28th August 2017

# HELP YOUR BUSY NEIGHBOURS
## DYNAMIC MULTICASTS
## OVER STATIC TOPOLOGIES

**Robert Kuban**, Randolf Rotta, Jörg Nolte

Distributed Systems / Operationg Systems

# OUR TARGET SCENARIO

**objective:** scalable multicasts
+ acknowledgement of completion
+ dynamic group membership (join/leave)

**applications:** cache invalidation, esp. TLB shootdown

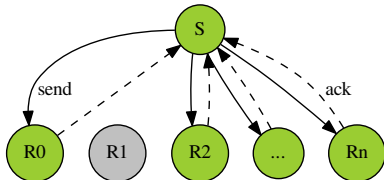**hardware:** many-cores like Intel XeonPhi, Tilera TilePro. . .
+ cache-coherent shared memory
+ point-to-point message passing

# EXAMPLE: LINUX TLB SHOOTDOWN
**Linux 4.11 x86 `smp_call_function_many()`**

## Initiator (Sender)

1. update page tables
2. enqueue invalidation tasklet at each thread
3. send IPI to each thread
4. wait on flag in each tasklet

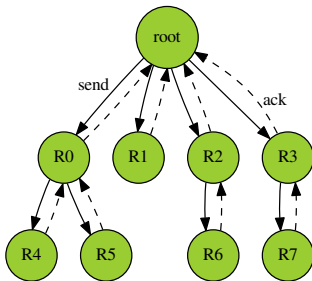## Other CPU Threads
IPI handler processes tasklet:

1. invalidate page(s) in TLB
2. set ACK flag in tasklet

$\Rightarrow$ flat topology

👍 fast join/leave via bit-mask

👎 $\mathcal{O}(n)$ latency

b·tu Brandenburg
University of Technology
Cottbus - Senftenberg



- propagate along a tree topology
- use constraint solver
  for optimized topology
- proposed for TLB shootdowns[1]

👎 expensive join/leave
   or interrupt ex-members

👍 $\mathscr{O}(\log n)$ latency

---

[1]Baumann et al., *The multikernel: A new OS architecture for scalable multicore systems*, 2009

# DESIGN SPACE

**b-tu** Brandenburg
University of Technology
Cottbus - Senftenberg

|  | Multicasts<br>(just members) | Broadcasts<br>(over all threads) |
|---|---|---|
| **Flat** | 👍 low latency<br>for small groups<br>👎 high latency<br>for large groups<br>👍 fast join/leave | 👎 always high latency<br>👎 interrupts non-members |
| **Tree** | 👍 always low latency<br>👎 costly join/leave | 👍 good latency<br>for large groups<br>👎 bad latency<br>for small groups<br>👎 interrupts non-members |

Brandenburg
University of Technology
Cottbus - Senftenberg

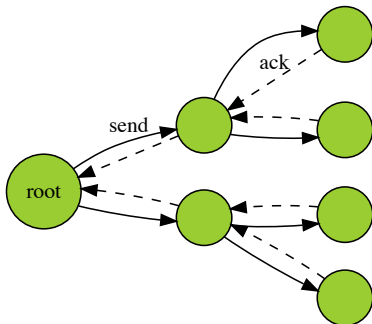## Problem Statement: Combine...

- **-** fast join/leave like with flat topology
- **-** low latency like in tree topologies (parallel propagation)

## Solution Idea

- **-** use static tree topology like in broadcasts
  (can be hand-crafted for the processor)
- **-** membership as bit-mask for fast join/leave
- **-** exploit shared memory to skip non-members,
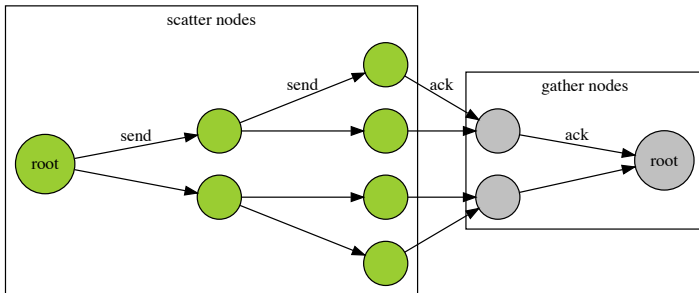  just message passing to actual members

# TREES WITH ACKNOWLEDGEMENT
## Logical nodes for larger design space & simpler code

Brandenburg
University of Technology
Cottbus - Senftenberg

scatter nodes

gather nodes

send

send

ack

ack

root

root

Brandenburg
University of Technology
Cottbus - Senftenberg

b·tu

Brandenburg
University of Technology
Cottbus - Senftenberg

Brandenburg
University of Technology
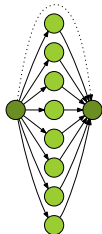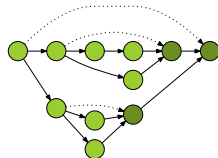Cottbus - Senftenberg

## Setup

## Flat Topology

## Binary Tree

- Intel XeonPhi Knights Corner (1.053 GHz)
- 60 cores
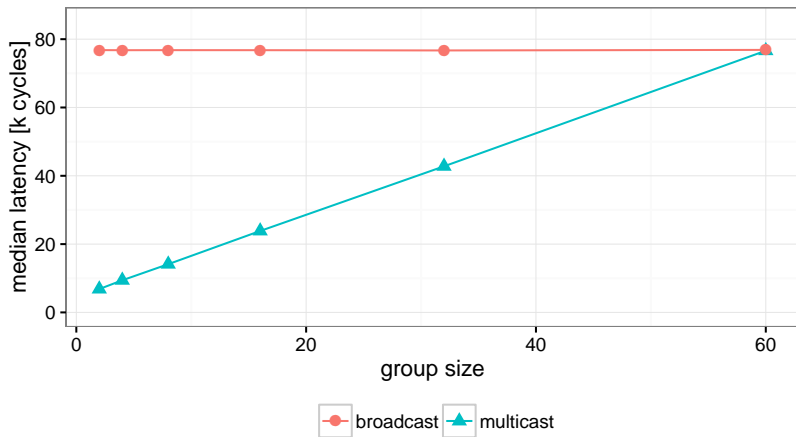- message passing via shared memory
- polling

