



# **ProtoRAIL :** A Risk-cognizant Imitation Agent for Adaptive vCPU Oversubscription In the Cloud

Lu Wang, Mayukh Das, Fangkai Yang, Bo Qiao, Hang Dong, Si Qin, Victor Ruehle, Chetan Bansal, Eli Cortez, Íñigo Goiri, Saravan Rajmohan, Qingwei Lin, Dongmei Zhang

Presented at MLSys 2025, Santra Clara, CA, USA, 12-15 May, 2025

## Outline

- 1. The Oversubscription Problem
- 2. Landscape
- 3. The ProtoRAIL Agent
- 4. Evaluations
- 5. Discussion

### The Oversubscription Problem

Can we bite more than we can chew? Why? How?



## **AIOps for Cloud Efficiency**



# **Capacity Efficiency**





Ioud: vCPU oversubscriptionIn Transport: Airline OverbookingAllocated vCPU < Ask vCPU</td>Passengers ticketed > Airline capacity == x<100% of 1 seat per passenger</td>

#### The Oversubscription problem:

- <u>What</u>? System offers more resources than its available capacity, assuming not all users would simultaneously utilize their allocated capacity
- <u>Why</u>? Diminish the sum of unutilized resources and increase gains.
- Impact? COGS (Cost of Goods Sold) reduction  $\rightarrow$  higher revenue margin
- CPU Oversubscription → minimize stranded cores/memory → Higher capacity efficiency

# The vCPU oversubscription problem





### vCPU Oversubscription problem:

• Why Oversubscription?

User requests 16 #vCore VMs | Physical Node has #60 physical CPU cores

1 vCore = 1.2 Physical Core

 $16 * 3 * 1.2 = 57.6 \approx 60$  physical

So node can fit at most 3 VMs

- VM1 99 %tile usage = 6 vCores
- VM2 99 %tile usage = 12 vCores
- VM3 99 %tile usage = 8 vCores
   (6 + 12 + 8 = 26) \* 1.2 = 31.2 ≈ 32 physical

#### 28 Stranded Cores

If we allocate less than requested we can pack more VMs into physical node.

Oversubscription Ratio

- Should we allocate **50% of requested vCore** to VM1 or is it too tight?
- What if usage peaks higher than usual? Will node have enough extra capacity?

Intractable

• Does the problem of deciding % and packing VMs be solved together

6

# The vCPU oversubscription problem





### Oversubscription problem:

- We de-couple for tractable solution.
- Find optimal Oversubscription Ratio  $\zeta = \frac{Allocated}{Requested} \rightarrow$  Assuming VM allocator is doing a fair job (even if not optimal).
- $\operatorname{argmin}_{\pi} \mathbb{E}_{\zeta \sim \pi} [\sum_{nodes} Risk_{\zeta}] \oplus \operatorname{argmax}_{\pi} \mathbb{E}_{\zeta \sim \pi} [\sum_{nodes} Saved \ \#VCores_{\zeta}]$ **Risk** Benefit

### Challenges:

- Competing objectives
  - Aggressive policy  $\rightarrow$  Contention/JITTER risk;
  - Conservative policy  $\rightarrow$  Stranded cores / memory i.e. Efficiency loss
- Varying demand patterns
- Granularity
- Safety

We need an adaptive solution cognizant of risks

### Landscape

Story so far ...



### [Efficiency Infra] Intelligent Capacity Efficiency Platforms

#### Azure Resource Central Platform [2]

A Data-driven Resource intelligence platform for all of Azure 1st party

- Data science pipelines | Statistical Models
- Robust VM allocators based on predictive heuristics
- Utilization Predictions
- Power Predictions

#### Current naïve solution for oversubscription

- Heuristics-based Binary Usage Classifier
- If Expected Usage <= Threshold, THEN Oversubscription = True/False
- Constant X% oversubscription rate for all VMs of a client across entire lifetime
- Trivial savings Not risk aware



<sup>[1]</sup>WWW '23; <sup>[2]</sup>SOSP '17; <sup>[3]</sup>IEEE Transactions on Sustainable Computing '24; <sup>[4]</sup>CIKM '24,

### [Efficiency Infra] Intelligent Capacity Efficiency Platforms

#### [Online] Multi-Agent RL based Oversubscription Control [1]

Multi-agent Reinforcement Learning for adaptive oversubscription control Each user/VM is agent – needs to factor in allocator/packing heuristics

#### CHALLENGES:

- I. How to train RL agents on real environments
- II. Suitable reward function reward sparsity, delayed feedback etc.
- III. Granularity of decisions(s) cluster, group/service, subscription, VMs
- IV. Safety against "risk" is not trivial convergence is tricky in constrained RL formulations.

