# Table of Contents

# Table of Contents

# What is Reachability Analysis?

- **Reachability Problem:** Given a graph $G = (V, E)$, initial states $I \subseteq V$ and goal states $F \subseteq V$, check if $F$ is *reachable* from $I$ via edges in $E$

- **In Model Checking:** Allows verification of most temporal safety properties *("something bad will never happen")*

- **Examples:**
  - Finding solutions in games (e.g. Chess, Sokoban, etc.)
  - Asserting mutual exclusion in parallel software
  - Asserting safety in traffic lights
  - etc.

# State Space Explosions

- **Reachability Problem:** Given a graph $G = (V, E)$, initial states $I \subseteq V$ and goal states $F \subseteq V$, check if $F$ is *reachable* from $I$ via edges in $E$

- $G$ is often *implicitly* described *(with a transition relation)*
- Set of reachable states is often determined *on-the-fly*
- Therefore, the size of $G$ is often initially unknown

- **State Space Explosions:** occur when $G$ does not longer fit into the available memory
    - Often happens in practice $\implies$ major limitation
    - State space Chess: $\sim 10^{43}$
    - Stars in universe: $\sim 10^{23}$

# Fighting State Space Explosions

- **Reduction Techniques:**
    - Partial Order Reduction *(exploit commutative transitions)*
    - Bisimulation Minimization *(merge "similar" states)*

- **Compression Techniques:**
    - Decision Diagrams *(e.g. BDDs, MDDs, LDDs, ZDDs)*
    - SAT-based approaches *(e.g. IC3)*

- **Adding Hardware Resources:**
    - More memory $\implies$ *larger* state spaces supported
    - More processors $\implies$ *faster* reachability analysis

# Parallel Symbolic Reachability

- **Symbolic Reachability:**
    - Represent the state space as a BDD
    - Represent initial states and the transition relation as BDDs
    - Perform reachability via BDD operations

- **Parallel Reachability:** Using a many-core cluster with a large amount of memory to perform reachability
    - Sylvan reaches speedups up to 38 with 48 cores

- **Disadvantages:**
    - Upgrading is *expensive*
    - Upgrading is *limited*

# Distributed Symbolic Reachability

- **Distributed Reachability:** Using a network of workstations, connected via a high-performance network.

- **Compared to a Many-core Cluster:**
  - *Cheaper* scalability
  - *Unlimited* scalability

- **Challenges:**
  - Only small amounts of computation per memory access
  - Many *remote* memory accesses required
  - Network latency *easily* becomes a bottleneck

- **Achievements:** Very large state spaces are supported, but no speedups are obtained...

# Improving Distributed Symbolic Reachability

- **Zhao et al (2009):** Most important design considerations for improvements are:
    - Data-distribution
    - Load-balancing maintenance
    - Reducing communication overhead
    - Exploting data-locality *(suggested by Chung et al)*

- **Contribution:** Employing *modern* techniques to implement these design considerations.

- **Reducing Communicational Overhead:** Infiniband and RDMA

- **Load-balancing:** Work stealing *(due to the success of Sylvan and Lace)*

- **Exploting data-locality:** *Hierarchical* work stealing

- **Data Distribution:** RDMA-based distributed hash table

# Research Questions

- **Main Question:** How efficient can RDMA-based distributed implementations of BDD operations scale along all processing units and available memory connected via a high-performance network?

- **Subquestions:**
  1. How can the storage and retrieval of data efficiently be managed to minimize their latencies?
  2. How can the total computational work be divided and distributed to maximize scalability along processors over a high-performance network?
  3. How can the idle-times of processes be minimized while performing network communication?

# Research Method

- **Project:** We split the project into three parts:
  1. Storing states *(Distributed hash table)*
  2. Load-balancing mechanisms *(Hierarchical work stealing)*
  3. Distributed BDD operations

- All parts have separately been designed, implemented, and experimentally evaluated
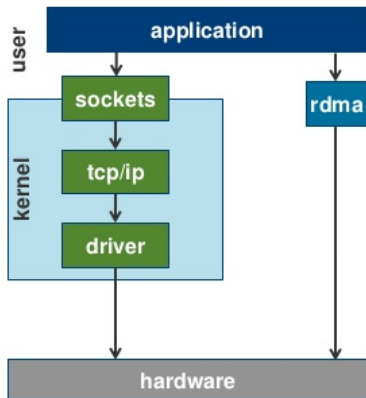
# Table of Contents

# The Infiniband Architecture

- **Infiniband:** Specialized hardware used to construct high-performance networks

- **Advantages:**
  1. Comparable in price to standard Ethernet hardware
  2. Supports up to 100 GB/s
  3. NICs can *directly* access main-memory via PCI-E bus
  4. End-to-end latencies of $1\mu s$ have been measured *(according to the IB website)*
  5. Supports RDMA

- **UTwente:** 10 Dell M610 machines, connected via a 20 GB/s Infiniband network

# Remote Direct Memory Access

- **RDMA:** Directly access memory of a remote machine, without invoking its CPU
  - one-sided vs. two-sided RDMA

- **Advantages:**
  1. Zero-copy
  2. Kernel bypassing
  3. CPU efficiency

- **Roundtrip Latency:** Within $3\mu s$ in Infiniband hardware, compared to $60\mu s$ with TCP on Ethernet hardware
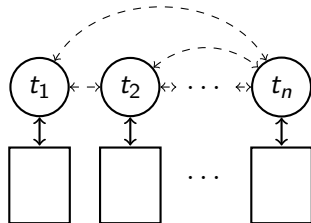
# Kernel Bypassing



(VMWorld 2013 - How Latency Destroys Performance... And What to Do About It)

