# Rethink the Sync

Ed Nightingale

Kaushik Veeraraghavan

Peter Chen

Jason Flinn

University of Michigan

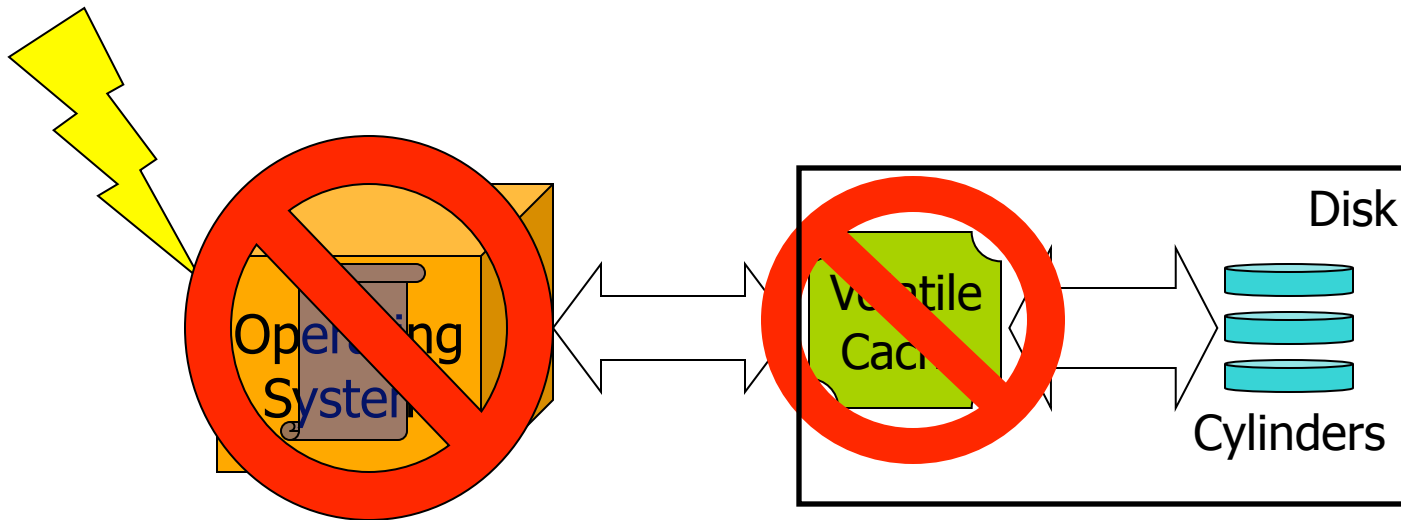# Problem

- Asynchronous I/O is a <span style="color:red">poor</span> abstraction for:
  - Reliability
  - Ordering
  - Durability
  - Ease of programming

- Synchronous I/O is superior but <span style="color:red">100x slower</span>
  - Caller blocked until operation is complete

# Solution

- Synchronous I/O can be fast

- New model for synchronous I/O
  - External synchrony
  - Same guarantees as synchronous I/O
  - Only 8% slower than asynchronous I/O

