

A Case for Scaling Applications to Many-core with OS Clustering

(reviewed by Oleg legorov and Jander Nascimento)

University Joseph Fourier

December 11, 2011

Process
management

Sharing file

The prototype

Evaluation

Our evaluation

Roadmap

Process management

Sharing file

The prototype

Evaluation

Our evaluation

A Case for Scaling
Applications to
Many-core with
OS Clustering

(reviewed by Oleg
legorov and Jander
Nascimento)

Process
management

Sharing file

The prototype

Evaluation

Our evaluation

Scaling Operating Systems For Many Cores

Nowaday's tendencies:

- ▶ the number of CPU cores in a system increases;
- ▶ parallel applications use fine-grained components;

A Case for Scaling
Applications to
Many-core with
OS Clustering

(reviewed by Oleg
legorov and Jander
Nascimento)

Process
management

Sharing file

The prototype

Evaluation

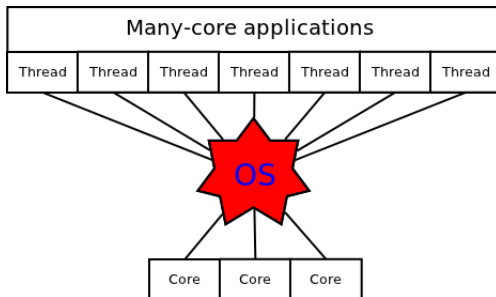
Our evaluation

Scaling Operating Systems For Many Cores

Nowaday's tendencies:

- ▶ the number of CPU cores in a system increases;
- ▶ parallel applications use fine-grained components;

OS becomes a bottleneck!



A Case for Scaling
Applications to
Many-core with
OS Clustering

(reviewed by Oleg
Igorov and Jander
Nascimento)

Process
management

Sharing file

The prototype

Evaluation

Our evaluation

Scaling Operating Systems For Many Cores

A Case for Scaling
Applications to
Many-core with
OS Clustering

(reviewed by Oleg
legorov and Jander
Nascimento)

Two main approaches:

1. designing new OSes

- ▶ distribute independent kernels on multiple cores;
- ▶ reduce resource sharing & improve data locality;

2. refining commodity kernels

- ▶ design new data structures to reduce the contention on shared data structures;
- ▶ *interesting fact*: Linux kernel can efficiently scale to 48 cores!

Process
management

Sharing file

The prototype

Evaluation

Our evaluation

OS clustering approach

New approach proposed by this paper: **OS clustering**

- ▶ cluster multiple OSes atop a VMM;
- ▶ several OSes can serve one application;

Motivation

1. commodity kernels scale well with small # of CPU cores;
2. one VMM can efficiently consolidate multiple OSes;

A Case for Scaling
Applications to
Many-core with
OS Clustering

(reviewed by Oleg
legorov and Jander
Nascimento)

Process
management

Sharing file

The prototype

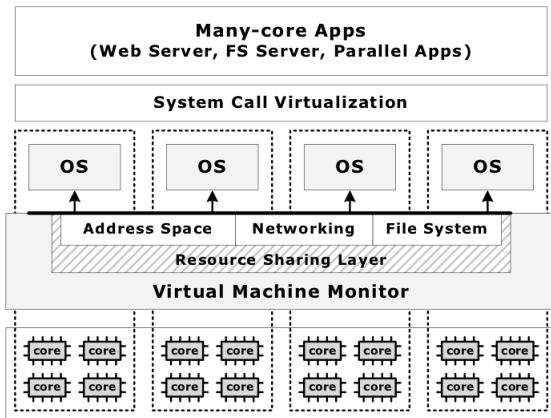
Evaluation

Our evaluation

OS clustering approach

New approach proposed by this paper: **OS clustering**

- ▶ cluster multiple OSES atop a VMM;
- ▶ several OSES can serve one application;



A Case for Scaling Applications to Many-core with OS Clustering

(reviewed by Oleg legorov and Jander Nascimento)

Process management

Sharing file

The prototype

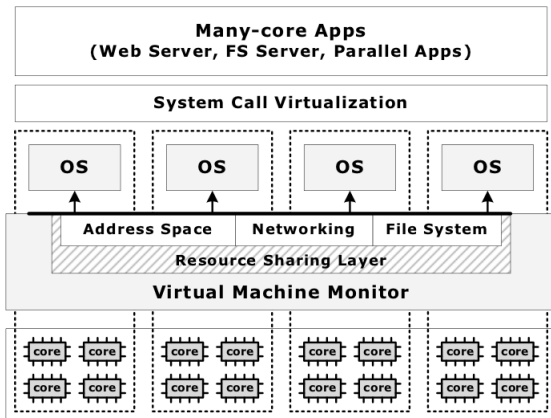
Evaluation

Our evaluation

OS clustering approach

Results

- ▶ single OS manages fewer cores;
- ▶ resource contention is mitigated.



