# UMEÅ UNIVERSITY

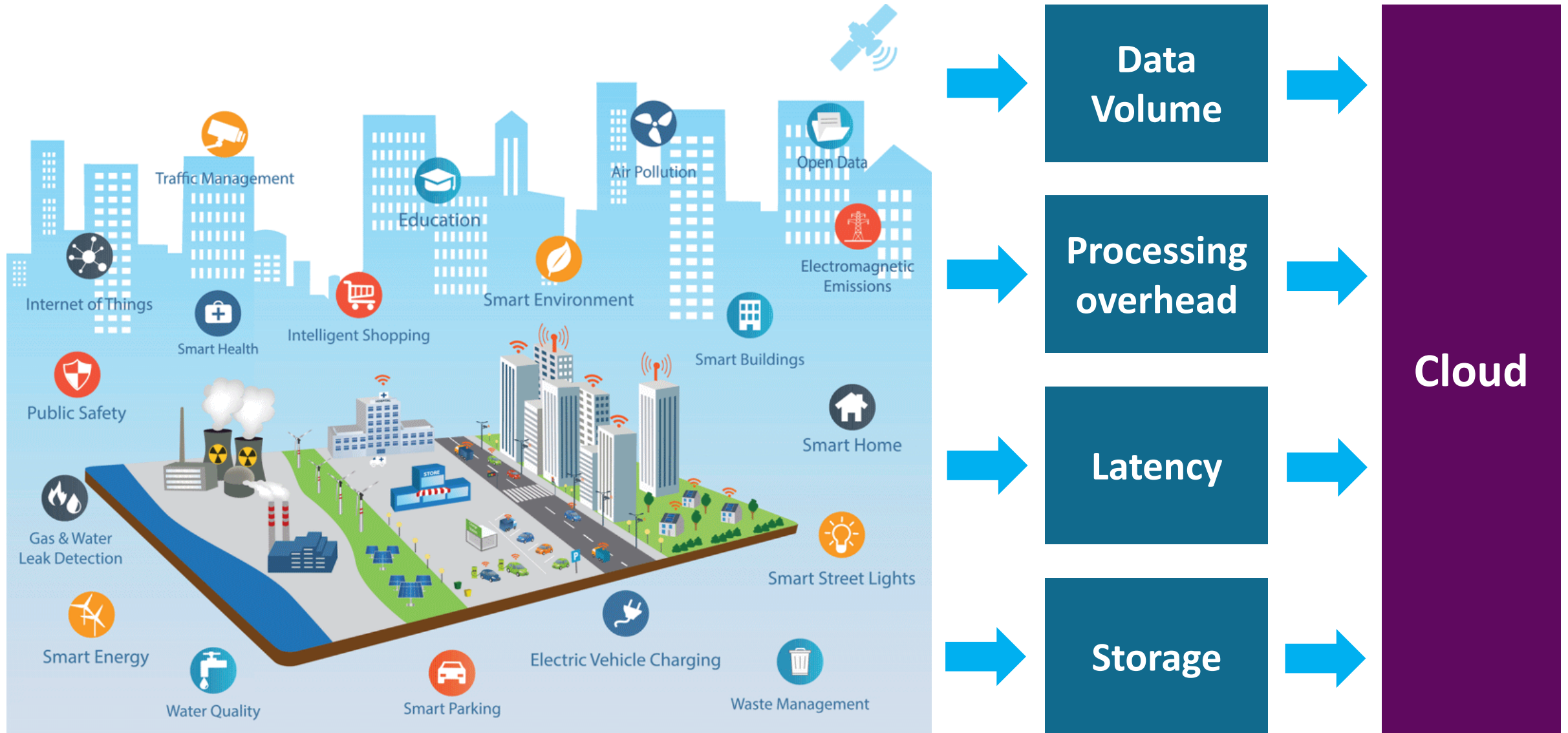# MODELLING AND MANAGING THE FUTURE CLOUD-EDGE CONTINUUM

## PAUL TOWNEND

UMEÅ UNIVERSITY, SWEDEN
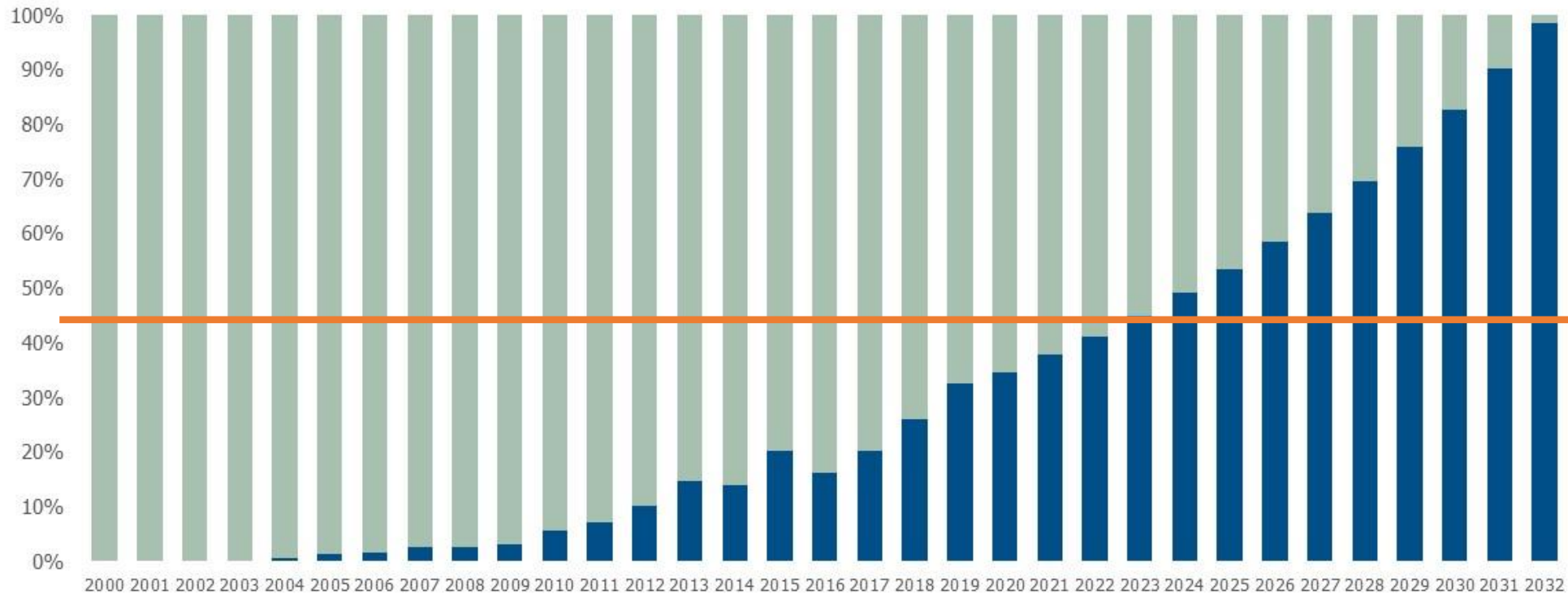
# Cloud is eating software
## Cloud will become majority of software market within 5 years



