# Adaptive Management in Extended Clouds

## Marin Litoiu

York University
Toronto, Canada
mlitoiu@yorku.ca
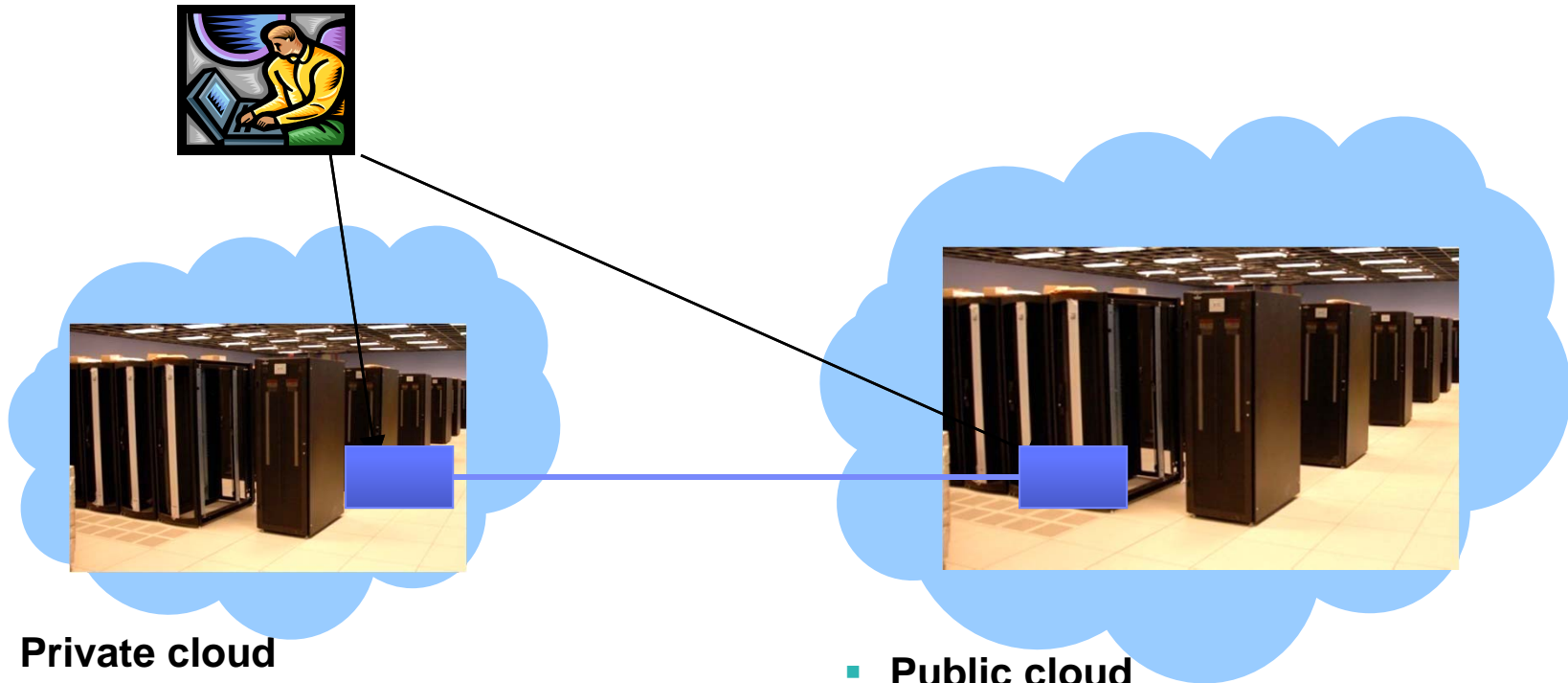
# Content

- **Extended Clouds**
  - Hybrid Clouds
  - SAVI Cloud

- **Extended Clouds Adaptive Management Platform**

- **Conclusions**

# Cloud Landscape

© Marin Litoiu

# Hybrid Clouds



- **Private cloud**
  - Limited capacity
  - Low latency
  - Privacy

- **Public cloud**
  - High capacity
  - Low cost
  - Lack of privacy
  - High latency

We are interested in applications that run/migrate seamlessly across private and public cloud

# Use Case 1: Disaster Mitigation

# Use Case 2: Cloud Bursting

- **Problem**

  - Applications run in private clouds

  - Few weeks/months a year, e-commerce applications experience high demand (think Black Friday in US, Boxing Day in Canada, etc..)