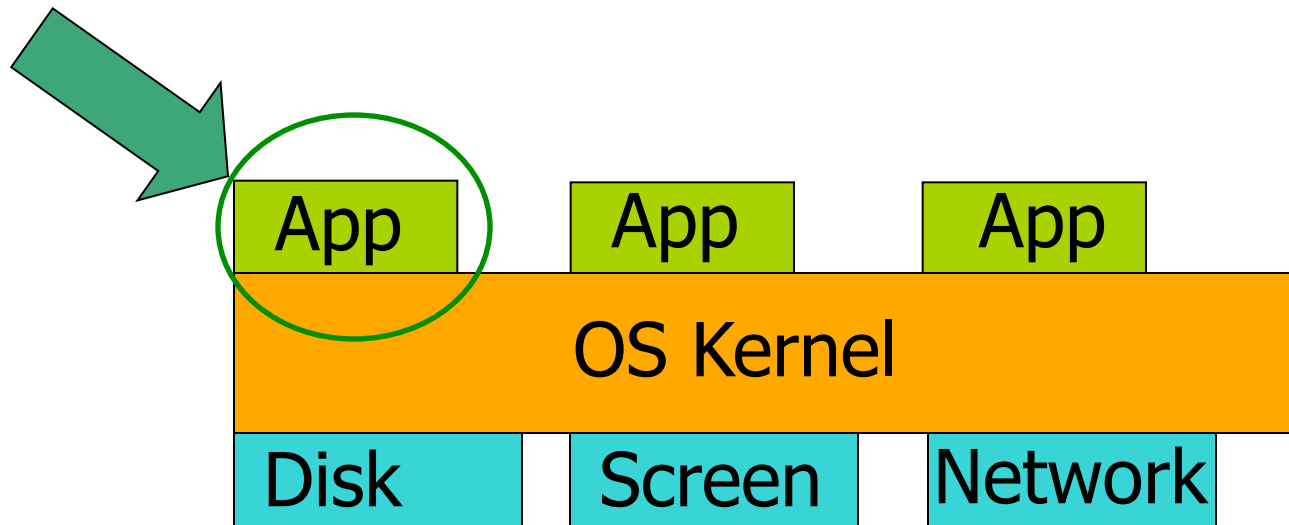# When a sync() is really async

- On sync() data written only to volatile cache
  - 10x performance penalty and data NOT safe



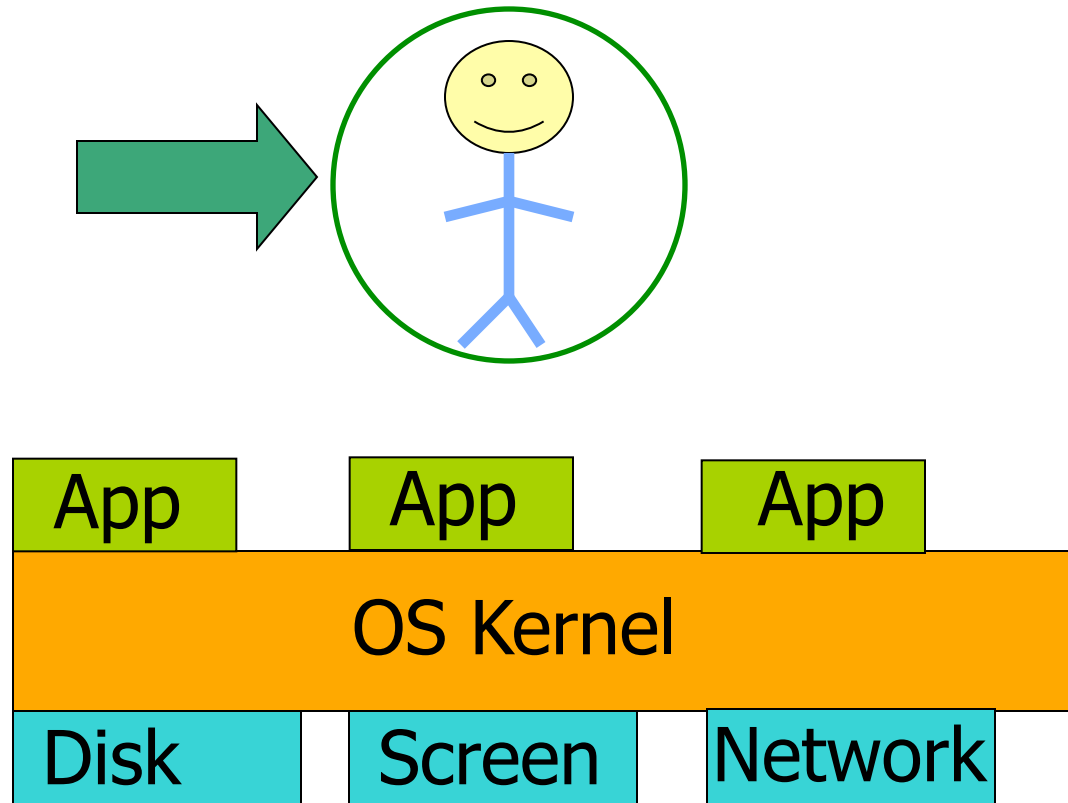- 100x slower than asynchronous I/O if disable cache

