

Solr sparse faceting

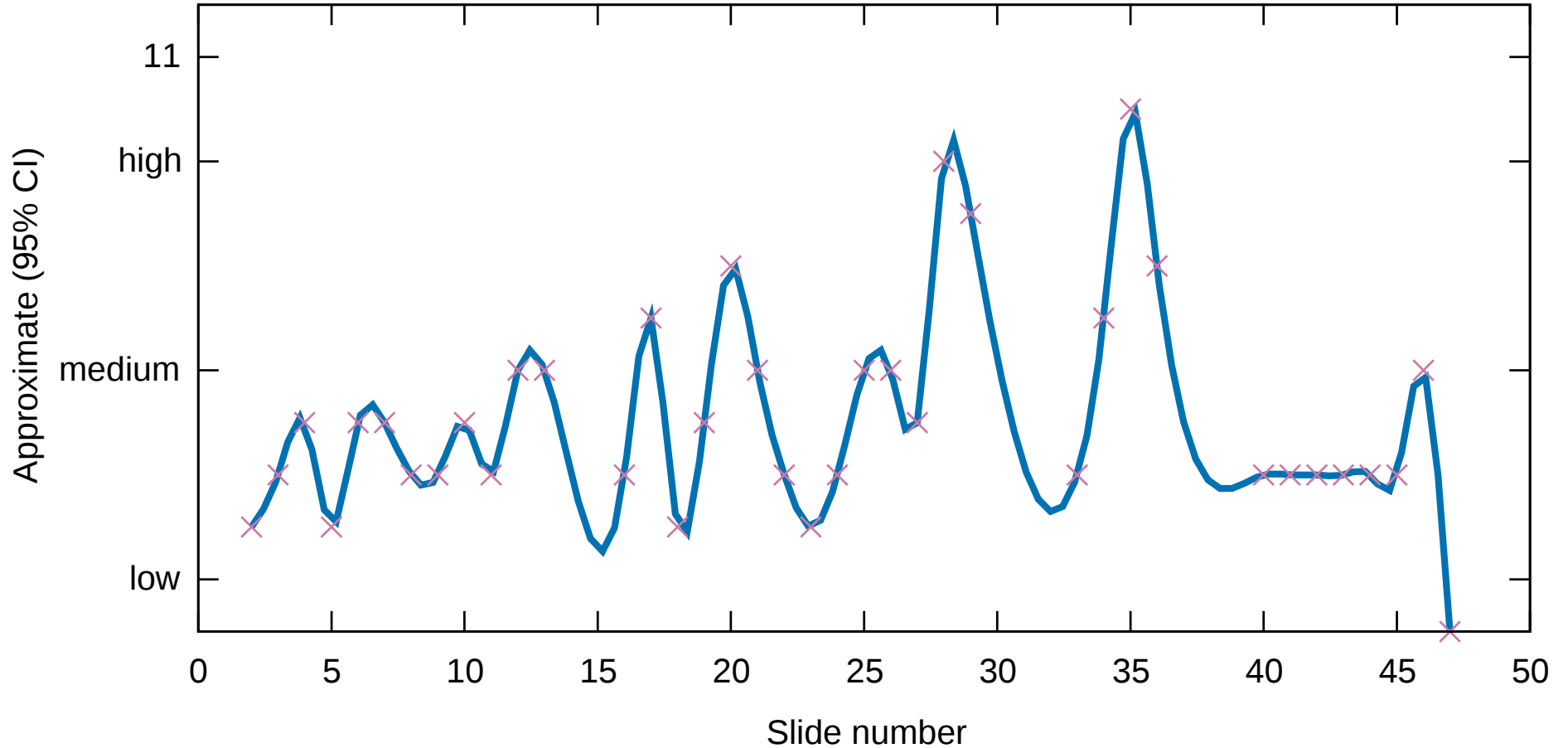
Everything counts in large amounts

@TokeEskildsen

State and University Library, Denmark

<https://tokee.github.io/lucene-solr/>

Presentation tech level



Nothing To Fear

- 500TB+ web resources from Danish Net Archive
- Estimated 50TB Solr index data when finished
- 3 machines of 16 CPU cores, 256GB RAM, 25 * 900GB SSD
 - Each machine holds: 25 Solrs
 - Each Solr holds: 1 optimized shard with 900GB / 250M docs
- Shards build externally, one at a time
- (Optimizations also relevant for smaller setups)

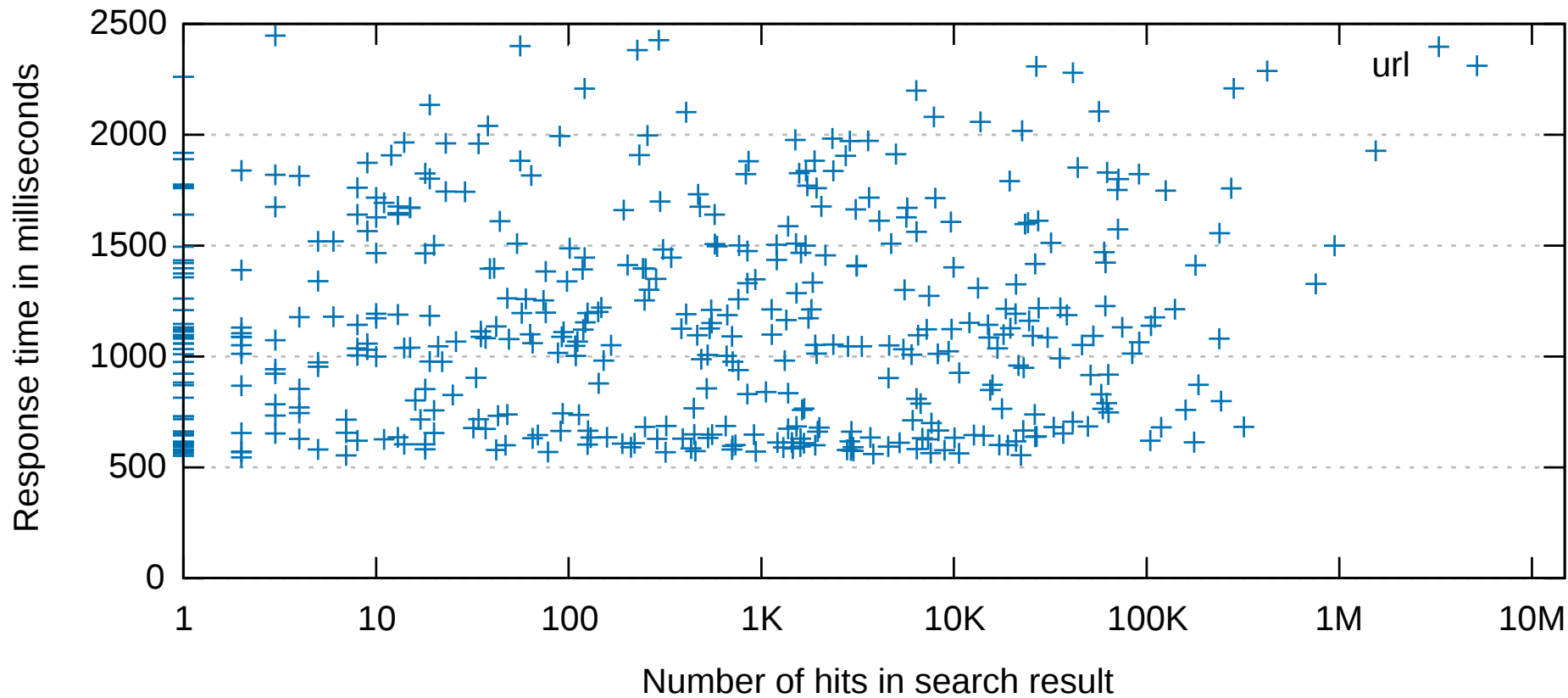
Pipeline

```
counter = new int[ordinals]
for docID: result.getDocIDs()
    for ordinal: getOrdinals(docID)
        counter[ordinal]++

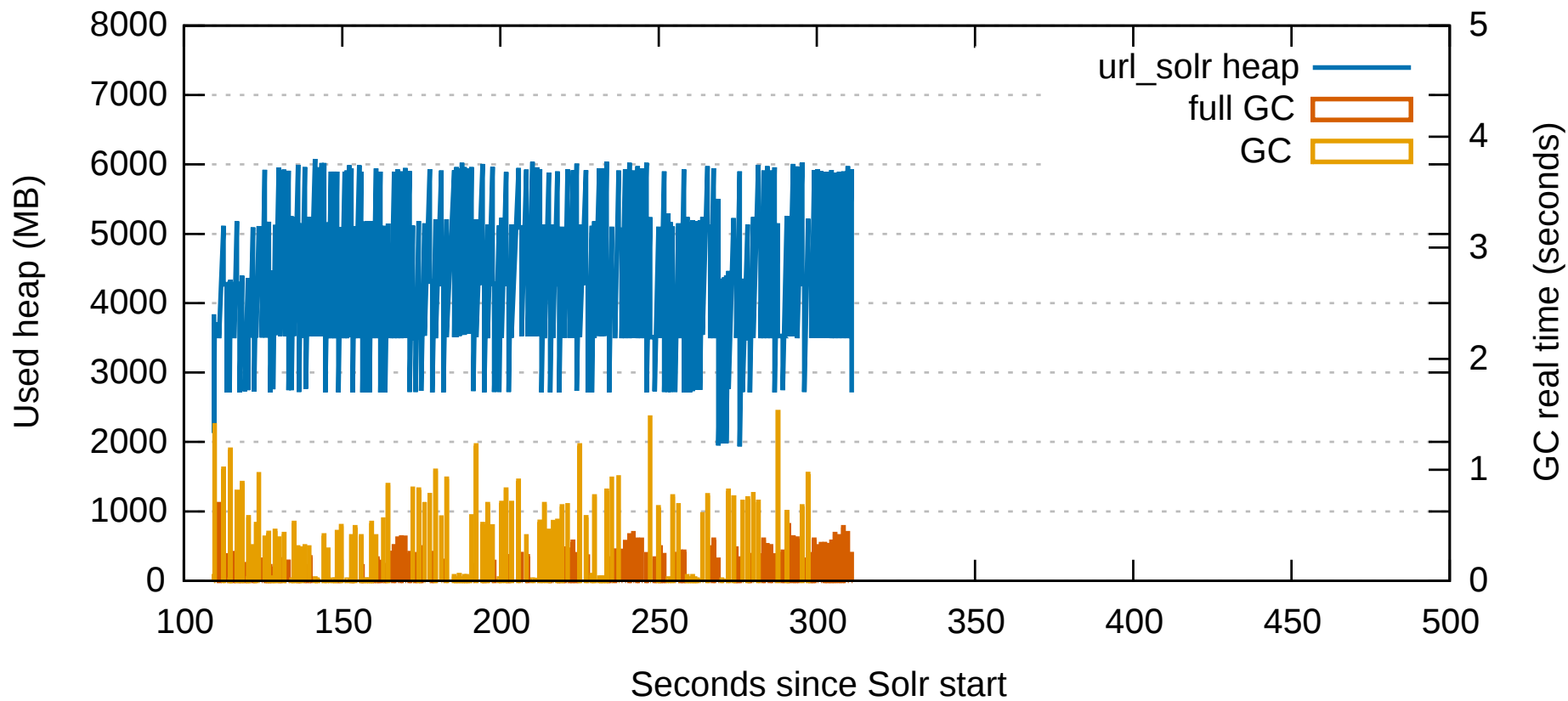
for ordinal = 0 ; ordinal < counters.length ; ordinal++
    priorityQueue.add(ordinal, counter[ordinal])
```

