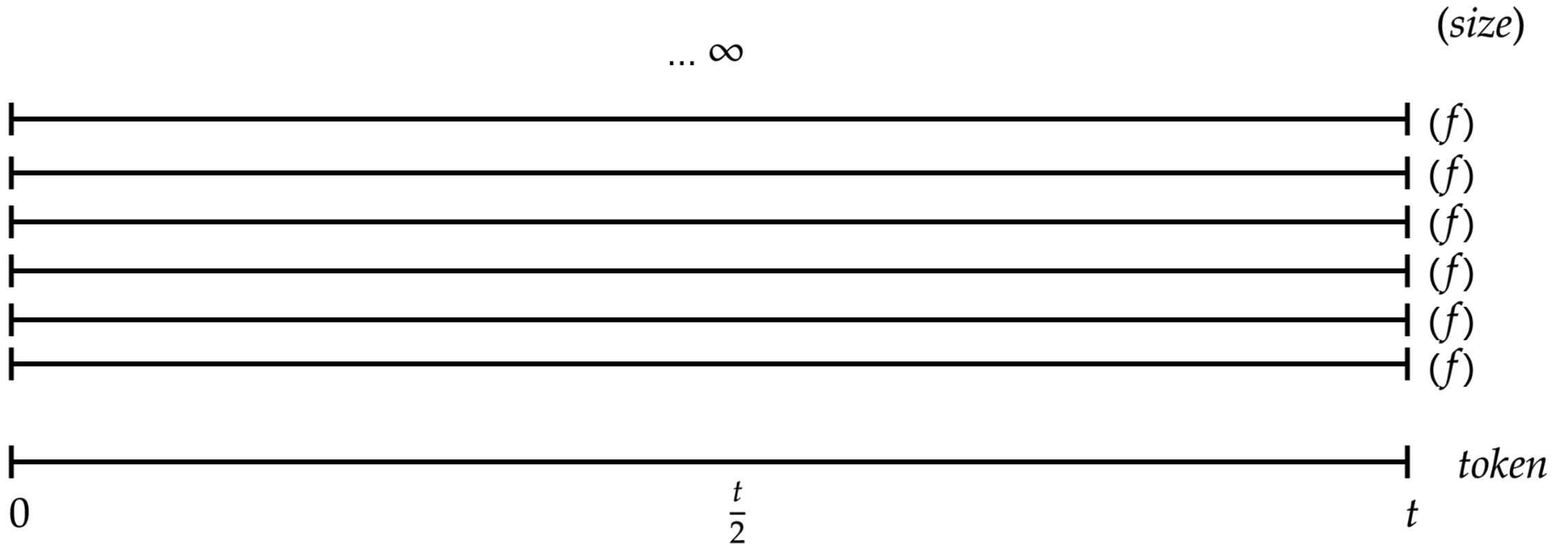
Writes are coalesced into a memtable, then flushed

# What is an “SSTable”

```
$ ls -latr | grep 79649
    2468 Jul  1 16:38 mc-79649-big-Summary.db
 1353995 Jul  1 16:38 mc-79649-big-Index.db
    9784 Jul  1 16:38 mc-79649-big-Filter.db
335548485 Jul  1 16:38 mc-79649-big-Data.db
     10 Jul  1 16:38 mc-79649-big-Digest.crc32
 189987 Jul  1 16:38 mc-79649-big-CompressionInfo.db
    92 Jul  1 16:38 mc-79649-big-TOC.txt
 11959 Jul  1 16:38 mc-79649-big-Statistics.db
```

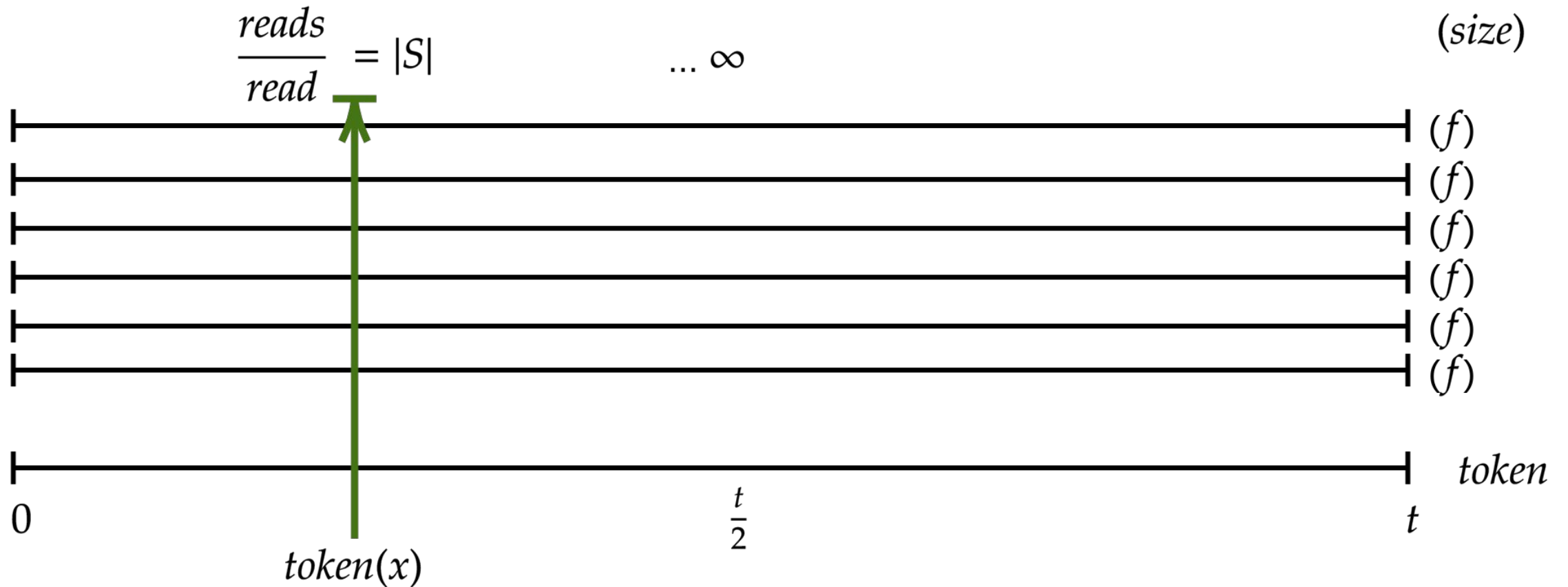
One “SSTable” actually multiple “components”

# Problem: Space Amplification



We have unbounded number of SSTables

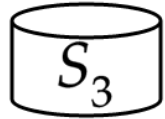
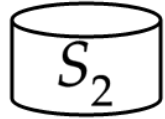
# Problem: Read Amplification



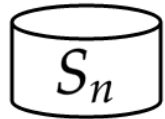
You wanted to *read* your data?



# Compaction Inputs



...



$size(S_i)$

=

Either compressed or uncompressed size

$range(S_i)$

=

What range of tokens does the SSTable span?

$level(S_i)$

=

Non-overlapping guarantee

Have some SSTables

Have some metrics

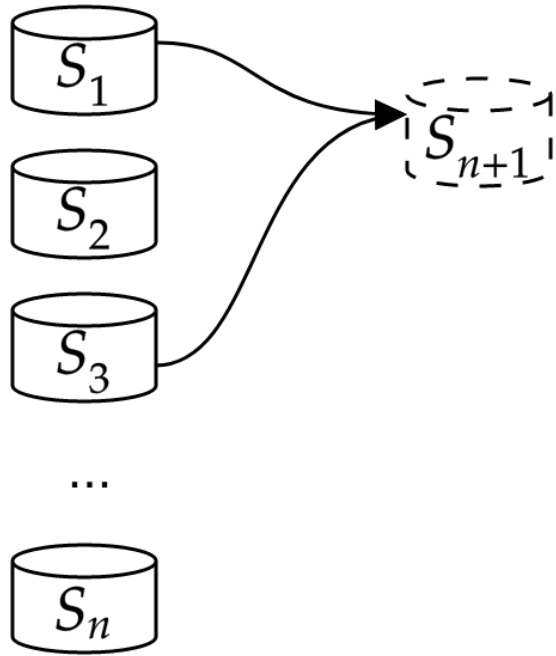
# Compaction Strategy Picks Candidates



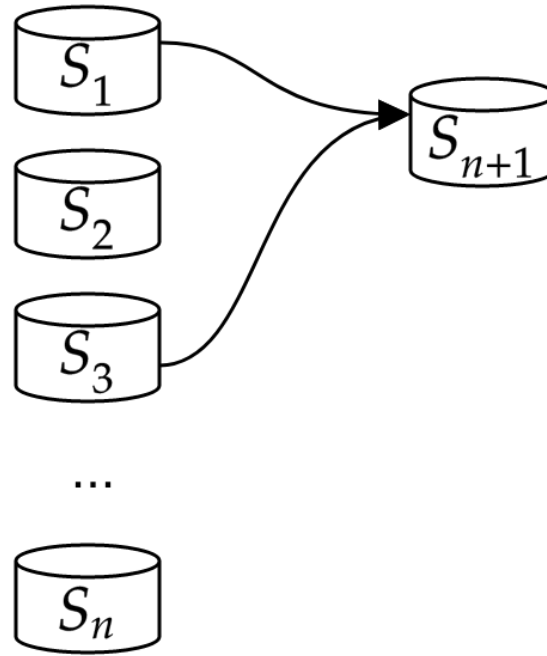
$$S = \{S_1, S_2, \dots, S_n\}$$

$$\begin{aligned} \text{let } C &= \text{candidates}(S) \\ \text{s.t. } C &\subseteq S \end{aligned}$$

