

MicroOpt: Model-driven Slice Resource Optimization in 5G and Beyond Networks

Muhammad Sulaiman, Mahdiah Ahmadi, Bo Sun, Mohammad A. Salahuddin, Raouf Boutaba, Aladdin Saleh



ROGERS

DAVID R. CHERITON SCHOOL OF COMPUTER SCIENCE

RESEARCH PROBLEM

BACKGROUND

- **5G Network:** Comprises Virtual Network Functions (VNFs) within the Radio Access Network (RAN), Transport, and Core
- **Network Slicing:** Involves creating isolated, virtual networks tailored to specific use-cases (e.g., eMBB slices tailored of VR/XR, URLLC slice for remote driving, telesurgery and mMTC for IOT applications)
- **Slice QoS & SLA:** Service Level Agreements (SLAs) define the minimum Quality of Service (QoS) 5G slices must receive
- **Resource Allocation vs. QoS:** The amount of resources allocated to VNFs directly impacts the slice's QoS

MOTIVATION

- **Dynamic Traffic Patterns:** Slices experience time-varying traffic
- **Resource Efficiency:** Resource allocation for peak traffic leads to resource over-provisioning

PROBLEM STATEMENT & CHALLENGES

- **Resource Optimization:** Minimize resource allocation while meeting QoS requirements
 - **Network Complexity:** Accurately modeling complex, multi-VNF slices
 - **Real-Time Adaptation:** Fast, efficient resource adjustment to varying slice traffic

SOLUTION

- **ML-Based VNF & Slice Modeling:** Leveraging machine learning for modeling VNFs that can be composed to create end-to-end 5G slice models
- **Resource Allocation Algorithm:** Using Lagrangian primal-dual algorithm coupled with gradient descent for fast, near-optimal resource scaling under QoS constraints

ML-BASED VNF & SLICE MODELING

VNF MODELING

- **Input/Output:** Predicts the egress traffic feature vector given input traffic trace (i.e., pcap) and VNF configuration & resource allocation
- **Dataset:** Uses data from in-lab 5G testbed and partner network operator's real network for training VNF models
- **Differentiability:** Leverages reparameterization trick to maintain VNF model differentiability
- **Composability:** Allows stacking VNF models to form end-to-end slice model