Source: CapIQ; Bessemer Venture Partners analysis;
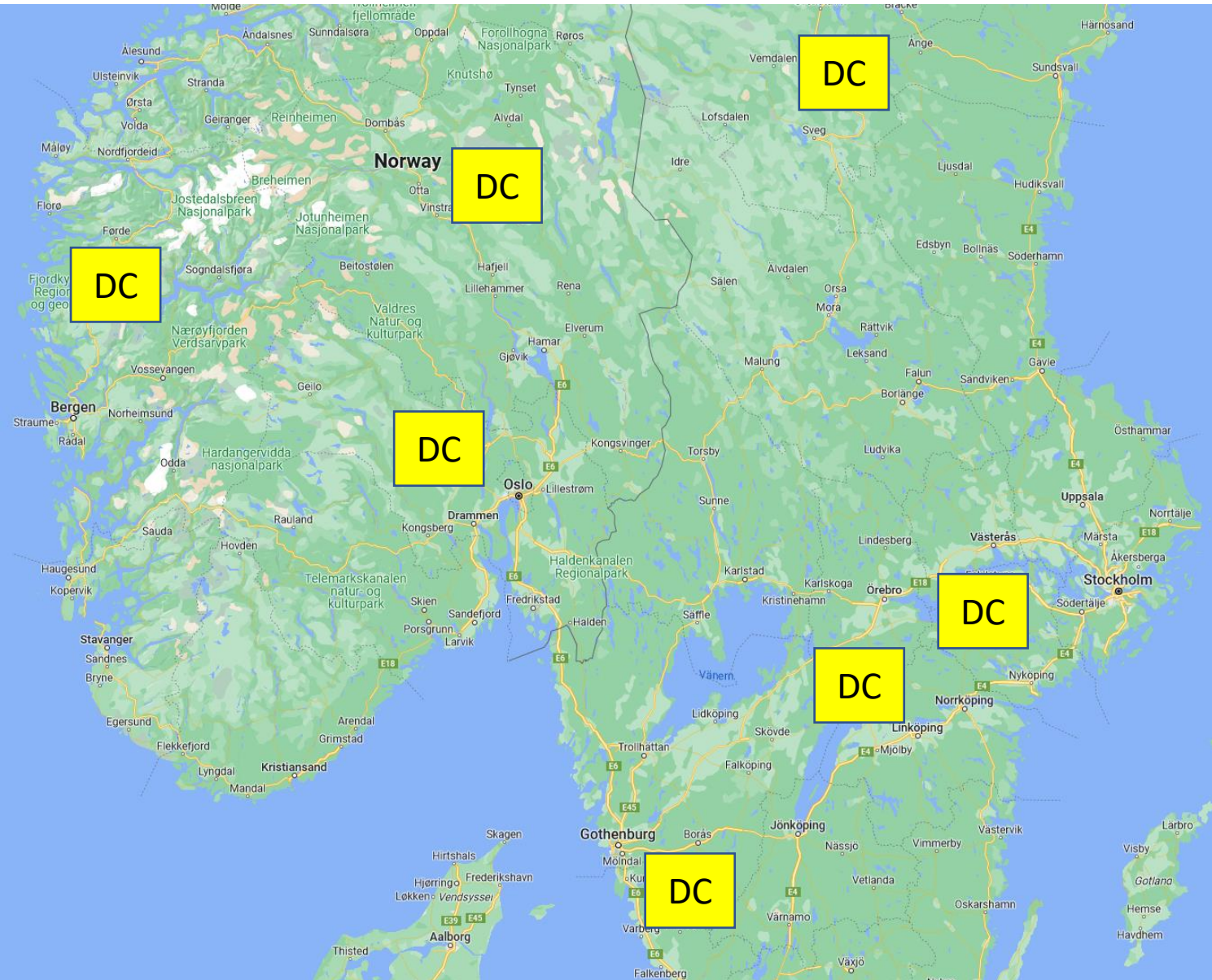Cloud CAGR – 20%, Software CAGR – 10%

