

# On Open Nesting in Distributed Transactional Memory

Alexandru Turcu\*

Binoy Ravindran\*

June 4, 2012

---

\*Virginia Tech

- Introduction
  - ◆ Distributed Transactional Memory
  - ◆ Nested Transactions
- Open Nesting
- TFA with Open Nesting (TFA-ON)
  - ◆ Transactional Forwarding Algorithm
  - ◆ TFA-ON
- Experimental analysis
- Conclusions

## Overview

Introduction

Open Nesting

TFA-ON

Experiments

Conclusions

# Introduction

Overview

**Introduction**

DTM

Transactions

Nested Transactions

Example

Open Nesting

TFA-ON

Experiments

Conclusions

# Distributed Transactional Memory

- Promising new model for programming distributed concurrency
  - ◆ Aims to replace distributed locks
  - ◆ Avoids their problems: distributed dead-locks, live-locks, difficult code composition
- Employs transactions
  - ◆ Successful abstraction originating in the database community
  - ◆ Provide atomicity, consistency, isolation

Overview

Introduction

**DTM**

Transactions

Nested Transactions

Example

Open Nesting

TFA-ON

Experiments

Conclusions

# Introducing Transactions

- We consider a redo-log approach:
  - ◆ Write ops are buffered to a write-set
  - ◆ Read ops first look at the write-set
  - ◆ Read ops are recorded in a read-set → used for validation (to make sure objects did not change since first seen)
- On commit, changes are propagated to the globally committed memory
- On abort, changes are discarded and the transaction is retried

Overview

Introduction

DTM

**Transactions**

Nested Transactions

Example

Open Nesting

TFA-ON

Experiments

Conclusions

# Nested Transactions

- Nesting is used to enable code composability
  - ◆ Transaction enclosed within another transaction
- Three types, based on parent/children interactions:
  - ◆ Flat nesting (monolithic transactions), conflict in child aborts parent
  - ◆ Closed nesting, children can abort independently
  - ◆ Open nesting, child releases isolation early

[Overview](#)

[Introduction](#)

[DTM](#)

[Transactions](#)

[Nested Transactions](#)

[Example](#)

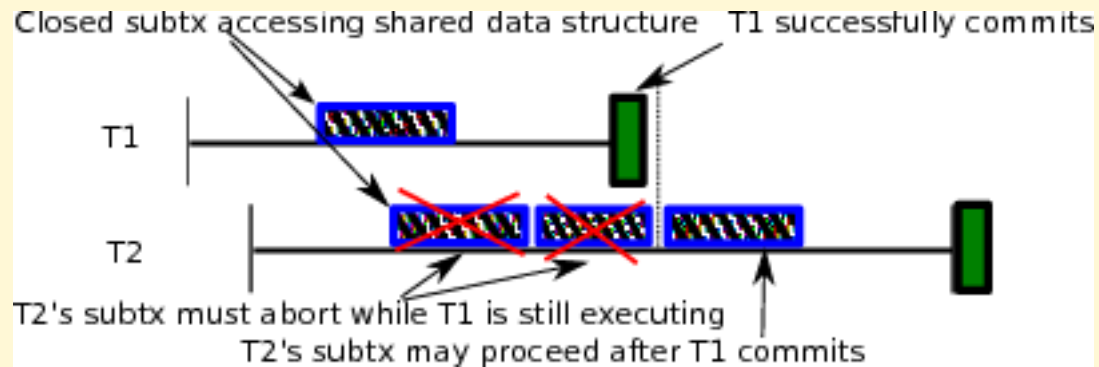
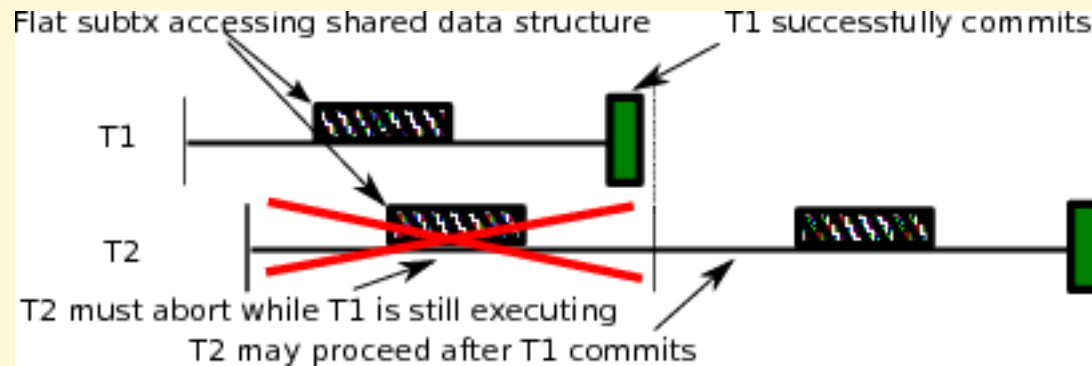
[Open Nesting](#)

[TFA-ON](#)

[Experiments](#)

[Conclusions](#)

# Flat vs Closed nesting



[Overview](#)

[Introduction](#)

[DTM](#)

[Transactions](#)

[Nested Transactions](#)

[Example](#)

