



## UvA-DARE (Digital Academic Repository)

### RAILS: Risk-Aware Iterated Local Search for Joint SLA Decomposition and Service Provider Management in Multi-Domain Networks

Hsu, Cyril Shih-Huan; Papagianni, Chrysa; Grosso, Paola

**DOI**

[10.1109/HPSR64165.2025.11038864](https://doi.org/10.1109/HPSR64165.2025.11038864)

**Publication date**

2025

**Document Version**

Final published version

**Published in**

2025 IEEE 26th International Conference on High Performance Switching and Routing

**License**

Article 25fa Dutch Copyright Act (<https://www.openaccess.nl/en/policies/open-access-in-dutch-copyright-law-taverne-amendment>)

[Link to publication](#)

**Citation for published version (APA):**

Hsu, C. S.-H., Papagianni, C., & Grosso, P. (2025). RAILS: Risk-Aware Iterated Local Search for Joint SLA Decomposition and Service Provider Management in Multi-Domain Networks. In *2025 IEEE 26th International Conference on High Performance Switching and Routing: HPSR 2025 : Suita, Osaka, Japan, 20-22 May 2025* (pp. 174-179). IEEE. <https://doi.org/10.1109/HPSR64165.2025.11038864>

**General rights**

It is not permitted to download or to forward/distribute the text or part of it without the consent of the author(s) and/or copyright holder(s), other than for strictly personal, individual use, unless the work is under an open content license (like Creative Commons).

**Disclaimer/Complaints regulations**

If you believe that digital publication of certain material infringes any of your rights or (privacy) interests, please let the Library know, stating your reasons. In case of a legitimate complaint, the Library will make the material inaccessible and/or remove it from the website. Please Ask the Library: <https://uba.uva.nl/en/contact>, or a letter to: Library of the University of Amsterdam, Secretariat, P.O. Box 19185, 1000 GD Amsterdam, The Netherlands. You will be contacted as soon as possible.

*UvA-DARE is a service provided by the library of the University of Amsterdam (<https://dare.uva.nl>)*

# RAILS: Risk-Aware Iterated Local Search for Joint SLA Decomposition and Service Provider Management in Multi-Domain Networks

Cyril Shih-Huan Hsu  
Informatics Institute  
University of Amsterdam  
Amsterdam, The Netherlands  
s.h.hsu@uva.nl

Chrysa Papagianni  
Informatics Institute  
University of Amsterdam  
Amsterdam, The Netherlands  
c.papagianni@uva.nl

Paola Grosso  
Informatics Institute  
University of Amsterdam  
Amsterdam, The Netherlands  
p.grosso@uva.nl

**Abstract**—The emergence of the fifth generation (5G) technology has transformed mobile networks into multi-service environments, necessitating efficient network slicing to meet diverse Service Level Agreements (SLAs). SLA decomposition across multiple network domains, each potentially managed by different service providers, poses a significant challenge due to limited visibility into real-time underlying domain conditions. This paper introduces Risk-Aware Iterated Local Search (RAILS), a novel risk model-driven meta-heuristic framework designed to jointly address SLA decomposition and service provider selection in multi-domain networks. By integrating online neural network (NN)-based risk modeling with iterated local search principles, RAILS effectively navigates the complex optimization landscape, utilizing historical feedback from domain controllers. We formulate the joint problem as a Mixed-Integer Nonlinear Programming (MINLP) problem and prove its NP-hardness. Extensive simulations demonstrate that RAILS achieves near-optimal performance, offering an efficient, real-time solution for adaptive SLA management in modern multi-domain networks.

**Index Terms**—network slicing, service level agreement, risk model, quality of service, deep neural network, optimization

## I. INTRODUCTION

The advent of 5G has transformed mobile networks into multi-service environments tailored to diverse industry needs. A key enabler of this shift is network slicing, which creates multiple End-to-End (E2E) logical networks over shared infrastructure, each customized per Service Level Agreements (SLAs). SLAs define expected Quality of Service (QoS) through Service-Level Objectives (SLOs), covering metrics like throughput, latency, reliability, and security. A single network slice may span across multiple segments of the network, including (radio) access, transport, and core networks, and it may involve collaboration between different operators and infrastructure providers. To ensure that the service meets the agreed-upon SLOs across these domains, it is essential to adjust the service parameters accordingly. As a result, the E2E SLA associated with a network slice must be partitioned into specific SLOs for each domain. This decomposition is crucial for effective resource allocation and remains a core challenge in network slicing. Several studies have discussed this issue. [1] highlights the complexity of mapping E2E requirements

to transport networks. [2] focuses on lifecycle automation, orchestration, and real-time monitoring for SLA compliance. [3] stresses the role of SLA parameters in E2E QoS and the need for appropriate transport resources. Additionally, [4] underscores the importance of SLA decomposition for resource allocation, while [5] explores AI-assisted SLA decomposition in automating 6G business processes.

