



# A ~~Green~~ (*Less brown*) Cloud

## How to use less resources?

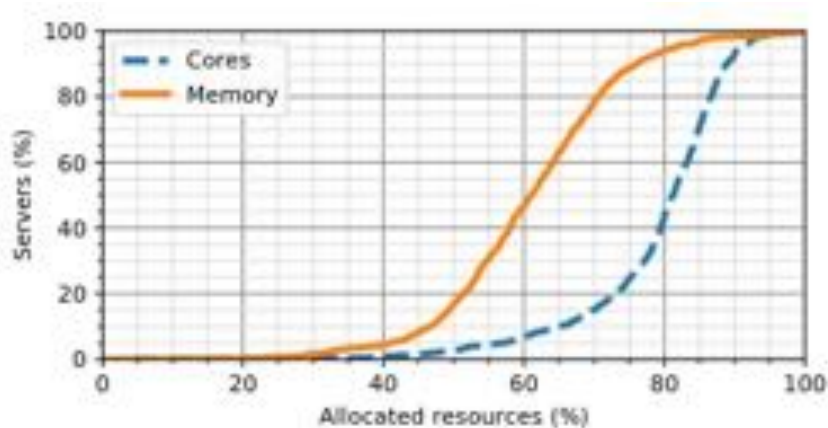
Pierre Jacquet

Under the supervision of Romain Rouvoy, Thomas Ledoux

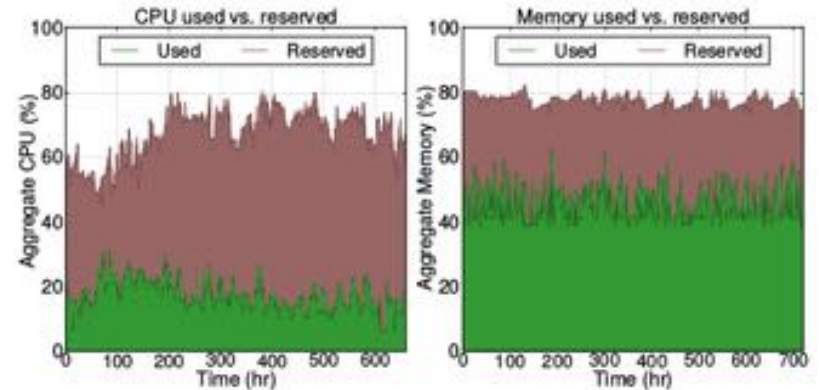


# State-of-the-Art of Cloud Resource Utilization

## Cloud infrastructure are under-used



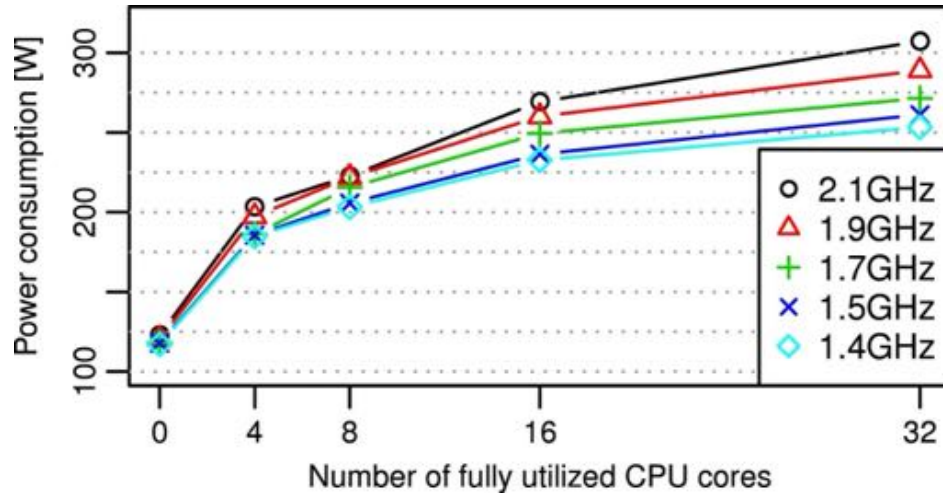
Not all resources are **allocated** [1]



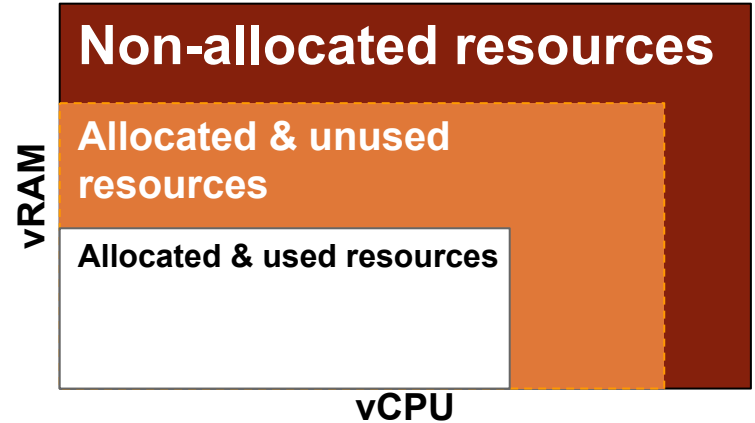
Not all allocated resources are **used** [2]

