
VMR²L: Virtual Machines Rescheduling Using Reinforcement Learning in Data Centers

Xianzhong Ding^{*1} Yunkai Zhang^{*2} Binbin Chen³ Donghao Ying² Tieying Zhang³
Jianjun Chen³ Lei Zhang³ Alberto Cerpa¹ Wan Du¹
¹UC Merced ²UC Berkeley ³ByteDance

Abstract

Modern industry-scale data centers receive thousands of virtual machine (VM) requests per minute. Due to the continual creation and release of VMs, many small resource fragments are scattered across physical machines (PMs). To handle these fragments, data centers periodically reschedule some VMs to alternative PMs. Despite the increasing importance of VM rescheduling as data centers grow in size, the problem remains understudied. We first show that, unlike most combinatorial optimization tasks, the inference time of VM rescheduling algorithms significantly influences their performance, causing many existing methods to scale poorly. Therefore, we develop a reinforcement learning system for VM rescheduling, VMR²L, which incorporates a set of customized techniques, such as a two-stage framework that accommodates diverse constraints and workload conditions, as well as an effective feature extraction module. Our experiments on an industry-scale data center show that VMR²L can achieve a performance comparable to the optimal solution, but with a running time of seconds.²

1 Introduction

Cloud service providers allow end-users to access computing resources, such as CPU and memory. They adopt resource virtualization to maximize hardware utilization, allocating Virtual Machines (VMs)[1, 2] with the requested resources to end-users [3, 4, 5]. An industry-scale data center typically has hundreds to thousands of Physical Machines (PMs), where each PM can host multiple VMs that run independently [6]. A central server manages all VM requests on PMs by performing two tasks, scheduling and rescheduling, in order to achieve different resource utilization goals, such as minimizing the overall fragment rate (FR) or maximizing the number of available PMs.³

VM Scheduling (VMS). When new VM requests arrive, there can be multiple available PMs to host them. Figure 1 shows the average number of VMs changes (VMs arriving and exiting) per minute over a 30-day period from our in-house data center. To ensure the system is robust, the VMS algorithm needs to meet the maximum number of VMs changes as indicated by the green line. The high queries per second (QPS) requirement deems only heuristic methods feasible for VMS. In fact, ByteDance uses best-fit [7, 8], which sorts all PMs that meet the requirements of the current VM according to the amount of FR reduction before and after this VM is added, and chooses the PM with the largest reduction.

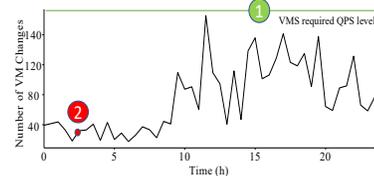


Figure 1: Avg. VM changes.

^{*}The authors contribute equally.

²We release our datasets and code here: <https://anonymous.4open.science/r/VMR2L-BEA6>.

³Due to space limits, we mainly explain our methods in terms of minimizing FR, whose formal definition is presented in Appendix A.1.

VM Rescheduling (VMR). However, simple heuristics are often far from optimal, and the continual exiting of completed VMs results in many fragments scattered across PMs. As a result, rescheduling is critical to optimize resource usage, which migrates VMs from their current PMs to new destination PMs. Due to the overhead of VM migrations, a migration number limit (MNL) is set to control the number of VMs to migrate. Note that while the VMR algorithm computes a solution, VMS is still handing new VM requests and completed VMs are also being deleted. The dynamic nature of VM states causing the computed VMR solution to no longer be optimal or even feasible⁴. Therefore, even though VMR mostly happens during off-peak hours where there are fewer VM changes⁵, VMR still needs to be very efficient. We conduct an experiment to quantify how the VMR inference time affects the achieved performance.

Motivation Experiment. We can formulate the VM rescheduling problem as a Mixed Integer Programming (MIP) problem, where the constraints come from the service expectations and the available hardware resources. An off-the-shelf MIP solver, such as Gurobi [9] and CPLEX [10], can achieve a near-optimal objective, but it suffers from an inference time that grows exponentially with the number of VMs and PMs, and thus fails to scale to large data centers.

To see how a near-optimal solution can result in a suboptimal achieved performance due to its poor inference time, we conduct an experiment on real traces from our in-house data center by selecting 20 random initial VM-PM mappings. For each mapping, we use Gurobi to compute a near-optimal solution to the MIP formulation of VMR, which takes 50.55 minutes. However, since VMs were dynamically arriving and exiting, most actions were no longer feasible and will fail to be deployed after 50 minutes. We then compute the final performance that could be achieved as if the near-optimal solution was instead returned in a shorter period of time, averaged over the 20 mappings. Figure 2 shows that the solution remains near-optimal if it could be computed within 5 seconds. However, FR reduction quickly diminishes when the inference time exceeds 20 seconds.

Limitations of the Current Methods. The experiment results reveal that different from bin-packing and other MIP applications, VM rescheduling requires an algorithm that has an inference time strictly under five seconds. To accelerate the running speed of the MIP approach, hand-tuned heuristics can be integrated into the process, e.g., adding constraints to limit the solver’s search space. These heuristics must trade off between the optimality of the solution and the tractability of the problem. Unfortunately, even highly-skilled experts need many iterations to manually find a proper trade-off, and no universal heuristics can achieve a good trade-off for all VM rescheduling scenarios.

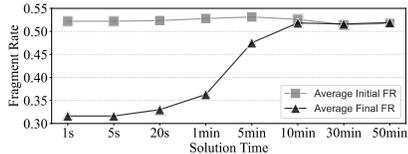


Figure 2: Effect of inference time on achieved performance.

In this work, we develop VMR²L, a deep Reinforcement Learning (RL) system for VM rescheduling. VMR²L trains a Deep Neural Network (DNN) as the rescheduling agent. RL is a great fit for VMR since VMR has no aleatoric uncertainties, i.e., the next state can be exactly simulated given the current state and action. This allows us to build a simulator that only requires the initial VM-PM mappings for training, without having to interact with a real data center. With our RL formulation and customized embedding techniques, inference can be done within one second and can scale easily to a large number of VMs and PMs. Extensive evaluation on two collected datasets demonstrates that VMR²L can generate a solution that is only 2.67% worse than the optimal solution. We summarize the contributions of this paper as follows:

- **RL for VM rescheduling.** We identify the unique characteristics of the VM rescheduling problem in terms of latency requirement and aleatoric uncertainties, which motivate its formulation as a RL problem.
- **Customized RL techniques for VM rescheduling.** We tackle three challenges in designing a RL framework for VM rescheduling with three customized techniques, including a two-stage framework and effective feature extraction module.

⁴A VM will not be rescheduled if it has exited or the destination PM no longer has enough resources.
⁵In less common cases, VMR is also performed if a high FR is observed that could potentially lead to insufficient resources for upcoming VM requests, which requires an even higher latency.

- **A VMR²L prototype and extensive evaluation.** We develop a prototype of VMR²L and conduct extensive experiments over two real datasets. We release the datasets and a custom Gym environment for RL training.

2 Design of VMR²L

2.1 VM Rescheduling as an RL Problem

In this paper, we adopt a deep reinforcement learning (RL) approach to address the VM rescheduling problem. A VMR request starts an episode, which involves migration number limit (MNL) steps. At each migration step, the agent reschedules one VM from its source PM to a new destination PM based on the current PM and VM states. Notably, the environment is deterministic – given the initial state and an action, we can directly simulate the change in objective and the next state. This allows us to build a simulator to train the agent, as detailed in Appendix B.1.

State Representation. The state input to the DRL agent contains two sets of features. The first set contains four PM features for each NUMA⁶, specifically the remaining CPU and memory resources, current FR, and fragment sizes. The second set contains 14 VM features including requested CPU and memory for each NUMA, fragment sizes, and the source PM details. If a single NUMA is requested, zeros are used as placeholders for the other NUMA. Features are min-max normalized.

Action Representation. The action at each step can be represented as a 2-tuple (k, i) . Specifically, the action is to reschedule a VM k from its source PM to a destination PM i . Note that the source PM can be retrieved once we select k .