Software    Cloud

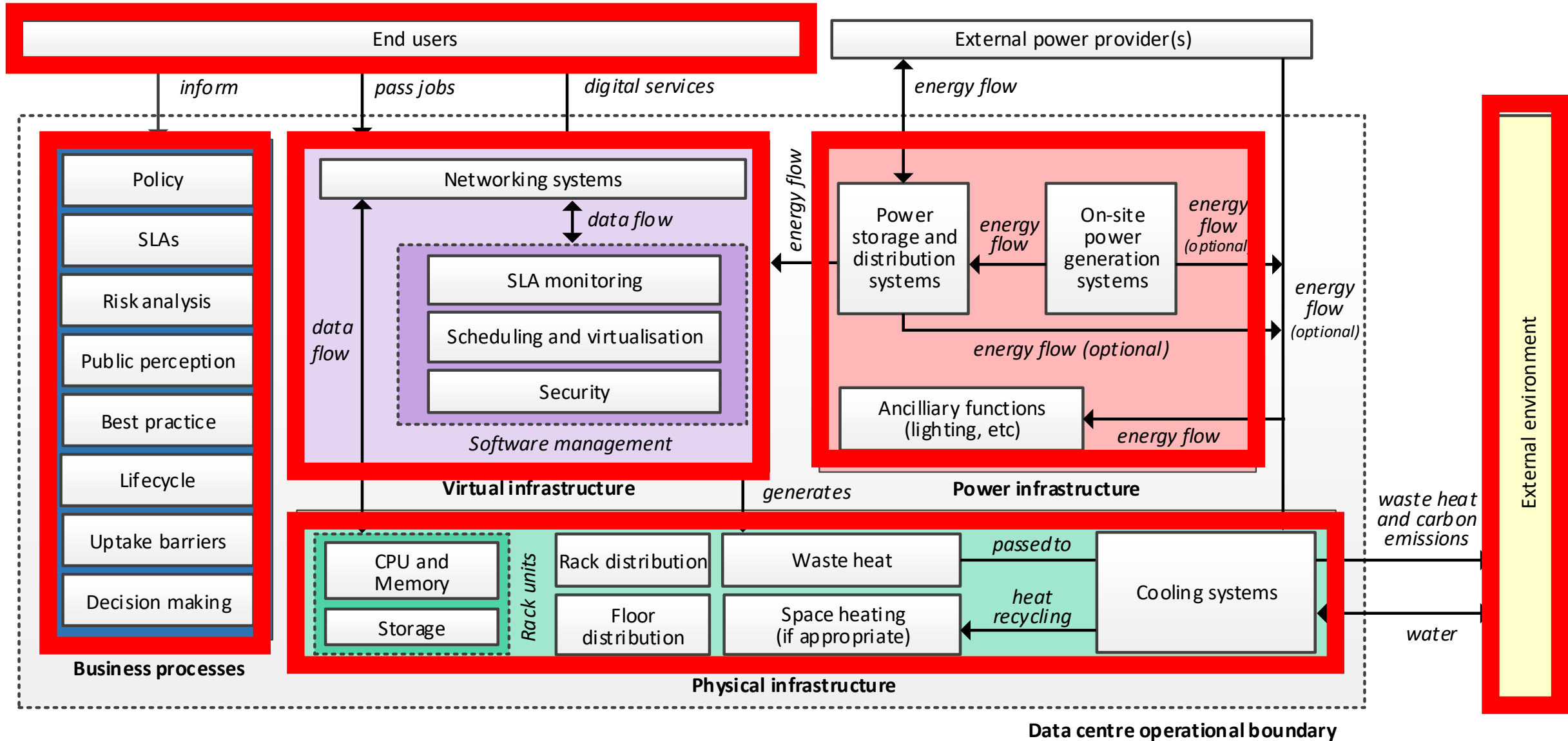**Users connect to large data centres**

(e.g. Facebook, Google, etc.)

**These are powerful resources**

Lots of servers, storage, bandwidth, etc.

**Expensive, slow to build, huge impact**
on power grids and the environment

GOOGLE DATA CENTRE, HAMINA, FINLAND

**Modelling and Managing the Future Cloud-Edge Continuum**  -  Paul Townend

**DataCenter Knowledge**™

**RECENT**

**The Biggest Problem in AI? Lying Chatbots**

MAY 30, 2023

**Biden's Former Tech Adviser on What Washington Is Missing About AI**

MAY 30, 2023

**HPE and Ampere Take Aim at Intel With Vision of Arm-Based Open RAN Server**

MAY 25, 2023

**Amazon's Answer to ChatGPT Seen as Incomplete**

MAY 24, 2023

**Are Data Centers Taking Over Oregon's**

COMPANIES  >  GOOGLE (ALPHABET)

