# Placement of Virtual Containers on NUMA Systems

Justin Funston\*, Maxime Lorrillere<sup>†</sup>, and Alexandra Fedorova, University of British Columbia Baptiste Lepers, EPFL David Vengerov and Jean-Pierre Lozi, Oracle Labs Vivien Quéma, IMAG

> \* Currently Huawei R&D † Currently Arista Networks

## Motivation



3% of global electricity usage 86% used by servers+cooling

Ö







VS.

- Half as many servers!
- Half as much energy!
- Half as much infrastructure!
- Performance?



## Background – "Module"



# Background – NUMA Node



# Background – Interconnect Topology



# Assumptions

- The number of vCPUs a container uses is fixed
- Max of one vCPU per core
- Containers packed together should not interfere with each other – containers will not share NUMA nodes





#### **Placements**







### Motivation – ua.B (scientific simulation)



## Motivation – Spark Pagerank



## Workload Placement Overview



## **Abstract Machine Model**

#### How to represent placements?

- Too many (e.g. trillions for 16 threads on 64 cores) for a naïve approach
- Need to exploit symmetry

# Scheduling Concerns



# Scheduling Concern Example – L3



# Scheduling Concern Example – L2



## Abstract Machine Model – Important Placements

- Scheduling concerns + Important placements:
- ~10<sup>14</sup> Placements  $\rightarrow$  12 Placements
  - 16 threads on 64 cores
- Can train on all important placements

# Performance Prediction Model – Features/Inputs

- Hardware Performance Events (HPEs)
  - Standard in existing work
  - Surprisingly, poor predictive performance!
  - Excessive training time
- Performance Measurements

# Performance Predictions, Online Inference



## Performance Predictions – ua.B (scientific sim.)



## Performance Predictions – Spark Pagerank



## Performance Predictions – is.D (sort)



## Machine Learning – Prediction Accuracy



# Conclusion

- Data centers: 3% world's electricity (86% servers+cooling)
- Packing containers onto servers
  - Increase efficiency
  - Maintain performance goals
- **2-4**× better utilization in many cases!

- Abstract machine model
- Performance prediction model





## Workload Placement – Related Work

	Predicts Performance	Multiple Hardware Resources	Easily Adapted	Deployable Online
Our Solution	$\checkmark$	$\checkmark$	$\checkmark$	$\checkmark$
Pandia (EuroSys '17)	$\checkmark$	$\checkmark$	×	×
SMiTe (Micro '14)	$\checkmark$	$\checkmark$	×	$\checkmark$
Bubble-Flux (ISCA '13)	$\checkmark$	$\checkmark$	×	$\checkmark$
Asymsched (ATC '15)	X	X	$\checkmark$	$\checkmark$
DINO (ASPLOS '10)	X	X	$\checkmark$	$\checkmark$
<b>Thread Clustering</b> (EuroSys '07)	X	X	$\checkmark$	$\checkmark$