# Persistent hardware transactions can scale

João Barreto

(and Daniel Castro, Alexandro Baldassin, Paolo Romano)
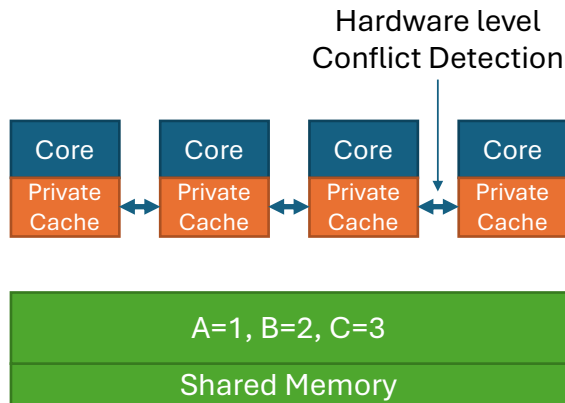
# The **transactional memory** abstraction

- Great for concurrency control
- One of the most popular abstraction for "failure-atomic critical sections" in PM literature

```
1  TX_BEGIN(pool) {
2      TX_ADD_DIRECT(&D_RW(FOO));
3      TX_ADD_DIRECT(&D_RW(BAR));
4      D_RW(FOO) = D_RO(FOO)-100;
5      D_RW(BAR) = D_RO(BAR)+100;
6  } TX_END
```

A persistent memory transaction with PMDK

# Hardware transactional memory (HTM)



Transactional Memory:
Architectural Support for Lock-Free Data Structures

Maurice Herlihy
Digital Equipment Corporation
Cambridge Research Laboratory
Cambridge MA 02139
herlihy@crl.dec.com

J. Eliot B. Moss
Dept. of Computer Science
University of Massachusetts
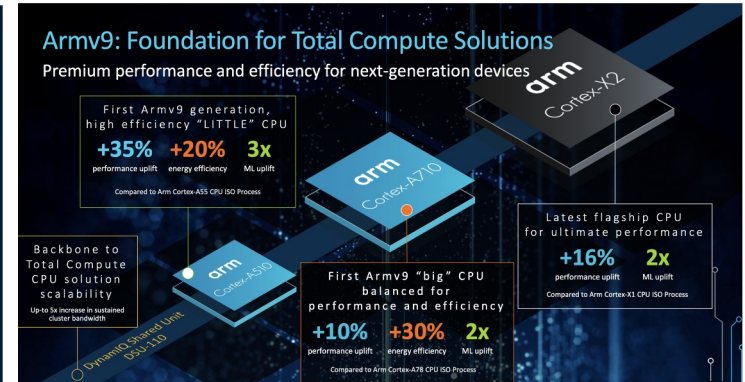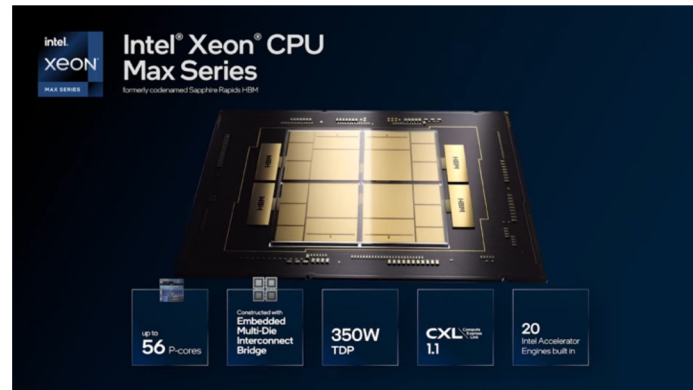Amherst, MA 01003
moss@cs.umass.edu

**Abstract**

A shared data structure is *lock-free* if its operations do not require mutual exclusion. If one process is interrupted in the middle of an operation, other processes will not be...

Herlihy and Moss, 1993

Hardware level
Conflict Detection



```
_xbegin;
foo = foo - 100;
bar = bar + 100;
_xend;
```

# After the multi-core revolution…
# … the persistent memory (PM) revolution



HTM provides **opacity**



hardware transactions can access persistent objects

However, HTM doesn't guarantee **durable opacity**

# How to help HTM support **persistent hardware transactions**?

# Why not just using write-ahead logging?

- Writes to PM also added to a durable redo log (in PM)
- However, we cannot flush redo log entries to PM before the HTM commits the transaction

Begin HW transaction

Log(X)   clflush   **Abort**

W(X,1)

time

On Cache

On Memory (NVM)