# To whom are guarantees provided?

- Synchronous I/O definition:
  - Caller blocked until operation completes



- Guarantee provided to application

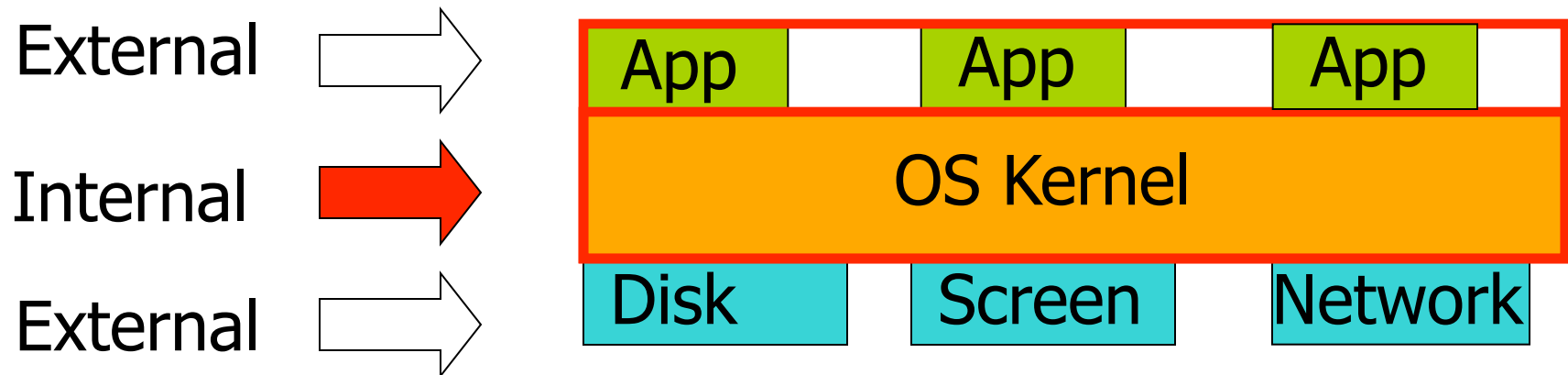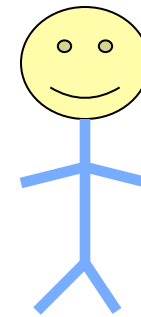# To whom are guarantees provided?

App   App   App

OS Kernel

Disk   Screen   Network

- Guarantee really provided to the user

# Providing the user a guarantee

- User *observes* operation has completed
  - User may examine screen, network, disk…

- Guarantee provided by synchronous I/O
  - Data durable when operation observed to complete

- To observe output it must be externally visible
  - Visible on external device

# Why do applications block?



External ⇒

Internal ⇒

External ⇒

| App | | App | | App | |
|---|---|---|---|---|---|
| OS Kernel | | | | | |
| Disk | | Screen | | Network | |