# Parallel Programming Models

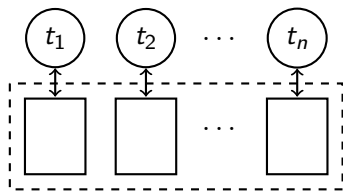

**Shared Memory**
- Single addr. space
- e.g. NUMA, SMP

**Distributed Memory**
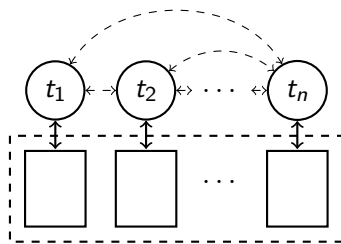- Only local memories
- Communication via Message Passing
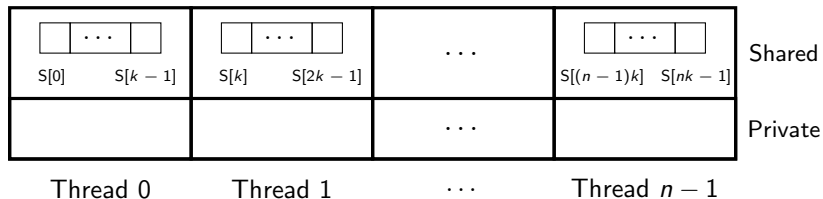
# Partitioned Global Address Space



**PGAS**
- Shared + Distributed
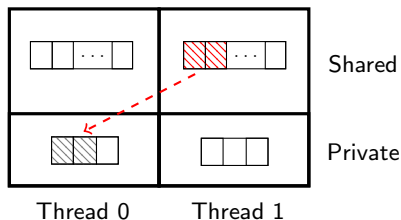- Data locality exploited

**Hybrid PGAS**
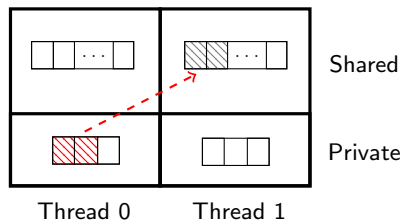- PGAS + message passing

# PGAS: Memory Model

# PGAS: `memget` and `memput`



memget($P, S$)
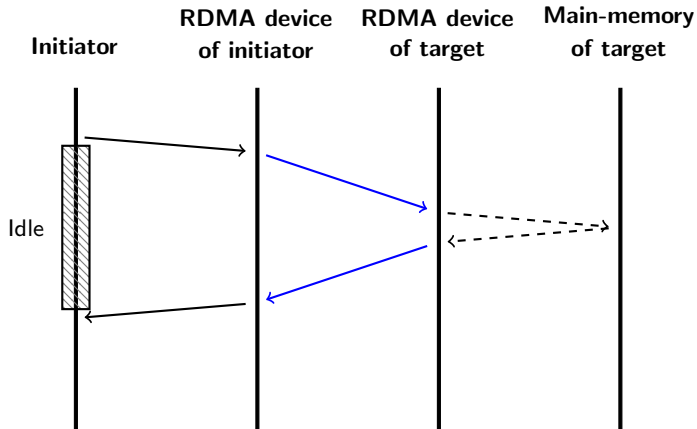
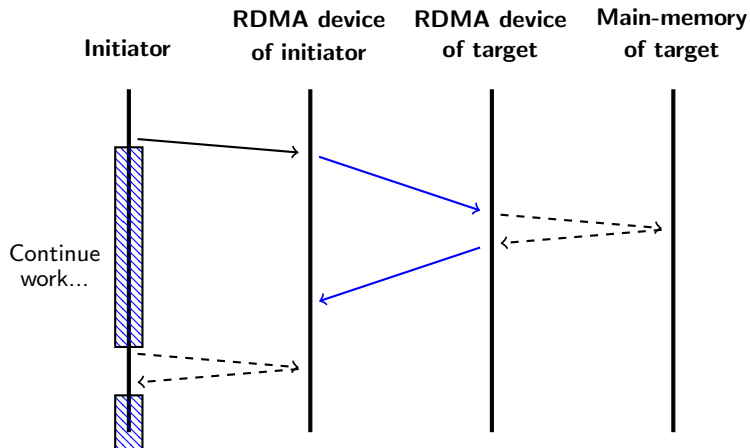- Copies block of *shared* memory $S$ into *private* memory $P$

memput($S, P$)

- Copies block of *private* $P$ memory into *shared* memory $S$

# PGAS: Synchronous Operations

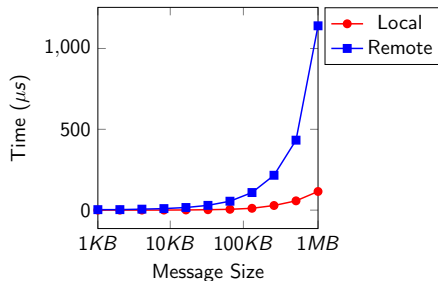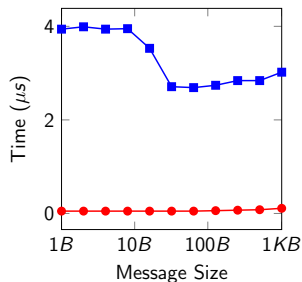# PGAS: Asynchronous Operations

# PGAS: Implementations

- **Many Implementations:**
  - Berkeley UPC
  - OpenSHMEM
  - Co-array Fortran
  - Titanium
  - X10
  - Chapel
  - etc.

