

15-410

Atomic Transactions

April 19, 2006

Jeffrey L. Eppinger

Professor of the Practice
School of Computer Science

So Who *Is* This Guy?

Jeff Eppinger (eppinger@cmu.edu, EDSH 229)

- Ph.D. Computer Science (CMU 1988)
- Asst Professor of Computer Science (Stanford 1988-1989)
- Co-founder of Transarc Corp. (Bought in 1994 by IBM)
 - Transaction Processing Software
 - Distributed File Systems Software
- IBM Faculty Loan to CMU eCommerce Inst. (1999-2000)
- Joined SCS Faculty in 2001
- Lecture Style: ¿Questioning?

What Do Transactions Do?

- They ensure the *consistency* of data
 - In the face of *concurrency*
 - In the face of *failure*
- They *improve performance*
 - In many cases
 - In many common cases
 - But not always

Do You Do ACID?

- What is ACID?
- The ACID properties are the guarantees provided by the transaction system:
 - Atomicity: **all or none**
 - Consistency: **if consistent before transaction, so too after**
 - Isolation: **despite concurrent execution, \exists serial ordering**
 - Durability: **committed transaction cannot be undone**

When Are Transactions Used?

- When you use:
 - Databases
 - File Systems
- Applications built on the above
 - Banking Applications
 - Web Applications
 - BeanFactory

Who Invented Atomic Transactions?

- The guys that built TP Monitors
- Most notable advocate: Jim Gray
 - The guru of transactions systems
 - Berkeley, Ph.D.
 - Famously worked at IBM
 - Now at Microsoft Research in San Francisco
 - Wrote the bible on transaction systems:

Transaction Processing: Concepts and Techniques, 1992

Outline

- ✓ *What* Do Transactions Do?
- ✓ *When* Are Transactions Used?
- ✓ *Who* Invented Atomic Transactions?
- *How*
 - How do you use transactions?
 - How do you implement them?

How do I use transactions?

```
public void deposit(int acctNum, double amount)
    throws RollbackException
{
    Transaction.begin();
    Acct a = acctFactory.lookup(acctNum);
    a.setBalance(a.getBalance()+amount);
    Transaction.commit();
}
```


Accounts are JavaBeans

```
public class Acct {
    private int    acctNum;
    private double balance;

    public Acct(int acctNum) { this.acctNum = acctNum; }

    public int    getAcctNum() { return acctNum; }
    public double getBalance() { return balance; }

    public void setBalance(double x) { balance = x; }
}
```

BeanFactory

```
public abstract class BeanFactory<B> {
    public abstract B      create(Object... priKeyValues) thr...
    public abstract void delete(Object... priKeyValues) thr...
    public abstract int  getBeanCount()                thr...
    public abstract B    lookup(Object... priKeyValues) thr...
    public abstract B[]  match(MatchArg... constaints)  thr...

    public abstract void createTable(String... priKeyNames);
    public abstract void deleteTable();
    ...
}
```

- BeanFactory uses introspection to obtain the bean properties
- Methods operating on beans throw RollbackException

Creating a BeanFactory

table name

```
BeanFactory<Acct> acctFactory =  
    BeanFactory.getInstance(Acct.class, "acct");
```

or

```
BeanFactory<Acct> acctFactory =  
    BeanFactory.getCSVInstance(Acct.class, "acct.csv");
```

- BeanFactory implementations use the Abstract Factory pattern
- There are multiple implementations of BeanFactory:
 - Using a relational database
 - Using files
- Each factory supports the same BeanFactory interface

Transactions

- Transactions are associated with threads
- When called in a transaction, beans returned by `create()`, `lookup()`, and `match()` are tracked and their changes are “saved” at commit time

```
public class Transaction {
    public static void begin() throws RollbackException {...}
    public static void commit() throws RollbackException {...}
    public static boolean isActive() {...}
    public static void rollback() {...}
}
```

The classic debit/credit example

```
public void xfer(int fromAcctNum,
                int toAcctNum,
                double amount) throws RollbackException {
    {
        Transaction.begin();
        Acct f = acctFactory.lookup(fromAcctNum);
        f.setBalance(f.getBalance()-amount);
        Acct t = acctFactory.lookup(toAcctNum);
        t.setBalance(t.getBalance()+amount);
        Transaction.commit();
    }
}
```

- Error cases not addressed (acct not found, low balance)

Remember the ACID Properties?