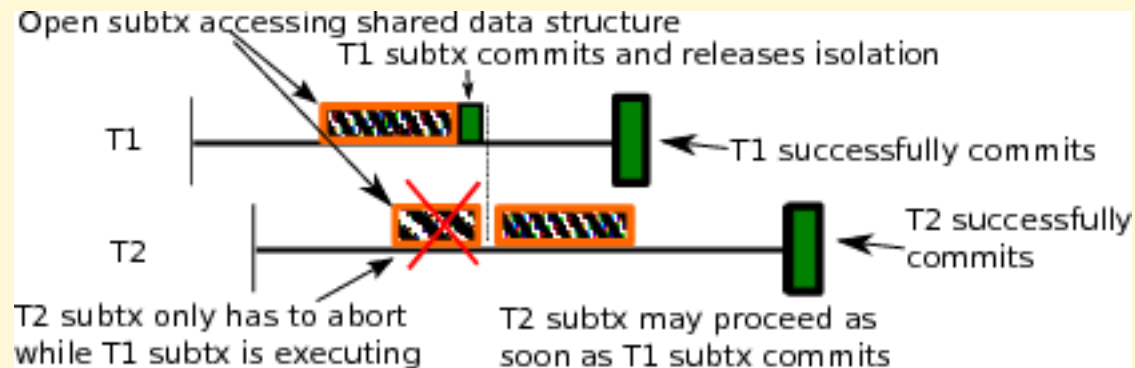
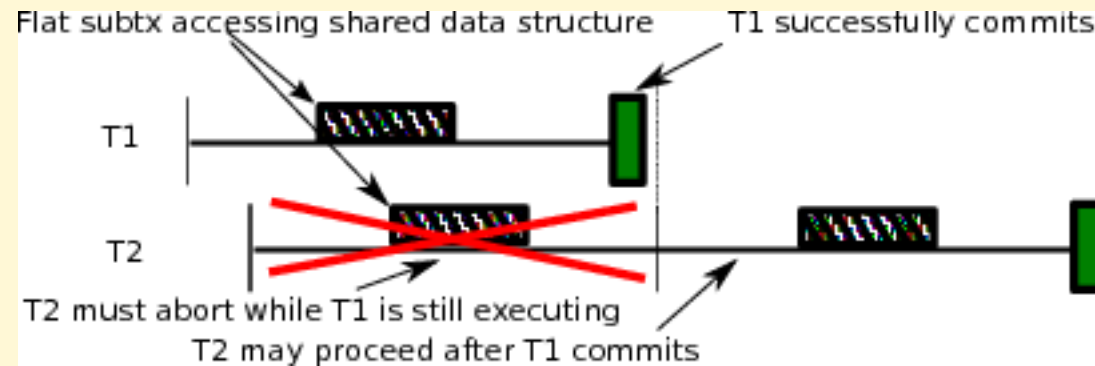
[Open Nesting](#)

[TFA-ON](#)

[Experiments](#)

[Conclusions](#)

# Flat vs Open nesting



Overview

Introduction

DTM

Transactions

Nested Transactions

Example

Open Nesting

TFA-ON

Experiments

Conclusions



# Open Nesting

Overview

Introduction

**Open Nesting**

Multilevel Txn

Safety

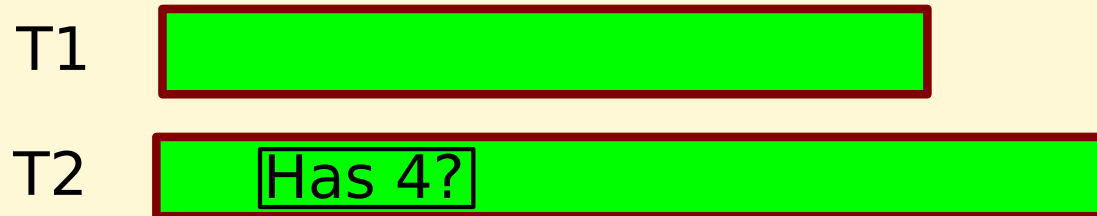
TFA-ON

Experiments

Conclusions

# Multilevel Transactions

- Why? Consider an example:
  - ◆ Set implemented using a skip-list
  - ◆ Operations accessing neighboring nodes conflict
  - ◆ For short transactions, this is OK
  - ◆ Long transactions would abort in vain.



Overview

Introduction

Open Nesting

**Multilevel Txn**

Safety

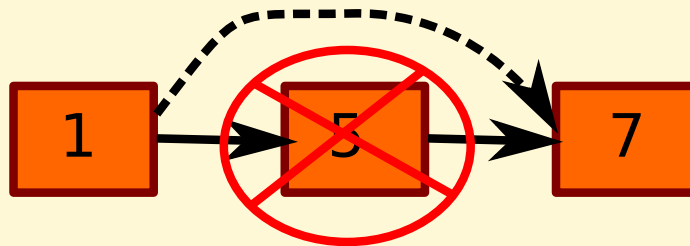
TFA-ON

Experiments

Conclusions

# Multilevel Transactions

- Why? Consider an example:
  - ◆ Set implemented using a skip-list
  - ◆ Operations accessing neighboring nodes conflict
  - ◆ For short transactions, this is OK
  - ◆ Long transactions would abort in vain.



T1 Del 5! Commit

T2 Has 4? Abort

Overview

Introduction

Open Nesting

**Multilevel Txn**

Safety

TFA-ON

Experiments

Conclusions

# Multilevel Transactions

- False conflicts:
  - ◆ Memory operations do conflict
  - ◆ Conceptually, no conflict
- Separate transactions into multiple levels of abstraction:
  - ◆ Make transactions at the lower (memory) level shorter, releasing isolation earlier
  - ◆ Preserve **fundamental conflicts** at the higher level
    - When two ops can not commute (e.g., *Has 5?* and *Del 5!*)
    - Detect using other mechanisms (e.g., abstract locks)

Overview

Introduction

Open Nesting

**Multilevel Txn**

Safety

TFA-ON

Experiments

Conclusions

# Multilevel transactions

- When parent needs to abort after child committed:
  - ◆ Revert the data structure to its original (abstract) state by applying a compensating action
  - ◆ Example: to compensate for adding 5 to a set → remove 5
- Abstract locks:
  - ◆ Do not lock physical memory locations, thus *abstract*
  - ◆ Can be used to implement mutual exclusion, R/W, etc.