A Case for Scaling Applications to Many-core with OS Clustering

(reviewed by Oleg Ilegorov and Jander Nascimento)

Process management

Sharing file

The prototype

Evaluation

Our evaluation

What's next?

The implemented system was called **Cerberus**

- ▶ overview of Cerberus;
- ▶ more implementation details;
- ▶ evaluation results

A Case for Scaling
Applications to
Many-core with
OS Clustering

(reviewed by Oleg
legorov and Jander
Nascimento)

Process
management

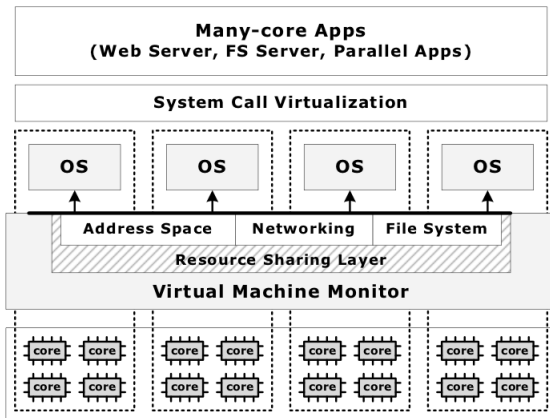
Sharing file

The prototype

Evaluation

Our evaluation

Overview of Cerberus architecture



A Case for Scaling
Applications to
Many-core with
OS Clustering

(reviewed by Oleg
legorov and Jander
Nascimento)

Process
management

Sharing file

The prototype

Evaluation

Our evaluation

System call virtualization

A Case for Scaling
Applications to
Many-core with
OS Clustering

(reviewed by Oleg
legorov and Jander
Nascimento)

Process
management

Sharing file

The prototype

Evaluation

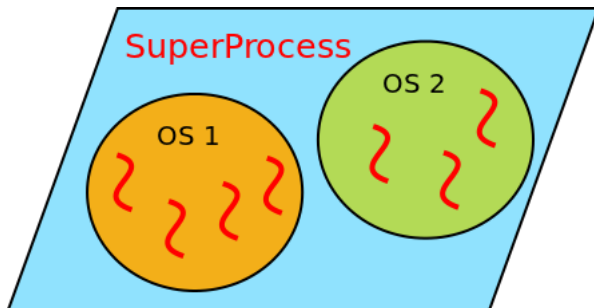
Our evaluation

Why it is needed?

- ▶ all syscalls should be intercepted;
- ▶ allow processes/threads of one application to run on several OSes;
- ▶ transparent to the user (POSIX API)

System call virtualization

The notion of **SuperProcess**



A Case for Scaling
Applications to
Many-core with
OS Clustering

(reviewed by Oleg
legorov and Jander
Nascimento)

Process
management

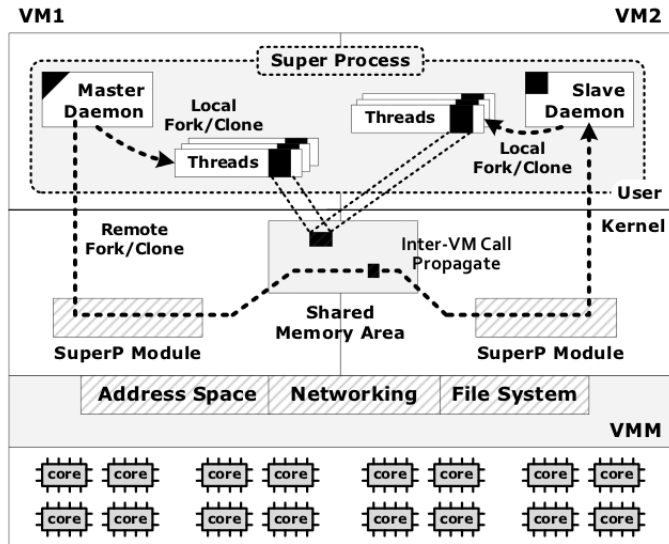
Sharing file

The prototype

Evaluation

Our evaluation

System call virtualization



A Case for Scaling Applications to Many-core with OS Clustering

(reviewed by Oleg Ilegorov and Jander Nascimento)