Reward Representation. The goal of VM rescheduling is to minimize the FR across all PMs. While we could return the FR of all PMs as a single final reward to the agent after finishing an entire episode, it corresponds to a form of sparse reward and it is known to be difficult for training [11]. Instead, we propose to generate dense rewards and use the change in FR on the source PM and the destination PM as an intermediate reward at each step. As such, the reward range is naturally scaled down to $[-2, 2]$, which we further normalize by dividing with a constant c [12]. We calculate the rescaled fragment size on each NUMA by $S_i = \sum_{j=0}^1 (\tilde{U}_{i,j} \% X) \div c$, and define reward as $R = (S_{\text{before, src}} - S_{\text{after, src}}) + (S_{\text{before, dest}} - S_{\text{after, dest}})$, where $S_{\text{before, \cdot}}$ and $S_{\text{after, \cdot}}$ are the fragment changes before and after the selected VM leaves (enters) the source (destination) PM.

2.2 A Two-Stage Framework

We use PPO [13] as the backbone of our DRL framework. To better accommodate a variety of constraints, we leverage the characteristics of the VMR problem and design a two-stage framework that allows the action tuple to be generated sequentially. As shown in Figure 3, in Stage 1, the VM actor selects the VM candidate to be rescheduled. Once a candidate VM is selected, we can efficiently mask out all the PMs that cannot host the candidate VM and then proceed to Stage 2, where the PM actor selects an appropriate destination PM from the remaining PMs. Additionally, when we select a VM to reschedule, a considerable portion of the PMs cannot meet its resource requirements. The exploration of RL agent on these PMs inevitably hinders the learning process. The proposed framework can effectively reduce the large action space and thus mitigate the exploration challenge.

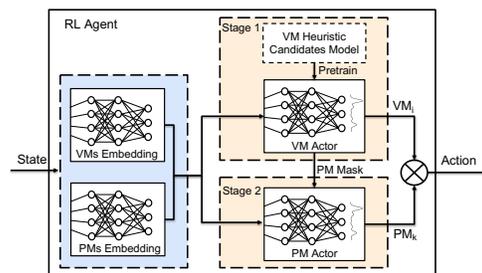


Figure 3: The two-stage RL agent.

Scalability. To make effective rescheduling decisions, VMR²L must extract meaningful representations of the state observation, which include features of each individual PM and VM as well as their affiliations. However, the number of VMs can vary drastically even in the same data center. To encode these features, we propose to share two small embedding networks across all VMs and PMs (Figure 4) — one to process each PM’s features and another one to process each VM’s features. As

⁶Each PM contains two non-uniform memory accesses (NUMAs).

such, the number of weight parameters is *independent* of the number of machines in the system. This is achieved via an attention-based transformer model [14] but tailored for rescheduling. Transformers have demonstrated strong performances in combinatorial optimization, such as in vector bin-packing [15, 16]. However, there is a notable difference in VM rescheduling: we must choose from a set of VMs that have already been assigned to PMs.

2.3 Feature Extraction with Sparse Attention

Tree-level Features. Consider the example shown in Appendix C.1. Not knowing which exact types of VMs are hosted on the same PM prevents the vanilla transformer from optimizing for such long-term rewards. We argue that each VM must be aware of which VMs are co-hosted on the same PM, resembling a tree structure with PM as the root and VMs as leaves. To enable VMs to recognize co-hosted VMs, we introduce an additional stage of *sparse local-attention* within each PM tree. This restricts PMs and VMs to attend to each other if and only if they belong to the same tree. Detailed network architecture can be found in Appendix C.

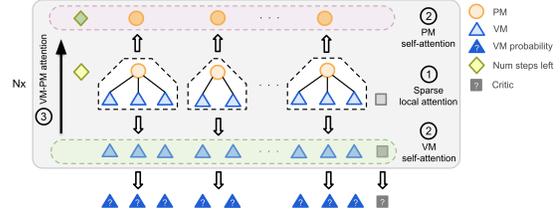


Figure 4: VM actor with sparse local-attention capturing tree-level features.

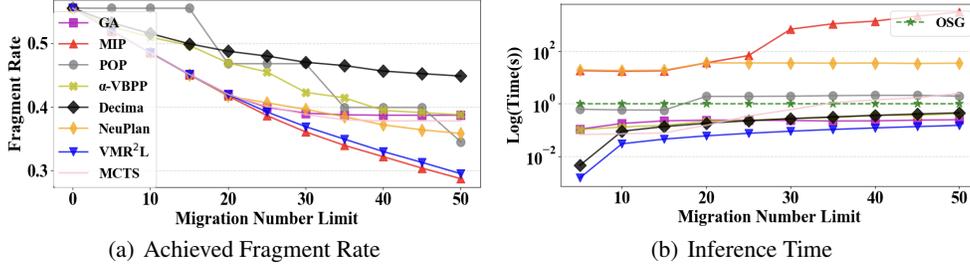


Figure 5: FR and inference time on the medium dataset at different MNLs.

3 Evaluation

We collect two datasets from an industry-scaled cloud data center – one medium dataset with 2089 VMs and 280 PMs, and one large dataset with 4546 VMs and 1176 PMs. We implement VMR²L based on the CleanRL framework [17] with PyTorch using PPO as the backbone [13]. We compare with seven baseline methods from six categories: heuristics (e.g., greedy, α -VBPP), optimization algorithms (e.g., MIP), approximate algorithms (e.g., POP), search-based algorithms (e.g., MCTS), deep learning-based (e.g., Decima), and hybrid methods (e.g., NeuPlan).⁷

Experiment. We compare VMR²L with the baselines on the dataset in terms of the optimality and inference latency under different migration number limits (MNLs). Figure 5 summarizes the results, which shows that VMR²L is able to consistently reduce FR given different MNLs at a rate closest to MIP compared to other baselines. For a more comprehensive analysis, refer to Appendix B.3. The VMR²L training details can be found in Appendix A.4.

4 Conclusion

Compared to conventional bin-packing applications, VM rescheduling presents unique challenges due to the expanding size of data centers. It needs to handle many VMs while meeting a strict inference speed requirement. To this end, we introduce VMR²L, a tailored deep RL solution featuring a two-stage framework for diverse service constraints and a sparse attention module for better feature extractions. We hope that our released datasets and RL environment will support future research in this area.

⁷See Appendix B.2 and B.1 for details on the baselines and datasets, respectively.

References

- [1] Jörg Thalheim, Peter Okelmann, Harshavardhan Unnibhavi, Redha Gouicem, and Pramod Bhatotia. Vmsh: hypervisor-agnostic guest overlays for vms. In *Proceedings of the Seventeenth European Conference on Computer Systems*, pages 678–696, 2022.
- [2] Christopher Clark, Keir Fraser, Steven Hand, Jacob Gorm Hansen, Eric Jul, Christian Limpach, Ian Pratt, and Andrew Warfield. Live migration of virtual machines. In *Proceedings of the 2nd conference on Symposium on Networked Systems Design & Implementation-Volume 2*, pages 273–286, 2005.
- [3] Xianzhong Ding, Le Chen, Murali Emani, Chunhua Liao, Pei-Hung Lin, Tristan Vanderbruggen, Zhen Xie, Alberto Cerpa, and Wan Du. Hpc-gpt: Integrating large language model for high-performance computing. In *Proceedings of the SC'23 Workshops of The International Conference on High Performance Computing, Network, Storage, and Analysis*, pages 951–960, 2023.
- [4] Mao Lin, Keren Zhou, and Pengfei Su. Drgpum: Guiding memory optimization for gpu-accelerated applications. In *Proceedings of the 28th ACM International Conference on Architectural Support for Programming Languages and Operating Systems, Volume 3*, pages 164–178, 2023.
- [5] Kang Yang, Yuning Chen, Xuanren Chen, and Wan Du. Link quality modeling for lora networks in orchards. In *Proceedings of the 22nd International Conference on Information Processing in Sensor Networks*, pages 27–39, 2023.
- [6] Ori Hadary, Luke Marshall, Ishai Menache, Abhisek Pan, Esaias E Greeff, David Dion, Star Dorminey, Shailesh Joshi, Yang Chen, Mark Russinovich, et al. Protean: Vm allocation service at scale. In *Proceedings of the 14th USENIX Conference on Operating Systems Design and Implementation*, pages 845–861, 2020.
- [7] Chi Trung Ha, Trung Thanh Nguyen, Lam Thu Bui, and Ran Wang. An online packing heuristic for the three-dimensional container loading problem in dynamic environments and the physical internet. In *Applications of Evolutionary Computation: 20th European Conference, EvoApplications 2017, Amsterdam, The Netherlands, April 19-21, 2017, Proceedings, Part II 20*, pages 140–155. Springer, 2017.
- [8] Francisco Parreño, Ramón Alvarez-Valdés, Jose Manuel Tamarit, and Jose Fernando Oliveira. A maximal-space algorithm for the container loading problem. *INFORMS Journal on Computing*, 20(3):412–422, 2008.
- [9] Gurobi solver. <https://www.gurobi.com/>.
- [10] Cplex optimizer. <https://www.ibm.com/analytics/cplex-optimizer>.
- [11] Desik Rengarajan, Gargi Vaidya, Akshay Sarvesh, Dileep Kalathil, and Srinivas Shakkottai. Reinforcement learning with sparse rewards using guidance from offline demonstration. *arXiv preprint arXiv:2202.04628*, 2022.
- [12] Peter Henderson, Riashat Islam, Philip Bachman, Joelle Pineau, Doina Precup, and David Meger. Deep reinforcement learning that matters. In *Proceedings of the AAAI conference on artificial intelligence*, volume 32, 2018.
- [13] John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017.
- [14] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30, 2017.
- [15] Jingwei Zhang, Bin Zi, and Xiaoyu Ge. Attend2pack: Bin packing through deep reinforcement learning with attention. *ArXiv*, abs/2107.04333, 2021.