Overview

Introduction

Open Nesting

**Multilevel Txn**

Safety

TFA-ON

Experiments

Conclusions

# Open Nesting Safety

- Because isolation is released early, maintaining correctness becomes the task of the programmer
  - ◆ Correct usage of abstract locks is paramount
  - ◆ Open nesting recommended for experts only (e.g., library devs)

Overview

Introduction

Open Nesting

Multilevel Txn

**Safety**

TFA-ON

Experiments

Conclusions

# TFA-ON

Overview

Introduction

Open Nesting

**TFA-ON**

TFA

TFA-ON

Locks

Defining Txn

Experiments

Conclusions

# Transactional Forwarding Algorithm

- TFA is an existing protocol for distributed STM
  - ◆ Based around Transactional Locking II and Lamport clocks
- Provides a way to establish "happens before" relationships
  - ◆ Each node holds a node-local clock
  - ◆ Clock value affixed to all messages
  - ◆ Clock incremented on local transactions' commits
  - ◆ When a message from a node with a higher clock is received, local clock is updated

Overview

Introduction

Open Nesting

TFA-ON

**TFA**

TFA-ON

Locks

Defining Txn

Experiments

Conclusions



# Transactional Forwarding Algorithm

- Each txn stores its starting time
- When a txn communicates with a node with a higher clock:
  - ◆ Attempt to update txn's starting time (i.e. *transactional forwarding*)
  - ◆ Must validate read-set before forwarding
    - Success → update txn starting time and continue
    - Failure → abort txn

Overview

Introduction

Open Nesting

TFA-ON

**TFA**

TFA-ON

Locks

Defining Txn

Experiments

Conclusions

- Our contribution: Transactional Forwarding Algorithm with Open Nesting (TFA-ON), implemented in HyFlow Java DTM framework.
- Open nested transactions in TFA-ON resemble the root transactions of TFA:
  - ◆ They record the local clock when they start
  - ◆ They increment the local clock upon commit
- For closed nested sub-transactions, the closest open-nested ancestor acts as a *local root*
  - ◆ Objects in all levels underneath this *local root* are considered for validation
  - ◆ When performing transactional forwarding, the starting time of the *local root* is updated

Overview

Introduction

Open Nesting

TFA-ON

TFA

**TFA-ON**

Locks

Defining Txn

Experiments

Conclusions

# Abstract Locks in TFA-ON

- Abstract lock acquisition:
  - ◆ With respect to the open-nested child → late acquisition (at commit time)
  - ◆ With respect to the parent transaction → early acquisition (encounter time)
- Deadlocks are possible
  - ◆ We chose to abort the whole transaction chain when an abstract lock acquisition fails.
  - ◆ This ensures all abstract locks are released, avoiding deadlocks.

[Overview](#)

[Introduction](#)

[Open Nesting](#)

[TFA-ON](#)

[TFA](#)

[TFA-ON](#)

[Locks](#)

[Defining Txn](#)

[Experiments](#)

[Conclusions](#)

# Defining Transactions

```
new Atomic<Boolean>(NestingModel.OPEN) {
  @Override boolean atomically(Txn t) {
    // ...
    t.acquireAbsLock(bst, lockName);
    return true;
  }
  @Override void onAbort(Txn t) {
    // ...
    t.releaseAbsLock(bst, lockName);
  }
  @Override void onCommit(Txn t) {
    t.releaseAbsLock(bst, lockName);
  }
}.execute();
```

[Overview](#)

[Introduction](#)

[Open Nesting](#)

[TFA-ON](#)

[TFA](#)

[TFA-ON](#)

[Locks](#)

[Defining Txn](#)

[Experiments](#)

[Conclusions](#)

# Experiments

Overview

Introduction

Open Nesting

TFA-ON

**Experiments**

Setup

Results

Conclusions

- Micro-benchmarks
  - ◆ Three distributed data structures (skip-list, hash-table, binary search tree)
  - ◆ Enhanced counter application
- Parameters:
  - ◆ Read-only ratio ( $r$ )
  - ◆ Number of calls ( $c$ )
  - ◆ Key domain size ( $k$ )
- Measured quantities (among others):
  - ◆ Throughput
  - ◆ Committed/aborted transactions, sub-transactions, commit/compensating actions
- 48-node AMD Opteron testbed running Linux

[Overview](#)

[Introduction](#)

[Open Nesting](#)

[TFA-ON](#)

[Experiments](#)

