

"Flash is good"  
-- Jim Gray, 2006

1M IOPS systems are within reach  
Page caching throttles IOPS  
But why cache?...Locality

Then, Forget locality  
Implementing locality is expensive and  
workloads exhibit low and low locality  
• graphs (linked data and web snapshots)  
• ozone (big data analysis)

### Forget Locality (Radio Edit)

Randal Burns  
IBM Research  
Almaden

Why does Randal want 1M IOPS?  
...to solve big graph problems

"There is no big graph problem"  
-- Martin Kersten, 2010

A less reactionary approach?  
Scrub for some cache performance to improve  
IOPS into the page cache.  
For graph analysis workloads that exhibit low  
locality and benefit from in-cache.

### I/O Challenges for Massive Graphs

Can we build data-intensive analysis for graphs?  
human brain: 10<sup>11</sup> vertices, 10<sup>15</sup> edges

n-body: 10<sup>11</sup> particles, 10<sup>15</sup> snapshots

Current computing paradigms:  
• supercomputing  
• custom hardware (Cray XMT)  
Require in-memory or semi-external  
memory representations of graphs  
• batch computing on cloud  
Not fast!

### What's Hard about Graphs?

- small and shrinking diameter
- no good cuts
- no locality
- bad I/O lower bounds

### I/O Tricks Don't Work

- Partition and parallelize (no good partitions)  
Localize and cache (no locality)  
Stream data (no natural ordering)
- I/O friendly (sort-like) algorithms
- hot tracking
  - connected components
  - minimum spanning forest
  - maximal matching

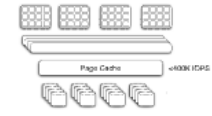
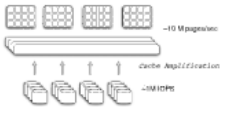
### I/O unfriendly (traversal) algorithms

- traversable list search
  - shortest paths
  - maximal clique
  - diameter
- and anything else you study want to know

Two things to debate  
• are workloads evolving to have less locality  
• are there big graph problems

\$400

\$7600  
1m Random Read IOPS

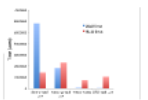
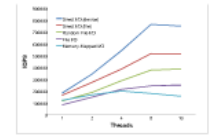


What went wrong?  
There's a lot of bookkeeping

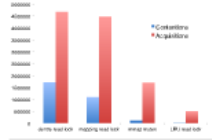
Address space: List, size, dirty, locked, and so on  
Page cache: Page lock table

Buffer cache: List, size, dirty, locked, and so on  
Buffer hash table

...and locks for all



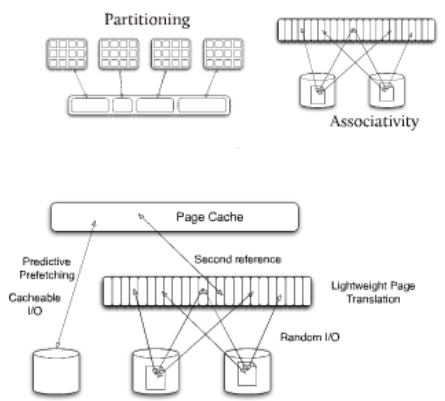
What we think?  
it's locking overhead  
What we can't show?  
it's locking overhead



Locks contention is moderate  
wait times are small  
But, locking needs expensive cache coherency  
when requested from multiple cores  
so we are performing "lock out" experiments  
to verify lock overheads

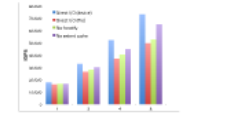
Both techniques used in  
processor caches.  
Something more  
interesting to do?

Directions  
• reduce lock intersection (decluster locks)  
• speculative lock elision (eliminate locks)  
• move away from global data structures



Investment Strategies  
Promote page on second hit  
Old data (e.g. HPI, LRU, ARC)  
but about investing cycles, not replacement

In this manner, we can  
accelerate direct file I/O  
Random interference (just overhead)



Thought for food  
• Are there large graphs? (that you care about?)  
• Is workload locality decreasing?  
• Are random read IOPS a performance  
limiting factor for analysis engines?

# Forget Locality

(Radio Edit)

Randal Burns

Department of Computer Science  
Institute for Data Intensive Engineering  
and Science (IDIES)  
Johns Hopkins University

"Flash is good"

-- Jim Gray, 2006

M IOPS systems are within reach

"Flash is good"

-- Jim Gray, 2006

1M IOPS systems are within reach

Page caching throttles IOPS

But why cache?....Locality

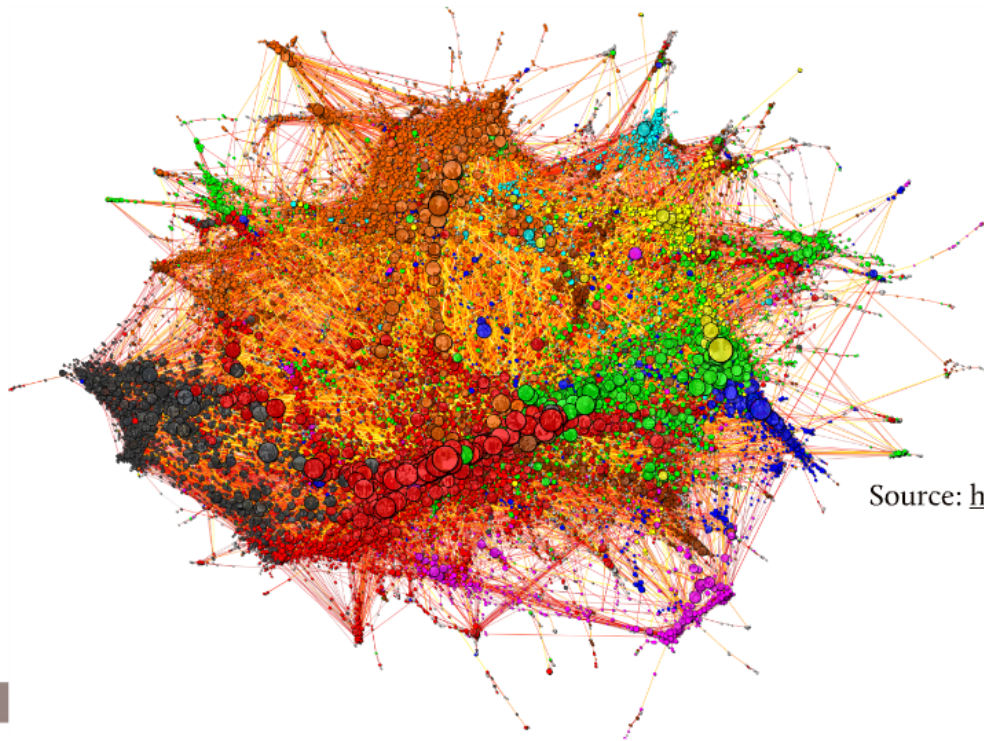
# But why cache?....Locality

## Then, F(orget) locality

Implementing locality is expensive and workloads exhibit less and less locality

- graphs (linked data and relationships)
- scans (big data analysis)

# Why does Randal want 1M IOPS?



....to solve big  
graph problems

Source: <https://http://sixdegrees.hu/last.fm/>



'There is no big graph problem'

-- Martin Kersten, 2010

## A less reactionary approach?

Sacrifice some cache performance to improve IOPS into the page cache.

For graph analysis workloads that exhibit little locality and benefit less from caching.