**Externalization of cache-lines while the transactions is running causes it to abort!**

# Durable hardware transactions based on a shadow copy



We decouple transaction isolation (via HTM) from durability (via WAL)

**Challenge: Persistence order** must be consistent with **happens-before order**

# But does it scale?...

## with NV-HTM [IPDPS'18]

T$_0$ starts

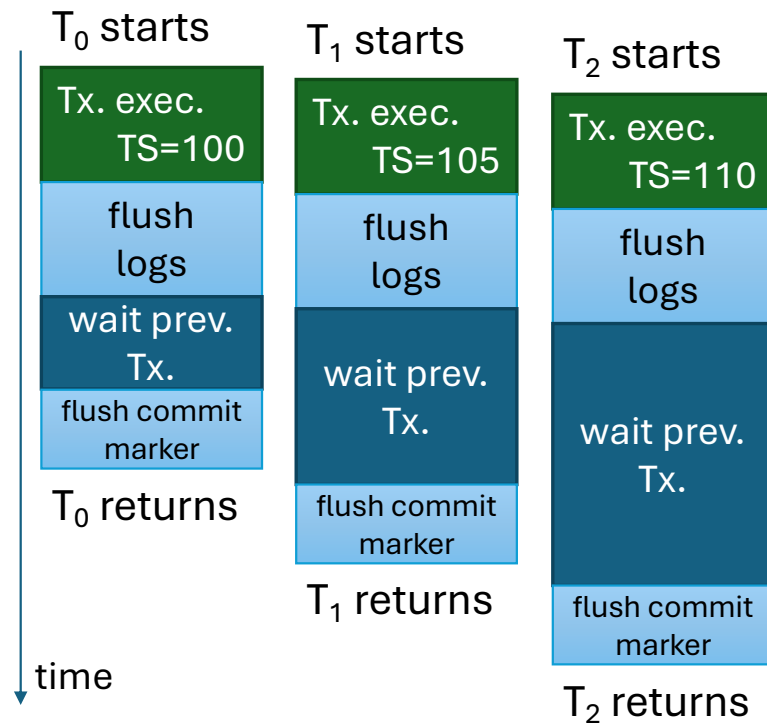| Tx. exec. TS=100 |
| flush logs |
| wait prev. Tx. |
| flush commit marker |

T$_0$ returns

T$_1$ starts

| Tx. exec. TS=105 |
| flush logs |
| wait prev. Tx. |
| flush commit marker |

T$_1$ returns

T$_2$ starts

| Tx. exec. TS=110 |
| flush logs |
| wait prev. Tx. |
| flush commit marker |

T$_2$ returns

time

## with SPHT [FAST'21]

T$_1$ starts

T$_0$ starts

| Tx. exec. TS=115 |
| flush logs |
| wait prev. Tx. |
| CAS |
| flush global marker |

T$_0$ returns

| Tx. exec. TS=105 |
| flush logs |
| wait prev. Tx. |
| wait global marker |

T$_1$ returns

T$_1$ detects that T$_0$ will flush the global marker

time

8

# STAMP



VACATION_LOW

# What about read-only transactions?

# Scalability issue #1



Legend:

htmBegin ●━━━━━━● htmCommit ● Flushing redo logs/
        Executing in HTM   commit marker to PM

# Scalability issue #2

- Hardware transactions have a (very) limited read set capacity
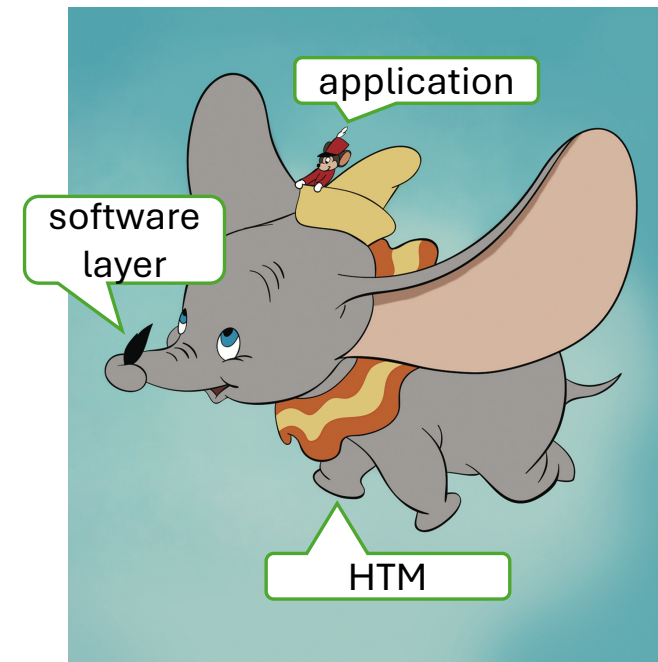- When a transaction exceeds its capacity, the HTM abort it
- Typically, the transaction must acquire a single global lock to execute