- [16] Dongda Li, Changwei Ren, Zhaoquan Gu, Yuexuan Wang, and Francis Lau. Solving packing problems by conditional query learning, 2020.
- [17] Shengyi Huang, Rousslan Fernand Julien Dossa, Chang Ye, Jeff Braga, Dipam Chakraborty, Kinal Mehta, and João G.M. Araújo. Cleanrl: High-quality single-file implementations of deep reinforcement learning algorithms. *Journal of Machine Learning Research*, 23(274):1–18, 2022.
- [18] Marco A Wiering and Martijn Van Otterlo. Reinforcement learning. *Adaptation, learning, and optimization*, 12(3):729, 2012.
- [19] Kai Arulkumaran, Marc Peter Deisenroth, Miles Brundage, and Anil Anthony Bharath. Deep reinforcement learning: A brief survey. *IEEE Signal Processing Magazine*, 34(6):26–38, 2017.
- [20] Raja Wasim Ahmad, Abdullah Gani, Siti Hafizah Ab Hamid, Muhammad Shiraz, Abdullah Yousafzai, and Feng Xia. A survey on virtual machine migration and server consolidation frameworks for cloud data centers. *Journal of network and computer applications*, 52:11–25, 2015.
- [21] Xianzhong Ding, Zhiyong Zhang, Zhiping Jia, Lei Ju, Mengying Zhao, and Huawei Huang. Unified nvtcam and stcam architecture for improving packet matching performance. *ACM SIGPLAN Notices*, 52(5):91–100, 2017.
- [22] Qingpeng Cai, Will Hang, Azalia Mirhoseini, George Tucker, Jingtao Wang, and Wei Wei. Reinforcement learning driven heuristic optimization. *Workshop on Deep Reinforcement Learning for Knowledge Discovery (DRL4KDD)*, abs/1906.06639, 2019.
- [23] Haoyuan Hu, Xiaodong Zhang, Xiaowei Yan, Longfei Wang, and Yinghui Xu. Solving a new 3d bin packing problem with deep reinforcement learning method, 2017.
- [24] Lu Duan, Haoyuan Hu, Yu Qian, Yu Gong, Xiaodong Zhang, Jiangwen Wei, and Yinghui Xu. A multi-task selected learning approach for solving 3d flexible bin packing problem. In *Proceedings of the 18th International Conference on Autonomous Agents and MultiAgent Systems*, AAMAS '19, page 1386–1394, Richland, SC, 2019. International Foundation for Autonomous Agents and Multiagent Systems.
- [25] Ye Xia, Mauricio Tsugawa, Jose AB Fortes, and Shigang Chen. Large-scale vm placement with disk anti-colocation constraints using hierarchical decomposition and mixed integer programming. *IEEE Transactions on Parallel and Distributed Systems*, 28(5):1361–1374, 2016.
- [26] Rina Panigrahy, Kunal Talwar, Lincoln Uyeda, and Udi Wieder. Heuristics for vector bin packing. *research. microsoft. com*, 2011.
- [27] A. Paul Davies and Eberhard E. Bischoff. Weight distribution considerations in container loading. *European Journal of Operational Research*, 114(3):509–527, May 1999.
- [28] Xijun Li, Mingxuan Yuan, Di Chen, Jianguo Yao, and Jia Zeng. A data-driven three-layer algorithm for split delivery vehicle routing problem with 3d container loading constraint. In *Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, KDD '18, page 528–536, New York, NY, USA, 2018. Association for Computing Machinery.
- [29] Hang Zhao, Qijin She, Chenyang Zhu, Yin Yang, and Kai Xu. Online 3d bin packing with constrained deep reinforcement learning. In *Thirty-Fifth AAAI Conference on Artificial Intelligence, AAAI 2021, Thirty-Third Conference on Innovative Applications of Artificial Intelligence, IAAI 2021, The Eleventh Symposium on Educational Advances in Artificial Intelligence, EAAI 2021, Virtual Event, February 2-9, 2021*, pages 741–749. AAAI Press, 2021.
- [30] Qianwen Zhu, Xihan Li, Zihan Zhang, Zhixing Luo, Xialiang Tong, Mingxuan Yuan, and Jia Zeng. Learning to pack: A data-driven tree search algorithm for large-scale 3d bin packing problem. In *Proceedings of the 30th ACM International Conference on Information Knowledge Management*, CIKM '21, page 4393–4402, New York, NY, USA, 2021. Association for Computing Machinery.

- [31] Ameer Haj-Ali, Qijing Jenny Huang, John Xiang, William Moses, Krste Asanovic, John Wawrzynnek, and Ion Stoica. Autophase: Juggling hls phase orderings in random forests with deep reinforcement learning. *Proceedings of Machine Learning and Systems*, 2:70–81, 2020.
- [32] Maxime Gasse, Didier Chételat, Nicola Ferroni, Laurent Charlin, and Andrea Lodi. Exact combinatorial optimization with graph convolutional neural networks. *Advances in Neural Information Processing Systems*, 32, 2019.
- [33] Marc Etheve, Zacharie Alès, Côme Bissuel, Olivier Juan, and Safia Kedad-Sidhoum. Reinforcement learning for variable selection in a branch and bound algorithm. In *International Conference on Integration of Constraint Programming, Artificial Intelligence, and Operations Research*, pages 176–185. Springer, 2020.
- [34] Prateek Gupta, Maxime Gasse, Elias Khalil, Pawan Mudigonda, Andrea Lodi, and Yoshua Bengio. Hybrid models for learning to branch. *Advances in neural information processing systems*, 33:18087–18097, 2020.
- [35] Haoran Sun, Wenbo Chen, Hui Li, and Le Song. Improving learning to branch via reinforcement learning. 2020.
- [36] Yunhao Tang, Shipra Agrawal, and Yuri Faenza. Reinforcement learning for integer programming: Learning to cut. In *International conference on machine learning*, pages 9367–9376. PMLR, 2020.
- [37] Quentin Cappart, Thierry Moisan, Louis-Martin Rousseau, Isabeau Prémont-Schwarz, and Andre A Cire. Combining reinforcement learning and constraint programming for combinatorial optimization. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 35, pages 3677–3687, 2021.
- [38] Lara Scavuzzo, Feng Yang Chen, Didier Chételat, Maxime Gasse, Andrea Lodi, Neil Yorke-Smith, and Karen Aardal. Learning to branch with tree mdps. *arXiv preprint arXiv:2205.11107*, 2022.
- [39] Jialin Song, Yisong Yue, Bistra Dilkina, et al. A general large neighborhood search framework for solving integer linear programs. *Advances in Neural Information Processing Systems*, 33:20012–20023, 2020.
- [40] Thomas Barrett, William Clements, Jakob Foerster, and Alex Lvovsky. Exploratory combinatorial optimization with reinforcement learning. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, pages 3243–3250, 2020.
- [41] Meng Qi, Mengxin Wang, and Zuo-Jun Shen. Smart feasibility pump: Reinforcement learning for (mixed) integer programming. *arXiv preprint arXiv:2102.09663*, 2021.
- [42] Deepak Narayanan, Fiodar Kazhamiaka, Firas Abuzaid, Peter Kraft, Akshay Agrawal, Srikanth Kandula, Stephen Boyd, and Matei Zaharia. Solving large-scale granular resource allocation problems efficiently with pop. In *Proceedings of the ACM SIGOPS 28th Symposium on Operating Systems Principles*, pages 521–537, 2021.
- [43] Sham M Kakade. A natural policy gradient. *Advances in neural information processing systems*, 14, 2001.
- [44] John Schulman, Sergey Levine, Pieter Abbeel, Michael Jordan, and Philipp Moritz. Trust region policy optimization. In *ICML*. PMLR, 2015.
- [45] Kyriakos G Vamvoudakis and Frank L Lewis. Online actor–critic algorithm to solve the continuous-time infinite horizon optimal control problem. *Automatica*, 46(5):878–888, 2010.
- [46] John Schulman, Philipp Moritz, Sergey Levine, Michael Jordan, and Pieter Abbeel. High-dimensional continuous control using generalized advantage estimation. *arXiv preprint arXiv:1506.02438*, 2015.
- [47] Vincent Mai, Kaustubh Mani, and Liam Paull. Sample efficient deep reinforcement learning via uncertainty estimation. In *International Conference on Learning Representations*, 2022.