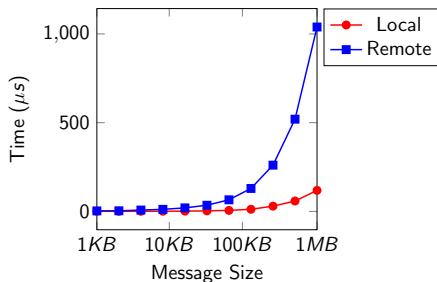- We chose UPC because it supports async operations

# PGAS: Latency of `memget`



- **For messages $\leq 16$ bytes:**
    - Local 76 times faster than remote
    - Local: $50ns$ on average
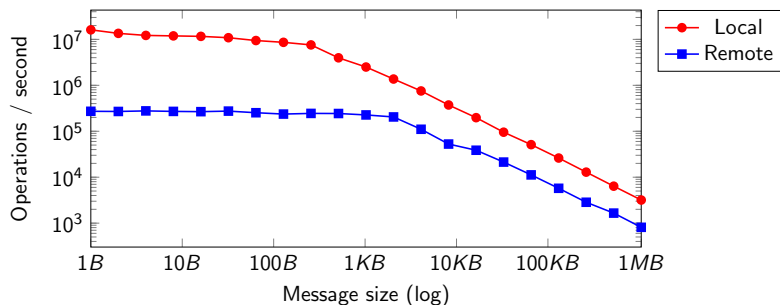    - Remote: $3.87\mu s$ on average

# PGAS: Latency of `memput`



- **For messages $\leq 16$ bytes:**
  - Local only 6 times faster than remote
  - Local: $0.44\mu s$ on average
  - Remote: $2.68\mu s$ on average

# PGAS: Throughput of `memput`



- **For messages $\leq 16$ bytes:**
  - Local throughput 48 times higher than remote throughput

# Table of Contents

# Hash Table: Notations

## Hash Table

$T = \langle b_0, \ldots, b_{n-1} \rangle$ as a sequence of buckets $b_i$, where:

- Hash table size: $n$, number of inserted elements: $m$
- Load factor: $\alpha = \frac{m}{n}$

## Hash function

$h : U \to R$, with:

- Range of keys: $R = \{0, \ldots, r-1\}$
- Universe: $\{0,1\}^w$ *(of all w-sized binary words)*
- Mapping words $x \in U$ to buckets $b_{h(x)}$ by letting $r < n$

## Notation

- For $x \in U$, we write $x \in b_i$ if bucket $b_i$ contains $x$
- For $x \in U$, we write $x \in T$ if $x \in b_i$ for some $0 \leq i \leq n-1$

# Distributed Hash Table: Requirements

## Requirements

1. Minimal number of roundtrips
2. Minimal memory overhead
3. CPU efficient (no polling)
4. Should support `find-or-put`
5. Should support PGAS

## `find-or-put`($d$)

- If $d \in T$, return `found`
- If $d \notin T$, insert $d$ and return `inserted`
- If $d \notin T$ and $d$ cannot be inserted, return `full`

# Resolving Collisions

## Hash collision

Occurs when $h(x) = h(y)$ for $x, y \in U$ with $x \neq y$

- Hashing strategy determines the number of roundtrips required by `find-or-put`

## Existing work (RDMA-based key/value stores)

- Pilaf *(Cuckoo hashing)*
- Nessie *(Cuckoo hashing)*
- FaRM *(Hopscotch hashing)*
- HERD *(high throughput, but CPU inefficient)*

We investigated hashing strategies and determined their performance

# Chained Hashing

## Chained Hashing

- Every bucket is a linked list
- Inserting $x \in U$ performed by adding it to $b_{h(x)}$
- Finding $x \in U$ performed by traversing $b_{h(x)}$

## Complexity of `find-or-put(d)`

- $\Theta(m)$ in worst case *(when all m elements are in $b_{h(d)}$)*
- $\Theta(1 + \alpha)$ on average when a *universal hash function* is used

## Universal hash function

$h : U \to R$ is called *universal* if $Pr[h(x) = h(y)] \leq \frac{1}{|U|}$ for every $x, y \in U$

- Good theoretical properties
- In practice: "cheaper" functions are often used

# Cuckoo Hashing

## Cuckoo Hashing

- Uses $k \geq 2$ *independent* hash functions $h_1, \ldots, h_k : U \to R$ with $h_i \neq h_j$ for every $i \neq j$ *(k-way Cuckoo hashing)*
- Nessie: 2-way Cuckoo hashing
- Pilaf: 3-way Cuckoo hashing

## Cuckoo Invariant

For every element $x \in U$ it holds that either $x \notin T$ or $x \in b_{h_i(x)}$ for exactly one $1 \leq i \leq k$

## Complexity

- Lookups require $k$ roundtrips
- Inserts may require many when all $k$ buckets are occupied

# Bucketized Cuckoo Hashing

## Bucketized Cuckoo Hashing

- Every bucket $b_i$ is subdivided into $l$ slots
- Every slot may contain an element from $U$
- Denoted by $(k, l)$-Cuckoo hashing

## Bucketized Cuckoo Hashing

- Same as Cuckoo hashing, but linearly reduced by $l$
- Efficient even when $\alpha > 0.9$ *(Andersen et al, 2013)*
- Pilaf: $(2, 4)$-Cuckoo hashing could be very effective

# Hopscotch Hashing

## Hopscotch Hashing

- Each bucket $b_i$ has a fixed-sized neighborhood $N(b_i)$ of constant size $H \geq 1$
- $N(b_i) = \langle b_i, \ldots, b_j \rangle$ with $j = (i + H - 1) \mod n$
- $N(b_i)$ thus contains $b_i$ itself and the next $H - 1$ buckets *(modulo n)*
- Neighborhoods are thus consecutive in memory

## Hopscotch Invariant

Let $x \in U$ and $N(b_{h(x)}) = \langle b_1, \ldots, b_H \rangle$. Then either $x \notin T$ or $x \in b_i$ for exactly one $1 \leq i \leq H$

## Complexity

- `find-or-put(d)` may obtain $N(b_{h(d)})$ in 1 roundtrip
- Inserts may require *many* more when $N(b_{h(d)})$ is full

# Linear Probing

## Linear Probing

- For $x \in U$, it examines $b_{h(x)+0}, b_{h(x)+1}, \ldots, b_{h(x)+t}$ (modulo $n$) with threshold $t > 0$
- Buckets are consecutive in memory
- Therefore, cache-line efficient

