

History-Based Harvesting of Spare Cycles and Storage in Large-Scale Datacenters

Yunqi Zhang, George Prekas, Giovanni Matteo Fumarola,
Marcus Fontoura, Inigo Goiri, Ricardo Bianchini.
University of Michigan and Microsoft Research

12th USENIX symposium on Operating System Design and Implementation

Presenter: Yi-Ning Chang

July 27, 2017

Outline

- 1 Introduction
- 2 Behavior Patterns
- 3 Co-location Techniques
- 4 Experiment
- 5 Conclusion

- 1 Introduction
- 2 Behavior Patterns
- 3 Co-location Techniques
- 4 Experiment
- 5 Conclusion

Introduction

- Overprovision resources in datacenters
 - Low tail latency requirement
 - Provisioned for peak load
 - Unexpected load spikes and failures
- A way to increase utilization and reduce costs in datacenters is to co-locate their latency-critical services and batch workloads.
- Harvest spare compute cycles and storage space for co-location purpose.

Challenges

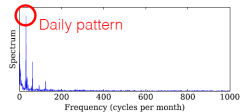
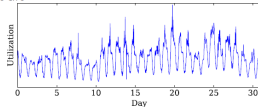
- Interactive services “own” the servers
- Resource availability
 - Interactive services performance
 - Resource availability dynamics - task killing
- Data storage co-location
 - Data availability
 - Data durability

- 1 Introduction
- 2 Behavior Patterns**
- 3 Co-location Techniques
- 4 Experiment
- 5 Conclusion

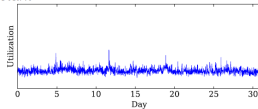
Resource Utilization

- Identify three main classes of primary tenants.

Periodic

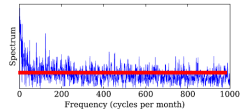


Constant

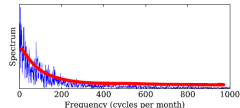
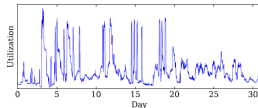


Fourier
Transform

A diagram showing a 3D box with a red outline, representing the Fourier Transform process. The box is tilted, and a red line connects it to the spectrum graph on the right.

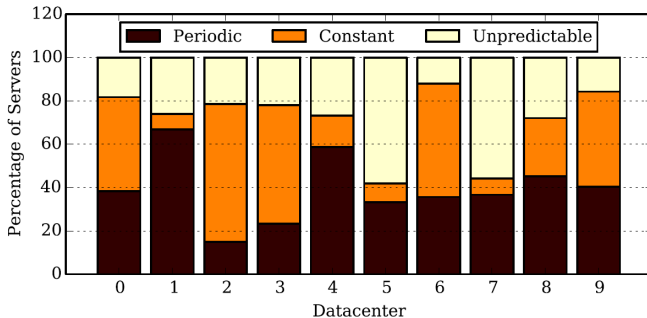


Unpredictable



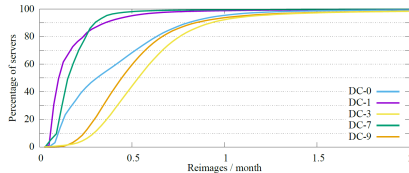
Resource Utilization

- Percentages of servers per class.

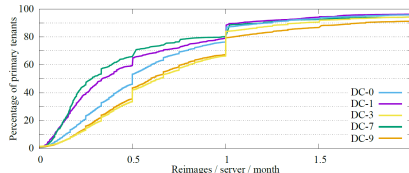


Disk Reimaging

- Per-server number of reimages in three years.



- Per-tenant number of reimages in three years.



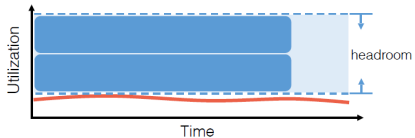
- 1 Introduction
- 2 Behavior Patterns
- 3 Co-location Techniques**
- 4 Experiment
- 5 Conclusion

History-Based Task Scheduling

Long Jobs



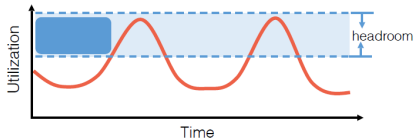
- Constant
- $1 - \text{MAX}(\text{Peak}, \text{Current})$



Medium Jobs



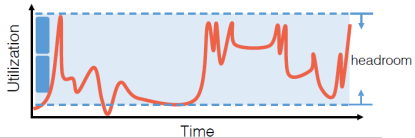
- Periodic
- $1 - \text{MAX}(\text{Average}, \text{Current})$



Short Jobs



- Unpredictable
- $1 - \text{Current}$



Task Scheduling

- Fast Fourier Transform (FFT)
 - Get 3 patterns.
- Clustering algorithm
 - K-means algorithm
 - Average and peak utilizations
- Class selection algorithm