In typical network slicing management architectures, a two-level hierarchy is employed [6]–[8]. This includes an E2E service orchestrator, responsible for overseeing the lifecycle management of network services, and local domain controllers, which manage the instantiation of network slices within their specific domains. The orchestrator determines how the E2E SLA is partitioned into domain-specific SLOs. However, a common constraint is that the orchestrator usually lacks real-time visibility into the state of each domain’s infrastructure at the moment of decomposition. Instead, it relies on historical data reflecting the outcomes of previous slice requests. Several studies [9]–[11] have introduced prediction-based approaches for SLA management, though they do not explicitly tackle the E2E SLA decomposition problem. In [9], the authors proposed a mapping layer that oversees the network within a service area, managing radio resource allocation to slices to ensure their target service requirements are met. The work in [10] presented an SLA-constrained optimization method leveraging Deep Learning (DL) to estimate resource requirements based on per-slice traffic. Similarly, [11] utilized a context-aware approach, employing graph representations to predict SLA violations in cloud computing environments. Additionally, heuristic-based SLA decomposition methods have been explored in prior research [4]. In [12], the authors introduced an E2E SLA decomposition system that applies supervised machine learning to partition E2E SLAs into access, transport, and core SLOs.

In our previous work [6], [7], we tackled the SLA decomposition problem using neural network (NN)-based risk models in a two-step approach that combined machine learning and optimization. Building on that, [8] introduced an online learning–decomposition framework for dynamic, multi-

domain SLA management. However, these studies assumed a preselected service provider per domain, focusing solely on optimizing E2E acceptance probabilities. In real-world network environments, multiple service providers are often available within each domain, offering varying performance characteristics and capabilities. For example, in a 5G network slice for autonomous vehicles, Ericsson provides high-capacity RAN for low-latency urban coverage, Nokia ensures reliable transport with energy-efficient networking, and AWS offers a scalable cloud-native core. This combination ensures stringent SLA requirements for real-time communication. As a result, the optimization process should consider both the decomposition of SLAs and the selection of providers across domains to ensure more flexible and efficient resource utilization in multi-domain networks. To address these limitations, this paper introduces Risk-Aware Iterated Local Search (RAILS), a novel risk model-driven meta-heuristic framework. RAILS extends the principles of Iterated Local Search (ILS) by integrating dynamic risk modeling and SLA decomposition techniques proposed in [8]. The main contributions of this paper are:

1. We formulate the joint SLA decomposition and service provider selection tasks as a Mixed-Integer Nonlinear Programming (MINLP) problem and demonstrate its NP-hardness.
2. We propose RAILS, a novel risk-aware meta-heuristic framework designed to jointly address the SLA decomposition and service provider selection problem.
3. We empirically show that RAILS achieves near-optimal performance within an analytic model-based simulation environment with low computational overhead.

The paper is organized as follows: Section II defines the system model and formulates the problem. Section III introduces the RAILS framework. Section IV details the simulation setup, while Section V presents and discusses the results. Section VI concludes the paper.

## II. SYSTEM MODEL

### A. Problem Formulation

Let  $N$  denote the number of domains that the service spans. For each domain  $i$  (with  $i = 1, \dots, N$ ), let  $\mathcal{J}_i$  denote the set of available service providers. We define the following decision variables:

- $x_{ij} \in \{0, 1\}$ : a binary variable that is 1 if provider  $j \in \mathcal{J}_i$  is selected for domain  $i$ , and 0 otherwise.
- $d_i \geq 0$ : the portion of the E2E delay budget allocated to domain  $i$ .

The acceptance probability of domain  $i$  using provider  $j$  when allocated a delay of  $d_i$  is given by the function  $p_{ij}(d_i)$ . Assuming that the decisions made in the domains are statistically independent, the E2E acceptance probability  $p_{e2e}$  is modeled as the product of the acceptance probabilities of all domains:

$$p_{e2e} = \prod_{i=1}^N \left( \sum_{j \in \mathcal{J}_i} x_{ij} p_{ij}(d_i) \right). \quad (1)$$

