Proactive Optimization of Virtual Machine Placement Using Predictive Models Based on Time Series

Naby Doumbouya © EPROAD, UR 4669 Université de Picardie Jules Verne Amiens, France e-mail: ndoumbouya@u-picardie.fr Mhand Hifi EPROAD,UR 4669 Université de Picardie Jules Verne Amiens, France e-mail: mhifi@u-picardie.fr

Abstract—To overcome the limitations of traditional reactive VM placement strategies, which often struggle with dynamic workload variations, this paper introduces an innovative proactive approach based on predictive time series analysis. By evaluating the ARIMA, LSTM, and Prophet models, we assess their effectiveness in accurately forecasting VM workload fluctuations, thereby minimizing unnecessary migrations, reducing energy consumption, and lowering operational costs. These predictions are then integrated into an advanced optimization algorithm to determine optimal VM placement in anticipation of workload spikes, leading to significant improvements in system performance, stability, and quality of service across distributed data centers.

Keywords-Cloud computing; virtual machine placement; time sequences; ARIMA; LSTM; Prophet; cloudsim; Time series forecasting; Proactive optimization.

I. INTRODUCTION

The proactive optimization of virtual machine (VM) placement is crucial for enhancing resource efficiency in cloud environments, where dynamic resource management is essential to ensure optimal performance. Cloud systems often face fluctuating demands, making predictive methods imperative for effective resource allocation.

Time series forecasting models such as ARIMA, Prophet, and LSTM offer various approaches to anticipate future resource needs. ARIMA is well suited for stationary datasets, while Prophet excels in handling seasonal data. LSTM provides flexibility for modeling complex patterns in dynamic contexts [1].

Cloud resource management relies on reactive and proactive approaches. The reactive approach adjusts resources based on current demand, while the proactive approach leverages historical data to predict future needs and optimize resource utilization. Proactive resource allocation has become a major research topic in cloud computing, aiming to optimize resource management and utilization [2].

Our approach combines ARIMA, LSTM, and Prophet to generate accurate workload forecasts. These forecasts are used in an optimization algorithm that minimizes energy consumption and reduces VM migrations while respecting server capacity constraints. By anticipating future workloads, our method enables more efficient resource allocation and improves system performance. The remainder of this paper is organized as follows. Section II reviews the state of the art in VM placement and predictive optimization. Section III formulates the problem as a multi-objective optimization model incorporating energy consumption and migration cost. Section IV describes the proposed methodology, combining time series forecasting with hybrid prediction weighting and linear programming. Section V presents the experimental setup and results obtained on the CloudSim dataset. Section VI highlights our main contributions in terms of prediction accuracy and optimization efficiency. Finally, Section VII concludes the paper and outlines future research directions.

II. BACKGROUND

Classical approaches to VM placement optimization, such as First Fit, Best Fit, and genetic algorithms, do not account for workload forecasting. Time series models (ARIMA, LSTM, Prophet) enable proactive decision making but struggle with sudden workload fluctuations [3].

Advanced optimization algorithms like Harris Hawk Optimization (HHO) outperform traditional methods, achieving a 27% reduction in energy consumption and a 17% increase in resource utilization [4]. Hybrid approaches combining optimization algorithms and machine learning are essential for efficient cloud resource allocation.

ARIMA excels with stationary series, and Prophet performs well with seasonal data, but both are limited in handling nonlinear variations and unexpected spikes [5]. LSTM captures long-term dependencies but requires significant computational resources [5]. Hybrid models like TempoScale enhance forecast accuracy and system stability [3].

Traditional optimization approaches face challenges related to resource heterogeneity and workload variability [6]. Hybrid solutions combining classical methods with artificial intelligence are promising for optimizing VM placement and ensuring proactive resource management [7].

This evolution highlights the need for integrated strategies to maximize cloud infrastructure efficiency. Future research should further explore hybrid approaches and assess their practical implementation for responsive and sustainable resource allocation.

III. PROBLEM MODELING

We model the VM placement problem as a **multi-objective optimization problem** with the following objectives:

A. Minimize Energy Consumption

The total energy consumption is calculated as the sum of the energy consumed by each server, weighted by the CPU usage of the VMs placed on it [8]:

$$E_{total} = \sum_{i=1}^{m} \sum_{j=1}^{n} x_{ij} \cdot E(M_i) \tag{1}$$

where x_{ij} is a binary variable indicating whether VM j is placed on server i, and $E(M_i)$ is the energy consumption of server i.