- ✓ Atomicity: **all or none**
 - ✓ Consistency: **if before than after**
 - ✓ Isolation: **serial ordering**
 - ✓ Durability: **cannot be undone**
-

```
public void xfer(int fromAcctNum,  
                int toAcctNum,  
                double amount) throws RollbackException {  
  
    Transaction.begin();  
    Acct f = acctFactory.lookup(fromAcctNum);  
    f.setBalance(f.getBalance() - amount);  
    Acct t = acctFactory.lookup(toAcctNum);  
    t.setBalance(t.getBalance() + amount);  
    Transaction.commit();  
  
}
```

How Are ACID Properties Enforced?

- A simple, *low-performance* implementation
 - One file holds contains all the data
 - *Atomicity* – write a new file and then use rename to replace old version
 - *Consistency* – app's problem
 - *Isolation* – locking, specifically one mutex
 - *Durability* – trust the file system (weak)

How Are ACID Properties Enforced?

- A *high-performance* implementation
 - Complex disk data structures (trees)
 - *Atomicity* – write-ahead logging
 - *Consistency* – app's problem
 - *Isolation* – two-phase locking
 - *Durability* – write-ahead logging

Write-ahead Logging

- Provides atomicity & durability
- Buffer database disk pages in a memory buffer cache
- Log all changes in a log before they are written to disk
 - When changing data pages, describe changes in log records
 - When committing, write commit-record into log, flush log
 - Before flushing cached pages, check ensure log was flushed
- Recover from the log
 - When restarting after a failure, scan the log:
 - (Case 1) Redo transactions with commit records, as necessary
 - (Case 2) Undo transactions without commit records, as necessary
 - When handling user or system initiated rollbacks:
 - (Case 3) Scan the log and undo all the work

How Do You Describe Changes?

- Value Logging
 - E.g., old value = 4, new value = 5
- Operation Logging
 - E.g., increment by 1,
 - E.g., insert file 436 into directory 123

Sample Log

*Disk
Storage*

<i><fromAcctNum></i>
balance: \$100

<i><toAcctNum></i>
balance: \$3

*Memory
Buffer Cache*

Log

⋮

*Green log
records have
been flushed to
disk*

Sample Log

```
Transaction.begin();  
Acct f = factory.lookup(fromAcctNum);  
...f.getBalance()...
```