The goal is to **choose a provider in each domain and allocate the delay budgets**  $\{d_i\}_{i=1}^N$  such that the **E2E acceptance probability is maximized**, subject to the constraint that the domain-specific partial delays sum up to the E2E delay budget  $d_{e2e}$ . Formally, the problem is formulated as follows:

$$\begin{aligned} \max_{\{x_{ij}, d_i\}} \quad & \prod_{i=1}^N \left( \sum_{j \in \mathcal{J}_i} x_{ij} p_{ij}(d_i) \right) \\ \text{s.t.} \quad & \sum_{i=1}^N d_i = d_{e2e}, \\ & d_i \geq 0, \quad \forall i = 1, \dots, N, \\ & \sum_{j \in \mathcal{J}_i} x_{ij} = 1, \quad \forall i = 1, \dots, N, \\ & x_{ij} \in \{0, 1\}, \quad \forall i = 1, \dots, N, \forall j \in \mathcal{J}_i. \end{aligned} \quad (2)$$

The presence of both integer and continuous variables, coupled with the nonlinear characteristics of the objective function, designates the problem as a canonical MINLP problem. However, due to the lack of knowledge about  $p_{ij}(d_i)$ , we leverage historical feedback data to construct a NN-based risk model for each domain [7]. These risk models serve as surrogates for  $p_{ij}(d_i)$  in the optimization process, providing an estimated acceptance probability based on past observations.

### B. NP-Hardness Analysis

We demonstrate that the joint optimization problem described in (2) is NP-hard by reducing from the well-known Multiple-Choice Subset Sum Problem (MCSSP). To this end, we consider a simplified version of the problem where only one provider is available per domain. In this case, the decision problem simplifies to the following objective:

$$\begin{aligned} \max_{d_i} \quad & \prod_{i=1}^N p_i(d_i) \\ \text{s.t.} \quad & \sum_{i=1}^N d_i = d_{e2e}, \\ & d_i \geq 0, \quad \forall i = 1, \dots, N. \end{aligned} \quad (3)$$

Now, we define an instance of MCSSP: Given  $N$  disjoint sets  $\{S_i\}_{i=1}^N$  of nonnegative integers and a target sum  $K$ , the goal is to determine whether it is possible to choose exactly one element from each set such that the sum of the selected elements equals  $K$ . We then construct a corresponding instance of the simplified problem described in (3) as follows:

- The number of domains is set to  $N$ .
- The total E2E delay budget is  $d_{e2e} = K$ .
- For the provider in each domain  $i$ , we define the acceptance probability function  $p_i(d_i)$  as:

$$p_i(d_i) = \begin{cases} 1, & \text{if } d_i \in S_i, \\ 0, & \text{otherwise.} \end{cases} \quad (4)$$

- The decision version of the simplified problem asks *whether there exists a feasible delay allocation*  $\{d_i\}_{i=1}^N$  such that  $p_{e2e} = 1$ , subject to the constraints in (3).

This construction ensures that  $p_i(d_i) = 1$  if and only if the allocated delay  $d_i$  is equal to at least one element in  $S_i$ . As a result,  $p_{e2e}$  is equal to 1 if and only if  $d_i \in S_i$  for all  $i$ . Moreover, since  $\{d_i\}_{i=1}^N$  sum up to  $K$ , any feasible solution that achieves  $p_{e2e} = 1$  corresponds to a valid solution to the MCSSP instance. Therefore, solving the decision version of the simplified problem is equivalent to solving MCSSP, which is known to be NP-complete [13]. Hence, the optimization version of simplified problem in (3) is NP-hard. Since the original problem in (2) generalizes this setting by incorporating continuous, nonlinear acceptance probability functions and allowing multiple providers per domain, it is also NP-hard.

Given the NP-hard nature of the problem, finding an exact solution is computationally intractable for large-scale systems. As a result, we resort to meta-heuristic approaches, leveraging the domain-specific risk models built from historical data to guide the search for near-optimal solutions.

### III. METHODOLOGY

#### A. Background

**ILS.** Iterated Local Search (ILS) [14] is a meta-heuristic approach that enhances local search algorithms by escaping local optima through iterative perturbations and refinements. ILS operates by first generating an initial solution, either randomly or via a heuristic method, and then refining it through a local search procedure to find a local optimum. Once a local optimum is identified, the algorithm introduces controlled randomness to perturb the solution and push it away from the identified optimum. This cycle of local search and perturbation continues until a stopping condition, such as reaching a maximum number of iterations or meeting a convergence criterion, is satisfied. ILS is used in networked cloud resource mapping to address the challenge of optimally partitioning and embedding virtual resources across multiple cloud providers [15]. ILS-based request partitioning has been shown to effectively balance cost and performance, leading to improved virtual network embedding outcomes.