#### ScroogeVM – LTSM based oversubscription manager [3] LSTM – based predictor from offline telemetry

#### COIN – Chance-Constrained Offline Learning from uncertain data [4]

Imitation Learning with stochastic constraints for handling uncertainty

Lot of usage telemetry | Can be trained offline | Safer

BUT: Data is noisy / uncertain. No Ground truth. Granularity of decisions(s) – temporal/spatial



#### <sup>[1]</sup>WWW '23; <sup>[2]</sup>SOSP '17; <sup>[3]</sup>IEEE Transactions on Sustainable Computing '24; <sup>[4]</sup>CIKM '24,

### [Resource Level Efficiency] **Key Takeaways**

#### Key problem: Oversubscription

• CPU Oversubscription  $\rightarrow$  minimize stranded cores/memory  $\rightarrow$  Higher capacity efficiency

[Problem] Less physical nodes  $\leftarrow \rightarrow$  More Clients

[Forecast] CPU Usage Forecast

[Optimization] Optimize/Decide oversubscription rate (allocated/requested) → better VM packing

[Objective] COGS (stranded resources) vs Overloading Risk

#### Solution: Smart Imitation Learning (offline mode RL) Input time Tempora

- Policy model from usage telemetry
- + Risk min. via domain knowledge
- + Prototype learning for multigranular decision.

\*\*Fine grained decision & Adaptive solution

**8X increase** in P95 sellable cores ~0% (Negligible) Risk of Jitter

all VMs for a client – for entire lifetime

• Trivial savings

• Oversub = True/False

Usage Classifier

In-Prod Heuristics solution

Const. X% oversub. rate for





#### Challenges

- Dynamic demand NO ground truth;
- Granularity VMs / Subscription / Service
- iii. Noisy uncertain data - risk

#### Impact: 1) 8X Savings 2) Multi-million \$ footprint

#### Status & Road Map:

- Model / Algorithmic development as well as POC complete results promising by deploying on subset of 1<sup>st</sup> party (internal clients)
- Pilot studies going on for all 1<sup>st</sup> party VMs
- Full deployment in pipeline
- Adapt to memory / cache and power oversubscription near future

#### WWW '23; CIKM '24

### The ProtoRAIL Agent



# [Resource Level Efficiency] Intelligent CPU Oversubscription

#### **PROTO**typical **R**isk-cognizant **A**ctive Imitation Learning (**PROTORAIL**)



Module 1

Policy model from usage telemetry

- <u>Prototype learning</u> for **multi-granular** decision.
- Prototypes<sup>1</sup> are "representative points" for equivalence classes of approximately symmetric patterns
- Temporal Usage trajectories  $\rightarrow$  Encode to embedding space  $h_t$
- Learn N prototype embeddings
- Customized similarity function
- Ensure diversity of learned prototypes

**LOSS**:  $L(\pi_D, \pi_\theta) = \left[\ell(h, C) + \ell(C_i, C_j)\right]$ 

<sup>[1]</sup> Molnar, C. (2020). Interpretable machine learning. Lulu. com.; Kim et al., NeurIPS 2016

# [Resource Level Efficiency] Intelligent CPU Oversubscription

#### **PROTO**typical **R**isk-cognizant **A**ctive Imitation Learning (**PROTORAIL**)



### Module (2) Policy model from <u>usage telemetry</u>

- Imitation Learning over prototypes
- New sample  $\rightarrow h_t^{new} \rightarrow sim(h_t^{new}, \{Prototypes\}) = \{w_1, \dots, w_N\}$
- Prediction = aggregate

$$\pi_{\theta}(\bar{\zeta}|.) = \{w_k\}_{k=1}^N \cdot \{\pi_{\theta}^k(\zeta|.)\}_{k=1}^N$$

