

Affinity-aware Dynamic Pinning Scheduling for Virtual Machines

Zhi Li

lizhi@cse.buaa.edu.cn

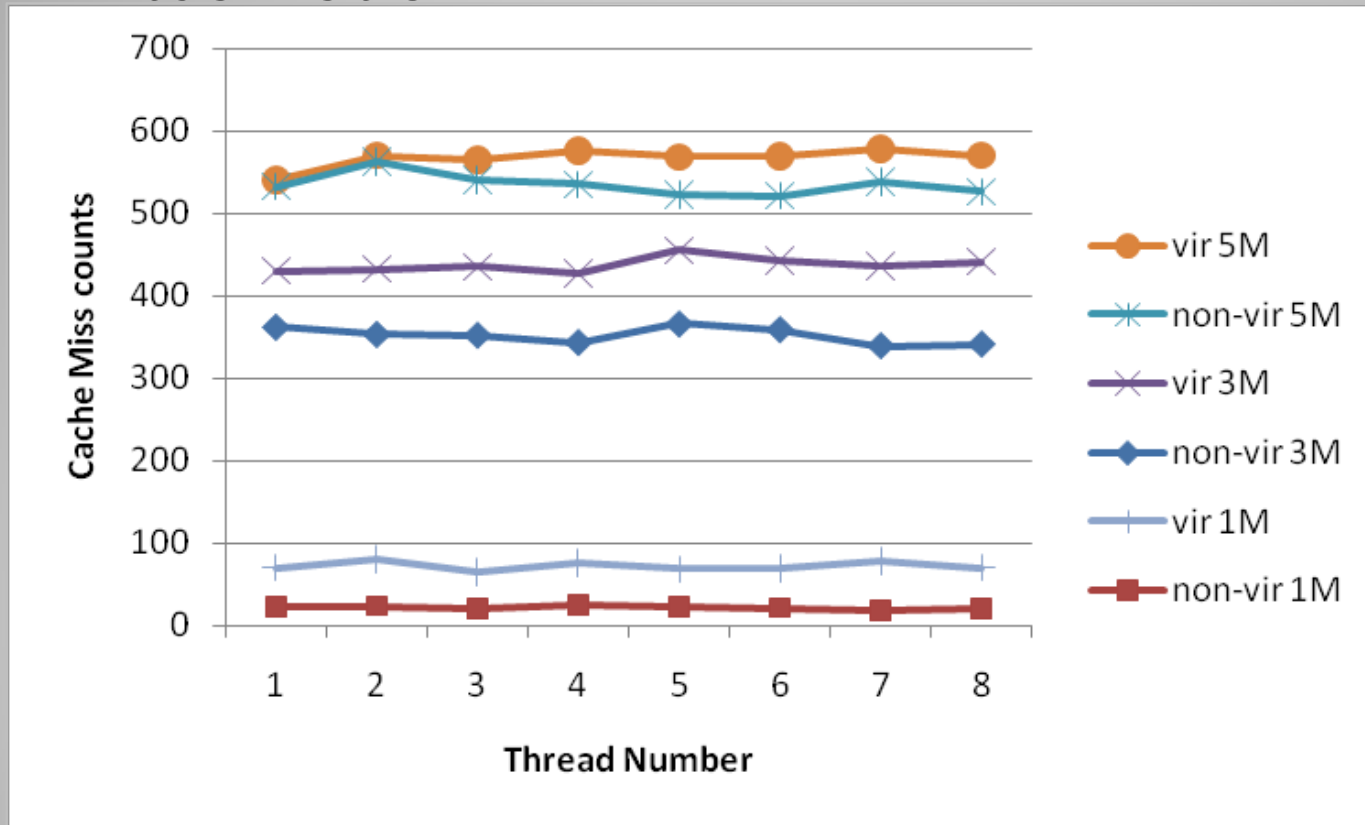
School of Computer Science, BeiHang University,
Beijing, China

Outline

- Motivation
- CPU Affinity-aware Method
- Dynamic Pinning Scheduling
- Performance Evaluation

Motivation

L2 Cache misses on both non-virtualization and virtualization



Virtualization leads to worse cache miss

Analysis of the Issue

- VCPUs from same Domain take turns to run in a same CPU runq
- Frequent migrations from this kind of VCPU