**RADE.** Real-time Adaptive DEcomposition (RADE) [8] is an advanced SLA decomposition framework that dynamically adjusts decomposition strategies based on real-time feedback of network conditions. Unlike static decomposition approaches, RADE employs online learning to enhance adaptability and accuracy. It utilizes a two-step decomposition approach [6], [7]. First, the orchestrator maintains domain-specific NN-based risk models trained on historical SLA acceptance and rejection feedback. Next, the E2E SLA is decomposed into domain-specific SLAs to maximize the overall acceptance probability, using a grid search followed by Sequential Least Squares Programming (SLSQP) algorithm. To adapt to evolving network conditions, these risk models are updated timely via Online Gradient Descent (OGD). A First In First Out (FIFO) memory buffer preserves recent observations, ensuring stable learning while mitigating overfitting caused by transient anomalies. RADE addresses key limitations of static decomposition methods by incorporating real-time adaptation. It also offers resilience against data corruption through its FIFO

memory buffer. Experimental results show that RADE consistently outperforms traditional methods in dynamic multi-domain environments, making it a promising solution for adaptive SLA management in modern network architectures.

#### B. Risk-Aware Iterated Local Search

In this work, we propose RAILS, a risk model-driven meta-heuristic method to solve the joint provider selection and SLA decomposition problem. The optimization problem involves two interconnected sets of decision variables (see Section II-A). On one hand, we have discrete variables that determine which provider is selected in each domain. On the other hand, we have continuous variables that specify how the E2E delay budget is decomposed among the domains to maximize the E2E acceptance probability. These two aspects of the problem are inherently intertwined because the acceptance probability in each domain is computed using risk models that depend on both the selected provider  $j$  and the assigned delay requests. Specifically, once provider  $j$  is chosen for domain  $i$ , the corresponding risk model serves as its surrogate, predicting  $p_{ij}(d_i)$  for a given delay budget  $d_i$ . In RAILS, the framework efficiently explores the complex search space by iteratively refining provider selections with risk models.

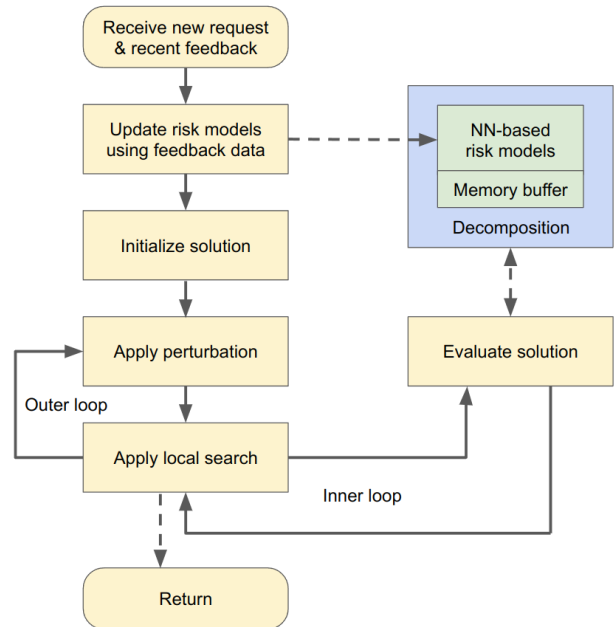


Fig. 1: A single-iteration workflow of RAILS.

Fig. 1 illustrates the workflow of the RAILS algorithm. The process begins with updating the risk models and memory buffer using the latest feedback data. Next, an initial solution for provider selection is generated. The algorithm then enters an iterative refinement phase with perturbation and local search steps. The perturbation step randomly alters the provider selection for each domain with a probability  $p_\mu$  to possibly escape local optima, while the local search step refines the solution by randomly selecting a domain and exhaustively checking all provider options within the domain to identify

the best one based on the risk models. This iterative process continues within an inner and outer loop structure until a predefined stopping condition is met, at which point the best solution identified is returned. RAILS operates within an ILS framework, where the core evaluation mechanism is powered by RADE. Specifically, ILS performs the repeated perturbation and local search steps, while RADE handles dynamic risk modeling and real-time decomposition. This synergy enables effective handling of the coupled discrete–continuous nature of provider selection and SLA decomposition.

#### IV. PERFORMANCE EVALUATION

##### A. Simulation Environment

In our simulation environment, we model the dynamic behavior of each provider’s system load and the resulting performance characteristics that affect SLA acceptance. In particular, we capture the temporal variations in load and their impact on the minimum delay that a provider can support, which in turn governs the acceptance probability of an SLA request. Because this subsection focuses on a single-domain provider, we omit the  $i$  and  $j$  subscripts for clarity.

