

The On-going Evolutions of Power Management in Xen

Kevin Tian

Open Source Technology Center (OTC)

Intel Cooperation



Software and Solutions Group

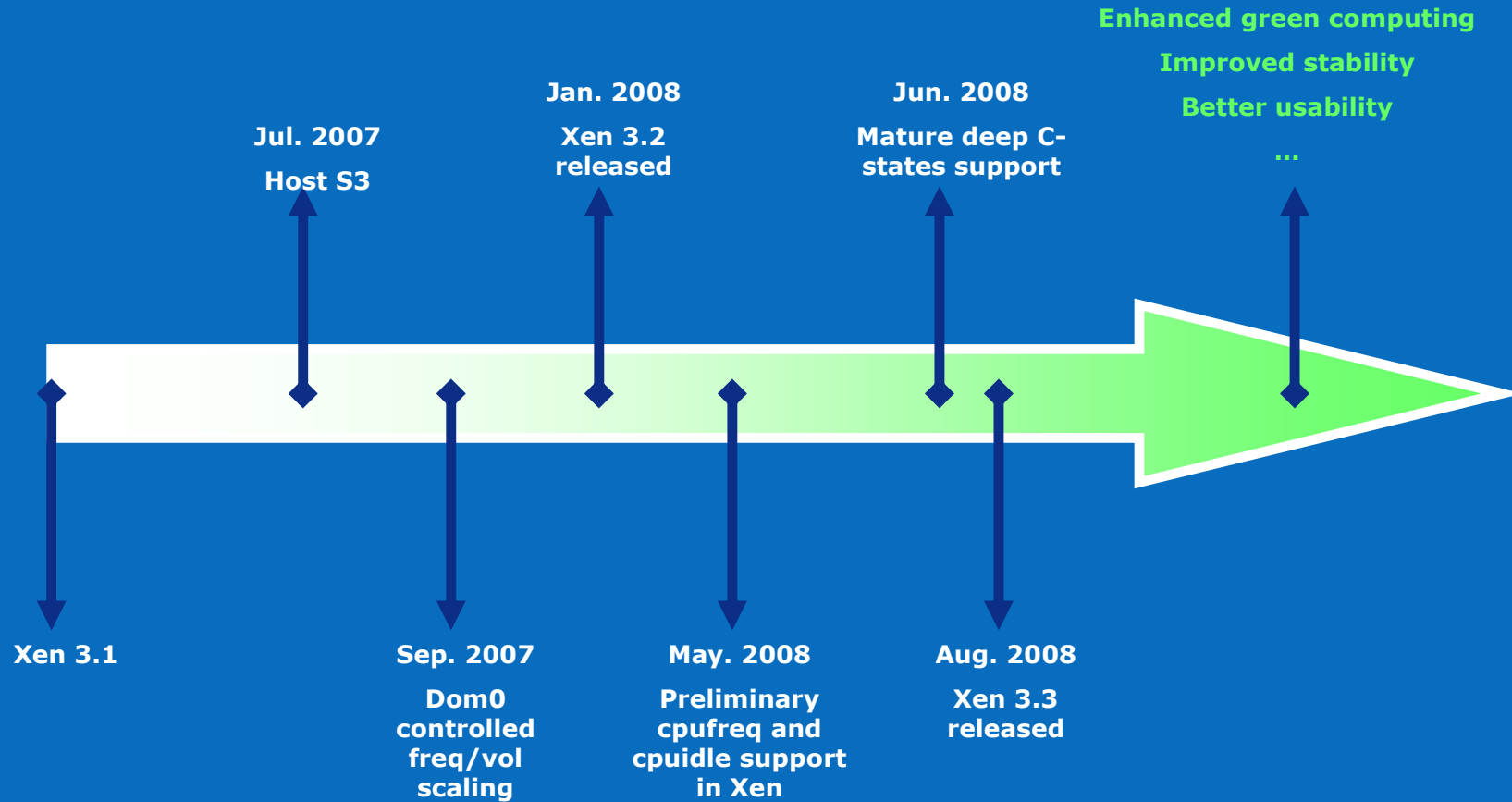


Agenda

- Brief history
- Evolutions of Idle power management
- Evolutions of run-time power management
- Tools
- Experimental data about Xen power efficiency
- Power impact from VM



Brief History



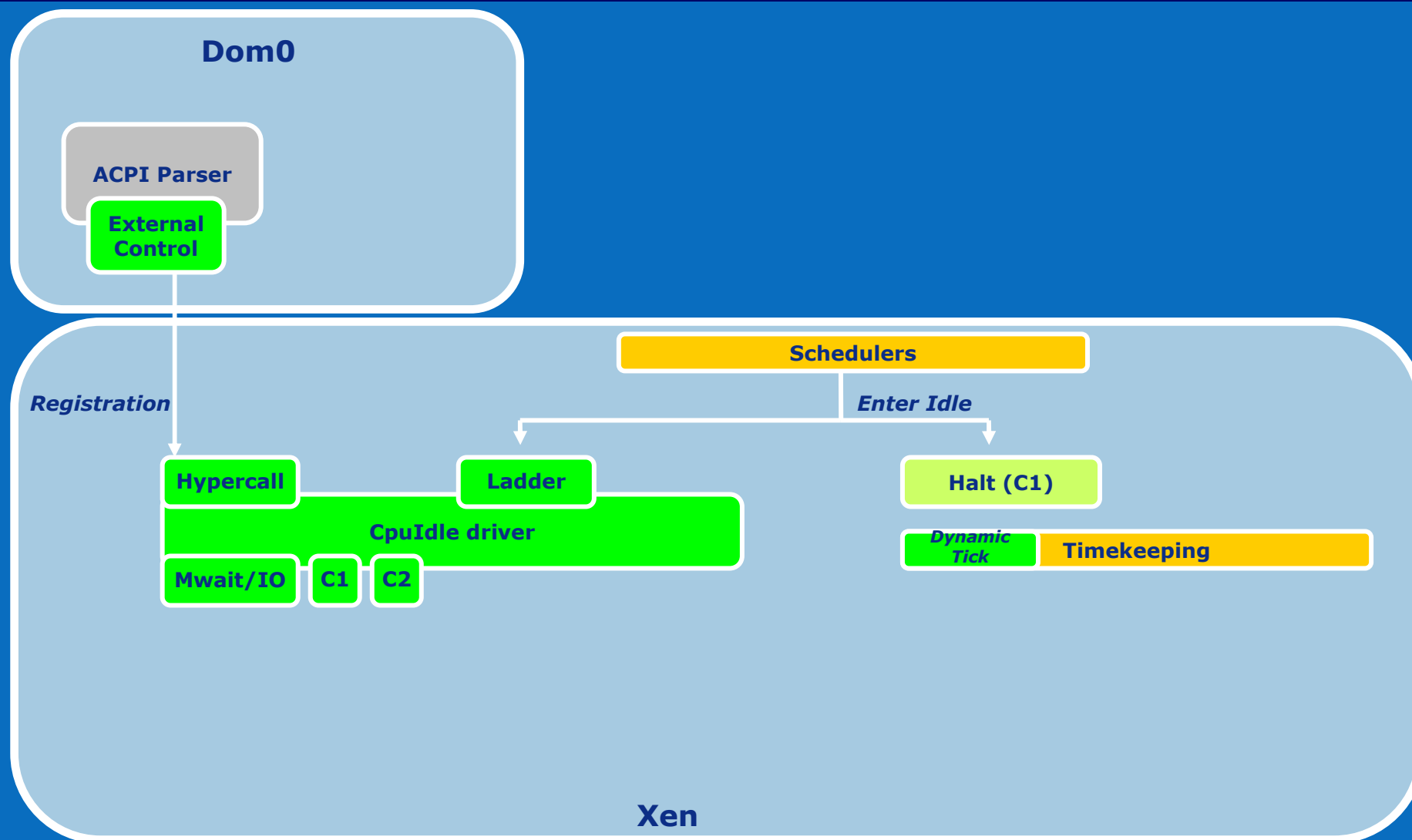
Evolutions of Idle power management (Xen cpuidle)



