

# 5G Network Slicing and Orchestration

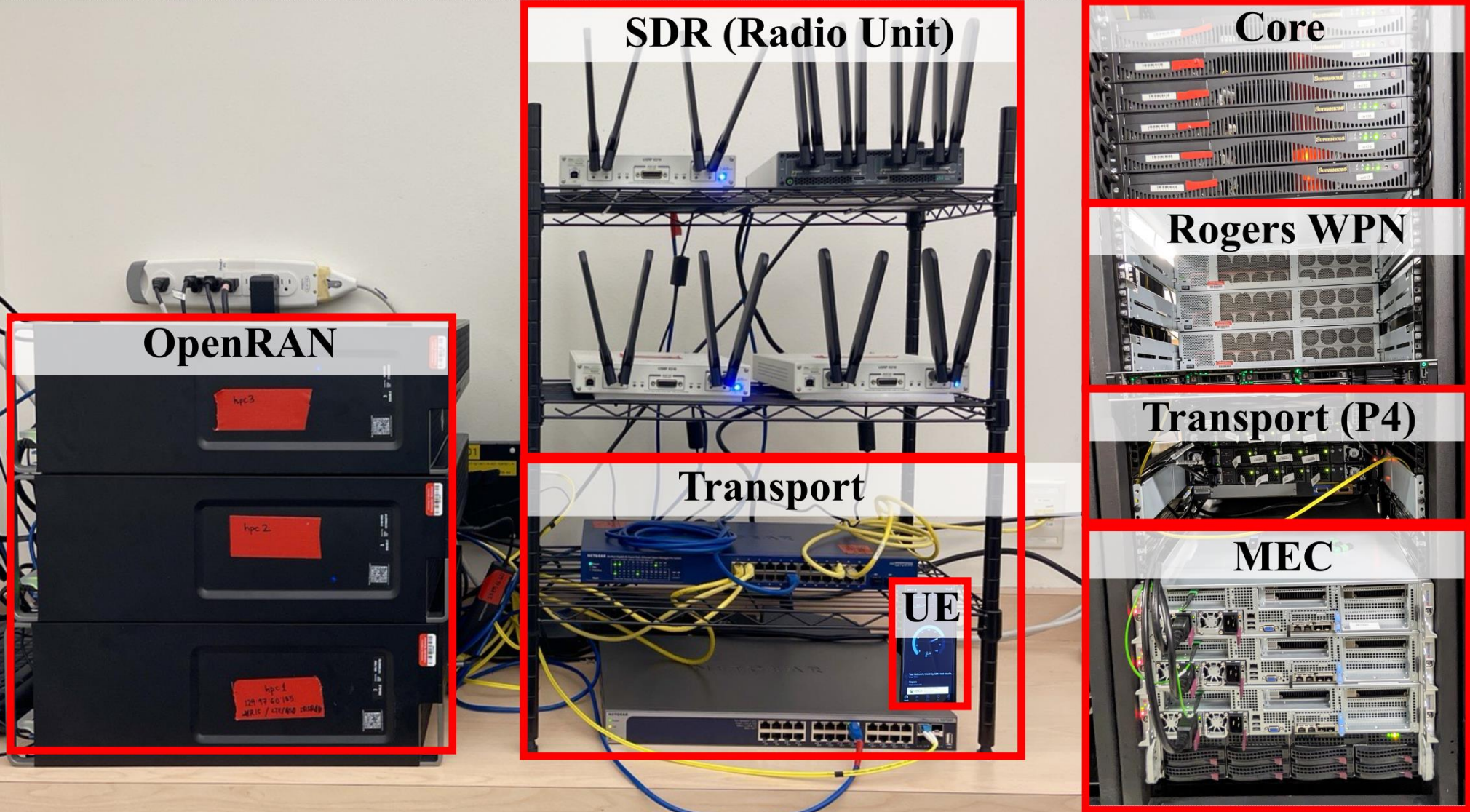
## Demo

**Raouf Boutaba**

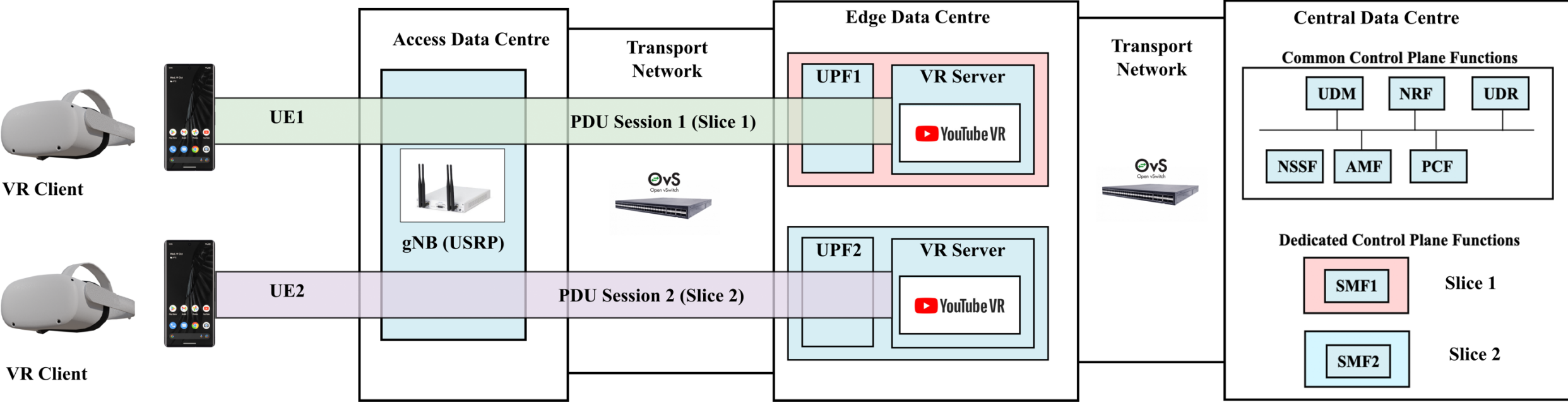
David R. Cheriton School of Computer Science  
University of Waterloo