Optimize with Imitation (behavior cloning)
 loss

$$LOSS: L(\pi_D, \pi_\theta) = \begin{bmatrix} \ell(\zeta, \hat{\zeta}) + \\ \ell(h, C) + \ell(C_i, C_j) \end{bmatrix}$$

# [Resource Level Efficiency] Intelligent CPU Oversubscription

#### **PROTO**typical **R**isk-cognizant **A**ctive Imitation Learning (**PROTORAIL**)



### Module (3)

- Risk min. via domain knowledge (KITL)
- ACTIVE feedback on quality of prototypes + predictions
- Human Domain Expert / LLM expert
- Feedback = upvote / downvote / merge / split
- Very Efficient Minimal queries ~6
- Loss scaling with  $e^{\alpha . [\sum votes]}$
- LLMs are equally good for generic pattern equivalence feedback – Humans better with nuanced **risk related feedback**

LOSS:  $L(\pi_D, \pi_\theta) = \left[\ell(\zeta, \hat{\zeta}) + \ell(h, C) + \ell(C_i, C_j)\right] \times e^{\alpha.Feedback}$ 

## Evaluation

Does it work, for real?



## **Evaluation**

#### Prototypical Risk-cognizant Active Imitation Learning

vCPU Oversubscription: Offline Replay Emulation on 2-week long observations of CPU usage for 1<sup>st</sup> Party workloads on the clusters of only 2 regions and on 300 clusters.

Flight Tickets data is collection from US DOT database for last 2 decades. Risk here is **involuntary offboarding risk**.

Approach	vCPU Over	Flight Tickets		
Approach	Hot Node/Risk↓	Core (Benefit)↑	Cost/Risk↓	Profit↑
Grid-search	0%	7450	0M	0M
Moving Average	1.39%	7628	0.96M	6.79M
DDPG	1.47%	5030	12.37M	2.35M
Behavior Cloning	1.19%	7870	1.47M	7.21M
GAIL	1.2%	6980	2.74M	4.56M
Dagger (20 time steps)	0.96%	7938	0.47M	6.95M
LSTM	1.27%	7749	1.82M	4.98M
Coop. Multi-Agent RI	L 0.89%	7897	0.59M	8.17M
PROTORAIL (w/o KITL	) 0%	8153	0.31M	8.79M
PROTORAIL	0%	8161	0.14M	13.65M

"Manual" indicates manually decided oversubscription, "Heuristics" indicates statistical non-adaptive prediction model based on collected heuristics in production.

Method	$MSE\downarrow$	VMs w/ $\beta \uparrow$	VMs at Risk $\downarrow$	BRR ↑
Manual	-NA-	109	14	7.79
Heuristics	0.065	3502	185	18.95
PROTORAIL	0.042	3542	113	31.34

	Heuristics (Binary)	ProtoRAIL	Total
% savings in core hours	5.2%	20.5%	100%

#### **Deployment Metrics**

- ~ 6-8X increase in P95 # sellable cores
- ~ 4x increase in saved (core hours) 1
- ~ 0% (Negligible) Risk of Jitter↓

~ 0.2% Underprediction Ratio compared to 7%  $\downarrow$ 

### [Resource Level Efficiency] Prototypes Learning / Refinement & Safety



Prototypes do capture approximate symmetries in temporal patterns across ANY granularity

Knowledge in the loop helped get rid of marginally useless prototypes.

Note how risk gets minimized, Often with Higher Savings, with Knowledge in the loop (KITL)

# **Conclusion and Discussion**

- Adaptive Risk-aware Intelligent Multi-granular oversubscription
- **~**
- Substantially higher sellable cores/ core\_hours savings, Minimized overloading risk
- Real tests in internal deployments show significant gains
- Cloud efficiency that is reliable making cloud operations sustainable
- Al contribution: Novel Prototypical Imitation Learning
  - Full Deployment for all 1<sup>st</sup> Party
  - Adapting to other oversubscription / overbooking problems
    - Workload characteristics aware policies



A representation of Mundell-Flemming's Trilemma in Macroeconomics applied on cloud business

# Thank you!





Our Research Group M365 Research, Microsoft: <u>https://www.microsoft.com/en-us/research/group/m365-research/</u>

Systems Innovation Effort at Microsoft: <u>https://www.microsoft.com/en-us/research/group/systems-innovation/</u>

Acknowledgements: DKI (Microsoft, Beijing), Azure Resource Central (Microsoft)