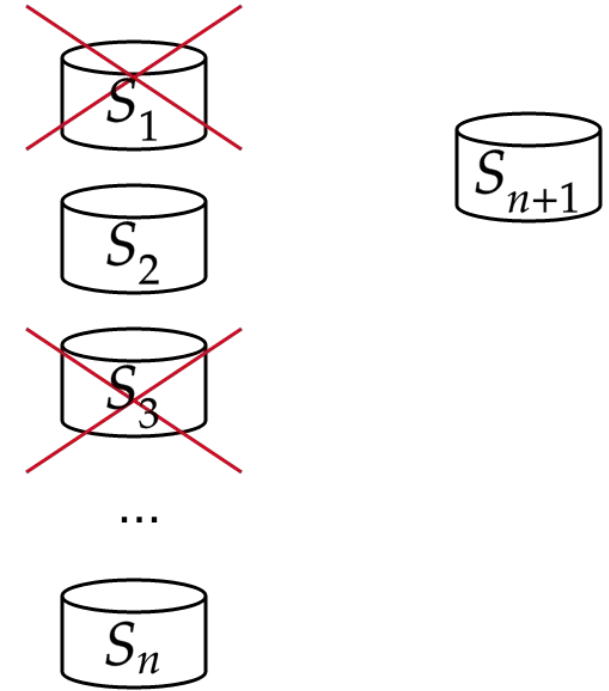
# Compaction Outputs



Phase 1



Phase 2



Phase 3

# Compaction Is Expensive

## Phase 1

Reads through a bunch of data, materializing it into heap

**This is very expensive**

## Phase 2

Re-builds view

This can be expensive

## Phase 3

Re-builds view,  
Drops old data, OS page cache drops

This can be expensive

## Problems?

Space Amplification

Read Amplification

Write Amplification

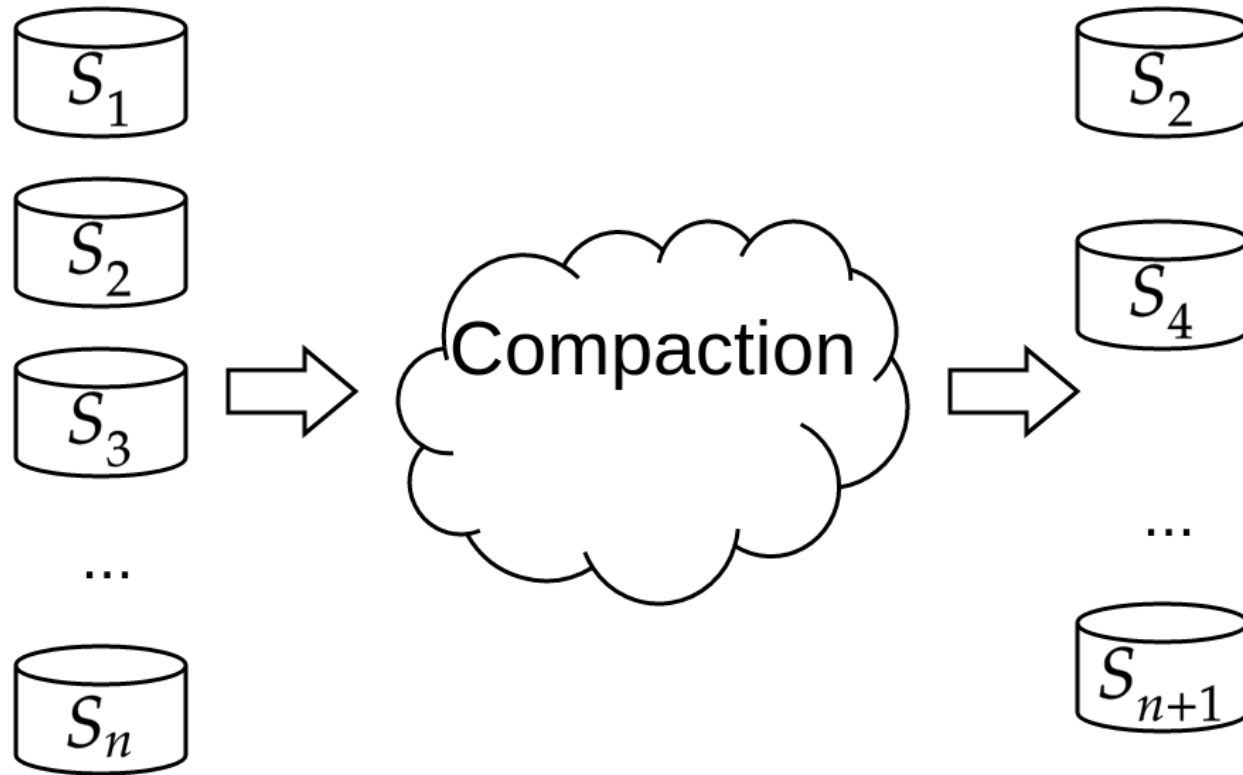
## Why?

Can't take infinite disk space ...

Reads should be fast

Compaction is expensive

# Compaction Strategy Goal



*choose candidates  $C$*

*let  $S' = S - C + S_n + 1$*

*s. t*

*$amp_r(S') < amp_r(S) \quad \wedge$*

*$amp_s(S') < amp_s(S) \quad \wedge$*

*$min(|C|)$*

$$S = \{S_1, S_2, \dots, S_n\}$$

$$S' = \{S_2, S_4, \dots, S_{n+1}\}$$

**“How much disk space did I reclaim?”**

**“Is the data really gone?”**

**“My reads are slow!!”**

## Problems?

Full compaction

Manual compaction

Handle tons of small  
SSTables

## Why?

Have to be able to make  
guarantees

Targeted incidents (hot  
partitions)

Repair is fun this way



## Lessons Learned

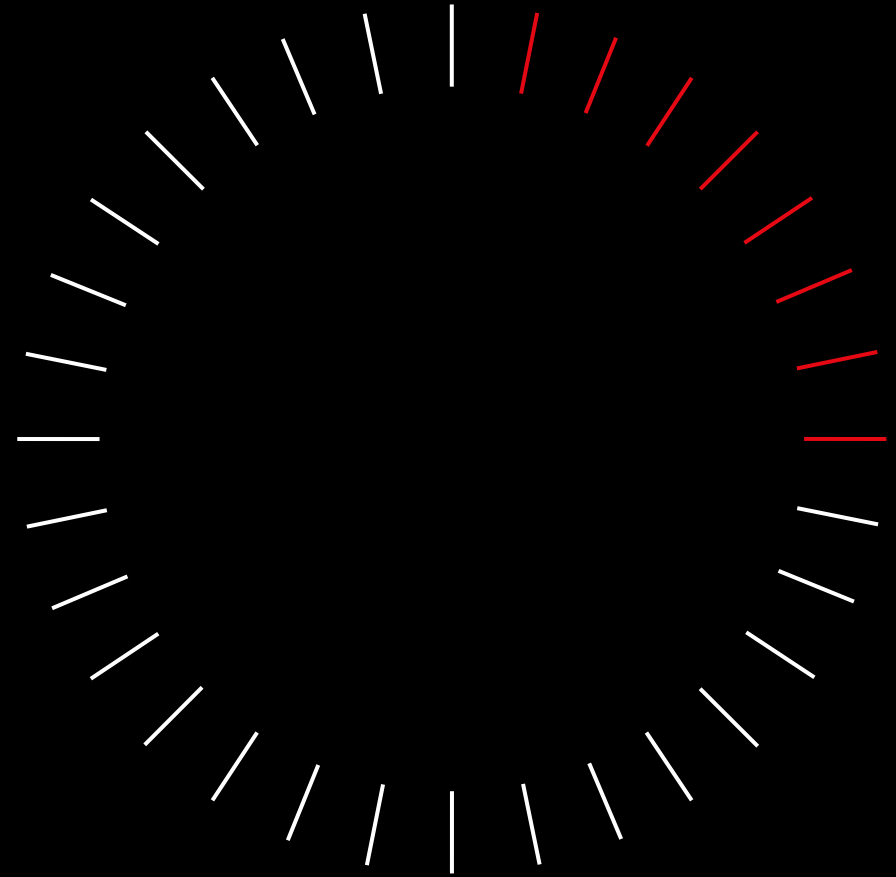
Compaction is an **optimization** problem