  - Private clouds cannot handle the demand

- **Solution**

  - Applications "burst" into public clouds during peak intervals

    - Applications are monitored
    - When performance degrades, components of the applications are migrated/instantiated in public clouds
    - Applications are scaled out in public cloud
    - Then they are scaled in back in private clouds, when the peak load is gone
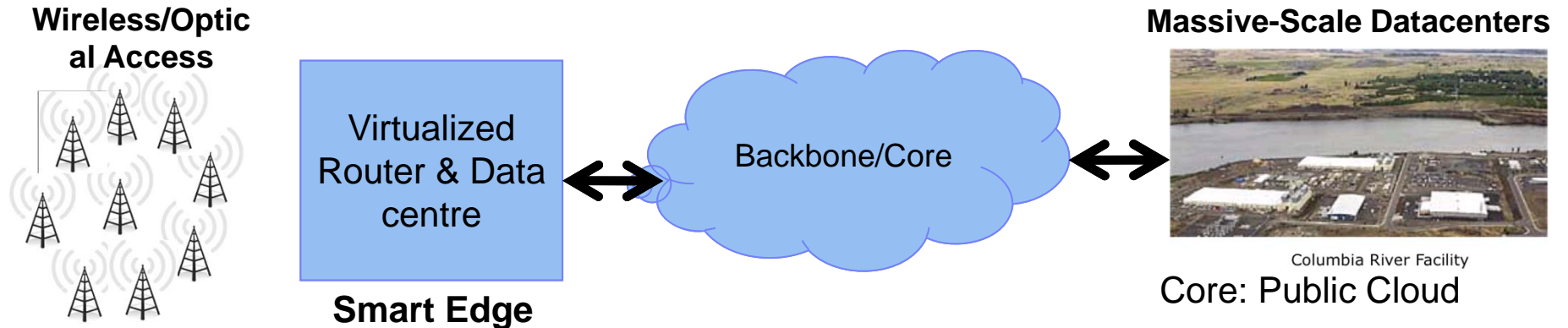
# Challenges

- **Not all parts of the application can be deployed in public clouds, due to**

  - Privacy

  - Regulation concerns

- **Need to partition the code into public and private portions**

- **Private data cannot be moved/accessed from public cloud unless is annonymized**