*Disk
Storage*

<i><fromAcctNum></i>
balance: \$100

<i><toAcctNum></i>
balance: \$3

*Memory
Buffer Cache*

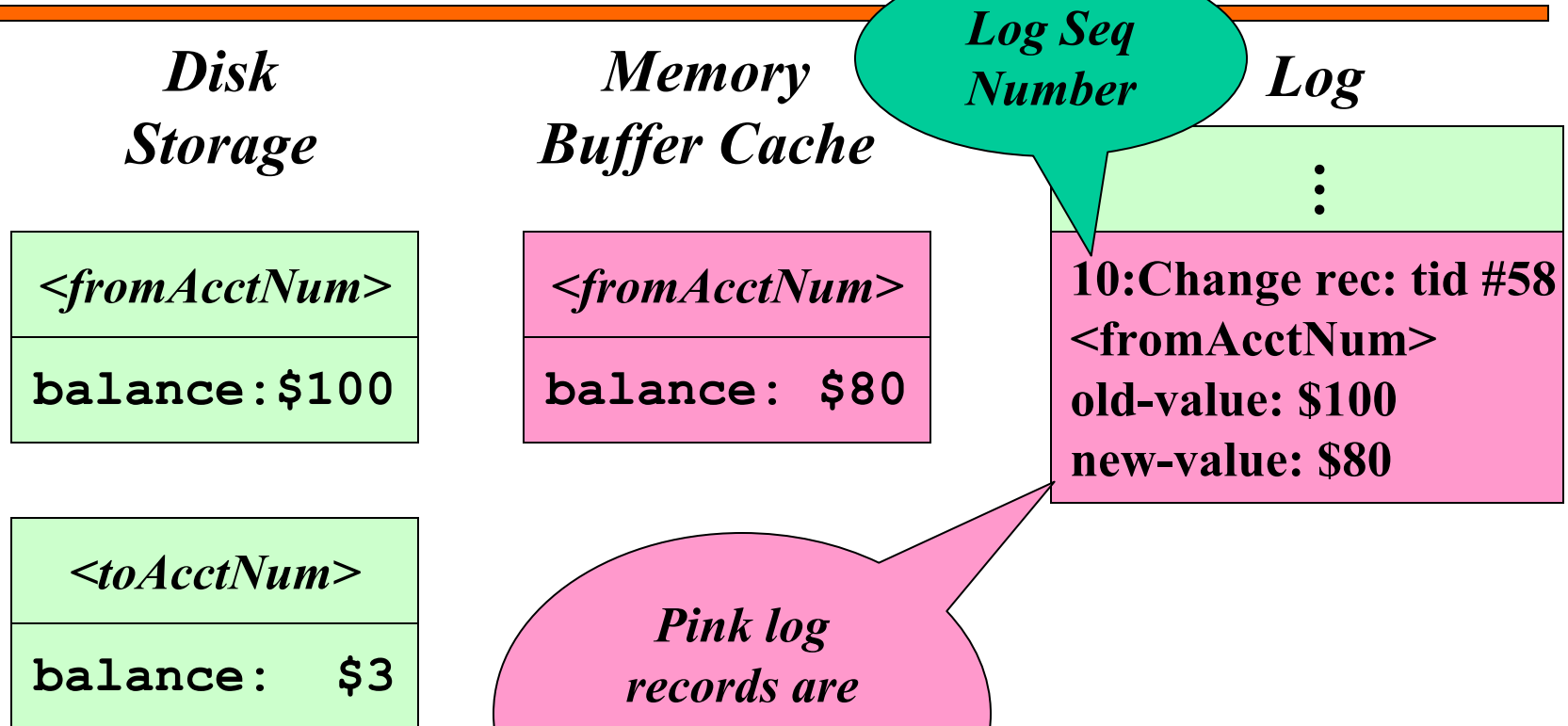
<i><fromAcctNum></i>
balance: \$100

Log

⋮

Sample Log

```
Transaction.begin();  
Acct f = factory.lookup(fromAcctNum);  
f.setBalance(f.getBalance() - amount);
```



Sample Log

```
Transaction.begin();  
Acct f = factory.lookup(fromAcctNum);  
f.setBalance(f.getBalance() - amount);  
Acct t = factory.lookup(toAcctNum);  
...t.getBalance()...
```

Disk Storage

<i><fromAcctNum></i>
balance: \$100

<i><toAcctNum></i>
balance: \$3

Memory Buffer Cache

<i><fromAcctNum></i>
balance: \$80

<i><toAcctNum></i>
balance: \$3

Log

⋮
10:Change rec: tid #58 <fromAcctNum> old-value: \$100 new-value: \$80

Sample Log

```
Transaction.begin();  
Acct f = factory.lookup(fromAcctNum);  
f.setBalance(f.getBalance() - amount);  
Acct t = factory.lookup(toAcctNum);  
t.setBalance(t.getBalance() + amount);
```

Disk Storage

<i><fromAcctNum></i>
balance: \$100

<i><toAcctNum></i>
balance: \$3

Memory Buffer Cache

<i><fromAcctNum></i>
balance: \$80

<i><toAcctNum></i>
balance: \$23

Log

⋮
10:Change rec: tid #58 <i><fromAcctNum></i> old-value: \$100 new-value: \$80
⋮
12:Change rec: tid #58 <i><toAcctNum></i> old-value: \$3 new-value: \$23

Sample Log

```
Transaction.begin();  
Acct f = factory.lookup(fromAcctNum);  
f.setBalance(f.getBalance()-amount);  
Acct t = factory.lookup(toAcctNum);  
t.setBalance(t.getBalance()+amount);  
Transaction.commit();
```

Disk Storage

<i><fromAcctNum></i>
balance: \$100

<i><toAcctNum></i>
balance: \$3

Memory Buffer Cache

<i><fromAcctNum></i>
balance: \$80

<i><toAcctNum></i>
balance: \$23

Log

⋮
10:Change rec: tid #58 <i><fromAcctNum></i> old-value: \$100 new-value: \$80
⋮
12:Change rec: tid #58 <i><toAcctNum></i> old-value: \$3 new-value: \$23
13:Commit: tid #58

