# Enabling System Transactions

## via Lightweight Kernel Extensions
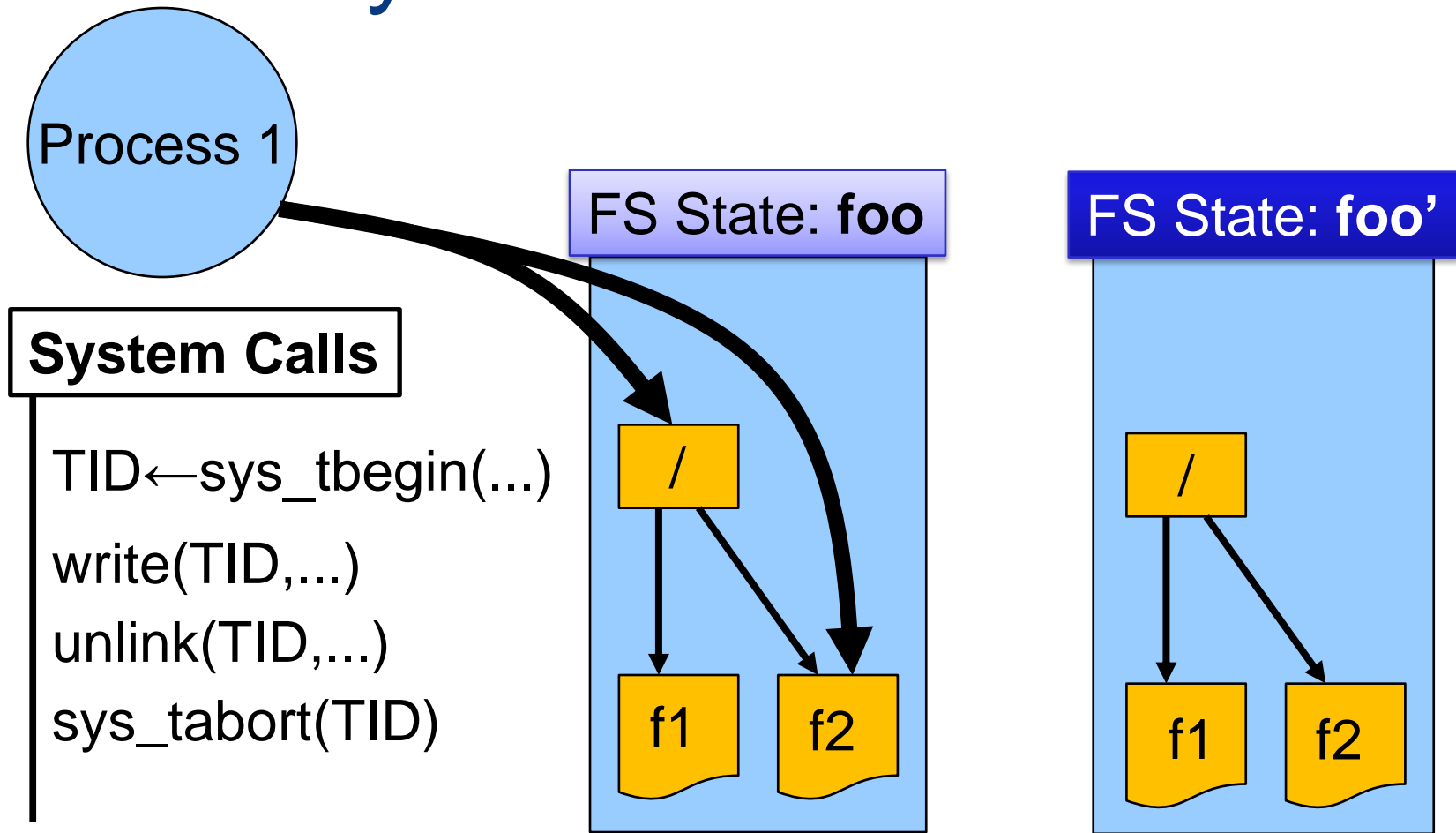
R.P. Spillane, S. Gaikwad. M. Chinni,
C.P. Wright, E. Zadok