- **What about the network, do we have any control?**

© Marin Litoiu

# Smart Applications on Virtual Infrastructure (SAVI)

www.savinetwork.org

- **A Canadian NSERC Strategic Network**

- **8 universities, 15 companies, over 50 graduate students**

- **Several research themes**

  – Future Internet Applications

  – Adaptive Management of Applications

  – Network Management

  – Integrated Wireless/Optical Access

  – Experimental Testbed

# SAVI Goals

**Wireless/Optical Access**



**Smart Edge**

Virtualized Router & Data centre

Backbone/Core

**Massive-Scale Datacenters**

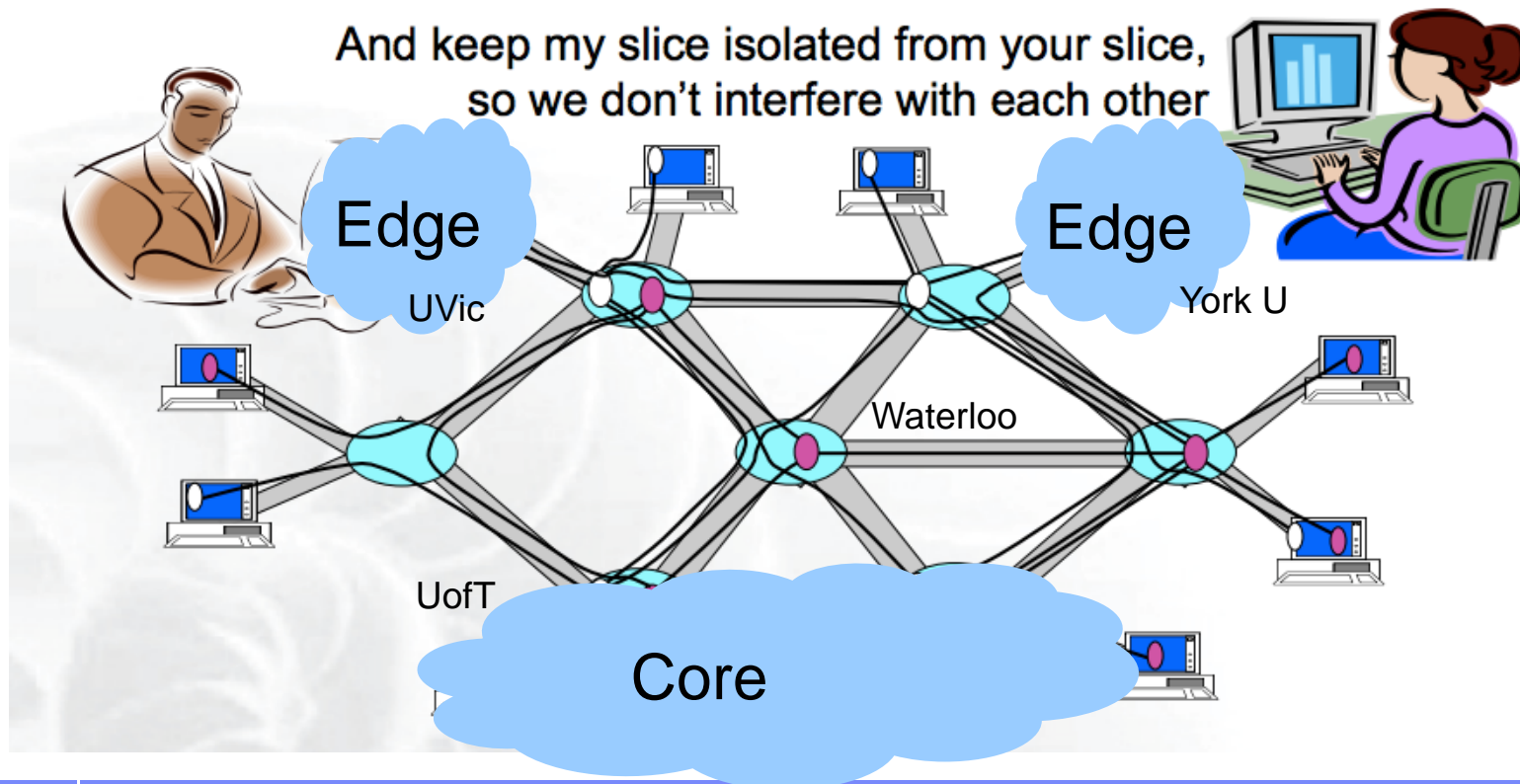

Columbia River Facility

Core: Public Cloud

- **Explore two tier clouds**
  - *Edge: low latency and high bandwidth; limited storage and computing*
  - *Core: infinite storage and computing capacity*
- **Integrated end to end adaptation (from wireless access to core cloud)**
- **Enable smart application development and deployment**
  - **Smart apps: sense the environment, analyze, predict and optimize their execution**

# SAVI Cloud: Software Defined Infrastructures (SDI)

- In SAVI, the network and the cloud converge, each cloud edge is both a cloud and a router (OpenStack and OpenFlow)
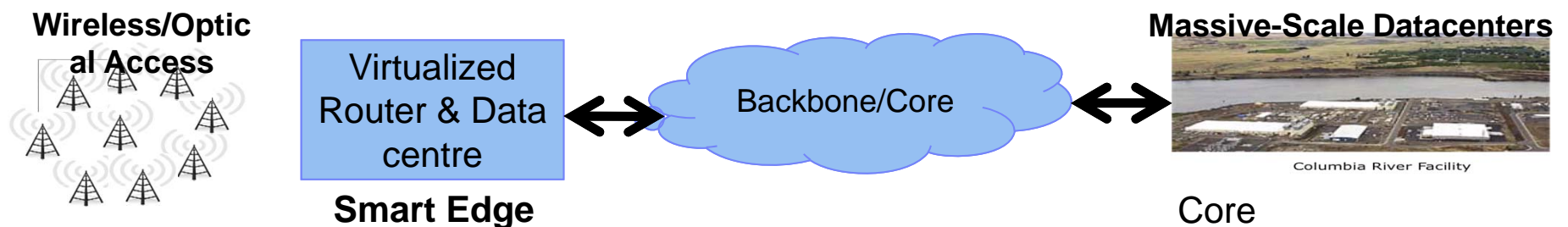


Install the software I want *throughout* my network slice
(into firewalls, routers, clouds, …)

And keep my slice isolated from your slice,
so we don't interfere with each other

Edge
UVic

Edge
York U

Waterloo

UofT

Core

© Marin Litoiu

# Use Cases for SAVI Clouds

- **Flash Crowds supporting applications**
  - 50000 people in a stadium/main square/emergency
  - 10000 people streaming video from mobiles

- **Sudden surge in demand for bandwidth, computation, storage**

- **Apps are "Smart" ( Instrumented, Interconnected, Intelligent)**
  - Monitor, Analyze, Plan and Execute loops
  - Provision/unprovision network, computing, storage