Process management

Sharing file

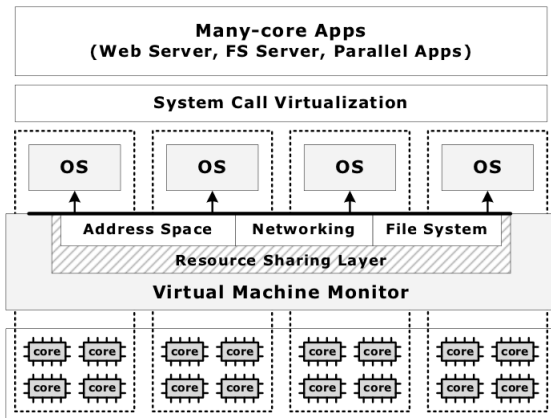
The prototype

Evaluation

Our evaluation

Resource Sharing

- ▶ processes/threads of one application must share some resources;
- ▶ VMMs are designed to isolate VMs from each other;
- ▶ *resource sharing layer* must modify VMM's code



A Case for Scaling
Applications to
Many-core with
OS Clustering

(reviewed by Oleg
legorov and Jander
Nascimento)

Process
management

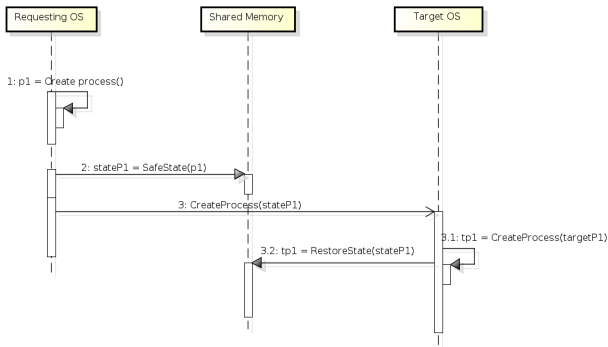
Sharing file

The prototype

Evaluation

Our evaluation

Remote process



- ▶ Create local process
- ▶ Saves the state of the process
- ▶ Request target OS to create a process
- ▶ Target OS restores the state

A Case for Scaling
Applications to
Many-core with
OS Clustering

(reviewed by Oleg
legorov and Jander
Nascimento)

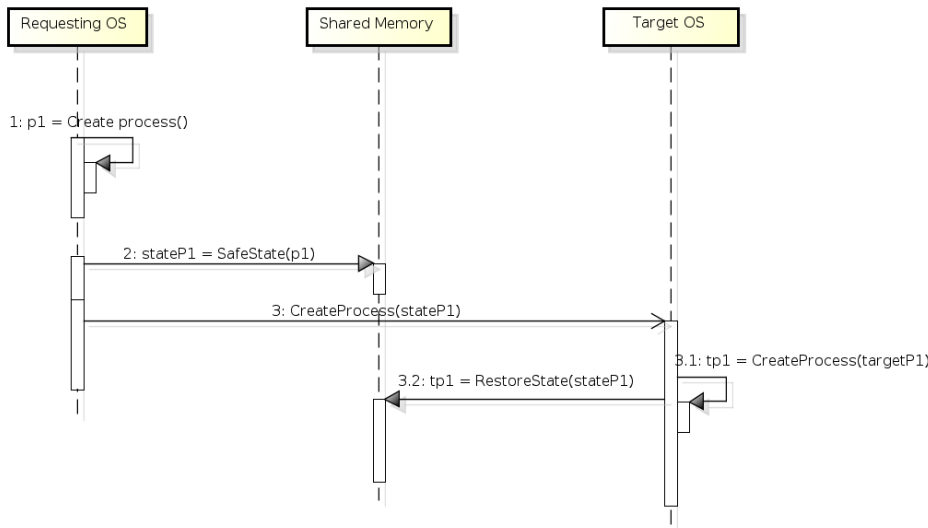
Process
management

Sharing file

The prototype

Evaluation

Our evaluation



Performance techniques

A Case for Scaling
Applications to
Many-core with
OS Clustering

(reviewed by Oleg
legorov and Jander
Nascimento)

Process
management

Sharing file

The prototype

Evaluation

Our evaluation

- ▶ Only certain levels of the page tables are shared
- ▶ Use domains to address page faults
- ▶ Avoid conflicting by enforcing update memory state

Shared Data Segment

A Case for Scaling
Applications to
Many-core with
OS Clustering

(reviewed by Oleg
legorov and Jander
Nascimento)

Process
management

Sharing file

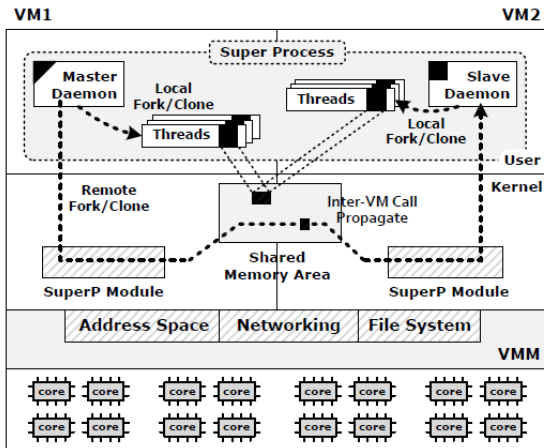
The prototype

Evaluation

Our evaluation

- ▶ Not common for multiprocess application
- ▶ Very common for multi-threaded

Cerberus arch



A Case for Scaling
Applications to
Many-core with
OS Clustering

(reviewed by Oleg
Ilegorov and Jander
Nascimento)

