Introduction

Clocks 00000 Synchronization

Performance

Syster 000 Précis

TSCCLOCK: A LOW COST, ROBUST, ACCURATE SOFTWARE CLOCK FOR NETWORKED COMPUTERS

Darryl Veitch

d.veitch@ee.unimelb.edu.au http://www.cubinlab.ee.unimelb.edu.au/~darryl Collaboration with **Julien Ridoux**

CUBIN, Department of Electrical & Electronic Engineering The University of Melbourne

Google Tech Talk, Googleplex CA, July 23, 2007





Introduction	Clocks	Synchronization	Performance	System	Précis		
000	00000	000000	0000000	000	00		

WE NEED CLOCKS

THE OBVIOUS

- · Clocks used everywhere in networking and computing
- Network measurement (active and passive)
- Real-time services

The Future

- Tighter server integration
- Coordination of simultaneous multiple connections
- Latency a fundamental constraint in distributed services/computing Example: Virtual world realism in distributed games

Introduction	Clocks	Synchronization	Performance	System	Précis		
000	00000	0000000	0000000	000	00		

WE NEED CLOCKS

THE OBVIOUS

- · Clocks used everywhere in networking and computing
- Network measurement (active and passive)
- Real-time services

THE FUTURE

- Tighter server integration
- Coordination of simultaneous multiple connections
- Latency a fundamental constraint in distributed services/computing Example: Virtual world realism in distributed games

		e system	Frecis
00000 000000	0000000	000	00

RELIABILITY/ROBUSTNESS

- Can't build tight systems unless have bounds
- Status quo (*ntpd*) can be 'good' (1ms), bad (20ms), outrageous (100's ms)
- Measurement community pushed to GPS synchronized capture cards

Accuracy

- Greater accuracy → wider range of applications
- Status quo (*ntpd*) limited to 1ms, if all is well..
- Constant error versus 'jitter', absolute time versus time differences

- Use existing hardware (oscillator(s) in PCs)
- Network based synchronization cheap and convenient
- But GPS also cheap, right?

duction	Clocks	Synchronization	Performance	System	Précis
)	00000	0000000	0000000	000	00

RELIABILITY/ROBUSTNESS

- Can't build tight systems unless have bounds
- Status quo (*ntpd*) can be 'good' (1ms), bad (20ms), outrageous (100's ms)
- Measurement community pushed to GPS synchronized capture cards

ACCURACY

Intro

- Greater accuracy \rightarrow wider range of applications
- Status quo (*ntpd*) limited to 1ms, if all is well..
- Constant error versus 'jitter', absolute time versus time differences

- Use existing hardware (oscillator(s) in PCs)
- Network based synchronization cheap and convenient
- But GPS also cheap, right?

luction	Clocks	Synchronization	Performance	System	Précis
	00000	0000000	0000000	000	00

RELIABILITY/ROBUSTNESS

- Can't build tight systems unless have bounds
- Status quo (*ntpd*) can be 'good' (1ms), bad (20ms), outrageous (100's ms)
- Measurement community pushed to GPS synchronized capture cards

ACCURACY

Intro

- Greater accuracy \rightarrow wider range of applications
- Status quo (*ntpd*) limited to 1ms, if all is well..
- Constant error versus 'jitter', absolute time versus time differences

- Use existing hardware (oscillator(s) in PCs)
- Network based synchronization cheap and convenient
- But GPS also cheap, right?

luction	Clocks	Synchronization	Performance	System	Précis
	00000	0000000	0000000	000	00

RELIABILITY/ROBUSTNESS

- Can't build tight systems unless have bounds
- Status quo (*ntpd*) can be 'good' (1ms), bad (20ms), outrageous (100's ms)
- Measurement community pushed to GPS synchronized capture cards

ACCURACY

Intro

- Greater accuracy \rightarrow wider range of applications
- Status quo (*ntpd*) limited to 1ms, if all is well..
- Constant error versus 'jitter', absolute time versus time differences

- Use existing hardware (oscillator(s) in PCs)
- Network based synchronization cheap and convenient
- But GPS also cheap, WRONG! Can cost 12 months and \$10,000 to instrument a machine room