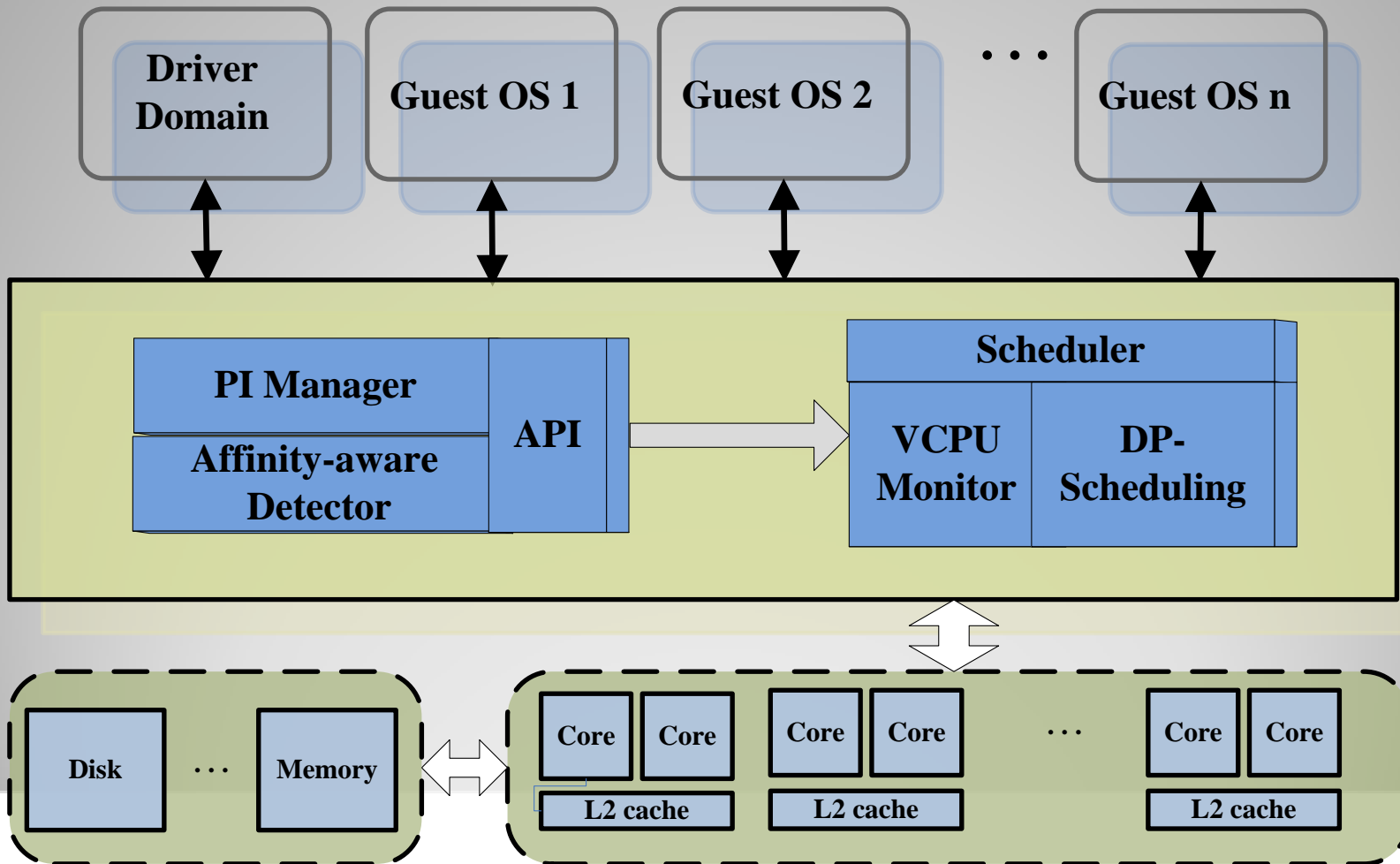
How to bridge the semantic gap between Guest OS and VMM