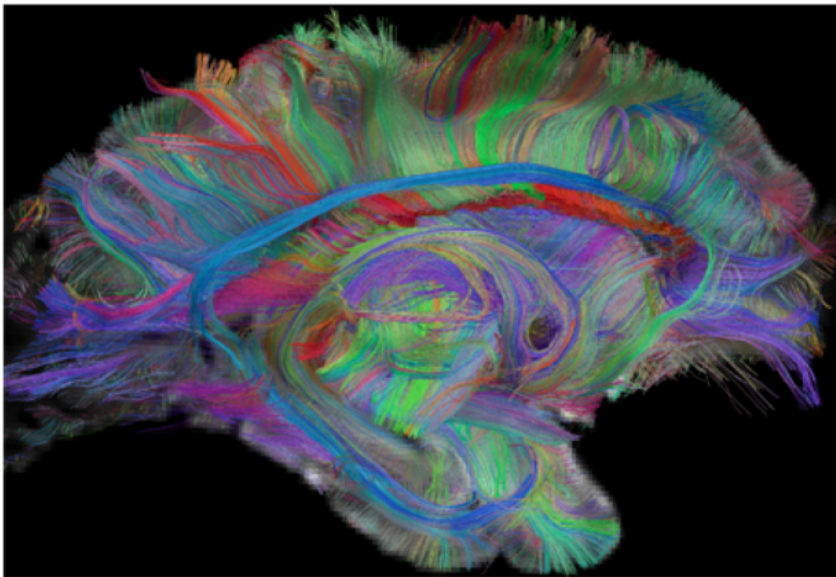


Redesign the page cache for  
concurrency rather than locality.

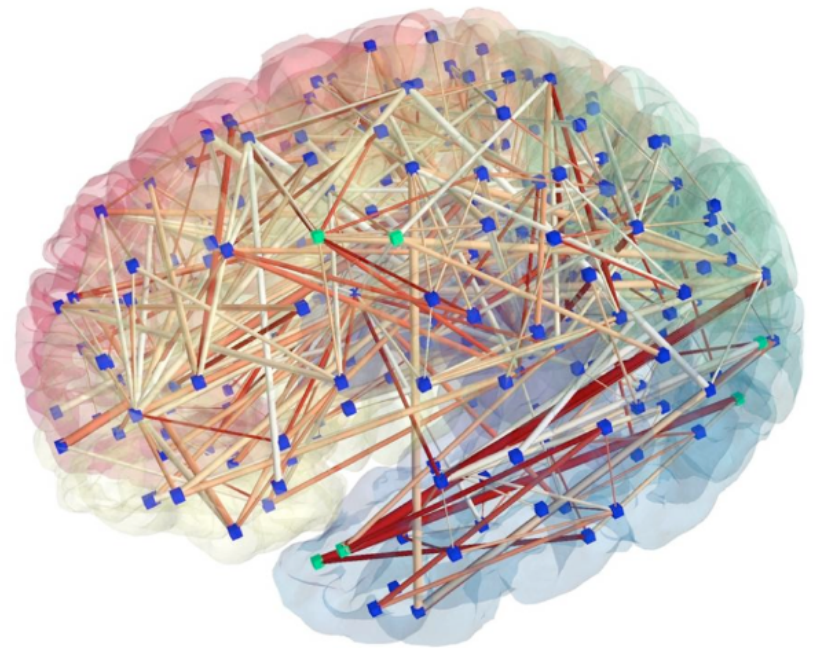
# I/O Challenges for Massive Graphs

Can we build data-intensive analysis for graphs?

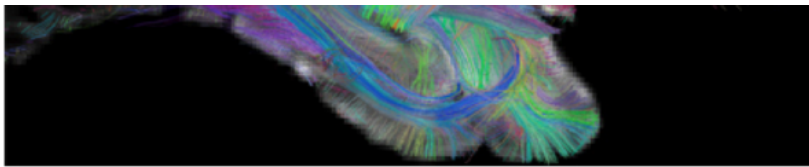
human brain:  $10^{11}$  vertices,  $10^{15}$  edges



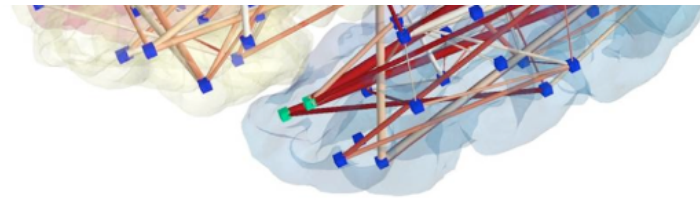
human dmri



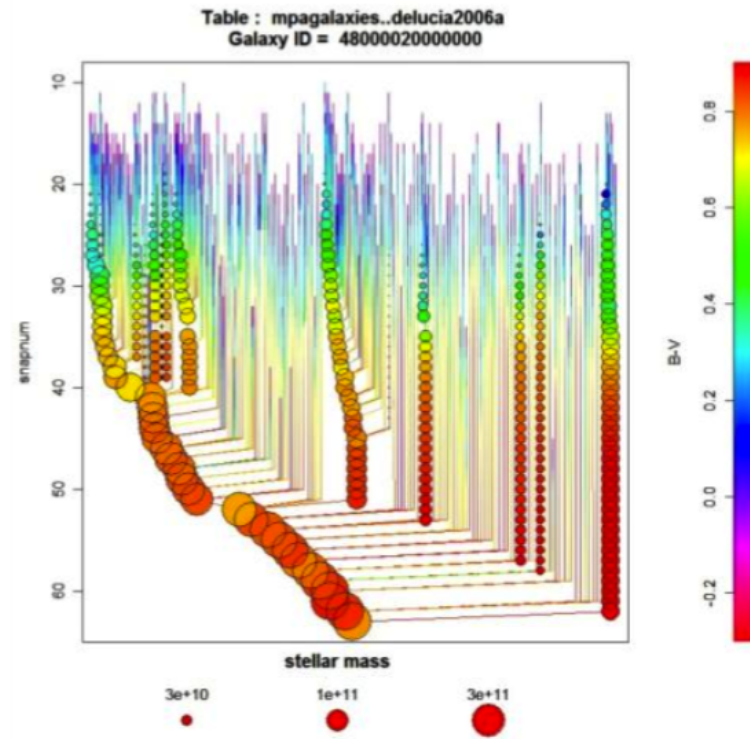
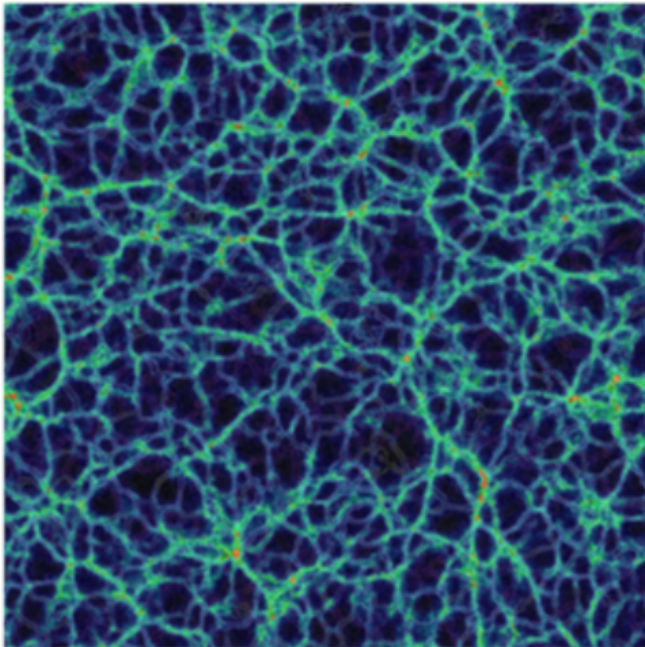
n-body: 1B particles, 10k snapshots