ord	term	counter
0	A	0
1	B	3
2	C	0
3	D	1006
4	E	1
5	F	1
6	G	0
7	H	0
8	I	3

1 shard / 900GB / 250M docs, facet url 200M values



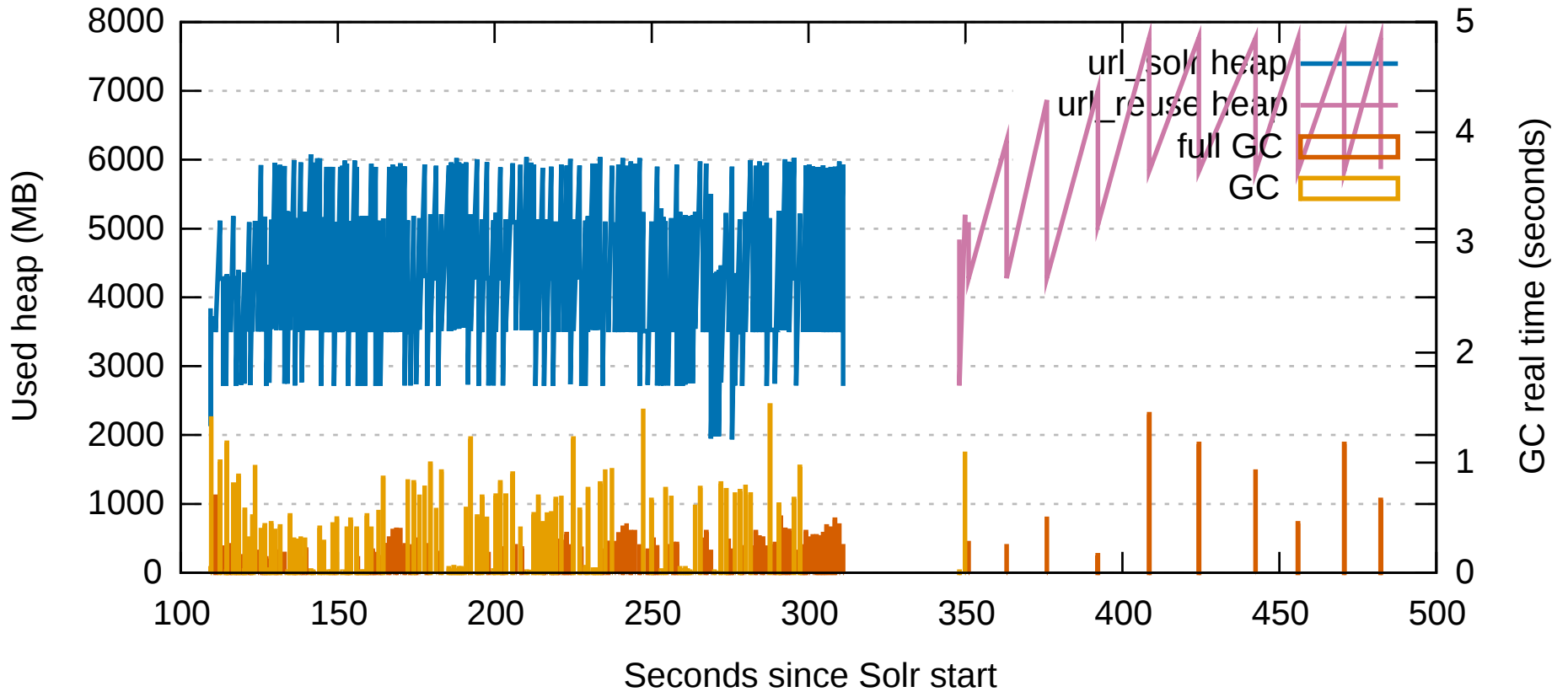
Garbage collections with 3 concurrent requests on field url



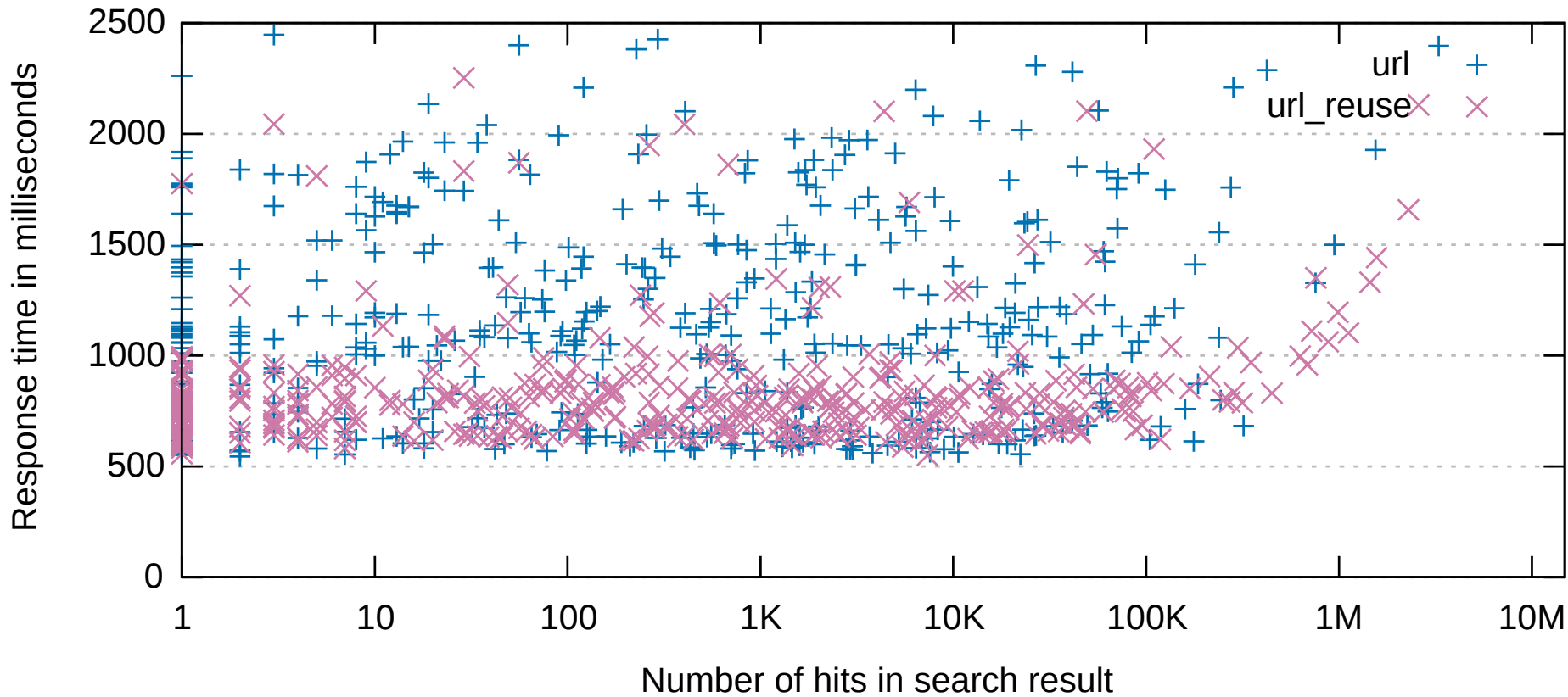
Recycle

```
counter = pool.getCounter()  
for docID: result.getDocIDs()  
    for ordinal: getOrdinals(docID)  
        counter[ordinal]++  
  
for ordinal = 0 ; ordinal < counters.length ; ordinal++  
    priorityQueue.add(ordinal, counter[ordinal])  
  
pool.release(counter)
```