# Network Slicing Testbed at UW

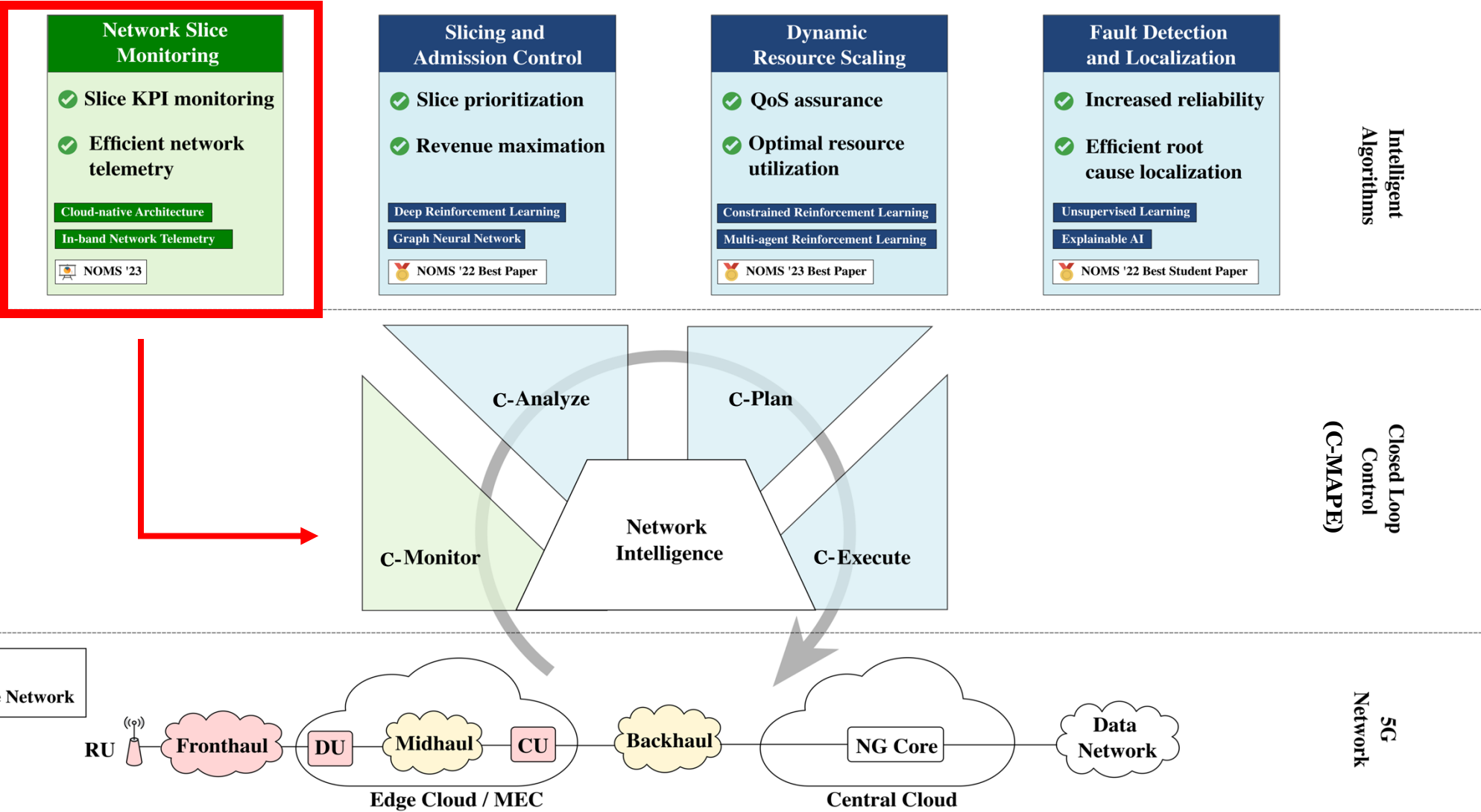


# Use-case: VR Streaming

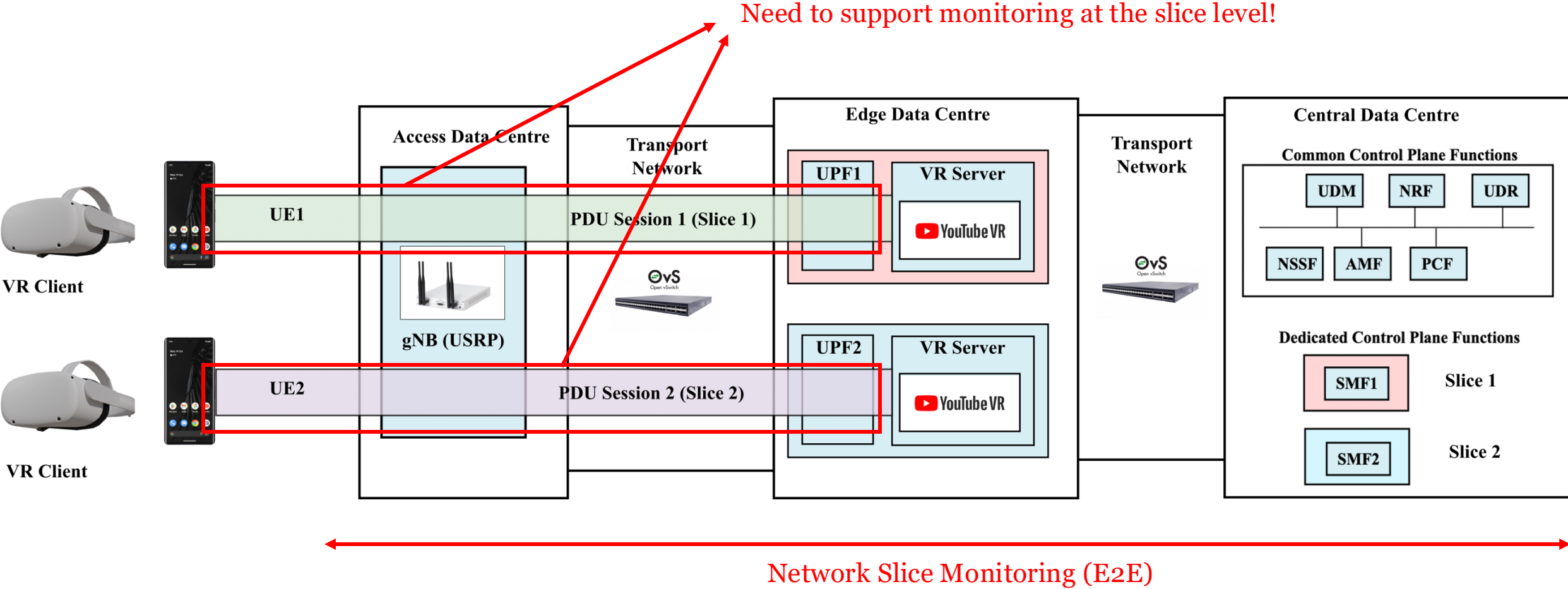