# State-of-the-Art of Cloud Resource Utilization

## Cloud infrastructure are under-used



Impacts **server consumption** [1]



Impacts the number of **provisioned servers**

# State-of-the-Art of Cloud Resource Utilization

- A problem addressed using different (complementary) approaches:
  - Fill with **heterogeneous workloads**, e.g. *Batch, FaaS, HarvestVM*
  - Fill with **homogeneous workloads** using oversubscription (n:1) ←

# From adopting oversubscription techniques

- Ratios at the cluster level
  - Statically defined
- Ratios at the server level
  - Statically defined
  - Dynamically defined
    - Generic formulas on VM percentile or host standard deviation

# To extending oversubscription techniques

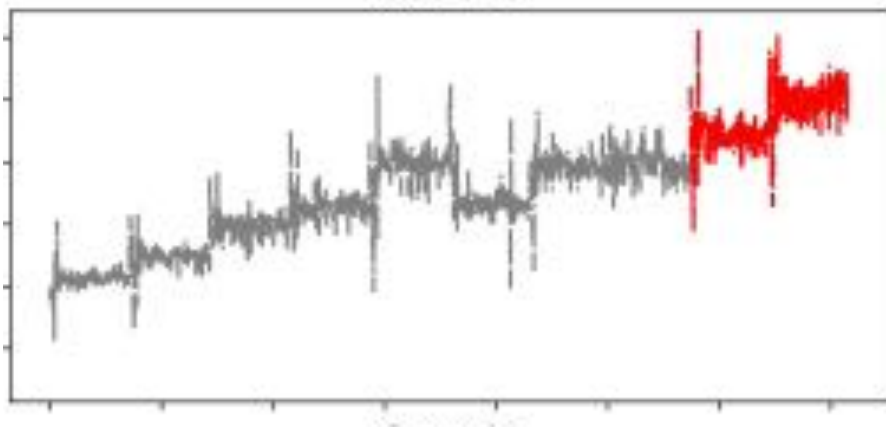
- Ratios at the cluster level
  - Statically defined
- Ratios at the server level
  - Statically defined
  - Dynamically defined
    - Generic formulas on VM percentile or host standard deviation
    - Per server formula customization
- Ratios at the VM level
- Ratios at the resource level (vCPU)

# Contribution 1: **ScroogeVM**

- Ratios at the cluster level
  - Statically defined
- Ratios at the server level
  - Statically defined
  - Dynamically defined
    - Generic formulas on VM percentile or host standard deviation
    - **Per server formula customization**
- Ratios at the VM level
- Ratios at the resource level (vCPU)

# Contribution 1: ScroogeVM

- **Adjust resources oversubscription ratios according to quiescent state**



adjust optimistic degree  
of next peak prediction

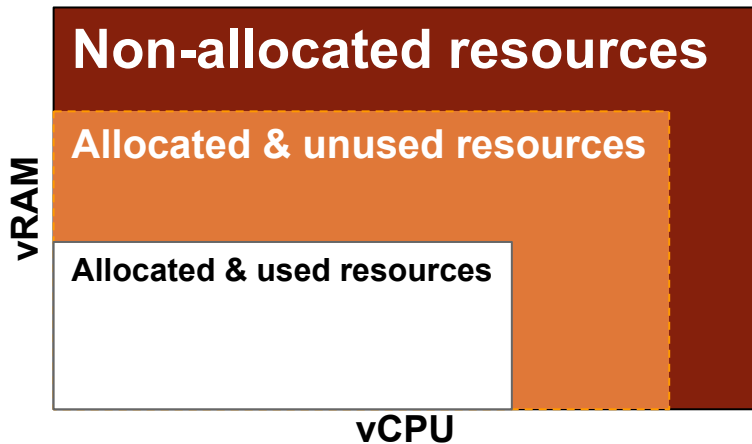
Being closer to the usage improves packing



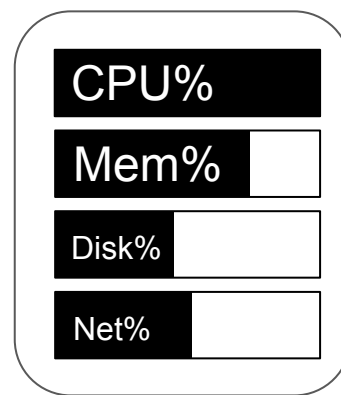
# Contribution 2: **SlackVM**

- Ratios at the cluster level
  - Statically defined
- Ratios at the server level
  - Statically defined
  - Dynamically defined
    - Generic formulas on VM percentile or host standard deviation
    - **ScroogeVM**: per server formula customization
- **Ratios at the VM level**
- Ratios at the resource level (vCPU)

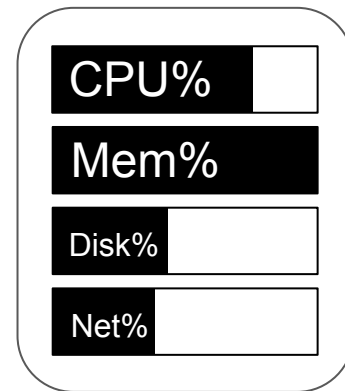
# Contribution 2: **SlackVM**



Oversubscription reduces  
unused resources...