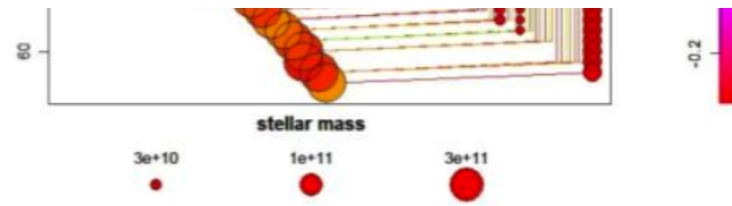
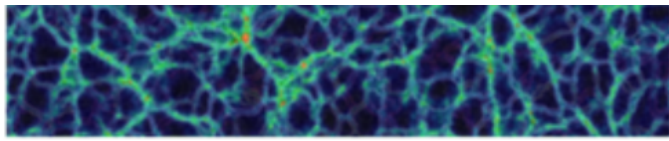
human dmri



n-body: 1B particles, 10k snapshots



Current computing paradigms:



## Current computing paradigms:

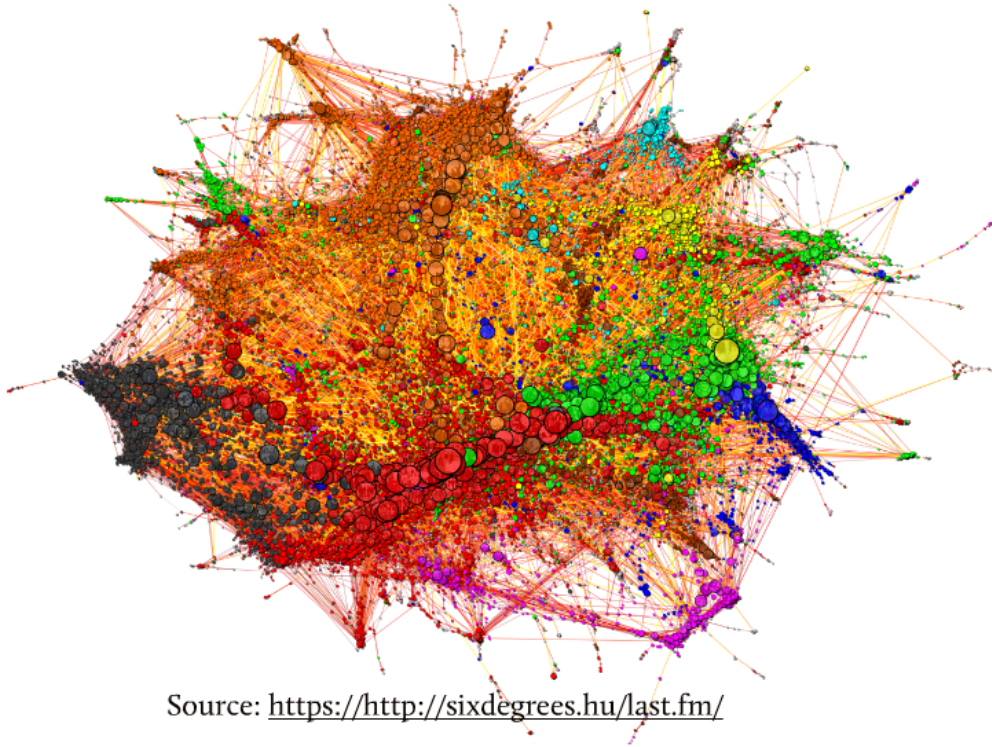
- supercomputing
- custom hardware (Cray XMT)

Require in-memory or semi-external memory representations of graphs

- batch computing on cloud

Not fast!

# What's Hard about Graphs??



Source: <https://http://sixdegrees.hu/last.fm/>

- small and shrinking diameter
- no good cuts
- no locality
- bad I/O lower bounds

# I/O Tricks Don't Work

Partition and parallelize (no good partitions)

Localize and cache (no locality)

Stream data (no natural orderings)

I/O friendly (sort-like) algorithms

Localize and cache (no locality)

Stream data (no natural orderings)

## I/O friendly (sort-like) algorithms

- list ranking
- connected components
- minimum spanning forest
- maximal matching