B. Minimize VM Migrations

When VMs change their host between time steps t-1 and t, it incurs a cost. We define $x_{ij}^{(t)}$ and $x_{ij}^{(t-1)}$ as binary variables indicating placement at time t and t-1, respectively.

The migration cost function is given by:

$$C_{mig} = \sum_{i=1}^{m} \sum_{j=1}^{n} \left| x_{ij}^{(t)} - x_{ij}^{(t-1)} \right| \cdot D_{mig}(V_j)$$
(2)

where $D_{mig}(V_j)$ denotes the migration cost (e.g., based on memory size or state size) of VM V_j [9], [10]. This penalty discourages unnecessary movement and ensures SLA stability.

Constraints The total resource usage of VMs on each server must not exceed the server's capacity:

$$\sum_{j=1}^{n} x_{ij} \cdot R(V_j) \le C(M_i), \quad \forall i$$
(3)

where $R(V_j)$ is the resource requirement of VM j, and $C(M_i)$ is the capacity of server i. In addition, a VM must be placed on exactly one server:

$$\sum_{i=1}^{m} x_{ij} = 1, \quad \forall j \tag{4}$$

IV. PROPOSED METHODOLOGY

Our methodology involves three steps: data preprocessing, workload prediction, and VM placement optimization.

A. Data preprocessing and Workload prediction

We use time series data from cloudsim, aggregated at 5 minutes intervals, and train predictive models (ARIMA, LSTM, Prophet) using a sliding window approach.

B. VM Placement Optimization

The predicted workloads are used as input to a **linear programming (LP)** optimization problem. The objectives are two fold: (1) minimize energy consumption and (2) minimize the cost of VM migrations, while respecting server capacity constraints.

Predictive-Aware Placement Strategy: The predicted workload for each VM is generated using a weighted combination of ARIMA, LSTM, and Prophet forecasts:

$$\hat{L}_j = \sum_{m \in \{\text{ARIMA, LSTM, Prophet}\}} w_m \cdot \hat{L}_{j,m}$$
(5)

with weights w_m based on the inverse of each model's error (RMSE + MAE), normalized:

$$w_m = \frac{1}{\text{RMSE}_m + \text{MAE}_m + \epsilon} \bigg/ \sum_{m'} \frac{1}{\text{RMSE}_{m'} + \text{MAE}_{m'} + \epsilon}$$
(6)

These predicted loads are injected into the optimization model to guide placement before overloads occur.

Solver: The combined multi-objective function is minimized using a weighted-sum scalarization:

$$\min\left(\alpha \cdot E_{total} + \beta \cdot C_{mig}\right) \tag{7}$$

where α and β are tunable coefficients reflecting the tradeoff between energy efficiency and migration stability, as recommended in [11], [12].

This LP model is implemented with PuLP in Python. The forecast-driven optimization allows proactive VM placement while balancing operational cost and performance constraints.

C. General scheme



Figure 1. Overview of the proposed method

V. EXPERIMENTS AND RESULTS

We evaluate our approach on the **CloudSim Dataset**. Key findings include:

- **Energy Efficiency**: Energy consumption is reduced by 15%, with an additional 5% improvement from the RMSE + MAE weighting.
- VM Migrations: Migrations are reduced by 20%, with a further 10% reduction due to the RMSE + MAE weighting.
- **Robustness**: The RMSE + MAE weighting improves stability under workload variations.

Courtesy of IARIA Board and IARIA Press. Original source: ThinkMind Digital Library https://www.thinkmind.org

VI. OUR CONTRIBUTION

Our work makes the following key contributions:

- Hybrid Prediction Model: We propose a novel approach that combines the strengths of three predictive models ARIMA, LSTM, and Prophet using a weighted average based on both RMSE and MAE. This hybrid approach improves prediction accuracy and robustness compared to using individual models.
- **Proactive Optimization Framework**: Unlike traditional reactive methods, our framework uses predicted workloads to proactively optimize VM placement, reducing energy consumption and unnecessary migrations.
- **RMSE + MAE Weighting Scheme**: We introduce a weighting scheme that balances the impact of large errors (RMSE) and average errors (MAE), leading to more stable and reliable predictions. This scheme significantly improves the robustness of the optimization process.
- Energy and Migration Efficiency: Our approach demonstrates a 15% reduction in energy consumption and a 20% reduction in VM migrations compared to traditional methods, with further improvements achieved through the RMSE + MAE weighting.