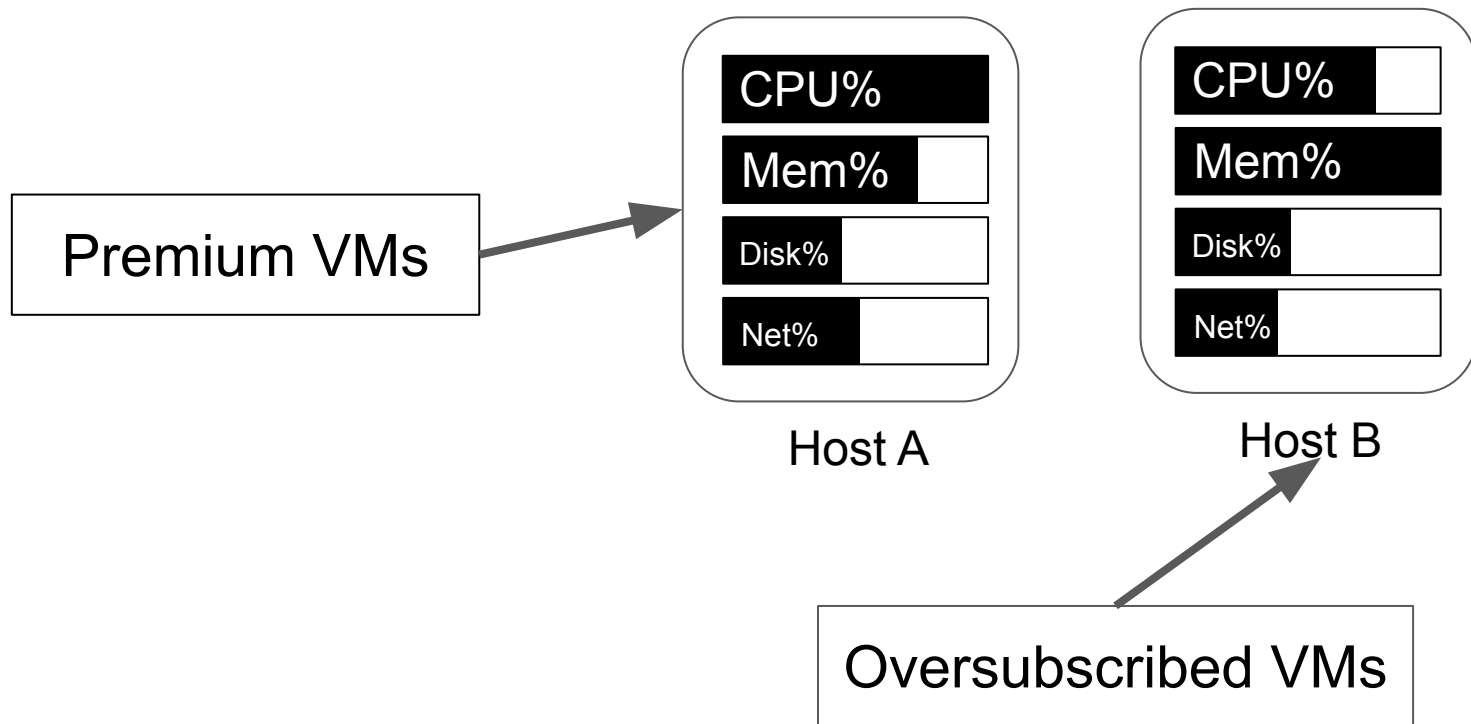
Host A



Host B

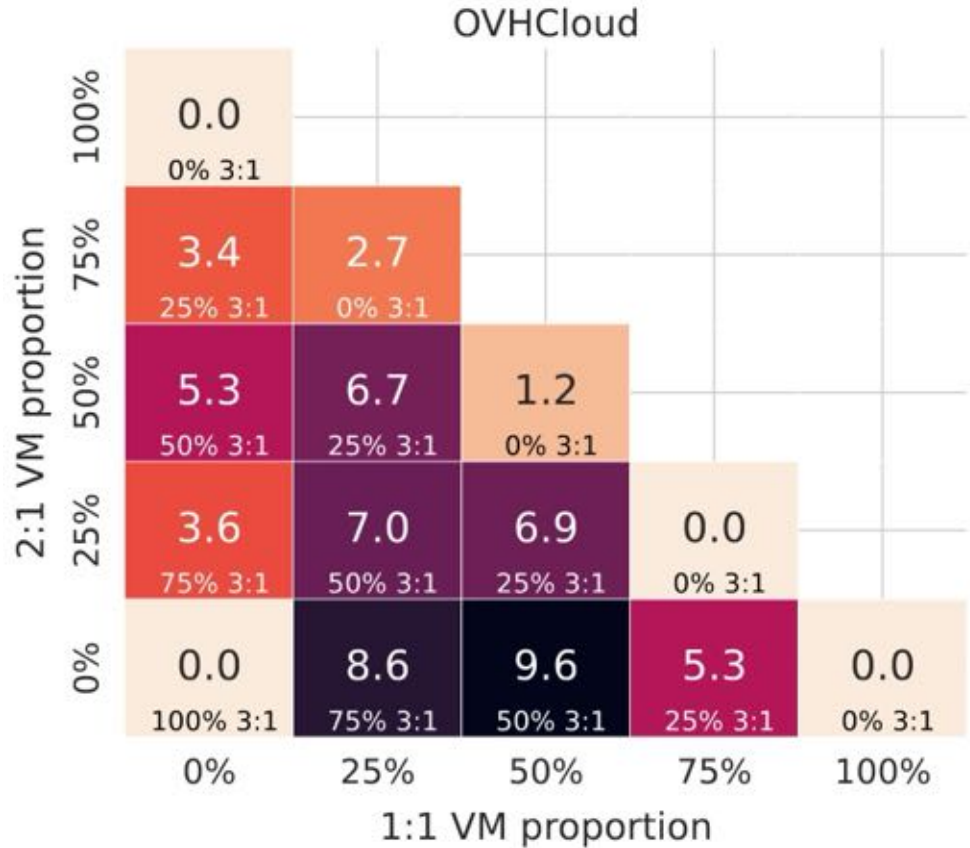
...But may also reduce  
unallocated resources