## Complexity

- Same as Hopscotch, but without relocation schemes
- Hopscotch invariant not maintained, lookups are more expensive
- But inserts are arguable cheaper *(amortized complexity)*

# Complexity of Linear Probing

### Theorem: Examining buckets *(Knuth, 1997)*

Assuming that a *universal hash function* is used, the *expected* number of buckets to examine until an empty bucket is found is at most:

$$\frac{1}{2}\Big(1 + \frac{1}{(1-\alpha)^2}\Big)$$

### Chunk retrieval

- Similar to Hopscotch, a fixed-sized range of buckets can be obtained with a *single* roundtrip, which we refer to as *chunks*
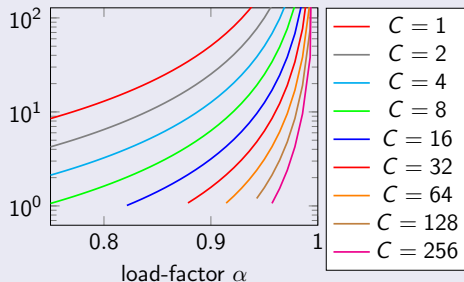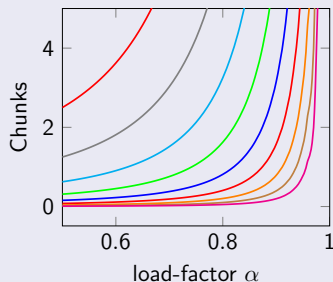- We denote the *chunk size* by $C \geq 1$

# Complexity of Chunk Retrievals

## Corollary: Number of chunks

The expected number of chunks to be inspected is at most:

$$\frac{1}{2C}\left(1 + \frac{1}{(1-\alpha)^2}\right)$$
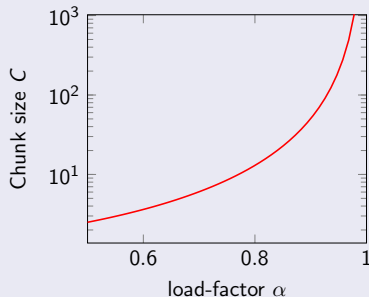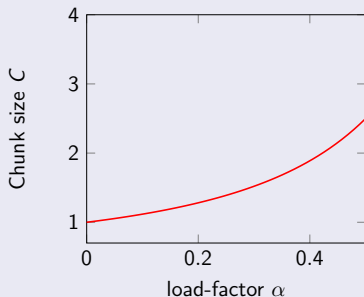
## Number of chunks to read

# Bounding Efficiency

## Theorem: Efficiency bound

A chunk of size $C \geq 1$ is expected to contain an empty bucket if:

$$\alpha \leq 1 - \sqrt{\frac{1}{2C - 1}}$$

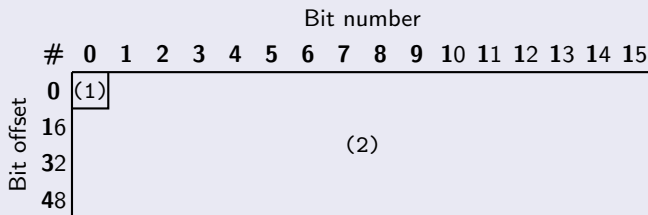## Expected load-factor at which chunk is full

# Designing `find-or-put`

## Memory layout

- Shared array $B[0], \ldots, B[kn-1]$ of buckets, so that each thread *owns* $k$ buckets
- 2D array $P[0][0], \ldots, P[M-1][C-1]$ per thread in *private* memory

## Bucket layout



Bit number

| # | **0** | **1** | **2** | **3** | **4** | **5** | **6** | **7** | **8** | **9** | 10 | 11 | 12 | 13 | 14 | 15 |

Bit offset: **0**, **16**, **32**, **48**

(1) is at bit offset 0

(2) fills the rest

(1) is a locking bit (1 bit)

(2) contains data (63 bits)

# Cache Efficiency

## Cache lines

- Typically 64 bytes in size
- So 8 buckets per cache line
- Therefore, we choose $C$ to be a multiple of 8

## Cache line alignment

- The arrays P are cache line aligned
- The array B is *not*, since it is shared (could not find support from UPC to align shared memory)
- But the IB verbs libary *has* support for shared memory alignment...

# Chunk Retrieval

## Asynchronous chunk retrievals

- Before iterating over a chunk, request the *next* consecutive chunk
- Done to overlap roundtrips with actual work (interleaving queries)

## The query-chunk($i, p$) operation

- Transfers the $i$th chunk, starting from $b_p$, from B into $P[i][0], \ldots, P[i][C-1]$ *asynchronously*
- Returns a *handle r* for synchronization
- Synchronization can be performed by calling sync($r$)