## I/O unfriendly (traversal) algorithms

- breadth-first search

- connected components
- minimum spanning forest
- maximal matching

## I/O unfriendly (traversal) algorithms

- breadth-first search
- shortest paths
- transitive closure
- diameter

and anything else you really want to know



\$400

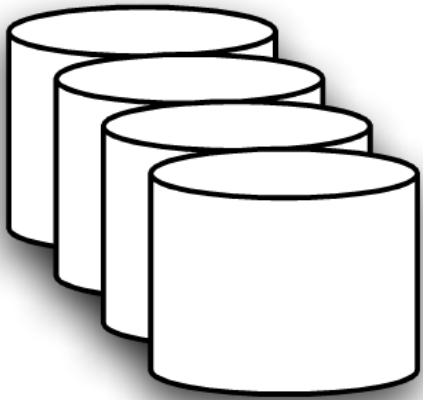


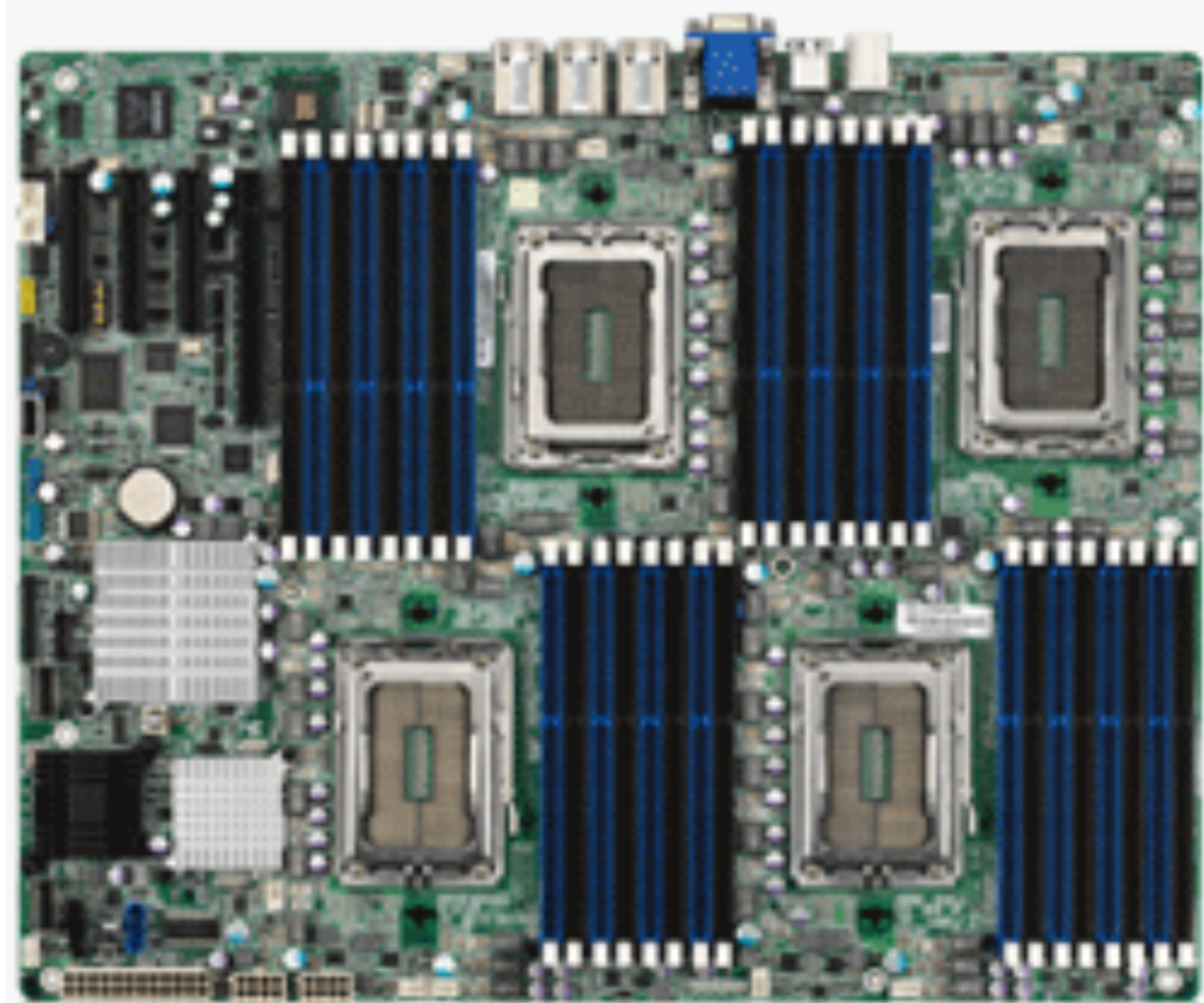
\$7600

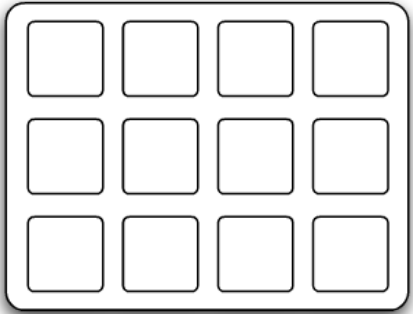
Im Random  
Read IOPS



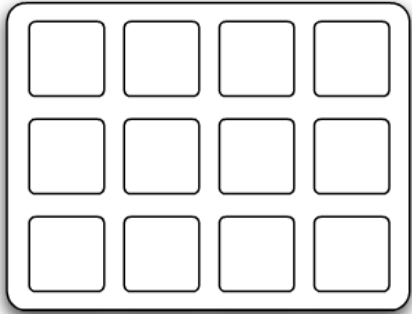
**>800K IOPS**  
**(random read)**



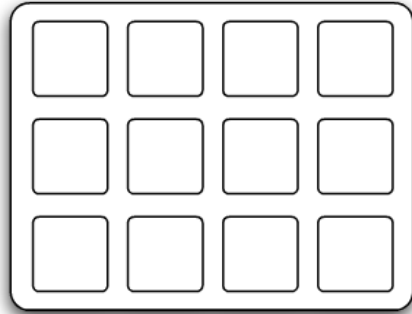




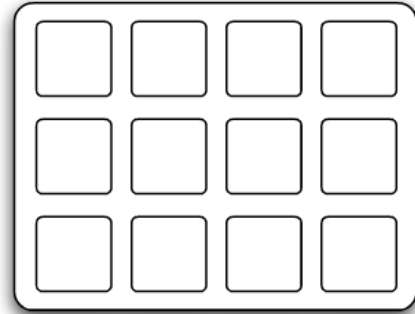