# Contribution 2: **SlackVM**



# Contribution 2: **SlackVM**

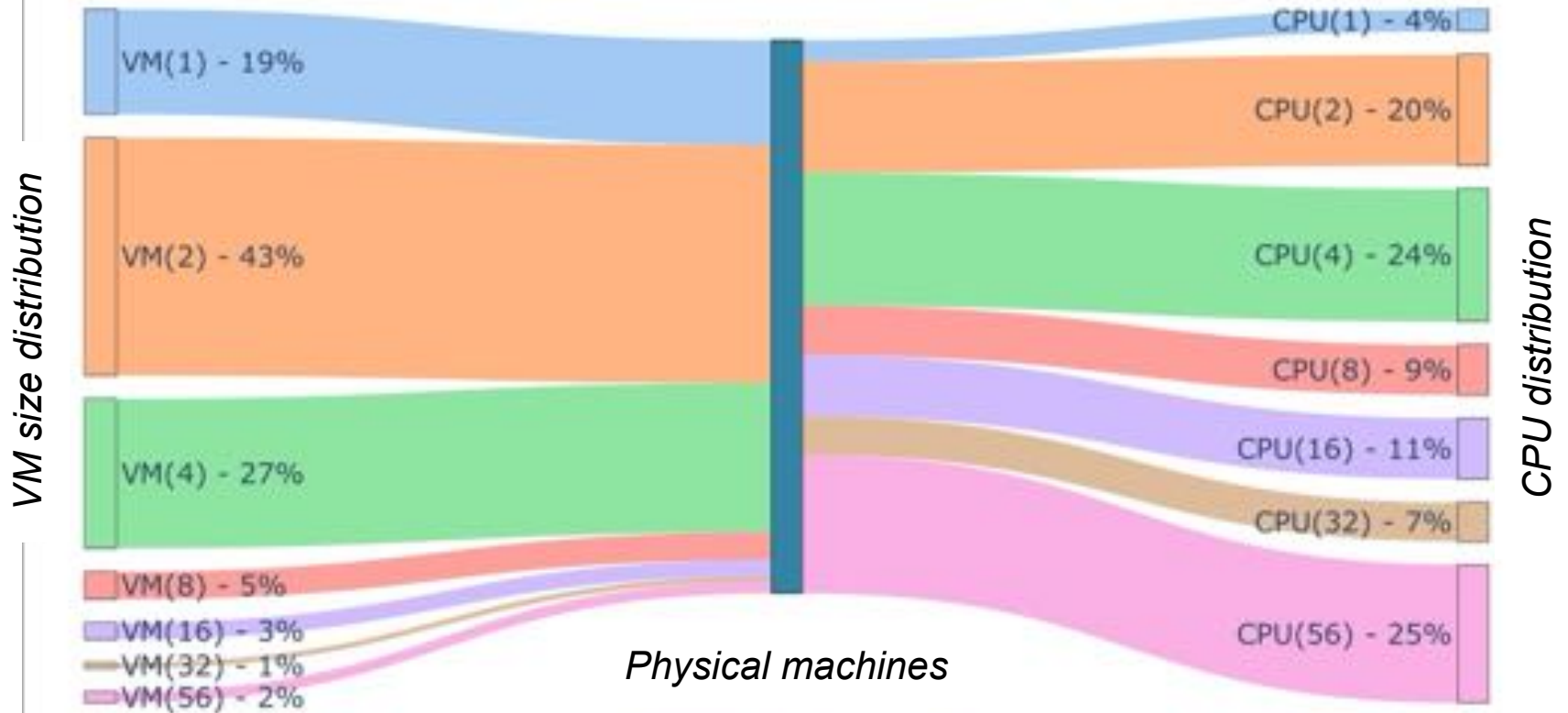
Combining oversubscription levels can **save up to 10% physical machines**