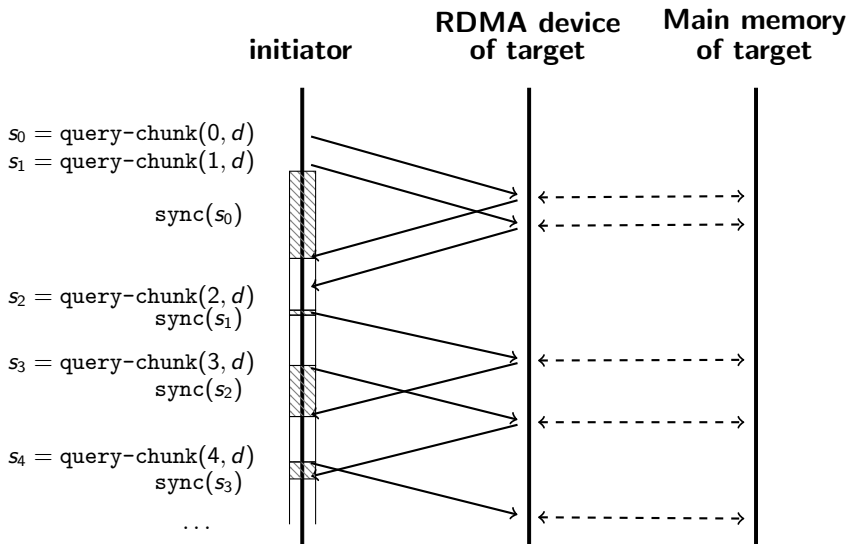# Design of `find-or-put`

```
1  def find-or-put(data):
2      h ← hash(data)
3      s₀ ← query-chunk(0, h)
4      for i ← 0 to M − 1:
5          if i < M − 1
6              s_{i+1} ← query-chunk(i + 1, h)
7          sync(s_i)
8          for j ← 0 to C − 1:
9              if (P[i][j] & OCCUPIED) = 0
10                 ▷ Try to claim bucket B[a]
11                 a ← (h + iC + j) mod kn
12                 d ← data(data) | OCCUPIED
13                 val ← cas(B[a], P[i][j], d)
14                 if val = P[i][j]
15                     return inserted
16                 elif data(val) = data
17                     return found
18             elif data(P[i][j]) = data
19                 return found
20     return full
```

```
1  def query-chunk(i, h):
2      ▷ Find start and end index
3      start ← (h + iC) mod kn
4      end ← (h + (i + 1)C − 1) mod kn
5      if end < start
6          return split(start, end)
7      else
8          S ← ⟨B[start], . . . , B[end]⟩
9          P ← ⟨P[i][0], . . . , P[i][C − 1]⟩
10         return memget-async(S, P)
```

```
1  def split(start, end):
2      ▷ Find the blocks in shared memory
3      S₁ ← ⟨B[start], . . . , B[kn − 1]⟩
4      S₂ ← ⟨B[0], . . . , B[end]⟩
5      ▷ Corresp. blocks in private memory
6      P₁ ← ⟨P[i][0], . . . , P[i][|S₁| − 1]⟩
7      P₂ ← ⟨P[i][|S₁|], . . . , P[i][C − 1]⟩
8      ▷ Retrieve the chunk
9      s₁ ← memget-async(S₁, P₁)
10     s₂ ← memget-async(S₂, P₂)
11     return ⟨s₁, s₂⟩
```

# Asynchronous Query Retrievals

# Experimental Evaluation
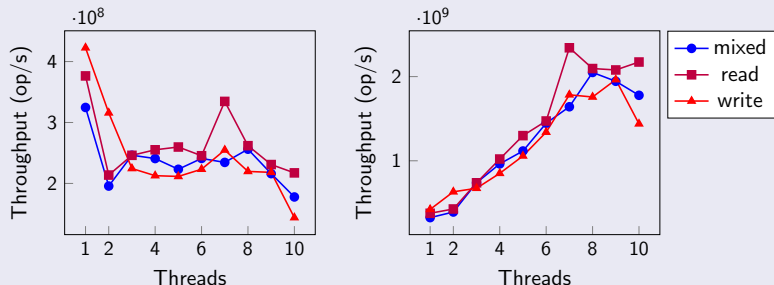
## Experimental setup (m610 partition)

- 10 Dell M610 machines
- 8 GPU cores and 24 GB main-memory (each)
- Ubuntu 14.04.2 LTS, kernel version 3.13.0
- 20 GB/s Infiniband network

## Benchmarks

- Throughput of `find-or-put`
- Latency of `find-or-put`
- Roundtrips required by `find-or-put`
- Under different workloads: mixed, read-intensive, and write-intensive
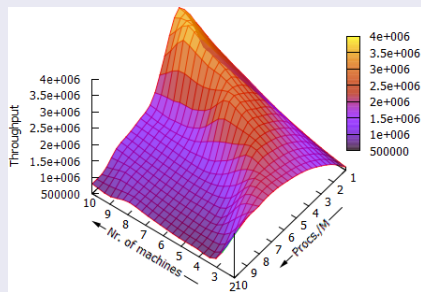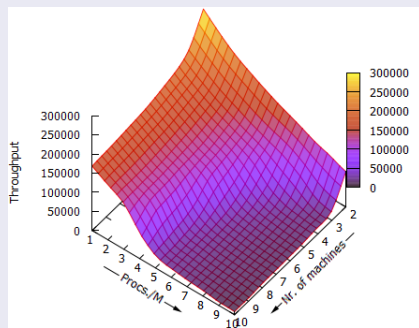
# Local Throughput of `find-or-put`

## TP per thread (left) and total TP (right)



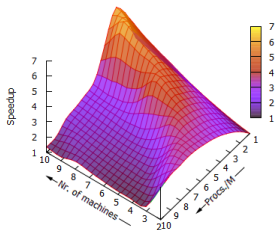| Workload | Base Througput | | Best Throughput | | Speedup |
|---|---|---|---|---|---|
| | Throughput | Procs. | Througput | Procs. | |
| Mixed | 324,676,333 | 1 | 2,049,388,900 | 8 | 6.31 |
| Read-intensive | 376,434,000 | 1 | 2,342,278,333 | 7 | 6.22 |
| Write-intensive | 422,593,000 | 1 | 1,963,570,267 | 9 | 4.65 |

# Remote Throughput of `find-or-put`
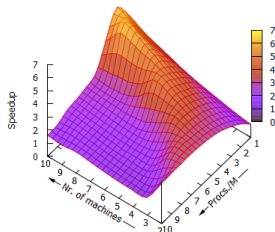
## Mixed workload: TP per thread (left) and total TP (right)