- [48] Xianzhong Ding, Wan Du, and Alberto E Cerpa. Mb2c: Model-based deep reinforcement learning for multi-zone building control. In *Proceedings of the 7th ACM international conference on systems for energy-efficient buildings, cities, and transportation*, pages 50–59, 2020.
- [49] Greg Brockman, Vicki Cheung, Ludwig Pettersson, Jonas Schneider, John Schulman, Jie Tang, and Wojciech Zaremba. Openai gym, 2016.
- [50] Kubernetes scheduler. <https://kubernetes.io/docs/concepts/scheduling-eviction/kube-scheduler/>.
- [51] Hongzi Mao, Malte Schwarzkopf, Shaileshh Bojja Venkatakrisnan, Zili Meng, and Mohammad Alizadeh. Learning scheduling algorithms for data processing clusters. In *Proceedings of the ACM special interest group on data communication*, pages 270–288. 2019.
- [52] Hang Zhu, Varun Gupta, Satyajeet Singh Ahuja, Yuandong Tian, Ying Zhang, and Xin Jin. Network planning with deep reinforcement learning. In *Proceedings of the 2021 ACM SIGCOMM 2021 Conference*, pages 258–271, 2021.
- [53] Jimmy Lei Ba, Jamie Ryan Kiros, and Geoffrey E. Hinton. Layer normalization, 2016.
- [54] Christopher M. Bishop. *Pattern Recognition and Machine Learning (Information Science and Statistics)*. Springer, 1 edition, 2007.
- [55] Tianyu He, Xu Tan, Yingce Xia, Di He, Tao Qin, Zhibo Chen, and Tie-Yan Liu. Layer-wise coordination between encoder and decoder for neural machine translation. In S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 31. Curran Associates, Inc., 2018.
- [56] Dan Hendrycks and Kevin Gimpel. Gaussian Error Linear Units (GELUs). 2016.
- [57] Fabio Pardo, Arash Tavakoli, Vitaly Levnik, and Petar Kormushev. Time limits in reinforcement learning, 2018.

Appendix Contents

A Appendix	10
A.1 VM Rescheduling Problem Formulation	10
A.2 Background on (Deep) Reinforcement Learning	11
A.3 Related Work	12
A.4 PPO	12
B Experiment Details	13
B.1 Simulator and Datasets	13
B.2 Existing Baseline Algorithms	14
B.3 Performance on the Medium Dataset	14
B.4 Scalability to the Large Dataset	15
B.5 Different Service Objectives	16
B.6 Generalizing to Different MNLs	16
C Architecture Overview	16
C.1 Sparse Attention Details	17
C.2 A Case Study on Sparse Attention	19

Table 1: The VM types we consider in our experiments.

VM Types	large	xlarge	2xlarge	4xlarge	8xlarge	16xlarge	22xlarge
Requested CPU	2	4	8	16	32	64	88
Requested Memory (GB)	4	8	16	32	64	128	256
Deploy NUMA	Single	Single	Single	Single	Double	Double	Double
Percentage	11.01%	43.51%	16.13%	17.62%	6.22%	5.31%	0.20%

A Appendix

A.1 VM Rescheduling Problem Formulation

In this section, we formulate the VM rescheduling problem as a Mixed-Integer Programming (MIP) problem.

In a data center, let \mathcal{V}, \mathcal{P} be the set of VMs and PMs, respectively. On the supply side, a PM $i \in \mathcal{P}$ has two NUMAs⁸, where NUMA j can provide $U_{i,j}$ CPU resources and $V_{i,j}$ memory resources. On the demand side, a VM $k \in \mathcal{V}$ requires u_k CPU resources and v_k memory resources and should be deployed on a single PM using $w_k \in \{1, 2\}$ NUMAs. w_k is the number of NUMAs required by VM k (1 for single-NUMA deployment, 2 for double-NUMA). After deploying several VMs on PM $i \in \mathcal{P}$, there remains $\tilde{U}_{i,j}$ spare CPU resources on NUMA j . We define **X-core fragment** of PM i as $\sum_j (\tilde{U}_{i,j} \% X)$, i.e., the remaining CPU resources cannot be further utilized by any additional X-core VMs.

Given an initial assignment of M VMs each onto one of the N PMs, the VM rescheduling task is to reassign a subset of deployed VMs and migrate them each onto a new PM. The maximum number of VMs that we can migrate for a given task is called *Migration Number Limit (MNL)*. We formulate the VM rescheduling as an optimization problem that searches for a reassignment of MNL selected VMs, in order to minimize the total X-core fragments across all PMs:

$$\text{Minimize:} \quad \sum_{i,j} \left(U_{i,j} - \sum_k \frac{x_{k,i,j} \cdot u_k}{w_k} - X y_{i,j} \right) \quad (1)$$

$$\text{Subject to:} \quad \sum_k \frac{x_{k,i,j} \cdot u_k}{w_k} + X y_{i,j} \leq U_{i,j}, \quad (2)$$

$$\sum_k \frac{x_{k,i,j} \cdot v_k}{w_k} \leq V_{i,j}, \quad (3)$$

$$\sum_{i,j} x_{k,i,j} = w_k, \quad (4)$$

$$\sum_k (1 - x_{k,i_k,j_k}) \leq MNL, \quad (5)$$

$$x_{k,i,0} = x_{k,i,1}, \quad \forall k \in \{k | w_k = 2\}, \quad (6)$$

$$x_{k,i,j} \in \{0, 1\} \text{ and } y_{i,j} \in \mathbb{Z}. \quad (7)$$

Here, $\{x, y\}$ are the decision variables, where $x_{k,i,j}$ represents whether VM k is deployed to the NUMA j of PM i in the new assignment (0 for No, 1 for Yes), and $y_{i,j}$ represents the maximum number of X-core VMs can be deployed on NUMA j of PM i using the remaining CPU resources. The objective in Equation 1 is to minimize the total X-core fragments.

Equation 2 and 3 enforce that the resource usage by VMs cannot exceed the total capacity of a PM. Equation 4 indicates that each VM must be deployed on exactly one PM. Equation 5, in which i_k and j_k are the initial PM id and NUMA id (0 for double-NUMA VMs) of VM k , means the total migration number should not exceed the limit. Lastly, Equation 6 restricts VMs with double NUMAs from deploying both of its NUMAs on the same PM.

Note (1) each PM has two NUMAs; (2) w_k is a constant for each VM as determined by their types (Table 1). Thus, $\sum_{i,j} x_{k,i,j} = w_k$ (Equation 4) enforces that the actual NUMA allocation number

⁸Non-uniform memory access.