# Google Using Sea Water to Cool Finland Project

Google will use cool sea water in the cooling system for its new data center in Hamina, Finland, which may be the first sea-cooled data center. The initiative continues Google's focus on data center efficiency and sustainability.
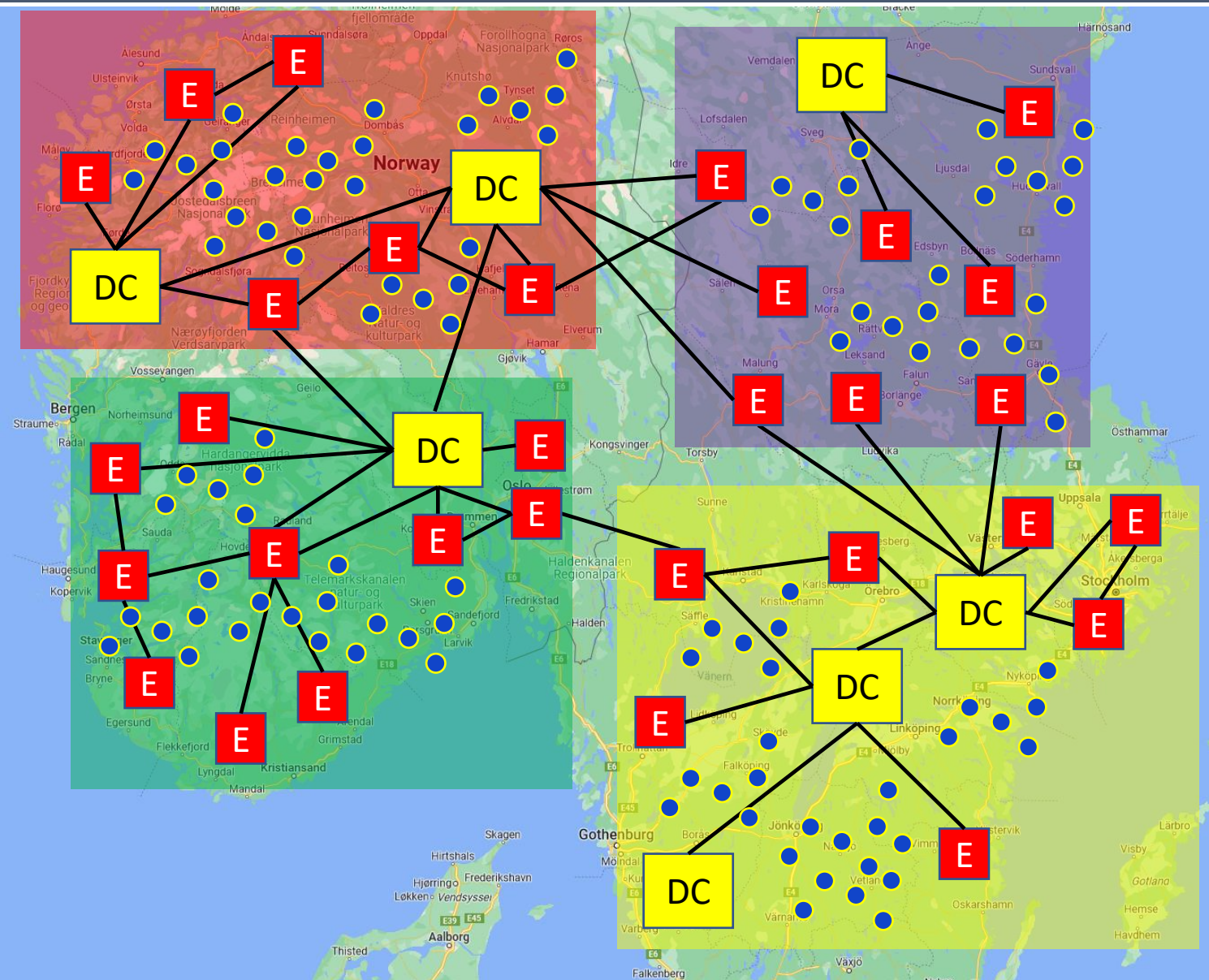
Rich Miller | Sep 15, 2010