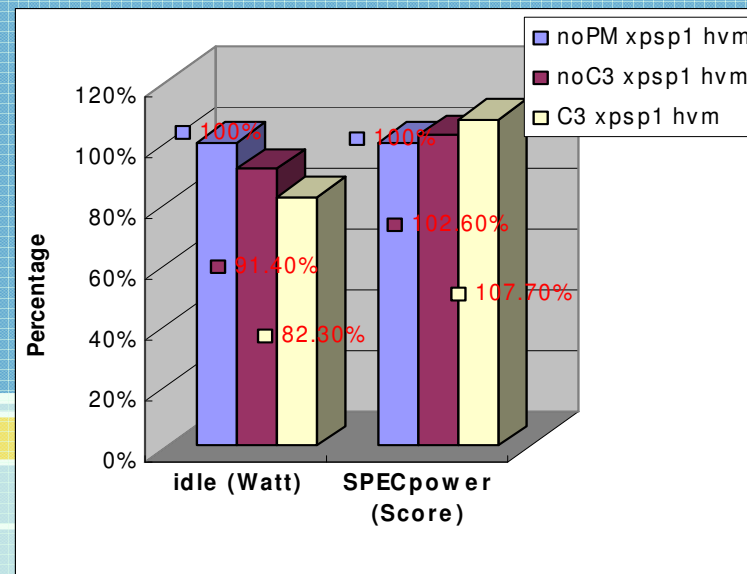
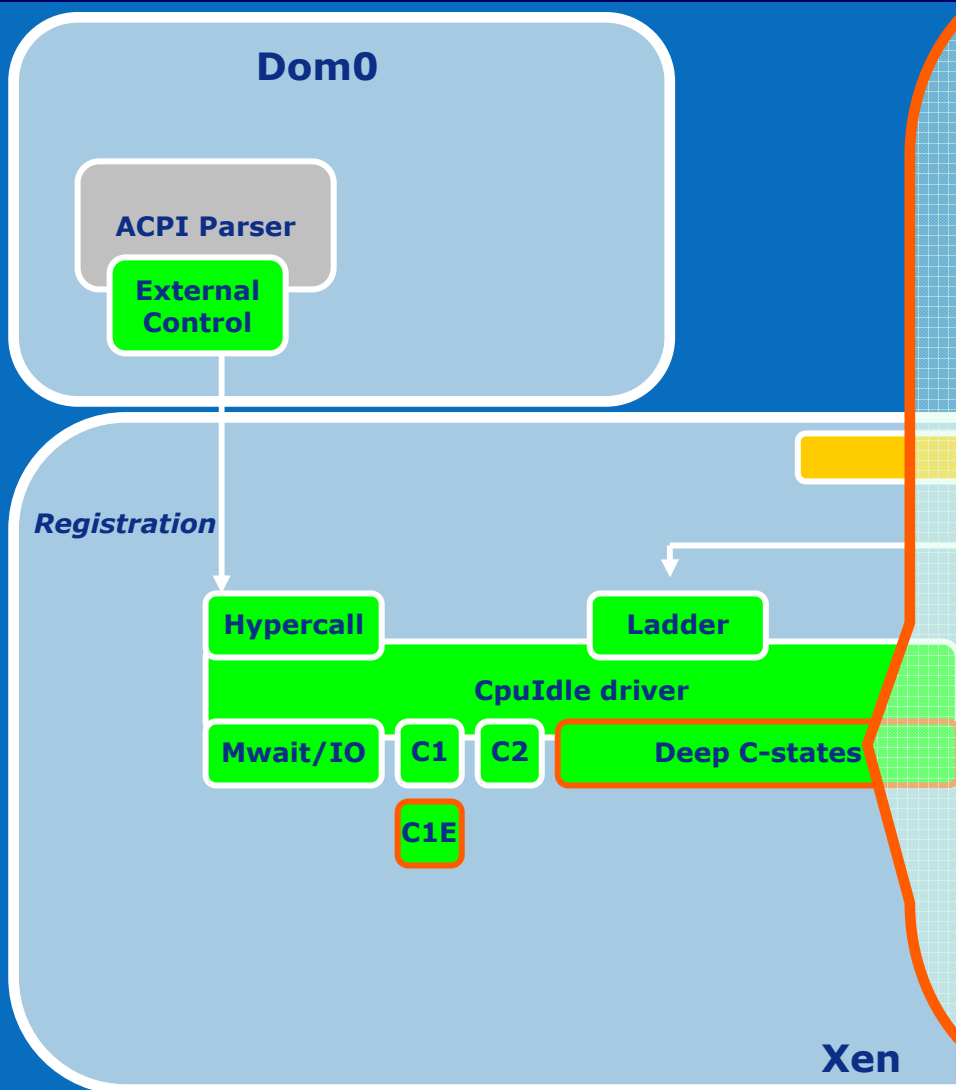
Software and Solutions Group