# **NETWORK SLICE MONITORING**

# Network Slice Monitoring



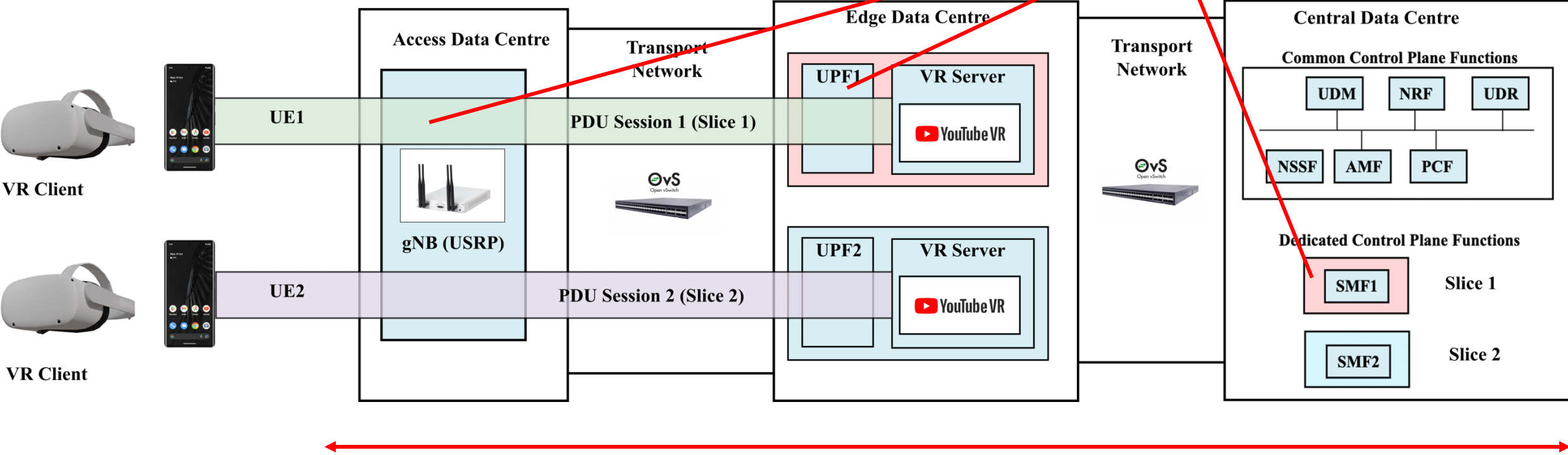
# Network Slice Monitoring





# Network