Trying to reduce space and read amplification

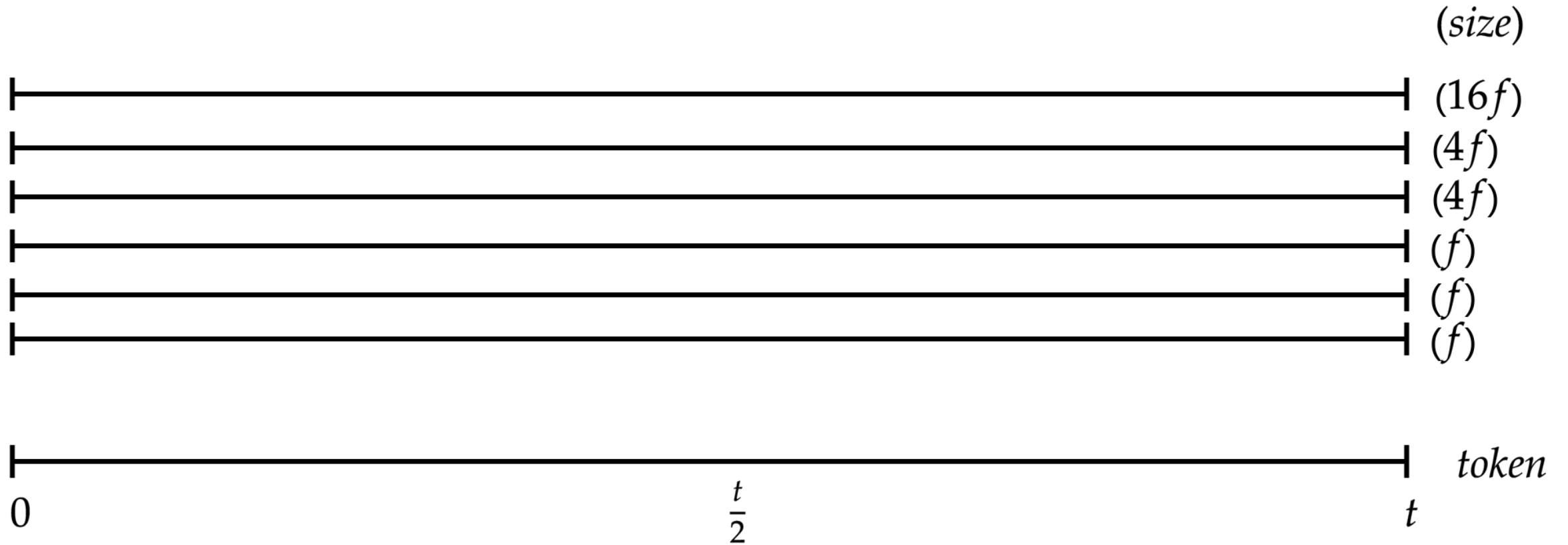
... while doing as little write amplification as we can

# Size Tiered Compaction

Let's do the most  
obvious thing



# Group SSTables by Size



## Main Tunables

min\_threshold

bucket\_low

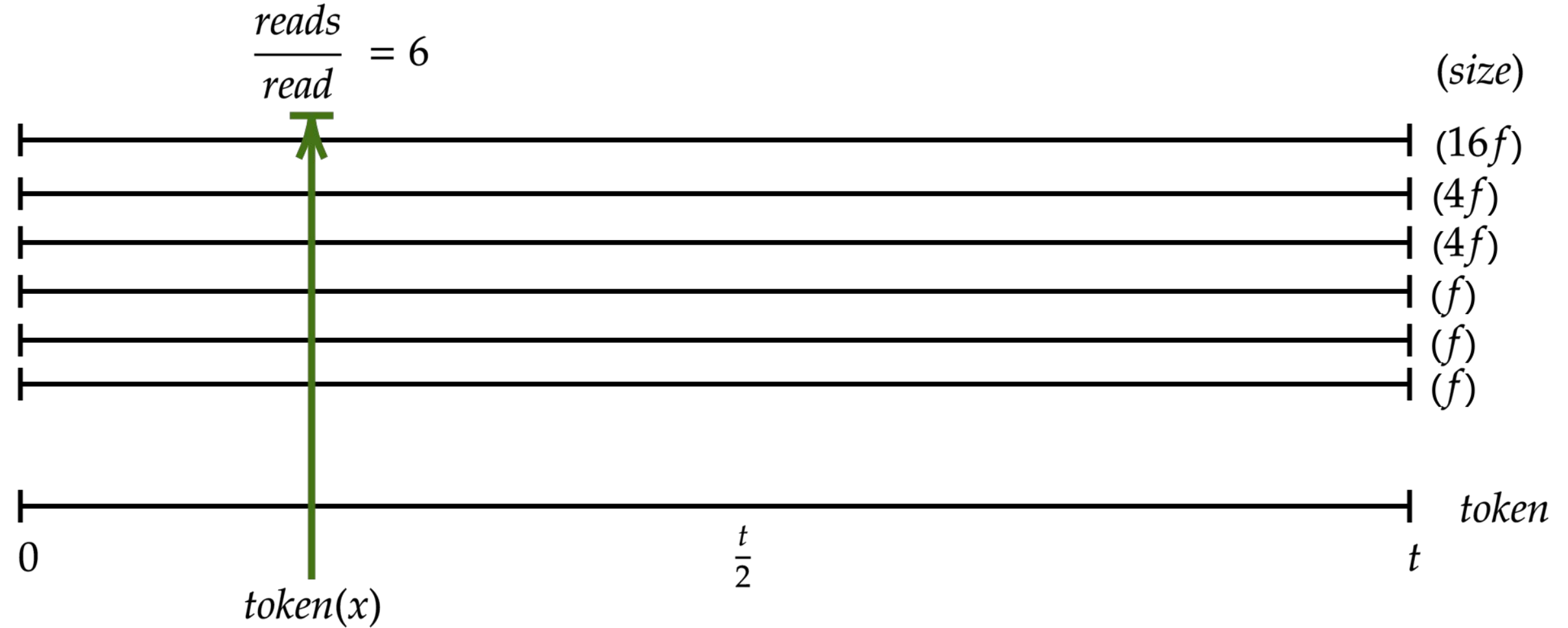
bucket\_high

## Meaning?

This dictates the level of tiering

What should be bucketed together

# Group SSTables by Size



# Advantages

Good for write only workloads

Full compaction works great

Super simple

# Why?

Exponential re-write curve

Small number of files

Fewer bugs

## Problems?

50% space overhead

Giant SSTables

Begat early re-open 🔥 🔥

Read amplification

## Why do we care?

Hard to provision for

Hard to backup / transfer

This had a **ton** of bugs

Variable read  
performance is bad

## Problems?

Needless write  
amplification

Have to do periodic full  
compactions to give any  
guarantees

## Why do we care?

Wasting CPU is bad

Re-writing the whole  
dataset all at once every  
few weeks seems like a  
suboptimal choice