**System Load Modeling.** For each provider, the system load is assumed to evolve periodically over time. Let  $t$  denote the current time. The system load  $\ell(t)$  is modeled using a sinusoidal function [16] as follows:

$$\ell(t) = \ell_{\text{base}} \cdot k + \ell_{\text{base}} \cdot (1 - k) \cdot \frac{1 + \sin\left(\frac{2\pi t}{T} + \phi\right)}{2}, \quad (5)$$

where  $\ell_{\text{base}}$  is a constant representing the baseline load of the provider,  $k \in [0, 1]$  is a parameter that determines the fraction of the load that is static,  $T$  is the period of the sinusoidal fluctuation, and  $\phi$  is the phase shift. This formulation ensures that the system load varies between the minimum load  $\ell_{\text{base}} \cdot k$  and the maximum load  $\ell_{\text{base}}$ . The parameter  $k$  allows for a mixture of a constant baseline load and a dynamic component.

**Minimum Supportable Delay.** Given the dynamic system load defined in (5), the minimum delay that a provider can support for an incoming request is assumed to depend on both a fixed latency component and an exponential function of the system load. Specifically, the minimum supportable delay  $d^{\min}(t)$  is defined as:

$$d^{\min}(t) = \alpha + \exp(\beta \cdot \ell(t)), \quad (6)$$

where  $\alpha$  represents the inherent latency of the system when the load is minimal (i.e., the baseline latency), and  $\beta$  is a parameter that characterizes how sensitive the delay is to changes against system load. This expression reflects that as the system load increases, the provider’s capability to handle requests with low delay diminishes, leading to a higher minimum supportable delay. Fig. 2 demonstrates this effect for different parameter values. The exponential relationship captures the non-linear increase in minimum supportable delay as the load intensifies, while  $\alpha$  sets the lower bound of latency.

**Acceptance Probability.** Based on the minimum supportable delay defined in (6), the acceptance probability of a service level request is defined as a function of the requested delay  $d$ .

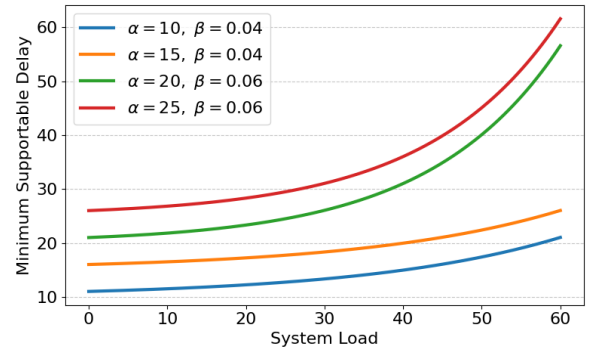


Fig. 2: Minimum supportable delay versus system load.

Specifically, if the requested delay is less than the minimum supportable delay  $d^{\min}(t)$ , the request is rejected. Otherwise, the acceptance probability follows an S-curve, increasing with the excess delay and saturating as the requested delay becomes much larger than  $d^{\min}(t)$  [6]. Formally, the acceptance probability  $p(d; t)$  is defined as:

$$p(d; t) = \begin{cases} 0, & \text{if } d < d^{\min}(t), \\ 1 - \exp(-\lambda(d - d^{\min}(t))), & \text{if } d \geq d^{\min}(t), \end{cases} \quad (7)$$

where  $\lambda > 0$  is a parameter that controls the rate at which the acceptance probability increases as the delay requirement exceeds the minimum supportable delay. The response curves with different parameter settings are illustrated in Fig. 3.

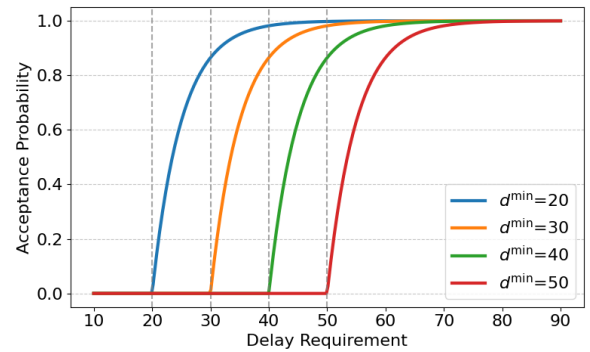


Fig. 3: Piecewise acceptance response curves.

The piecewise function ensures that any request with a delay requirement less than the system’s minimum supportable delay is rejected. For requests above this threshold, the probability of acceptance increases rapidly at first and gradually levels off, reflecting the typical behavior of admission control in real-world service provider systems.