SLICE MODELING

- **Construction:** Composed using individual VNF models

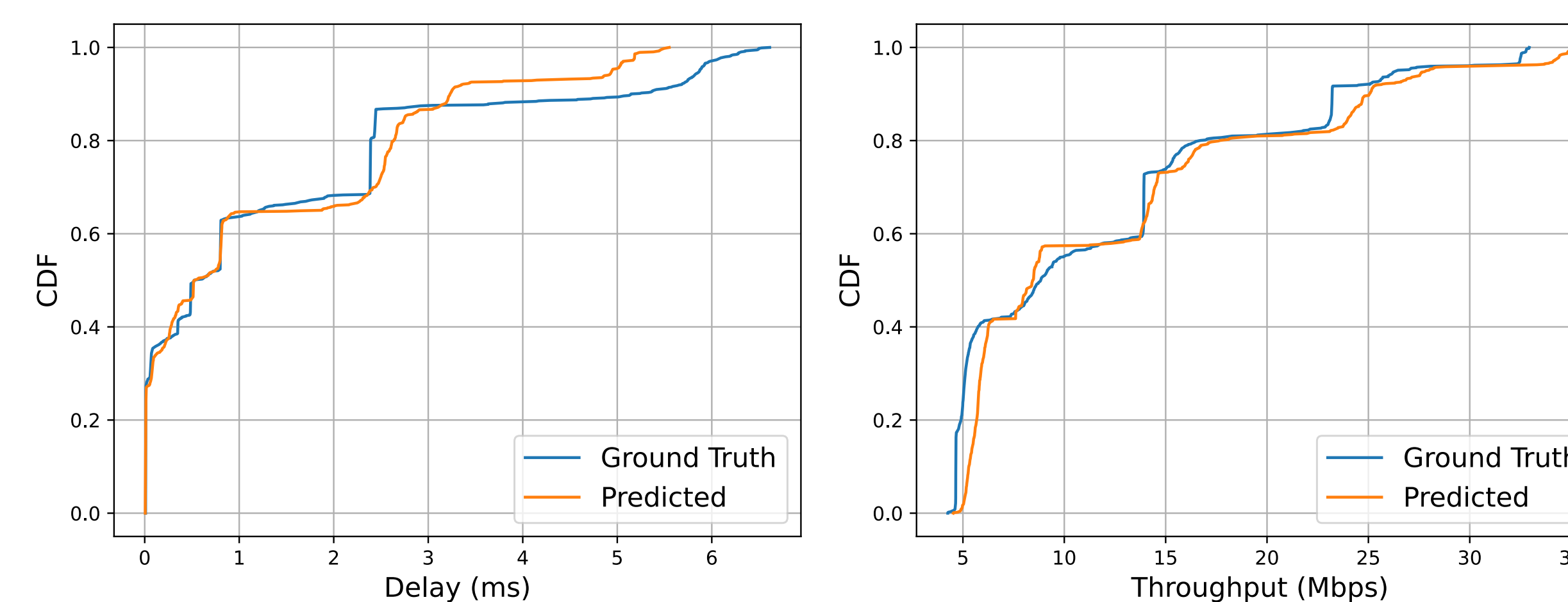
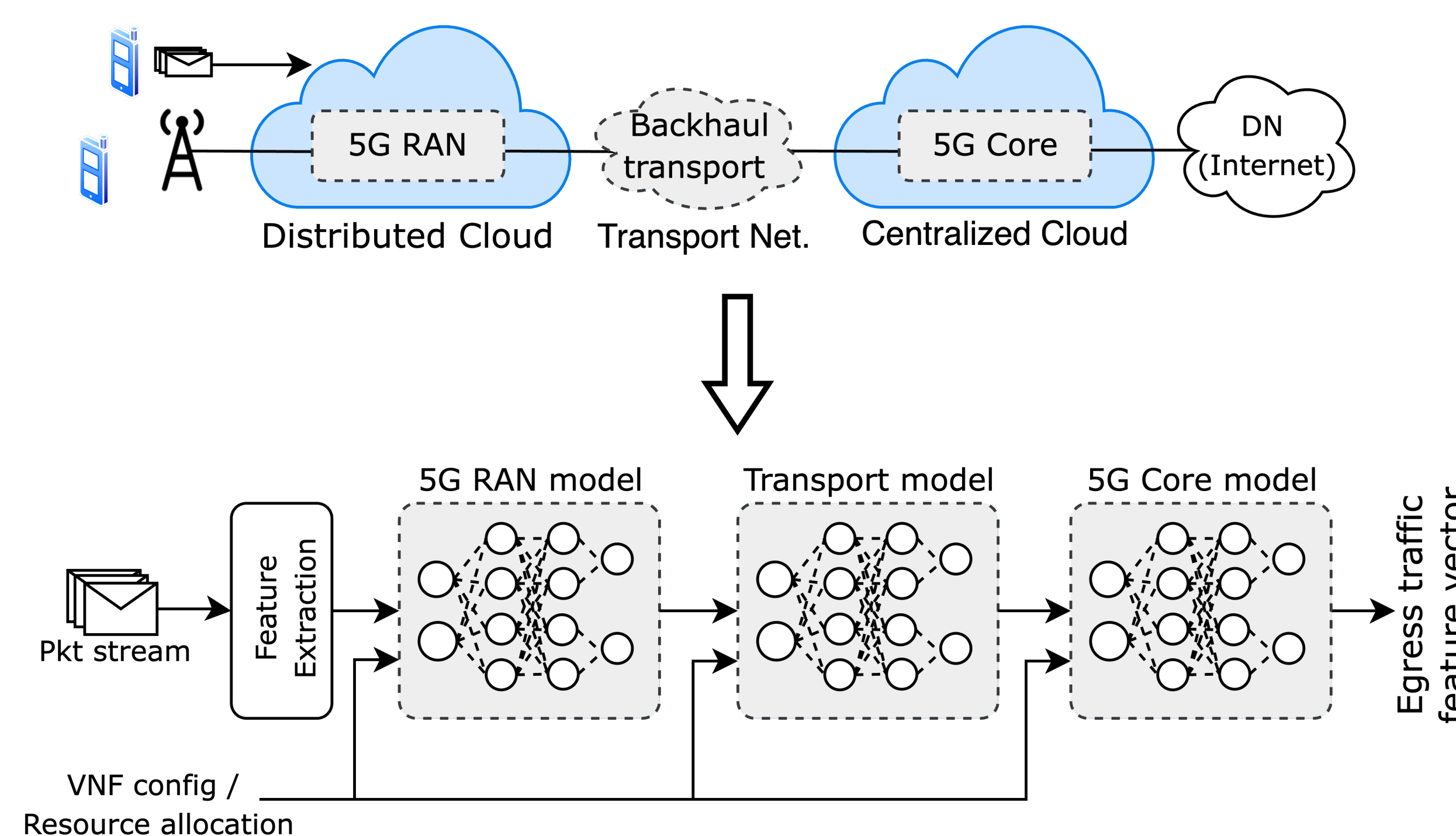