**6M pages/sec**



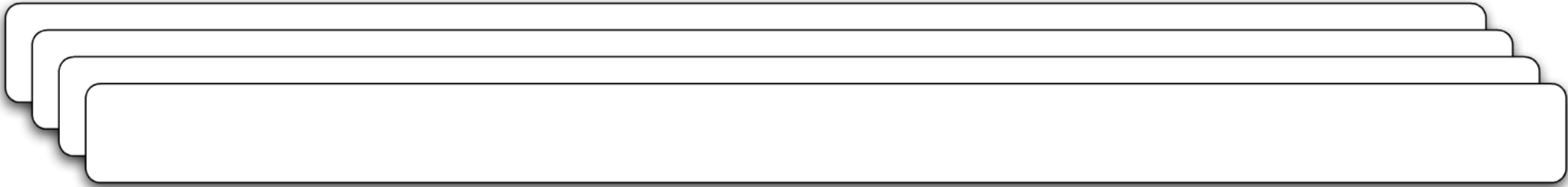
**6M pages/sec**

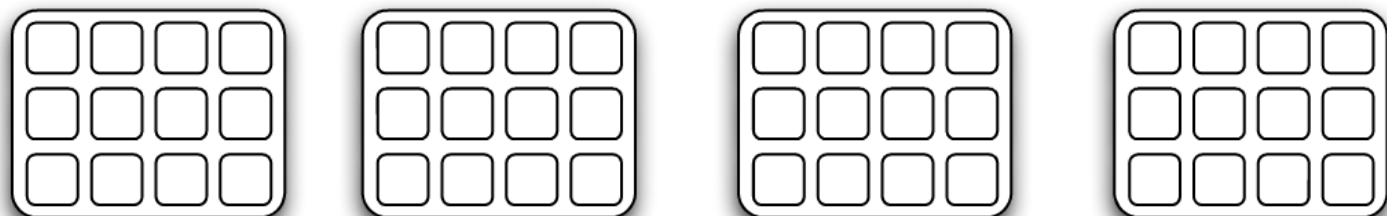


**6M pages/sec**

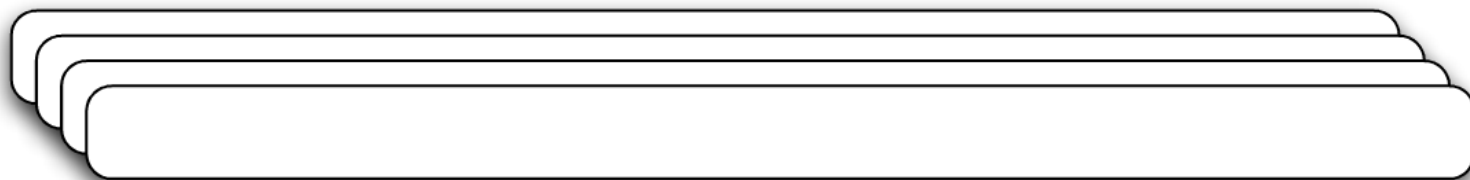


**6M pages/sec**

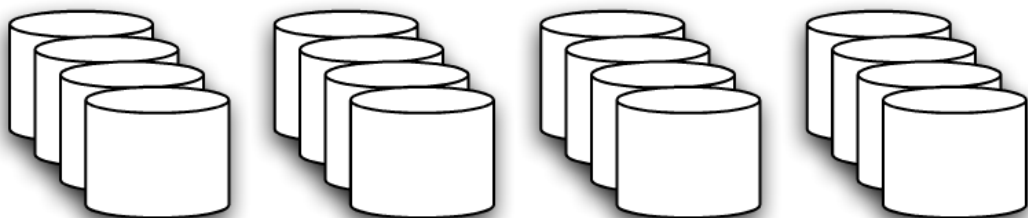




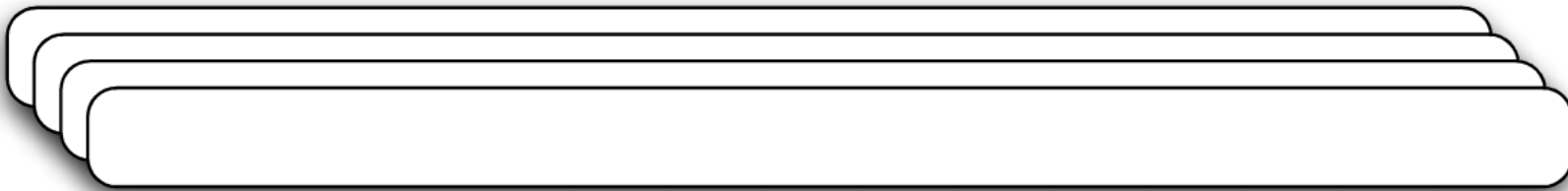
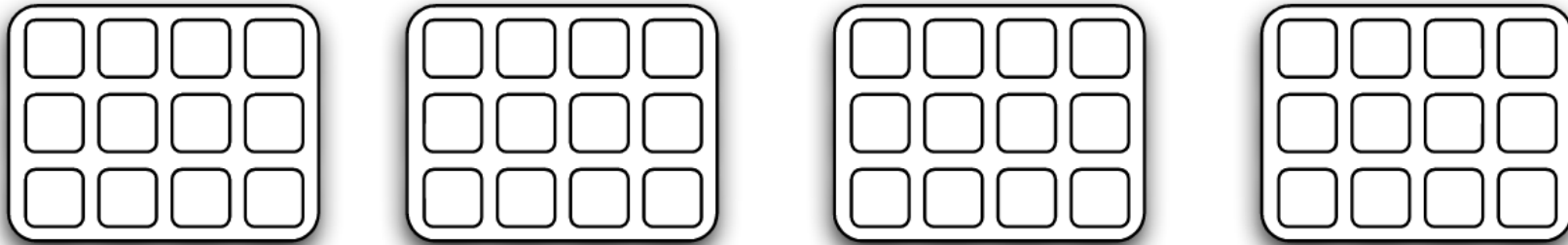
~10 M pages/sec



*Cache Amplification*

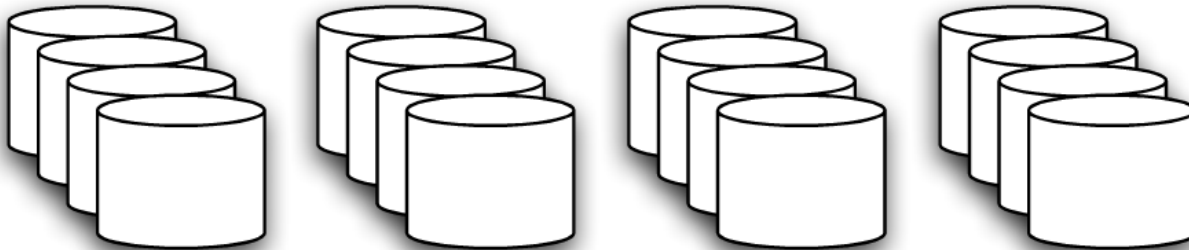


~1M IOPS



Page Cache

<400K IOPS



# What went wrong?

There's a lot of bookkeeping

Address space:

Lists:

clean  
dirty  
locked

RB Tree:  
page maps

Page cache:

Page hash table

Buffer cache:

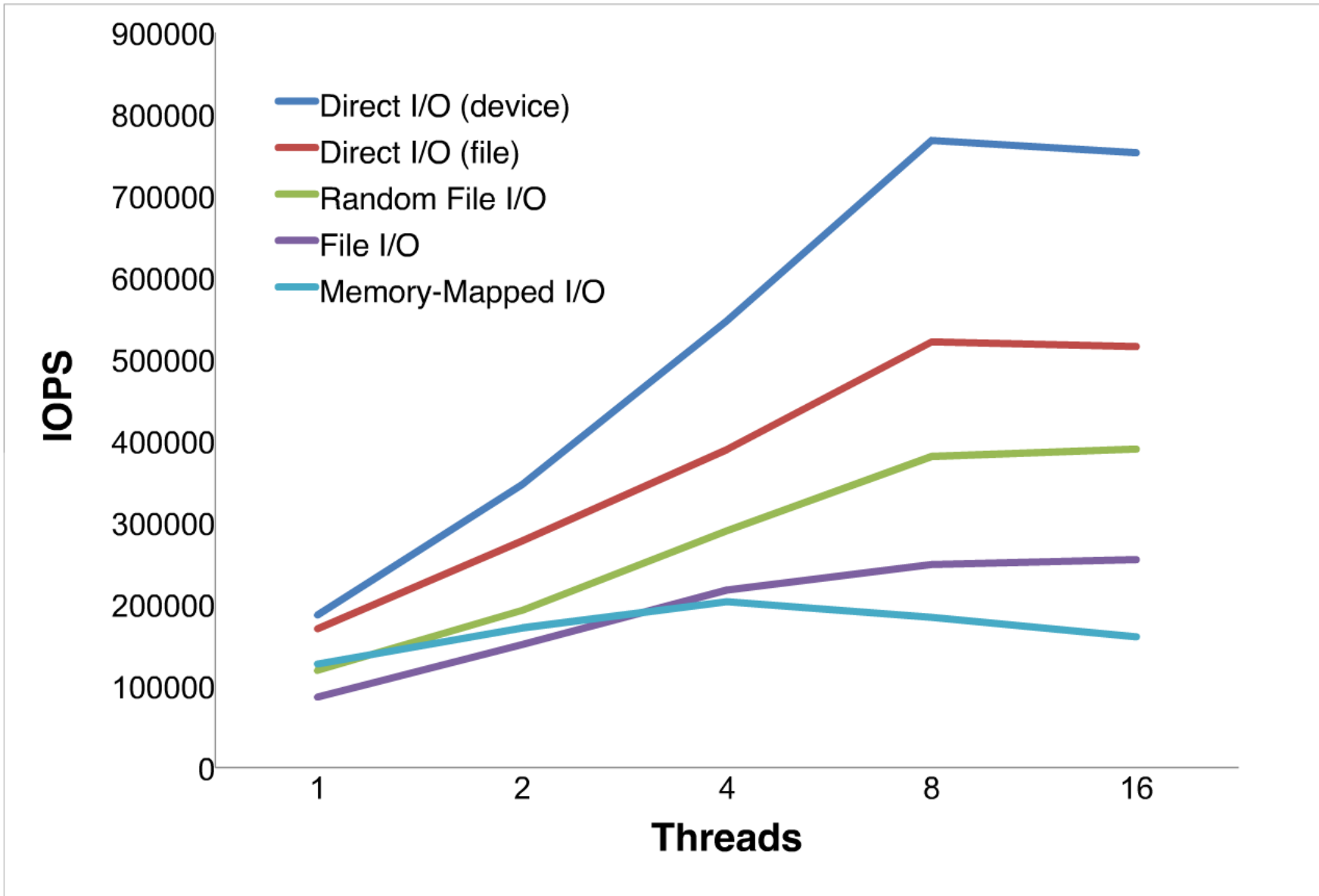
Lists:

clean  
dirty  
locked  
shared  
unshared

Buffer hash table

....and locks for all





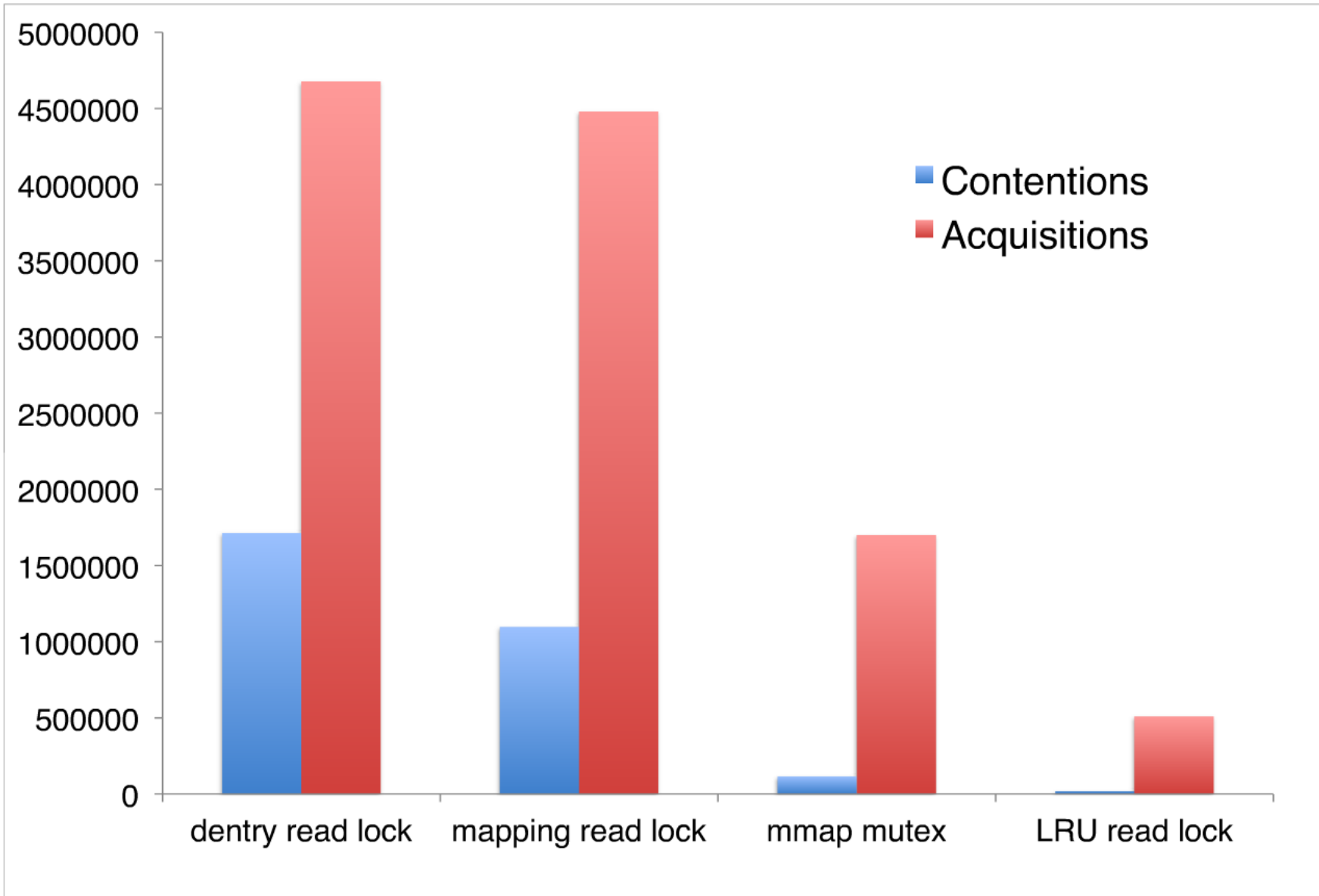


What we think?

it's locking overhead

What we can't show?