##### B. Evaluation Scenarios

To assess the long-term performance of the proposed RAILS framework, we conduct simulations over 100 discrete time steps within the designed simulation environment. At each time step, a new service request arrives with an E2E delay budget of 100 ms. The RAILS is then applied to select a service

provider for each domain and to determine the corresponding delay decomposition that maximizes the E2E acceptance probability according to the risk models built from historical data. Specifically, once the RAILS provides a provider selection and delay decomposition, this assignment is evaluated using the ground-truth acceptance probability models (described in (7)) to compute the actual E2E acceptance probability. This process is repeated at every time step, and the resulting E2E acceptance probabilities are collected to compute an overall average performance metric, namely:

$$\bar{P}_{e2e} = \frac{1}{T} \sum_{t=1}^T P_{e2e}^{(t)} = \frac{1}{T} \sum_{t=1}^T \prod_{i=1}^N \left( \sum_{j \in \mathcal{J}_i} x_{ij}^{(t)} p_{ij}^{(t)}(d_i^{(t)}) \right), \quad (8)$$

where  $T$  denote the total number of simulation time steps, and  $N$  represents the total number of involved domains.

At each time step, we assume that a set of historical requests and their associated feedback are available from each provider. A historical sample is represented by a pair  $(d^{(t)}, a^{(t)})$ , where  $d^{(t)}$  is the delay request and  $a^{(t)} \in \{0, 1\}$  is the binary decision outcome from the admission control. For generating the feedback data, we simulate requests by sampling delay requirements uniformly from the interval [10 ms, 100 ms]. Each request is then processed through the corresponding ground-truth acceptance probability model to obtain its actual acceptance probability, and a coin-flipping process is used to determine whether the request is accepted or not by the admission control. For performance comparison, we consider two baselines:

**Non-Risk-Aware (NRA).** In this baseline, a provider is selected at random for each domain, and the E2E delay budget is evenly decomposed among the domains as a heuristic guess, without leveraging risk models.

**Optimal (OPT).** This benchmark uses an exhaustive search with full access to ground-truth acceptance probability models, representing the theoretical upper bound on performance.

### C. Experimental Setup

In our experiments, we consider a network slicing scenario involving 3 domains, each comprising 10 service providers. The ground-truth models for each provider are generated using a set of randomly selected parameters to emulate realistic and heterogeneous operational conditions. For each provider, the baseline latency  $\alpha$  is drawn uniformly from [0, 2]; The parameter  $\lambda$  in (7) is set to 0.2, and  $k$  is set to 0.5 in (5); A domain-wise additional latency, randomly chosen from the set  $\{0, 10, 20\}$ , is added to the baseline latency  $\alpha$  to reflect inter-domain behavioral shifts; the load-sensitivity parameter  $\beta$  is sampled uniformly from [0.04, 0.06]; the baseline system load  $\ell_{\text{base}}$  is drawn uniformly from [30, 50]; the period of the sinusoidal load fluctuation is selected as an integer uniformly from the range [30, 60]; and the phase shift in the load function is chosen uniformly from the interval  $[0, \pi]$ .

At each time step, we assume that the number of recent feedback samples available from each provider is proportional to its current system load defined in (5). We use the integer

part of the load as the number of samples. For the RAILS framework, the number of iterations is set empirically to the total number of providers across all domains (i.e.,  $3 \times 10 = 30$ ), and the perturbation probability  $p_\mu$  is 0.8. Results are averaged over 10 independent runs to ensure statistical significance. Each risk model is implemented as a 2-layer monotonic Multi-Layer Perceptron (MLP) with a hidden dimension of 16, similar to that described in [8]. AdamW optimizer is used with a learning rate of 0.01. A memory buffer of size 300 is maintained, and each risk model update involves 10 iterations.

## V. RESULTS AND DISCUSSION

Fig. 4a presents the average E2E acceptance probability for the considered approaches. Each bar reflects the average performance metric defined in (8). The NRA approach achieves an average acceptance probability of approximately 0.71. The large error bars reveal significant performance fluctuations, suggesting that the absence of strategic decision-making leads to suboptimal performance. In contrast, the proposed RAILS method demonstrates a substantial improvement, achieving an average acceptance probability of around 0.89. The error bars are smaller compared to the NRA approach, indicating more stable outcomes. Moreover, we include results for a RAILS variant, denoted RAES, where the ILS search component is replaced with an exhaustive search. RAES's performance reflects RAILS running with an effectively infinite number of iterations. Both RAILS and RAES achieve comparable performance, while RAILS requires only 31.6% of RAES's run time, demonstrating its computational efficiency, as shown in Fig. 4b. The OPT approach, which represents the theoretical performance upper bound, achieves an average acceptance probability of about 0.95. While RAILS does not fully reach the OPT benchmark, it comes remarkably close, balancing computational efficiency with SLA acceptance rates.