Xen Summit Boston 2008



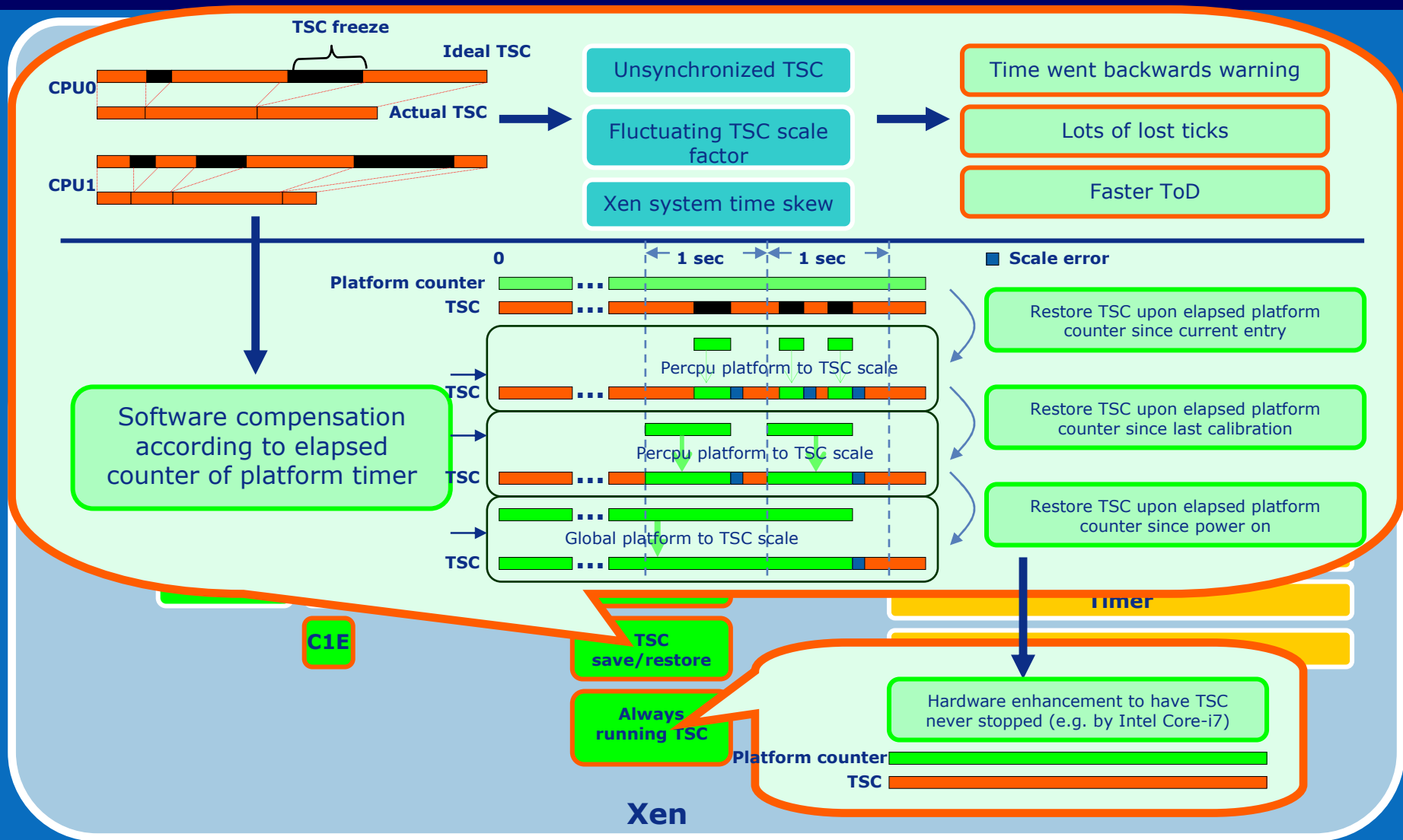
Enhanced C-states support



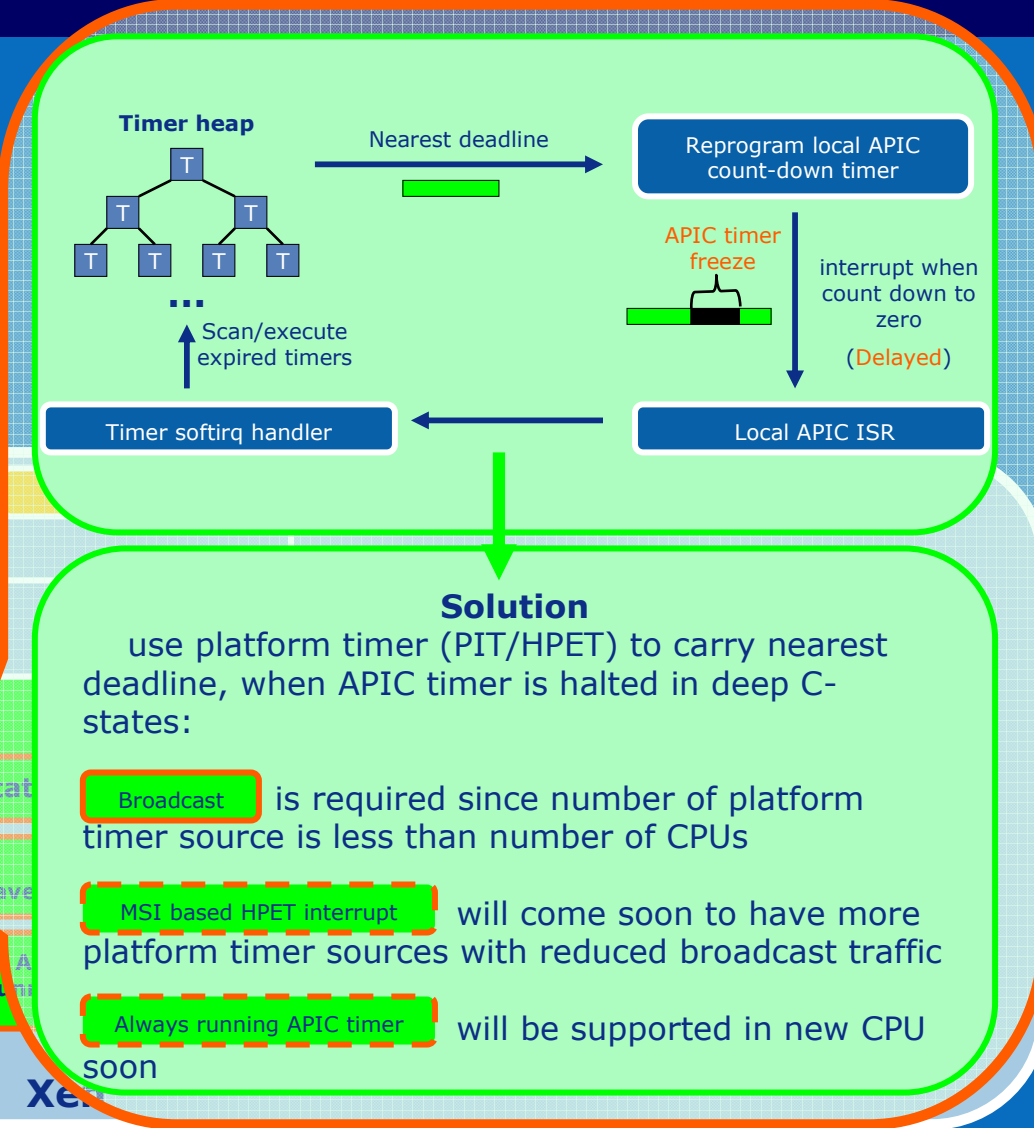
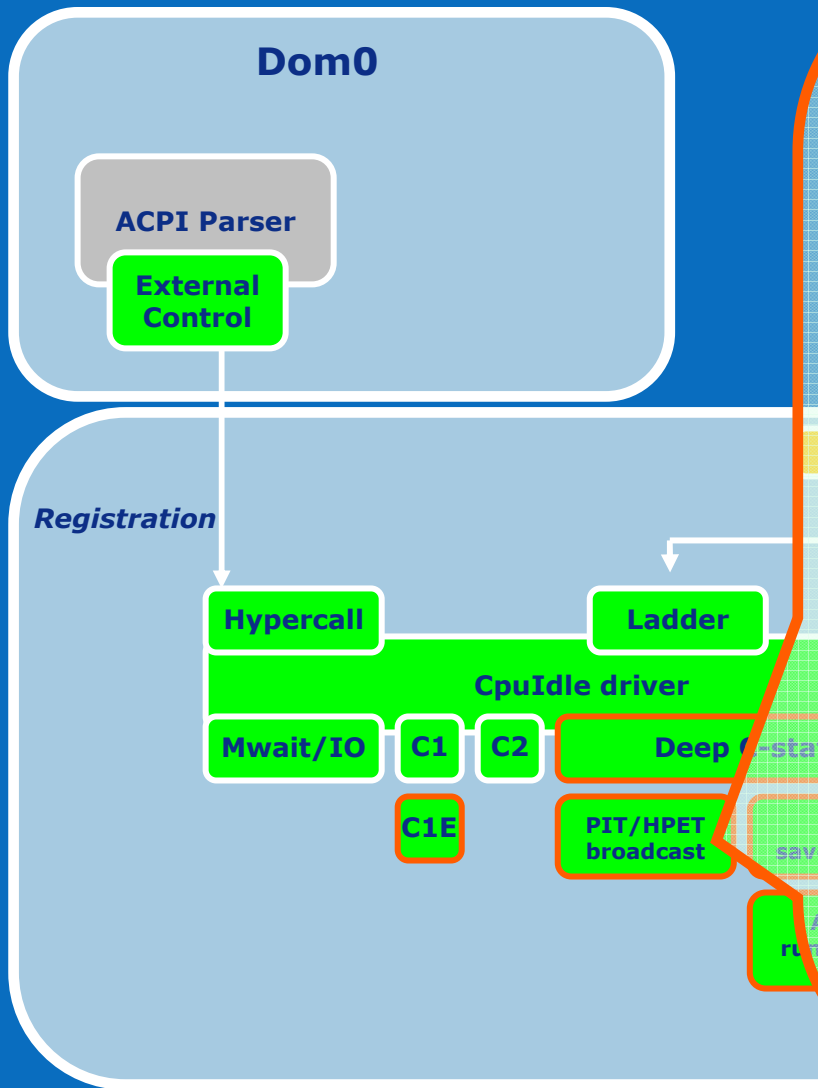
'noPM' has both cpuidle and cpufreq disabled, and vice versa for other two cases. Compared to 'C3', noC3 has maximum C-state limited to C2

For idle watt, lower value means greener. For SPECpower score, higher value indicates more power efficient

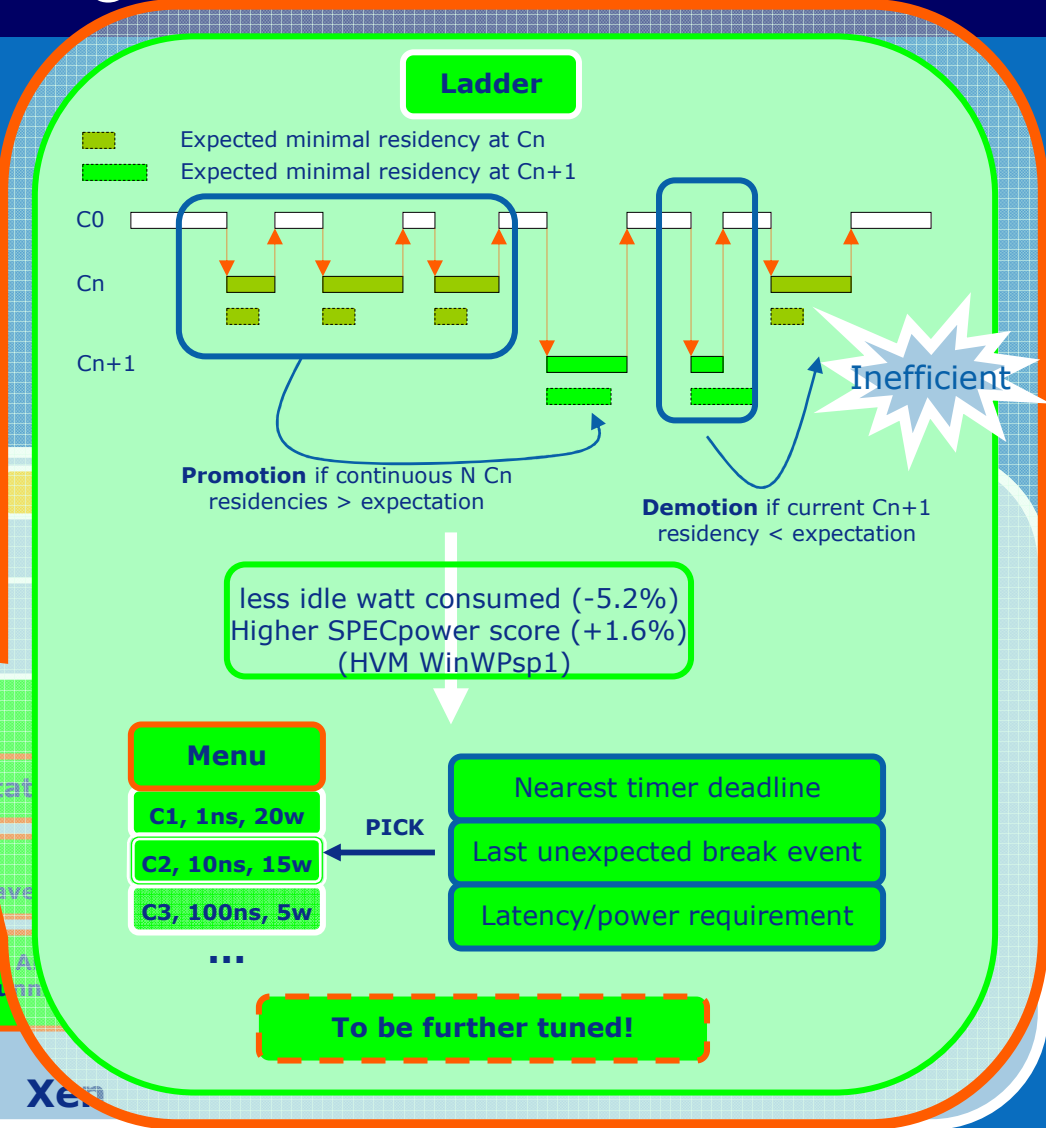
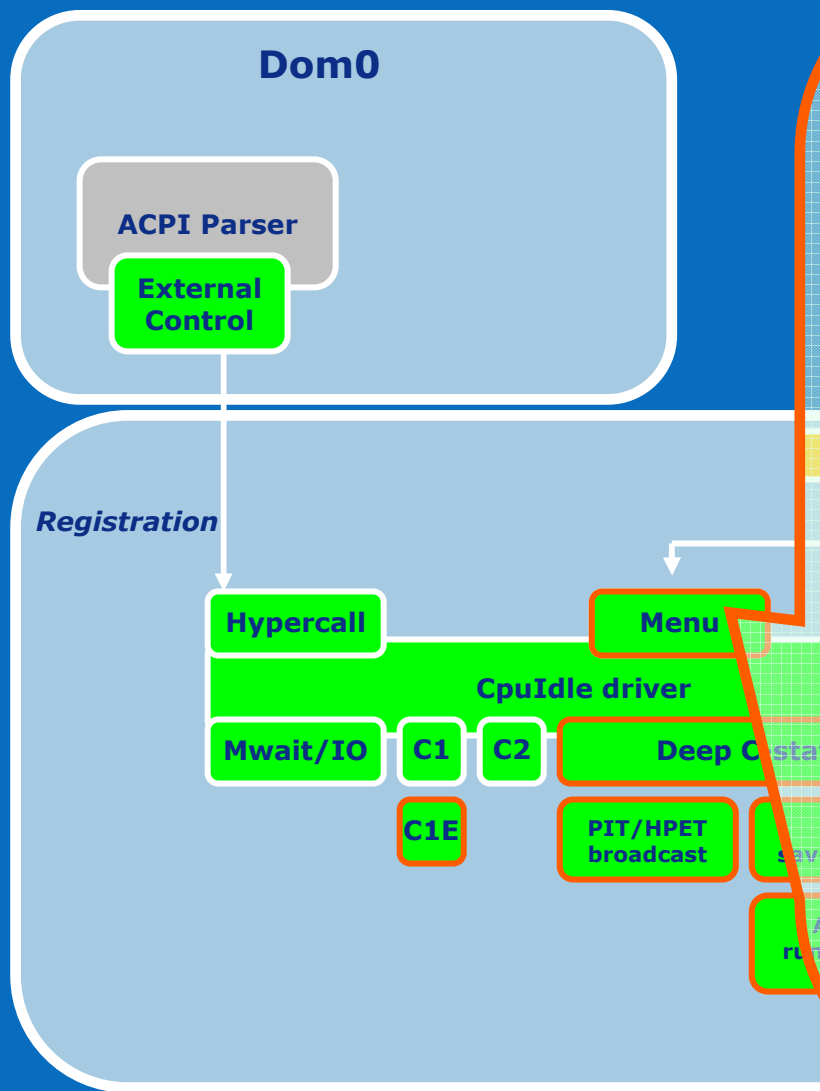
TSC freeze