# Contribution 3: **SweetspotVM**

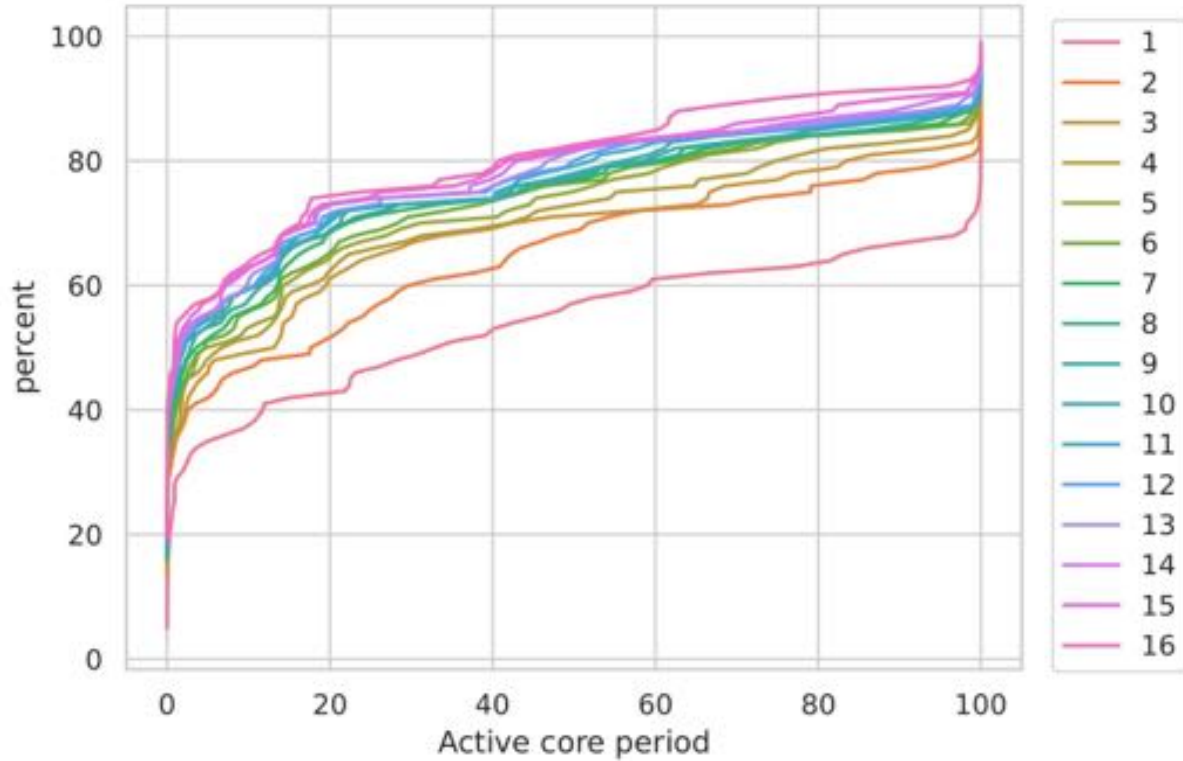
- Ratios at the cluster level
  - Statically defined
- Ratios at the server level
  - Statically defined
  - Dynamically defined
    - Generic formulas on VM percentile or host standard deviation
    - **ScroogeVM**: per server formula customization
- **SlackVM**: ratios at the VM level
- **Ratios at the resource level (vCPU)**