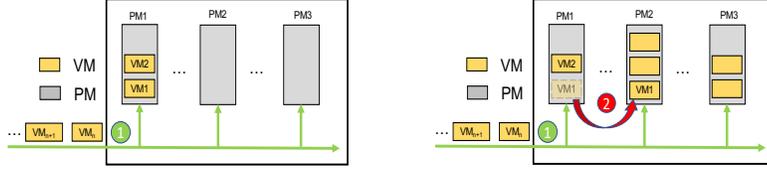


Figure 6: VM scheduling procedure. Figure 7: VM rescheduling procedure.

of VM k matches the desired configuration. When $w_k = 1$, Equation 4 constraints VM k to be deployed on one NUMA of a PM; when $w_k = 2$, Equation 4 constraints VM k to be deployed on both NUMAs of a PM. Note that deploying VM k on two NUMAs of two different PMs (each PM hosting a NUMA) violates $x_{k,i,0} = x_{k,i,1}, \forall k \in \{k | w_k = 2\}$ (Equation 6). Because $w_k \neq 0$, it guarantees each VM is deployed.

A.2 Background on (Deep) Reinforcement Learning

In this section, we give a brief overview on (deep) reinforcement learning while referring the readers to [18, 19] for a more detailed introduction.

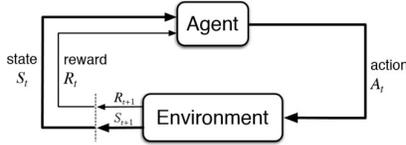


Figure 8: Illustration Diagram of Reinforcement Learning

Reinforcement Learning A standard Reinforcement Learning (RL) problem is typically characterized by a Markov Decision Process (MDP), where an agent continuously interacts with its environment, as depicted in Fig. 8. At each step (or period), the agent observes a state s_k from the environment and responds with an action a_k in accordance with the current policy π , i.e., $a_k \sim \pi(\cdot | s_k)$. The agent subsequently receives a reward $r_k = r(s_k, a_k)$ as feedback, based on the state and action. The transition to the next state depends solely on the current state and action, i.e., $s_{k+1} \sim \mathbb{P}(\cdot | s_k, a_k)$, where $\mathbb{P}(\cdot | \cdot, \cdot)$ represents the transition probability. In a tabular RL setting, both the state space S and the action space A are assumed to be finite and discrete.

Commencing with the initial state s_0 , the primary goal of the agent is to learn a policy that optimizes the so-called value function. This function is defined as the expected cumulative rewards received over K periods, i.e.,

$$J^\pi(s_0) := \mathbb{E} \left[\sum_{k=0}^{MNL} \gamma^k r(s_k, a_k) \middle| \pi, s_0 \right]. \quad (8)$$

Here, the expectation is taken over all possible trajectories under policy π . γ is a discount factor that affects how much weight is given to future rewards. Additionally, for any initial state-action pair (s_0, a_0) , the corresponding state-action value function (Q-function) is defined as

$$Q^\pi(s_0, a_0) := \mathbb{E} \left[\sum_{k=0}^K \gamma^k r(s_k, a_k) \middle| \pi, (s_0, a_0) \right]. \quad (9)$$

Deep Reinforcement Learning Deep RL, a subset of RL, integrates RL and deep learning. In tabular RL, the policy is directly optimized as a table of dimensions $|S| \times |A|$. However, the exponential scaling of $|S|$ and $|A|$ in high-dimensional state and action spaces makes traditional RL algorithms impractical. Deep RL addresses this issue by employing (deep) neural networks to represent the policy function (or other learned functions) and developing specialized algorithms for scalability.

A.3 Related Work

Connections to Bin Packing. The use of optimized placement mechanisms proved to be successful in a broad set of use cases, including production quality scenarios [20, 21] as well as transportation logistics [22, 23, 24, 25]. A typical solution exploits heuristics based on bin packing [26]. In fact, VM placement can be modeled as a bin-packing problem, where VMs and PMs are objects and bins, respectively. Bin packing typically involves packing a set of items into fixed-sized bins such that the number of bins required [22] or the total surface area is minimized [23, 24]. However, there are two notable differences. First, the problem of VM rescheduling concerns adjusting an initial assignment of VMs to PMs and has practical value in the industry, since some VMs might terminate causing the problem state to change. On the other hand, rebalancing items that are already packed in bins has received little attention in the context of other traditional bin-packing applications. Second, the total number of VMs and PMs in a data center can easily go into the range of several thousands or more [25], and is far more than the typical scale of bin packing problems, which typically involve no more than a few hundred items [27, 28, 29, 30].

RL for Optimization Problems. RL has been recently introduced to solve optimization problems, e.g., building ML compilers and optimizing neural network architectures [31]. In particular, RL is used to select branching variables or find cutting planes in the Branch-and-cut method [32, 33, 34, 35, 36, 37, 38]. Besides, RL can also be applied to existing heuristics for MIPs to further increase the quality of solutions [39, 40, 41]. However, the above approaches are not directly appropriate for the VM rescheduling task due to their poor computation complexity. Although they are designed to accelerate MIPs, even a state-of-the-art technique as POP [42] fails to deliver a satisfying solution within the second-level inference time limit required for the VM rescheduling task.

A.4 PPO

The algorithm VMR²L is developed based on the Proximal Policy Optimization (PPO) algorithm [13]. Recognized for its stability and robustness to various hyperparameters and network architectures [13], PPO has exhibited superior performance compared to Natural Policy Gradients (NPG) [43] and Trust Region Policy Optimization (TRPO) [44], and exhibited less bias compared to Q-learning [45]. Despite PPO’s high sample complexity, it becomes less of a concern due to our cost-effective simulator as detailed in Section B.1. In VMR²L, we employ two actors — a VM actor $\pi_{vm}(a_{vm}|s; \theta_{vm})$ and a PM actor $\pi_{pm}(a_{pm}|s, VM; \theta_{pm})$. Given the current state encoding s , the VM actor samples the next-step vm action, denoted as a_{vm} . Based on both s and a_{vm} , the PM actor samples a destination PM from the pm distribution, denoted as a_{pm} . Then, the critic $V(s; \theta_v)$ outputs an evaluation for the current state embedding s .

Below, we provide a brief overview for the training algorithm, with the complete pseudocode shown in Algorithm 1. In Lines 1-3, the algorithm first initializes the parameters of the VM actor, the PM actor, and the critic all with random weights. To simplify the implementation, we incorporate the critic network into the VM actor by appending a special VM with all ‘-1’ features and using its output as the critic score.

Each episode is made up of MNL steps, with each step representing a rescheduling action. During every rescheduling time step, VMR²L commences vm and pm selection with the original vm - pm mapping. It then generates the new vm - pm mapping by iteratively executing an action, as computed by the vm and pm actors, on the current vm - pm mapping until the rescheduling time step concludes (Line 6-16). The episode is terminated if the rescheduling time step length exceeds a pre-determined threshold.

Upon the completion of each episode, we compute the gradient for both actors and critics using their loss functions (Line 17). The gradient loss of the vm and pm actors is defined as the mean error between the advantage estimate and the logit of the corresponding sampled vm and pm actions ($\log p_{vm}$ and $\log p_{pm}$) across the episode. The Generalized Advantage Estimate for step i (GAE_i) is calculated as per the GAE-Lambda advantage [46] by

$$GAE_i = r_i + \gamma \cdot v_{i+1} - v_i + \gamma \cdot \lambda \cdot GAE_{i+1}, \quad (10)$$

where r_i , v_i represent the reward and the critic’s output at step i respectively. γ is the discount factor and λ is a smoothing parameter for variance reduction. The critic gradient loss is defined as the mean-square error between the reward to go and the critic’s output, v , across the epoch. The

Algorithm 1: Learning a VM rescheduling policy using a two-stage actor-critic algorithm.