VII. CONCLUSION AND PERSPECTIVES

Our proactive VM placement strategy, which integrates hybrid time series forecasting with a weighted integer linear optimization model, has proven effective in reducing energy consumption and eliminating unnecessary migrations. The weighted combination of ARIMA, LSTM, and Prophet based on RMSE and MAE metrics significantly enhances forecasting robustness.

The results obtained demonstrate that our method offers a reliable and efficient trade-off between resource optimization and service stability. The approach remains scalable, and the low computation time allows real-time or near-real-time decision-making.

As future work, we plan to:

Extend the optimization model to a distributed and federated cloud environment, where coordination among data centers is required;

Incorporate network-related constraints such as bandwidth, latency, and routing cost;

Integrate adaptive dynamic weights based on SLA policies and QoS priorities;

Experiment on real-world traces such as the **Google Cloud Cluster Trace Dataset** to evaluate the model at scale.

REFERENCES

- J. Chen, Y. Wang, and T. Liu, "A proactive resource allocation method based on adaptive prediction of resource requests in cloud computing", *EURASIP Journal on Wireless Communications and Networking*, vol. 2021, p. 24, 2021. DOI: 10.1186/ s13638-021-01912-8.
- [2] T. Kamble, S. Deokar, V. S. Wadne, D. P. Gadekar, and H. B. Vanjari, "Predictive resource allocation strategies for cloud computing environments using machine learning", *Journal of Electrical Systems*, vol. 19, no. 2, pp. 68–77, 2023.
- [3] L. Wen, M. Xu, A. N. Toosi, and K. Ye, "Temposcale: A cloud workloads prediction approach integrating short-term and long-term information", in 2024 IEEE 17th International Conference on Cloud Computing (CLOUD), 2024.
- [4] H. S. Madhusudhan, T. S. Kumar, P. Gupta, and G. McArdle, "A harris hawk optimisation system for energy and resource efficient virtual machine placement in cloud data centers", *PLOS ONE*, vol. 18, no. 8, e0289156, 2023. DOI: 10.1371/ journal.pone.0289156.
- [5] S. Yadav, "A comparative study of arima, prophet and lstm for time series prediction", *Journal of Artificial Intelligence*, *Machine Learning and Data Science*, vol. 1, no. 1, pp. 1813– 1816, 2022. DOI: 10.51219/JAIMLD/sandeep-yadav/402.
- [6] A. Abdelaziz, M. Anastasiadou, and M. Castelli, "A parallel particle swarm optimisation for selecting optimal virtual machine on cloud environment", *Applied Sciences*, vol. 10, p. 6538, 2020. DOI: 10.3390/app10186538.
- [7] A. Poghosyan *et al.*, "An enterprise time series forecasting system for cloud applications using transfer learning", *Sensors*, vol. 21, p. 1590, 2021. DOI: 10.3390/s21051590.
- [8] X. Li, Z. Qian, S. Lua, and J. Wu, "Energy efficient virtual machine placement algorithm with balanced and improved resource utilization in a data center", *Mathematical and Computer Modelling*, vol. 58, pp. 1222–1235, 2013. DOI: 10.1016/j.mcm.2013.02.003.
- [9] M. Masdari and H. Khezri, "Efficient vm migrations using forecasting techniques in cloud computing: A comprehensive review", *Cluster Computing*, vol. 23, pp. 2629–2658, Jan. 2020. DOI: 10.1007/s10586-019-03032-x.
- [10] M. T. et al, "Borg: The next generation", in *Fifteenth European Conference on Computer Systems (EuroSys '20)*, 2020, pp. 1–14. DOI: 10.1145/3342195.3387517.
- [11] C. Vijaya and P. Srinivasan, "Multi-objective meta-heuristic technique for energy efficient virtual machine placement in cloud computing data centers", *Informatica*, vol. 48, pp. 1–18, Jun. 2024. DOI: 10.31449/inf.v48i6.5263.
- [12] R. Keshri and D. P. Vidyarthi, "Energy-efficient communicationaware vm placement in cloud datacenter using hybrid aco-gwo", *Cluster Computing*, vol. 27, pp. 13047–13074, 2024. DOI: 10.1007/s10586-024-04623-z.