Stony Brook University
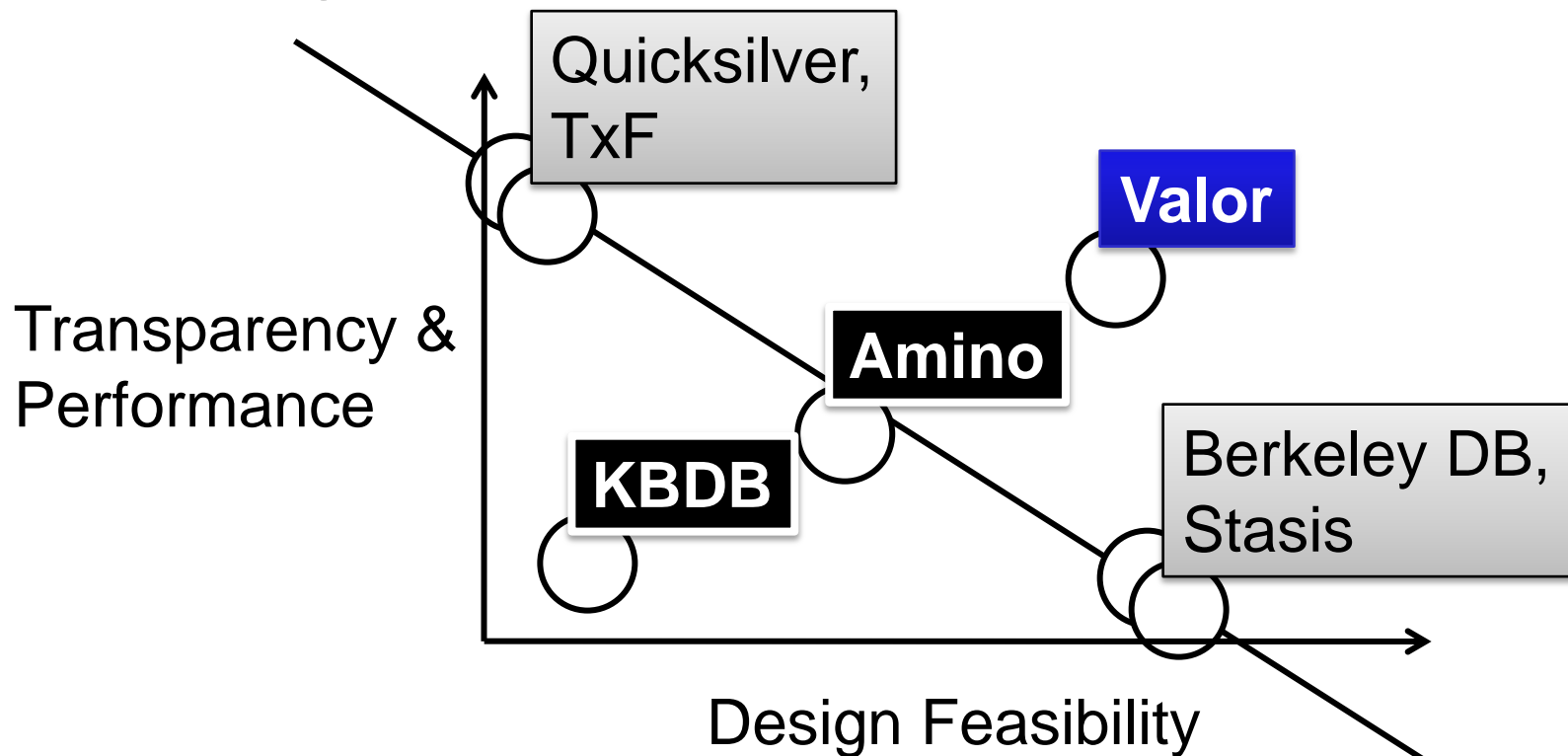
http://www.fsl.cs.sunysb.edu/

# Summary

- What is the design complexity of system transactions implemented in the VFS?
  - Low
    - 100 lines of code added to page writeback
    - 4000 lines of module code (log implementation)
- What is the performance?
  - Valor: 35% overhead on top of theoretical best, compared to…
  - 104% overhead for an efficient user-level alternative

STONY BROOK
STATE UNIVERSITY OF NEW YORK

# System Transaction

Process 1

**System Calls**

TID←sys_tbegin(...)

write(TID,...)
unlink(TID,...)
sys_tabort(TID)

FS State: **foo**

/

f1    f2

FS State: **foo'**

/

f1    f2

STONY BR🔆🔆K
STATE UNIVERSITY OF NEW YORK

# The Design Spectrum

- Valor side-steps the traditional trade-off by working with the Kernel's page cache in a general way.