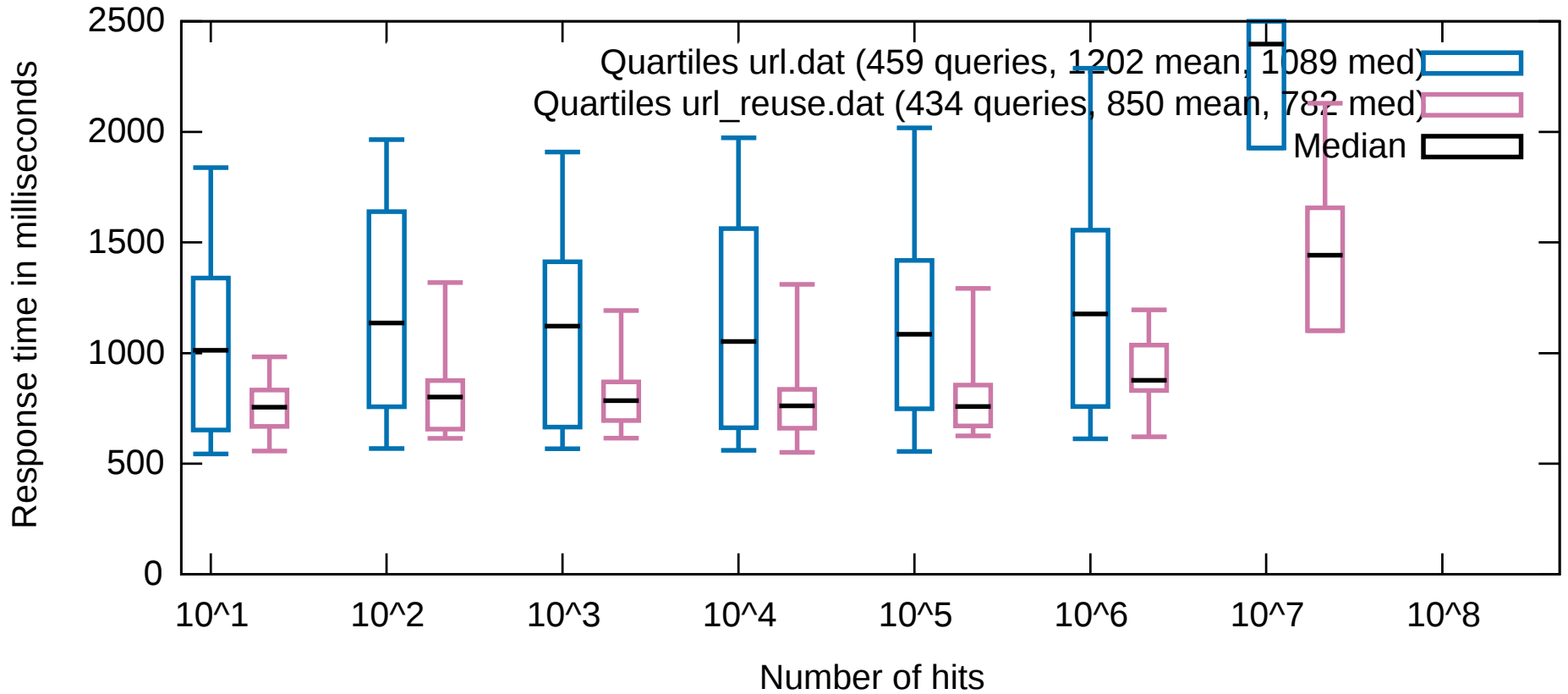
Garbage collections with 3 concurrent requests on field url



1 shard / 900GB / 250M docs, facet url 200M values



Quartiles for response times with top as 95 percentile, 893 samples, 2015-05-29



Counting

```
counter = pool.getCounter()
for docID: result.getDocIDs()
    for ordinal: getOrdinals(docID)
        counter[ordinal]++

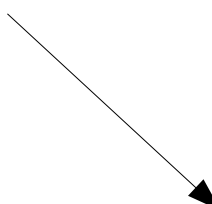
for ordinal = 0 ; ordinal < counters.length ; ordinal++
    priorityQueue.add(ordinal, counter[ordinal])

pool.release(counter)
```

ord	term	counter
0	A	0
1	B	3
2	C	0
3	D	1006
4	E	1
5	F	1
6	G	0
7	H	0
8	I	3

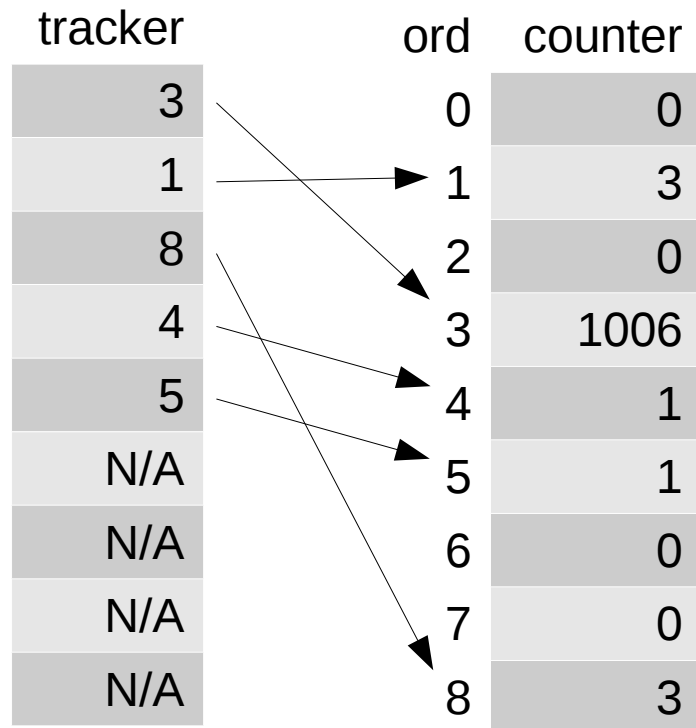
tracker	ord	counter
N/A	0	0
N/A	1	0
N/A	2	0
N/A	3	0
N/A	4	0
N/A	5	0
N/A	6	0
N/A	7	0
N/A	8	0

tracker	ord	counter
3	0	0
N/A	1	0
N/A	2	0
N/A	3	1
N/A	4	0
N/A	5	0
N/A	6	0
N/A	7	0
N/A	8	0



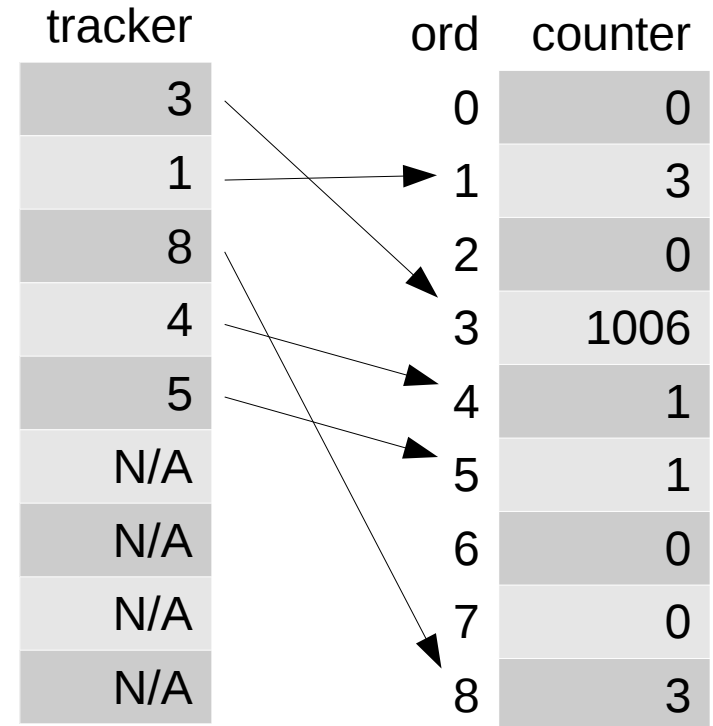
tracker	ord	counter
3	0	0
1	1	1
N/A	2	0
N/A	3	1
N/A	4	0
N/A	5	0
N/A	6	0
N/A	7	0
N/A	8	0

tracker	ord	counter
3	0	0
1	1	1
N/A	2	0
N/A	3	2
N/A	4	0
N/A	5	0
N/A	6	0
N/A	7	0
N/A	8	0