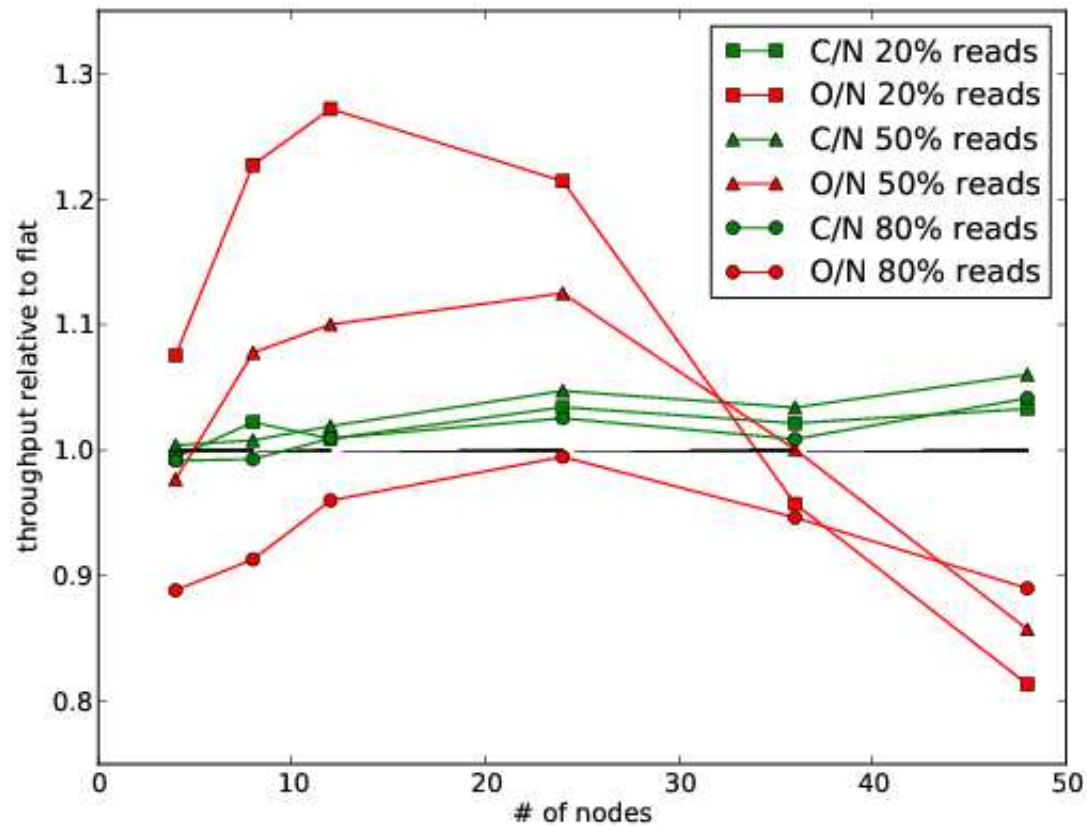
[Setup](#)

[Results](#)

[Conclusions](#)

# Results: Throughput

Throughput relative to flat, on skip-list, with  $c=3$



[Overview](#)

[Introduction](#)

[Open Nesting](#)

[TFA-ON](#)

[Experiments](#)

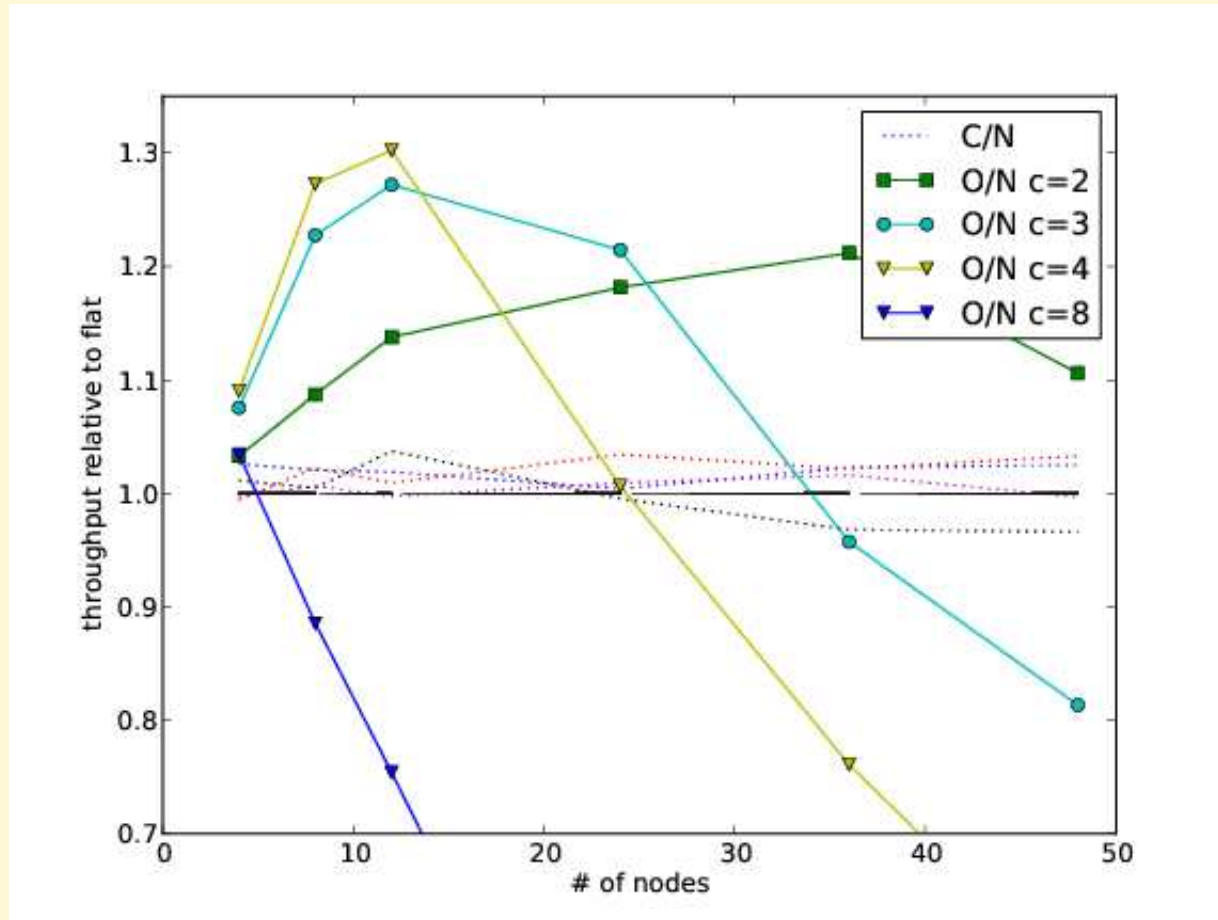
[Setup](#)

[Results](#)

[Conclusions](#)

# Results: Throughput

Throughput relative to flat, on skip-list, with 20% reads



[Overview](#)

[Introduction](#)

[Open Nesting](#)

[TFA-ON](#)

[Experiments](#)

[Setup](#)