Introduction	Clocks	Synchronization	Performance	System	Précis
000	00000	000000	0000000	000	00
		T TOO			

DESIGN

- Re-engineered from scratch
- Built on
 - time-scale aware abstraction of oscillator performance
 - separate and decoupled treatment of rate and absolute time
 - RTT based delay filtering
 - feedforward not feedback
- Use oscillator driving CPU, accessible via TSC register (commonly available, high resolution, hardware updating, fast read)

- Very high robustness
- Accuracy an order of magnitude higher than *ntpd* (or more)
- Separate Absolute and Difference clocks

Introduction	Clocks	Synchronization	Performance	System	Précis
000	00000	000000	0000000	000	00
		T TOO			

DESIGN

- Re-engineered from scratch
- Built on
 - time-scale aware abstraction of oscillator performance
 - separate and decoupled treatment of rate and absolute time
 - RTT based delay filtering
 - feedforward not feedback
- Use oscillator driving CPU, accessible via TSC register (commonly available, high resolution, hardware updating, fast read)

- Very high robustness
- Accuracy an order of magnitude higher than *ntpd* (or more)
- Separate Absolute and Difference clocks

Introduction	Clocks	Synchronization	Performance	System	Précis		
000	00000	0000000	0000000	000	00		

DESIGN

- Re-engineered from scratch NOT just the fact that TSC is used!
- Built on
 - time-scale aware abstraction of oscillator performance
 - separate and decoupled treatment of rate and absolute time
 - RTT based delay filtering
 - feedforward not feedback
- Use oscillator driving CPU, accessible via TSC register (commonly available, high resolution, hardware updating, fast read)

- Very high robustness
- Accuracy an order of magnitude higher than *ntpd* (or more)
- Separate Absolute and Difference clocks

Introduction	Clocks	Synchronization	Performance	System	Précis		
000	00000	0000000	0000000	000	00		

DESIGN

- Re-engineered from scratch NOT just the fact that TSC is used!
- Built on
 - time-scale aware abstraction of oscillator performance
 - separate and decoupled treatment of rate and absolute time
 - RTT based delay filtering
 - feedforward not feedback
- Use oscillator driving CPU, accessible via TSC register (commonly available, high resolution, hardware updating, fast read)

- Very high robustness
- Accuracy an order of magnitude higher than *ntpd* (or more)
- Separate Absolute and Difference clocks

Introduction	Clocks	Synchronization	Performance	System	Pré
000	•0000	0000000	000000	000	00

OFFSET, SKEW AND DRIFT



Offset: Skew: Drift: rror $\theta(t) = C(t) - t$ of clock C(t) at time t rror in rate. E.g.: $\theta(t) = C + \gamma t$ (Simple Skew Model (SKM)) on-linear evolution of $\theta(t)$

Introduction	Clocks	Synchronization	Performance	System	Préc
000	00000	000000	0000000	000	00

OFFSET, SKEW AND DRIFT



Offset: error $\theta(t) = C(t) - t$ of clock C(t) at time t

- Skew: error in rate. E.g.: $\theta(t) = C + \gamma t$ (Simple Skew Model (SKM))
- Drift: non-linear evolution of $\theta(t)$

Clocks 0000

THE ROLE OF TIME SCALE

Laboratory:

 $\bar{p} = 1.82263812 * 10^{-9}$ (548.65527 Mhz) Machine Room: $\bar{p} = 1.82263832 * 10^{-9}$ (548.65521 Mhz)



Clocks 0000

THE ROLE OF TIME SCALE

Laboratory:

 $\bar{p} = 1.82263812 * 10^{-9}$ (548.65527 Mhz) Machine Room: $\bar{p} = 1.82263832 * 10^{-9}$ (548.65521 Mhz)



Short timescales: Large timescales: Simple Skew Model applies unpredictable drift must be tracked tion Clocks Synchronization Performance System COSCILLATOR STABILITY Allan deviation: scale dependent rate errors: $y_{\tau}(t) = \frac{\theta(t+\tau) - \theta(t)}{\theta(t+\tau) - \theta(t)}$