Input: All VMs and PMs’ states
Output: A VM rescheduling agent

- 1 Initialize the VM actor parameters θ_{vm} with random weights ;
- 2 Initialize PM actor parameters θ_{pm} with random weights ;
- 3 Initialize critic parameters θ_v with random weights ;
- 4 **for** $episode = 0, 1, \dots, M$ **do**
- 5 Randomly sample the initial state s_0 from all available PM-VM mappings in the train set;
- 6 buffer.clear() ;
- 7 **for** *rescheduling step* $t = 0, 1, \dots, MNL$ **do**
- 8 $\log p_{t,vm} \leftarrow \pi(a|s_t, \theta_{vm})$;
- 9 $a_{t,vm} \leftarrow p_{t,vm}.sample()$;
- 10 $\log p_{t,pm} \leftarrow \pi(a|s_t, a_{t,vm}, \theta_{pm})$;
- 11 $a_{t,pm} \leftarrow p_{t,pm}.sample()$;
- 12 $v_t \leftarrow V(s_t, \theta_v)$;
- 13 buffer.append($\log p_{t,vm}, \log p_{t,pm}, a_{t,vm}, a_{t,pm}, v_t$) ;
- 14 // Compute gradients wrt. *vm* actor, *pm* actor, and critic loss ;
- 15 $d\theta_{vm}, d\theta_{pm}, d\theta_v \leftarrow ComputeLoss(buffer)$;
- 16 Perform update of θ_{vm} using $d\theta_{vm}$, θ_{pm} using $d\theta_{pm}$, and θ_v using $d\theta_v$;

reward-to-go is calculated by applying the discount factor to the intermediate rewards. Finally, in Line 18, we update the parameters of the actor and critic networks with the calculated gradients.

B Experiment Details

We conduct extensive experiments to answer:

- How far is VMR²L from the optimal solution? (§ B.3)
- How well does VMR²L generalize to larger data centers? (§ B.4)
- How well does VMR²L generalize to different objectives? (§ B.5)
- How well does VMR²L generalize to different MNLs? (§ B.6)

B.1 Simulator and Datasets

While DRL can be very powerful, its main drawback is the amount of training data required [47, 48]. In light of this, we design a simulator for the VM rescheduling task. The simulator follows the OpenAI Gym environments [49] including specific file hierarchy and function abstractions. Given an existing VM-PM mapping and a rescheduling action, we can directly calculate the change in FR caused by the action. Thus, during training, VMR²L only needs to interact with the simulator instead of with the real environment, which drastically lowers the amount of real-world data required to train the agent.

As for the datasets, we have two seed initial mappings from an industry-scaled real cloud data center – one medium dataset with 2089 VMs and 280 PMs, and one large dataset with 4546 VMs and 1176 PMs. Note that the RL agent must be able to generalize to VM-PM mappings unseen during training and a dynamic number of VMs at deployment time. To better evaluate the agent’s performance under various initial mappings, we generate 4400 initial mappings with different numbers of VMs for both the medium and the large datasets. Each mapping is generated by removing the existing VMs on each PM and randomly scheduling some of them to any PM that can fit them. We generate three versions of the middle dataset with low, middle, and high workloads (different remaining CPU resources). We leverage 4000 datasets for training, 200 datasets for both validation and test. Both the simulator and datasets are available to the research community.

B.2 Existing Baseline Algorithms

Existing baselines can be summarized into six categories: heuristics (e.g., greedy, α -VBPP), optimization algorithms (e.g., MIP), approximate algorithms (e.g., POP), search-based algorithms (e.g., MCTS), deep learning-based (e.g., Decima), and hybrid methods (e.g., NeuPlan). We compare with at least one representative algorithm from each category.

MIP Algorithm: We formulate the VM rescheduling problem as an optimization problem (Appendix A.1) and solve this problem with an off-the-shelf MIP solver such as Gurobi [9] and CPLEX [10], which can find an optimal solution through algorithms of branch & bound, cutting planes, etc. In our experiments, we use the former.

Greedy Algorithm: To obtain a feasible solution within a short time frame, greedy algorithms [50] are often used. They normally include two stages: filtering and scoring. In the filtering stage, we calculate the change in FR for each VM if it is removed from its source PM, and only select the VM candidate that corresponds to the most significant drop in FR. In the scoring stage, we calculate the change in FR if the selected VM is migrated to each of the eligible PMs. We then greedily assign the selected VM to the PM that leads to the largest drop in FR. The above two stages are repeated until the migration number limit is reached.

Vector Bin Packing Problem (α -VBPP): We generalize the VBPP [26] algorithm for initial scheduling to rescheduling. We first divide the entire episode into MNL/α stages. During each stage, we greedily remove α number of VMs that lead to the most fragments, and then apply VBPP to treat them as incoming VMs. We carefully tune α (10 in our case) to achieve the best FR reduction.

Partitioned Optimization Problems (POP): The POP method [42] solves the optimization problem formulated in Section A.1 by randomly splitting the large-scale VM rescheduling problem into smaller subproblems (each contain a subset of VMs and PMs) and coalescing the resulting sub-rescheduling solution into a global rescheduling solution for all VMs. We perform a grid search for the best POP parameter that balances the FR and time on both medium and large datasets.

Monte-Carlo Tree Search (MCTS) [30]: As traditional search-based methods need to perform multiple rollouts during inference time to achieve a good performance, we use DDTs [30] to prune the search space.

Decima [51]: Decima uses RL and neural networks to learn the VM rescheduling algorithm. A graph neural network is leveraged to encode the VM and PM information in a set of graph embedding vectors to process a large amount of state information. Decima balances the size of the action space and the number of actions required by decomposing VM rescheduling decisions into a two-dimensional action, which outputs (i) the VM that needs to migrate, and (ii) an upper number of PM subsets to choose as the destination.

NeuPlan [52]: NeuPlan uses a two-stage hybrid approach to address MIP’s scalability issue. In the first stage, an RL-based method takes in the problem in the form of a graph and generates the first few steps of VM rescheduling to prune the MNL search space. In the second stage, it uses a MIP solver to find the optimal VM migration given the remaining MNL. A relax factor α is used to control the size of the MNL space to explore by MIP.

B.3 Performance on the Medium Dataset

Fragment Rate. The results in Fig. 5(a) reveal that VMR²L can achieve +14.48%, +17.60%, +22.37%, +23.71%, 23.84%, and +34.17% performance gain when compared to POP, NeuPlan, MCTS, GA, α -VBPP and Decima respectively, when the task is to perform 50 migrations on the medium dataset. Notably, VMR²L is merely 2.67% behind the optimal MIP solution (0.2953 vs. 0.2859). It is worth noting that the performance gap between VMR²L and the near-optimal MIP does not increase with more MNLs.

Although α -VBPP can reduce FR with more MNL, its performance is inferior to that of GA. This is because α -VBPP only removes the α worst VMs for each stage based on a single timestep, failing to consider future opportunities to replace them back and achieve even higher FR reduction. POP fails to achieve a good performance since it still relies on MIPs to solve each subproblem. To meet the second-level latency requirement, we must divide the problem into many subproblems, causing its solutions to be only locally optimal. On the other hand, Decima reduces the huge action space by

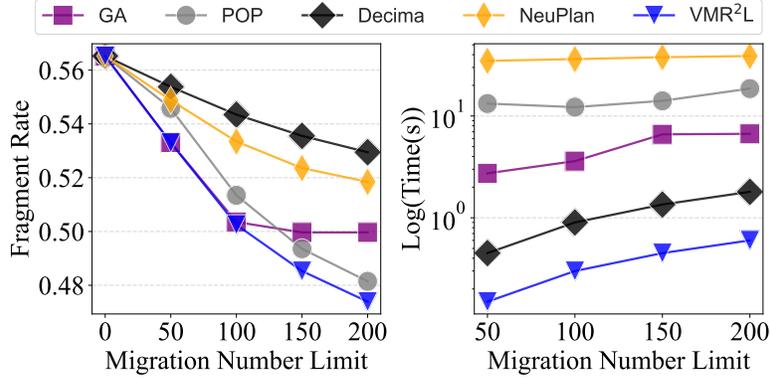


Figure 9: Fragment rate and inference time between different methods on the large dataset.

limiting the PM actor to only select from a subset of PMs, but the subsampling of PMs is completely random, as opposed to our solution. While MCTS with DDTS uses neural networks to prune the search space, it still requires a significant number of rollouts to achieve stable performance. Lastly, NeuPlan is able to achieve a low FR, since it solves a large subproblem entirely with MIP. We implement NeuPlan as follows: if MNL is less than 20 steps, MIP is used to solve the optimization problem. If MNL is larger than 20, the first 20 VMs are migrated by MIP and the remaining VMs will be handled by RL. Notice that NeuPlan’s FR is higher than VMR²L after MNL = 20, since after this step NeuPlan relies entirely on its RL agent, which fails to learn a good policy given such a huge state and action space.