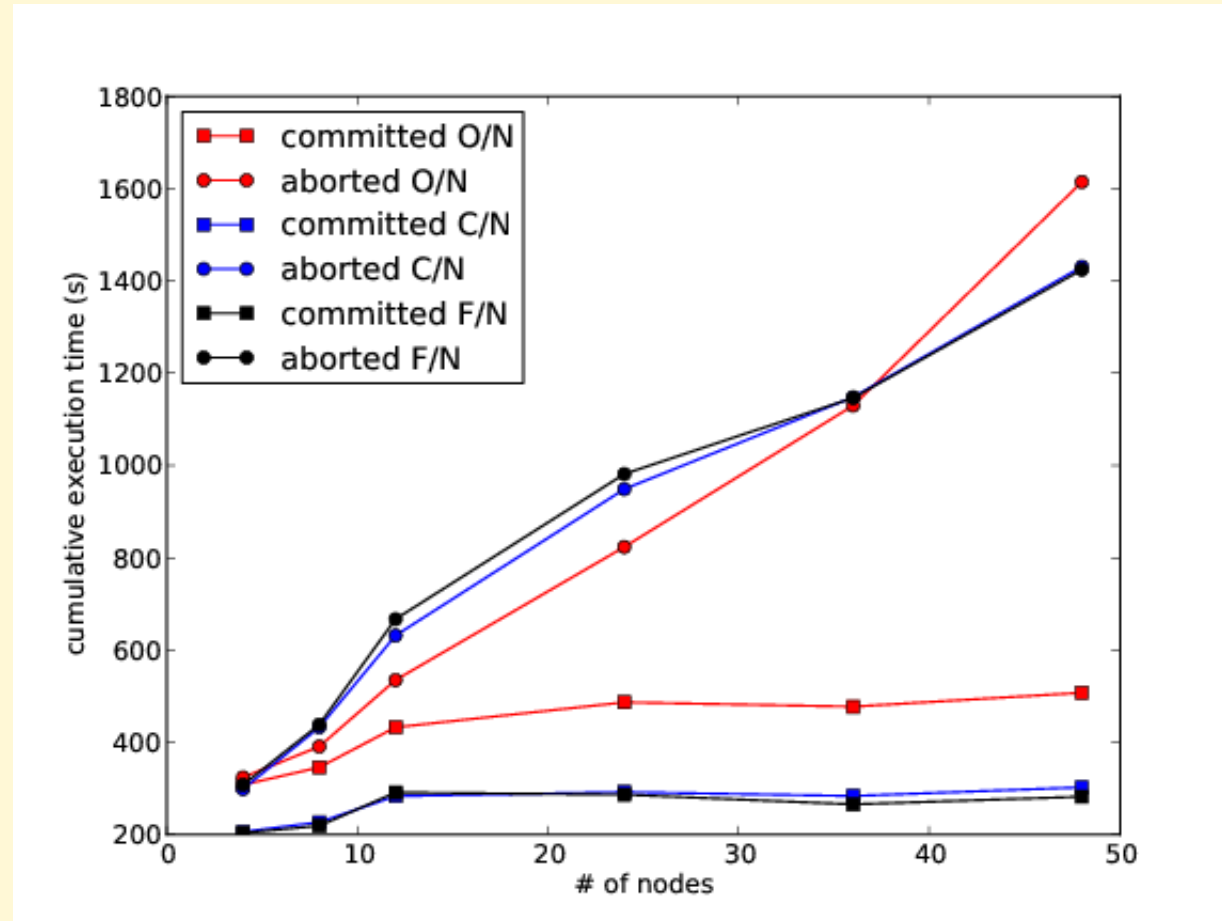
[Results](#)

[Conclusions](#)



# Results: abort vs. commit time

Time in committed vs aborted, hash-table,  $r=20$  and  $c=4$



[Overview](#)

[Introduction](#)

[Open Nesting](#)

[TFA-ON](#)

[Experiments](#)

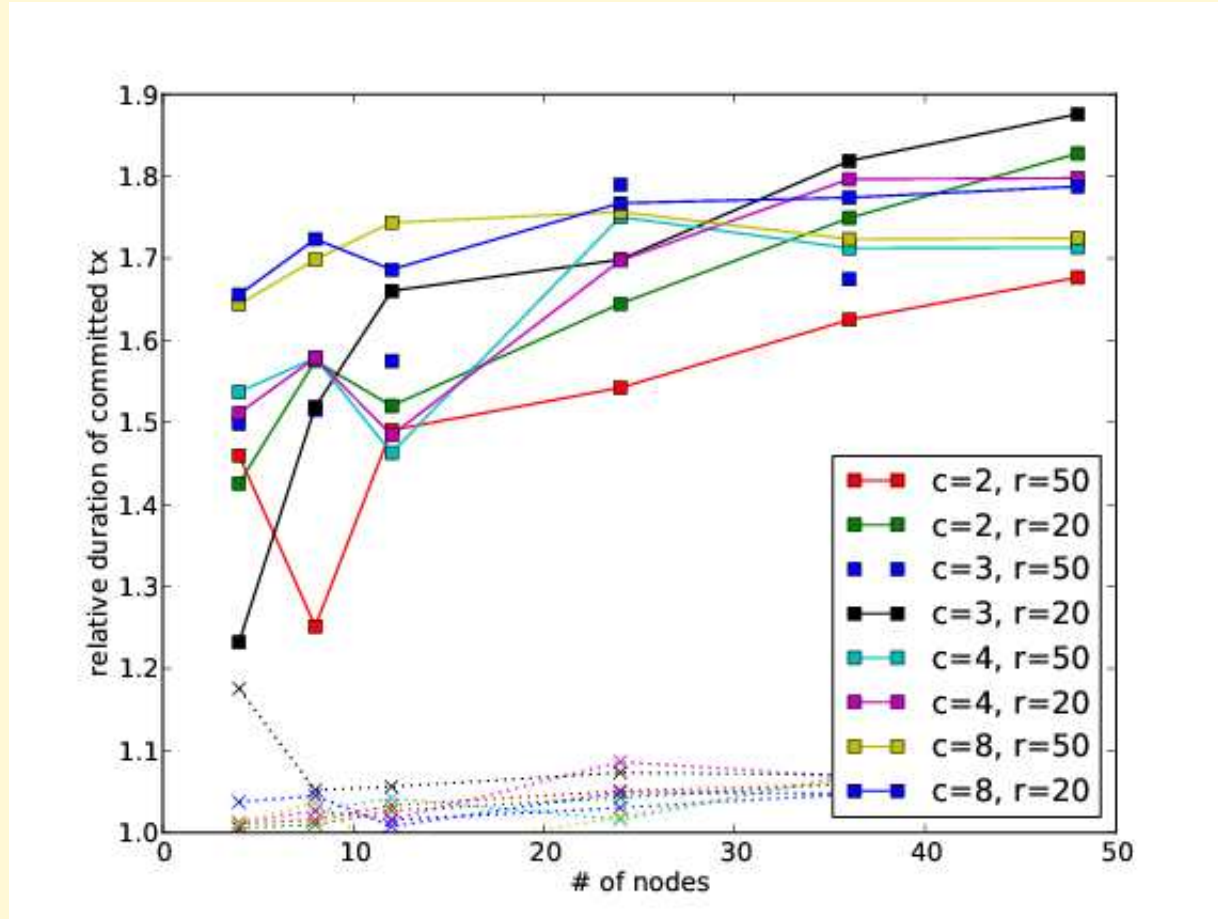
[Setup](#)

