

Smelt: Machine-aware Atomic Broadcast Trees for Multicores

Stefan Kaestle, **<u>Reto Achermann</u>**, Roni Haecki, Moritz Hoffmann, Sabela Ramos, Timothy Roscoe

Systems Group, Department of Computer Science, ETH Zurich



Large number of trees: topologies and send orders



There is no globally optimal tree structure

AMD Interlagos (4 Socket x 4 Cores x 2 Threads)



Intel Xeon Phi (61 Cores x 1 Thread) Execution Time [kCycles] 120 Binary tree or Fibonacci win 100 80 60 40 20 0 Reduction Barrier Broadcast 2PC Fibonacci Sequential MST

3



<u>Smelt</u>: Automatic optimization of broadcast and reduction trees



Example: Building fast and simple barriers



Barrier Benchmark on Intel Sandy Bridge 4x8x2

Dramatic improvement through automatic optimization of communication patterns Broadcasts and reductions are central building blocks for parallel programs

Performance

Atomic broadcasts

Replication for **data locality**

e.g. Shoal, Carrefour, SMMP OS, FOS **Fault-Tolerance**

Agreement protocols, atomic broadcasts

Replication for **failure resilience**

e.g. 1Paxos

Execution Control

Reductions, broadcast, barriers

Thread **synchronization**, data gathering

e.g. OpenMP

ETHzürich

Multicore hardware is complex



Smelt is based on peer-to-peer message passing



- Works well for our approach.
- Clear concept: Enables reasoning about send and receive costs

Message-passing on multicores is different

Classical Network Multicore interconnect Machine 4 Core 4 treceive Machine 3 Core 3 t_{receive} Machine 2 Core 2 t_{receive} On multicores send and receive times dominate propagation time **Goal:** Minimize total time of the broadcast

Minimizing the total time of a broadcast

- $t_{broadcast} = t_{last} t_{start}$
- Minimize the longest path from the root to the leaves.



$$t_{path} = \sum (t_{send} + t_{receive})$$

We need to know the send and receive cost between any pair of cores

¢ 10000

Information obtained from hardware discovery

AMD Interlagos 4x4x2

ֆ_ISCpu	
C NUMA distanc	e: abstract value
C Doesn't distin N send() a	nguish between and recv()
L ^{1d} cache:	161
L Symmetric:	$A \rightarrow B == B \rightarrow A$
L2 cache:	2048K
L3 cache:	6144K
NUMA node0 CPU(s):	0,4,8,12,16,20,24,28
NUMA node1 CPU(s):	32,36,40,44,48,52,56,60
NUMA node2 CPU(s):	2,6,10,14,18,22,26,30
NUMA node3 CPU(s):	34,38,42,46,50,54,58,62
NUMA node4 CPU(s):	3,7,11,15,19,23,27,31
NUMA node5 CPU(s):	35,39,43,47,51,55,59,63
NUMA node6 CPU(s):	1,5,9,13,17,21,25,29
NUMA node7 CPU(s):	33,37,41,45,49,53,57,61

\$ numactl -hardware node distances: node 0 1 2 3 4 5 6 7

ioue	0	T	Z	2	4	C	O	/
0:	10	16	16	22	16	22	16	22
1:	16	10	22	16	16	22	22	16
2 —	16	-22	10	16	-16	16	16	16
3:	22	16	16	10	16	16	22	22
4 –	16	16	-16	16	10	16	16	22
5:	22	22	16	16	16	10	22	16
6:	16	22	16	22	16	22	10	16
7:	22	16	16	22	22	16	16	10
			2		4			

Complement with microbenchmarks: pairwise send and receive

AMD Interlagos 4x4x2



Cost of Send Operation [Cycles]



Complement with microbenchmarks: pairwise send and receive

AMD Interlagos 4x4x2



Smelt

Using Smelt for group communication

#include <smelt/smelt.h>

void main() {

}

smelt_init();

smelt_topology_create();

smelt_broadcast(msg);



Smelt's tree generator heuristics

Remote Cores FirstAvoid Expensive
CommunicationImage: Cores FirstImage: Cores First<td

Maximize Parallelism



No redundancy







ETH zürich





Smelt Tree for Intel Sandy Bridge 4 Sockets x 8 Cores x 2 Threads

Evaluation Testbed

Intel

Architecture	Sockets	Cores / Socket	Threads / Core
Ivy Bridge	2	10	2
Nehalem	4	8	2
Knights Corner	1	61	4
Sandy Bridge	4	8	2
Sandy Bridge	2	10	2
Bloomfield	2	4	2

AMD

Architecture	Sockets	Cores / Socket	Threads / Core	
Magny Cours	4	12	1	
Barcelona	8	4	1	
Shanghai	4	4	1	
Interlagos	4	4	2	
Istanbul	4	6	1	

Full set of results online. http://machinedb.systems.ethz.ch

Smelt produces good trees across architectures



AMD Interlagos (4 Socket x 4 Threads)



Intel Xeon Phi (61 Threads)

Smelt produces good trees across architectures

slow	down										speedup
0.8		0.9	1.0		1.1		1.2		1.3		1.4
broadcast	- 1.24	1.06	1.11	1.37	1.13	1.10	1.16	1.15	1.07	1.22	1.01 -
barrier	- 1.12	1.07	1.30	1.41	1.09	1.08	1.13	1.03	1.09	1.38	1.02 -
2PC	- 1.17	1.09	1.22	1.35	1.10	1.13	1.11	1.07	1.17	1.33	1.01 -
reduction	- 1.18		1.08	1.27	1.24	1.01	1.24	1.09	1.18	1.53	1.21 -
1 KNC 1×61×A 2×4×2 2×8×2 2×10×2 4×4×1 4×4×1 4×4×2 4×12×1 4×8×2 4×8×2 4×8×2 8×4×1 A IS A IL A×4×2 15B 4×8×2 8×4×1 A IS A×6×1 A×4×2 12×1 A×8×2 4×8×2 8×4×1 A×1 A×1 A×1 A×1 A×1 A×1 A×1 A×1 A×1 A											

ETH zürich

Fast broadcast trees are good for reductions in most cases



Smelt provides simple and fast barriers

Barrier Benchmark on Intel Sandy Bridge 4x8x2

Execution Time [kCycles]



Barriers based on reduction and broadcast



Simple barrier implementation

OpenMP: EPCC OpenMP Benchmark Collection



Replaced GOMP barrier with Smelt

→ Remaining results on the website





Agreement Protocols: 1Paxos



4 clients to generate load N replicas executing 1Paxos

1Paxos Benchmark on AMD Interlagos 4x4x2



Summary

- Broadcasts and reductions are central building blocks
- No globally optimal tree topology
- Information from hardware discovery is not sufficient
- Smelt's produces good trees

Talk to us at

the first poster

session



machinedb.systems.ethz.ch