# Contribution 3: SweetspotVM



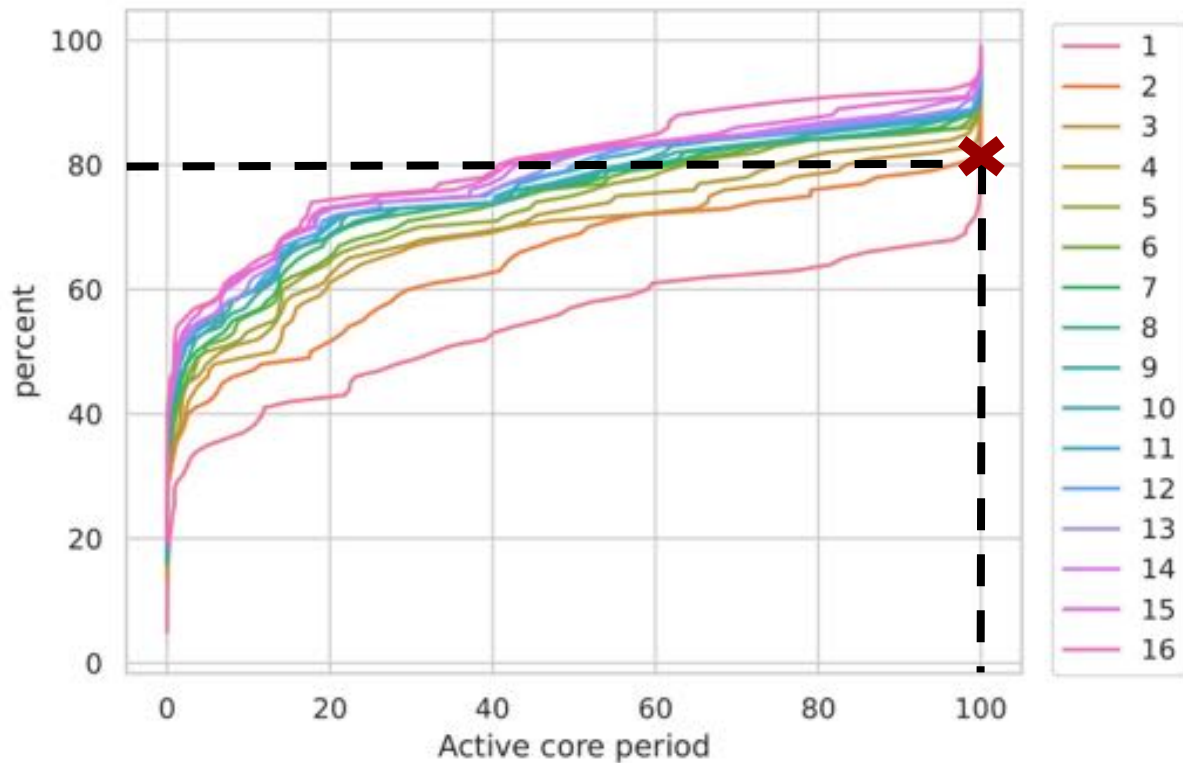
Most of CPUs are provisioned by a **small subset of VMs**

