# State of the Art Meeting

Multi-level Multi-core Distributed Trace Synchronization

Benjamin Poirier

# What is Distributed Tracing?

# What is Distributed Tracing?

# How is distributed tracing done?

# How is distributed tracing done?

# Clock parameters

# Clock parameters

# Clock parameters

# Clock parameters

# Clock parameters

# Clock parameters

# Packet ordering

# Packet ordering

# Packet ordering

# Packet ordering

# Packet ordering

# Packet ordering

# Correction factors

# Linear Regression Method

# Convex Hull Method

# Convex Hull Method



Optimize:
{minimize, maximize} β

Subject to:
$\alpha + \beta \, x1 \leq y1$
$\alpha + \beta \, x2 \leq y2$
...
$\alpha + \beta \, x5 \geq y5$
...

With bounds:
$-\infty \leq \alpha \leq \infty$
$0 \leq \beta \leq \infty$

# Convex Hull Method

# Convex Hull Method

# Convex Hull Method



Optimize:
{minimize, maximize} $\alpha + \beta\ t$

Subject to:
$\alpha + \beta\ x1\ \leq\ y1$
$\alpha + \beta\ x2\ \leq\ y2$
...
$\alpha + \beta\ x5\ \geq\ y5$
...

With bounds:
$-\infty \leq \alpha \leq \infty$
$0 \leq \beta \leq \infty$

# Convex Hull Method

# Convex Hull Method

# Synchronization Accuracy

# Results

# Results

# Results

# Synchronization Accuracy



30s.                              60s.                              120s.

# Synchronization Accuracy



Fast Ethernet (100Mbs)

Gigabit Ethernet

# Synchronization Accuracy



120s.

960s.
(16 minutes)

15360s.
(4:16 hours)

# Statistics

```
Event count (2 traces):
    2.1-2nodes-1024/noeud1: 1140075
    2.1-2nodes-1024/noeud2: 1140285
    total events: 2280360
LTTV processing stats:
    received frames: 922900
    received frames that are IP: 922888
    received and processed packets that are TCP: 553952
    received and processed packets that are UDP: 245766
    sent packets that are TCP: 553813
TCP matching stats:
    Message traffic:
        0 - 1  : sent 215138     received 215076
Broadcast matching stats:
    total broadcasts datagrams received: 245762
    total broadcast groups for which all receptions were identified:
        122881
```

# Statistics

```
Linear regression analysis stats:
    Individual synchronization factors:
        0 - 1  : a0=  3.81325e+08 a1= 1 + 6.16834e-06 accuracy 23092.8
        1 - 0  : a0= -3.81251e+08 a1= 1 - 6.18234e-06 accuracy 23876.0

Convex hull analysis stats:
    out of order packets dropped from analysis: 0
    Number of points in convex hulls:
        0 - 1  : lower half-hull 16    upper half-hull 14
    Individual synchronization factors:
        0 - 1  : Middle
                a0=  3.81288e+08 a1= 1 + 6.17755e-06 accuracy 5.54264e-07
                a0:  3.81197e+08 to  3.81378e+08 (delta= 181440)
                a1: 1 +5.90042e-06 to +6.45468e-06 (delta= 5.54264e-07)
```

# Statistics

```
Linear regression analysis stats:
     Individual synchronization factors:
         0 - 1   : a0=   3.81325e+08 a1= 1 + 6.16834e-06 accuracy 23092.8
         1 - 0   : a0= -3.81251e+08 a1= 1 - 6.18234e-06 accuracy 23876.0


Convex hull analysis stats:
     out of order packets dropped from analysis: 0
     Number of points in convex hulls:
         0 - 1   : lower half-hull 16     upper half-hull 14
     Individual synchronization factors:
         0 - 1   : Middle
                    a0=   3.81288e+08 a1= 1 + 6.17755e-06 accuracy 5.54264e-07
                    a0:   3.81197e+08 to   3.81378e+08 (delta= 181440)
                    a1: 1 +5.90042e-06 to +6.45468e-06 (delta= 5.54264e-07)
```

# Statistics

```
Synchronization evaluation analysis stats:
    sum of broadcast differential delays: 1.71
    average broadcast differential delay: 1.39159e-05
    Individual evaluation:
        Trace pair  Inversions              Too fast            Total
        0 - 1                               47842 (23%)         215138
        1 - 0                               0 (0%)              215076
        total                               47842 (12%)         430214

Synchronization evaluation analysis stats:
    sum of broadcast differential delays: 1.11123
    average broadcast differential delay: 9.04313e-06
    Individual evaluation:
        Trace pair  Inversions              Too fast            Total
        0 - 1       0 (0%)                  18243 (9%)          215138
        1 - 0       0 (0%)                  6951 (4%)           215076
        total       0 (0%)                  25194 (6%)          430214
```