# Making read-only transactions scale with DUMBO

- Scalability issue #1:
Read-only transactions *practically* **never need to wait** for writer transactions to persist

- Scalability issue #2:
Read-only transactions **run HTM-free** (without read instrumentation), so have **unlimited reads and never abort**

*Durable Unlimited-reads Memory transactions on Best-effOrt HTM (DUMBO)*

**Update transactions**

BeginTX() →

`htmBegin(...)`

Execute
transaction

CommitTX() →

`htmCommit`

Ensure tx is durable
(e.g. using SPHT)

return

$T_R$  BeginTx()    r(X)=0                                    r(X)=1   CommitTx()

BeginTx()

CommitTx()

Consistency
violation

$T_W$                              w(X,1)

**htmBegin**                **htmCommit**

## Read-only transactions

BeginTX(RO) →

| State[tid] = activeRO |
|:---:|
| Execute transaction |
| State[tid] = inactiveRO |

CommitTX() →

## Update transactions

BeginTX() →

| **htmBegin(...)** |
|:---:|
| Execute transaction |

CommitTX() →

Wait until
State[i] == inactiveRO
for every other thread

**htmCommit**

isolation wait

Ensure tx is durable
(e.g. using SPHT)

| Existing HTMs | HTM can suspend access tracking? |
|---|---|
| ARM TME | no |
| Intel TSX | loads |
| PowerHTM | loads or full |

## Read-only transactions

BeginTX(RO)

State[tid] = activeRO

Execute
transaction

CommitTX()

State[tid] = inactiveRO

Not ready yet.
We need to fix the "The persistence bug"
(see Erez Petrank's previous talk)

## Update transactions

BeginTX()

**htmBegin(...)**

Execute
transaction
**htmSusTrack**

Wait until
State[i] == inactiveRO
for every other thread

CommitTX()

**htmResTrack**

**htmCommit**

isolation
wait

Ensure tx is durable
(e.g. using SPHT)

## Read-only transactions

BeginTX(RO) →

| State[tid] = activeRO |
| :---: |
| Execute transaction |
| State[tid] = inactiveRO |

CommitTX() →

**Wait until previous or concurrent update txs durable**

durability wait

## Update transactions

BeginTX() →

| `htmBegin(...)` |
| :---: |
| Execute transaction `htmSusTrack` |
| Wait until State[i] == inactiveRO for every other thread |
| `htmResTrack` `htmCommit` |

CommitTX() →

Ensure tx is durable (e.g. using SPHT)

isolation wait

# Revisiting the durability wait of read-only txs



Any concurrent writer will HTM-commit after R. Therefore, R no longer needs to wait for them!