- Simple application need the for block more sync Application is external world need to block
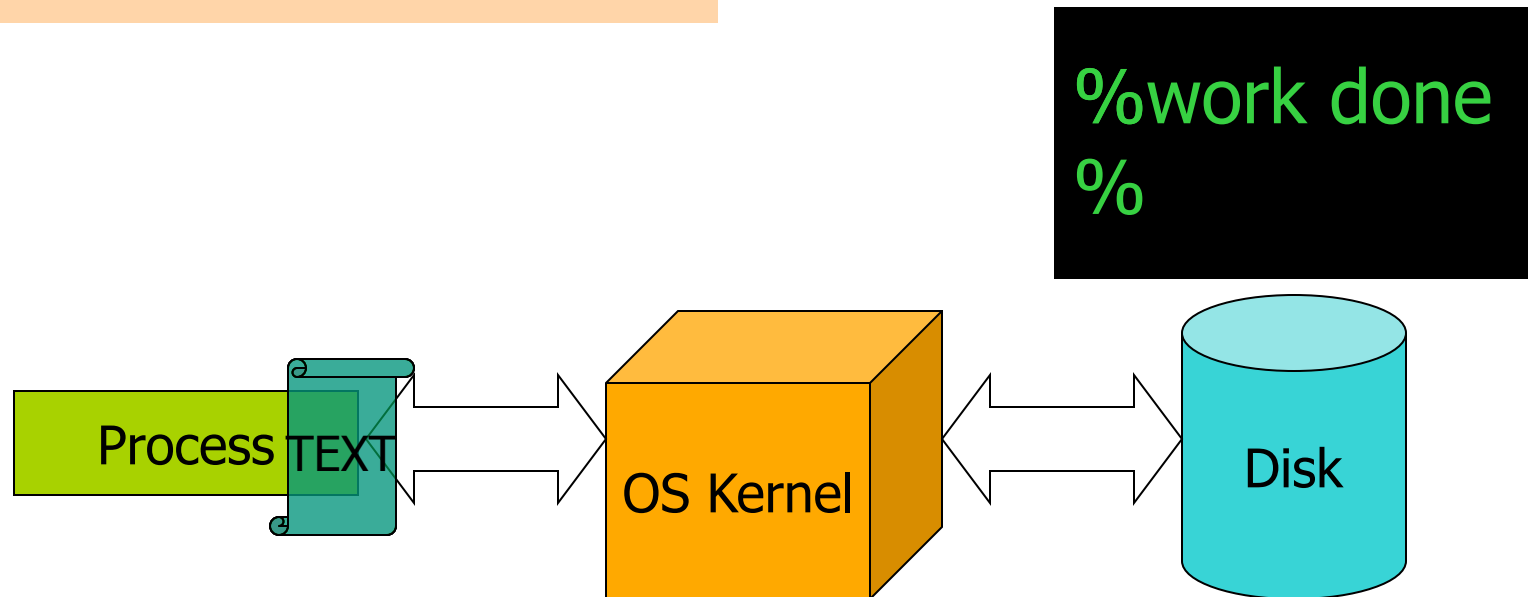Applicaiton is in external world call

# A new model of synchronous I/O

- Provide guarantee directly to user
  - Rather than via application

- Called externally synchronous I/O
  - Indistinguishable from traditional sync I/O
  - Approaches speed of asynchronous I/O

# Example: Synchronous I/O

101　write(buf_1);

102　write(buf_2);

103　print("work done");

104　foo();

← Application blocks
Application blocks

%work done
%

Process TEXT → OS Kernel ⟷ Disk

# Observing synchronous I/O

```
101   write(buf_1);
102   write(buf_2);
103   print("work done");
104   foo();
```
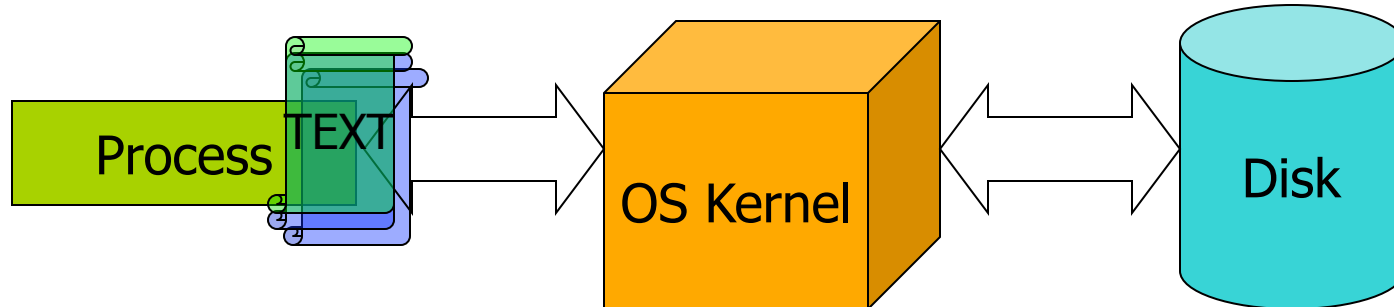
Depends on 1st write

Depends on 1st & 2nd write

- Sync I/O externalizes output based on causal ordering
  - Enforces causal ordering by blocking an application

- Ext sync: Same causal ordering without blocking applications

# Example: External synchrony

101   write(buf_1);

102   write(buf_2);

103   print("work done");

104   foo();

%work done %

Process  TEXT → OS Kernel ⟷ Disk

# Tracking causal dependencies

- Applications may communicate via IPC
  - Socket, pipe, fifo etc.