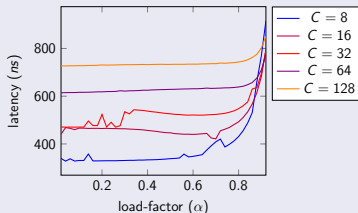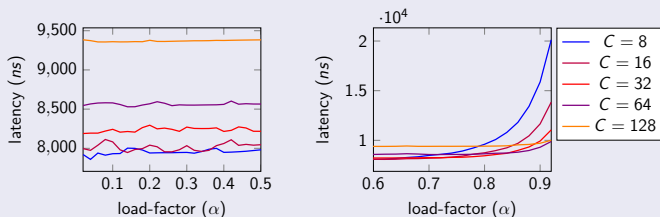| Workload | Base Throuput | | | Best Throughput | | | Speedup |
|----------|------|------|----------|------|------|----------|---------|
| | TP. | M. | Procs./M. | TP. | M. | Procs./M. | |
| Mixed | 592,929 | 2 | 1 | 3,607,003 | 10 | 3 | 6.08 |
| Read | 742,728 | 2 | 1 | 4,620,752 | 10 | 3 | 6.22 |
| Write | 495,370 | 2 | 1 | 2,999,234 | 10 | 3 | 6.05 |

Mixed



Read-intensive



Write-intensive

# Latency of `find-or-put`

## Local latency
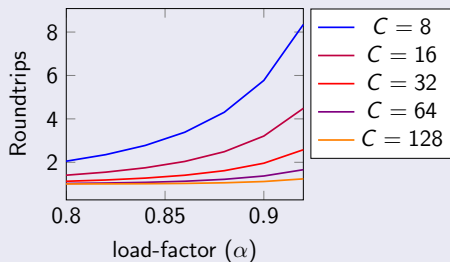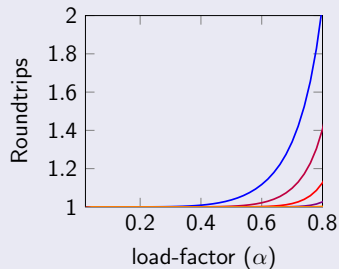


## Remote latency

## Number of roundtrips

# Conclusions

## General

- Minimizing roundtrips *critical* for increased performance
- Overlapping queries reduces waiting-times and increases latency
- Linear probing requires less roundtrips than Hopscotch and Cuckoo

## Performance

- `find-or-put` takes $9.3\mu s$ on average with $\alpha = 0.9$ and $C = 64$
- Peak-throughput of $3.6 \times 10^6$ op/s obtained with 10 machines

## Future work

- Use adaptive chunk sizes (based on efficiency bounds Theorem)
- In addition, update asynchronous queries to prevent unused retrievals

# Table of Contents

# Load Balancing

## Task-based parallelism

- Dividing computational problems into smaller *tasks*
- Task is a basic unit of work and only depend on intermediate *subtasks*
- All threads maintain *task pools*

## Load-balancing tasks

- Ideally tasks are perfectly distributed *(infeasible)*
- Instead: mapping tasks dynamically to threads

## Task granularity

The relation between the computational workload and the amount of communication required between threads

- Fine-grained: large number of small tasks
- coarse-grained: small number of large tasks

# Work Stealing and Sharing

## Work stealing

- Efficient technique for fine-grained task parallelism
- Threads are either idle or working
- When idle, threads *steal* from remote task pools
- Stealing thread is *thief*, targetted thread is *victim*
- Termination when all threads are idle

## Work sharing

- Threads communicate their status
- When idle, other threads *share* work
- Communication of work stealing is *more* efficient (Blumofe, 1999)

# Work Stealing Operations

## Operations

- **spawn**: push new task to task pool
- **call**: execute given task
- **sync**: pull task from pool and execute

## Fibonacci example

```
1  int fib(n):
2      if n < 2 return n
3      a ← fib(n − 1)
4      b ← fib(n − 2)
5      return a + b
```
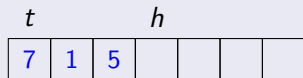
```
1  int par-fib(n):
2      if n < 2 return n
3      spawn(par-fib, n − 1)
4      r ← call(par-fib, n − 2)
5      return r + sync
```

# Implementing the Task Pool

## Double ended queue (deque)

- Similar to queue, but has two ends: *head* and *tail*
- Items can be pushed or popped from both ends
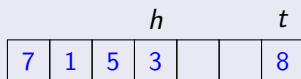- Implemented as a fixed-sized array

## Example



| $t$ | | $h$ | | | | |
|---|---|---|---|---|---|---|
| 7 | 1 | 5 | | | | |

Initial

| $t$ | | | $h$ | | | |
|---|---|---|---|---|---|---|
| 7 | 1 | 5 | 3 | | | |

$\mathtt{push}(3, h)$

| | | | $h$ | | $t$ | |
|---|---|---|---|---|---|---|
| 7 | 1 | 5 | 3 | | | 8 |

$\mathtt{push}(8, t)$

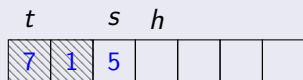| | | | $h$ | | $t$ | |
|---|---|---|---|---|---|---|
| 7 | 1 | 5 | 3 | | | 8 |

$\mathtt{pop}(h)$

# Implementing the Task Pool

## Split deque

- Deque with a *split point s*
- *s* determines what sections belongs to *head* and *tail*
- Used to denote a *public* and *private* region
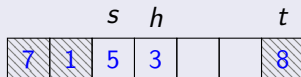- *s* can be relocated to increase/decrease public region

## Example



Initial

push(8, t)

push(3, h)

pop(h)

# Implementing the Task Pool

## Performance of split deques

- Modifying $s$ may conflict with steal operations
- Either locks or memory fences required
- Expensive in distributed environment!