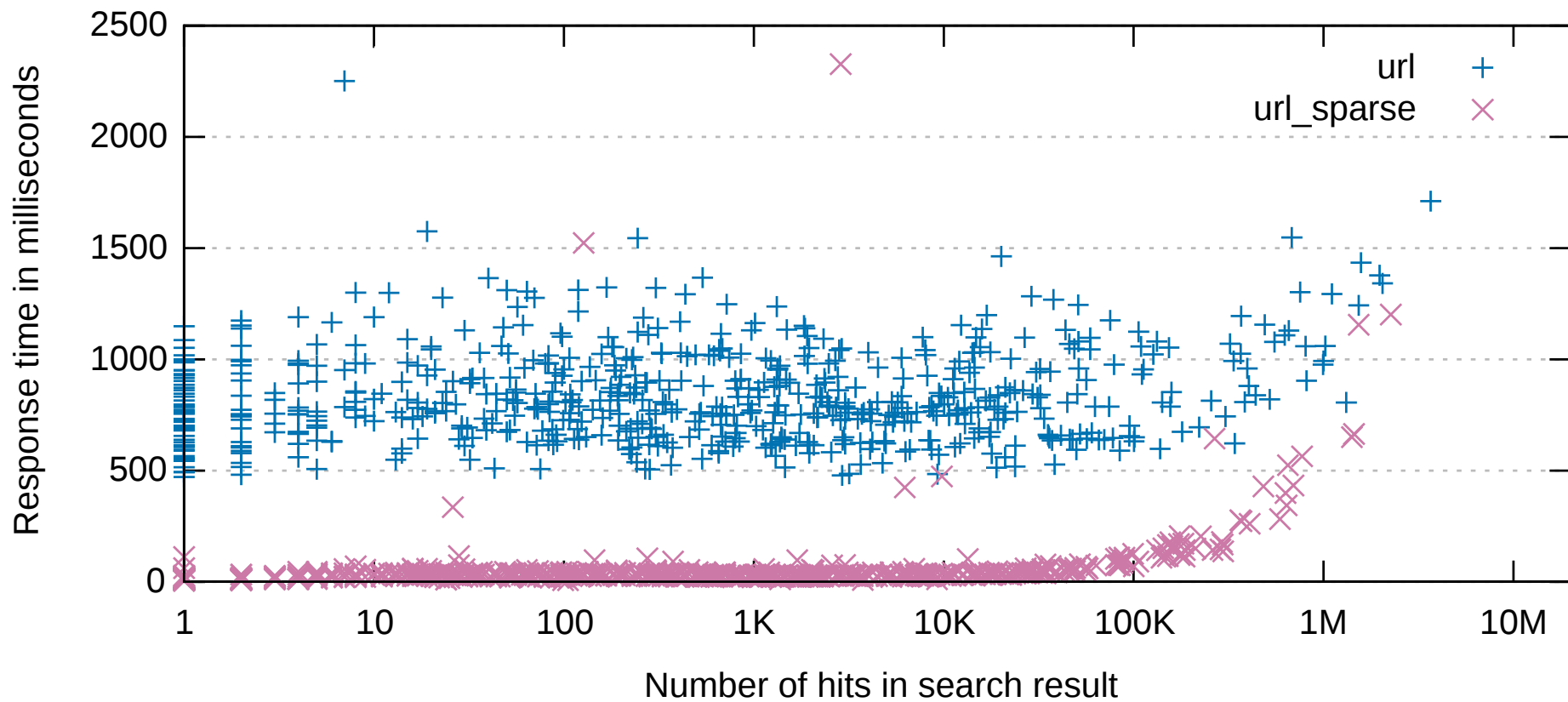


Sparse counting

```
counter = pool.getCounter()
for docID: result.getDocIDs()
  for ordinal: getOrdinals(docID)
    if counter[ordinal]++ == 0 && tracked < maxTracked
      tracker[tracked++] = ordinal
if tracked < maxTracked
  for i = 0 ; i < tracked ; i++
    priorityQueue.add(tracker[i], counter[tracker[i]])
else
  for ordinal = 0 ; ordinal < counter.length ; ordinal++
    priorityQueue.add(ordinal, counter[ordinal])
```



1 shard / 900GB / 250M docs, facet url 200M values



Get the Balance Right

Phase 1) All shards perform faceting

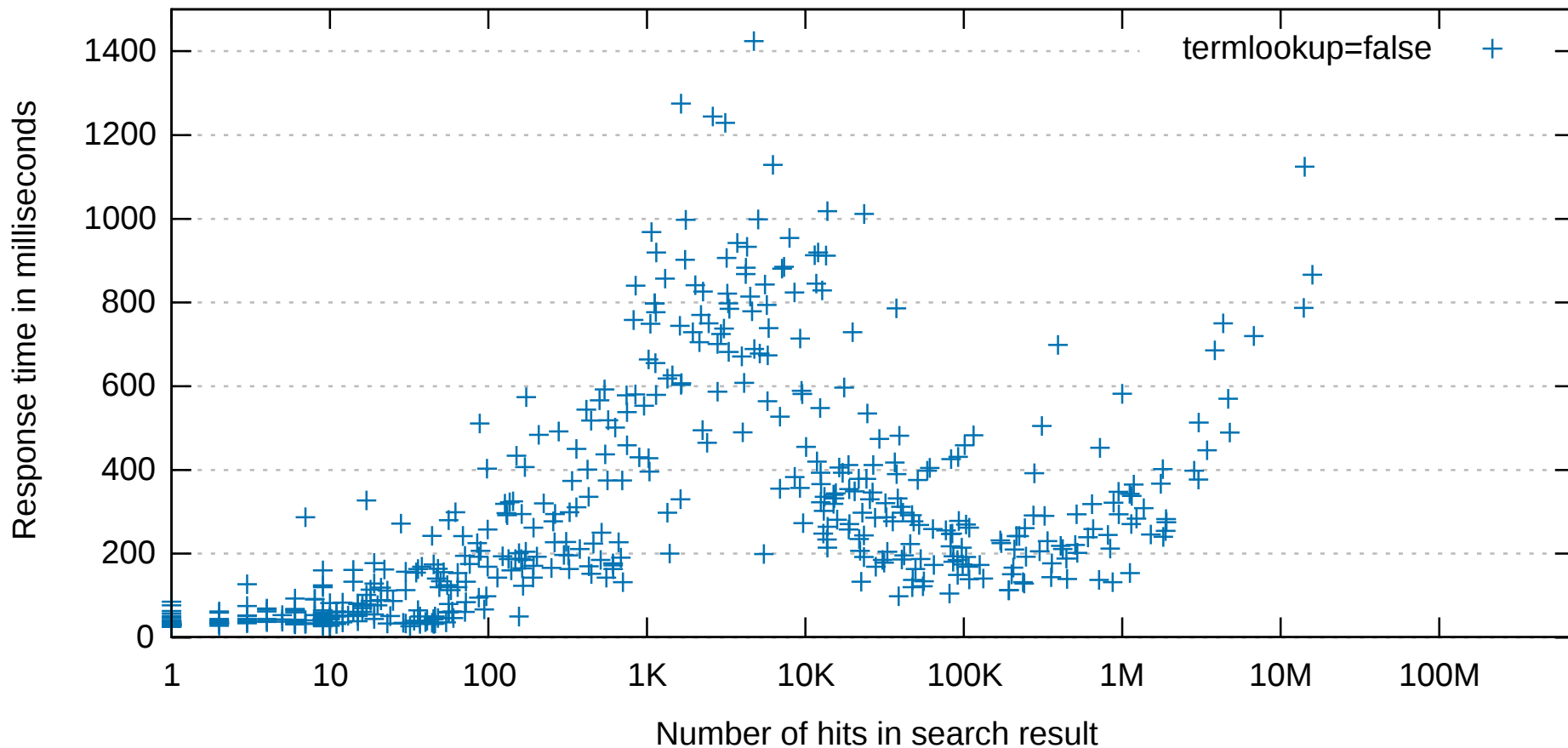
The Merger calculates top-X terms

Phase 2) The term counts are requested from the shards that did not return them in phase 1

```
for term: query.getTerms()  
    result.add(term, searcher.numDocs(  
        query(field:term), base.getDocIDs()  
    ).hitCount)
```

Pit of Pain™

9 shards / 7TB / 2.3G docs, 256GB RAM, 16 cores, facet 1.1M values/shard



Alternative fine counting

```
counter = pool.getCounter()
for docID: result.getDocIDs()
    for ordinal: getOrdinals(docID)
        counter.increment(ordinal)
```

} Same as phase 1

```
for term: query.getTerms()
    result.add(term, counter.get(getOrdinal(term)))
```