Quicksilver, TxF

Valor

Amino

KBDB

Berkeley DB, Stasis

Transparency & Performance

Design Feasibility

# Valor's Process Txn Model

- Transactional Model
  - ◆ Supported Operations:
    - dirtying a page
    - appending to a file, modifying an inode
    - modifying a directory
  - ◆ Locking:
    - directory locks, inode locks
    - page range locks for overwrites
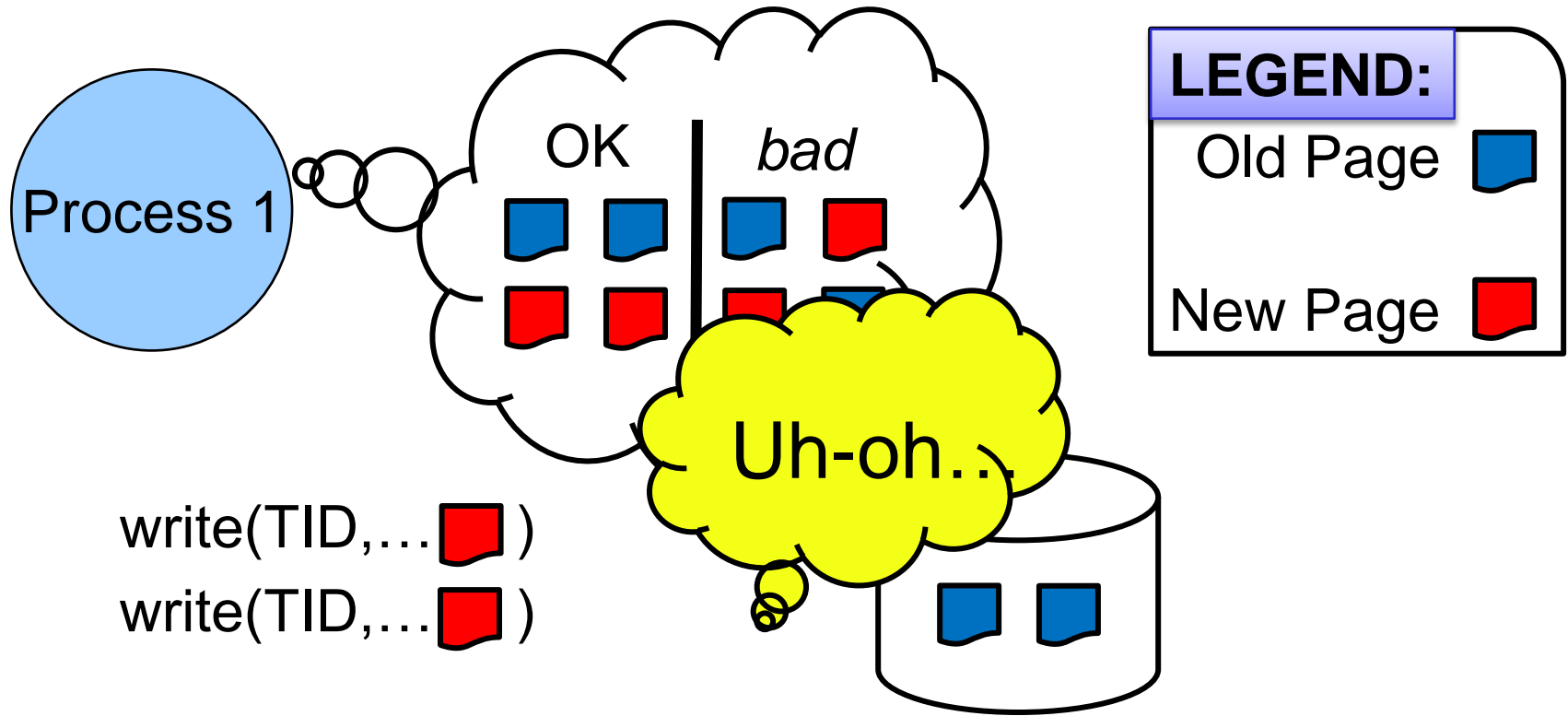    - intent locks for directory renames

# Asynchronous By Default

- ACI (no D w/o tsync)
- Similar to asynchronous write(2) with fsync(2)
- Same purpose (performance increase)
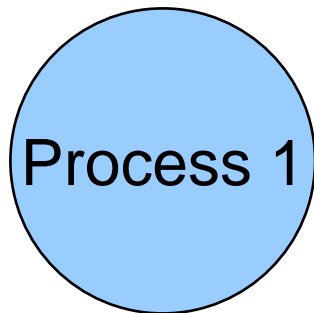- Requires page cache for files updated transactionally