## **Lessons Learned**

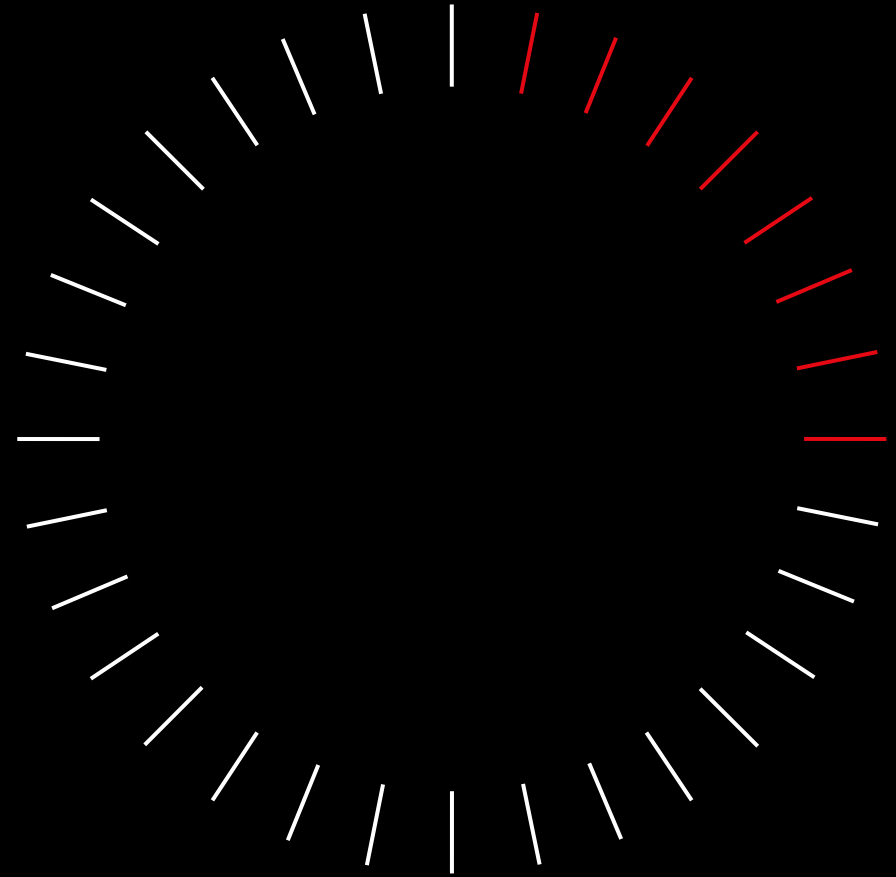
Large single files are problematic

Need to wait until we have enough work from flushes to do.

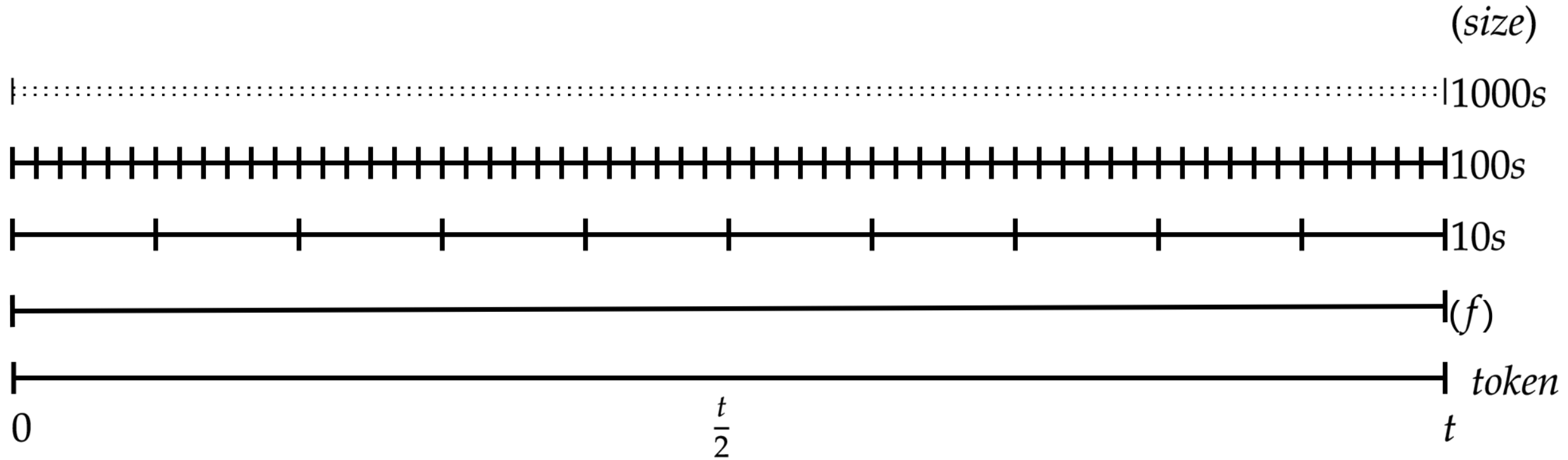
Time bounds on full compaction is nice

# Leveled Compaction

What if we only did  
useful work?