Process
management

Sharing file

The prototype

Evaluation

Our evaluation

Syscall virtualization

A Case for Scaling
Applications to
Many-core with
OS Clustering

(reviewed by Oleg
legorov and Jander
Nascimento)

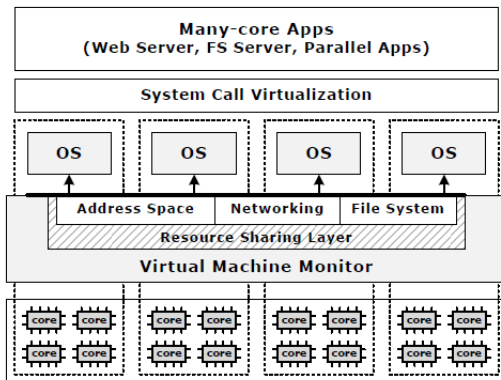
Process
management

Sharing file

The prototype

Evaluation

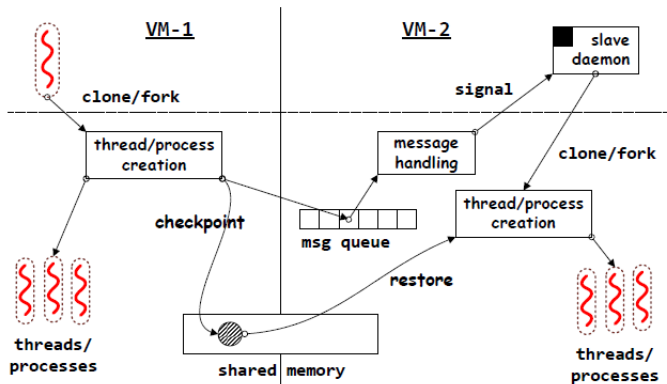
Our evaluation



Fork and Clone

A Case for Scaling
Applications to
Many-core with
OS Clustering

(reviewed by Oleg
legorov and Jander
Nascimento)



Process
management

Sharing file

The prototype

Evaluation

Our evaluation

Common approach

A Case for Scaling
Applications to
Many-core with
OS Clustering

(reviewed by Oleg
legorov and Jander
Nascimento)

Use existent network file systems:

- ▶ NFS

Issues?

- ▶ Centralized
- ▶ Rather slow

Solution?

Process
management

Sharing file

The prototype

Evaluation

Our evaluation

Sharing file

A Case for Scaling
Applications to
Many-core with
OS Clustering

(reviewed by Oleg
legorov and Jander
Nascimento)

Process
management

Sharing file

The prototype

Evaluation

Our evaluation

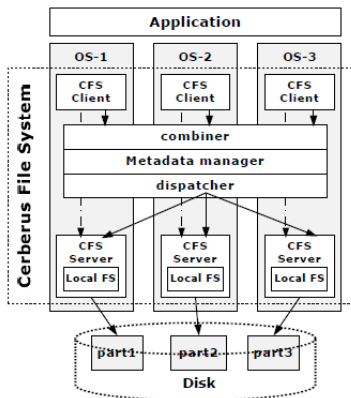
Question

How share data files efficiently?

The answer might be:

- ▶ Hybrid solution: local and remote file access.
- ▶ Basic security scheme: public or private file.

CFS arch



A Case for Scaling
Applications to
Many-core with
OS Clustering

(reviewed by Oleg
legorov and Jander
Nascimento)

Process
management

Sharing file

The prototype

Evaluation

Our evaluation

Prototype

A Case for Scaling
Applications to
Many-core with
OS Clustering

(reviewed by Oleg
legorov and Jander
Nascimento)

Tech overview:

- ▶ based on Xen
- ▶ Uses compare and swap to look for gaps
- ▶ using shadow mode page table
- ▶ re-implement a subset of POSIX calls
- ▶ Using P2M

Hardware overview

- ▶ AMD 48 cores
- ▶ 4 network cards

Process
management

Sharing file

The prototype

Evaluation

Our evaluation

Some negative points of the paper:

- ▶ Do not show the performance of CFS in operations like: writing and seeking
- ▶ Do not explain how solve the problem of centralization for load balance
- ▶ Do not show the numbers for acceptable performance
- ▶ Code is not available