- Affinity-aware DP-Scheduling:
 - Affinity-aware method: providing the task affinity information to VMM.
 - DP-Scheduling: implementing that VCPU can be pinned or unpinned dynamically

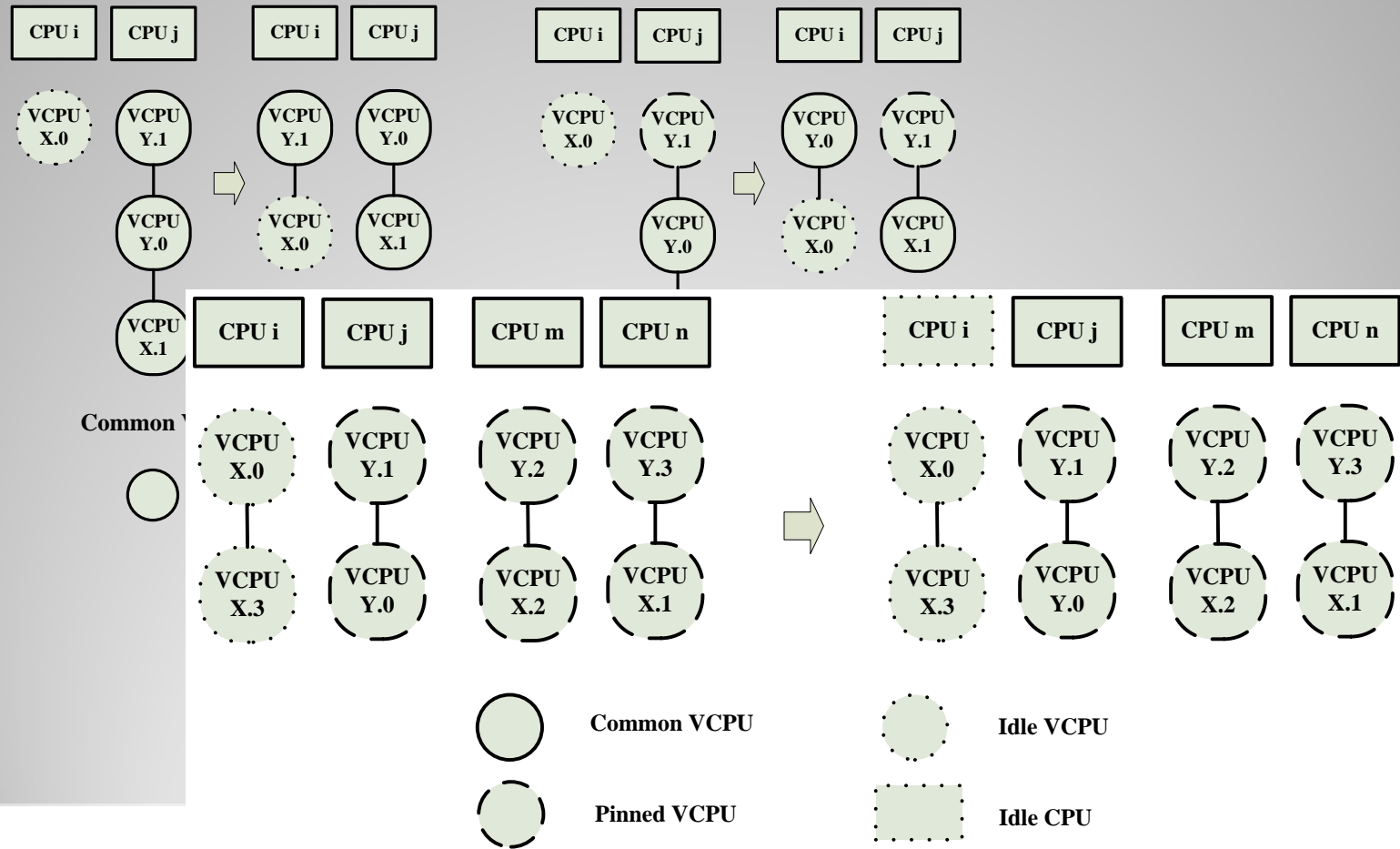
CPU Affinity-aware Method

- Timing Control
 - When CR3 is changing
- Methodology for Capture
 - Affinity Coefficient (AC)
- API
 - Provide AC for the scheduler