• SKM holds for $\tau^* = 1000$ [sec], (here TSC period p meaningful)

Average rate error upper bounded by 0.1 PPM no matter the scale





• SKM holds for $\tau^* = 1000$ [sec], (here TSC period *p* meaningful)

• Average rate error upper bounded by 0.1 PPM no matter the scale

Introduction	Clocks	Synchronization	Performance	System	Précis
000	00000	0000000	0000000	000	00

THE NEED FOR DIFFERENCE CLOCKS

Stable rate (p good to 10^{-7}) implies accurate $\Delta(t)$ measurement: Example: error in RTT of 100ms just 10ns

However an absolute clock $C_a(t)$ requires constant correction to negate drift:

- To synchronize $C_a(t)$, could
 - continuously modulate rate (*ntpd* uses ± 500 PPM band)
 - regularly add corrective jumps
- Either way, rate is disturbed
- Effect large! since drift estimation inherently difficult

Result: high native stability degraded by unbounded amount!

Introduction	Clocks	Synchronization	Performance	System	Précis
000	00000	0000000	0000000	000	00

THE NEED FOR DIFFERENCE CLOCKS

Stable rate (p good to 10^{-7}) implies accurate $\Delta(t)$ measurement: Example: error in RTT of 100ms just 10ns

However an absolute clock $C_a(t)$ requires constant correction to negate drift:

- To synchronize $C_a(t)$, could
 - continuously modulate rate (*ntpd* uses \pm 500 PPM band)
 - regularly add corrective jumps
- Either way, rate is disturbed
- Effect large! since drift estimation inherently difficult

Result: high native stability degraded by unbounded amount!

Introduction	Clocks	Synchronization	Performance	System	Précis
000	00000	0000000	0000000	000	00

THE NEED FOR DIFFERENCE CLOCKS

Stable rate (p good to 10^{-7}) implies accurate $\Delta(t)$ measurement: Example: error in RTT of 100ms just 10ns

However an absolute clock $C_a(t)$ requires constant correction to negate drift:

- To synchronize $C_a(t)$, could
 - continuously modulate rate (*ntpd* uses \pm 500 PPM band)
 - regularly add corrective jumps
- Either way, rate is disturbed
- Effect large! since drift estimation inherently difficult

Result: high native stability degraded by unbounded amount!

	4				43		
				c			
5							

Clocks

Synchronization

Performance 0000000 o O Précis

A DUAL CLOCK ARCHITECTURE

Foundation is the *uncorrected clock*: $C_u(t) = \bar{p} \cdot \text{TSC}(t) + K$

DIFFERENCE CLOCK

- Used for time differences below $au^* \sim 1000\,{
 m sec}$
- $C_d(t) = C_u(t)$ Example: $C_d(t_2) C_d(t_1) = \bar{p} \cdot (\text{TSC}(t_2) \text{TSC}(t_1))$
- Immune from errors in drift correction
- Use: RTTs, delay jitter, execution time, local event ordering ..

Absolute Clock

- Absolute timestamps (and time differences above τ^*)
- $C_a(t) = C_u(t) \hat{\theta}(t)$
- Drift correction estimate $\hat{\theta}(t)$ only applied when clock read
- Use: latency, global event ordering and scheduling ..

Require robust, accurate algorithms for \bar{p} and $\bar{\theta}$

duction	Clocks	Synchronization 1
)	00000	0000000

Performance

0

A DUAL CLOCK ARCHITECTURE

Foundation is the *uncorrected clock*: $C_u(t) = \bar{p} \cdot \text{TSC}(t) + K$

DIFFERENCE CLOCK

- Used for time differences below $\tau^* \sim 1000 \, {\rm sec}$
- $C_d(t) = C_u(t)$ Example: $C_d(t_2) C_d(t_1) = \bar{p} \cdot (\text{TSC}(t_2) \text{TSC}(t_1))$
- Immune from errors in drift correction
- Use: RTTs, delay jitter, execution time, local event ordering ..