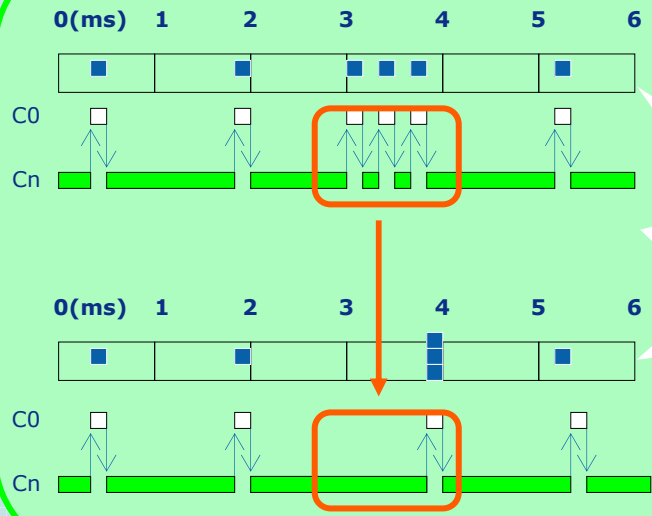
APIC timer freeze



Menu governor



Range timer

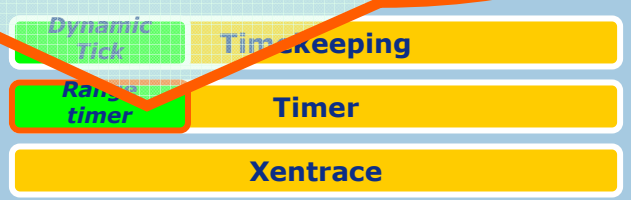
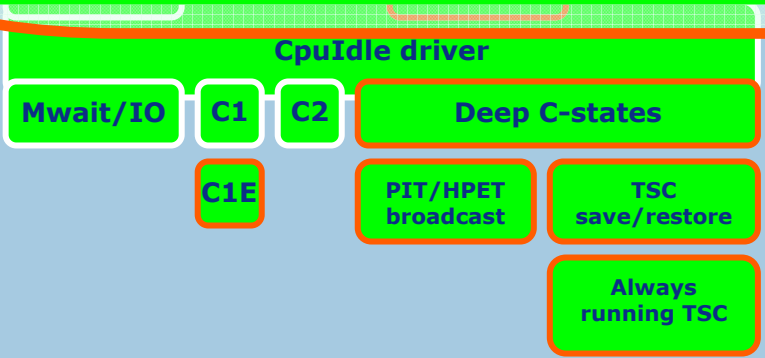


Frequent C-state Entry/exit may Instead consume More power

For each timer, it accepts a range for expiration now:

$[expiration, expiration + timer_slop]$
(default 50us for timer_slop)

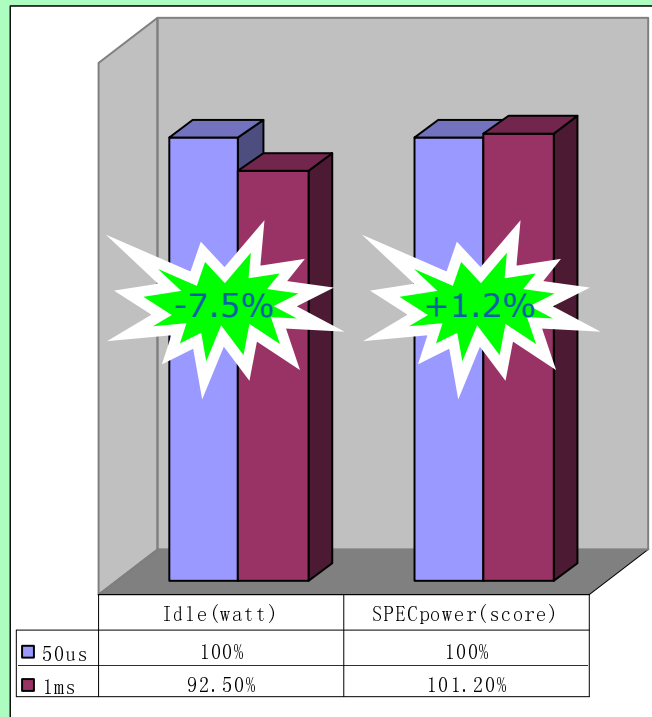
Overlapping ranges can be merged to reduce timer interrupt count