- Need to propagate dependencies through IPC

- We build upon Speculator [SOSP '05]
  - Track and propagate causal dependencies
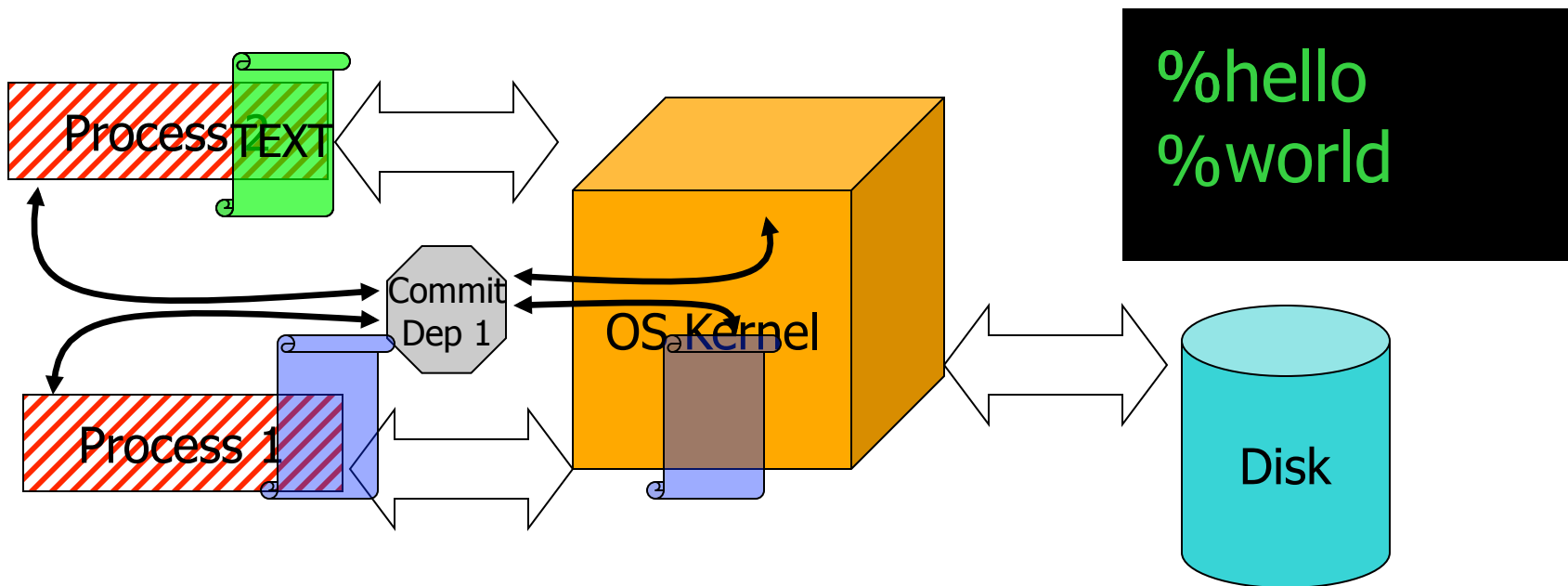  - Buffer output to screen and network

# Tracking causal dependencies

## Process 1

| | |
|---|---|
| 101 | write(file1); |
| 102 | do_something(); |

## Process 2

| | |
|---|---|
| 101 | print ("hello"); |
| 102 | read(file1); |
| 103 | print("world"); |