**Google** will use cool sea water in the cooling system for its new data center in Hamina,

# The Cloud-Edge Continuum and its characteristics

**Modelling and Managing the Future Cloud-Edge Continuum** - Paul Townend

**50 billion connected devices by 2025**
This would overwhelm centralised DCs

**Proliferation of "intermediate" edge DCs**

**Complex federations of devices**

**Regional considerations**

UMEÅ UNIVERSITY

Multiple disparate providers

Platform heterogeneity

Resource constrained devices

Infrastructural dynamicity
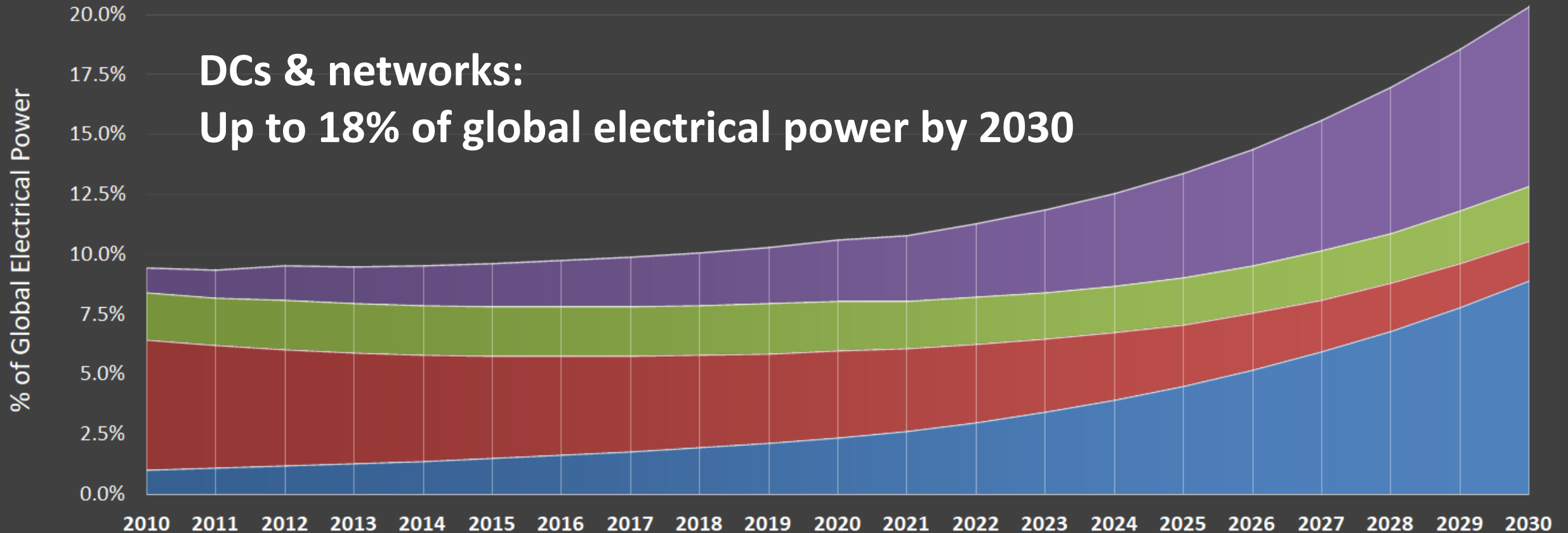
Secure orchestration over public networks
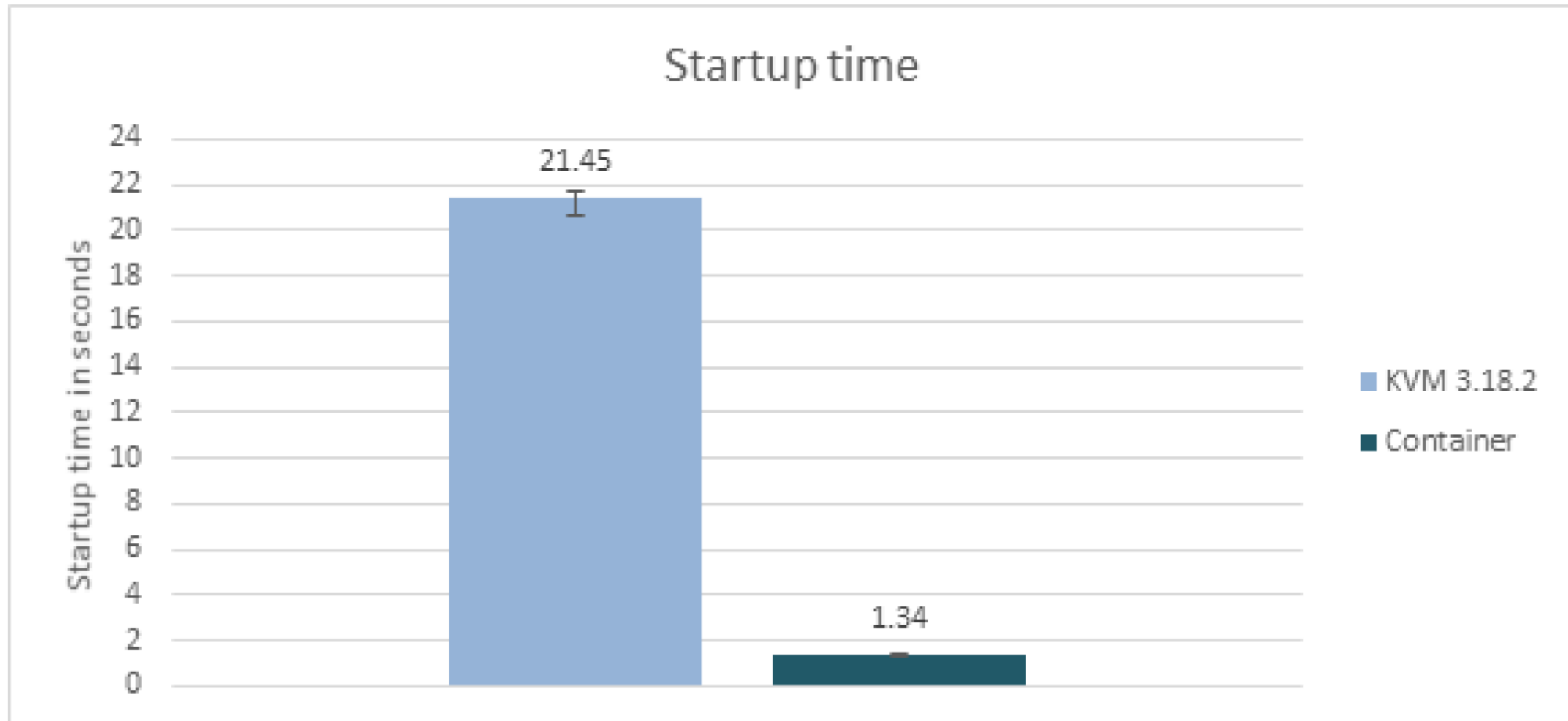
Massive complexity

High energy consumption

**DCs & networks:**
**Up to 18% of global electrical power by 2030**

A. Anders, T. Edler, "On global electricity usage of communication technology: trends to 2030.", Challenges 6, no. 1 (2015): 117-157

# Managing complexity:

## Containers and Serverless computing

UMEÅ UNIVERSITY

**Average Startup Time (Seconds) for a KVM Linux Virtual Machine and a Container Over Five Measurements**

# VIRTUALISED NETWORK INFRASTRUCTURE

We have an infrastructure to instantiate, run, and manage containers almost instantly