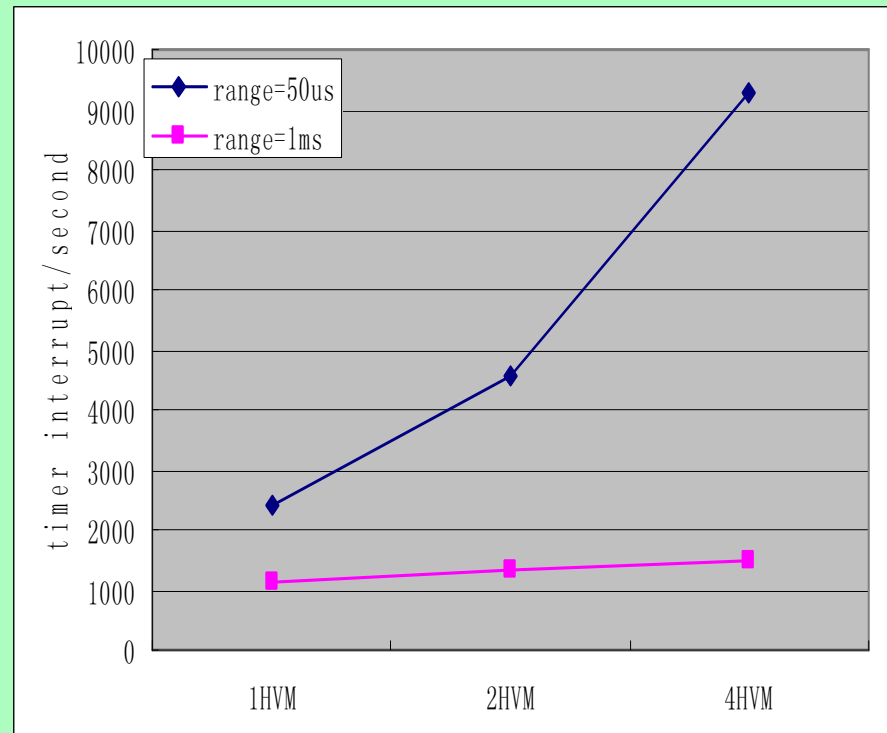
Xen

Range timer effect

One UP HVM RHEL5u1

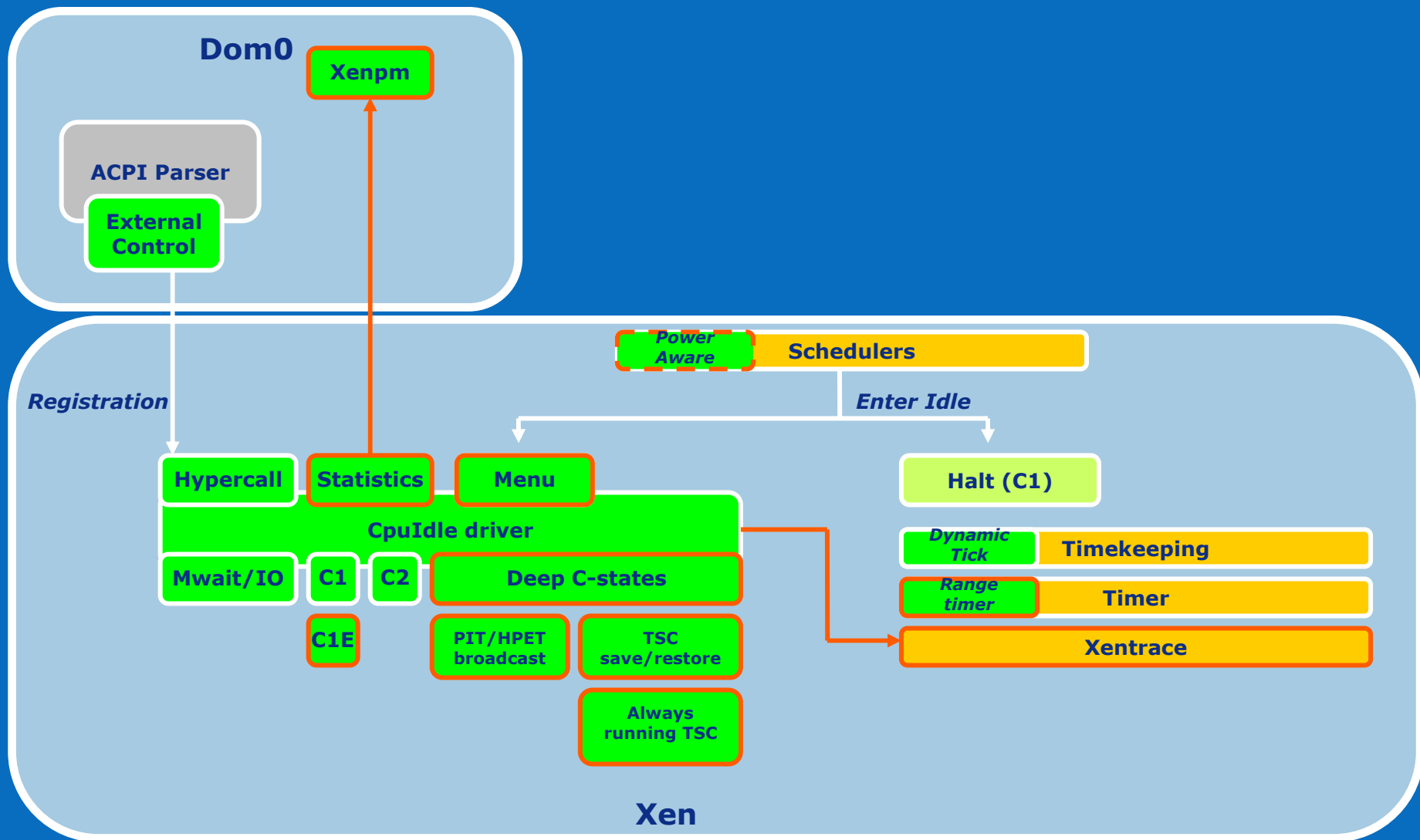


Multiple idle UP HVM RHEL5u1



Collected on a two-cores mobile platform

Current picture



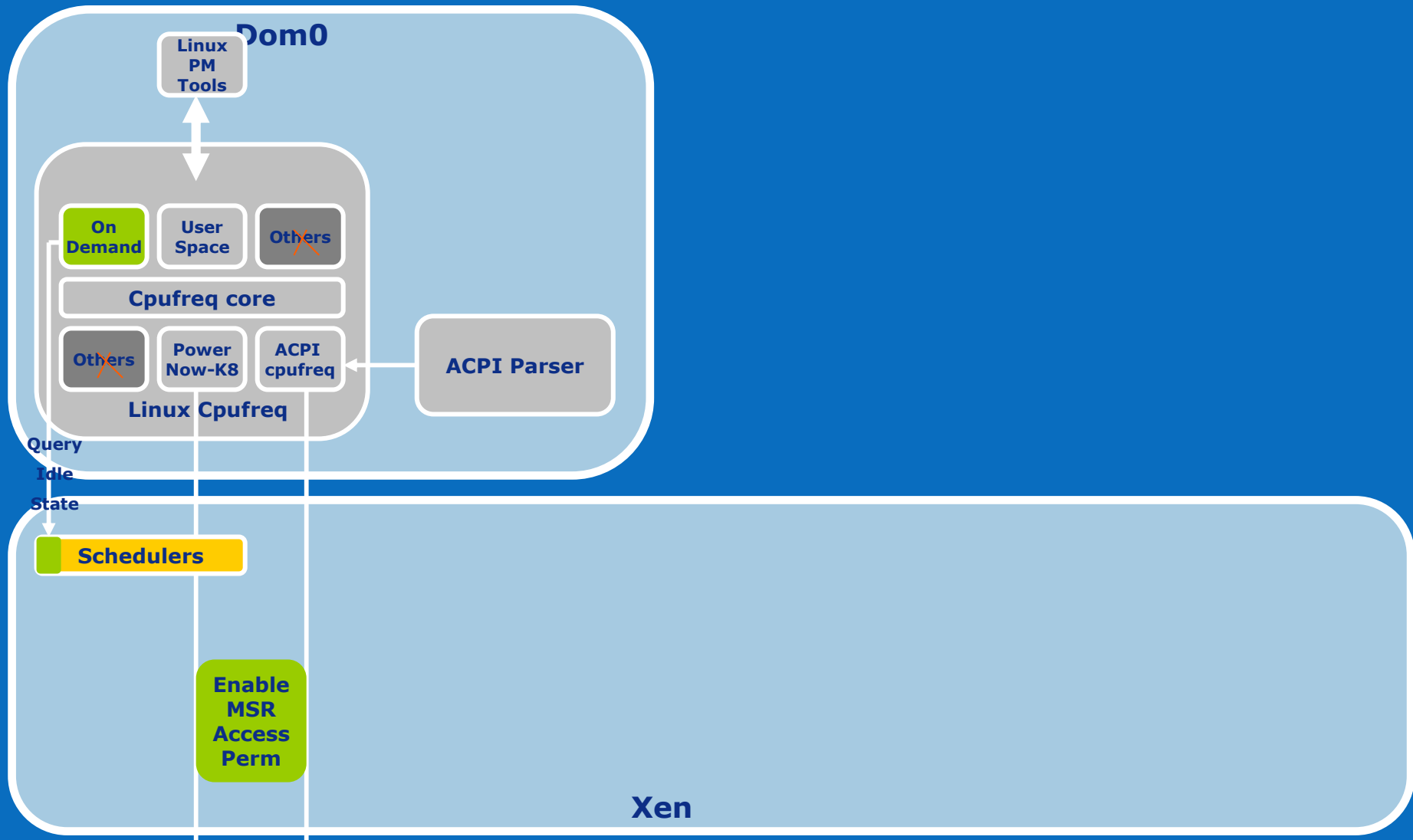
Evolutions of run-time power management (Xen cpufreq)