Absolute Clock

- Absolute timestamps (and time differences above τ^*)
- $C_a(t) = C_u(t) \hat{\theta}(t)$
- Drift correction estimate $\hat{\theta}(t)$ only applied when clock read
- Use: latency, global event ordering and scheduling ..

Require robust, accurate algorithms for \bar{p} and $\hat{\theta}$

ction	Clocks	
	00000	

Synchronization

Performance

em O

A DUAL CLOCK ARCHITECTURE

Foundation is the *uncorrected clock*: $C_u(t) = \bar{p} \cdot \text{TSC}(t) + K$

DIFFERENCE CLOCK

- Used for time differences below $au^* \sim 1000\,{
 m sec}$
- $C_d(t) = C_u(t)$ Example: $C_d(t_2) C_d(t_1) = \bar{p} \cdot (\text{TSC}(t_2) \text{TSC}(t_1))$
- Immune from errors in drift correction
- Use: RTTs, delay jitter, execution time, local event ordering ..

ABSOLUTE CLOCK

- Absolute timestamps (and time differences above τ^*)
- $C_a(t) = C_u(t) \hat{\theta}(t)$
- Drift correction estimate $\hat{\theta}(t)$ only applied when clock read
- Use: latency, global event ordering and scheduling ..

Require robust, accurate algorithms for \bar{p} and $\hat{\theta}$

ction	Clocks	
	00000	

Synchronization

Performance

stem 00

A DUAL CLOCK ARCHITECTURE

Foundation is the *uncorrected clock*: $C_u(t) = \bar{p} \cdot \text{TSC}(t) + K$

DIFFERENCE CLOCK

- Used for time differences below $au^* \sim 1000\,{
 m sec}$
- $C_d(t) = C_u(t)$ Example: $C_d(t_2) C_d(t_1) = \bar{p} \cdot (\text{TSC}(t_2) \text{TSC}(t_1))$
- Immune from errors in drift correction
- Use: RTTs, delay jitter, execution time, local event ordering ..

ABSOLUTE CLOCK

- Absolute timestamps (and time differences above τ^*)
- $C_a(t) = C_u(t) \hat{\theta}(t)$
- Drift correction estimate $\hat{\theta}(t)$ only applied when clock read
- Use: latency, global event ordering and scheduling ..

Require robust, accurate algorithms for \bar{p} and $\hat{\theta}$





Obtain timestamps $\{T_{a,i}, T_{b,i}, T_{e,i}, T_{f,i}\}$ from *i*-th exchange

 ${T_{a,i}, T_{f,i}}:$ host timestamps in TSC counter units ${T_{b,i}, T_{e,i}}:$ server timestamps in seconds

Introduction	Clocks	Synchronization	Performance	System	Précis
000	00000	000000	000000	000	00

FILTERING NETWORK DELAYS

Choose RTT based filtering, not one-way (using same clock good!) Round–Trip Times r_i of packet *i*



Model for RTT:

 $r_i = r + \text{positive random noise}$ Filter using point error: excess over minimum RTT

÷							
L	2		u	c	5		
C		C)				

Clocks 00000 Synchronization

Performance 0000000 oo

Précis

NAIVE RATE SYNCHRONIZATION

Wish to exploit the relation $\Delta(t) = \Delta(\text{TSC}) * \bar{p}$

Naive estimate based on widely separated packets *j* and *i*:

$$\hat{p}_{i,j}^{\uparrow} \equiv \frac{T_{b,i} - T_{b,j}}{T_{a,i} - T_{a,j}}$$



Network delay and timestamping noise $\sim \frac{1}{\Delta(\text{TSC})}$, but errors not bounded.