%hello
%world

ProcessTEXT

Process 1

Commit Dep 1

OS Kernel

Disk

# Output triggered commits

- Maximize throughput until output buffered
- When output buffered, trigger commit
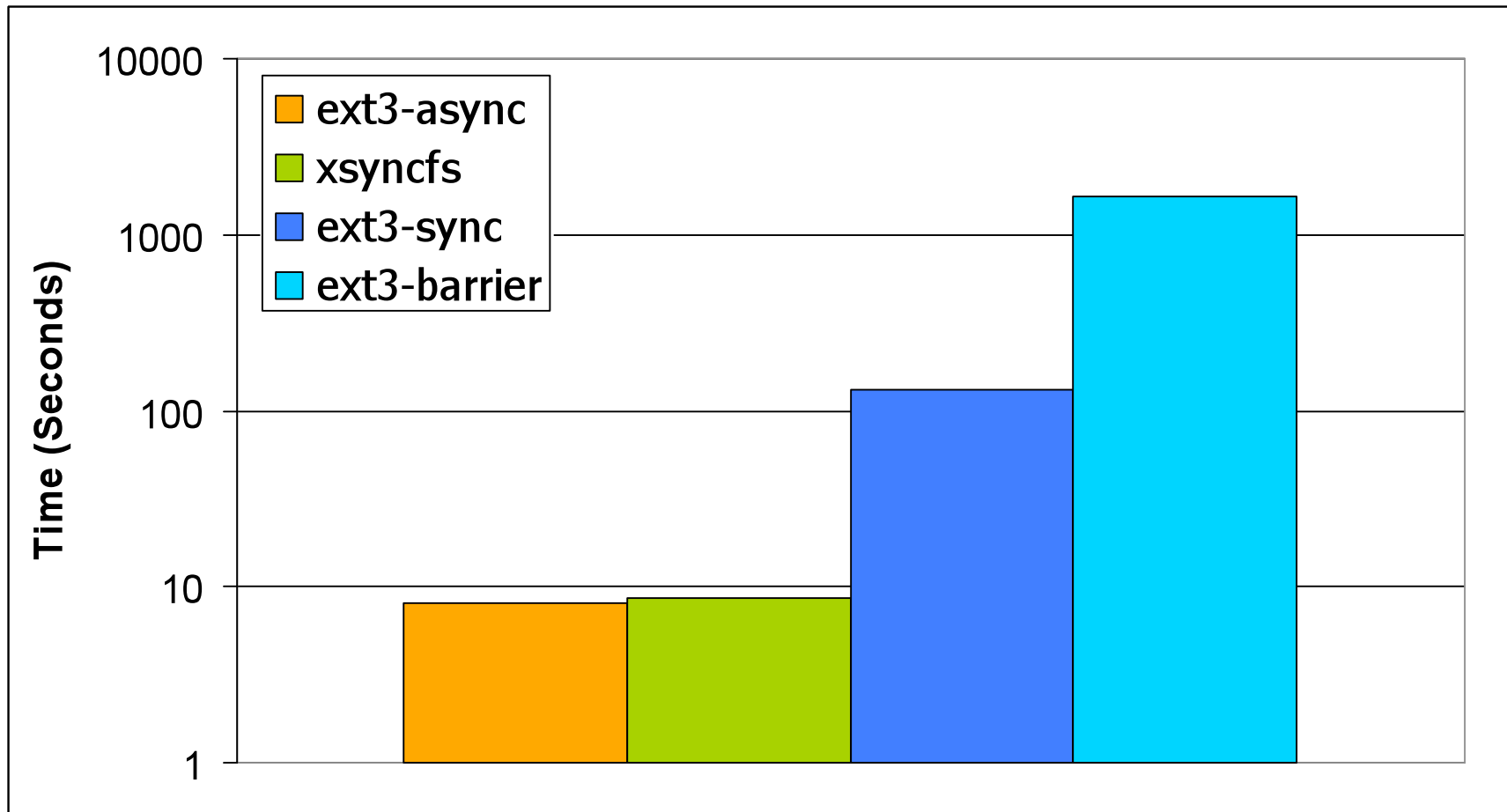  - Minimize latency only when important

%work done %

Process TEXT

OS Kernel

Disk

# Evaluation

- Implemented ext sync file system <u>Xsyncfs</u>
  - Based on the ext3 file system
  - Use journaling to preserve order of writes
  - Use write barriers to flush volatile cache