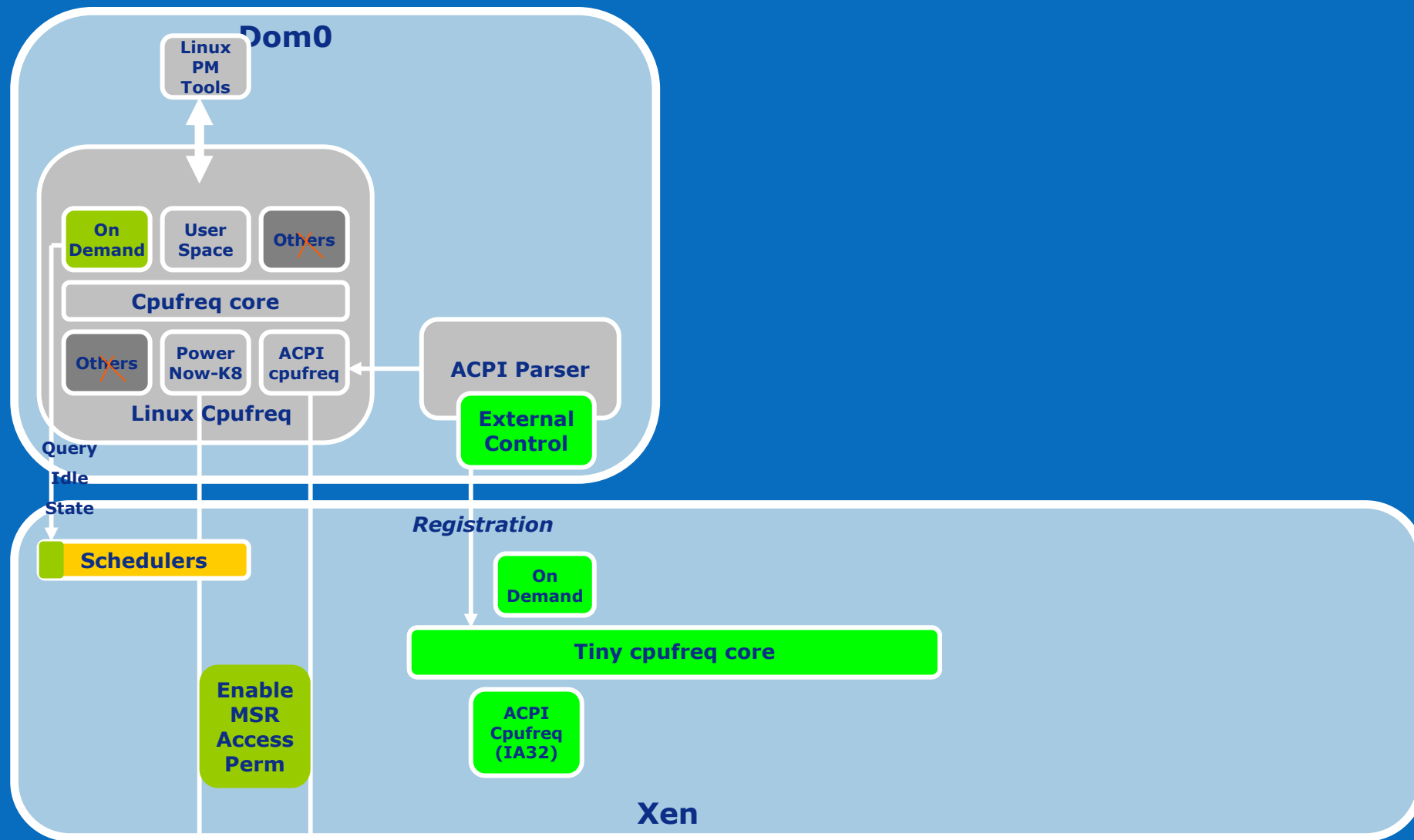
Software and Solutions Group



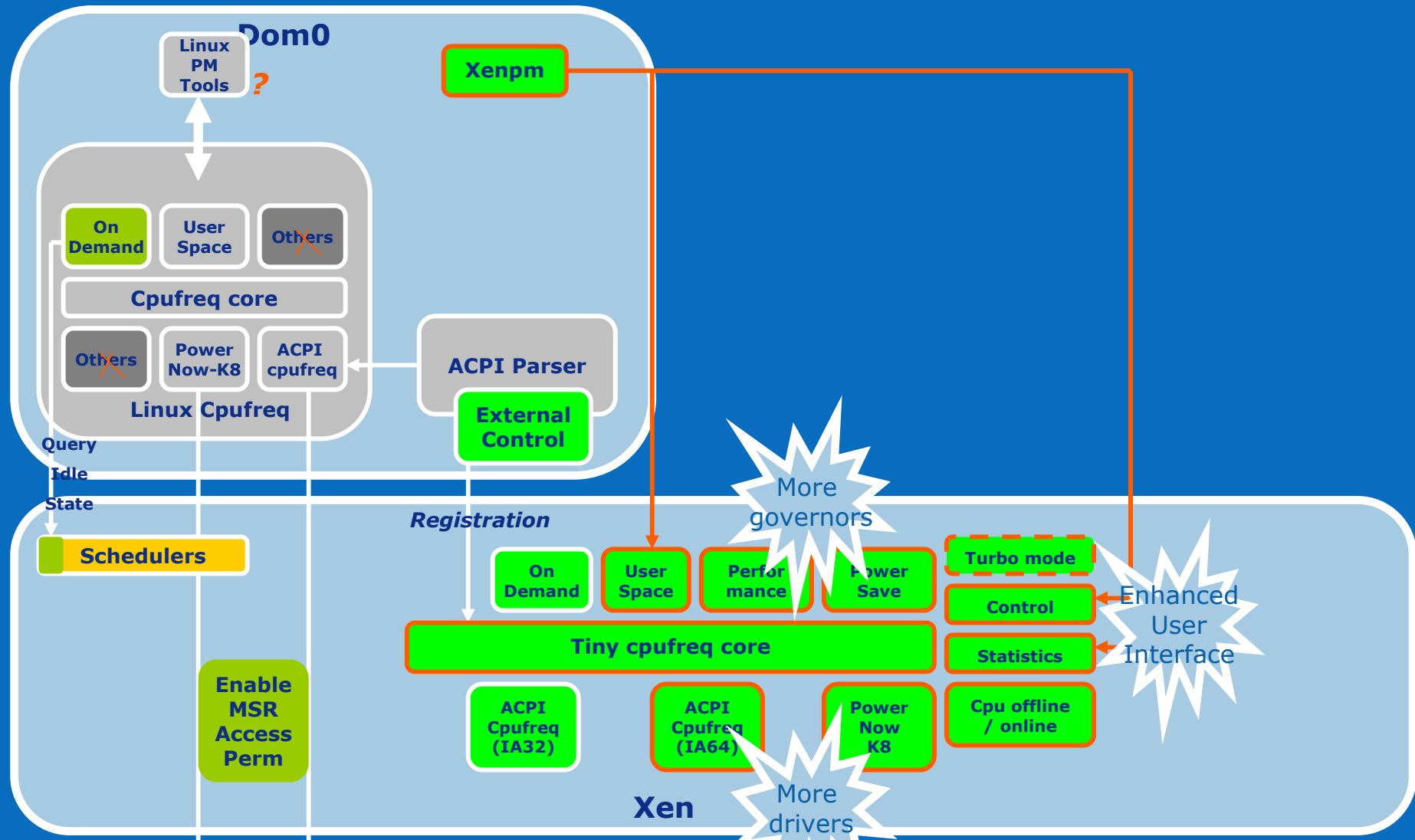
Xen 3.2



Xen Summit Boston 2008



Current picture



Open Source
Technology
Center

Software and Solutions Group



Tools



Software and Solutions Group



Dom0

Xenpm

Retrieve run-time statistics about Xen cpuidle and cpufreq

Apply user policy on exposed control knobs of Xen cpufreq (governor, set freq, etc)

More capabilities to be added later, e.g. profile...

```
[root@ktian1-dev ~]# xenpm start
Elapsed time (ms): 2883
CPU0:
  residency    percentage
C0      1 ms      0.06%
C1      0 ms      0.00%
C2     2882 ms  99.94%
P0      0 ms      0.00%
P1      0 ms      0.00%
P2      0 ms      0.00%
P3      0 ms     100.00%
CPU1:
  residency    percentage
C0      4 ms      0.16%
C1      0 ms      0.00%
C2     2879 ms  99.84%
P0      0 ms      0.00%
P1      0 ms      0.00%
P2      0 ms      0.00%
P3      2 ms     100.00%
```

Log every state change for Xen cpuidle and cpufreq:

```
CPU0 391365842416 (+ 21204) cpu_idle_entry [ idle to state 2 ]
CPU0 391375951050 (+10108634) cpu_idle_exit [ return from state 2 ]
```