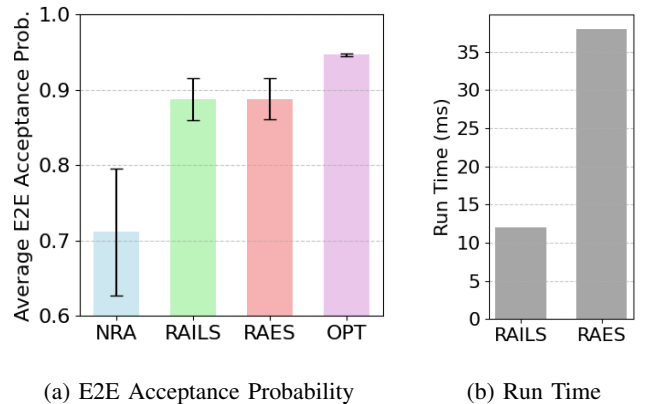


Fig. 4: Performance comparison: (a) E2E SLA acceptance probability over time; (b) computational run time.

Table I provides a partial trace of a simulation run, showcasing the E2E acceptance probabilities, delay decompositions, and the indices of the selected providers across different time steps. The OPT method shows that the optimal provider

TABLE I: Partial trace of a simulation run.

Step	Method	$p_{e2e}$	$d_1, d_2, d_3$	Providers
0	NRA	0.78	33.33, 33.33, 33.33	7, 9, 6
	RAILS	0.92	22.00, 33.00, 45.00	8, 6, 0
	OPT	0.94	24.23, 32.50, 43.27	3, 1, 0
5	NRA	0.81	33.33, 33.33, 33.33	3, 4, 9
	RAILS	0.88	22.63, 32.08, 45.30	4, 6, 1
	OPT	0.94	24.31, 32.49, 43.20	3, 1, 6
10	NRA	0.75	33.33, 33.33, 33.33	5, 4, 1
	RAILS	0.93	26.22, 30.30, 43.48	8, 6, 6
	OPT	0.94	24.17, 32.98, 42.85	5, 1, 3
15	NRA	0.73	33.33, 33.33, 33.33	7, 3, 9
	RAILS	0.90	20.00, 30.00, 50.00	1, 6, 7
	OPT	0.94	24.11, 33.00, 42.90	2, 4, 3

selections and delay decompositions change frequently over time, reflecting dynamic network conditions. For instance, optimal providers shift from (3, 1, 0) at step 0 to (2, 4, 3) at step 15, highlighting the need for adaptive risk models to maintain high SLA acceptance rates. Furthermore, RAILS achieves near-optimal performance across time steps, even without always selecting the optimal providers.

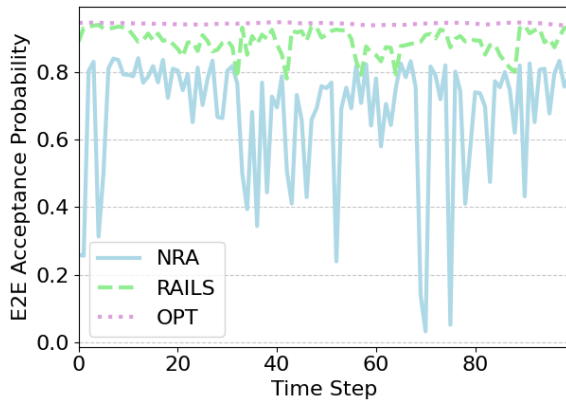


Fig. 5: E2E acceptance probability over time for a single run.

Fig. 5 illustrates the temporal dynamics of E2E acceptance probabilities for a single representative run. The OPT method consistently achieves near-perfect acceptance. RAILS closely tracks the optimal performance, with slight fluctuations, indicating its effectiveness in adapting to dynamic network conditions. In contrast, the NRA method experiences significant variability and frequent drops in acceptance probability, revealing its limitations in predicting and responding to environmental changes.

## VI. CONCLUSION

This paper introduced Risk-Aware Iterated Local Search (RAILS), a meta-heuristic framework driven by NN-based risk models, for SLA decomposition and service provider selection in multi-domain networks. We formulated the problem as a Mixed-Integer Nonlinear Programming (MINLP) problem and demonstrated its NP-hardness. RAILS integrates dynamic risk modeling with iterated local search, effectively handling the complex optimization landscape of interdependent decisions.