**Why not virtualise network functions?**

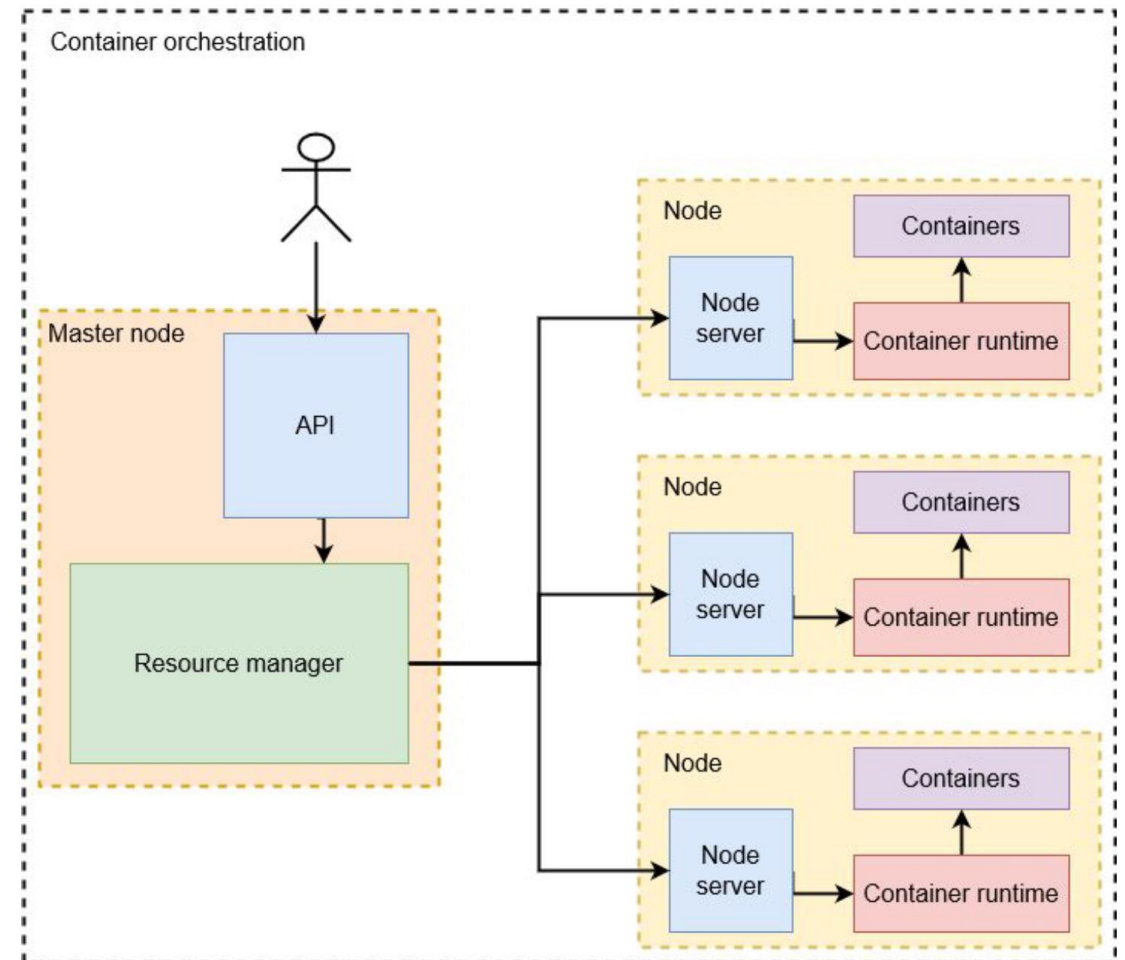Network address translation
Firewalls
Routing services
Etc.

Run and manage the network on cheap commodity servers

**The challenge then becomes:**

running and managing the network on cheap commodity servers

UMEÅ UNIVERSITY

**Container orchestrators are crucial for deploying, managing, and monitoring container systems**

Container engines deploy container images, running container runtimes

Container orchestrators manage the runtimes and the live system as a whole

**Azure Kubernetes Service (AKS)**



IBM Cloud
Kubernetes Service



Google Kubernetes Engine



Amazon ECS

# KUBERNETES

By far the leading container orchestration platform in the world
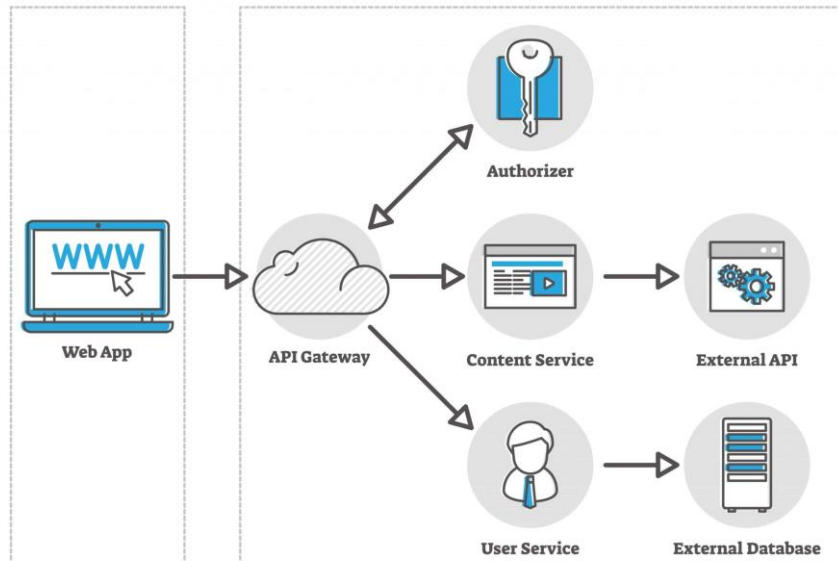
Portable

Extensible

Open-Source

Huge ecosystem

KubeCon 2022

17,000 attendees

3,000 companies

# CAN SERVERLESS HELP?

SERVERLESS

Web App → API Gateway → Authorizer / Content Service → External API / User Service → External Database

**Users / developers:**

Application functionality composed of invokable services (FaaS)

**Cloud providers:**

Auto-provision, deploy, and scale the services based on range of criteria

Abstract infrastructure from applications

**How does this work in a Cloud-Edge continuum, given the characteristics discussed?**

# Serverless at the Edge: Research challenges

# SERVERLESS CHALLENGES IN THE CONTINUUM

**Cognitive Cloud-Edge management:**
disaggregated, highly distributed, SDNs, etc.

**Deployment and migration of functions:**
Overlays to abstract heterogeneity etc.

**Seamless user access and programming models:**
automatic offload with specific requirements

**Secure and trusted execution:**
Fine-grained control in multi-provider networks

**New serverless functionalities:**
Long-term execution, location-awareness, etc.

**Energy efficiency and adaptation to green energy:**
Placement, renewables, models, etc.

**AI models and optimisation techniques:**
How to accurately manage and predict without incurring excessive latency or overhead

**50 billion connected devices by 2025**
This would overwhelm centralised DCs

**Proliferation of "intermediate" edge DCs**

**Complex federations of devices**

**Regional considerations**

# MODELLING THE CONTINUUM

A Simple Generalized Stochastic Petri Net Model

**Conceptual models:**

Stochastic process algebra, Discrete event simulation, Queueing theory, Approximation theory, Game theory, etc.