Introduction	Clocks	Synchronization	Performance	System	Préc
000	00000	000000	0000000	000	00

RATE SYNCHRONIZATION ALGORITHM

Use selected naive estimates based on point error threshold



PROPERTIES

- Error quickly < 0.1 PPM, In 10mins, better than GPS!
- Error reduction (in timestamping, latency) guaranteed by $\Delta(t)$
- Inherently robust to packet loss, congestion, loss of server..
- Based on \bar{p} , no local rate estimates

Introduction	Clocks	Synchronization	Performance	System	Préc
000	00000	000000	0000000	000	00

RATE SYNCHRONIZATION ALGORITHM

Use selected naive estimates based on point error threshold



PROPERTIES

- Error quickly < 0.1 PPM, In 10mins, better than GPS!
- Error reduction (in timestamping, latency) guaranteed by $\Delta(t)$
- Inherently robust to packet loss, congestion, loss of server..
- Based on \bar{p} , no local rate estimates

Introduction	Clocks	Synchronization	Performance	System	Précis
000	00000	0000000	0000000	000	00

NAIVE ABSOLUTE SYNCHRONIZATION

Wish to exploit SKM over small scales to measure $\theta(t)$

Naive estimate again ignores network congestion, exploits steady rate over RTT

$$\hat{\theta}_i = \frac{1}{2}(C(t_{a,i}) + C(t_{f,i})) - \frac{1}{2}(T_{b,i} + T_{e,i})$$



Introduction	Clocks	Synchronization	Performance	System	Précis
000	00000	0000000	0000000	000	00

ABSOLUTE SYNCHRONIZATION ALGORITHM

Must track, so use all naive estimates, but carefully

Algorithm for $\hat{\theta}(t)$

- Weighted estimate of naive θ_i 's over SKM window
- Weights very strict, based on RTT quality (if quality very bad, freeze)
- Meaningful sanity check: ignore if hardware rate bound exceeded



Introduction	Clocks	Synchronization	Performance	System	Précis
000	00000	000000	000000	000	00

THE PATH ASYMMETRY

FUNDAMENTAL AMBIGUITY

Asymmetry $A \equiv d^{\uparrow} - d^{\downarrow}$ and $2\theta(t)$ non-unique up to a constant.

IMPACT ON ABSOLUTE CLOCK

- A unknown: generally forced to assume A = 0
- However, bounded by minimum RTT: $A \in (-r, r)$
- Create constant errors from 5μ s to 100's ms
- Causes jumps when server changed
- \rightarrow Important to use a single, close, server.

IMPACT ON DIFFERENCE CLOCK

- None
- Difference clock can be used to measure *r*



- GPS synchronized DAG card for external validation
- GPS synchronized SW and modified kernels for internal validation



- GPS synchronized DAG card for *external* validation
 - timestamps accurate to 100ns, but
 - comparison polluted by 'system noise'
 - splits asymmetry: $A = A_n + A_h$
 - allows network component A_n to be measured
 - host component A_h can only be bounded, can be >200 μ s!
- GPS synchronized SW and modified kernels for *internal* validation



- GPS synchronized DAG card for external validation
- GPS synchronized SW and modified kernels for *internal* validation
 - side by side timestamps cancels noise, but
 - only relative performance measurable, not absolute



A QUICK COMPARISON WITH ntpd



ntpd: sync'd to stratum-1 NTP server on LAN (broadcast mode) TSCclock: sync'd to stratum-1 NTP server outside LAN



Server: Polling Period: System Noise: Asymmetry: Stratum-1 NTP on LAN 256 sec $\sim 20\mu s$ Measured at $36\mu s$ and removed





INTERNAL VALIDATION: TSCCLOCK VS SW-GPS



Server: Polling Period: System Noise: Asymmetry: Stratum-1 outside LAN 16 sec $\ll 1\mu s$

As before for TSCclock, but SW-GPS component?



duction	Clocks	Synchronization	Performance	System	Précis
)	00000	0000000	0000000	000	00

PARAMETER DEPENDENCE



ntroduction	Clocks	Synchronization	Performance	System	Précis
	00000	0000000	0000000	000	00
	•				

COMPARISON WITH *nptd*

Server on LAN





Asymmetry: Same for each clock

Introduction	Clocks 00000	Synchronization	Performance ○○○○○●	System 000	Précis 00
	DIFFER	ENCE CLOC	K VERSUS	GPS	







Now compare if *no connectivity* with server

Introduction	Clocks 00000	Synchronization	Performance ○○○○○●	System 000	Précis 00
	DIFFER			GPS	





Introduction	Clocks	Synchronization	Performance	System	Précis
000	00000		0000000	•OO	00
		T			