Simulation results showed that RAILS achieves near-optimal performance against the theoretical optima while maintaining low computational overhead. Overall, RAILS offers a robust and efficient solution for adaptive network slicing management in modern network systems. Future work will explore the long-term impact of each decision on subsequent ones by formulating the problem as a Markov Decision Process (MDP) and applying Deep Reinforcement Learning (DRL) techniques.

## ACKNOWLEDGMENT

This research was partially funded by the HORIZON SNS JU DESIRE6G project (grant no. 101096466) and the Dutch 6G flagship project “Future Network Services”.

## REFERENCES

- [1] X. Geng, L. M. Contreras, R. Rokui, J. Dong, and I. Bykov, “IETF Network Slice Application in 3GPP 5G End-to-End Network Slice,” Internet Engineering Task Force, Internet-Draft draft-ietf-teas-5g-network-slice-application-03, Jun. 2024.
- [2] R. Swamy and S. K. M., “5G network slicing,” HCL Technologies, Tech. Rep., 2023.
- [3] P. Iovanna, M. Svensson, A. Shapin, G. Bottari, F. Ubaldi, F. Ponzini, and M. Puleri, “End-to-end network slicing orchestration,” *Ericsson Technology Review*, vol. 2, 2022.
- [4] R. Su, D. Zhang, R. Venkatesan, Z. Gong, C. Li, F. Ding, F. Jiang, and Z. Zhu, “Resource allocation for network slicing in 5G telecommunication networks: A survey of principles and models,” *IEEE Network*, vol. 33, no. 6, pp. 172–179, 2019.
- [5] J. Wang, J. Liu, J. Li, and N. Kato, “Artificial intelligence-assisted network slicing: Network assurance and service provisioning in 6G,” *IEEE Vehicular Technology Magazine*, vol. 18, no. 1, pp. 49–58, 2023.
- [6] D. De Vleeschauwer, C. Papagianni, and A. Walid, “Decomposing SLAs for network slicing,” *IEEE Communications Letters*, vol. 25, no. 3, pp. 950–954, March 2021.
- [7] C. S.-H. Hsu, D. De Vleeschauwer, and C. Papagianni, “SLA decomposition for network slicing: A deep neural network approach,” *IEEE Networking Letters*, pp. 1–1, 2023.
- [8] C. S.-H. Hsu, D. De Vleeschauwer, C. Papagianni, and P. Grosso, “Online SLA decomposition: Enabling real-time adaptation to evolving network systems,” in *2025 Joint European Conference on Networks and Communications & 6G Summit (EuCNC/6G Summit)*, 2025, to be published.
- [9] B. Khodapanah, A. Awada, I. Viering, D. Oehmann, M. Simsek, and G. P. Fettweis, “Fulfillment of service level agreements via slice-aware radio resource management in 5G networks,” in *2018 IEEE 87th Vehicular Technology Conference (VTC Spring)*, 2018, pp. 1–6.
- [10] H. Chergui and C. Verikoukis, “Offline SLA-constrained deep learning for 5G networks reliable and dynamic end-to-end slicing,” *IEEE Journal on Selected Areas in Communications*, vol. 38, no. 2, pp. 350–360, 2020.
- [11] A.-C. Maroudis, T. Theodoropoulos, J. Violos, A. Leivadreas, and K. Tserpes, “Leveraging graph neural networks for SLA violation prediction in cloud computing,” *IEEE Transactions on Network and Service Management*, vol. 21, no. 1, pp. 605–620, 2024.
- [12] M. Iannelli, M. R. Rahman, N. Choi, and L. Wang, “Applying machine learning to end-to-end slice SLA decomposition,” in *2020 6th IEEE Conference on Network Softwarization (NetSoft)*, 2020, pp. 92–99.
- [13] H. Kellerer, U. Pferschy, and D. Pisinger, *The Multiple-Choice Knapsack Problem*. Berlin, Heidelberg: Springer Berlin Heidelberg, 2004, pp. 317–347.
- [14] H. R. Lourenço, O. C. Martin, and T. Stützle, *Iterated Local Search*. Boston, MA: Springer US, 2003, pp. 320–353.
- [15] A. Leivadreas, C. Papagianni, and S. Papavassiliou, “Efficient resource mapping framework over networked clouds via iterated local search-based request partitioning,” *IEEE Transactions on Parallel and Distributed Systems*, vol. 24, no. 6, pp. 1077–1086, 2013.
- [16] S. Islam, K. Lee, A. Fekete, and A. Liu, “How a consumer can measure elasticity for cloud platforms,” in *Proceedings of the 3rd ACM/SPEC International Conference on Performance Engineering*, ser. ICPE ’12. New York, NY, USA: Association for Computing Machinery, 2012, p. 85–96.