R TS=11 TS=19 Durability wait

$W_1$ durTS=9 durTS=NA

$W_2$ durTS=7 durTS=15 durTS=NA

$W_3$ durTS=13 durTS=17 durTS=NA

$W_4$ durTS=9 durTS=24

$W_5$ durTS=13 durTS=NA

## Read-only transactions

BeginTX(RO) →  State[tid] = activeRO

Execute
transaction

CommitTX() →  State[tid] = inactiveRO

Wait until **previous**
update txs durable

durability
wait

## Update transactions

BeginTX() →  `htmBegin(...)`

Execute
transaction
`htmSusTrack`

CommitTX() →

Wait until
State[i] == inactiveRO
for every other thread

`htmResTrack`
`htmCommit`

isolation
wait

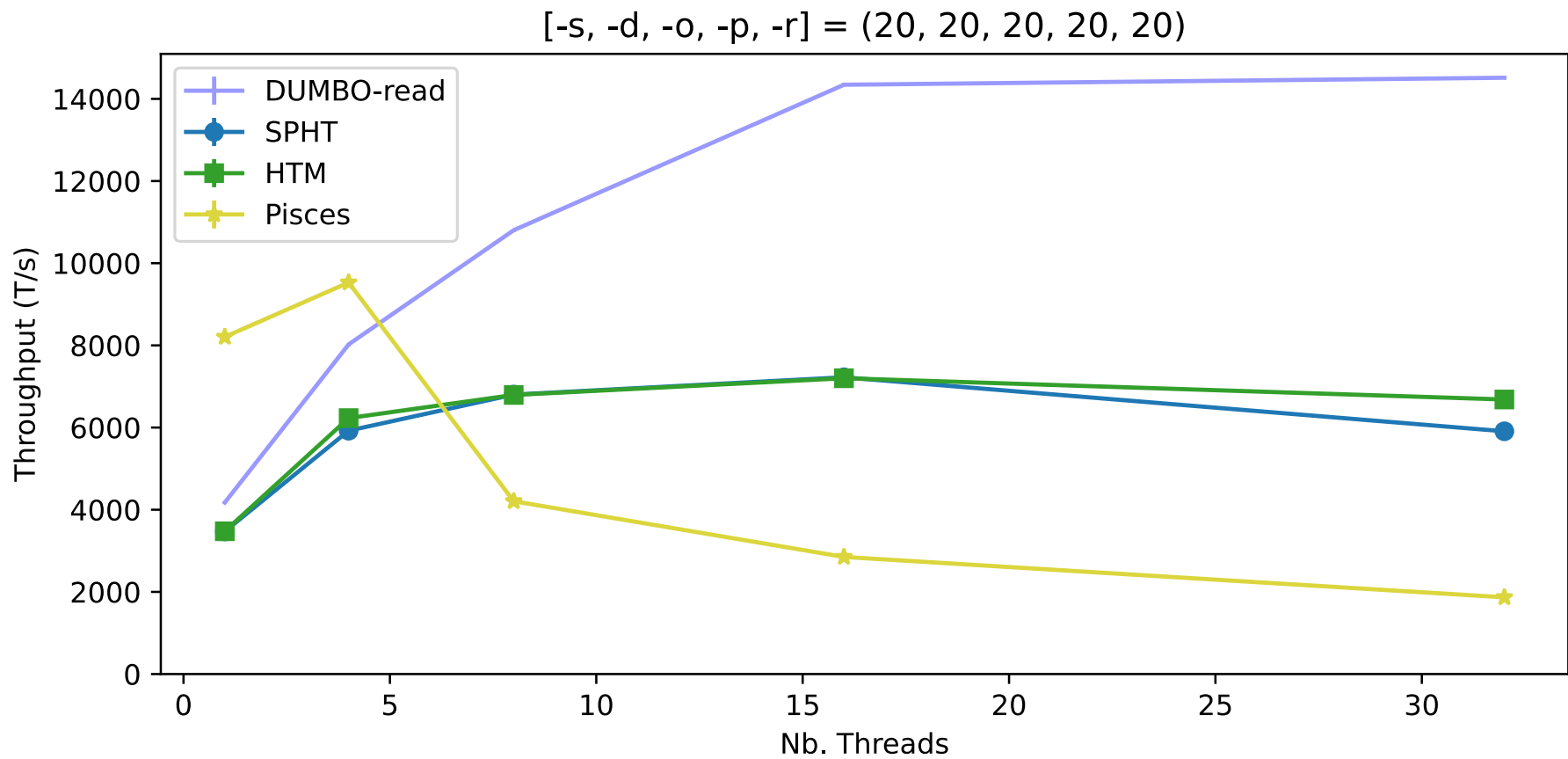Ensure tx is durable
(e.g. using SPHT)

# Evaluation

- IBM POWER9

- OS-assisted suspend/resume (HTM issues a trap) => high suspend/resume latency

- DUMBO, SPHT [FAST'21], HTM, Pisces (SI STM) [ATC'19]

- TPC-C benchmark with an even mix of operations:
  - "Stock level" and "Order status" (large read-only transactions)
  - "Payment" (small writer transactions)
  - "Delivery", "New order" (large writer transactions)

**Oregon State**
University

Machine provided by
OSU Open Source Lab

# Results



[-s, -d, -o, -p, -r] = (20, 20, 20, 20, 20)

# Summary

- We can have persistent hardware transactions with HTM+"software glue", however scalability is a huge challenge

- **SPHT** [FAST'21] accelerates writer transactions with a new commit logic that mitigates scalability bottlenecks of previous alternatives

- **DUMBO** [wip] boosts read-only transactions by granting them unlimited reads and a reduced durability wait

- Suspend/resume tracking support in HTM is useful, even if through expensive OS-assisted mechanisms

João Barreto | joao.barreto@tecnico.ulisboa.pt | https://www.dpss.inesc-id.pt/~jpbarreto/

# More DUMBO (not in this talk)

- On IBM POWER9, we can also suspend load+store tracking

- This also enables DUMBO to:
    - Hide redo log flush latency
    - Generlize unlimited reads to writer transactions
    - Improve durability&log replay logic