Inference Time. From Fig. 5(b), we can see that the solution time of GA, α -VBPP, Decima and VMR²L is less than OSG. VMR²L can generate one trajectory within 0.15s when MNL = 50. In comparison, MIP requires 50.55 minutes to provide the optimal solution. The running time of POP can be adjusted by setting how many subproblems to divide into. To meet the latency requirement of the VM rescheduling task, we set this number to 16 so POP can deliver a solution within 1.94s. Both GA and α -VBPP are greedy algorithms and can provide the solution within 1s. Decima requires 0.45s, which is at a roughly same scale as VMR²L, since both use end-to-end deep RL. Lastly, NeuPlan takes 34.8s to yield the final solution since it still needs MIP to solve 20 steps.

B.4 Scalability to the Large Dataset

We conduct experiments on a large dataset with 4546 VMs and 1176 PMs to analyze the scalability of VMR²L. Fig. 9 shows the FR and time performance of different baseline methods with the MNL from 50 to 200. The MIP is not included in this experiment since we cannot get a solution within 50 minutes. As it turns out, VMR²L again achieves better performance than the baselines on both FR and solution time.

From the left subfigure in Fig. 9, we can see that VMR²L can achieve average 2.01%, 2.08%, 6.14% and 7.87% and max 2.34%, 5.15%, 8.6%, and 10.5% performance gain compared with POP, GA, NeuPlan, and Decima respectively. Decima and NeuPlan both cannot achieve good FR performance on such a large dataset. NeuPlan mostly relies on RL’s solution since MIP cannot work with MNL larger than 20. GA gradually stops reducing FR after 100 steps. POP and VMR²L continue to reduce FR even at MNL 200. POP achieves worse FR than GA before 100 and better FRs for larger MNLs.

The right subfigure in Fig. 9 shows the running time to generate a migration solution with MNL set from 50 to 200. GA, POP, Decima, NeuPlan, and VMR²L can generate a solution within 4.91s, 14.53s, 1.125s, 37.23s, and 0.375s on average. GA increases the time when MNL is less than 150. GA calculates the score of all the VMs and PMs as a metric to evaluate their migration value. The calculation time increases with the number of VMs and PMs. POP costs more time with a bigger MNL. VMR²L increases time linearly from MNL 50 to 200. For VMR²L, the neural network inference time will only increase minimally with the number of PMs and VMs. Decima needs three times than VMR²L since the GNN needs to encode the VM-PM information. NeuPlan needs MIP to solve 20 MNL.

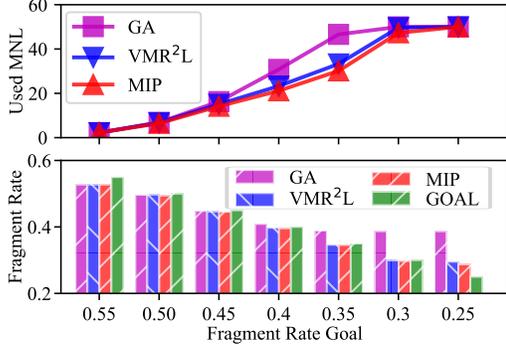


Figure 10: MNL Performance under Different FR Goals.

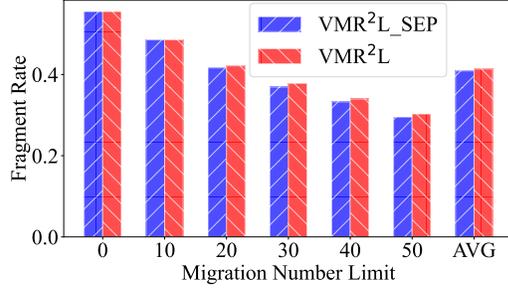


Figure 11: Fragment Rate Gap between VMR²L and VMR²L_{SEP}.

B.5 Different Service Objectives

VMR²L’s flexibility enables it to learn different policies depending on the high-level objective. We now consider a new objective: given FR goals, we would like to minimize the number of migrations needed to reach the FR goals. As seen in Fig. 10, the top subfigure displays the used MNL, while the bottom subfigure shows the FR, both sharing the x-axis with the FR goals. In general, as the FR goal decreases, GA requires significantly more migration steps than both MIP and VMR²L. Note that none of the three methods can achieve the FR goal of 0.25, since the optimal FR is 0.2859 at MNL 50. To summarize, the used MNLs of GA, VMR²L, and MIP increase with lower FR goals. Across all FR goals, MIP and VMR²L achieve 14.77% and 11.11% less MNL than GA, respectively. VMR²L performs similarly to MIP, with a difference of only 3.66%, but with a millisecond-level inference time. On the other hand, GA is stuck at a FR goal of 0.4 since it only optimizes for the next-step performance, instead of the long-term performance.

B.6 Generalizing to Different MNLs

In practice, the required migration number limit (MNL) can constantly vary due to changing business requirements such as target fragment rates. We show that we can readily achieve good performance by only training one VMR²L agent with MNL = 50, and deploying it for MNLs = {10, 20, 30, 40, 50}. For comparison, we train a separate VMR²L agent for each MNL, which we denote as VMR²L_{SEP}. As shown in Fig. 11, VMR²L performs only marginally worse than VMR²L_{SEP} with an average FR performance gap of 1.16%. This suggests that the VMR²L agent trained with a large MNL can be readily applied to the tasks with smaller MNLs. It avoids the overhead of maintaining a separate VMR²L agent for each MNL.

C Architecture Overview

To better incorporate the proposed sparse local-attention module, we modify the vanilla transformer architecture as follows. The model is composed of several attention blocks, where each block includes three stages as shown in Fig. 4:

1. All PMs and VMs exchange information if they belong to the same tree via sparse local-attention.
2. Each PM attends to other PMs’ updated embeddings and each VM attends to other VMs’ updated embeddings with self-attention.
3. The new VM embeddings are allowed to attend to the new PM embeddings through VM-PM attention.

At the end of the three stages, each machine is then further processed by two point-wise dense layers with ReLU activation and layer norm [53]. The updated embeddings for each machine is then feed into the next block and the process repeats. Finally, the VM embeddings from the final block are linearly projected into a set of logits followed by the Softmax operation[54] to generate the probability of each selected VM.

We also add the number of MNL steps left as an additional input to be processed together with the PMs. As for the critic, we add a node with all -1’s to be processed together with the VMs. The output embedding from the final block is linearly projected into the critic score. Note that the PM embeddings are updated block-by-block together with the VM embeddings, as it encourages the PMs and VMs to better coordinate and update gradually from low level to high level [55].

As for the PM actor architecture, we adapt the vanilla transformer encoder-decoder module, since we can directly inject the graph information by including the VM and PM embeddings from the VM actor as input. We only feed in the selected VM candidate to the encoder, while the decoder still takes in all the PMs. Additionally, we also add the VM-PM attention score from stage 3 for the selected VM, since the score implies why the VM actor selects the chosen VM and which PMs it attends to in the process. The output embeddings of each PM is linearly projected into a logit. Based on the selected VM, we mask out all the illegal PMs by setting their logits to be $-\infty$. The remaining logits are translated into the probability of selecting each PM as the destination PM.

C.1 Sparse Attention Details

In this section, we delve into the detailed formulation of *sparse attention*.