**Wireless/Optical Access**

Virtualized Router & Data centre

**Smart Edge**

Backbone/Core

**Massive-Scale Datacenters**
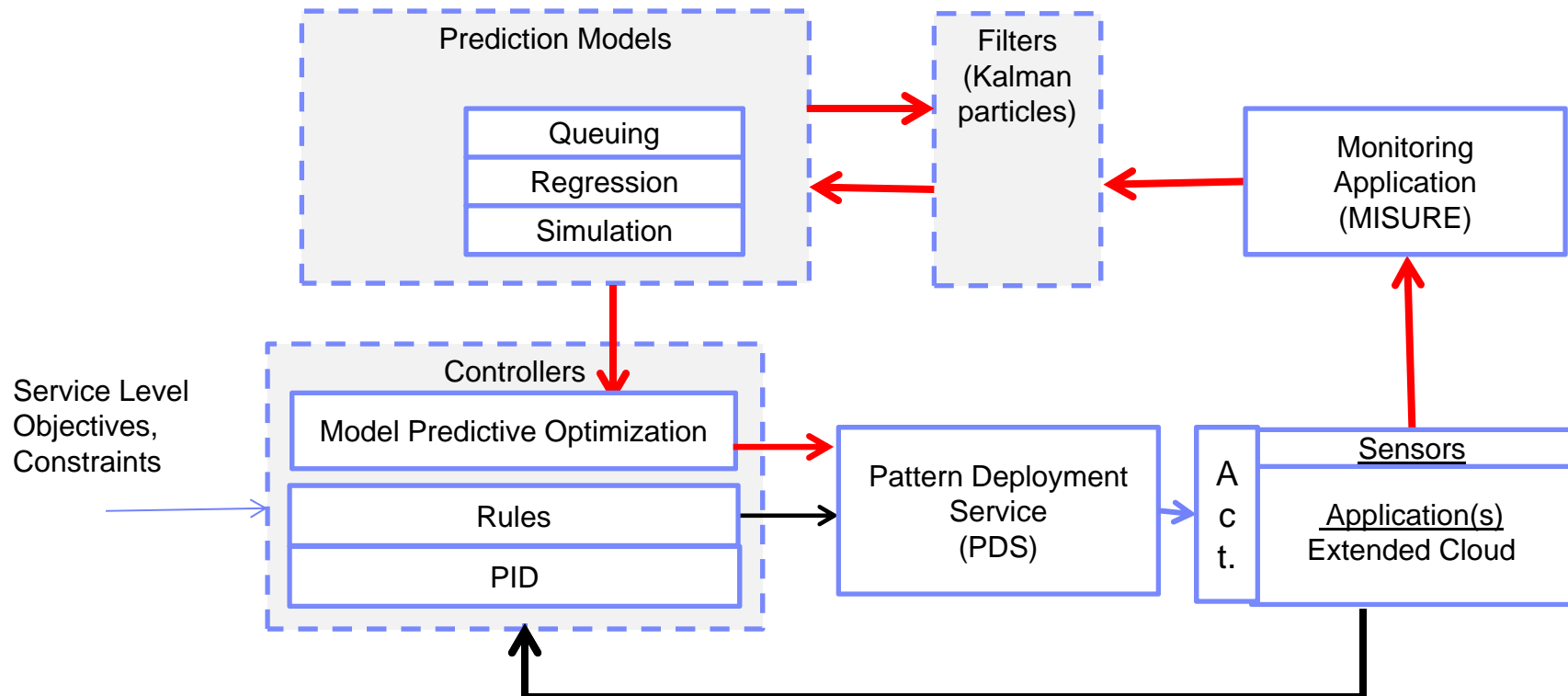
Columbia River Facility

Core

# Challenges

- **Need to partition the code into edge and core portions (performance driven)**

- **Integrate different adaptation layers**
  - Application
  - Platform
  - Network