THE SYSTEM

• API

- absolute and difference clock reading
- mode setting/reading
- diagnostics
- Timestamping solution
 - better with kernel support
- Synchronization algorithm:
 - runs as daemon or on command line
 - can store and replay log files
- Server
 - no server side solution, yet
 - client compatible with existing NTP servers
 - designed (and recommended) for use with a single server

Clocks	Synchronization	Performance	System	Précis
00000	0000000	0000000	000	00
	Turser			
	Clocks	Clocks Synchronization	Clocks Synchronization Performance	Clocks Synchronization Performance System

KERNEL

USER

45/52

Introduction	Clocks 00000	Synchronization	Performance	System ○●○	Précis 00
		Тилгота			

TIMESTAMPING

Kernel

- Packet timestamping:
 - Normal mode: TSCclock works in parallel with SW
 - TSCclock mode: SW also returns $C_a(t)$ transparently
- Other timestamping:
 - TSCclock works in parallel with SW

USER

Introduction	Clocks 00000	Synchronization	Performance	System OOO	Précis 00
		Тилгота			

TIMESTAMPING

Kernel

- Packet timestamping:
 - Normal mode: TSCclock works in parallel with SW
 - TSCclock mode: SW also returns $C_a(t)$ transparently
- Other timestamping:
 - TSCclock works in parallel with SW

USER

- Packet timestamping:
 - · Kernel packet timestamps inferred from userland
 - TSCclock works in parallel with SW

Introduction	Clocks	Synchronization	Performance	System	Précis		
000	00000	0000000	000000	000	00		
PACKAGING							

- Ubuntu 6.10 (Edgy)
- Ubuntu 7.04 (Feisty)
- Debian 4.0 (Etch)
- Fedora Core 6
- and soon Fedora Core 7 ...

Introduction	Clocks	Synchronization	Performance	System	Précis
000	00000	0000000	0000000	000	●○
			NC		

CONCLUSIONS

- TSCclock: for synchronization over networks
- Currently based on CPU oscillator accessible via TSC register
- Absolute Clock:
 - far more robust than *ntpd*
 - order of magnitude more accurate
- Difference Clock:
 - exceptionally robust
 - not available under *ntpd*
 - more accurate than standard GPS solution for small time intervals
- Kernel and userland packet timestamping solutions
- Low computational requirements
- Runs as daemon in parallel with *ntpd*
- Works with existing NTP server network
- Packages written for BSD and popular Linux distributions

Introduction	Clocks	Synchronization	Performance	System	Précis
000	00000		0000000	000	●○
		Conorna			

CONCLUSIONS

- TSCclock: for synchronization over networks
- Currently based on CPU oscillator accessible via TSC register
- Absolute Clock:
 - far more robust than *ntpd*
 - order of magnitude more accurate
- Difference Clock:
 - exceptionally robust
 - not available under *ntpd*
 - more accurate than standard GPS solution for small time intervals
- Kernel and userland packet timestamping solutions
- Low computational requirements
- Runs as daemon in parallel with *ntpd*
- Works with existing NTP server network
- Packages written for BSD and popular Linux distributions

Introduction	Clocks	Synchronization	Performance	System	Précis
000	00000		0000000	000	●○

CONCLUSIONS

- TSCclock: for synchronization over networks
- Currently based on CPU oscillator accessible via TSC register
- Absolute Clock:
 - far more robust than *ntpd*
 - order of magnitude more accurate
- Difference Clock:
 - · exceptionally robust
 - not available under *ntpd*
 - more accurate than standard GPS solution for small time intervals
- Kernel and userland packet timestamping solutions
- Low computational requirements
- Runs as daemon in parallel with *ntpd*
- Works with existing NTP server network
- Packages written for BSD and popular Linux distributions

Introduction	Clocks	Synchronization	Performance	System	Précis		
	00000	0000000	0000000	000	○●		
Links							

- Publications: http://www.cubinlab.ee.unimelb.edu.au/articles
- TSCclock page: http://www.cubinlab.ee.unimelb.edu.au/tscclock