# Workshop 1: 5G Core Slicing
## Slice Modeling and Dynamic Resource Scaling

**Raouf Boutaba**
David R. Cheriton School of Computer Science
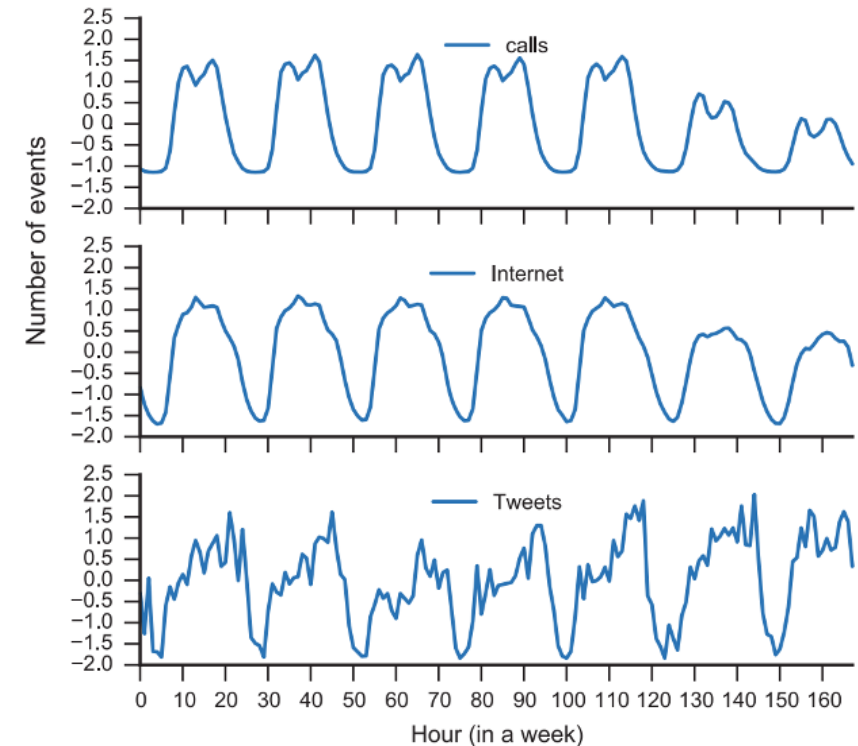University of Waterloo

UNIVERSITY OF
**WATERLOO**

# Table of Contents

- Introduction

- Challenges

- Solutions

- Session#2 structure

UNIVERSITY OF
WATERLOO

# INTRODUCTION

# Introduction: Dynamic Resource Scaling

- **Slice Traffic**: Number of active slice users

    - Varies throughout the week

- **Resource allocation and QoS**

    - Peak allocation vs. average allocation

- **Dynamic resource scaling**:

    - Dynamically Scaling resources allocated to slices based on current or predicted traffic

    - Objective:

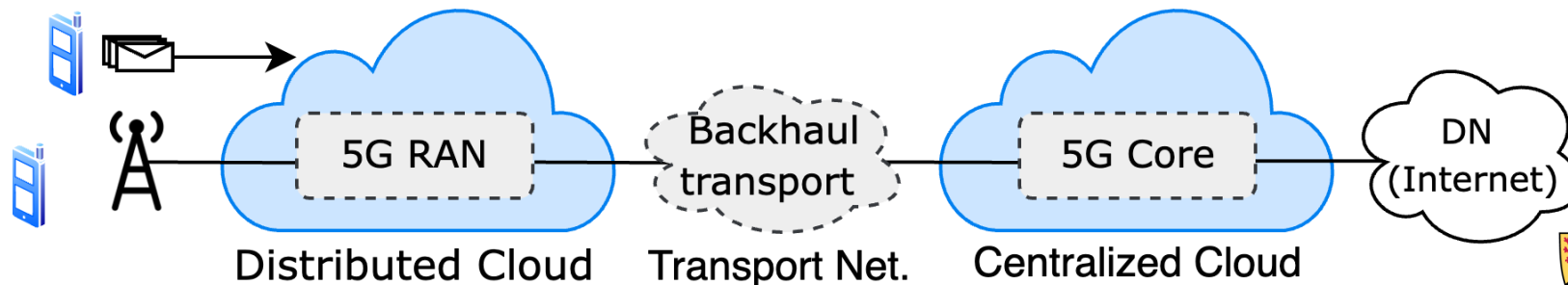        - Minimize resource allocation

        - Satisfy QoS requirements



Scaled weekly behavior of Calls, Internet, Tweets in MILAN [1]

[1] G. Barlacchi et al., "A multi-source dataset of urban life in the city of Milan and the province of Trentino," Scientific Data, 2015.

# CHALLENGES

# Dynamic Resource Scaling: Challenges (1/2)

## Slice Modeling:

- Modeling relationship between resource allocation and QoS

- Slices span multiple network segments (RAN, transport, core) networks

- Heterogeneous QoS and resource requirements

  - Latency, packet loss, reliability, jitter

  - CPU, PRBs, bandwidth, memory

- Traditional modeling approaches are slow and lack real-time application feasibility

# Dynamic Resource Scaling: Challenges (2/2)
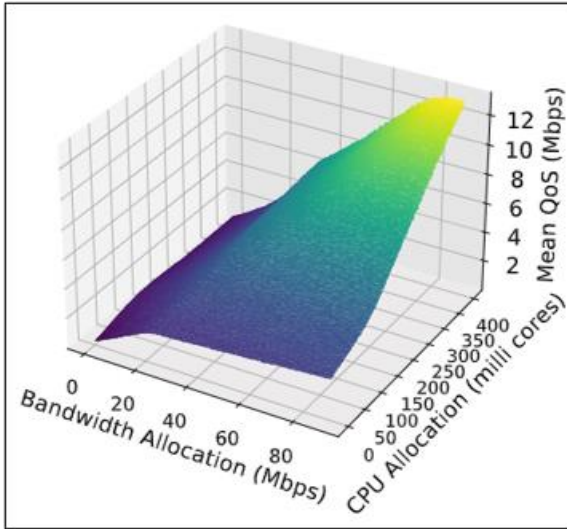
## Constrained Optimization:

- Finding minimum resource allocation that satisfies QoS requirements

- Needs to integrate neural network-based slice model

- Must be fast and efficient

$$\min_{\mathbf{r}} \quad \frac{1}{|T|} \sum_{t \in T} \sum_{s \in \mathbf{S}} \eta^\top \mathbf{r}_t^s \quad \rightarrow \quad \text{Minimize resource allocation to slices, subject to}$$

$$\text{s.t.} \quad \mathbb{E}\left(\beta_{max(T)}^s\right) \leq \beta_{s,thresh}, \quad \forall s \in S \quad \rightarrow \quad \text{QoS degradation threshold constraint}$$

$$\sum_{s \in S} \mathbf{r}_t^s \leq \mathbf{R}, \quad \forall t \in T, \quad \rightarrow \quad \text{Resource capacity constraint}$$

UNIVERSITY OF
WATERLOO

# SOLUTIONS

# Solution Overview

### Network Modelling (**vNetRunner**)



### Slice Traffic



### Resource Scaling Algorithm (**MicroOpt**)



**Algorithm 1** Resource Allocation Algorithm

**Input:** Traffic $x_{\tau_i}^s$, Network Model $f_{QoS}^s(x_{\tau_i}^s, r_{\tau_i}^s)$, QoS threshold $q_{thresh}^s$, QoS degradation threshold $\beta_{thresh}^s$, $\tau_{1,max}, \tau_{2,max}, \alpha_1, \alpha_2, \alpha_3, \epsilon_1, \epsilon_2$