Raw data could be further processed by other scripts

Xentrace

Xen



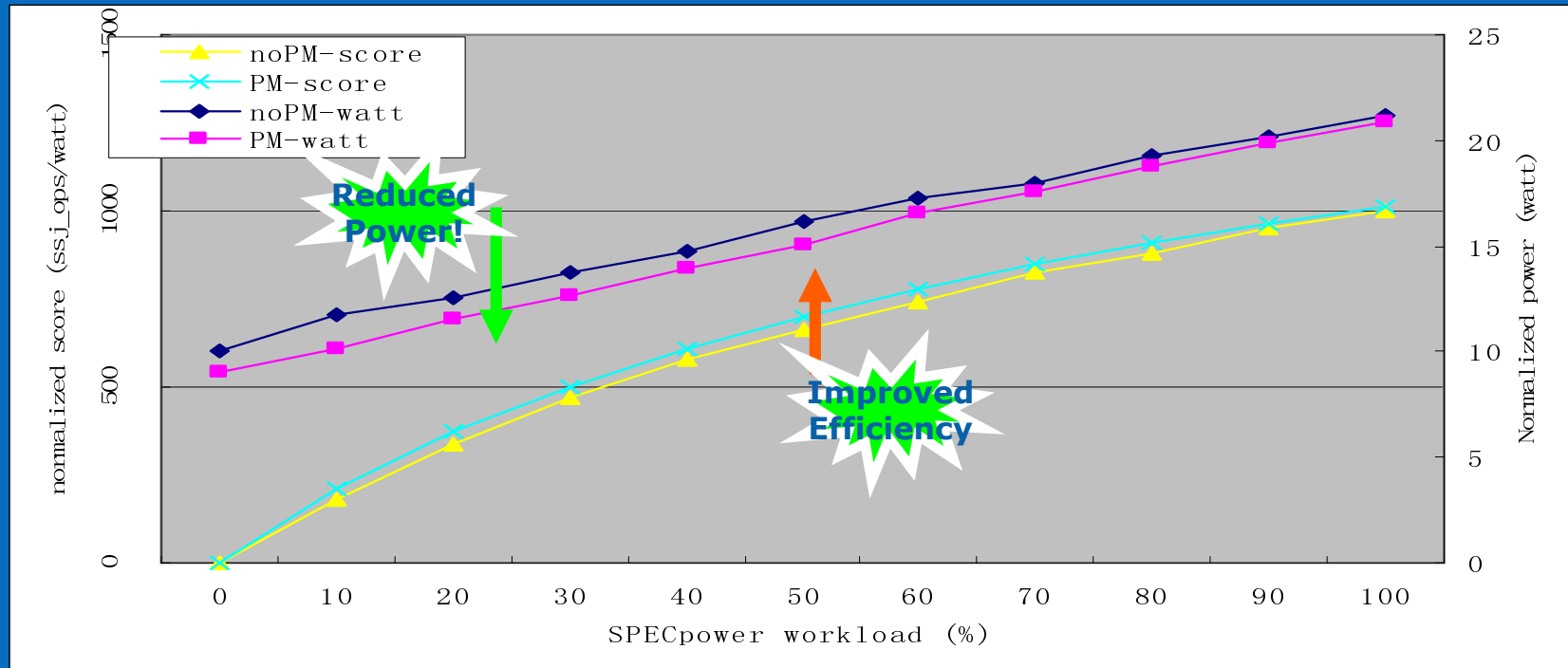
Experimental data about Xen power efficiency



Software and Solutions Group

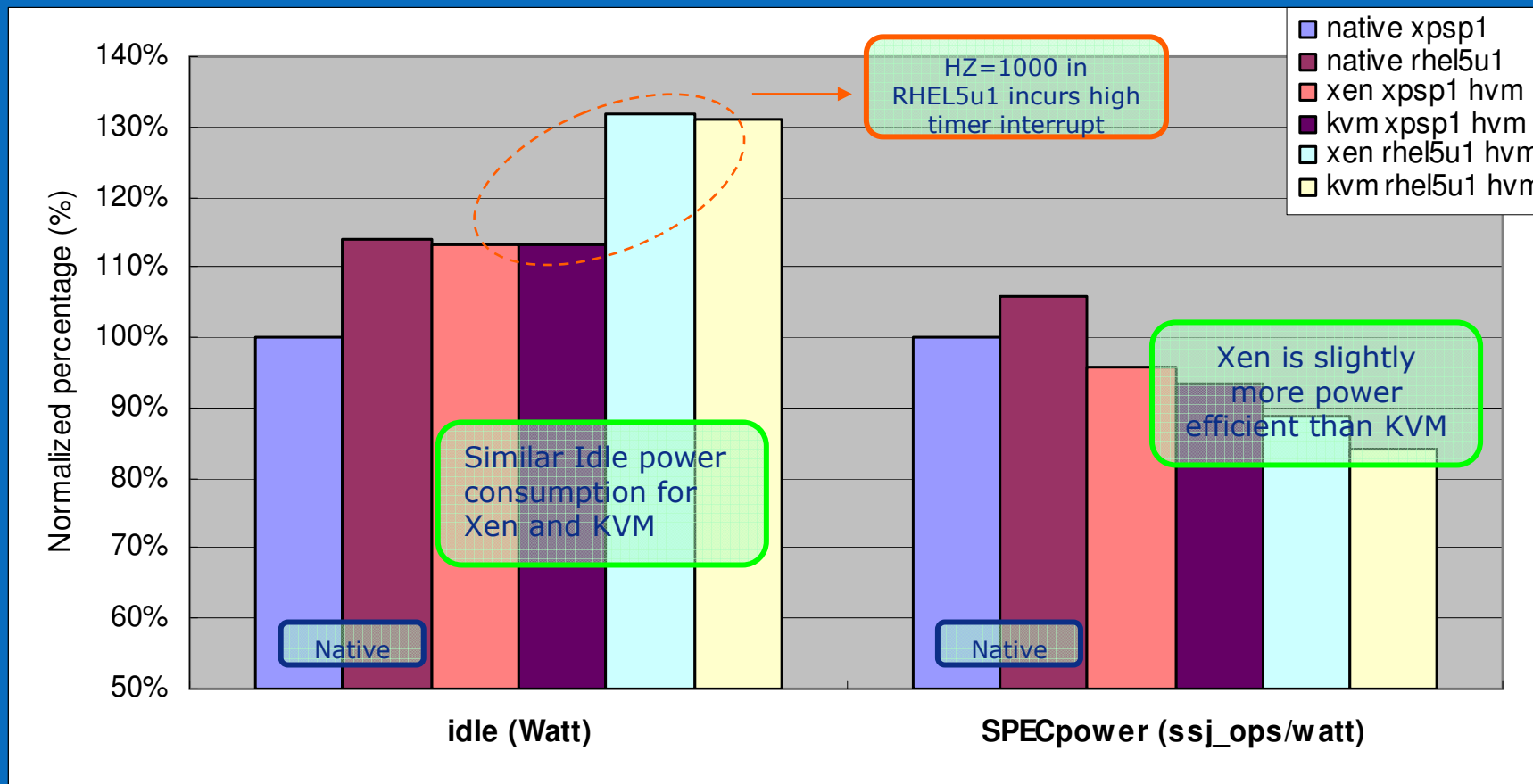


- All data shown in this section:
 - For reference only and not guaranteed
 - On a two-core mobile platform, with one HVM guest created
 - Server consolidation effect with multiple VMs/workloads are in progress
- Improvement when Xen cpuidle and cpufreq are enabled
 - SPECpower score is normalized (100% noPM score is 1000 ssj_ops / watt)
 - Similarly, consumed watt is also normalized (idle noPM watt is 10w)



- Below is a more attracting comparison

- Native WinXPsp1's idle watt and SPECpower score are both normalized to 100% as the base



Power impact from VM



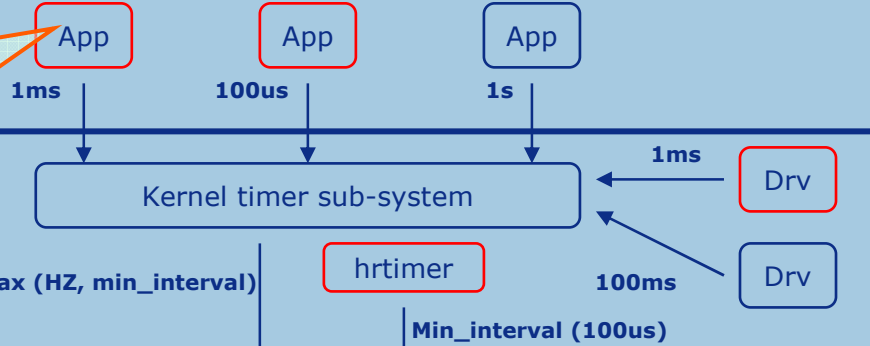
Software and Solutions Group



Dom0

DomN

Bad App eats power like in native!



Timekeeping
Schedule
Accounting

HZ

Idle tickles

Clock event

Kernel implementation matters!

Green features in VM helps!

VMM shouldn't Be only Blamed!