- **Geographical location of servers and clients need to be considered**

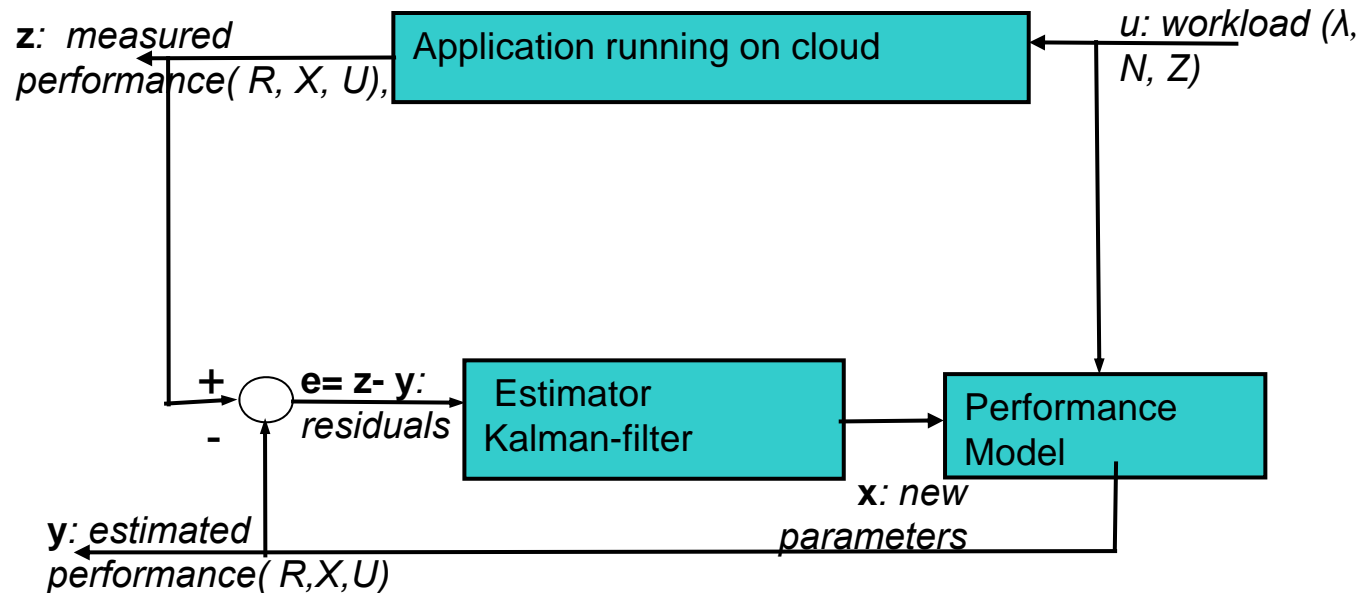# Summary so far…

- **Extended Clouds**

  – Have two tiers

    • Private-public for hybrid clouds

    • Edge-core for SAVI clouds

  – Network is programmable and part of the cloud

  – Expose many control actuators

    • E.g. application specific parameters, placement of application components, middleware parameters, network (flows and bandwidth), platform (VM migration, size), storage size and speed

  – Applications need to

    • *maintain SLOs (e.g Response_time < 100ms)*

    • *subject to constraints: cost, surging workloads, cloud topology, etc..*

    • *using an adaptive architecture (see next slide)*

# Extended Cloud Application Management Platform (XCAMP)



- **Reactive(black arrows):reacts to current load; implements simple controls(PID=proportional, integrative, derivative); fast but imprecise**

- **Predictive (red arrows): anticipates future load, performance, cost**
  – Uses prediction models, filters and predictive optimization. It is slow and effortful but efficient

# Parameter Estimation and Tracking



z: measured performance( R, X, U),

Application running on cloud

u: workload (λ, N, Z)

$+$ $e= z- y$: residuals

$-$

Estimator Kalman-filter

Performance Model

x: new parameters

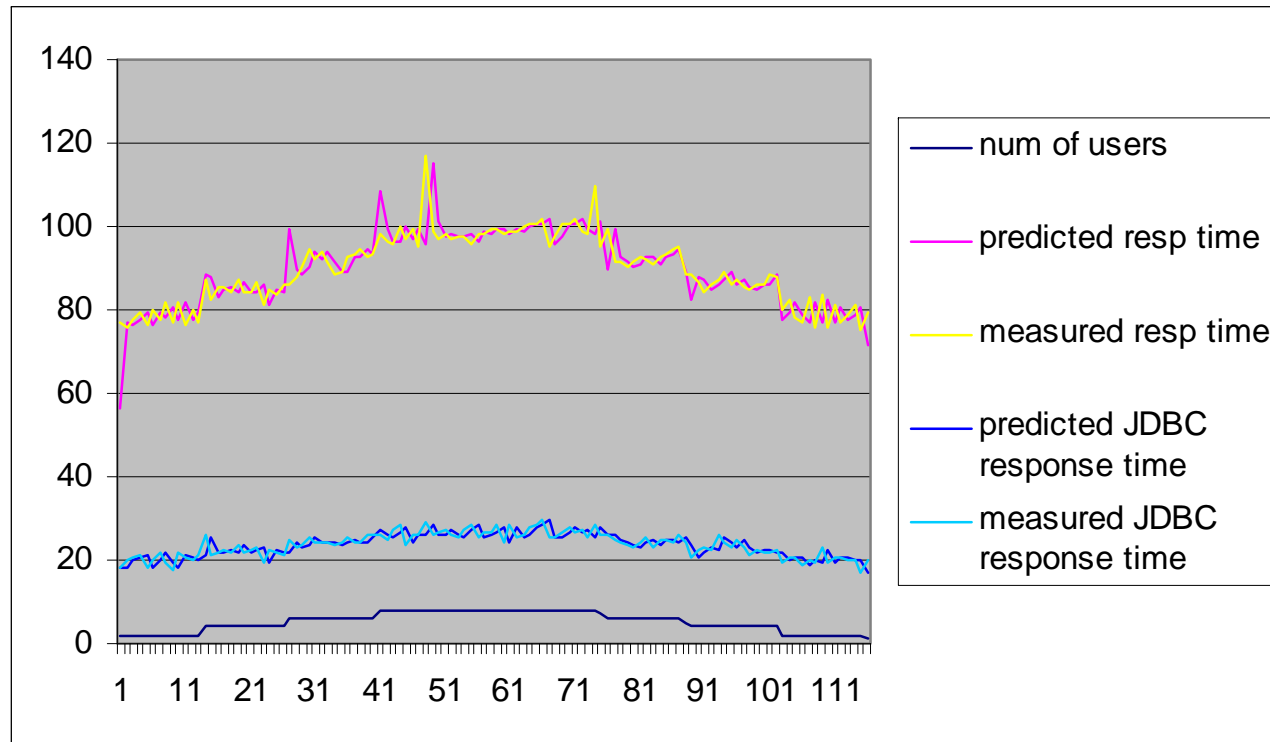y: estimated performance( R,X,U)

Parameter estimator (Kalman filter): a feedback based system, based on past and current data from the system

Continuously updates the parameters:

      - compares the measured and estimated performance metrics (e)

      - adjusts the parameter (state) of the model such that e~0.