STONY
BROOK
STATE UNIVERSITY OF NEW YORK

# Valor Design

- Modify page writeback to support simple write ordering

- Implement an ARIES style undo/redo log module for FS-operations

STONY BROOK
STATE UNIVERSITY OF NEW YORK

# Page Dirtying: No Txns

# Page Dirtying: With Txns

Process 1

**LEGEND:**
Old Page 🟦
U/R Page 🟨
New Page 🟥
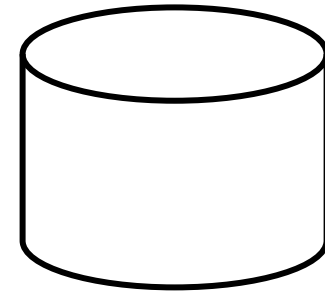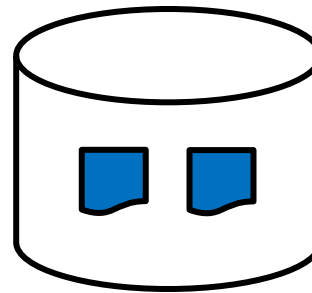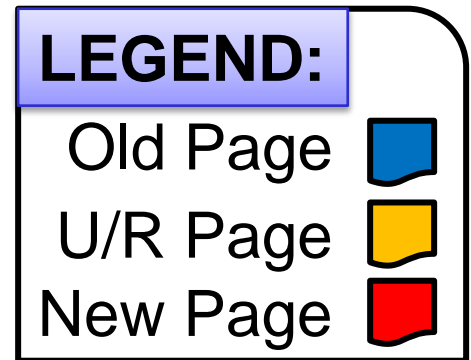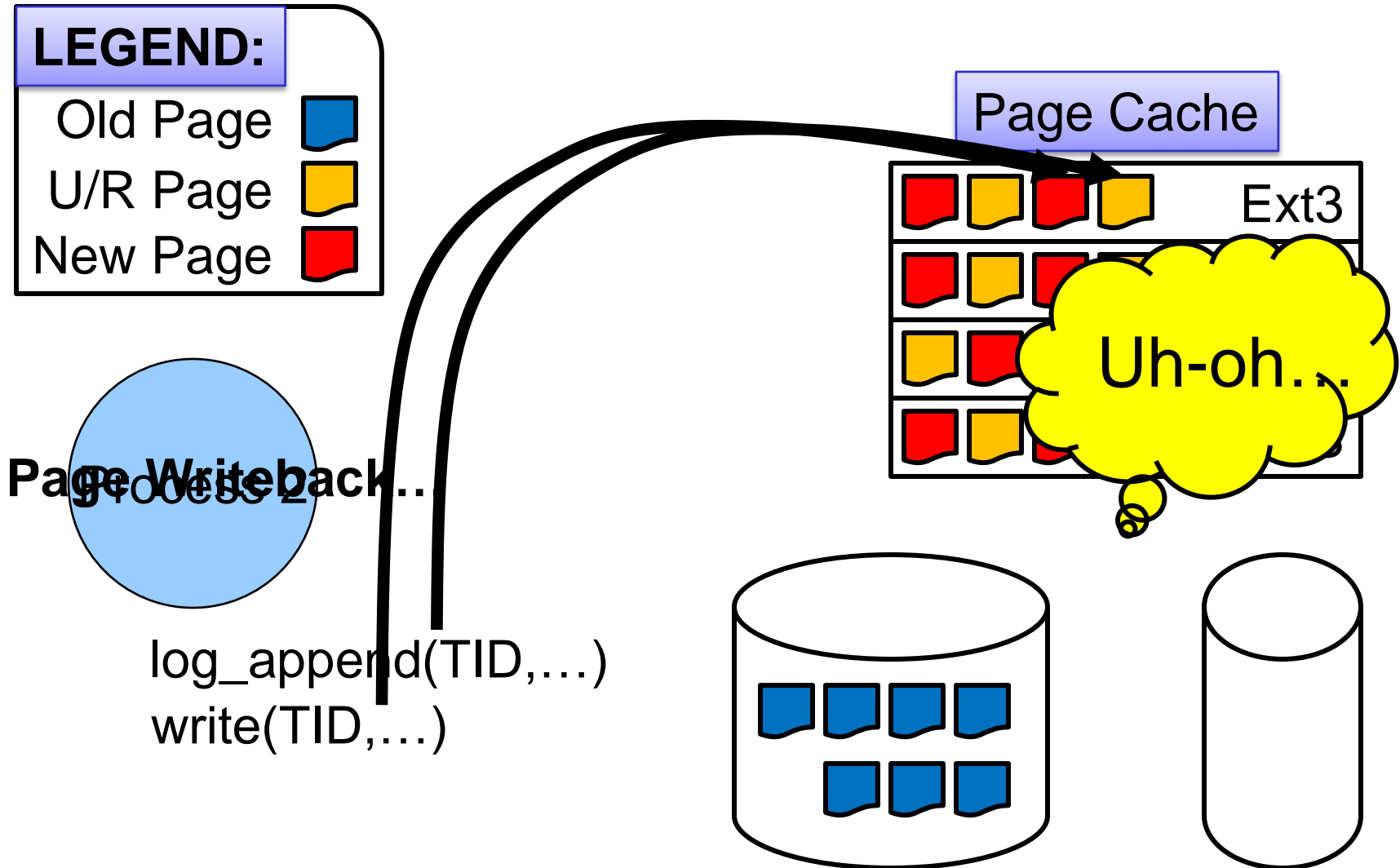
log_append(TID,… 🟨 )
log_append(TID,… 🟨 )
write(TID,… 🟥 )
write(TID,… 🟥 )