Performance Improvement

- You do not need to flush the buffer cache to commit a transaction
 - Only need to flush the buffered log records
 - Great locality...all those disparate buffer cache data pages can be written out later...writes of hot pages will contain changes from many transactions
- All transactions share one log
- The log is append only and rarely read
 - So it's very efficient to write...great locality
 - Optimizations abound for increasing throughput

Recovery after System Failure:

Crash after commit (Case 1)

Disk Storage

<i><fromAcctNum></i>
balance: \$100

<i><toAcctNum></i>
balance: \$3

Memory Buffer Cache

Log

⋮
10:Change rec: tid #58 <fromAcctNum> old-value: \$100 new-value: \$80
⋮
12:Change rec: tid #58 <toAcctNum> old-value: \$3 new-value: \$23
13:Commit: tid #58

Recovery after System Failure:

Redo committed transactions (Case 1)

<i>Disk Storage</i>	<i>Memory Buffer Cache</i>	<i>Log</i>		
<table border="1"><tr><td><i><fromAcctNum></i></td></tr><tr><td>balance: \$100</td></tr></table>	<i><fromAcctNum></i>	balance: \$100	80	⋮
<i><fromAcctNum></i>				
balance: \$100				
<table border="1"><tr><td><i><toAcctNum></i></td></tr><tr><td>balance: \$3</td></tr></table>	<i><toAcctNum></i>	balance: \$3	23	10:Change rec: tid #58 <i><fromAcctNum></i> old-value: \$100 new-value: \$80
<i><toAcctNum></i>				
balance: \$3				
		⋮		
		12:Change rec: tid #58 <i><toAcctNum></i> old-value: \$3 new-value: \$23		
		13:Commit: tid #58		

Buffer Cache Can Be Flushed Mid-Transaction

```
Transaction.begin();
Acct f = factory.lookup(fromAcctNum);
f.setBalance(f.getBalance() - amount);
Acct t = factory.lookup(toAcctNum);
t.setBalance(t.getBalance() + amount);
```

Disk Storage

<i><fromAcctNum></i>
balance: \$80

<i><toAcctNum></i>
balance: \$3

Memory Buffer Cache

<i><toAcctNum></i>
balance: \$23

Log

⋮
10:Change rec: tid #58 <i><fromAcctNum></i> old-value: \$100 new-value: \$80
⋮
12:Change rec: tid #58 <i><toAcctNum></i> old-value: \$3 new-value: \$23

Be sure the relevant portion of the log is flushed before the buffer cache is flushed

Recovery after System Failure:

partial work of

Undo_^uncommitted transactions (Case 2)

*Disk
Storage*

<i><fromAcctNum></i>
balance: \$80

100

<i><toAcctNum></i>
balance: \$3

*Memory
Buffer Cache*

Log

⋮
10:Change rec: tid #58 <fromAcctNum> old-value: \$100 new-value: \$80

Rollback using the log (Case 3)

```

Transaction.begin();
Acct f = factory.lookup(fromAcctNum);
f.setBalance(f.getBalance() - amount);
Acct t = factory.lookup(toAcctNum);
t.setBalance(t.getBalance() + amount);
Transaction.rollback();

```

Disk Storage

<fromAcctNum>
balance: \$100

<toAcctNum>
balance: \$3

Memory Buffer Cache

<fromAcctNum>
balance: \$80 100

<toAcctNum>
balance: \$23 3

Log

⋮
10:Change rec: tid #58 <fromAcctNum> old-value: \$100 new-value: \$80
⋮
12:Change rec: tid #58 <toAcctNum> old-value: \$3 new-value: \$23

What else is in the log?

- You cannot afford to process the whole log at system restart
 - You need to come up quickly
- Many optimizations and special cases
 - Periodically checkpoint record are written describing the state of the buffer cache
 - Rollback records written to the log
 - Log running transactions are rolled back
 - Storing Log Sequence Numbers (LSNs) on data pages
 - Page flush records written to the log

How Are ACID Properties Enforced?