[Results](#)

[Conclusions](#)

# Results: commit overheads

## Successful transactions relative overheads, on hash-table



Overview

Introduction

Open Nesting

TFA-ON

Experiments

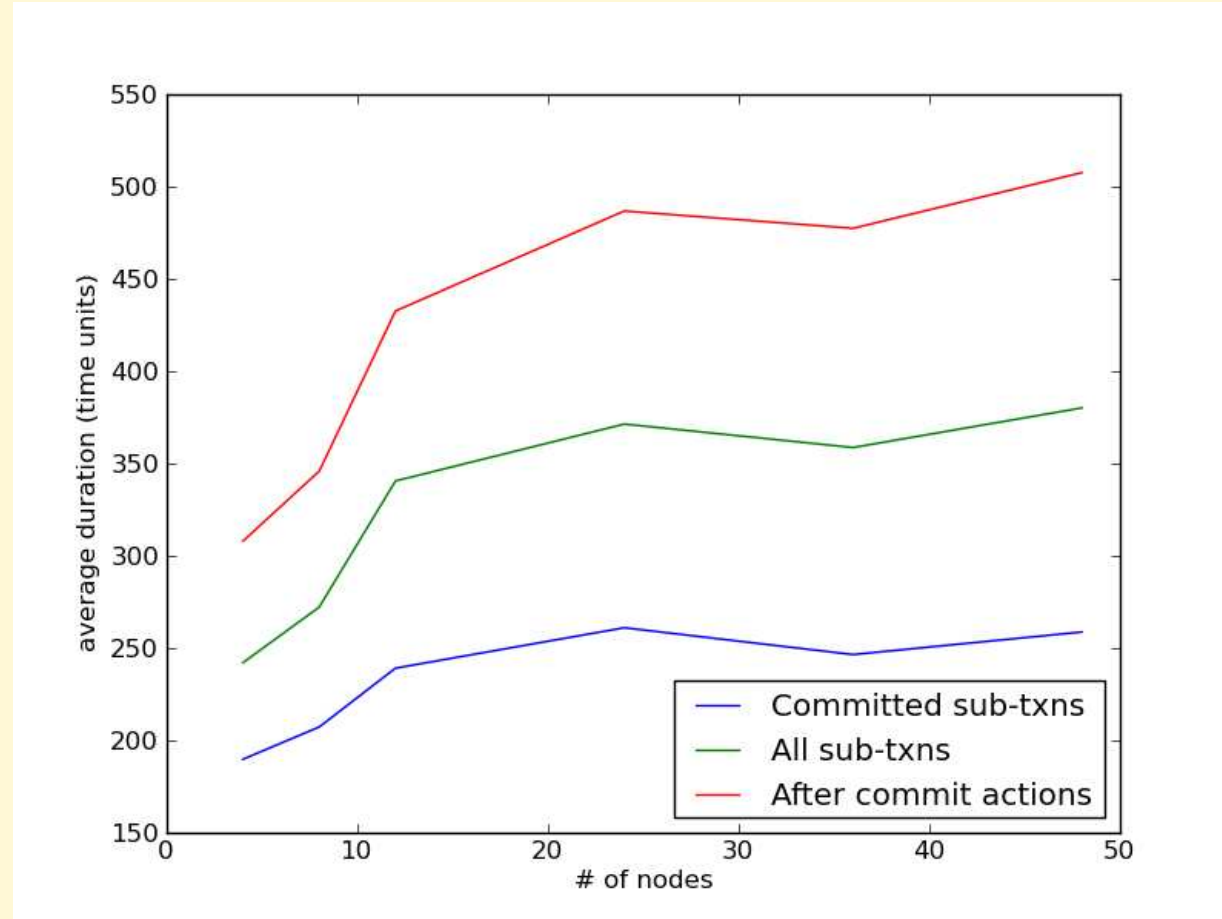
Setup

Results

Conclusions

# Results: breakdown

Breakdown time in successful txn, hash-table,  $r=20$ ,  $c=4$



[Overview](#)

[Introduction](#)

[Open Nesting](#)

[TFA-ON](#)

[Experiments](#)