# Sorted runs to the rescue!



# Advantages

Wicked fast reads

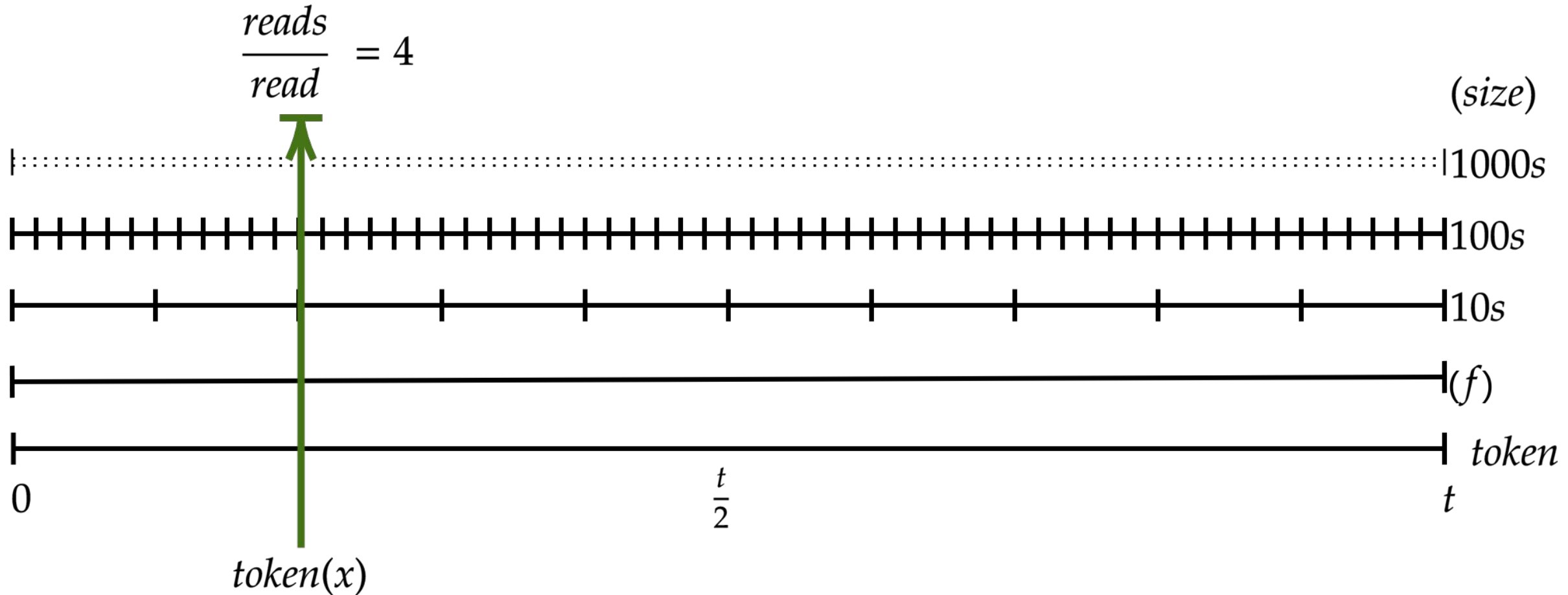
Relatively fixed disk  
usage

# Why?

Bounded read  
amplification

Bounded space  
amplification

# Sorted runs to the rescue!



## Main Tunables

`sstable_size_in_mb`

## Meaning?

The small, fixed size that LCS should output

Each sorted run or “level” will contain 10x the previous sorted run

# Advantages

Small files easy to work with

Good for TTL data

Good for updates  
Good for reads

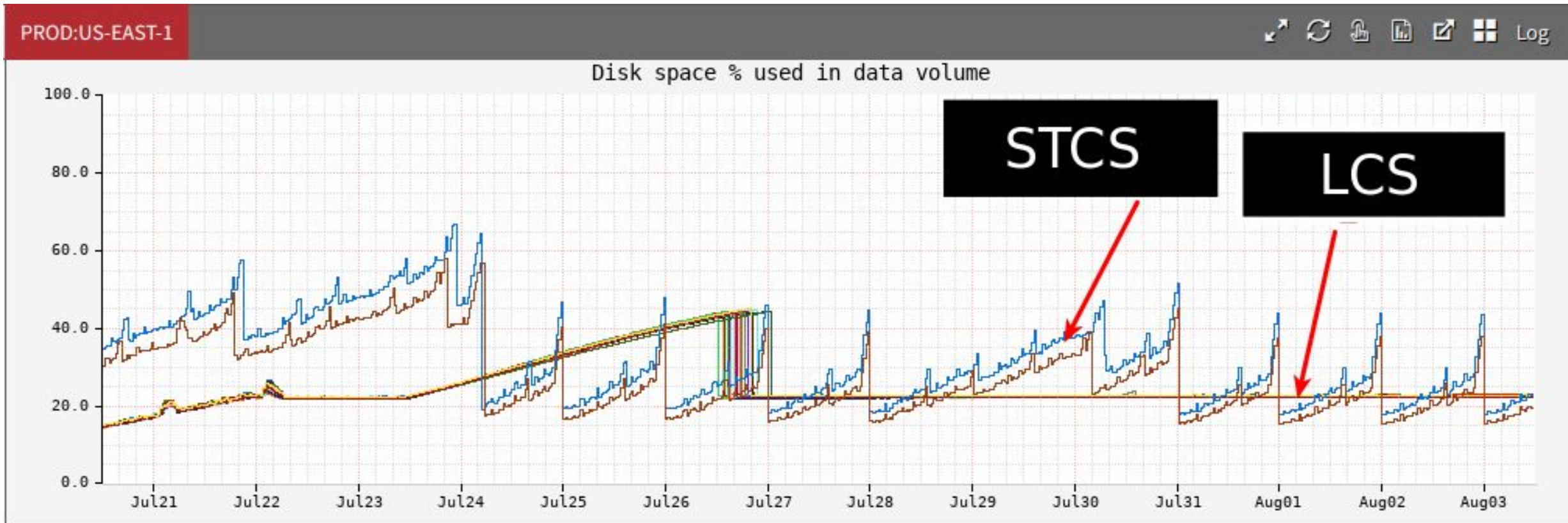
# Why?

File parallelism is generally easy

Self compactions viable

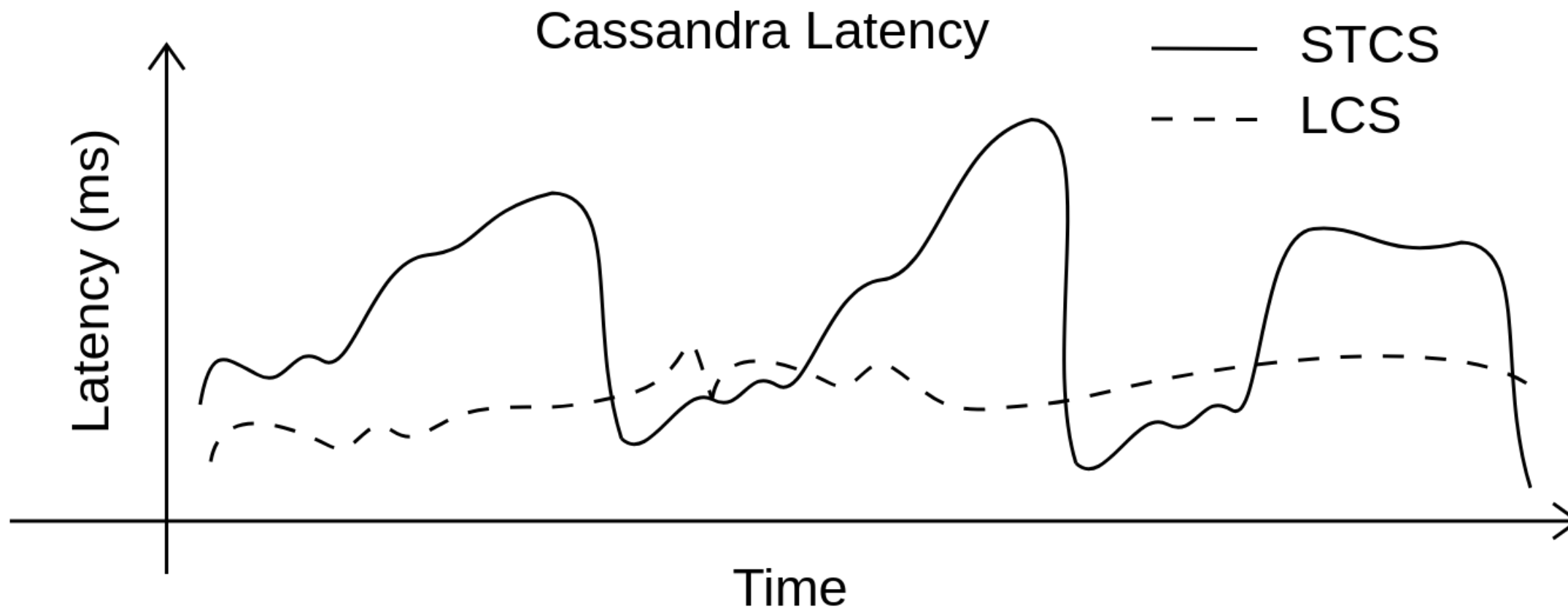
Levels are really nice ...

# It Really Works





# It Really Works



## Problems?

Data must go through L1

Does a lot of compaction

Really complicated

## Why do we care?

Can “fall behind”  
incoming writes

For spinning disks this  
can be a problem ...

A large number of bugs

## Problems?

Default 160 MiB can  
create a *lot* of sstables

SSTable size is  
underconstrained

## Why do we care?

Upstream processors  
may not be able to  
handle 20k files

Have to raise SSTable  
size when all you want is  
larger L1

## Problems?

Hard to get data from the upper levels “down”

Full compaction is currently very slow

## Why do we care?

Getting updates (deletes) to find each other

Pretty huge operability issue

# Four day compactions ...



Cassandra / CASSANDRA-14605

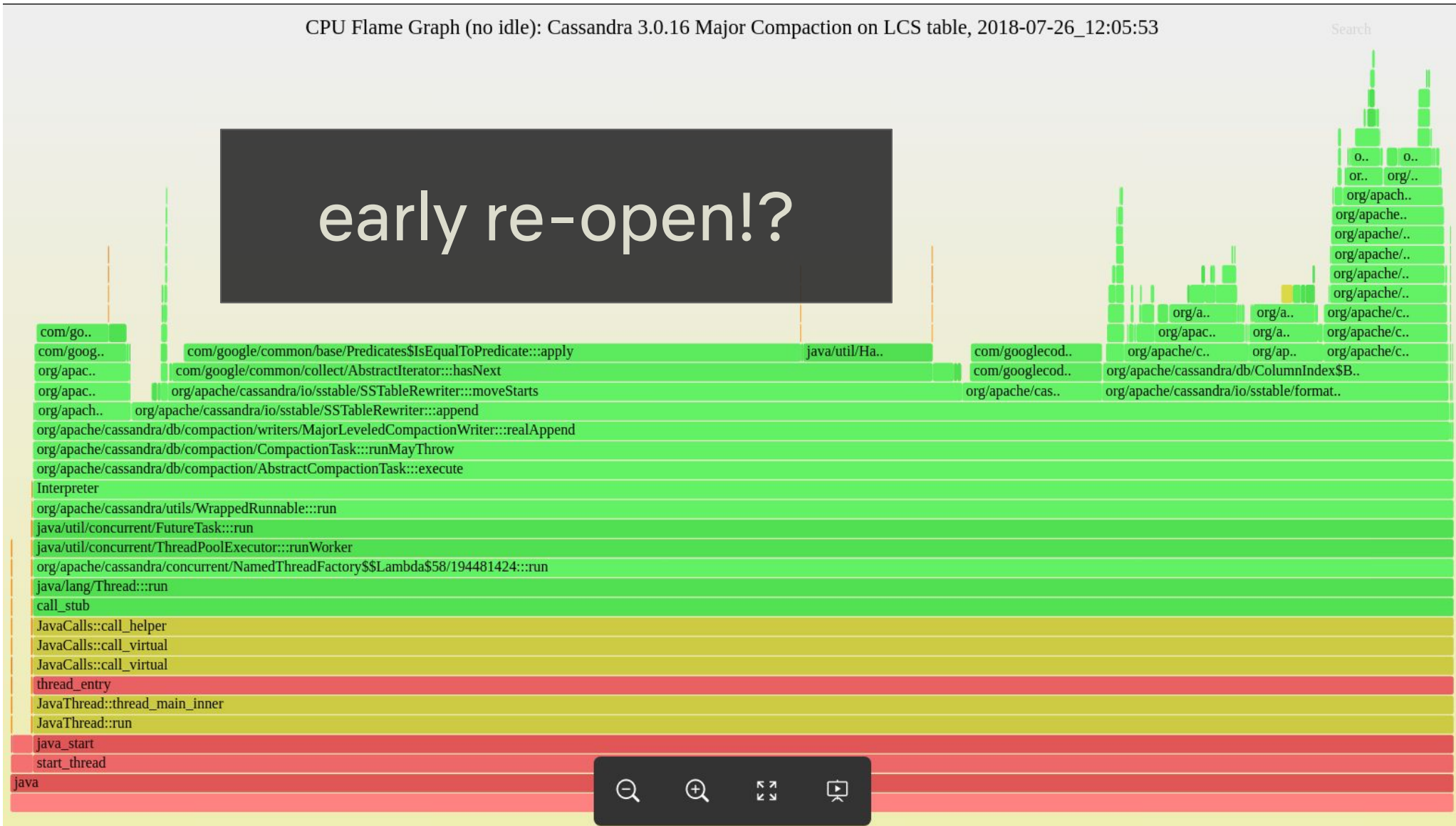
## Major compaction of LCS tables very slow

### > Details

### ▼ Description

We've recently started deploying 3.0.16 more heavily in production and today I noticed that full compaction of LCS tables takes a much longer time than it should. In particular it appears to be faster to convert a large dataset to STCS, run full compaction, and then convert it to LCS (with re-leveling) than it is to just run full compaction on LCS (with re-leveling).