Kalman estimators used in radar/missile tracking, autopilot, computer vision,etc.
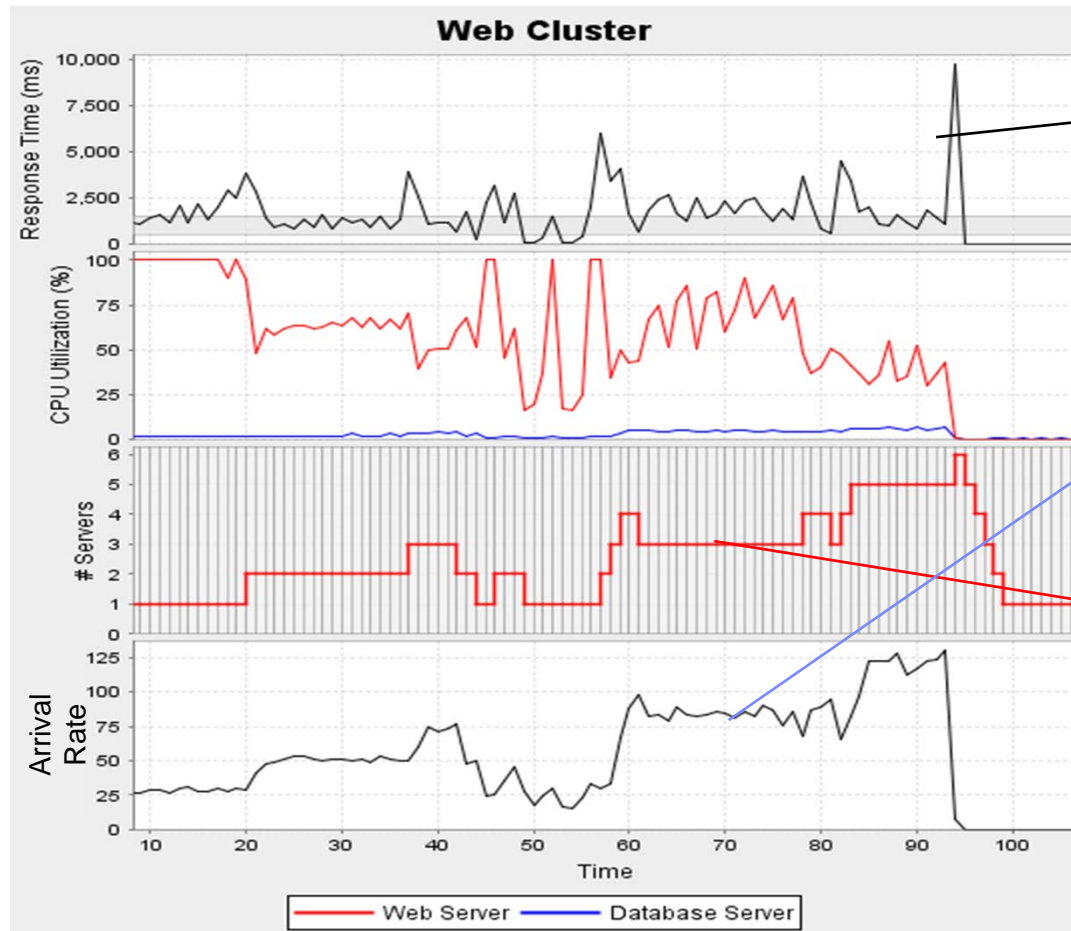
# Model + Estimator: Accuracy



- Measured: servlet response times and CPU utilizations on both tiers, throughput
- Estimated: transaction demands at each tier, no of invocations