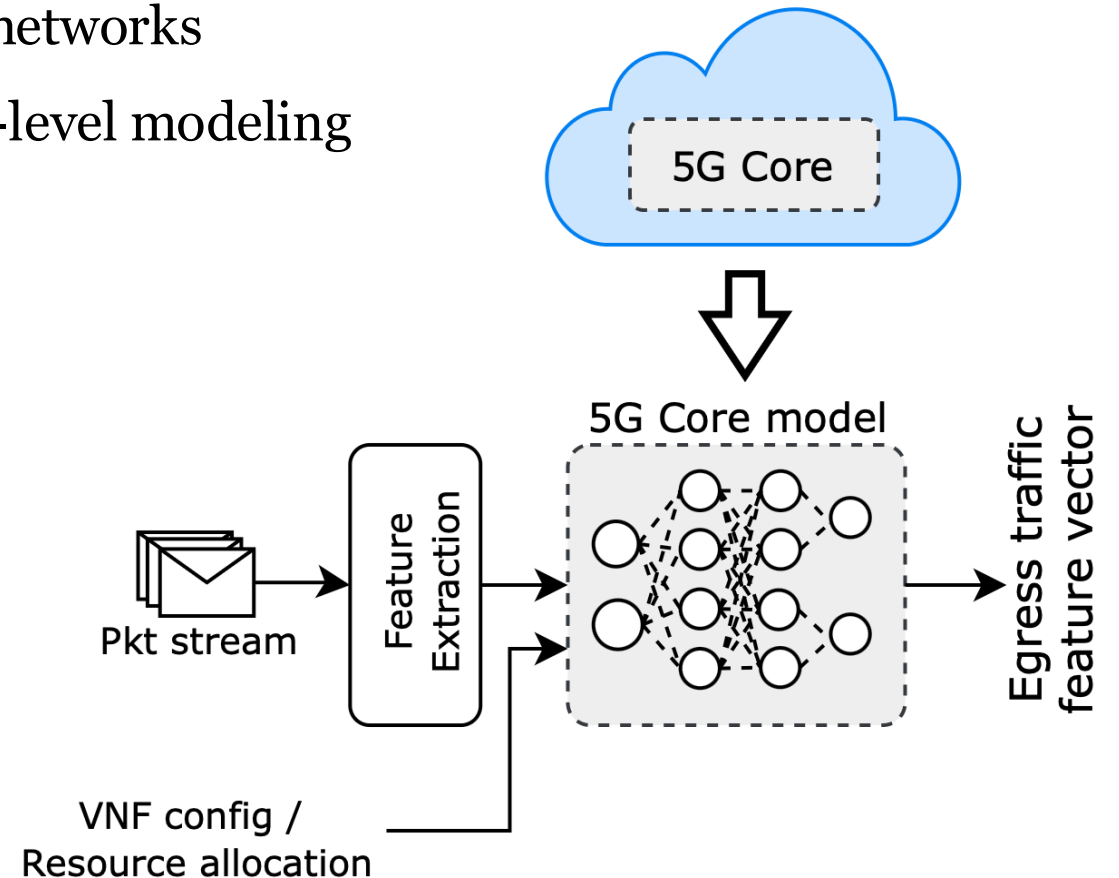**Output:** Optimal resource allocation vector $r_{\tau_i}^s$

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:    $r \leftarrow$ Gridsearch$(x_{\tau_i}^s, f_{QoS}(x_{\tau_i}^s, r))$
4:    **while** $|\nabla_r \hat{\mathcal{L}}| > \epsilon_2$ **or** $\tau_2 < \tau_{2,max}$ **do**
5:       $r \leftarrow [r - \alpha_1 \nabla_r \hat{\mathcal{L}}]^+$
6:       $\tau_2 \leftarrow \tau_2 + 1$
7:    **end while**
8:    $\lambda_s \leftarrow [\lambda_s + \alpha_2(\beta^s - \beta_{thresh}^s)]^+, \forall s$
9:    $\mu_k \leftarrow [\mu_k + \alpha_3(\sum_{s \in S} r^{s,k} - R^k)]^+, \forall k$
10:   LB = max(LB, $\mathcal{L}(r, \mu, \lambda)$)
11:   UB = min(UB, $\sum_{s \in S} \eta^\top r^s$)
12:   $\tau_1 \leftarrow \tau_1 + 1$
13: **end while**
14: **return** r