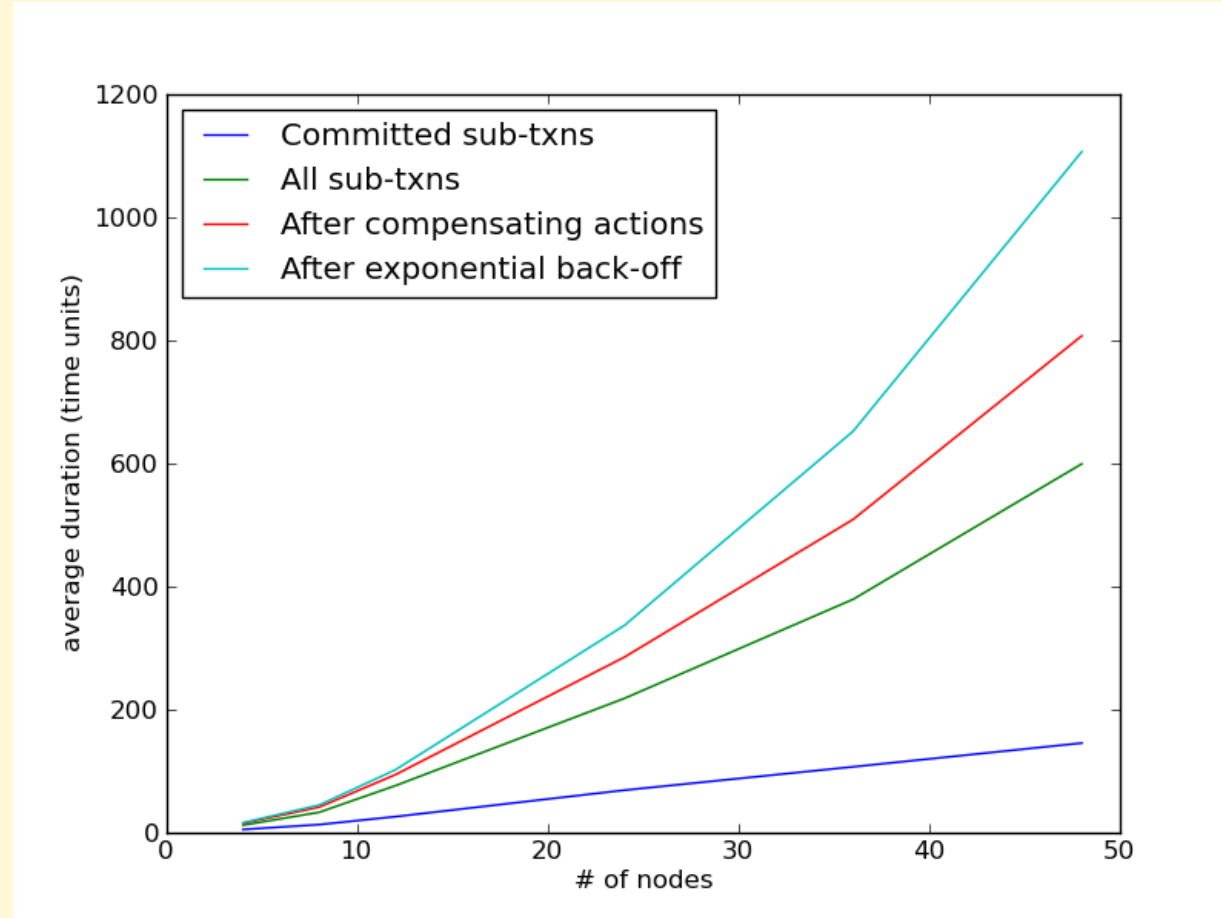
[Setup](#)

[Results](#)

[Conclusions](#)

# Results: breakdown

Breakdown time in failed txn, hash-table,  $r=20$ ,  $c=4$



[Overview](#)

[Introduction](#)

[Open Nesting](#)

[TFA-ON](#)

[Experiments](#)

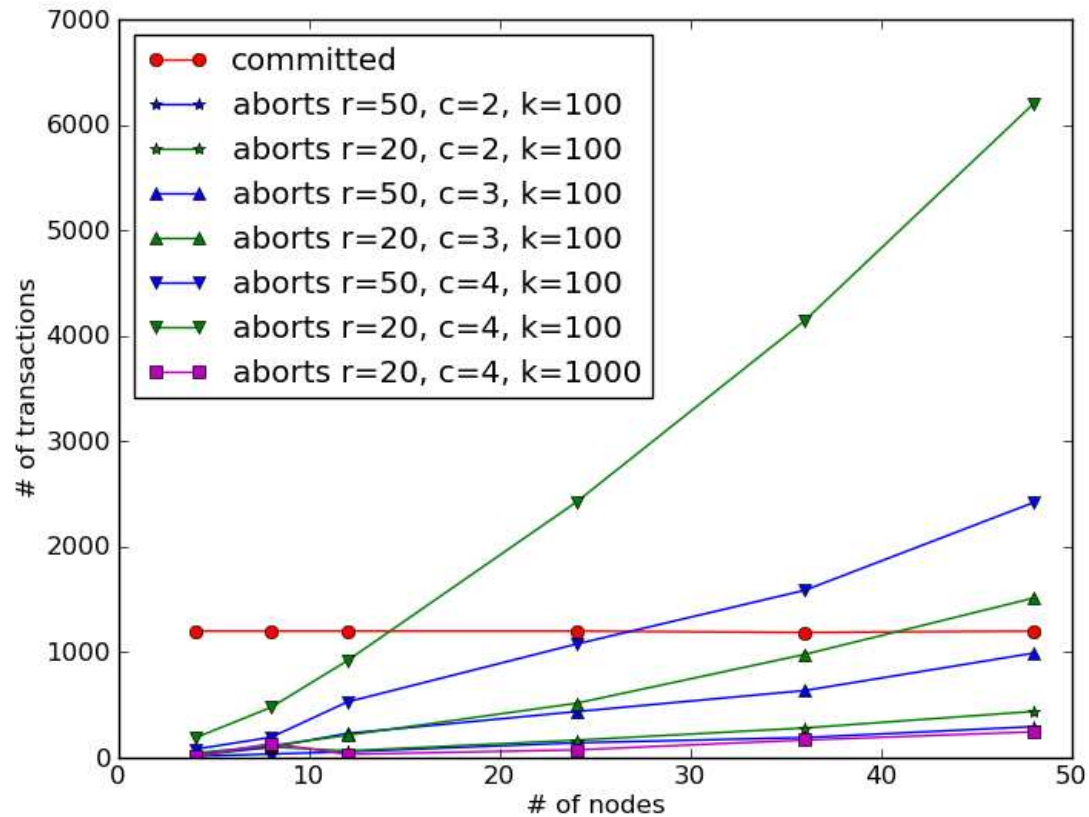
[Setup](#)