Why We Need Tree-level Features? Consider a PM with 2 CPUs left. It contains a VM with 4 CPUs and another VM with 2 CPUs. Suppose a second PM has a fragment size of 8 while hosting a VM with 8 CPUs. In order to minimize the total 16-core fragments, an ideal approach would be to first remove the two VMs with 2 and 4 CPUs from the first PM, and then reassign the VM with 8 CPUs from the second PM to the first PM. However, if we merely include the source PM’s features in each of the VM’s features and feed them into the vanilla transformer model, there will not be sufficient information for the two actors to take the above actions. Instead, each VM must also be aware of the other VMs that are hosted on the same PM, which is not possible in the vanilla transformer model. In fact, each PM can be viewed as a tree of depth one, where the PM acts as the root node and the VMs it hosts act as the leaf nodes. In order to allow every VM to recognize which other VMs are hosted on the same PM, we propose to include an additional stage of *sparse local-attention* within each PM tree, i.e., we only allow PMs and VMs to attend to each other if and only if they belong to the same tree.

Conceptual Overview. Fundamentally, attention can be understood as a trainable dictionary involving queries, keys, and values. Given an embedding vector (a VM or a PM) that requires an update, and a set of reference vectors (selected VMs and/or PMs), we project the vector-to-be-updated into a query vector. Concurrently, we project all reference vectors into key and value vectors. We then compute the similarity score between the query vector and each key vector. Based on these scores, we update the target embedding vector as the weighted sum of the corresponding values.

The term “VM self-attention” refers to the process of updating each VM’s embedding vector using all VM’s embedding vectors (including its own) as a reference. “PM self-attention” operates similarly for PMs. “VM-PM attention” involves updating each VM’s embedding vector using all PM’s embedding vectors as references. Lastly, the proposed tree-level *sparse attention* involves updating each VM or PM using only other VMs or PMs within the same tree—that is, those affiliated with the same PM.

Attention Formulation. Formally, let $\mathcal{V} = \{v_i\}$ and $\mathcal{P} = \{p_i\}$ denote the set of feature vectors for each VM and each PM, respectively. Consider an embedding vector $x_i \in \mathbb{R}^{d_{\text{model}}}$ that we aim to update, which is projected linearly⁹ from $v \in \mathcal{V}$ or $p \in \mathcal{P}$. Let (y_1, \dots, y_n) be a set of reference vectors that could be a combination of embeddings of $v \in \mathcal{V}$ and $p \in \mathcal{P}$, including x_i itself. Instead of directly operating in the feature space $\mathbb{R}^{d_{\text{model}}}$, we project these vectors into an embedding space \mathbb{R}^{d_k} .

An attention function updates the target vector x_i by first projecting it into a query $Q_i = x_i W^Q$ using a linear transformation, where $W^Q \in \mathbb{R}^{d_{\text{model}} \times d_k}$ is a learnable weight matrix. Similarly, each reference vector y_j is projected into a key $K_j = y_j W^K$ and a value $V_j = y_j W^V$, where

⁹A linear projection is necessary since v and p can have different number of features.

$W^K, W^V \in \mathbb{R}^{d_{\text{model}} \times d_k}$ ¹⁰. Note that we can perform the attention function on a set of queries in parallel by stacking Q_i 's into a larger matrix Q ¹¹.

The compatibility function, measuring the similarity between Q_i and K_j , is typically implemented as the dot product of these two vectors. However, when d_k is large, the magnitude of dot products tends to increase, leading to a high similarity for only a few keys and a marginal similarity for the rest of the references. This has been known to cause extremely small gradients [14]. To address this, we scale the dot product by $\frac{1}{\sqrt{d_k}}$,

$$\text{Attention}(Q, K, V) = \text{softmax} \left(\frac{QK^T}{\sqrt{d_k}} \right) V. \quad (11)$$

For the proposed *sparse attention*, we introduce an additional mask M . Here, $M_{i,j} = -\infty$ if x_i and y_j are not part of the same PM tree and zero otherwise.

$$\text{Sparse-Attention}(Q, K, V) = \text{softmax} \left(\frac{QK^T}{\sqrt{d_k}} + M \right) V. \quad (12)$$

Attention Block. In actual applications, we often utilize multiple sets of W^Q, W^K , and W^V . Each set, referred to as a single attention head, enables the model to extract pertinent information from a distinct representation space. We combine the resulting matrices from each attention head in multi-head attention by concatenation.

$$\text{Multi-Head}(Q, K, V) = \text{Concat}(\text{head}_1, \dots, \text{head}_h) W^O, \quad (13)$$

where head_i represents the output of each attention head from Equation 12 using W_i^Q, W_i^K , and W_i^V . $W^O \in \mathbb{R}^{h \cdot d_k \times d_{\text{model}}}$ is an additional learnable weight parameter.

Within each attention block, we consecutively execute the process for tree-level *sparse attention*, VM/PM self-attention, and VM-PM attention. Following these three attention submodules, we further process each updated embedding vector with a fully connected feed-forward network, which operates independently on each embedding vector.

$$\text{FFN}(x_i) = \text{GELU}(xW_1 + b_1) W_2 + b_2, \quad (14)$$

where GELU refers to the Gaussian Error Linear Units activation function [56]. The output of this process is then supplied to the next attention block, and the process is iteratively repeated. For more in-depth details regarding the attention module, we encourage readers to consult the paper by Vaswani et al. [14].

¹⁰ W^V can have different dimensions than W^K . We keep them the same here for simplicity.

¹¹As we use the same W^Q, W^K, Q^V for all vectors, the total number of weight parameters remains independent of the number of VMs or PMs in the problem, making this approach suitable for scaling to large data centers.

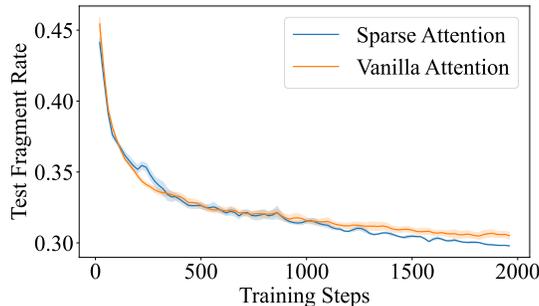


Figure 12: FR Performance of VMR²L with Sparse versus Vanilla Attention.

Table 2: A list of hyperparameters.

RL Parameters	Value	General Parameters	Value
Clip_coefficient	0.1	Total_iterations	4e6
Discount_factor	0.99	Learning_rate	2.5e-4
PPO update_epochs	8	Attention_blocks	2
PPO minibatches	128	Attention_heads	2
GAE_Lambda	0.95	Transformer d_{ff}	32
Risk-Seeking K	3	Transformer d_{hidden}	64

Overall Architecture. The architecture incorporates the remaining MNL steps as supplementary input for both the VM and PM actors, processed concurrently with the PMs. This inclusion is pivotal for reinforcement learning in scenarios with fixed-length episodes as demonstrated in prior studies[57]. As for the critic, a node containing all -1’s is incorporated and processed alongside the VMs. The output embedding from the terminal block undergoes a linear projection to yield the critic score. It’s important to note that the PM embeddings undergo iterative updates block-by-block in tandem with the VM embeddings. This approach fosters better coordination between PMs and VMs, enabling a gradual transition from low to high-level updates [55].

Ablation on Sparse Attention. In Fig. 12, we present an ablation study performed on a medium dataset with $MNL = 50$ to evaluate the impact of the *sparse attention* module in terms of FR. As a consequence of the supplementary parameters involved, *sparse attention* exhibits a longer convergence time but successfully achieves a lower FR. This observation validates that *sparse attention* empowers VMR²L to extract tree-level information, which is absent in vanilla attention. A visual illustration of this concept is provided in Section C.2.

C.2 A Case Study on Sparse Attention

We build a visualization tool that allows end-users to directly input a trained agent and receive the detailed migration actions at each step in the form of a gif file. This allows end-users to better interpret and analyze the actions taken by the agent. In Fig. C.1, we provide a case study on how the proposed sparse attention allows the agent to optimize for long-term rewards.



Figure 13: At step 38, VMR²L removes a VM with 4 CPUs to achieve zero fragments on the destination NUMA. However, this also creates four fragments on the source NUMA, leading to a net reward of zero at this step. Then at step 39, it reschedules another VM with 4 CPUs to achieve zero fragments on both the source and the destination NUMAs. Although the agent receives a reward of zero at step 38, *sparse attention* allows the agent to realize that there were two VMs with 4 CPUs on the source NUMA, so the agent can reassign the second VM with 4 CPUs elsewhere in order to achieve zero fragments on the source NUMA as well.