## Existing work (current state-of-the-art)

- HotSLAW: access to public region requires locking
- Scioto: *whole* split deque locked when stealing
- Lace: non-blocking, but shrinking public region requires memory fence

# Implementing the Task Pool

## Private deques

- Implemented as a stack
- Do not have a public region (completely private)

## Private deque work stealing

- When stealing, idle workers *explicitly* ask for work
- **Advantage:** No locking required
- **Disadvantage:** Requires participation from both victim and thief

# Victim Selection Protocols

## Selecting victims

- Random victim selection
- Hierarchical victim selection (Min et al, 2011)
- Leapfrogging

## Contribution and motivation
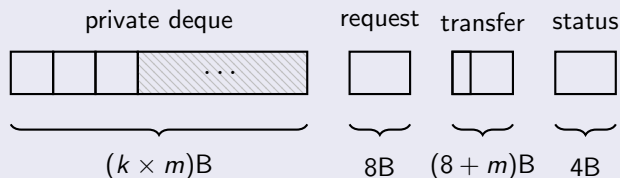
Private-deque work stealing operations:

- Minimal number of roundtrips
- Uses all three victim selection protocols
- Similar approach by Olivier et al. (2008), but requires *more* roundtrips and does *not* exploit network hierarcy

# Designing Private-Deque Work Stealing

## Memory layout

- Shared 2D-array: deque[0][0], . . . , deque[THREADS − 1][k − 1]
- Request cells: request[0], . . . , deque[THREADS − 1]
- Transfer cells: transfer[0], . . . , transfer[THREADS − 1]
- Status cells: status[0], . . . , status[THREADS − 1]

## Schematically



private deque        request   transfer   status

$(k \times m)$B      8B    $(8 + m)$B    4B

Where $m$ is the task size (in Bytes) and $k$ the deque size

# Performing Steals

## Request cell

- Contains either `BLOCKING`, `EMPTY`, or a thread ID:
- **blocking**: no tasks can be stolen
- **empty**: no pending steal requests
- **identifier**: pending steal request

## Transfer cell

- Contains either `EMPTY` or a task + location:
- **empty**: no task received
- **task**: task received + corresponding location in deque

## Status cell

- Contains either `IDLE` or `WORKING`

# Hierarchical Work Stealing

## Domain levels

- Berkeley UPC provides thread-distance($i, j$) function
- Which returns: verynear, near, far, or veryfar
- We use an array *domain*, so that *domain*[$i$] contain all thread IDs on the $i$th level
- We use a shuffle function that randomly *shuffles* a domain level.

## Hierarchical stealing

1. Threads start by performing leapfrogging
2. Threads perform count(*domain*[$i$]) steal attempts before moving to level $i + 1$
3. If all levels have been tried, perform termination detection

# Designing **spawn**, **call**, and **sync**

```
1  def sync():
2     task ← deque[MY-ID][head − 1]
3     if task.stolen
4        communicate()
5        while ¬task.completed:
6           ▷ Perform leapfrogging
7           if steal(task.owner) continue
8           if task.completed break
9           ▷ Perform hierarchical stealing
10          for i ← 0 to HIERARCHY-LVLS - 1:
11             shuffle(domain[i])
12             foreach victim ∈ domain[i] do
13                if steal(victim) goto line 5
14                if task.completed goto line 16
15       ▷ Return result from stolen task
16       head ← head − 1
17       tail ← tail − 1
18       return task.result
19    else
20       head ← head − 1
21       return call(task)
```

```
1  def spawn(desc, params):
2     ▷ Build a new task
3     task ← deque[MY-ID][head]
4     task.desc ← desc
5     task.stolen ← false
6     task.completed ← false
7     task.params ← params
8     ▷ Write new task to deque
9     deque[MY-ID][head] ← task
10    head ← head + 1
```

```
1  def call(desc, params):
2     communicate()
3     ▷ Find the intended function
4     func ← function-of(desc)
5     ▷ Invoke that function
6     return func(params)
```

# Initiating a Computation

```
1  def initiate(desc, params):
2      ▷ Wait for all workers to start
3      status[MY-ID] ← WORKING
4      barrier()
5      ▷ Perform task
6      result ← call(desc, params)
7      ▷ Wait for all workers to complete
8      status[MY-ID] ← IDLE
9      barrier()
10     return result

1  def compute(desc, params):
2      if MY-ID = 0
3          initiate(desc, params)
4      else
5          participate()
```

```
1  def participate():
2      ▷ Wait for all workers to start
3      status[MY-ID] ← IDLE
4      barrier()
5      ▷ Perform hierarchical stealing
6      while true:
7          status[MY-ID] ← IDLE
8          for i ← 0 to HIERARCHY-LVLS - 1:
9              shuffle(domain[i])
10             foreach victim ∈ domain[i] do
11                 if steal(victim) goto line 5

12         ▷ No worker had tasks to steal..
13         if termination-detection()
14             break
15     barrier()
```

# Communicate Work

```
1  def communicate():
2     if head − tail < 2
3        if request[MY-ID] ≠ BLOCKED
4           ▷ Not enough stealable tasks, block further requests
5           if request[MY-ID] ≠ EMPTY
6              reject-and-block()
7           elif cas(request[MY-ID], EMPTY, BLOCKED) ≠ EMPTY
8              reject-and-block()
9     elif request[MY-ID] = BLOCKED
10       request[MY-ID] ← EMPTY
11    elif request[MY-ID] ≠ EMPTY
12       thief ← request[MY-ID]
13       request[MY-ID] ← EMPTY
14       ▷ Prepare task to be stolen
15       deque[MY-ID][tail].stolen ← true
16       deque[MY-ID][tail].owner ← thief
17       ▷ Construct the transfer message
18       msg ← new TransferMessage
19       msg.index ← tail
20       msg.task ← deque[MY-ID][tail]
21       memput-async(transfer[thief], msg)
```