- ✓ *Atomicity* – write-ahead logging
- ✓ *Consistency* – app's problem
- *Isolation* – two-phase locking
- ✓ *Durability* – write-ahead logging

Isolation via Locking

- Multiple transactions can hold read locks (concurrently)
- Only one transaction can hold a write lock

		Lock Held by other Trans		
		None	R	W
Lock Requested	R	✓	✓	
	W	✓		

Fancy Locks (1)

- Read-intent-write...allows one transaction to lock-out other writers, but allow readers (until it upgrades to a write lock)

		Lock Held by other Trans			
		None	R	W	RIW
Lock Requested	R	✓	✓		✓
	W	✓			
	RIW	✓	✓		

Fancy Locks (2)

- Increment locks allow concurrent writes. Example: increment by x . Use operation logging: Redo: $+x$. Undo: $-x$.

		Lock Held by other Trans				
		None	R	W	RIW	Incr
Lock Requested	R	✓	✓		✓	
	W	✓				
	RIW	✓	✓			
	Incr	✓				✓

Two-phase Locking

- Grab locks and keep them until end-of-transaction, so others won't see uncommitted changes

		Lock Held by other Trans				
		None	R	W	RIW	Incr
Lock Requested	R	√	√		√	
	W	√				
	RIW	√	√			
	Incr	√				√

Avoiding Lock-out

- Locks are held on specific portions of the data
- Avoid dead-lock: E.g., ordering: if all transactions (threads) grab locks in “alphabetical” order (or any specific ordering)
- Avoid live-lock: E.g., waiting writers prevent new transactions from getting read locks

		Lock Held by other Trans				
		None	R	W	RIW	Incr
Lock Requested	R	√	√		√	
	W	√				
	RIW	√	√			
	Incr	√				√

How Does Data Get Written to Disk?

- Does the OS buffer the writes?
- Does the disk write happen atomically?

What is the Atomicity of Disk Writes?

- When you write to the disk, does it all go out?
 - Sector = 512 bytes
 - Track = n Sectors
 - Block (or page) = m Sectors
- OS writes blocks/pages
- Disk has ECC codes...can detect partial sector
 - Often there is hardware support (NV memory buffer)
- We steal a few bits on each sector to detect partial blocks / pages
 - Often there are extra bits in the sector header
 - Often we will store LSN in the sector/header or block

Bad blocks

- A block is bad if it's partially written
 - ECC detects sector error
 - Our tags on the sectors don't match
- If a log block is bad...it had better be part of the last write...good idea: mirror the log
- If data block (page) is bad...restore from backup and apply all committed changes

Remind You of Something?

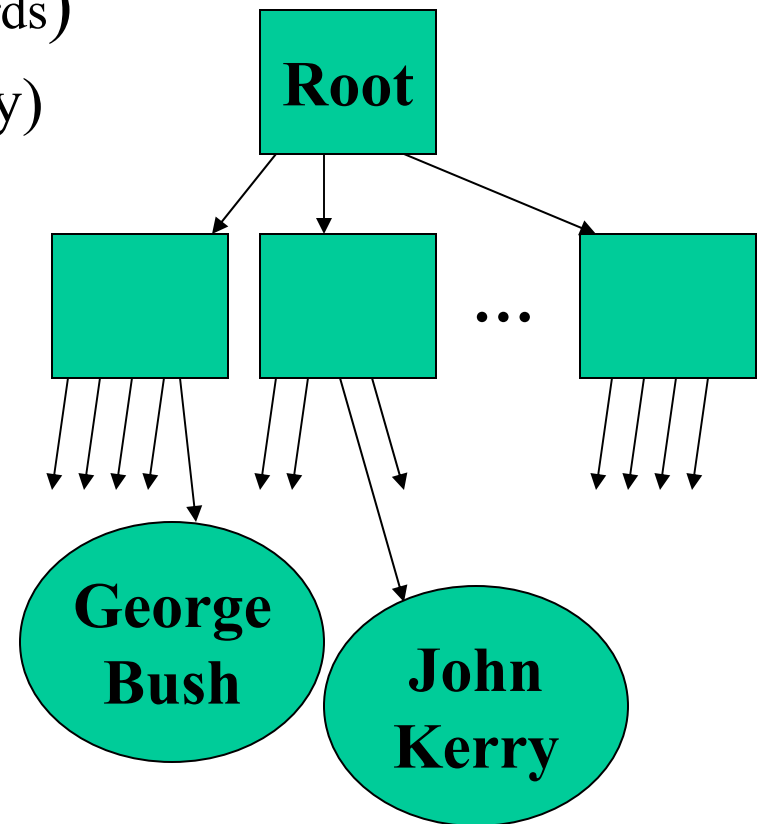
- A Relational Database
 - Any database

Why Is This Relevant to OS?