# Four day compactions ...



# There is a solution!

✓  **Marcus Eriksson** added a comment - 27/Jul/18 05:17

could you try setting `ssstable_preemptive_open_interval_in_mb` to -1 and see how it performs?

✓  **Joseph Lynch** added a comment - 30/Jul/18 20:25 

**Marcus Eriksson** I've set the parameter and taken another flamegraph, all of the `moveStarts` usage is gone (nice!) and it appears to be proceeding about twice as fast (instead of taking 4 days it will take 2 days). I've attached the new flamegraph

## **Lessons Learned**

Sorted runs are powerful  
Eliminates a lot of the need for early  
re-open

Bookkeeping complexity should be  
avoided if possible

Rigid structure (L0  $\rightarrow$  L1) is not so good

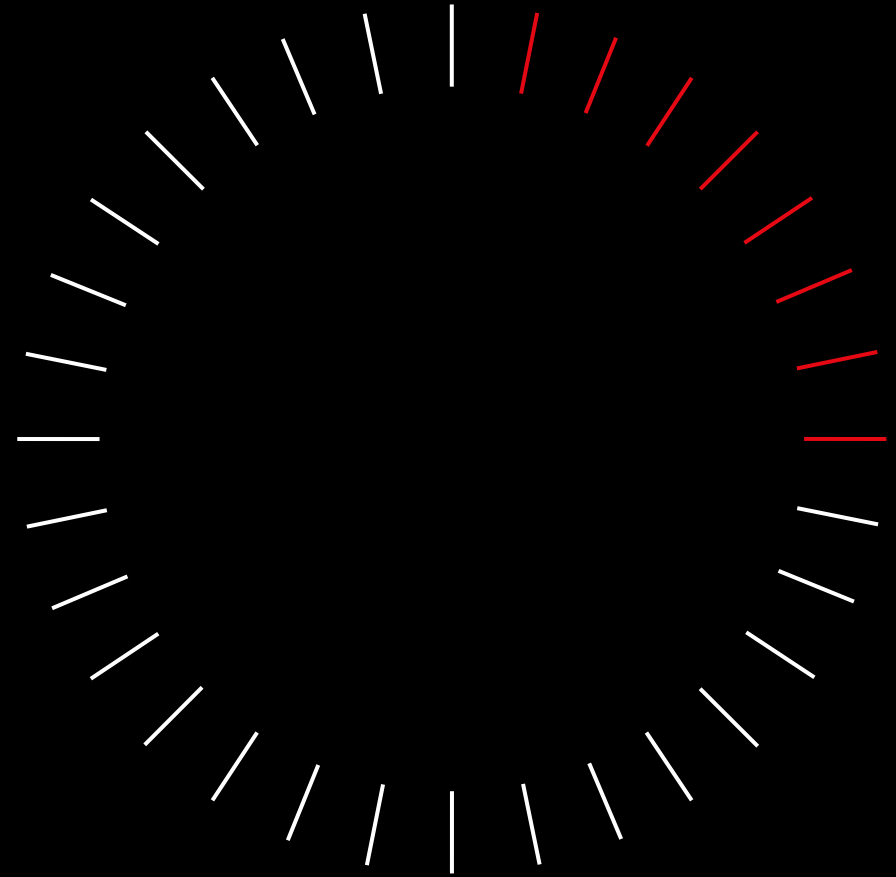
Need more levels!



**We need something better.**

# Target Overlap Compaction Strategy

Only do useful  
work, fewer rules



# Goals

Sorted runs

Configurable read vs write amplification with one tunable

Capable of “full” compaction without significant extra disk space

# Terminology

$size(S_i)$  = Either compressed or uncompressed size

$range(S_i)$  = What range of tokens does the SSTable span?

$level(S_i)$  = Non-overlapping guarantee

Sorted runs of SSTables

$SS_i \leftarrow level(S_i)$

$density(SS_i) = \frac{\sum size(S_i)}{\sum range(S_i)}$

# High Level Idea

Aim for target number of reads per read (overlap)

Compaction only operates across sorted runs

Select candidates using `IntervalTree`, bucket by density. Once dense enough start de-overlapping

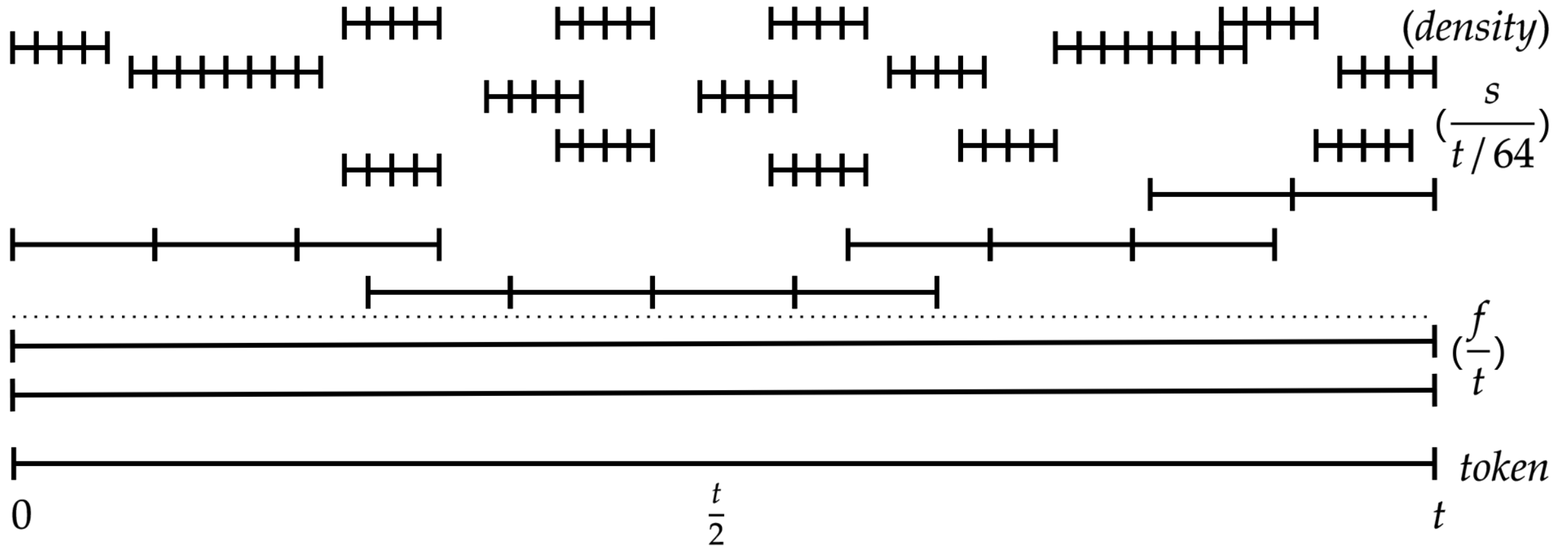
# Target Overlap Configuration

```
CREATE TABLE ... (  
    ...  
) WITH compaction = {  
    'class' : 'TargetOverlapCompactionStrategy',  
    # What is your target "read/read"?  
    'target_overlap': 4  
};
```

# Target Overlap Configuration

```
compaction = {  
    # Standard STCS things  
    'min_threshold': 2,  
    'max_threshold': 32,  
    'min_sstable_size_in_mb': 50,  
    # Target read per read  
    'target_overlap': 4,  
    # How big should my segments be  
    'target_sstable_size_in_mb': 256,  
    # How many reads per read can I get?  
    'max_overlap': target_overlap * min_threshold,  
    # How many sstables can I handle max per node (e.g. due to backup)  
    'max_sstable_count': 2000  
    # How old should SSTables be allowed to get? 0 for unlimited  
    'target_deadline': '10d',  
};
```