Stripped

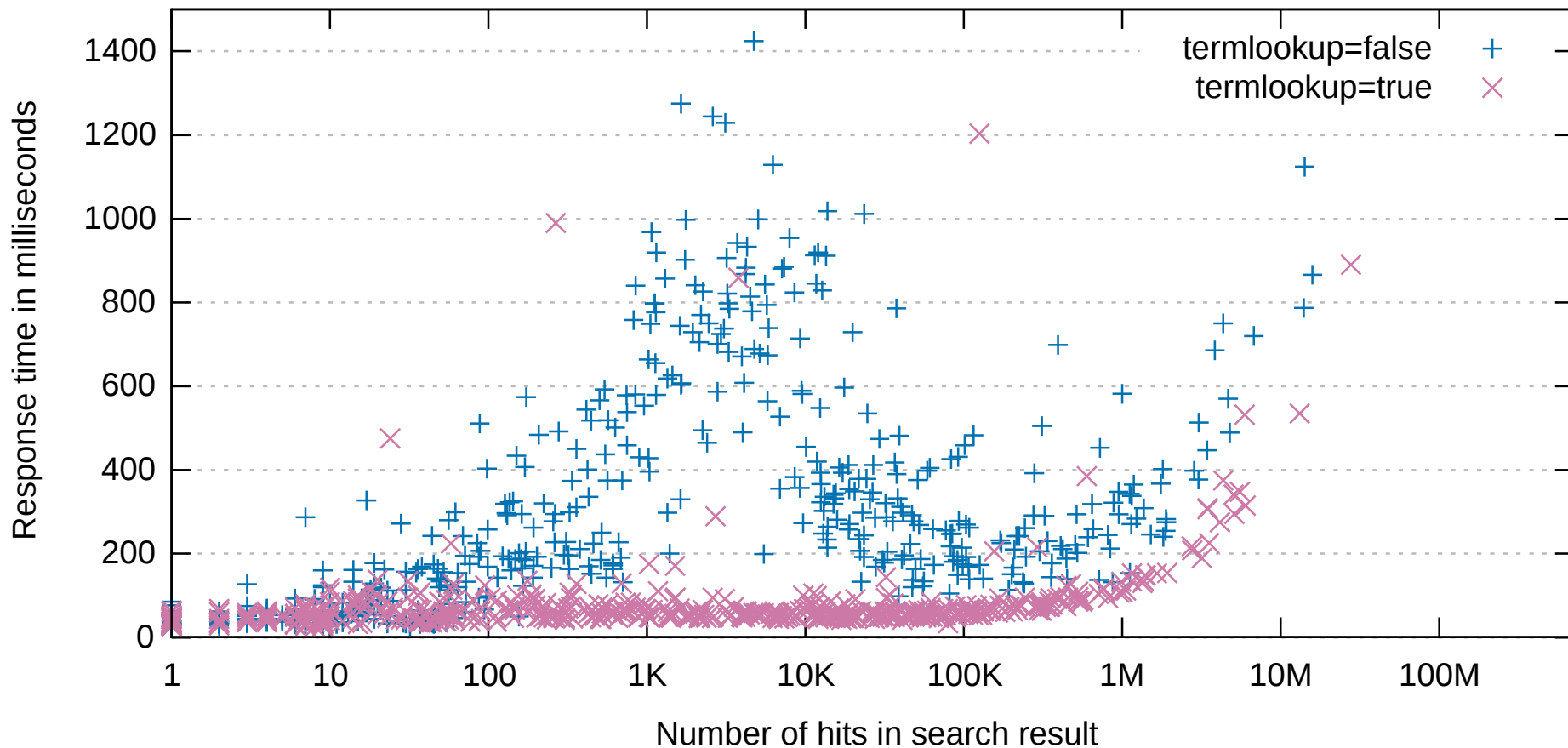
```
counter = pool.getCounter(key)
```

```
for term: query.getTerms()
```

```
    result.add(term, counter.get(getOrdinal(term)))
```

Plain of Platitude™

9 shards / 7TB / 2.3G docs, 256GB RAM, 16 cores, facet 1.1M values/shard

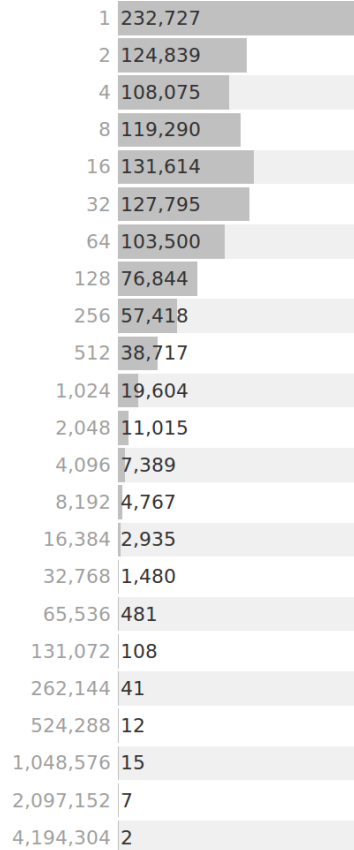


250,000,000 docs / 900GB, optimized

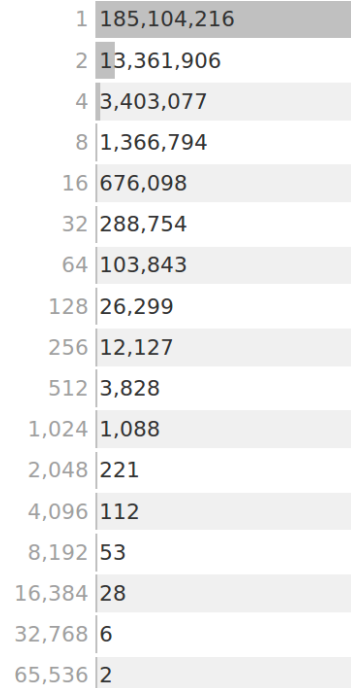
Field	References	Max docs/term	Terms
domain	250,000,000	3,000,000	1,100,000
url	250,000,000	56,000	200,000,000
links	5,800,000,000	5,000,000	610,000,000

Term distributions

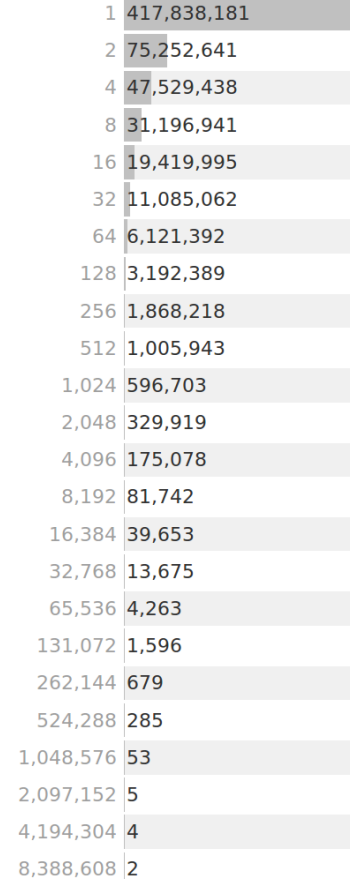
domain 1.1M



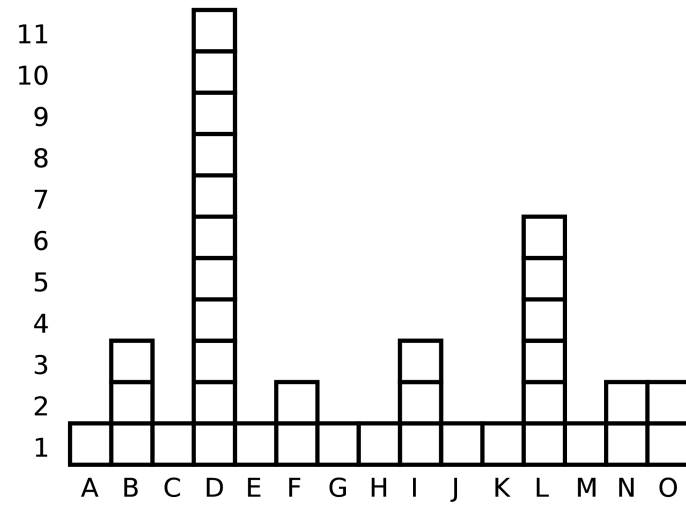
url 200M



links 600M

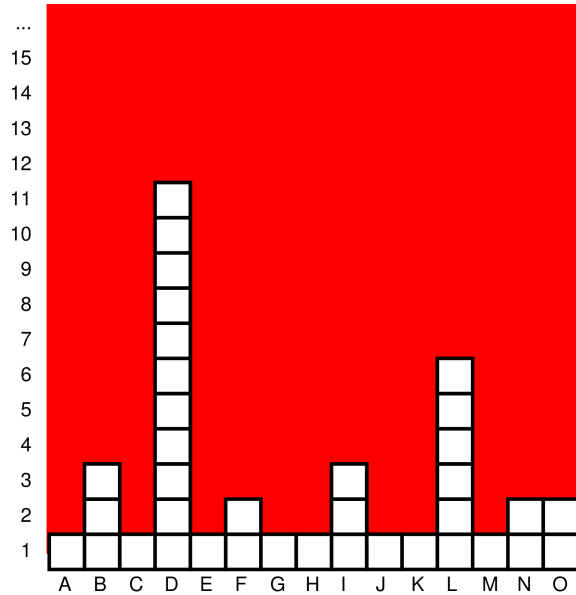


Clean



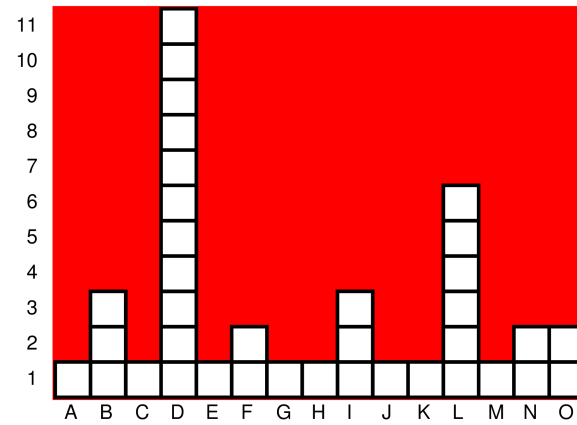
World Full Of Nothing

int[ordinals]