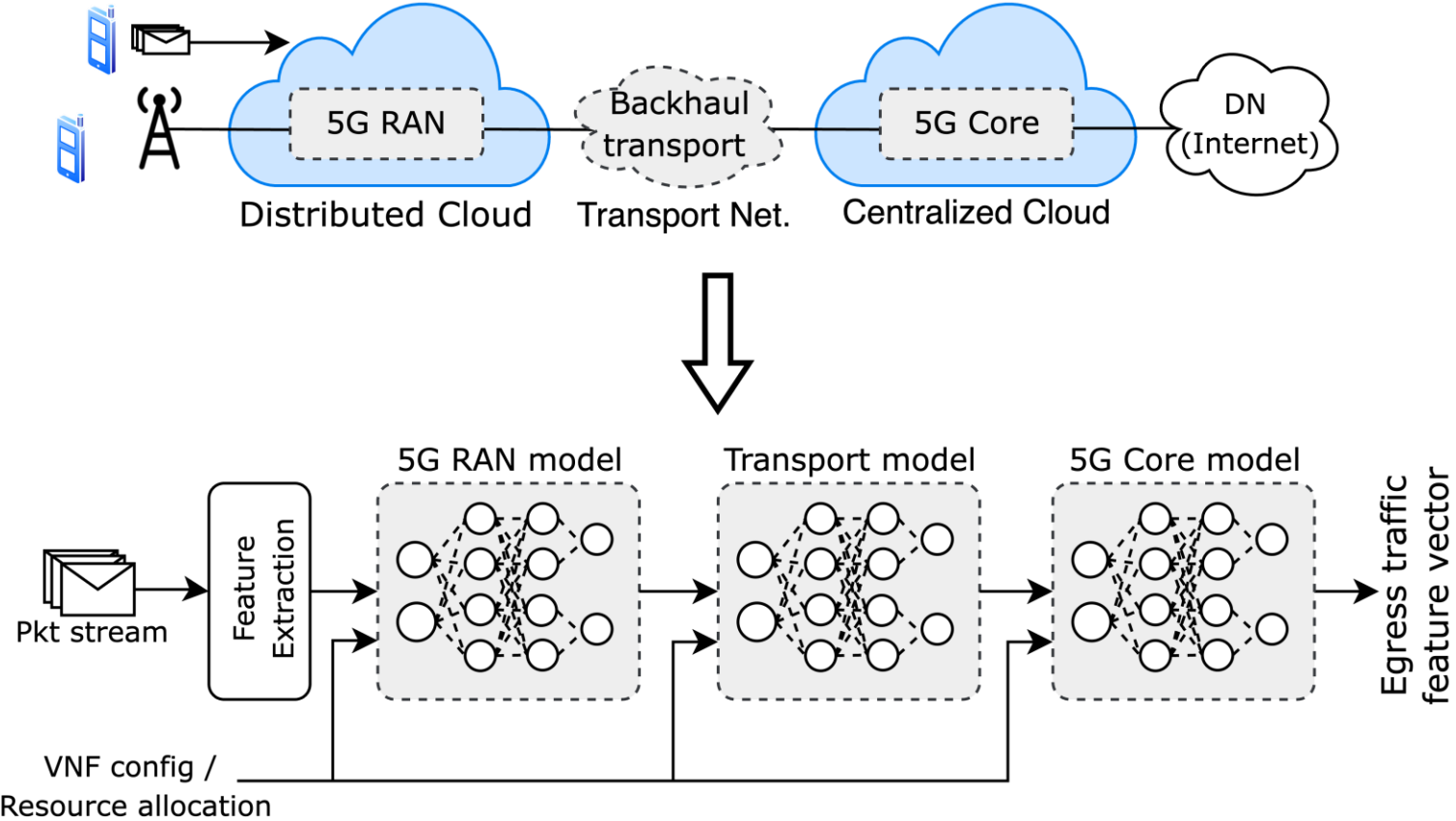### Resource Allocation

9

UNIVERSITY OF **WATERLOO**

# Network Modeling: vNetRunner (1/2)

- Slice Modeling using neural networks

- Two steps slice modeling: VNF modeling, Slice modeling using VNF models

- Step 1: Individual VNF slice modeling using neural networks

  - Reduces dataset requirement compared to slice-level modeling

  - Allows for composable VNF models

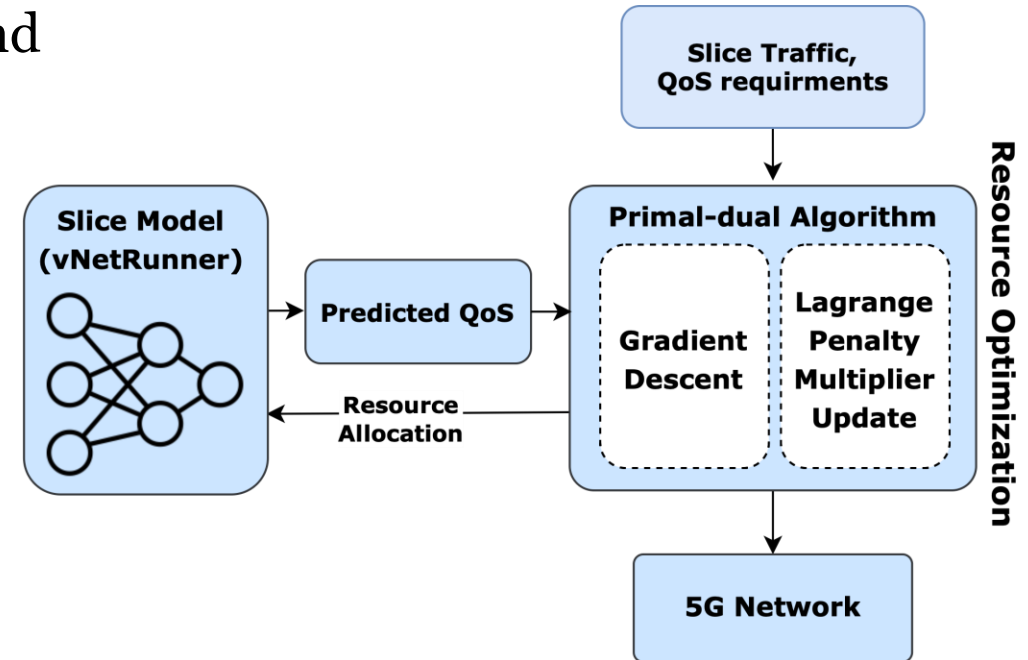  - Fast inference (milliseconds)

# Network Modeling: vNetRunner (2/2)

▪ Step 2: Composing slice model from VNF models

# Constrained Optimization: MicroOpt

- Primal-dual optimization and gradient descent for fast and efficient resource optimization

- User vNetRunner for QoS estimation

- Gradient Descent:
  - Adjusts resource allocation to minimize the overall resource usage while paying QoS violation penalty

- Lagrange Multiplier Update:
  - Ensures that QoS constraints are met by adjusting QoS violation penalties in each iteration.

# WORKSHOP SESSION#2 STRUCTURE

# Workshop Session#2 Structure

- **Part1: Data exploration and visualization**

  - Explore and visualize the resource allocation dataset gathered from in-lab 5G testbed.

- **Part2: vNetRunner**

  - Train and visualize VNF models using machine learning.

  - Use trained VNF models to compose end-to-end slice model.

- **Part3: MicroOpt**

  - Implement dynamic resource scaling with the MicroOpt framework.