- Databases stole all this from operating systems and transaction systems
- Some OS services are better implemented using ACID properties
 - Journalling file systems
- Let's start in the beginning...

In the Old Days: OS provided

- Structured files (containing records)
 - Entry-sequenced (append-only)
 - Relative (array)
 - B-tree clustered (hash table)
- Secondary access methods
- Many field types
 - Character data
 - Integers
 - Floats
 - Dates



In the Old Days: TPM provided

- ACID properties for the OS files
 - Transactions
 - Logging
 - Recovery

Today: Relational Databases

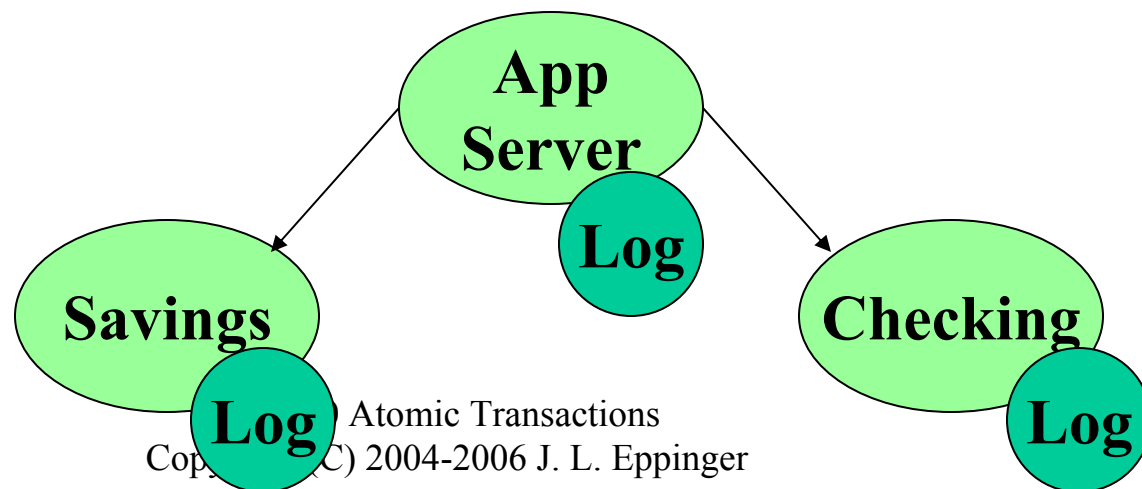
- Structured files
- ACID properties
- SQL Interface

History

- First, atomic transactions were added on at application-level (in TP Monitors)
- Then they were added to OS (mostly research OSs)
- Then they were back in the app with RBDs
- Then they were generalized to create DTP

Distributed Two-Phase Commit

- You can have distributed transactions
 - RPC, access multiple databases, etc
 - DTP: Prepare Phase (subs flush), Commit Phase (coord flush)



Why Do You Care?

- RDBs are happy to manage whole disks
- There is more to life than relational data
 - HTML, Images, Office Docs, Source, Binaries
- If you don't otherwise need a RDB, put your files in a file system

File Systems & Transactions

- If you don't allow user-level apps to compose transactions, implementation is easier
- FS Ops that require ACID properties:
 - For sure: create, delete, rename, modify properties
 - Often: write

How File Systems Implement ACID?

- Older/cheaper file systems are not log-based
 - Carefully writing to the disk
 - scandisk, chkdsk, fsck
- Newer file systems are log-based
 - E.g., NTFS, Network Appliance's NFS, JFS
 - Transactions are specialized
 - Not running general, user provided transactions
 - creat(), rename()
 - Allows specialized locking and logging

Any Questions?