Dynamic Pinning Scheduling



Dynamic Pinning Scheduling



Dynamic Pinning Scheduling

- strategies :
 - (a).pin VCPU to the CPU with no pinned VCPU at this time.
 - (b).pin VCPU to the CPU with lower workload.
 - (c).pin VCPU to the local CPU if both (a) and (b) do not happen.
 - (d).do not migrate the VCPU actively when it is unpinned.
 - (e).unpin the VCPU with the lowest value of AC when the number of CPU equals.
 - (f).unpin the VCPU pinned before if it goes to the state of OVER or IDLE.

DP-Scheduling Algorithm

Let Φ stand for VCPU set from all domains, that is, $\Phi = \{VCPU_1, VCPU_2, \dots, VCPU_{|\Phi|}\}$. Let φ_p be the set of VCPUs that were pinned, whereas φ_{up} is the set of unpinning VCPUs. Hence, $\Phi = \varphi_p \cup \varphi_{up}$, and $|\Phi| = |\varphi_p| + |\varphi_{up}|$. Besides, let $Y = \{\gamma_1, \gamma_2, \dots, \gamma_{|\Phi|}\}$, where γ_i is the value of AC for VCPU_{*i*}.

For each VCPU_{*i*} $\in \Phi$ {

 CPU_{*x*} = VCPU_{*i*} \rightarrow processor;

 /*VCPU_{*i*} has been pinned by its high value of AC*/

 If ($\gamma_i \geq THRESHOLD$ and VCPU_{*i*} $\in \varphi_p$) Return;

 /*VCPU_{*i*} with its high value of AC was not pinned */

 Else If ($\gamma_i \geq THRESHOLD$ and VCPU_{*i*} $\in \varphi_{up}$

 and VCPU_{*i*} $\rightarrow pri > OVER$){

 peer_cpu = select_peer(VCPU_{*i*});

 set_unmigrateable(peer_cpu, vcpu);

 del(φ_{up} , VCPU_{*i*});

 add(φ_p , VCPU_{*i*});

 }

 /*Pinned VCPU_{*i*} contains low value of AC currently
 according to strategy (d)/

 Else If ($\gamma_i < THRESHOLD$ and VCPU_{*i*} $\in \varphi_p$){

 set_migrateable(CPU_{*x*}, VCPU_{*i*});

 del(φ_p , VCPU_{*i*});

 add(φ_{up} , VCPU_{*i*});

 }

 /*The total number of VCPUs should not over

 CPU number -1/

 If ($|\varphi_p| == |CPU|$)

 set_migrateable_by_ac(φ_p , φ_{up});