domain: 4 MB
url: 780 MB
links: 2350 MB

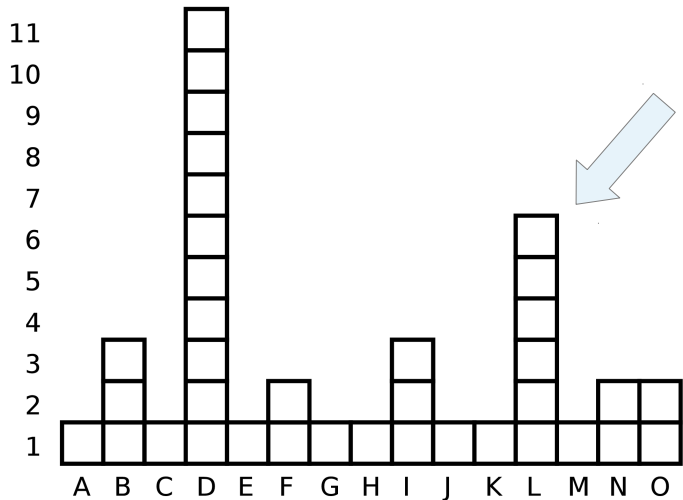
PackedInts(ordinals, maxBPV)



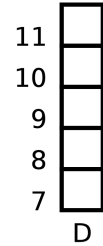
domain: 3 MB (72%)
url: 420 MB (53%)
links: 1760 MB (75%)

Construction Time Again

Platonic ideal

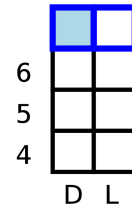


Plane 4

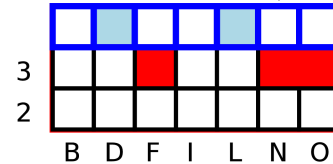


Harsh reality

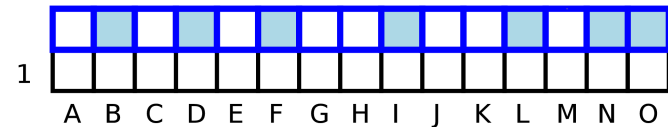
Plane 3



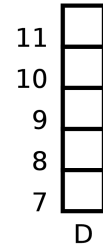
Plane 2



Plane 1

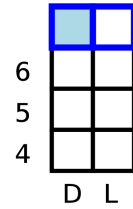


Plane 4

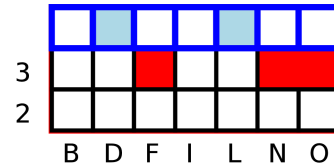


L: 0 ≡ 000000

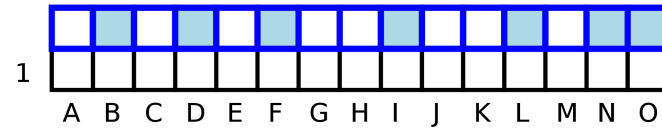
Plane 3



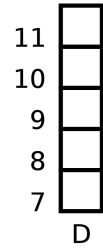
Plane 2



Plane 1

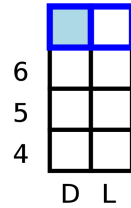


Plane 4

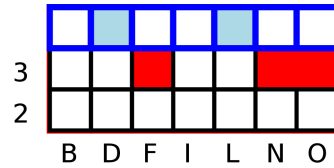


L: 0 ≡ 000000
L: 1 ≡ 000001

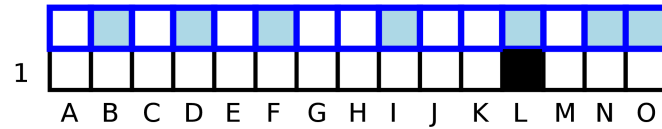
Plane 3



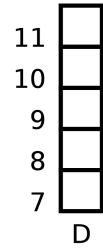
Plane 2



Plane 1

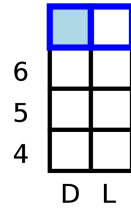


Plane 4

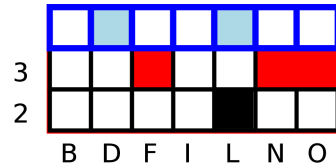


L: 0 ≡ 000000
L: 1 ≡ 000001
L: 2 ≡ 000010

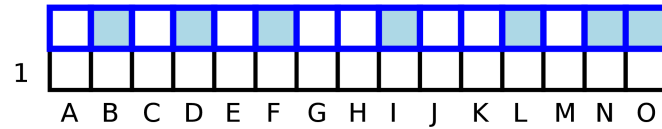
Plane 3



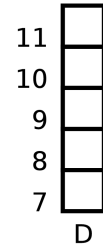
Plane 2



Plane 1

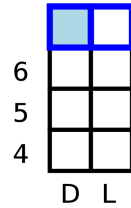


Plane 4

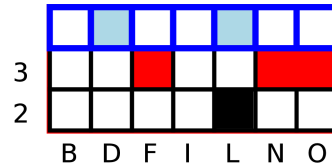


L: 0 ≡ 000000
L: 1 ≡ 000001
L: 2 ≡ 000010
L: 3 ≡ 000011

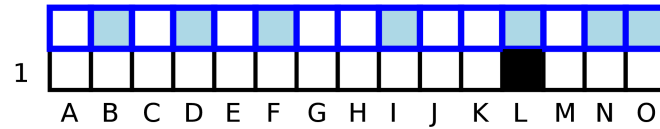
Plane 3



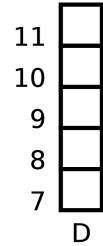
Plane 2



Plane 1

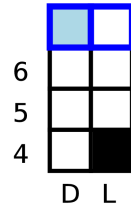


Plane 4



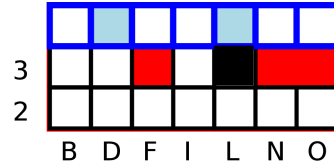
L: 0 ≡ 000000
 L: 1 ≡ 000001
 L: 2 ≡ 000010
 L: 3 ≡ 000011
 L: 4 ≡ 000100
 L: 5 ≡ 000101
 L: 6 ≡ 000110
 L: 7 ≡ 000111

Plane 3

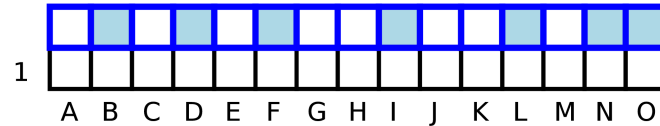


...
 L: 12 ≡ 001100

Plane 2



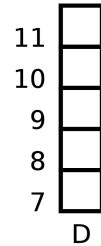
Plane 1



```
if counter[ordinal]++ == 0 && tracked < maxTracked
```

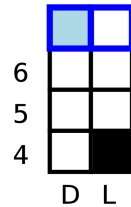
```
tracker[tracked++] = ordinal
```

Plane 4

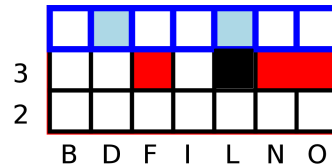


L: 0 ≡ 000000
L: 1 ≡ 000001
L: 2 ≡ 000010
L: 3 ≡ 000011
L: 4 ≡ 000100
L: 5 ≡ 000101
L: 6 ≡ 000110
L: 7 ≡ 000111
...
L: 12 ≡ 001100

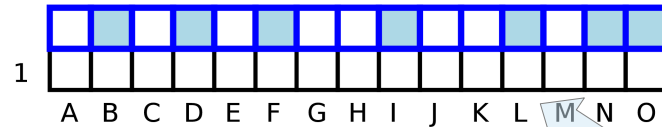
Plane 3



Plane 2

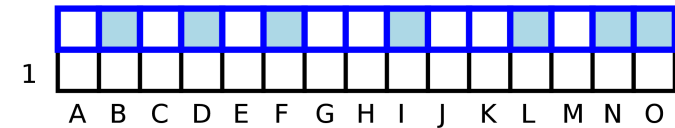
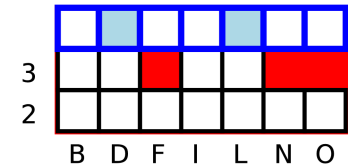
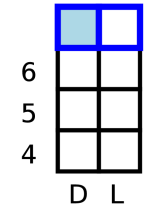
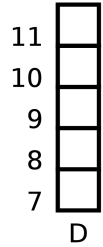


Plane 1

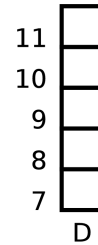


?