Xen



Open Source
Technology
Center

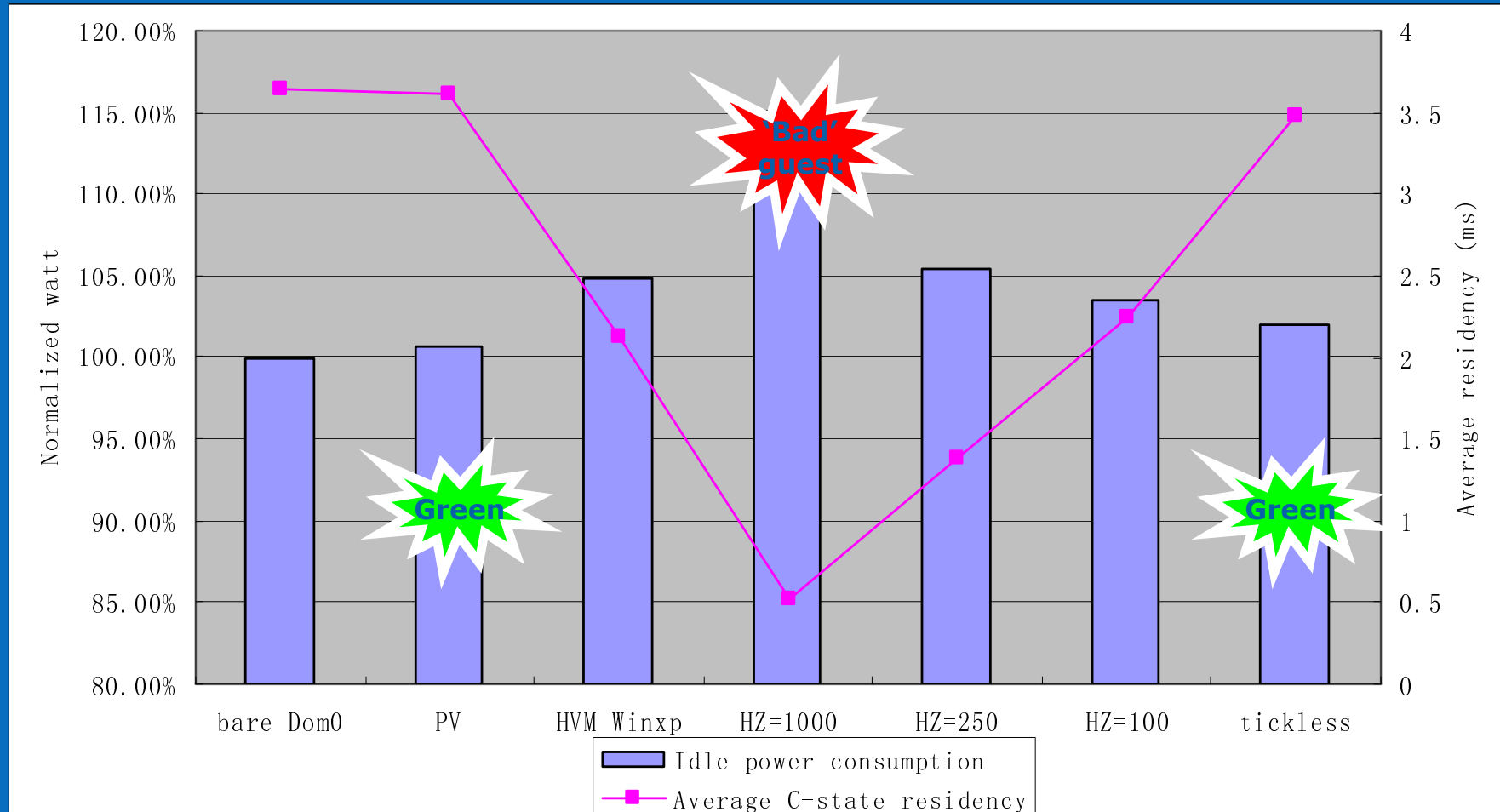
Software and Solutions Group



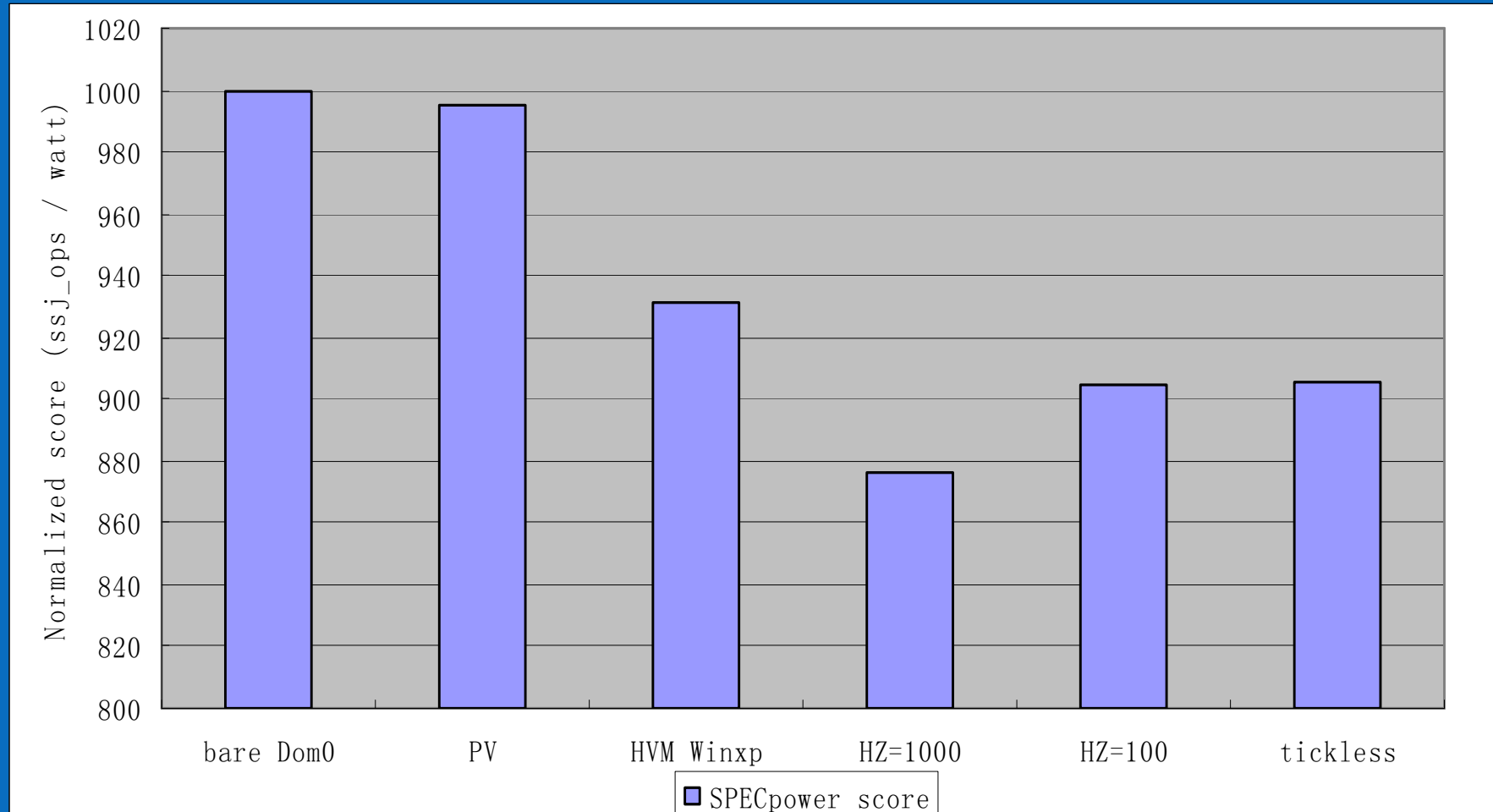
- VMM shouldn't be blamed as only reason for high power consumption!
- A 'bad' VM could eat power
 - Just like what 'bad' application could do in a native OS
 - Cause high break events (e.g. timer interrupts) with short C-state residency
- Which parts in VM could draw high power?
 - 'Bad' applications hurt just as what they could on native
 - How guest OS is implemented also matters
 - Periodic tick frequency – HZ
 - Timer usage in drivers
 - Time sub-system implementation
 - ...
- Green guest OS wins!
 - Smaller HZ, idle tickles or fully dynamic tick, range timer, etc.



Idle power consumption



Power efficiency



Software and Solutions Group



Legal Information

- INFORMATION IN THIS DOCUMENT IS PROVIDED IN CONNECTION WITH INTEL® PRODUCTS. EXCEPT AS PROVIDED IN INTEL'S TERMS AND CONDITIONS OF SALE FOR SUCH PRODUCTS, INTEL ASSUMES NO LIABILITY WHATSOEVER, AND INTEL DISCLAIMS ANY EXPRESS OR IMPLIED WARRANTY RELATING TO SALE AND/OR USE OF INTEL PRODUCTS, INCLUDING LIABILITY OR WARRANTIES RELATING TO FITNESS FOR A PARTICULAR PURPOSE, MERCHANTABILITY, OR INFRINGEMENT OF ANY PATENT, COPYRIGHT, OR OTHER INTELLECTUAL PROPERTY RIGHT.
- Intel may make changes to specifications, product descriptions, and plans at any time, without notice.
- All dates provided are subject to change without notice.
- Intel is a trademark of Intel Corporation in the U.S. and other countries.
- *Other names and brands may be claimed as the property of others.
- Copyright © 2007, Intel Corporation. All rights are protected.





Software and Solutions Group