# Let the Levels\* Run Free



\* Now known as “Sorted Runs”



# Algorithm: $SS_0$ ( $\approx$ Level Zero)

Flushes are low quality sorted runs

Pure density tiering ( $\approx$  size tiering)

Once we have  $\text{target\_size} * \text{max\_count}$  data,  
“promote” to a sorted run  $SS_{n+1}$

# Algorithm: $SS_n$

*let  $O_t$  = target overlap*  
*let  $O_m$  = max overlap*  
*let  $D$  = deadline*  
*let  $F$  = tier factor*

$C \leftarrow \{\}$   
 $\forall$  interval  $I \in \text{range}()$ :  
     $S \leftarrow \text{overlapping}(I)$   
     $SS \leftarrow \text{group}(S, \text{level}(S))$   
    if  $|SS| > O_m$ :  
         $C.\text{insert}(\text{overlapping}(SS))$   
    elif  $\text{age}(s \in SS) > D$ :  
         $C.\text{insert}(\text{overlapping}(s))$   
    elif  $|SS| > O_t$ :  
         $B_d \leftarrow \text{bucket}(SS, \text{density})$   
         $\forall b \in B_d$  s. t.  $|d| > F$   
             $C.\text{insert}(*b, \text{score}(d))$

# Algorithm: $SS_n$

Iterate intervals in range

Group by sorted run

$C \leftarrow \{\}$

$\forall \text{ interval } I \in \text{range}():$   
     $S \leftarrow \text{overlapping}(I)$   
     $SS \leftarrow \text{group}(S, \text{level}(S))$

    if  $|SS| > O_m$ :

$C.\text{insert}(\text{overlapping}(SS))$

    elif  $\text{age}(s \in SS) > D$ :

$C.\text{insert}(\text{overlapping}(s))$

    elif  $|SS| > O_t$ :

$B_d \leftarrow \text{bucket}(SS, \text{density})$

$\forall b \in B_d \text{ s.t. } |d| > F$

$C.\text{insert}(*b, \text{score}(d))$

# Algorithm: $SS_n$

If we have enough overlapping ranges, force a compaction

```
C ← {}  
∀ interval I ∈ range():  
    S ← overlapping(I)  
    SS ← group(S, level(S))  
    if |SS| > Om:  
        C.insert(overlapping(SS))  
    elif age(s ∈ SS) > D :  
        C.insert(overlapping(s))  
    elif |SS| > Ot:  
        Bd ← bucket(SS, density)  
        ∀ b ∈ Bd s. t. |d| > F  
            C.insert(*b, score(d))
```

# Algorithm: $SS_n$

If any SSTables are too old, compact all overlapping with that SSTable

Probably with jitter...

```
C ← {}  
∀ interval I ∈ range():  
    S ← overlapping(I)  
    SS ← group(S, level(S))  
    if |SS| > Om:  
        C.insert(overlapping(SS))  
    elif age(s ∈ SS) > D :  
        C.insert(overlapping(s))  
    elif |SS| > Ot:  
        Bd ← bucket(SS, density)  
        ∀ b ∈ Bd s. t. |d| > F  
            C.insert(*b, score(d))
```

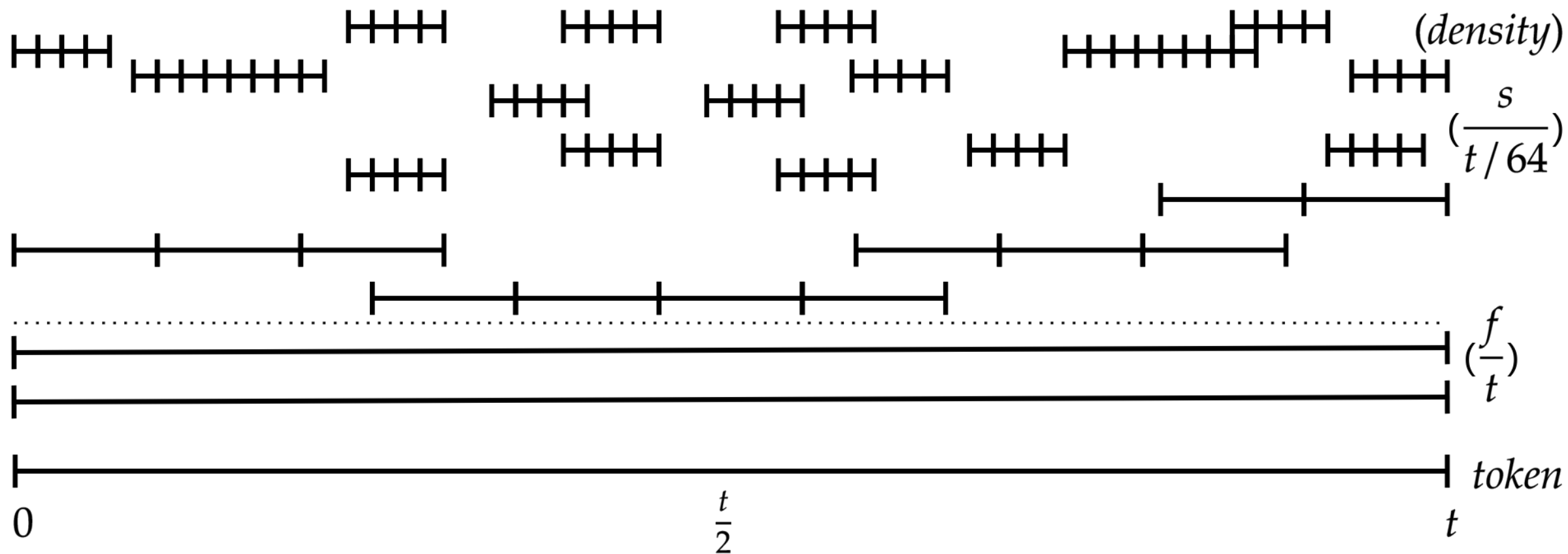
# Algorithm: $SS_n$

Density tiering within the sorted runs.

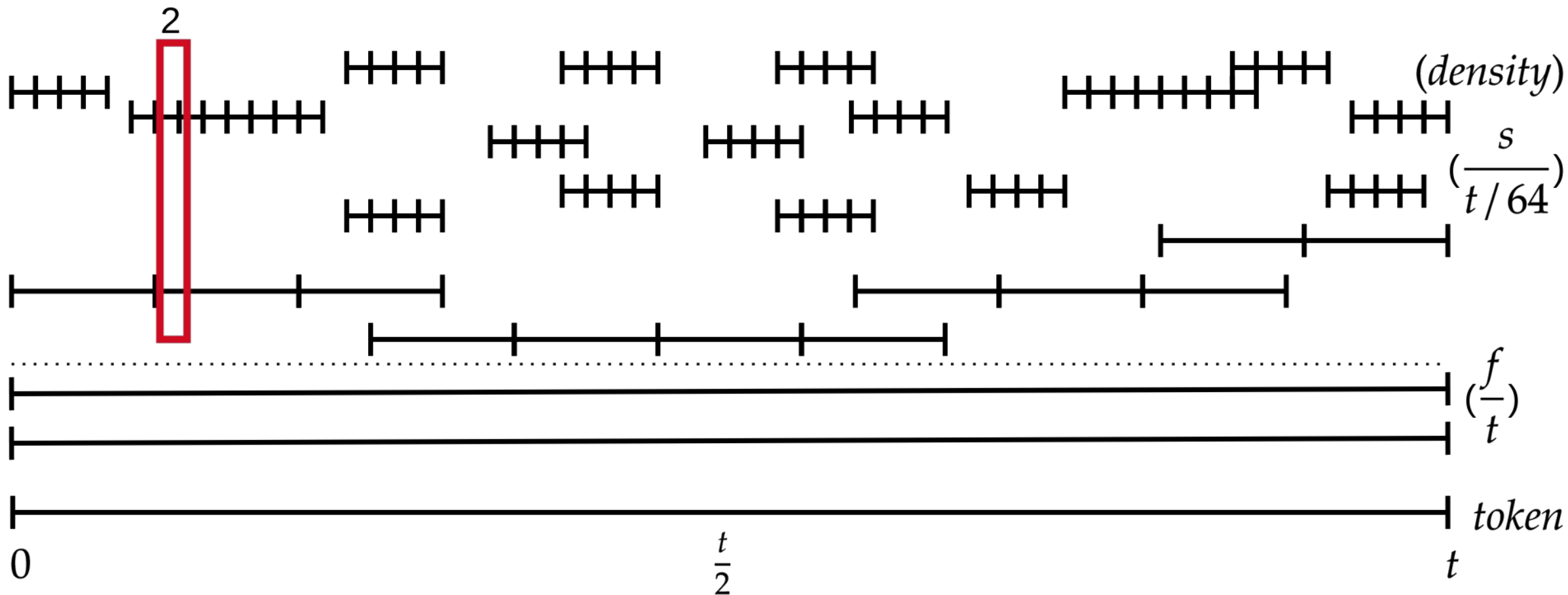
Prefer larger buckets first, break ties with density

```
C ← {}  
∀ interval I ∈ range():  
    S ← overlapping(I)  
    SS ← group(S, level(S))  
    if |SS| > Om:  
        C.insert(overlapping(SS))  
    elif age(s ∈ SS) > D :  
        C.insert(overlapping(s))  
    elif |SS| > Ot:  
        Bd ← bucket(SS, density)  
        ∀ b ∈ Bd s. t. |d| > F  
            C.insert(*b, score(d))
```