- Compare Xsyncfs to 3 other file systems
  - Default asynchronous ext3
  - Default synchronous ext3
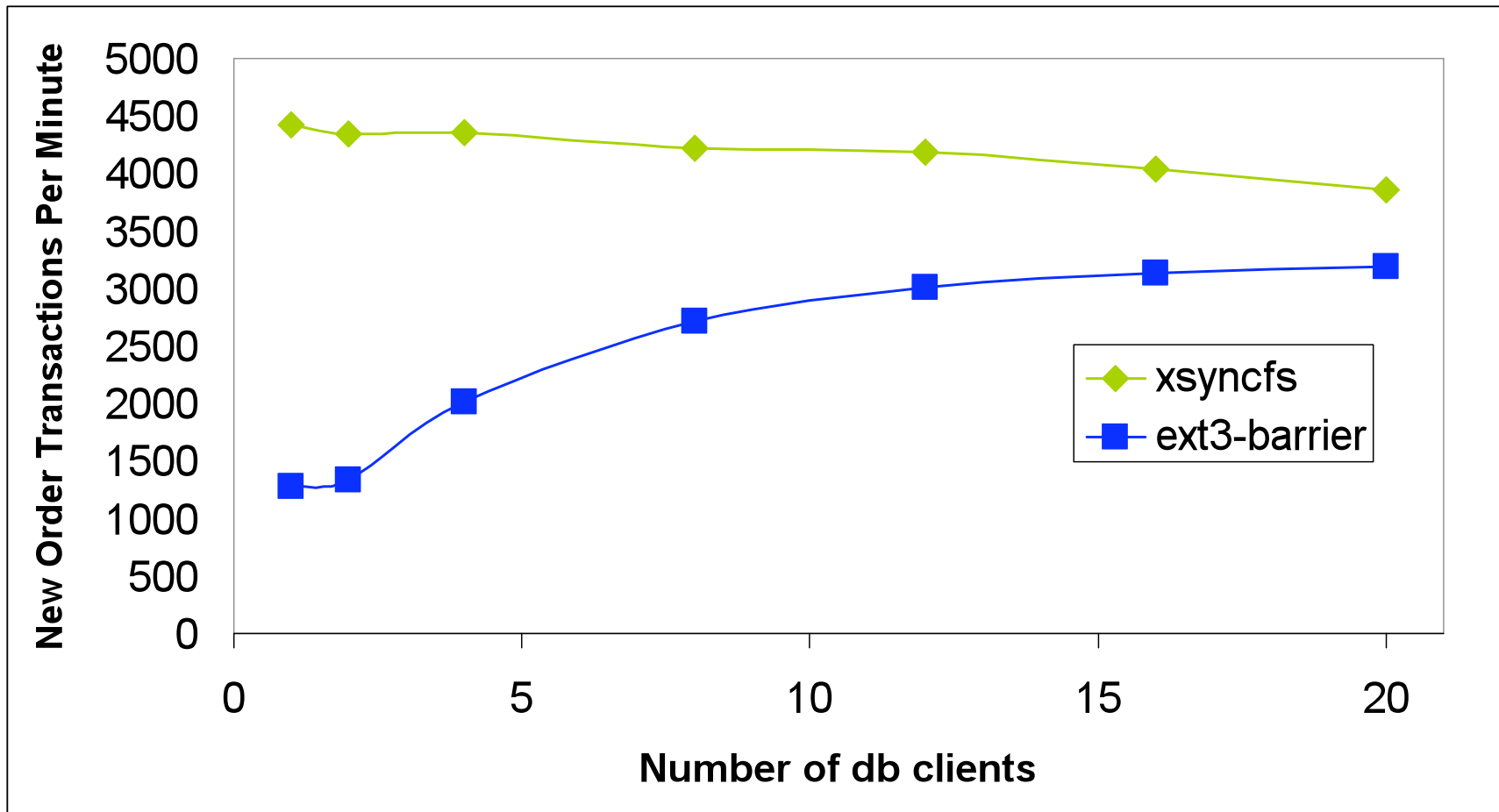  - Synchronous ext3 with write barriers