Slice Monitoring

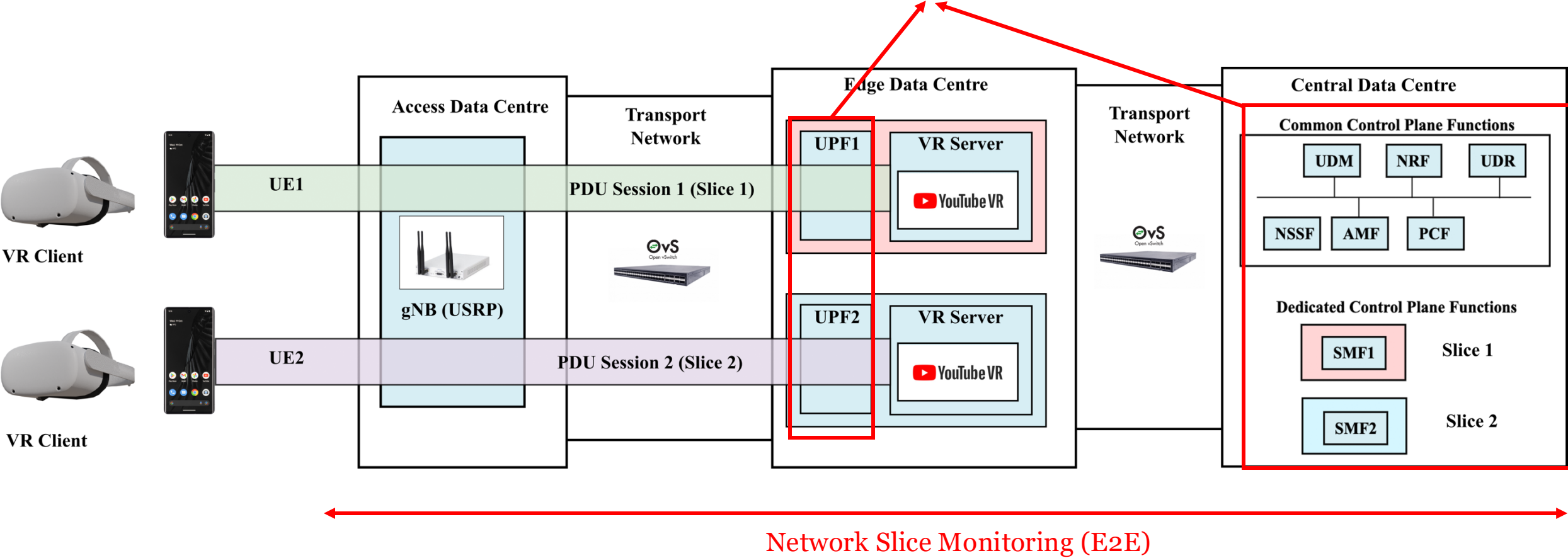
Collect and correlate monitoring data from network functions in different segments!