# Statistics

```
Synchronization evaluation analysis stats:
    sum of broadcast differential delays: 1.71
    average broadcast differential delay: 1.39159e-05
    Individual evaluation:
        Trace pair  Inversions          Too fast        Total
        0 - 1       156 (1%)            47842 (23%)     215138
        1 - 0       0 (0%)              0 (0%)          215076
        total       156 (1%)            47842 (12%)     430214

Synchronization evaluation analysis stats:
    sum of broadcast differential delays: 1.11123
    average broadcast differential delay: 9.04313e-06
    Individual evaluation:
        Trace pair  Inversions          Too fast        Total
        0 - 1       0 (0%)              18243 (9%)      215138
        1 - 0       0 (0%)              6951 (4%)       215076
        total       0 (0%)              25194 (6%)      430214
```

# Statistics

```
Synchronization evaluation analysis stats:
    sum of broadcast differential delays: 1.71
    average broadcast differential delay: 1.39159e-05
    Individual evaluation:
```

| Trace pair | Inversions | Too fast | Total |
|---|---|---|---|
| 0 - 1 | 156 (1%) | 47842 (23%) | 215138 |
| 1 - 0 | 0 (0%) | 0 (0%) | 215076 |
| total | 156 (1%) | 47842 (12%) | 430214 |

```
Synchronization evaluation analysis stats:
    sum of broadcast differential delays: 1.11123
    average broadcast differential delay: 9.04313e-06
    Individual evaluation:
```

| Trace pair | Inversions | Too fast | Total |
|---|---|---|---|
| 0 - 1 | 0 (0%) | 18243 (9%) | 215138 |
| 1 - 0 | 0 (0%) | 6951 (4%) | 215076 |
| total | 0 (0%) | 25194 (6%) | 430214 |

# Runtime Performance

# Future Developments

Synchronization across more than two nodes

# Future Developments

## Streaming trace synchronization

# Future Developments

Tracing on future multi-core architectures

# Future Developments

### Tracing on future multi-core architectures



AMD K10 (Opteron, Phenom)
- 2 to 6 cores
- Three-level cache
- Integrated memory controller
- HyperTransport bus

# Future Developments

## Tracing on future multi-core architectures



Tilera TILE64
- 64 cores
- Integrated memory controllers
- On-chip five lane "mesh" network

# Future Developments

Tracing on future multi-core architectures



STI Cell Broadband Engine
- 9 heterogenous cores
- Integrated memory controller
- Multi-lane ring bus

# Future Developments

Tracing on future multi-core architectures



Intel SCC
- 48 cores
- Integrated memory controllers
- Independant voltage and frequency "islands"

# Future Developments

core architectures

el SCC
48 cores
Integrated memory controllers
Independant voltage and frequency "islands"

# Future Developments

Tracing on future multi-core architectures



Intel SCC
- 48 cores
- Integrated memory controllers
- Independant voltage and frequency "islands"

# References

**Linear regression and convex hull synchronization algorithms**
- Duda, A., Harrus, G., Haddad, Y., and Bernard, G.: Estimating global time in distributed systems, Proc. 7th Int. Conf. on Distributed Computing Systems, Berlin, volume 18, 1987
- Ashton, P.: Algorithms for Off-line Clock Synchronisation, University of Canterbury, December 1995

**Streaming trace synchronization**
- Sirdey, R., and Maurice, F.: A linear programming approach to highly precise clock synchronization over a packet network, 4OR: A Quarterly Journal of Operations Research 6(4), volume 6, Springer, 393-401, 2008

**Systems of more than 2 nodes**
- Jezequel, J.M., and Jard, C.: Building a global clock for observing computations in distributed memory parallel computers, Concurrency: Practice and Experience 8(1), volume 8, John Wiley & Sons, Ltd Chichester, 1996
- Scheuermann, B., Kiess, W., Roos, M., Jarre, F., and Mauve, M.: On the Time Synchronization of Distributed Log Files in Networks With Local Broadcast Media, Networking, IEEE/ACM Transactions on 17(2), volume 17, 431-444, April 2009

**Multi-core architectures**
http://www.amd.com/us/products/desktop/processors/phenom/Pages/AMD-phenom-processor-X4-X3-at-home.aspx
http://www.tilera.com/products/TILE64.php
http://www.ibm.com/developerworks/power/cell/docs_documentation.html
http://www.intel.com/pressroom/archive/releases/20091202comp_sm.htm

# State of the Art Meeting

## Multi-level Multi-core Distributed Trace Synchronization

Benjamin Poirier