**Need of modelling solutions:**

Understanding the behavior, performance, and resource orchestration in cloud-edge systems.

Few formal models of federated Cloud-Edge systems exist

**None adequately represent and integrate energy and network considerations**

# MODELLING CHALLENGES

**How to model the system**

Stochastic features, network, pricing, energy distribution, policies, etc.

**How to maintain energy-perf trade-offs**

Best practices, Conflicting SLOs, local & global optima, etc.

**How to combine multiple models**

Model types, granularities, scaling, interactions

**How to model energy-driven systems**

Power grids capacities, renewables, service levels

**How to model different system regions**

Optimization, heterogeneity, monitoring, scheduling

**How to develop validation models**

Scalability scenarios, Iterative testing, mobility, topology, network behavior, energy, etc.

# SovereignEdge.COGNIT

**Modelling and Managing the Future Cloud-Edge Continuum** - Paul Townend

# PROJECT PARTICIPANTS

## Umea University, Sweden

**Paul Townend**

Monowar Bhuyan

P-O Ostberg

Erik Elmroth

## Ikerlan, Spain

Idoia de la Iglesia

Marco González

Iván Valdés

Aritz Brosa

Martxel Lasa

Goiuri Peralta

Samuel Pérez

## Open Nebula, Spain

**Alberto P. Martí**

Constantino Vázquez

Marco Mancini

Ignacio M. Llorente

Michael Abdou

## ACISA, Spain

Joan Iglesias

Antonio Lalaguna

Behnam Ojaghi

## SUSE, Germany

Torsten Hallmann

Holger Pfister

## Nature 4.0, Italy

Riccardo Valentini

Francesco Renzi

Micaela Onorati

## RISE, Sweden

Thomas Timoudas

Daniel Olsson

Johan Kristiansson

Shuai Zhu

## Atende, Poland

Dominik Bocheński

Tomasz Piasecki

Grzegorz Gil

## CETIC, Belgium

Nikolaos Matskanis

Philippe Massonet

Sébastien Dupont

Malik Bouhou

## Phoenix Systems, Poland

Kaja Swat

Tomasz Korniluk

Marek Białowąs

Gerard Świderski

Rafał Jurkiewicz

**10+ organisations
40 researchers**

# COGNIT architecture and use cases

**Stateless component responsible for managing the life cycle of the Serverless Runtimes**

**AI reasoning. First version will look at workload placement based on renewable power source availability**

UMEÅ UNIVERSITY

**1**

**Smart Cities**

Coordinated by ACISA

**2**

**Wildfire Detection**

Coordinated by Nature 4.0

**3**

**Energy**

Coordinated by Phoenix Systems & Atende Industries

**4**

**Cybersecurity**

Coordinated by CETIC and SUSE

UMEÅ UNIVERSITY

Commercial and Research Data Center

DCD Best Data Center Initiative 2017

IEEE
Scale Award 2017

Building 2000+ Node Container Facility

**RISE SICS, NORTHERN SWEDEN**

COGNIT testbed front-end deployment configuration

# COGNIT

**An open source reference implementation for serverless Continuum computing**

First version released in **September 2023**, more advanced version in **March 2024**

All versions of COGNIT will be tested in physical + virtual testbeds and use cases

**Happy to collaborate, incorporate interesting new technologies, etc.**

# Going forward

**Modelling and Managing the Future Cloud-Edge Continuum** - Paul Townend

**Energy-aware
Autonomous management**
Intelligently allocate resources:
e.g. target renewable energy, etc.

**Horizon Europe** SovereignEdge
(2023-2025)

**Efficiently monitor, predict, and
audit at massive scale**
How to adaptively do this to avoid huge overhead?

**WASP** WARA-Ops
(2023-2025)

**A formal model for energy-aware Cloud-Edge Systems**
Integrate energy providers, pricing, renewables, etc.
into existing Cloud-Edge models

**WASP** Academic PhD
(2023-2026)

# WARA-OPS

How do we model and integrate energy and network into Cloud-Edge?

How do we monitor at massive-scale without being overwhelmed with data?

How do we deal with conflicting demands between DCs and energy providers?

How do we optimise and negotiate in near real-time?

How do we store so much information for later audit?

Where does ML and other AI fit into this?

UMEÅ UNIVERSITY

What are the key components we need to introduce into our models?

How to integrate network models into our existing Cloud-Edge models?