# Stealing and Termination Detection

```
1  def steal(victim):
2      communicate()
3      transfer[MY-ID] ← EMPTY
4      res ← cas(request[victim], EMPTY, MY-ID)
5      if res = EMPTY
6          ▷ Wait for response from victim
7          while transfer[MY-ID] = EMPTY:
8              communicate()
9          if transfer[MY-ID] = EMPTY
10             return false
11         else
12             status[MY-ID] ← WORKING
13             i ← transfer[MY-ID].index
14             task ← transfer[MY-ID].task
15             task.result ← call(task)
16             ▷ Write back the task result
17             task.completed ← true
18             memput-async(deque[victim][i], task)
19             return true
20     return false
```

```
1  def reject-and-block():
2      ▷ Block further requests
3      thief ← request[MY-ID]
4      request[MY-ID] ← BLOCKED
5      ▷ Send a negative response
6      msg ← new TransferMessage
7      msg.index ← 0
8      msg.task ← EMPTY
9      memput-async(transfer[thief], msg)

1  def termination-detection():
2      ▷ Send status requests
3      for i ← 0 to THREADS − 1:
4          s_i ← memget-async(
5              res[i], status[i])
6      ▷ Wait for status responses
7      for i ← 0 to THREADS − 1:
8          sync(s_i)
9          if res[i] = WORKING
10             return false
11     return true
```

# Experimental Evaluation

## Benchmarks

- Performed a number of microbenchmarks
- Determined speedup when scaling along machines and threads per machine
- Compared speedup with HotSLAW

| Benchmark | Nr. of Tasks | Avg. Task Time | Input/Output Size | Input/Output Size Hotslaw |
|-----------|--------------|----------------|-------------------|---------------------------|
| fib(45) | 3,672,623,805 | 0.154 $\mu s$ | 16/8 bytes | 4/8 bytes |
| nqueens(15) | 171,127,071 | 4.14 $\mu s$ | 20/8 bytes | 28/8 bytes |
| uts(T2L) | 96,793,510 | 0.986 $\mu s$ | 20/8 bytes | 32/0 bytes |
| uts(T3L) | 111,345,631 | 0.722 $\mu s$ | 20/8 bytes | 32/0 bytes |
| matmul(512) | 32,767 | 188.30 $\mu s$ | 20/8 bytes | 28/0 bytes |

# Speedup Graphs



fib(45)

matmul(512)

uts(T3L)

nqueens(15)

uts(T2L)

# Computation Times

## Sequential- and best times

| Benchmark | Sequential Time | Best Configuration | | | Speedup |
|-----------|-----------------|------|----------|----------|---------|
|           |                 | Time | Machines | Procs./M. |         |
| `fib(45)`     | 563.87 | 8.85  | 10 | 8 | 63.69 |
| `matmul(512)` | 6.17   | 1.07  | 1  | 9 | 5.76  |
| `uts(T2L)`    | 90.48  | 1.90  | 10 | 8 | 47.62 |
| `uts(T3L)`    | 73.60  | 3.48  | 10 | 5 | 21.15 |
| `nqueens(15)` | 707.64 | 10.29 | 10 | 8 | 68.74 |

## Comparison with HotSLAW

| Benchmark | Our Implementation | | HotSLAW | |
|-----------|--------------------|--|---------|--|
|           | Seq. Time | Best Time | Seq. Time | Best Time |
| `fib(45)`      | 563.87 | 8.85  | 938.49 | 13.13 |
| `nqueens(15)`  | 707.64 | 10.29 | 387.53 | 5.68  |
| `uts(T2L)`     | 90.48  | 1.90  | 81.22  | 1.53  |
| `uts(T3L)`     | 73.60  | 3.48  | 67.64  | 5.48  |
| `uts(T3L)`*    | 73.60  | 3.48  | 51.69  | 1.32  |

# Comparison with HotSLAW (Computation Time)



fib(45)

nqueens(15)
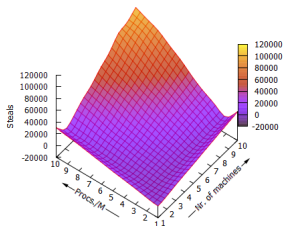
uts(T3L)

uts(T3L) with steal-10
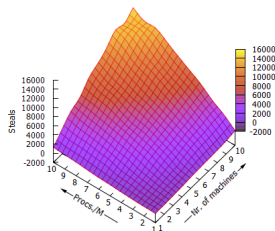
uts(T2L)

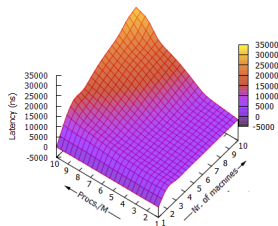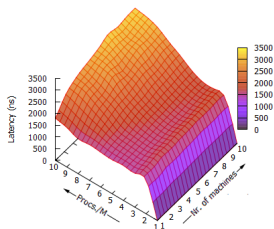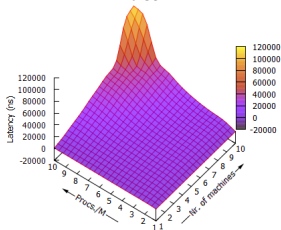# Average Number of Steals



fib(45)



nqueens(15)



uts(T3L)


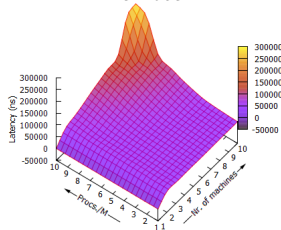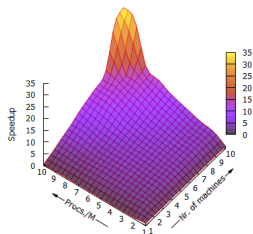
uts(T3L) with steal-10



uts(T2L)

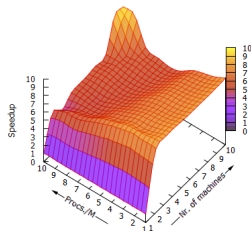# Latency of `steal`



local



remote



HotSLAW - local



HotSLAW - remote

# Speedup of `steal` (vs HotSLAW)



local

remote