# Contribution 3: **SweetspotVM**



Large VMs hosted by OVHCloud **do not use all vCPUs equally**

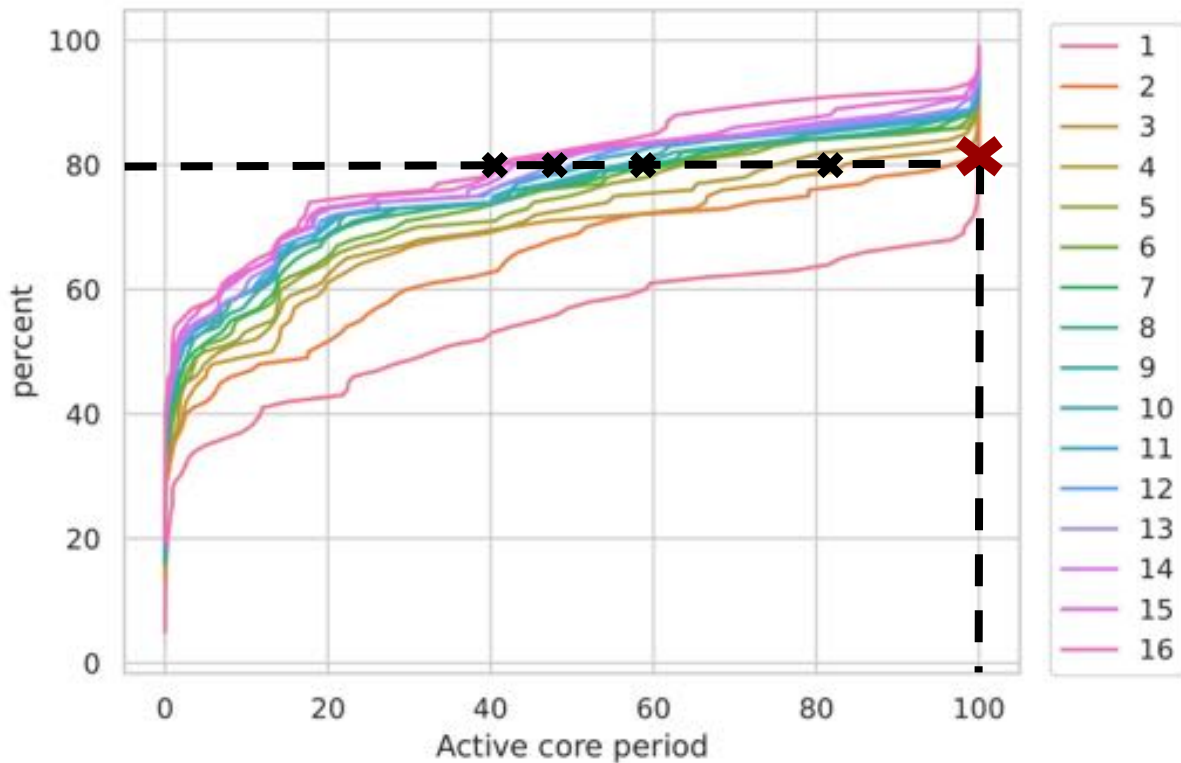
# Contribution 3: SweetspotVM



Large VMs hosted by OVHCloud **do not use all vCPUs equally**

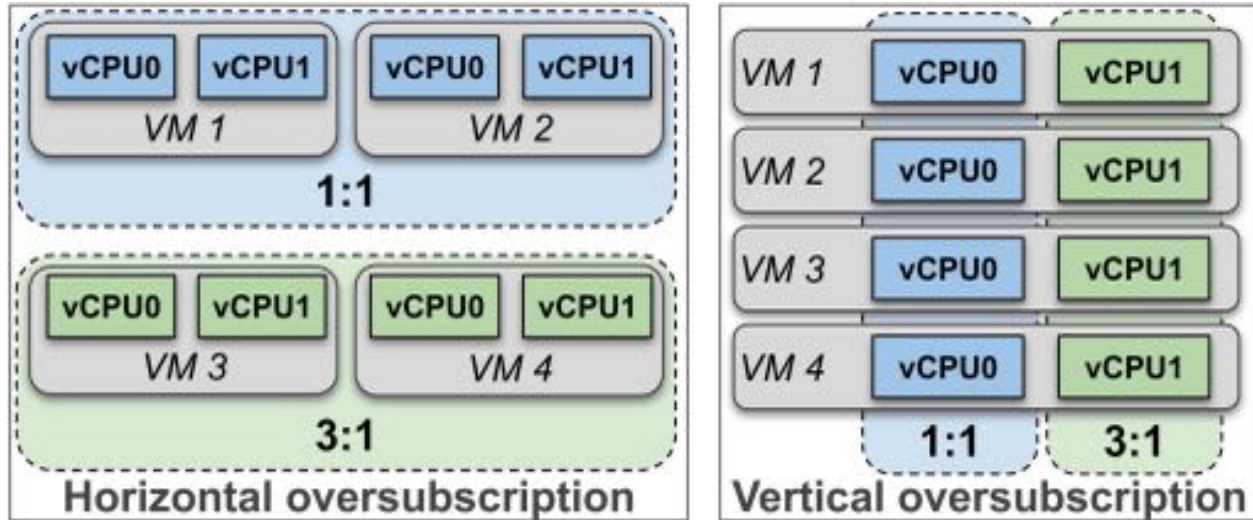


# Contribution 3: SweetspotVM



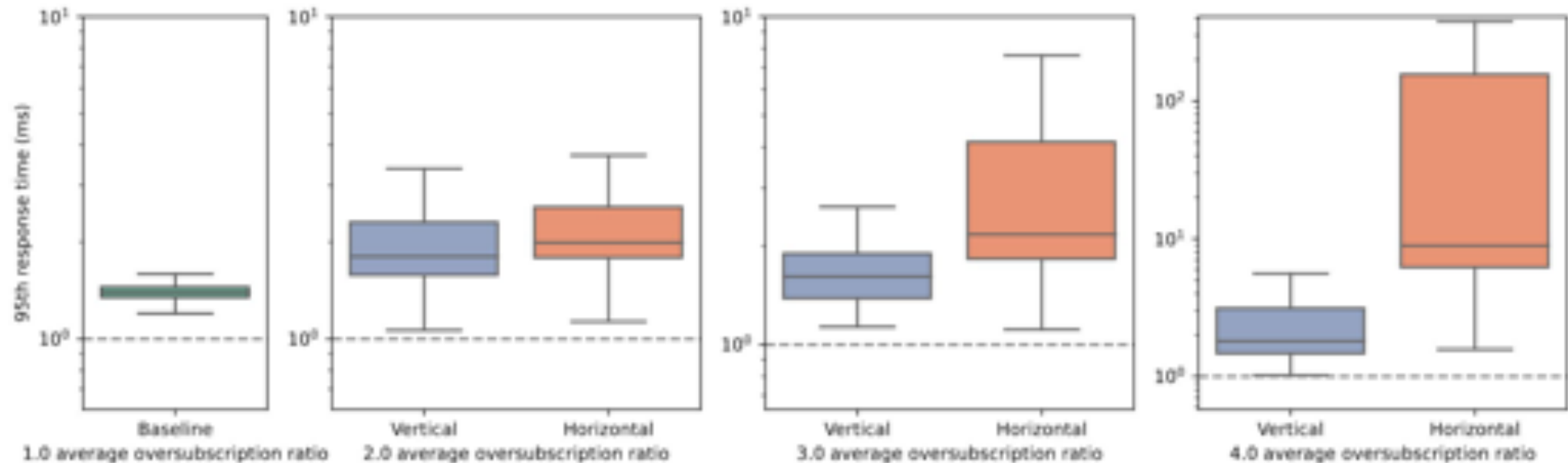
Large VMs hosted by OVHCloud **do not use all vCPUs equally**