Class Selection Algorithm

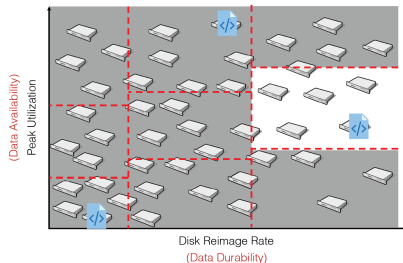
Algorithm 1 Class selection algorithm.

```

1: Given: Classes  $C$ , Headroom( $type, c$ ), Ranking Weights  $W$ 
2: function SCHEDULE(Batch job  $J$ )
3:    $J.type =$  Length (short, medium, or long) from its last run
4:    $J.req =$  Max amount of concurrent resources from DAG
5:   for each  $c \in C$  do
6:      $c.weightedroom = \text{Headroom}(J.type, c) \times W[J.type, c.class]$ 
7:   end for
8:    $F = \{\forall c \in C \mid \text{Headroom}(J.type, c) \geq J.req\}$ 
9:   if  $F \neq \emptyset$  then
10:    Pick 1 class  $c \in F$  probabilistically  $\propto c.weightedroom$ 
11:    return  $\{c\}$ 
12:   else if Job  $J$  can fit in multiple classes combined then
13:    Pick  $\{c_0, \dots, c_k\} \subseteq C$  probabilistically  $\propto c.weightedroom$ 
14:    return  $\{c_0, \dots, c_k\}$ 
15:   else
16:    Do not pick classes
17:    return  $\{\emptyset\}$ 
18:   end if
19: end function
  
```

History-Based Data Placement

- Data availability
 - Diverse in utilization pattern.
- Data durability
 - Diverse in reimaging pattern.



Replica Placement Algorithm

Algorithm 2 Replica placement algorithm.

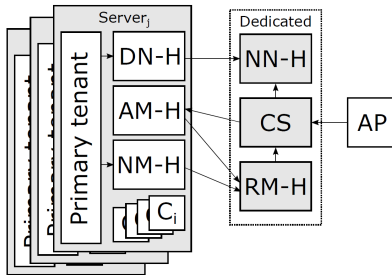
```

1: Given: Storage space available in each server, Primary reimaging
2:   stats, Primary peak CPU util stats, Desired replication  $R$ 
3: function PLACE REPLICAS(Block  $B$ )
4:   Cluster primary tenants wrt reimaging and peak CPU util
5:   into 9 classes, each with the same total space
6:   Select the class of the server creating the block
7:   Select the server creating the block for one replica
8:   for  $r = 2; r \leq R; r = r + 1$  do
9:     Select the next class randomly under two constraints:
10:    No class in the same row has been picked
11:    No class in the same column has been picked
12:    Pick a random primary tenant of this class as long as
13:    its environment has not received a replica
14:    Pick a server in this primary tenant for the next replica
15:    if  $(r \bmod 3) == 0$  then
16:      Forget rows and columns that have been selected so far
17:    end if
18:  end for
19: end function
  
```

- 1 Introduction
- 2 Behavior Patterns
- 3 Co-location Techniques
- 4 Experiment**
- 5 Conclusion

System Implementation

- Overview of YARN-H, Tez-H, and HDFS-H in a co-location scenario.

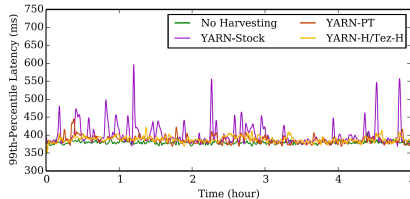


Evaluation

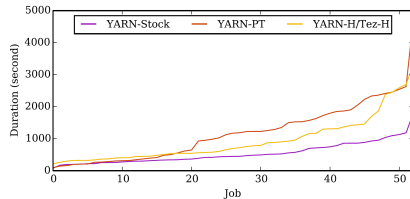
- Environment:
 - 102-server cluster
- Primary tenant (interactive service):
 - Apache Lucene search engine with utilization trace
- Batch task:
 - TPC-DS benchmark

Batch Task Scheduling

- Interactive services performance

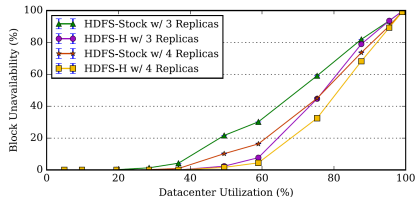


- Job duration - reducing task killing

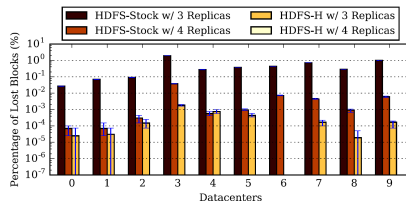


Data Placement

- Data availability



- Data durability



- 1 Introduction
- 2 Behavior Patterns
- 3 Co-location Techniques
- 4 Experiment
- 5 Conclusion**

Conclusion

- Embody knowledge of the existing primary workloads, and leverage historical utilization.
- Improve batch job performance while protecting primary workloads.
- Eliminate data loss and unavailability in many scenarios.

Future Work

- Paper
 - Isolation and security in public cloud.
- Project
 - Network utilization.
 - Constraint of Ethernet bandwidth.