# When is data safe?

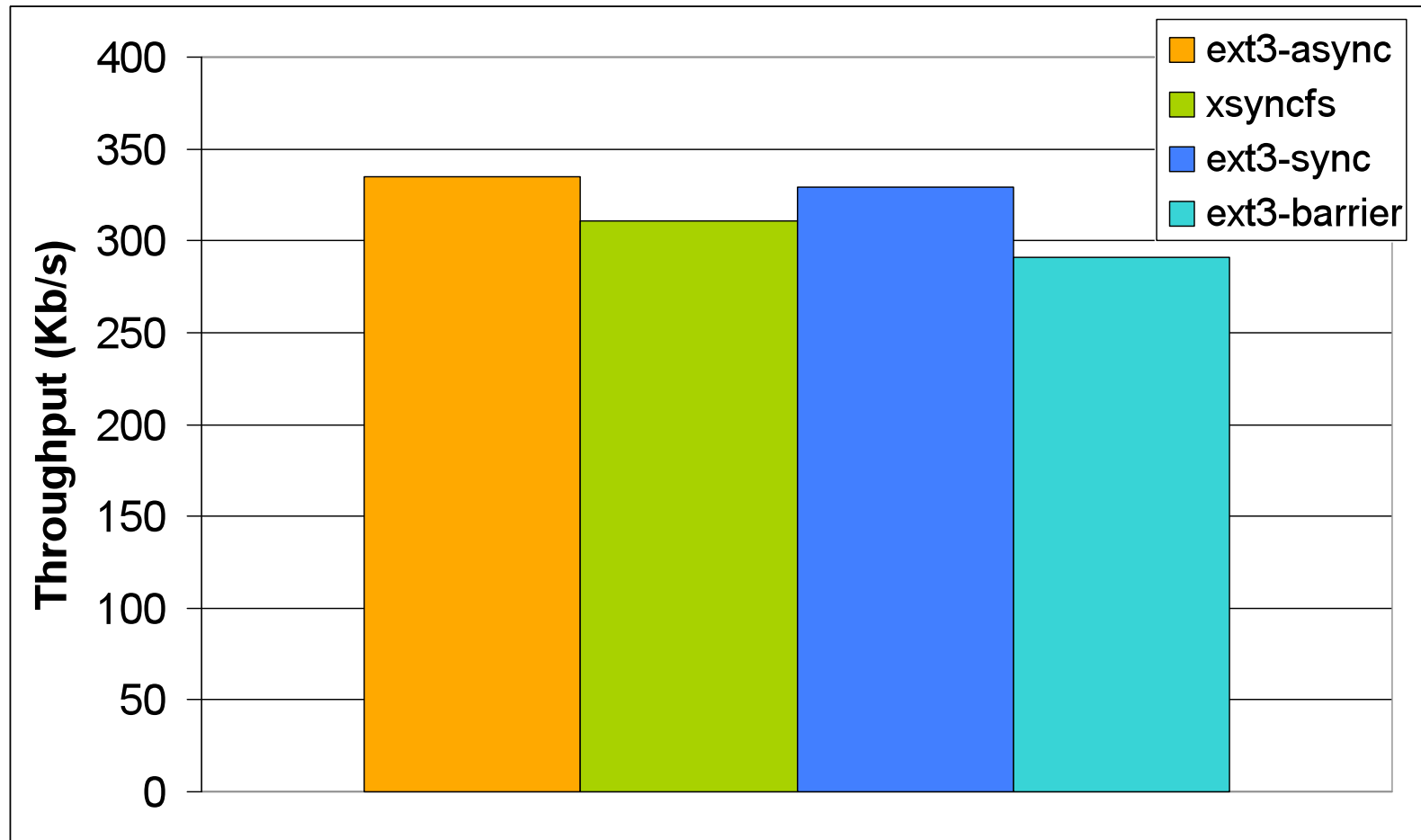| File System Configuration | Data durable on write() | Data durable on fsync() |
|---|---|---|
| Asynchronous | No | Not on power failure |
| Synchronous | Not on power failure | Not on power failure |
| Synchronous w/ write barriers | Yes | Yes |
| External synchrony | Yes | Yes |

# Postmark benchmark



- **Xsyncfs within 7% of ext3 mounted asynchronously**

# The MySQL benchmark



- **Xsyncfs can group commit from a single client**

# Specweb99 throughput



- **Xsyncfs within 8% of ext3 mounted asynchronously**

# Specweb99 latency

| Request size | ext3-async | xsyncfs |
|---|---|---|
| 0-1 KB | 0.064 seconds | 0.097 seconds |
| 1-10 KB | 0.150 second | 0.180 seconds |
| 10-100 KB | 1.084 seconds | 1.094 seconds |
| 100-1000 KB | 10.253 seconds | 10.072 seconds |

- Xsyncfs adds no more than 33 ms of delay

# Conclusion

- Synchronous I/O can be fast

- External synchrony performs with 8% of async

- Questions?