**FLAT TOPOLOGY**
multicast similar to Linux TLB shootdown

# FLAT TOPOLOGY WITH HELPING
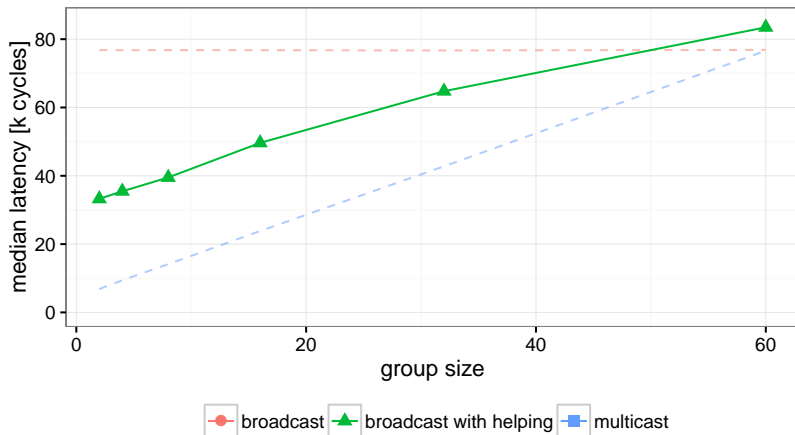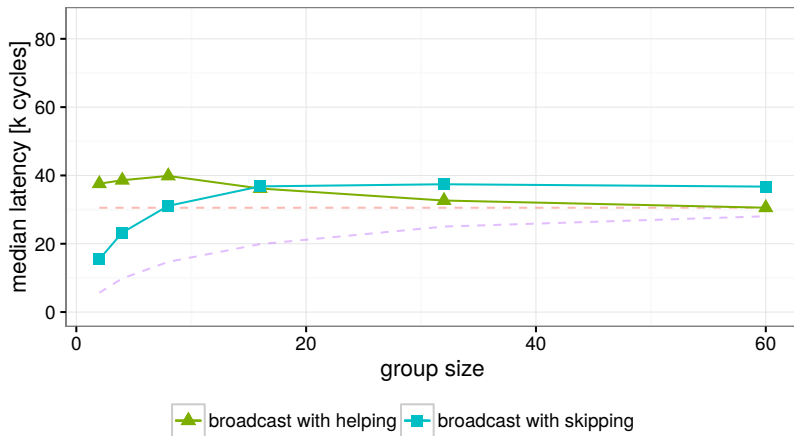**Overhead from membership tests and graph traversal**

b-tu Brandenburg
University of Technology
Cottbus - Senftenberg

# BINARY TREE WITH HELPING, SKIPPING



broadcast with helping   broadcast with skipping

Brandenburg
University of Technology
Cottbus - Senftenberg

Scalable, acknowledged, dynamic multicasts for manycores:

**Challenges:** generating good topologies is costly,
flat topology not scalable,
non-members should not be interrupted

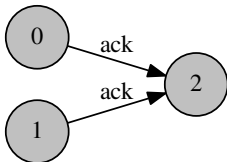**Solution:** static optimized broadcast topology,
help and skip non-member cores

**Result:** success for large groups, alright for small

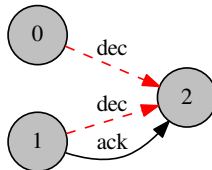**Implications:** improve Linux TLB shootdown for Many-Core HPC apps

Only message passing:

Using shared memory:

Brandenburg
University of Technology
Cottbus - Senftenberg

## Many-core systems

- **-** (incoherent) shared
  caches in cores/tiles
- **-** abundance of cores
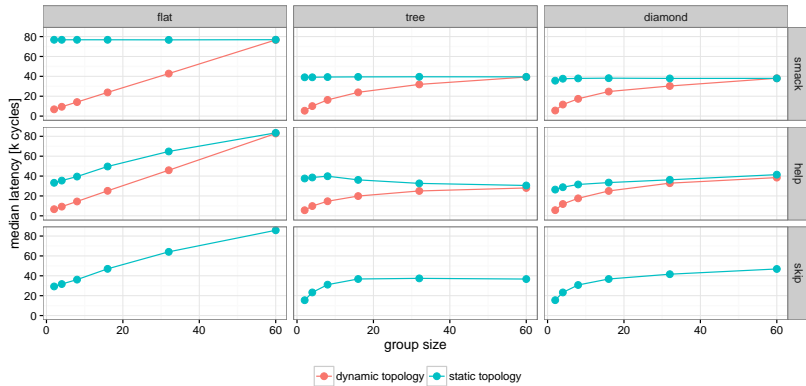- **-** symmetric multithreading

## Assumption

- **-** shared TLB can be
  invalidated by any sharer
- **-** superfluous hardware
  threads per shared cache

## Algorithm

- **-** one dedicated hardware thread per core/tile → still many cores
- **-** polling/waiting for notifications → core can not sleep
- **-** TLB invalidation for all threads in core/tile

# EVALUATION