# Managing Web Applications Deployed  SAVI Cloud

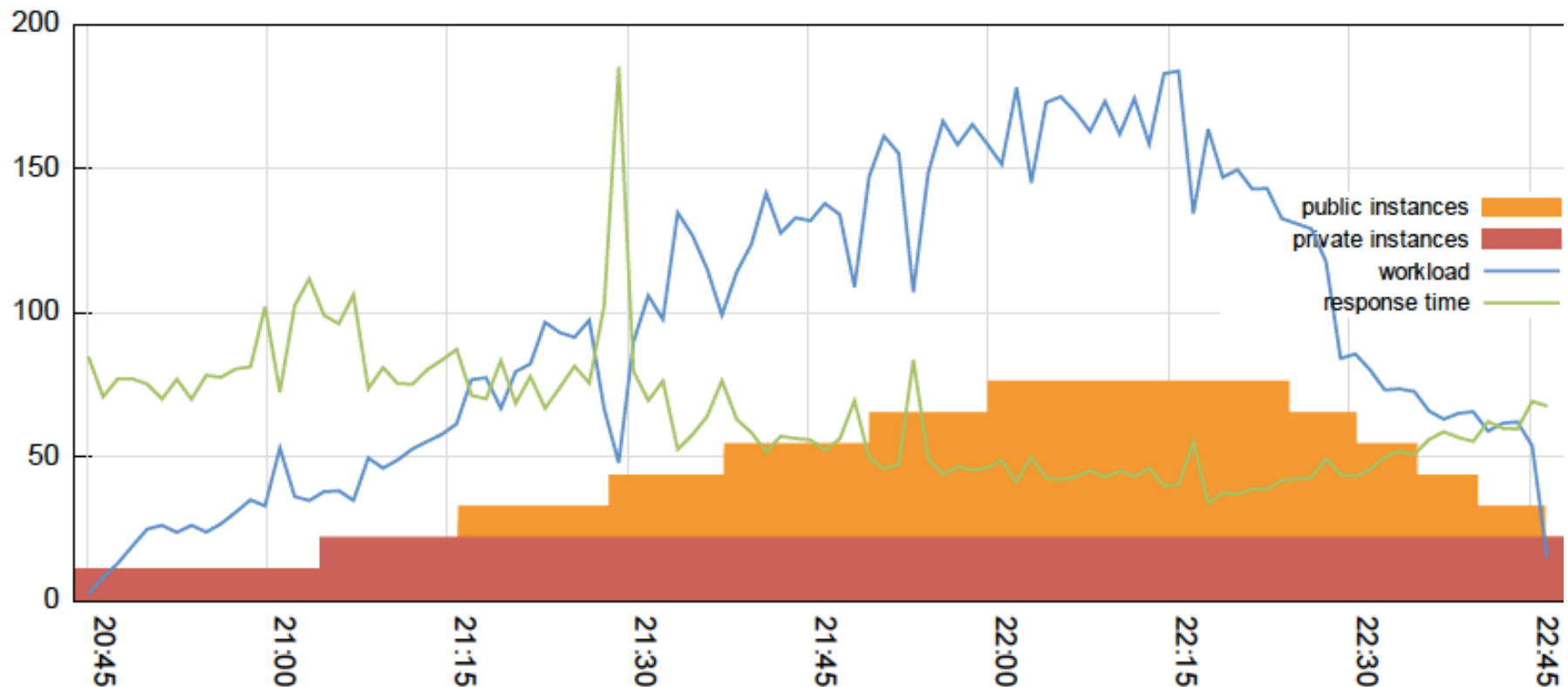**Cost and performance**



SLO: response time <1500ms

…

As workload changes

…

…

VM are added and removed

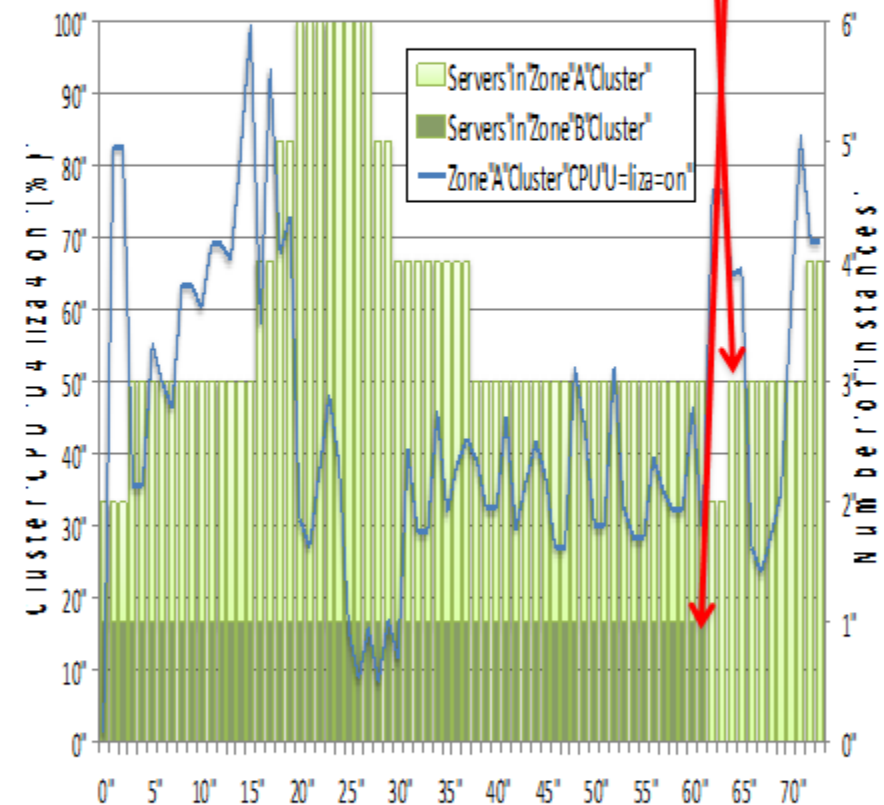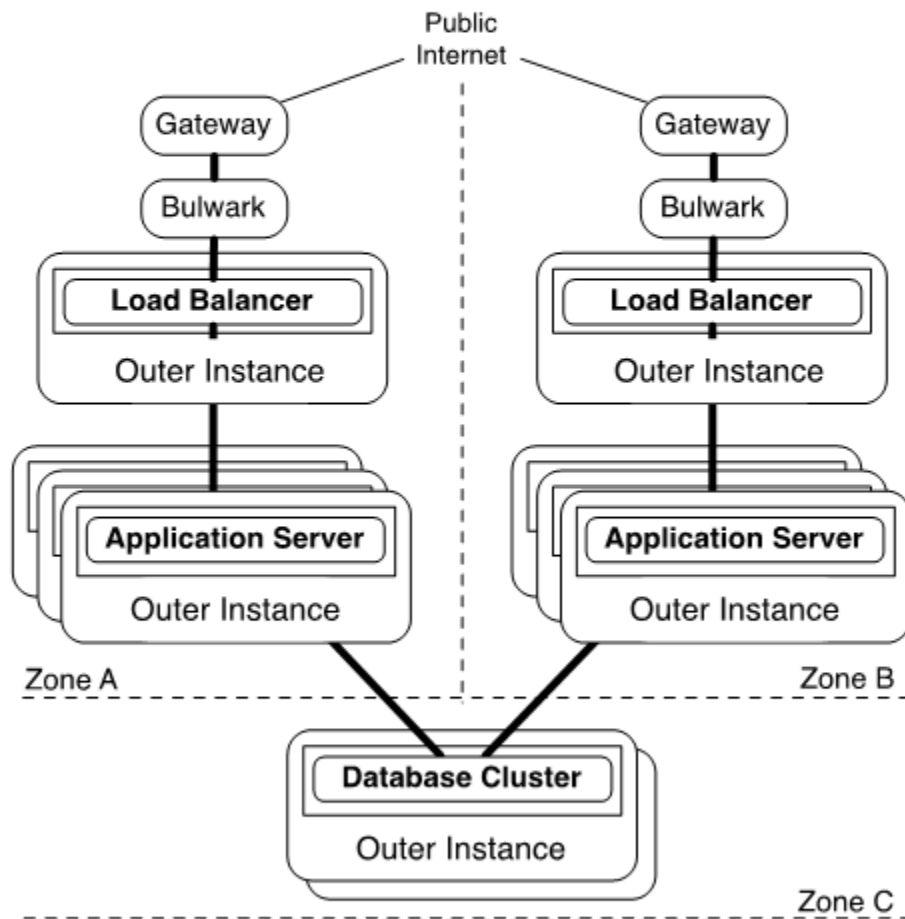# Cloud Bursting with XCAMP: Results



Private cloud: York Univ (IBM Blade Centre running OpenStack)
Public cloud: Amazon
Bursting architecture: Monitor (Misure),   Deployer (PDS),
Controller(XCAMP)

# Achieving Availability with XCAMP*

QEMU emulator version 0.14.1, TrueCrypt 7.1a, OpenVPN 2.2.0, OSSEC version 2.6, lbcd version 3.30, Snort version 2.8.5.2, and both mod-security version 2.6.0 (using default/s- tandard community rules) and mod-evasive version 1.10.1 for Apache.

# Summary

- **Extended Clouds**
  - Have two (or more) tiers
    - Private-public for hybrid clouds
    - Edge-core for SAVI clouds
  - Network is integral part of the cloud

- **Expose many sensors and actuators**
  - E.g application specific parameters, placement of application components, middleware parameters, network (flows and bandwidth), platform (VM migration, size), storage size and speed