Network Slice Monitoring (E2E)

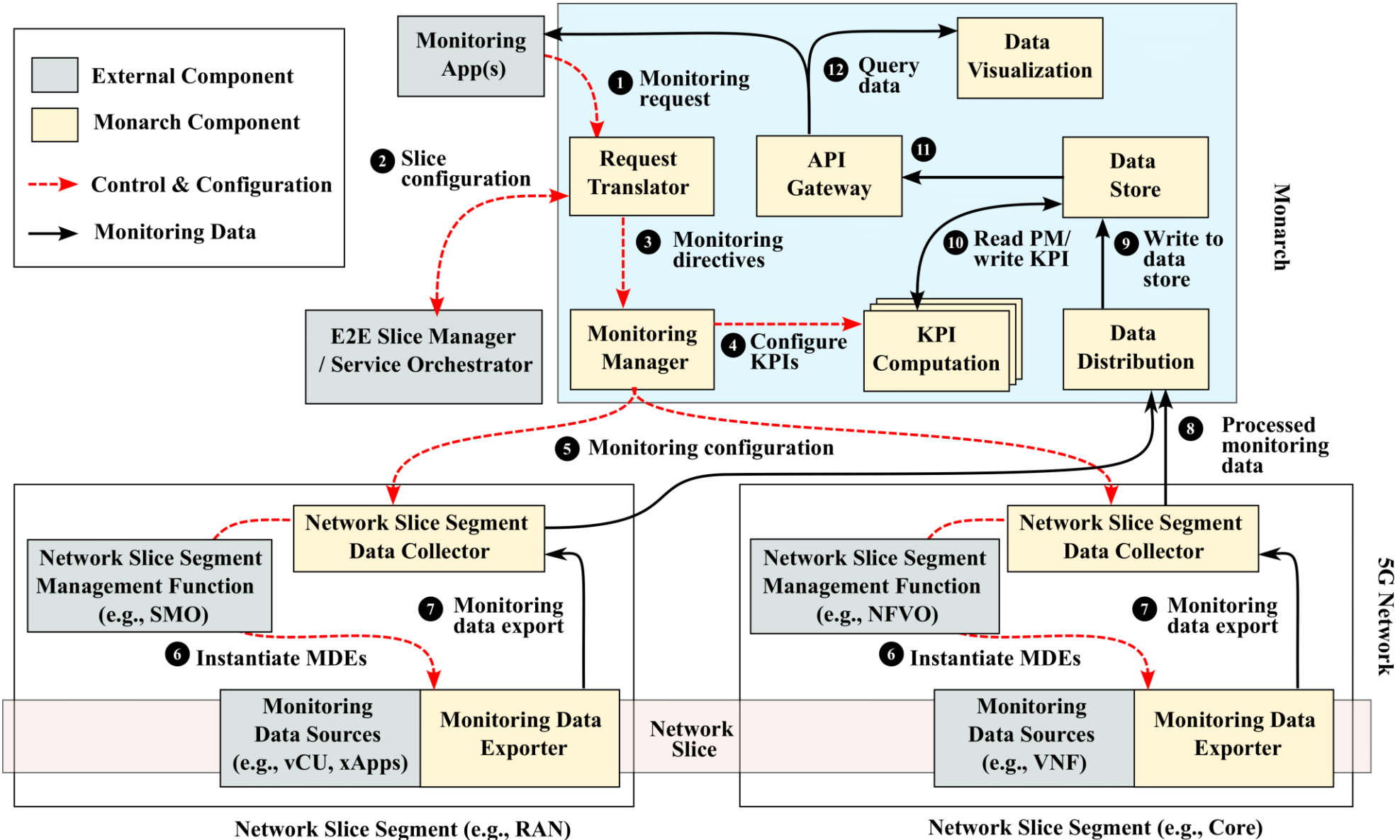
# Network Slice Monitoring

Monitoring must seamlessly integrate with cloud-native network functions!