[Results](#)

[Conclusions](#)

# Results: number of aborts

## Number of aborts vs commits, on hash-table



Overview

Introduction

Open Nesting

TFA-ON

Experiments

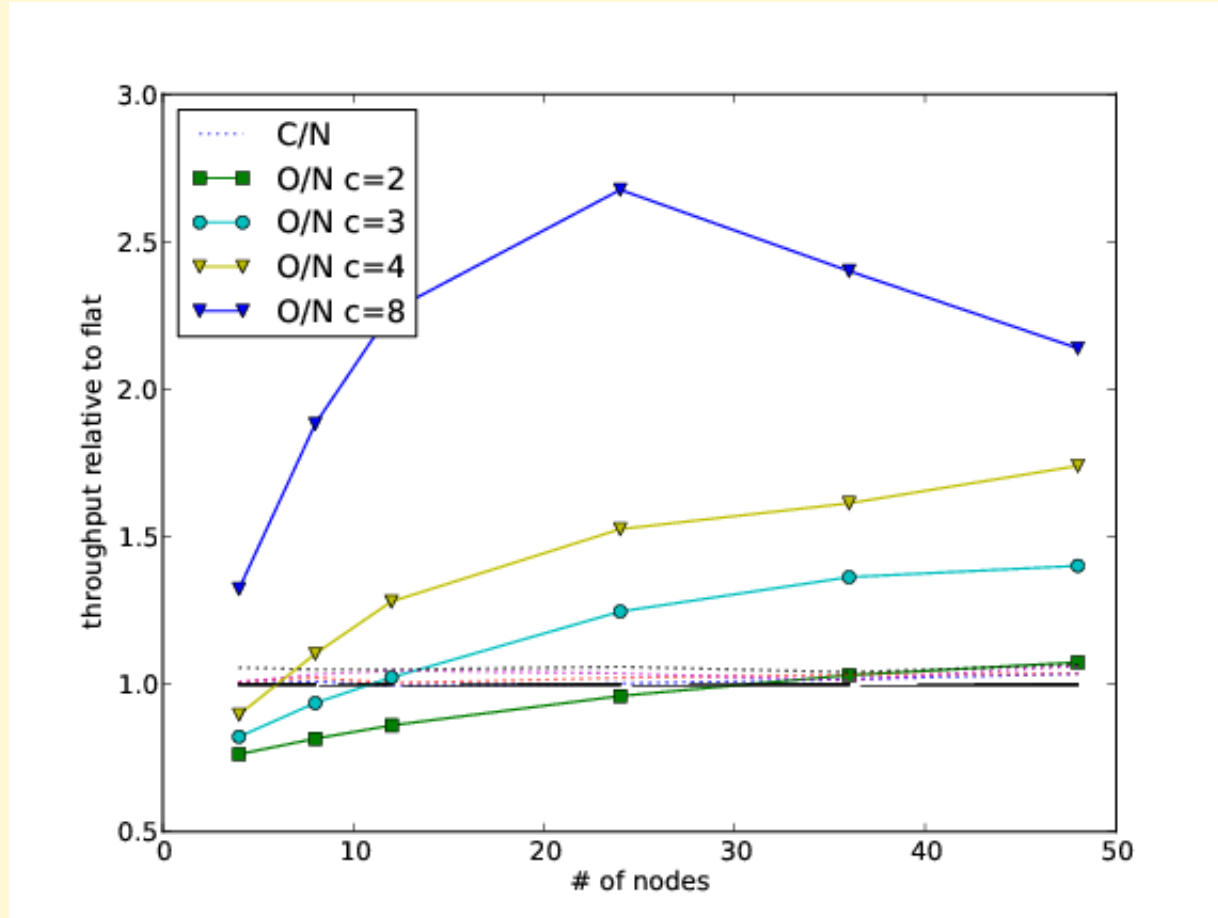
Setup

Results

Conclusions

# Results: increased key-space

Throughput, on hash-table,  $r=20$ ,  $k=1000$



[Overview](#)

[Introduction](#)

[Open Nesting](#)

[TFA-ON](#)

[Experiments](#)

[Setup](#)

[Results](#)

[Conclusions](#)

# Conclusions

Overview

Introduction

Open Nesting

TFA-ON

Experiments

**Conclusions**

Conclusions

Factors limiting performance:

- Commit overhead at low node-count
  - ◆ Significant in read-dominated workloads
- Increased fundamental conflicts at high node-count
  - ◆ Depends on available key-space for abstract locking

Open nesting optimistically assumes the parent will commit. Contention in parent after open-nested child reduces benefit.

[Overview](#)

[Introduction](#)

[Open Nesting](#)

[TFA-ON](#)

[Experiments](#)

[Conclusions](#)

[Conclusions](#)