  - Constraints: cloud topology, geographic distribution of clients and servers, cost, etc..

- **Multiple complementary feedback loops might be needed**

- **Predictive adaptation mitigates delays and long term goals such as cost and revenue**

  - Performance models

  - Filters (Kalman, particle)

  - Model Predictive Optimization

# Acknowledgements

- **ASRL Team**
  - Post Doctoral Fellows: Brad Simmons, Mike Smit , Mark Shtern
  - Graduate Students: Hamoun Ghanbari, Cornel Barna, Przemyslaw Pawluk, Mona Yousefian , Parisa Zoghi, Vasileios Theodorou, Hongbin Lu, Mihai Iacob , Mircea Constantinescu

IBM

ca

SAVI

amazon
web services™

JUNIPER
NETWORKS

NSERC
CRSNG  People. Discovery. Innovation.

YORK U
UNIVERSITÉ
UNIVERSITY
redefine THE POSSIBLE.

Ontario Centres of Excellence

# 2013 IEEE TCSE - Vote by September 15

http://www.cs-tcse.org/

**Candidates:** http://www.computer.org/portal/web/tandc/tcse

**TCSE membership:** https://supportcenter.ieee.org/app/answers/detail/a_id/353/session/L3RpbWUvMTM3NTc0NDYxMi9zaWQvdXlxcVAxeGGw%3D

**Vote by Sep 15:** http://www.surveymonkey.com/s/77HW8TD

© Marin Litoiu