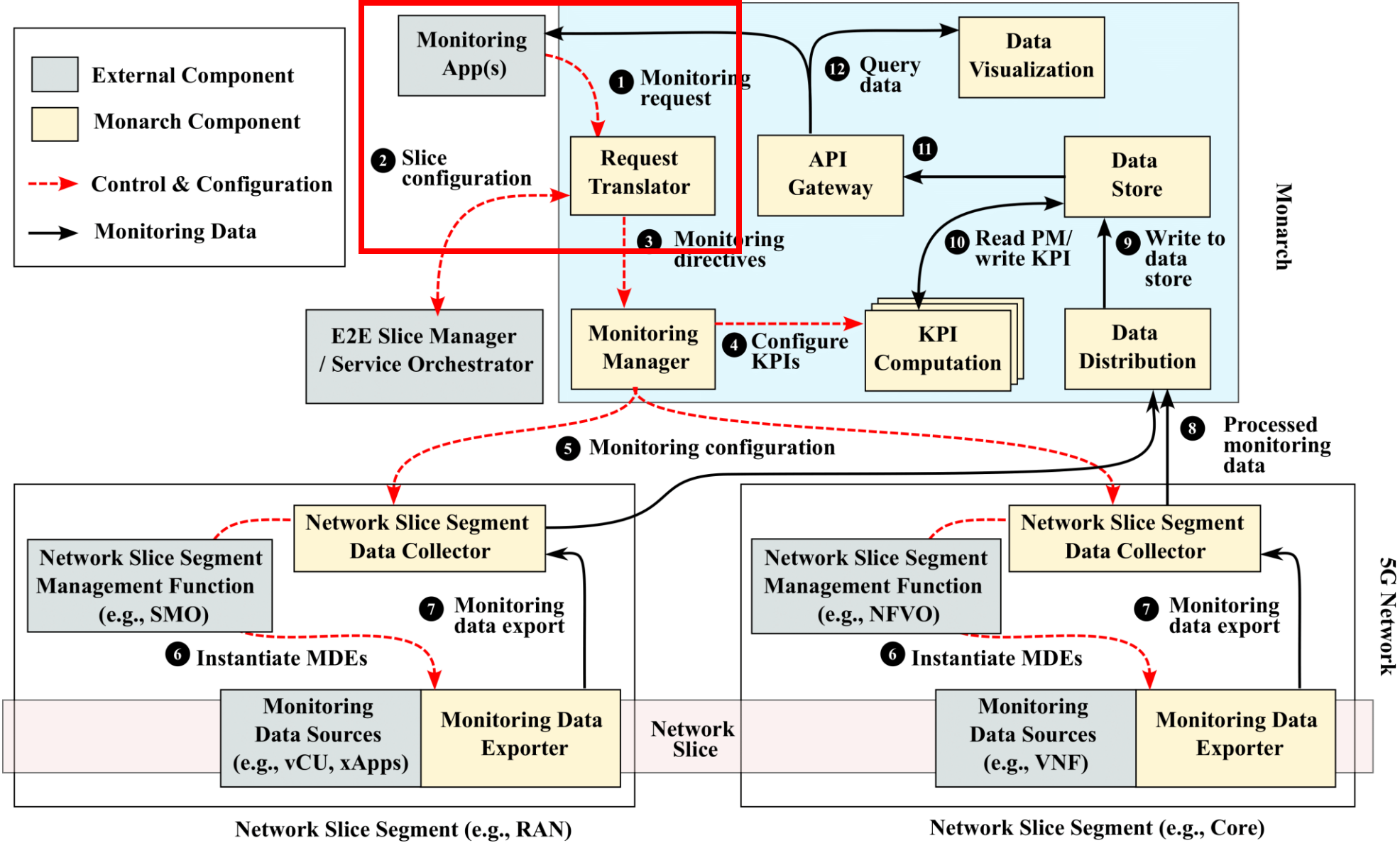




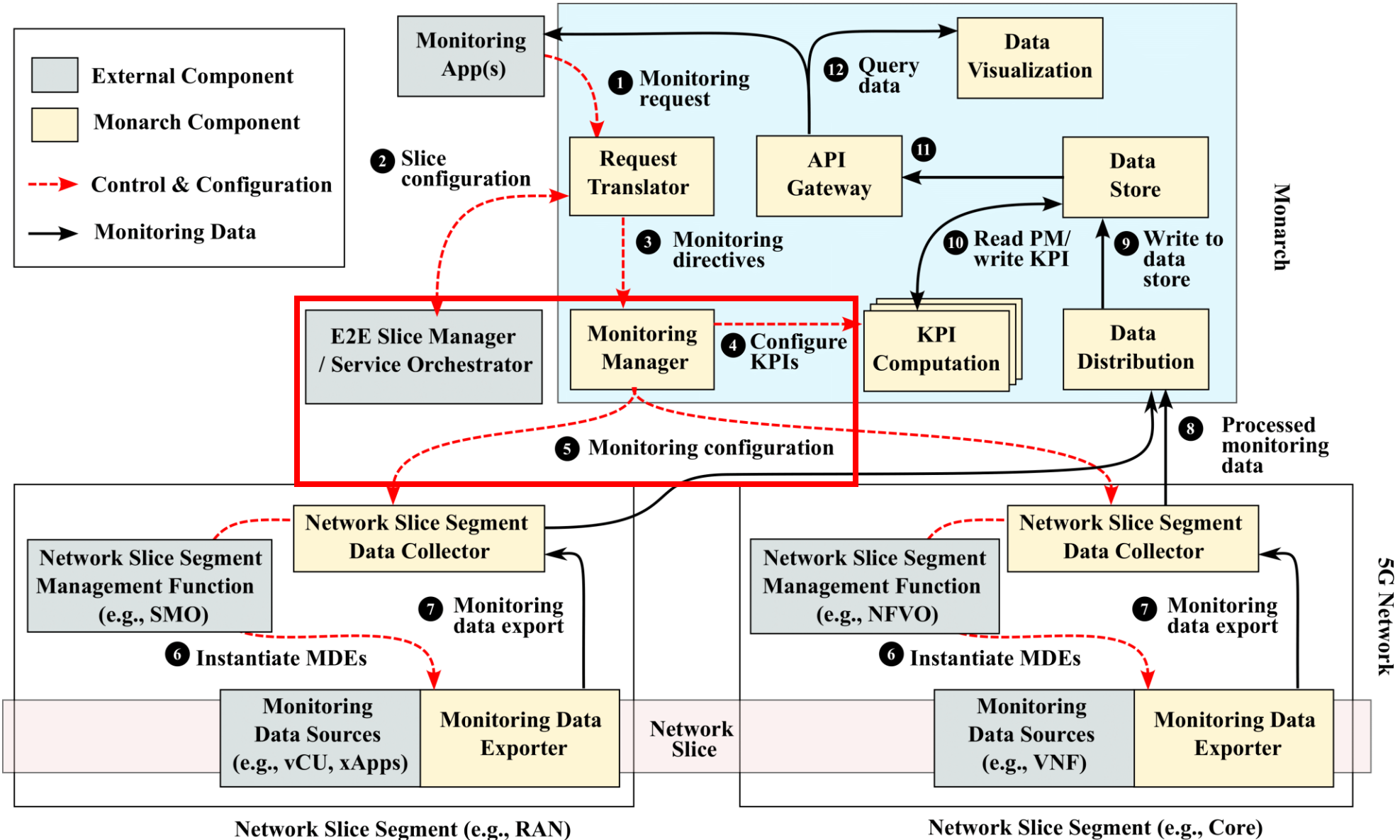
# Network Slice Monitoring with Monarch



# Network Slice Monitoring with Monarch

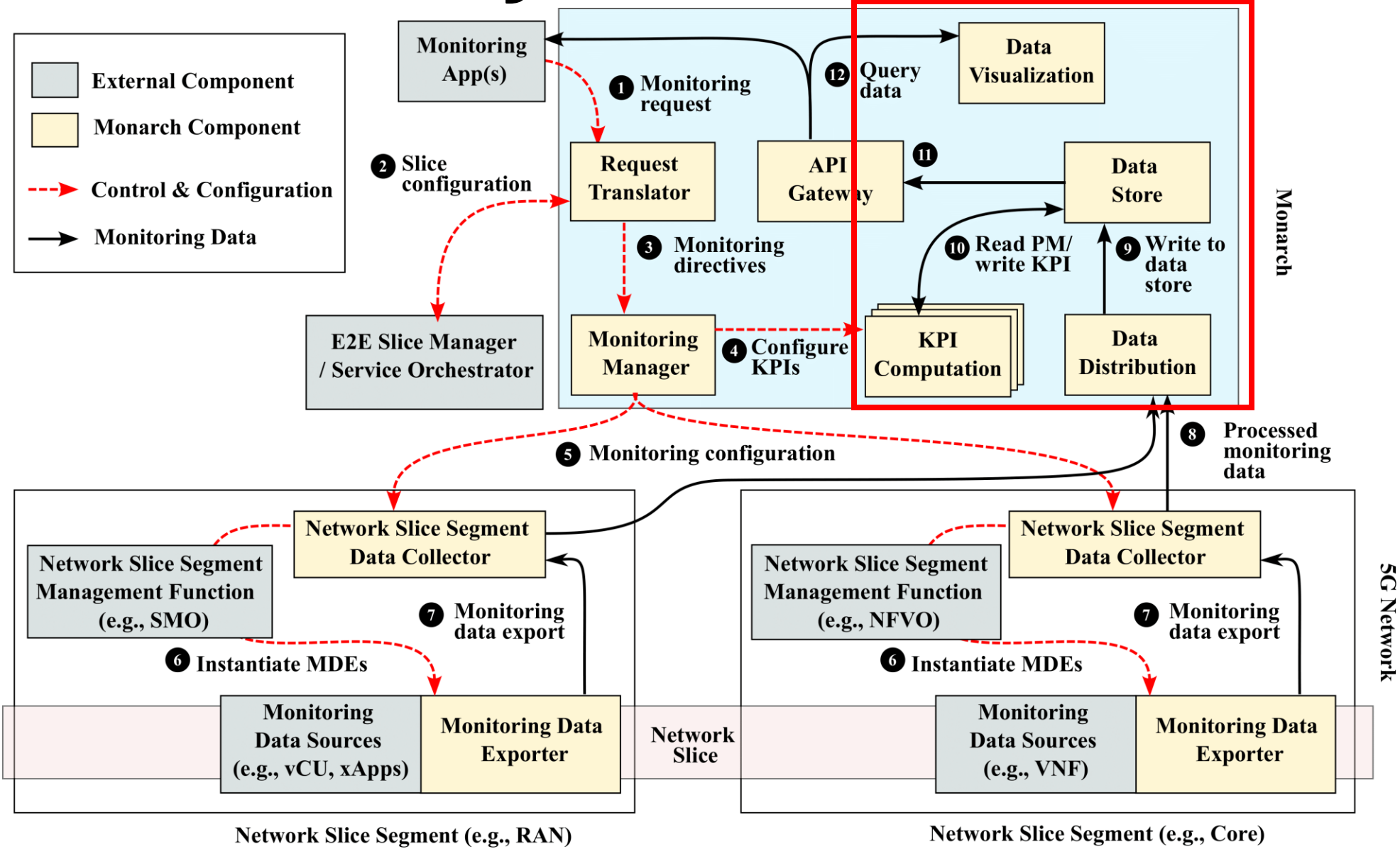


# Network Slice Monitoring with Monarch





# Network Slice Monitoring with Monarch



# **USE-CASE: VR STREAMING**



# **DYNAMIC RESOURCE SCALING**

# Dynamic Resource Scaling (1/2)

## DRS: Main Idea

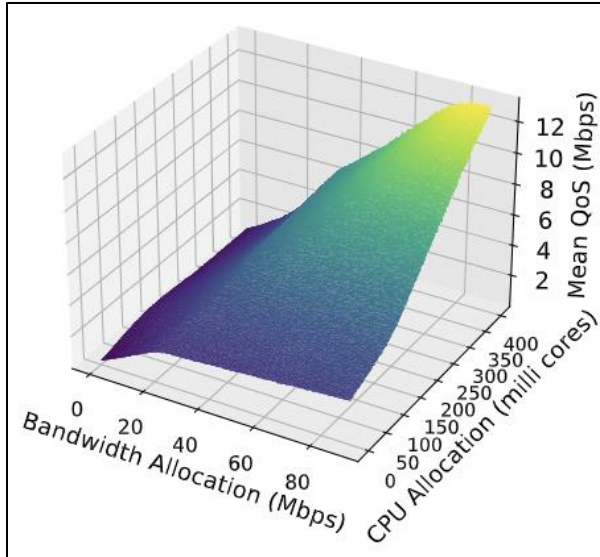
Scaling resources allocated to slices based on predicted traffic

## DRS: Objective

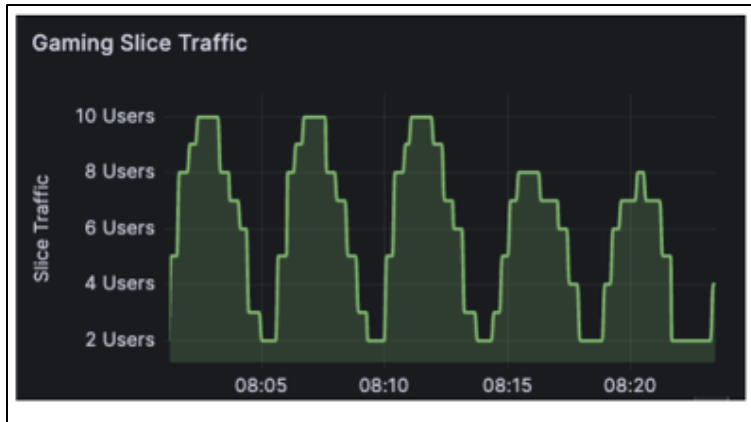