it's locking overhead



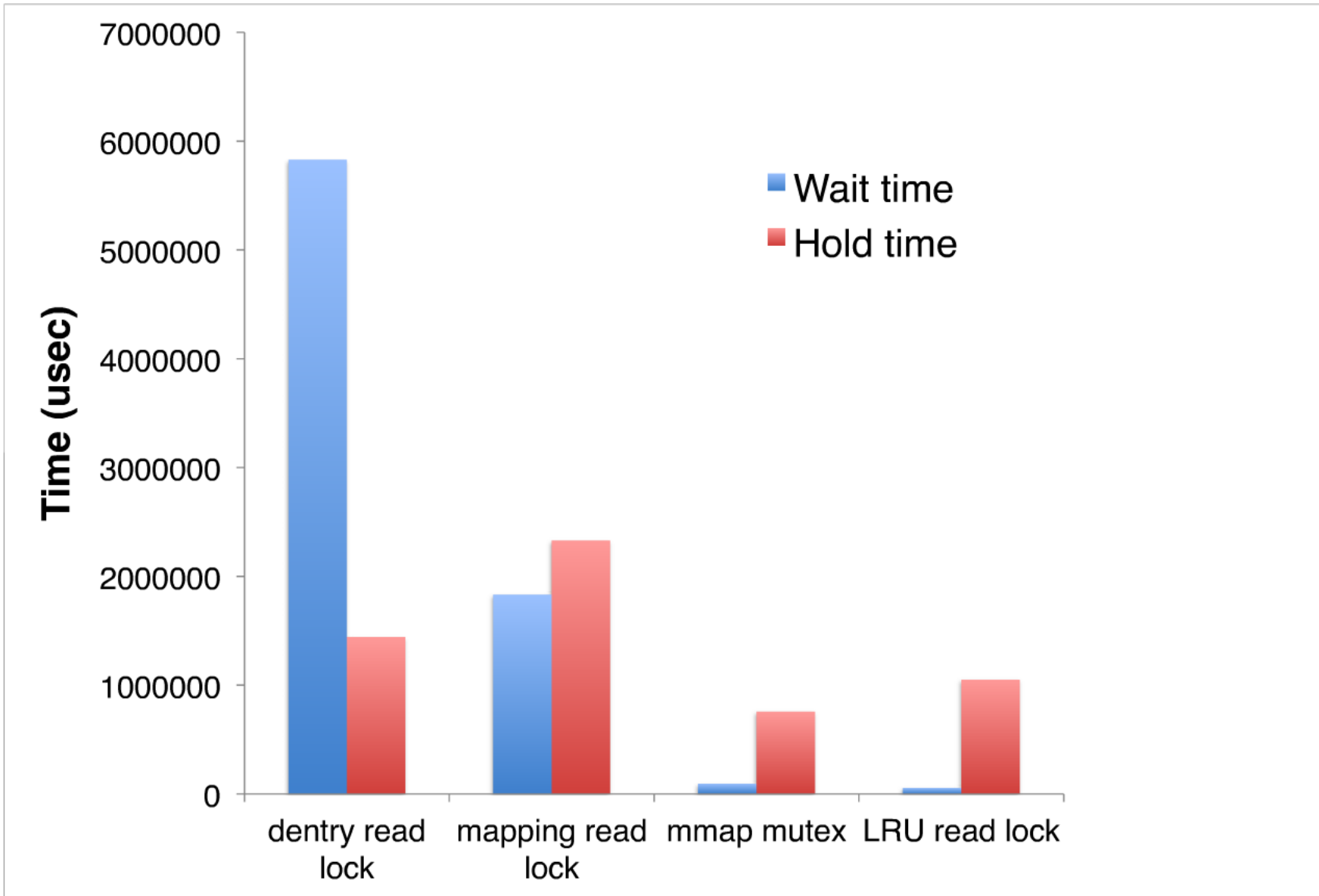
Locks contention is moderate

wait times are small

But, locking needs expensive cache coherency

when requested from multiple cores

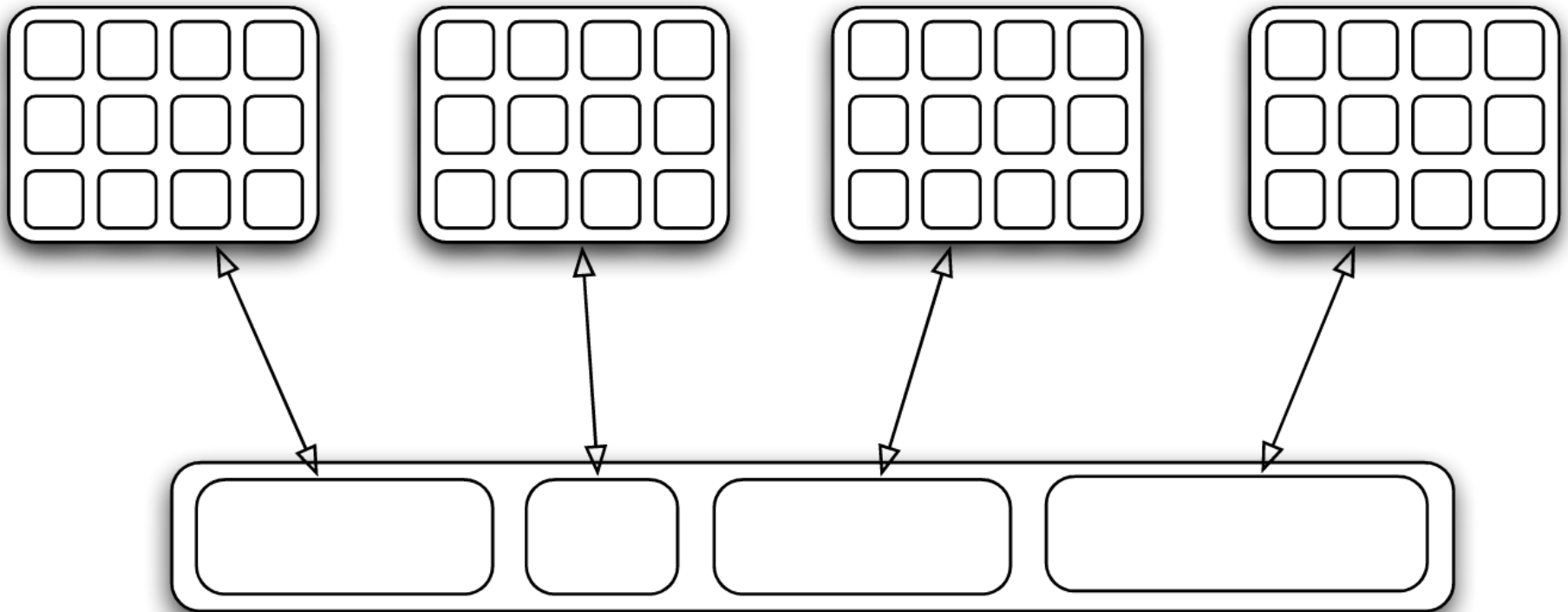
we are performing "knock out" experiments  
to verify lock overheads

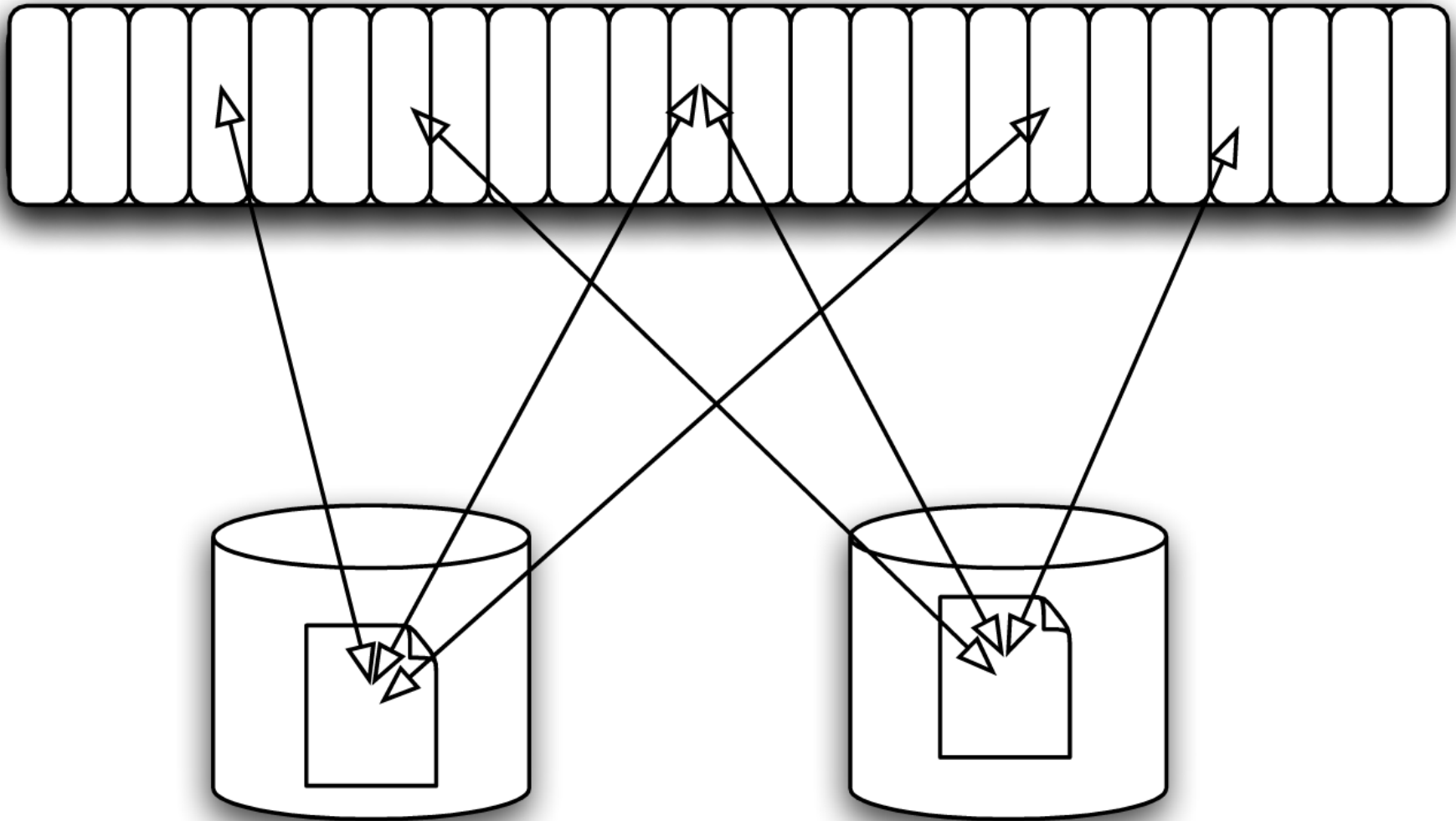


## Directions

- reduce lock intersection (decluster locks)
- speculative lock elision (eliminate locks)
- move away from global data structures