# Current Kernel Design

# What DBs Do

# Simple Write Ordering

**LEGEND:**

Old Page ▮

U/R Page ▮

New Page ▮

Page Cache

| | |
|---|---|
| ▮ ▮ | FS1 |
| ▮ ▮ | FS2 |
| ▮ | FS3 |
| ▮ ▮ | FS4 |
| ▮ ▮ ▮ ▮ | Valor |

STONY BROOK
STATE UNIVERSITY OF NEW YORK

# Log Module

Process 2

**Valor Module**

1  tbegin(TID,...)
2  tlog(TID,...)
3  ...
4  **page writeback**
5  tlog(TID,...)
6  write(TID,...)
7  tresolve(TID,...)
8  **page writeback**
9  **page writeback**

State File

Log File

Disk

U/R,1   U/R,1
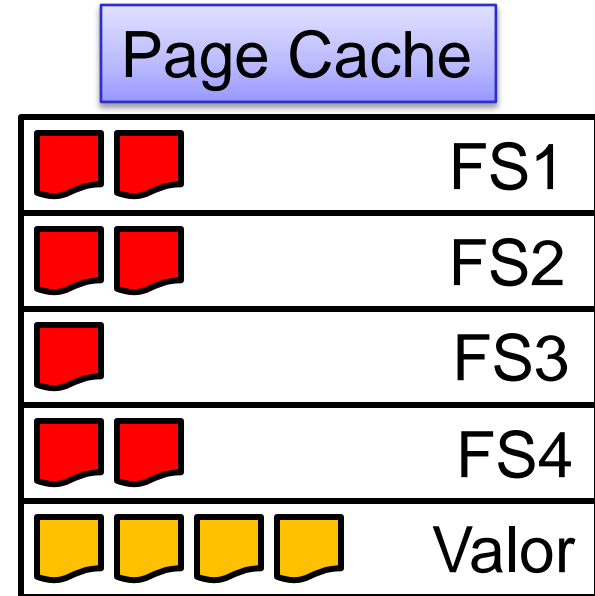| 1 | 2 | 3 |

U/R,1   U/R,1   C,1
| 4 | 5 | 6 |

Record Maps

U/R,1   U/R,1   U/R,1
| 1 | 2 | 3 |

U/R,1   U/R,1   C,1
| 4 | 5 | 6 |

U/R Page 1   1
U/R Page 2   2
U/R Page 3   3
U/R Page 4   4
U/R Page 5   5
             6

STONY BROOK
STATE UNIVERSITY OF NEW YORK

# Atomicity Argument

- Transition from pre-writeback to post-writeback disk state atomically ***iff***
  - All writes preceded by **sys_log_append**
  - Simple write ordering is implemented
  - writes to a single sector are atomic
- Valor satisfies the top 2 constraints
- A supported hard disk satisfies the third