# Contribution 3: SweetspotVM




Changing the perspective on vCPU oversubscription

# Contribution 3: SweetspotVM

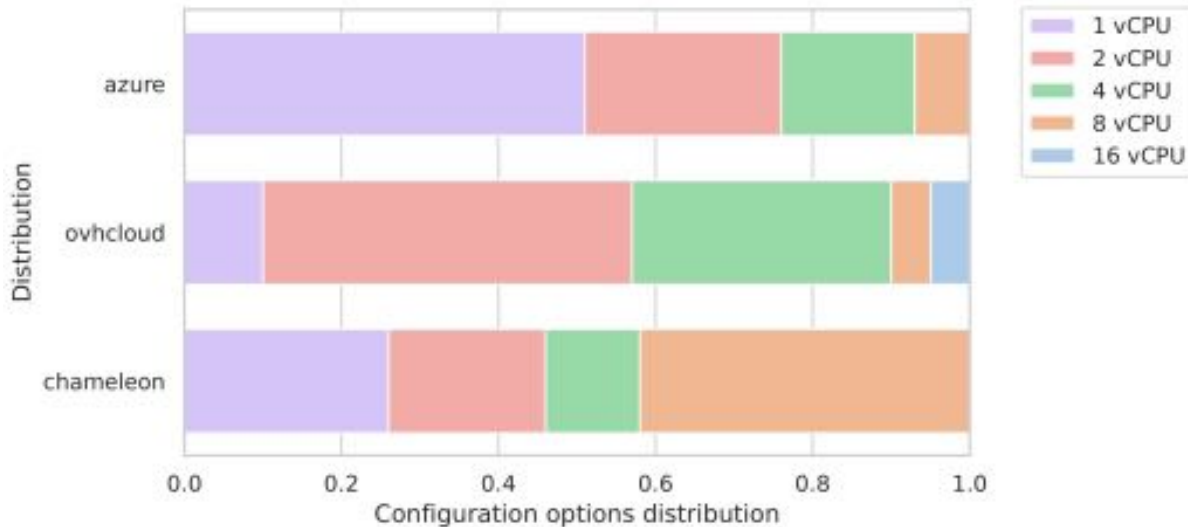
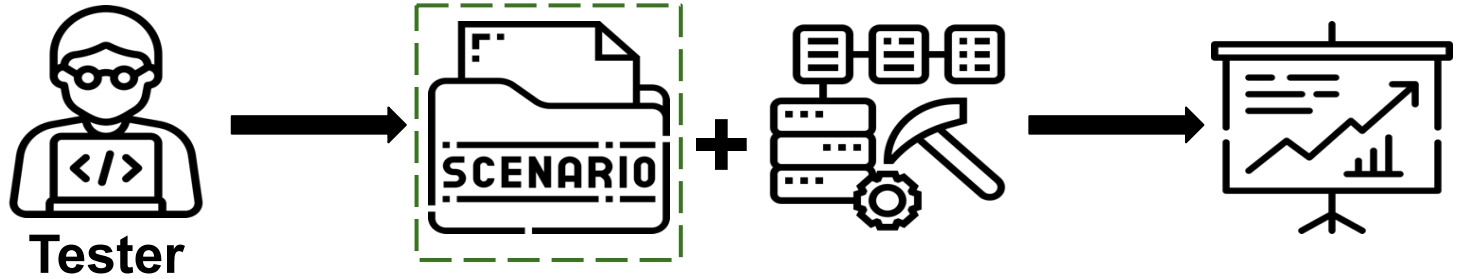


Oversubscribing differently vCPUs **reduces the performance overhead**

# Contribution 4: **CloudFactory**

- Ratios at the cluster level
    - Statically defined
  - Ratios at the server level
    - Statically defined
    - Dynamically defined
      - Generic formulas on VM percentile or host standard deviation
      - **ScroogeVM**: per server formula customization
  - **SlackVM**: ratios at the VM level
  - **SweetspotVM**: ratios at the resource level (vCPU)
- 
- Validation**

# Contribution 4: CloudFactory



# Takeaways

- Cloud resources are underused
  - “All clients do not need all their resources all the time”
- Oversubscription is a way to improve server packing
- Defining oversubscription ratios closer to individual usage is promising

*Papers & code here!*



# Backup

---

## Algorithm 1 Quiescent state detection algorithm

---

**Input** Historical dataset, last window

**Output** Quiescent state

- 1:  $\mathcal{M}_{LSTM} \leftarrow$  Generate model from historical dataset
  - 2:  $forecasted \leftarrow predict(\mathcal{M}_{LSTM}, historical\_set)$
  - 3:  $predicted \leftarrow predict(\mathcal{M}_{LSTM}, last\_window)$
  - 4:  $\delta \leftarrow |RMSE(forecasted) - RMSE(predicted)|$
  - 5: **if**  $\delta > threshold$  **then**
  - 6:     **return** UNSTABLE
  - 7: **end if**
  - 8: **return** QUIESCENT
- 

## QUANTITATIVE EVALUATION OF QUIESCENT STATE CLASSIFIERS

Classifier	Accuracy	Precision	Recall	F-score
<b>Average</b>	0.68	0.8	0.52	0.63
<b>Percentile</b>	0.57	0.55	<b>0.87</b>	0.67
<b>P-value</b>	0.52	<b>1.0</b>	0.06	0.12
<b>LSTM</b>	<b>0.88</b>	0.93	0.84	<b>0.88</b>

Using SotA Google N-Sigma peak predictor:

$$\overline{cpu} + N \times \sigma$$