# TOCS In Action

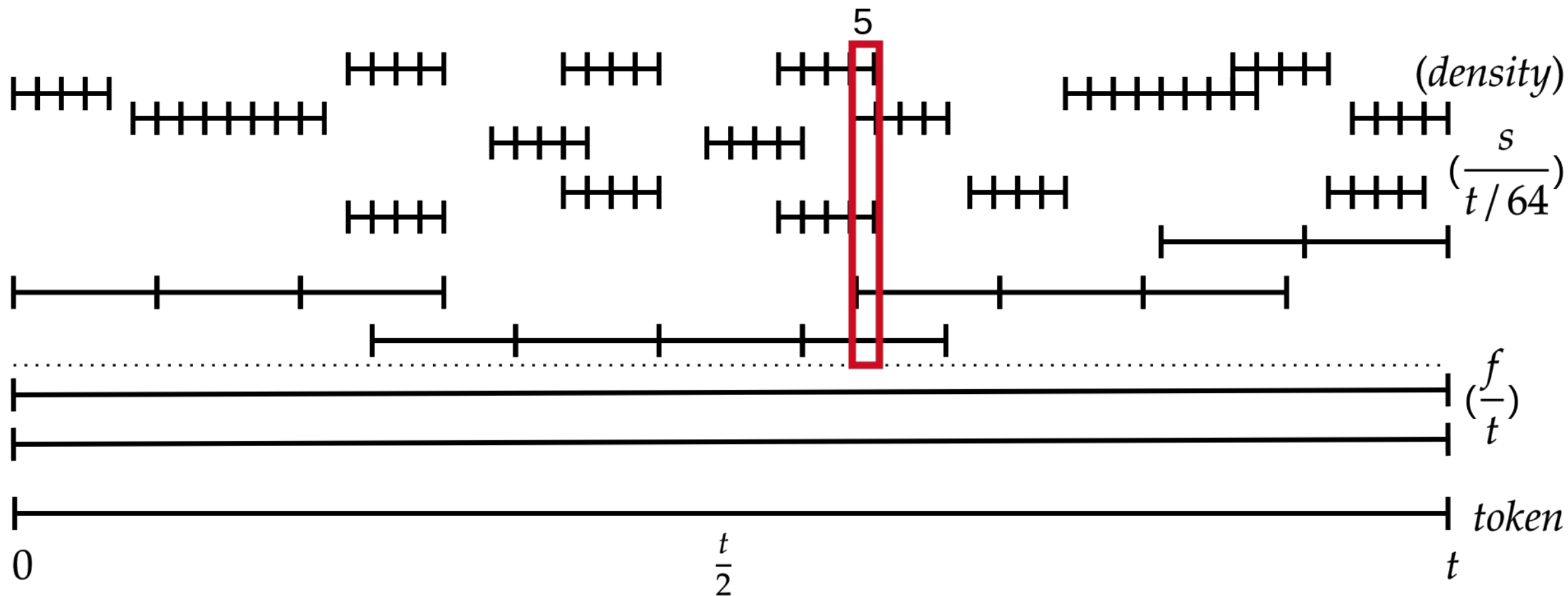


# TOCS In Action

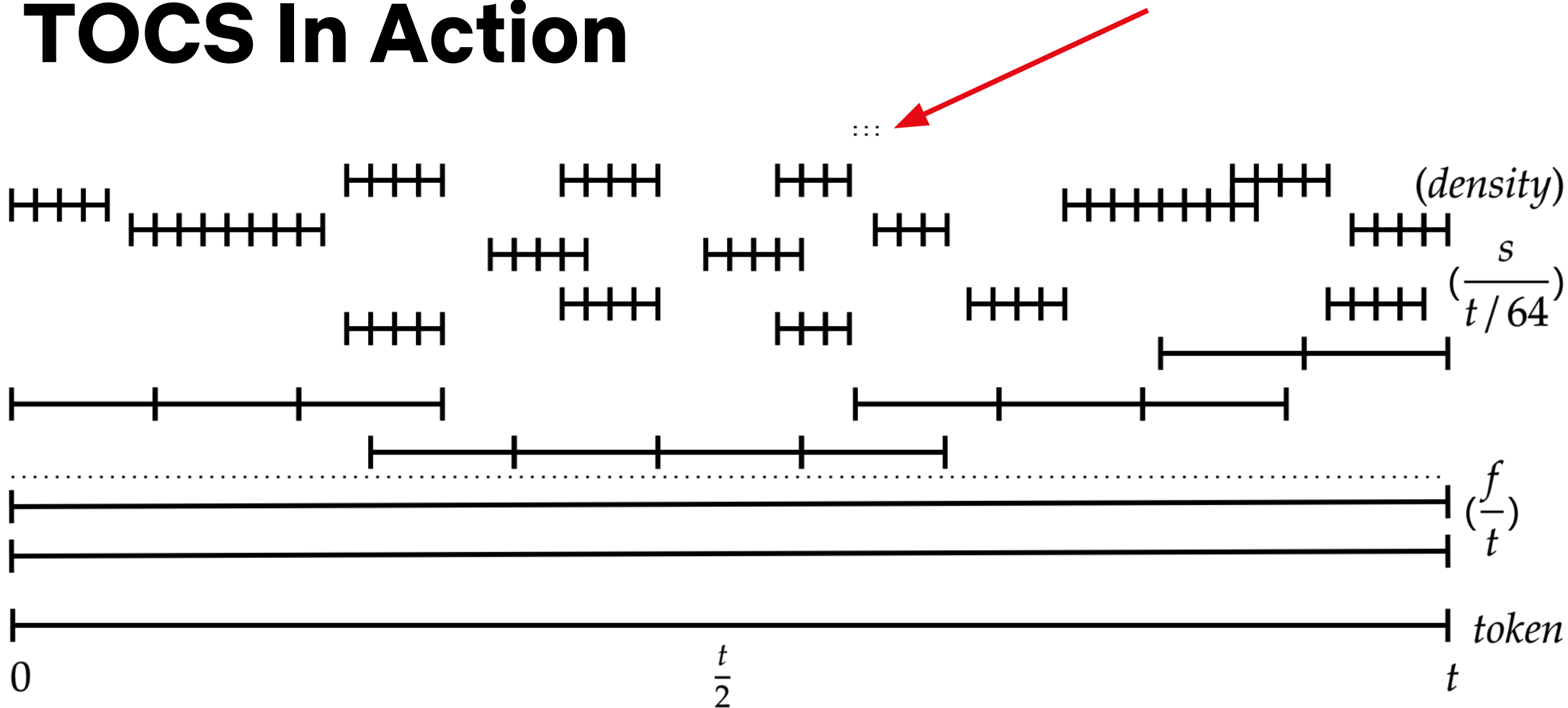




# TOCS In Action



# TOCS In Action



# Property: Easy to Use

Single main tunable, `target_overlap`, allows easy to understand read vs write amplification tradeoff

Secondary tunable, `target_deadline`, allows easy to understand tombstone / update properties

# Property: Low Overhead

Largest compaction will be  $\sim 1/32$  of data size

Full compactions implemented through declarative desire rather than imperative command.

# Property: Small predictable files

Since the strategy only produces sorted runs of tables about 256 MiB large, it is easy on:

1. Backup
2. Streaming (full sstable streaming)
3. Page cache (no more early re-open)

# Still TODO

It's not clear density tiering in the top levels is worth the complexity. Just do the whole vertical every time?

Should we split sorted runs into 32 ranges on promotion?

# Take Away

Target Overlap (name?) is a potential way to unify our learnings from LCS and STCS.

# Thank You.



Joey Lynch  
josephl@netflix.com