 }

Performance Evaluation

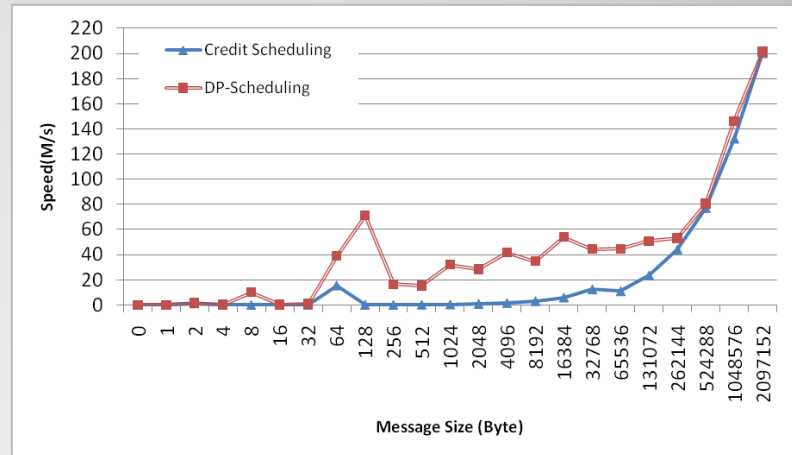
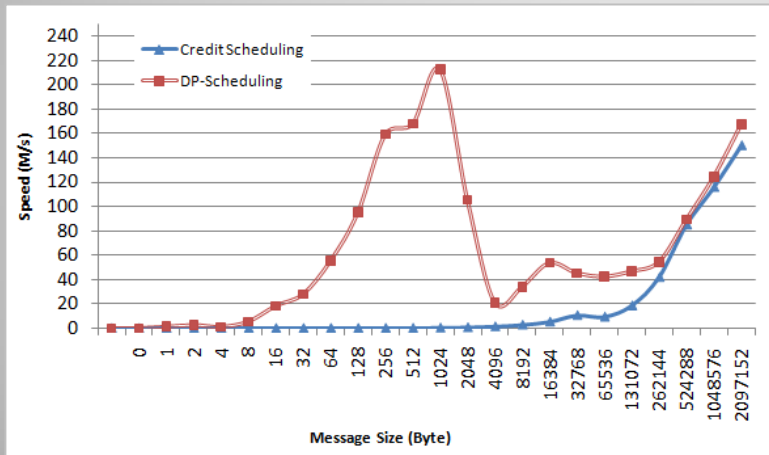
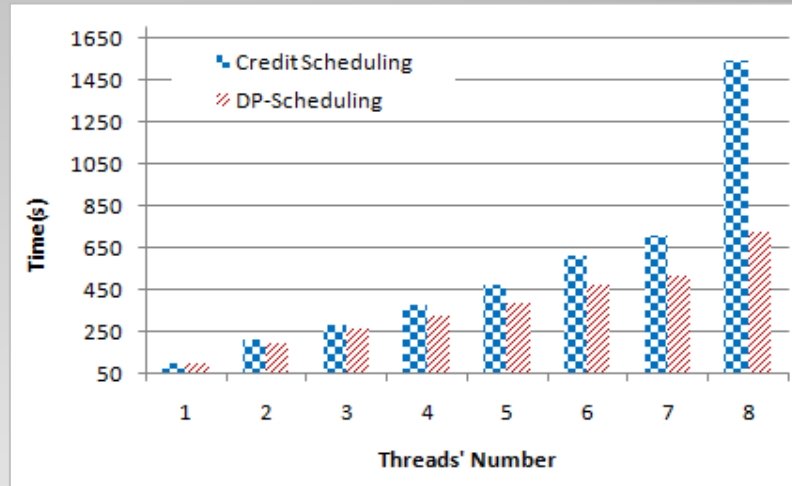
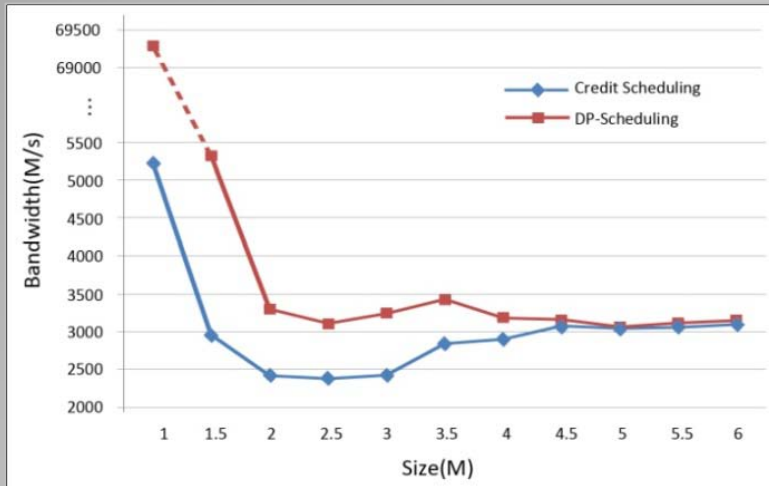
Platform

- Xeon 5405 (two quad-core)
- L2 cache: 6M
- RAM: 4G
- Xen: 3.4.3

Benchmark

Benchmark Category	Code Name	Variable	Measurement
HPCC	STREAM	Array size	Memory Bandwidth
EPCC	OpenMP Micro-benchmark suite	Thread Number	Time
IMB	Sendrecv	Message Size	Transfer Speed
	Exchange		

Performance Evaluation



- Conclusion

DP-Scheduling outperforms Credit scheduling for kinds of CPU-bound tasks, without interfering the load balance

Thank You !