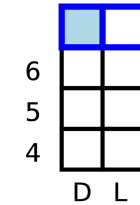
Now This is Fun



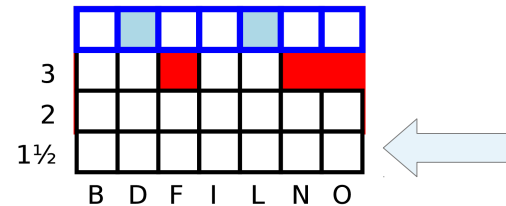
Plane 4



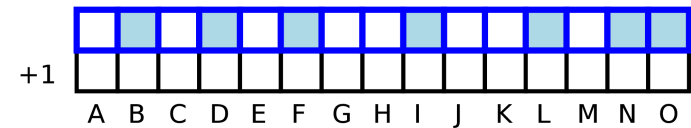
Plane 3



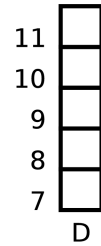
Plane 2



Plane 1

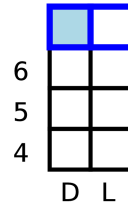


Plane 4

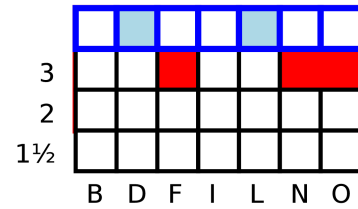


L: 0 ≡ 000000

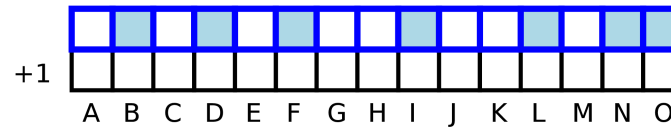
Plane 3



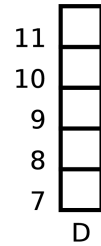
Plane 2



Plane 1

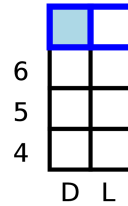


Plane 4

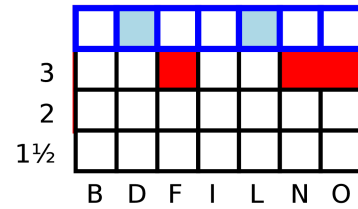


L: 0 ≡ 000000
L: 1 ≡ 000001

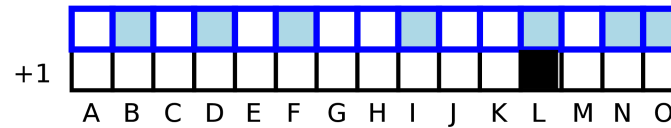
Plane 3



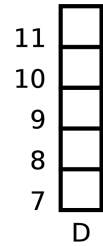
Plane 2



Plane 1

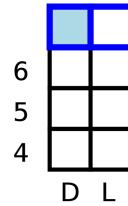


Plane 4

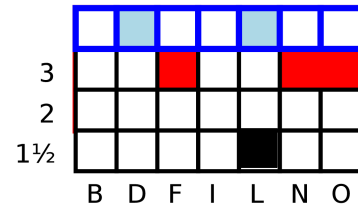


L: 0 ≡ 000000
L: 1 ≡ 000001
L: 2 ≡ 000011

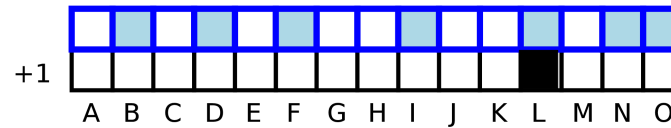
Plane 3



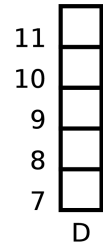
Plane 2



Plane 1

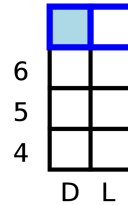


Plane 4

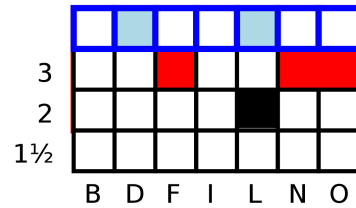


L: 0 ≡ 000000
L: 1 ≡ 000001
L: 2 ≡ 000011
L: 3 ≡ 000101

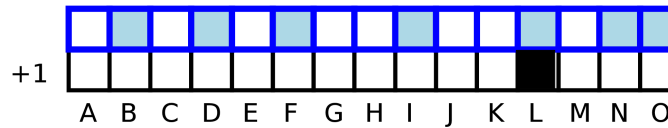
Plane 3



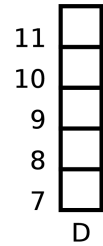
Plane 2



Plane 1

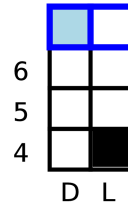


Plane 4



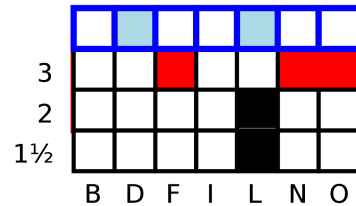
L: 0 ≡ 000000
 L: 1 ≡ 000001
 L: 2 ≡ 000011
 L: 3 ≡ 000101
 L: 4 ≡ 000111
 L: 5 ≡ 001001
 L: 6 ≡ 001011
 L: 7 ≡ 001101

Plane 3

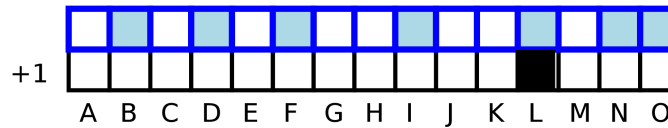


⋮
 L: 12 ≡ 010111

Plane 2



Plane 1



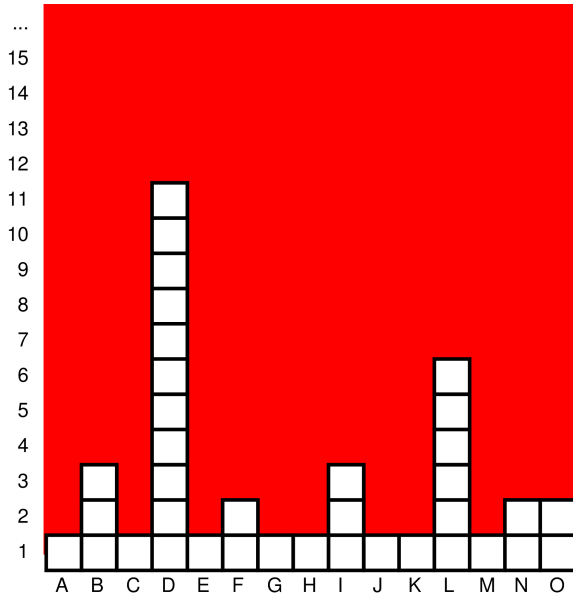
The Bottom Line

domain: 4 MB
 url: 780 MB
 links: 2350 MB

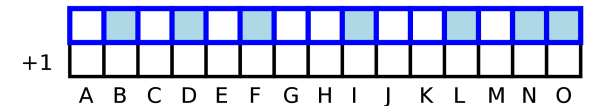
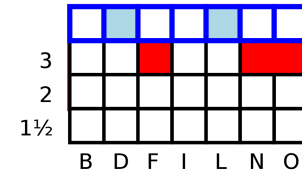
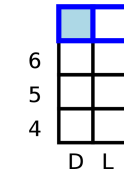
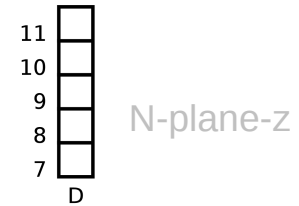
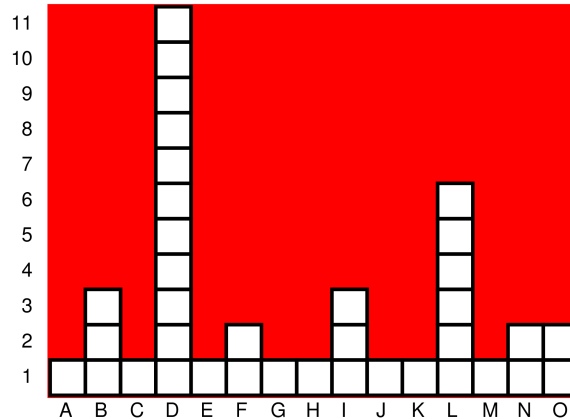
domain: 3 MB (72%)
 url: 420 MB (53%)
 links: 1760 MB (75%)

domain: 1 MB (30%)
 url: 66 MB (8%)
 links: 311 MB (13%)