Fig: Ground truth vs. predicted throughput and delay for end-to-end slice model

RESOURCE ALLOCATION ALGORITHM

ALGORITHM

- **Lagrangian primal dual algorithm coupled with Gradient descent**
 - **Relaxed Lagrangian formulation:**

$$\mathcal{L} = \mathbf{r} + \lambda \cdot (QoS_{required} - QoS_{predicted})$$
 - **Outer loop:** Updates Lagrangian variables (λ) based on QoS constraint violation (i.e., $QoS_{required} - QoS_{predicted}$)
 - **Inner loop:** Utilizes gradient descent to minimize the relaxed Lagrangian, with gradients from the differentiable slice model

Algorithm 1 MicroOpt Algorithm

Input: Traffic \mathbf{x}_i^s , Slice Model $f_{QoS}^s(\mathbf{x}_i^s, \mathbf{r}_i^s)$, QoS threshold q_{thresh}^s , QoS degradation threshold β_{thresh}^s , $\tau_{1,max}, \tau_{2,max}, \alpha_1, \alpha_2, \alpha_3, \epsilon_1, \epsilon_2$

Output: Optimal resource allocation vector \mathbf{r}_i^s

- 1: Initialize $\lambda, \mu, LB = 0, UB = \infty, \tau_1 = 0, \tau_2 = 0$
- 2: **while** $\frac{UB-LB}{UB} > \epsilon_1$ **or** $\tau_1 < \tau_{1,max}$ **do**
- 3: $\mathbf{r} \leftarrow \text{Initialization}(\mathbf{x}_i^s, f_{QoS}^s(\mathbf{x}_i^s, \mathbf{r}))$
- 4: **while** $|\nabla_r \mathcal{L}| > \epsilon_2$ **or** $\tau_2 < \tau_{2,max}$ **do**
- 5: $\mathbf{r} \leftarrow [\mathbf{r} - \alpha_1 \nabla_r \mathcal{L}]^+$
- 6: $\tau_2 \leftarrow \tau_2 + 1$
- 7: **end while**
- 8: $\lambda_s \leftarrow [\lambda_s + \alpha_2 (\beta_i^s - \beta_{thresh}^s)]^+, \forall s$
- 9: $\mu_k \leftarrow [\mu_k + \alpha_3 (\sum_{s \in S} r^{s,k} - R^k)]^+, \forall k$
- 10: $LB = \max(LB, \mathcal{L}(\mathbf{r}, \mu, \lambda))$
- 11: $UB = \min(UB, \sum_{s \in S} \eta^T \mathbf{r}^s)$
- 12: $\tau_1 \leftarrow \tau_1 + 1$
- 13: **end while**
- 14: **return** \mathbf{r}

RESULTS

- **Resource saving:** 14.60% and 20.74% improvement over previous SOTA and peak resource allocation
- **Significantly faster resource scaling:** 2-3 orders of magnitude faster resource scaling compared to previous SOTA