# Partitioning





*Associativity*



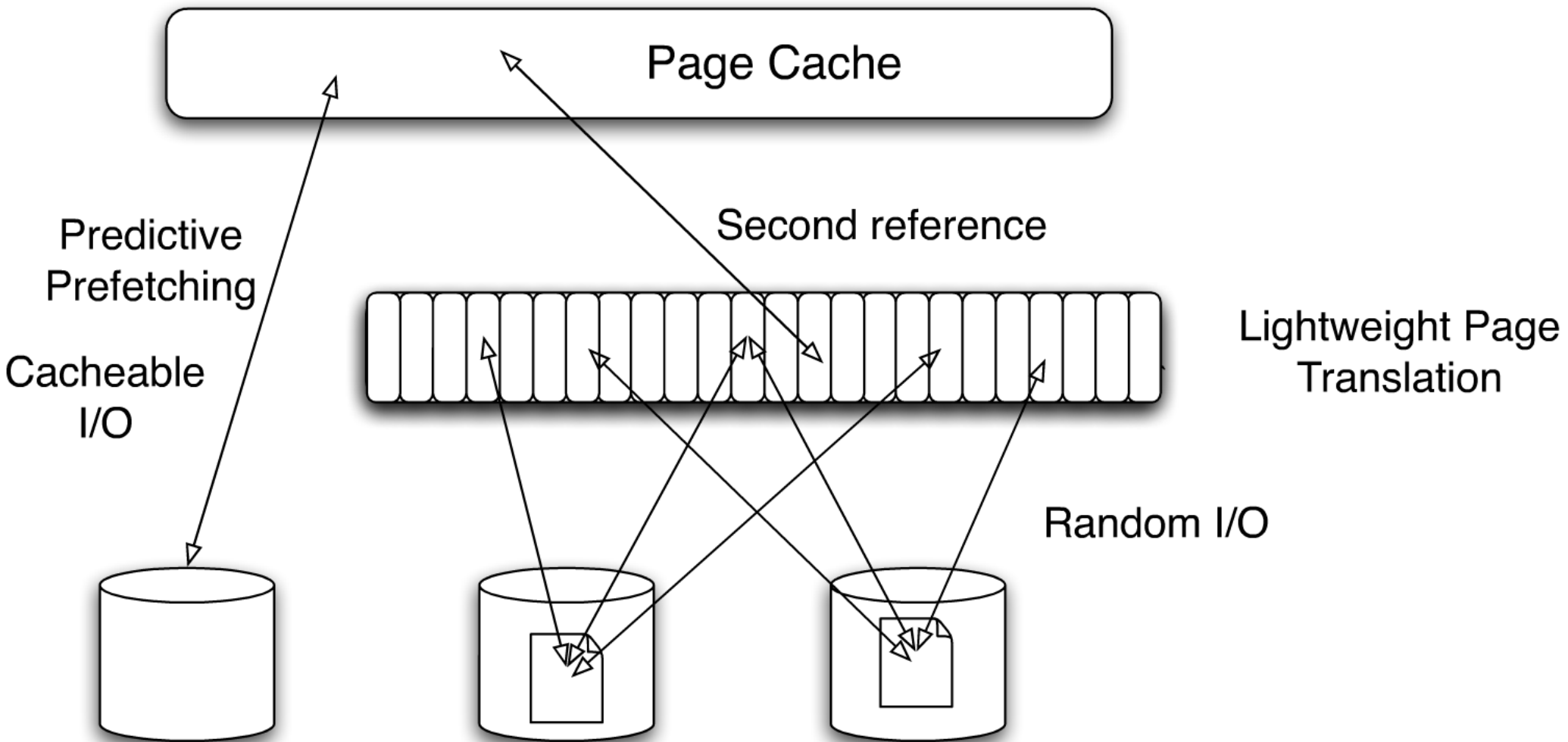
Both techniques used in  
processor caches.

Something more  
interesting to do?

# Investment Strategies

Promote pages on second hit  
Old idea (2Q-LRU, LIRS, ARC)  
but about investing cycles, not replacement





Page Cache

Predictive Prefetching

Second reference

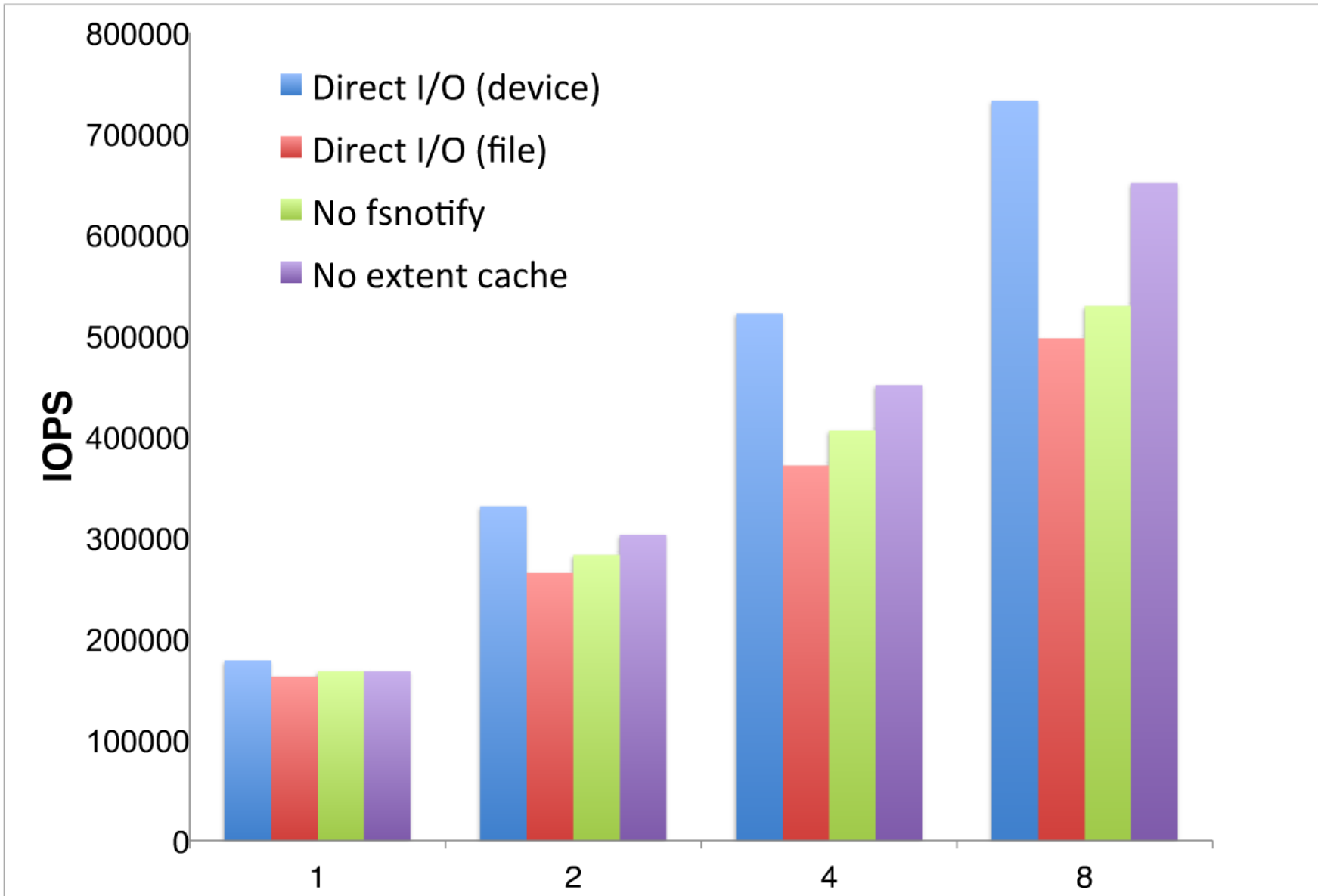
Lightweight Page Translation

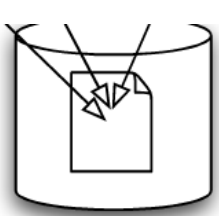
Cacheable I/O

Random I/O

In this manner, we can  
accelerate direct file I/O

Reduces interference (not overhead)





## Thought for food

- Are there large graphs? (that you care about?)
- Is workload locality decreasing?
- Are random read IOPS a performance limiting factor for analysis engines?