STONY
BR●●K
STATE UNIVERSITY OF NEW YORK

# Performing Recovery

- Two kinds of recovery are supported:
  - ◆ System Recovery
  - ◆ Application Recovery (per-process abort)
- Standard recovery process:
  - ◆ Reconstruct RAM state from log
  - ◆ In reverse LSN order commit/abort landed transactions
  - ◆ Perform a page writeback

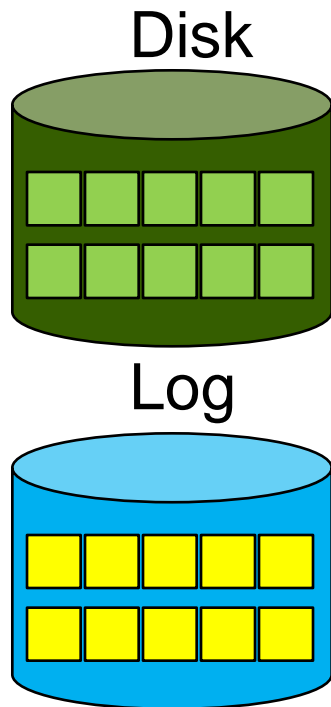STONY BROOK
STATE UNIVERSITY OF NEW YORK

# Evaluation

- We must compare against traditional asynchronous FSes
  - ◆ benchmark against asynchronous ext3
  - ◆ do serial transfer benchmarks for large files
- We turn *off* synchronous transactions for two other controls (for fairness)
  - ◆ FS built on top of Stasis
  - ◆ FS built on top of Berkeley DB

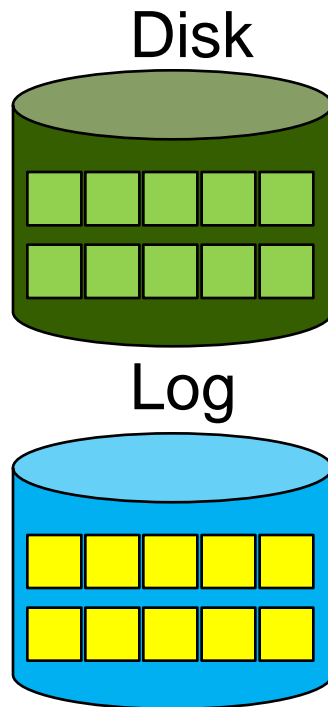STONY BROOK
STATE UNIVERSITY OF NEW YORK

# Mock ARIES Benchmark

- Important lower bound (not tight)

**MT-ow-noread**    **MT-ow**    **MT-ow-finite**

Disk    Disk    Disk

Log    Log    Log

STONY BROOK
STATE UNIVERSITY OF NEW YORK

# Mock ARIES Benchmark

# Serial Overwrite

Transaction size: 16 pages



22.75 x Ext3

5.0 x Ext3

2.75 x Ext3

# Transaction Throughput

# Conclusions

- System transactions are feasible
- Valor achieves good overhead
- Minimal changes to existing kernels

STONY
BROOK
STATE UNIVERSITY OF NEW YORK

# Limitations/Future Work

- Limitations
  - ◆ Locking slows interleaved writes to the same page
  - ◆ Some FSes/Disks do not fsync() when asked to
- Future Work
  - ◆ Explore use of logging device as a coordinator in a transactional disk array

# Q&A

## Enabling System Transactions via Lightweight Kernel Extensions

R.P. Spillane, S. Gaikwad. M. Chinni, C.P. Wright, E. Zadok

Stony Brook University

http://www.fsl.cs.sunysb.edu/

STONY BROOK
STATE UNIVERSITY OF NEW YORK

# TxF

- TxF is Microsoft's transactional file system
  - ◆ Motivation: program installation, system updates, website updates
- Pros
  - ◆ Backed by Microsoft
- Cons
  - ◆ Specific to NTFS

STONY
BR K
STATE UNIVERSITY OF NEW YORK

# Isolation

- Extended mandatory locking
  - ◆ Allows locking of directories
  - ◆ Do not have to set group exec/setgid bits
- Locking permissions
  - ◆ Let users decide if a file can be locked
- All processes acquire locks
  - ◆ Regular processes hold only for the syscall
- Lock inheritance
  - ◆ Allow multi-process transactions

# Valor != Journaling

- Journaling FSes good at **fast recovery**

- …but are too **special-purpose:**
  - *No-Steal Caching*
    - all state modified by a txn. must remain in memory until commit/abort
  - *Non-Modular Design*
    - does not handle rollback of VFS and page caches, just disk-state on boot