- Minimize resource allocation
- Maintain minimum QoS threshold

# Dynamic Resource Scaling (2/2)

Network Modelling



Slice Traffic Forecast



Resource Scaling Algorithm

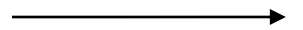
## Algorithm 1 Resource Allocation Algorithm

**Input:** Traffic  $\mathbf{x}_{\tau_1}^s$ , Network Model  $f_{QoS}^s(\mathbf{x}_{\tau_1}^s, \mathbf{r}_{\tau_1}^s)$ , 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  $\mathbf{r}_{\tau_1}^s$

- 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:    $\mathbf{r} \leftarrow \text{Gridsearch}(\mathbf{x}_{\tau_1}^s, f_{QoS}(\mathbf{x}_{\tau_1}^s, \mathbf{r}))$
- 4:   **while**  $|\nabla_r \hat{\mathcal{L}}| > \epsilon_2$  **or**  $\tau_2 < \tau_{2,max}$  **do**
- 5:      $\mathbf{r} \leftarrow [\mathbf{r} - \alpha_1 \nabla_r \hat{\mathcal{L}}]^+$
- 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,  $\mathcal{L}(\mathbf{r}, \mu, \lambda)$ )
- 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**  $\mathbf{r}$

Resource Allocation



# **USE-CASE: CLOUD GAMING**

# Use-case: Cloud Gaming



Cloud Gaming Client



Cloud Gaming Client