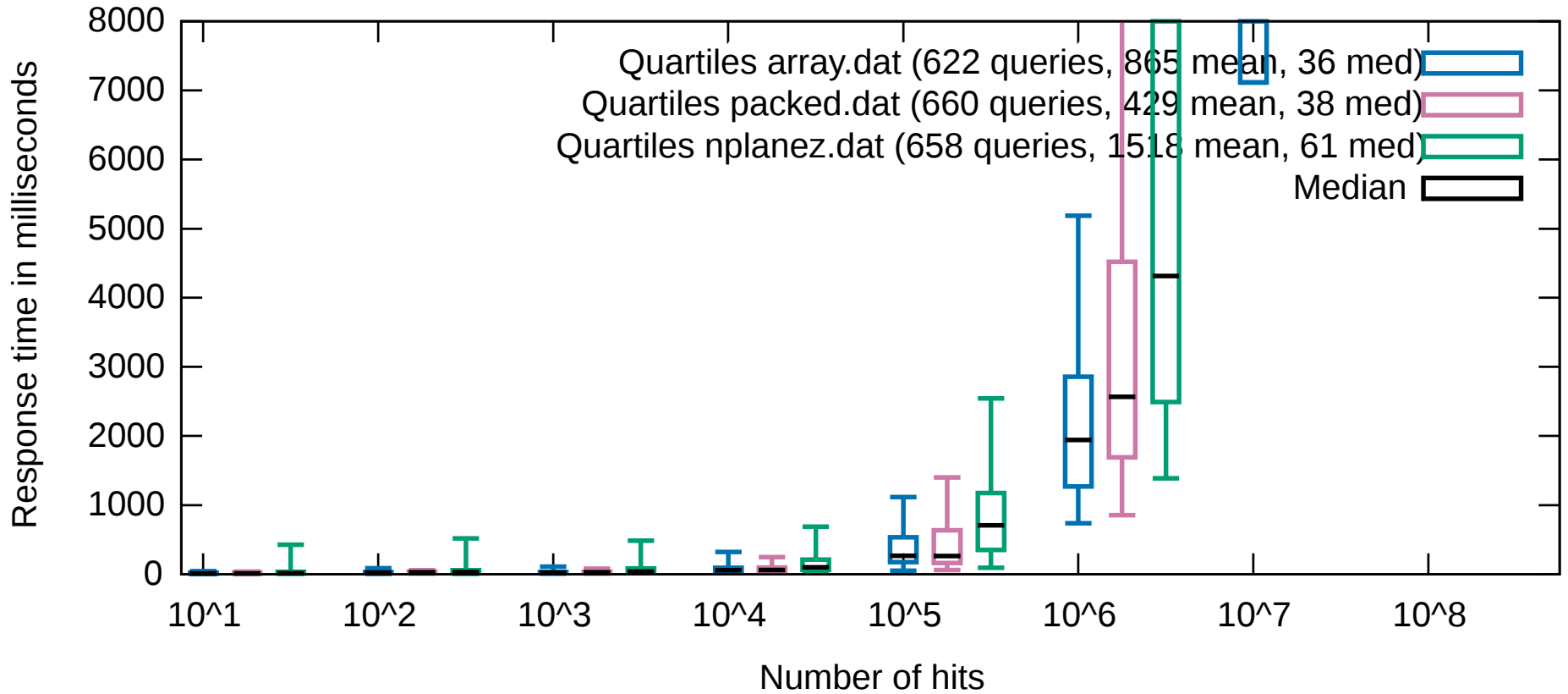
int[ordinals]



PackedInts(ordinals, maxBPV)



1 shard / 900GB / 250M docs, 256GB RAM, 16 cores, facet links 620M values



Kitchen sink

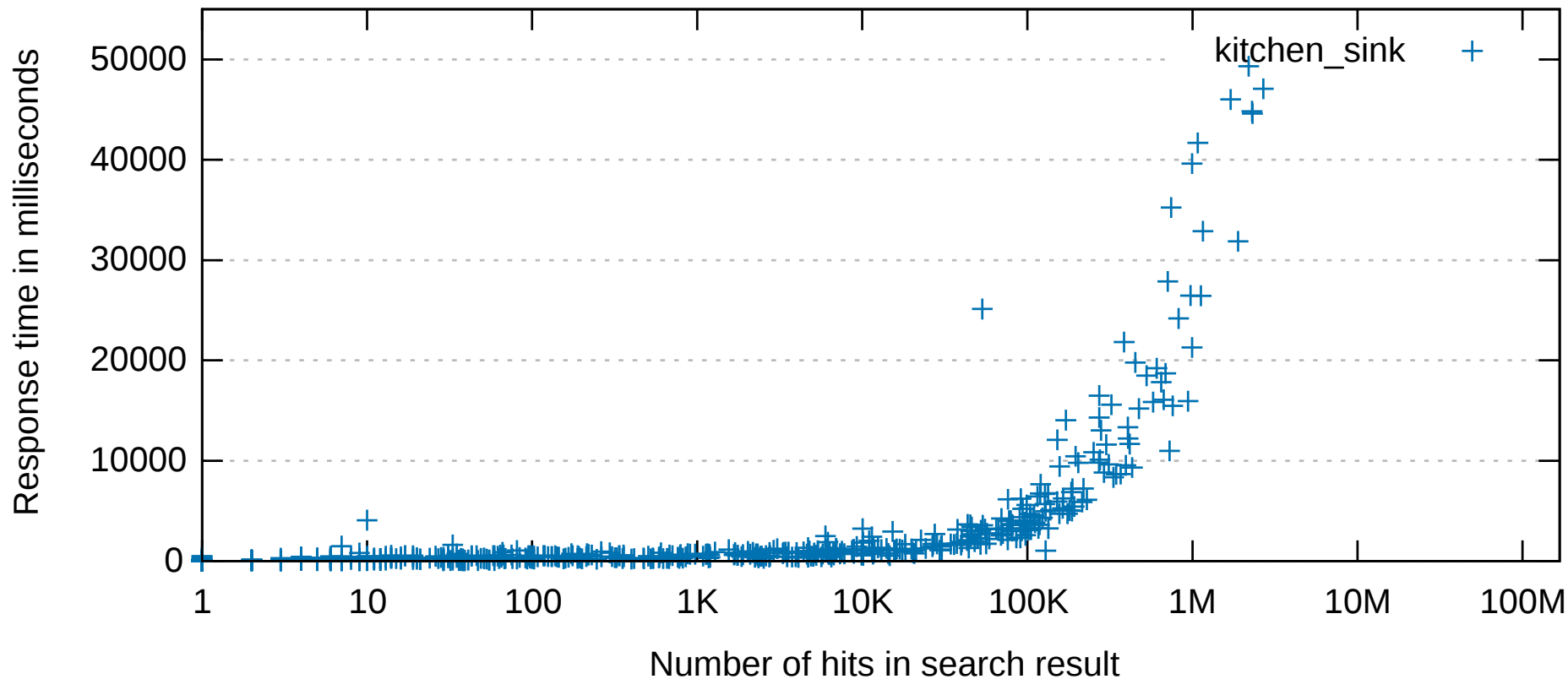
250,000,000 docs / 900GB, optimized			
Field	References	Max docs/term	Terms
domain	250,000,000	3,000,000	1,100,000
url	250,000,000	56,000	200,000,000
links	5,800,000,000	5,000,000	610,000,000

x 9 shards

x 3 concurrent requests

Shouldn't Have Done That

9 shards / 7TB / 2.3G docs, 256GB RAM, 16 cores, facet 1.1M/250M/610M values/shard



Some Great Reward

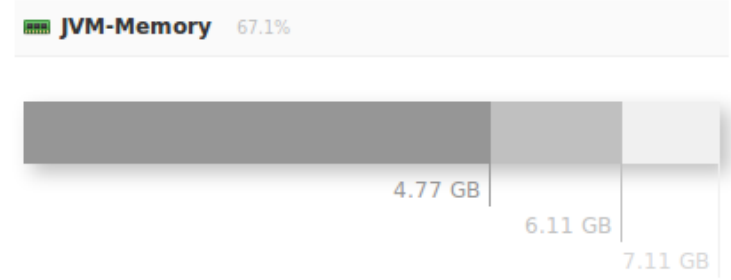
JVM

Runtime Oracle Corporation Java HotSpot(TM) 64-Bit Server VM (1.7.0_76 24.7...)

Processors 32

Args

- Djava.io.tmpdir=/home/summanet/tomcat-solr8a/temp
- Dcatalina.home=/home/summanet/tomcat-solr-master
- Dcatalina.base=/home/summanet/tomcat-solr8a
- Djava.endorsed.dirs=/home/summanet/tomcat-solr-master/endorsed
- Dlog4j.configuration=file:///home/summanet/services/conf/tomcat-so...
- Dsb.tomcat.maxthreads=200
- DzkHost=localhost:52001,localhost:52002,localhost:52003
- Dsb.solr.httpport=52308
- Dsb.solr.shutdownport=52508
- Dsb.solr.instanceid=8a
- Xmx8192m
- Djava.security.egd=file:/dev/./urandom
- Dsun.net.inetaddr.ttl=300
- Djava.util.logging.manager=org.apache.juli.ClassLoaderLogManager
- Djava.util.logging.config.file=/home/summanet/tomcat-solr8a/conf/lo...



8GB heap per
900GB shard

Dream On

- Threaded counting
- Monotonically increasing tracker for nplane-z
- Regexp filtering
- Fine count skipping
- Counter capping

Nothing's Impossible