How to integrate (lightweight) AI/ML for 6G?
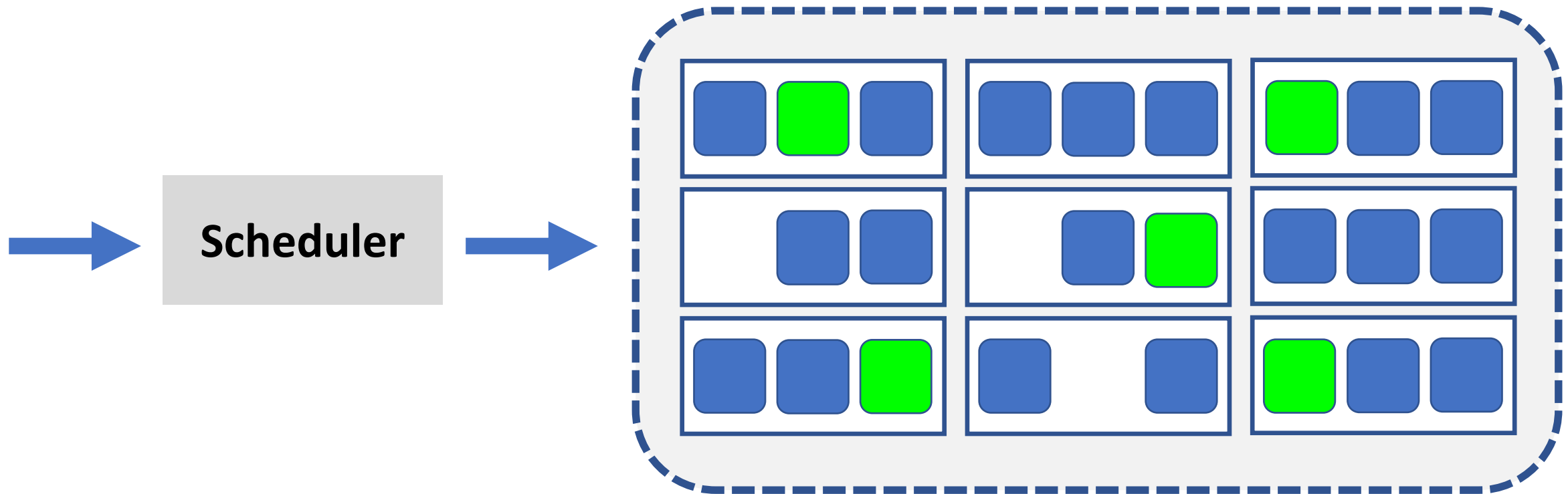
Rapid anomaly detection

Behavioural prediction

What orchestration technology is appropriate?

# An example of container orchestration for energy

Virtualise resources – and **schedule workloads** in a more effective manner

# WHAT CAN WE DO WITH SMARTER SCHEDULING?

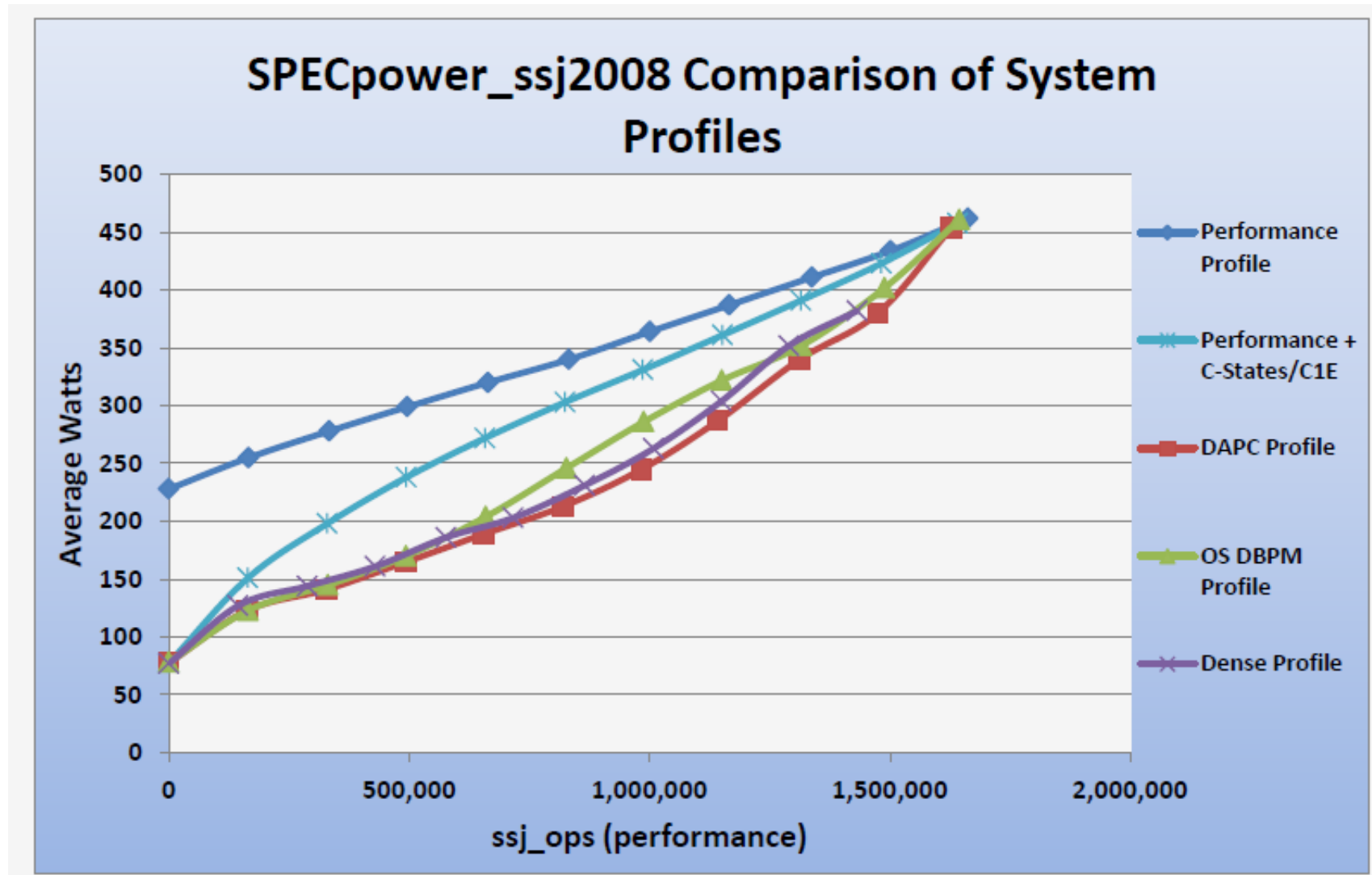| Over-allocation | Interference | Optimise hardware |
|---|---|---|
| Allocate more work on the same nodes | Avoid contention between co-located workloads | Allocate work until nodes are at "optimum" efficiency |
| Use less machines | Reduce power, improve performance | Reduce power, improve performance |

Commercial and Research Data Center

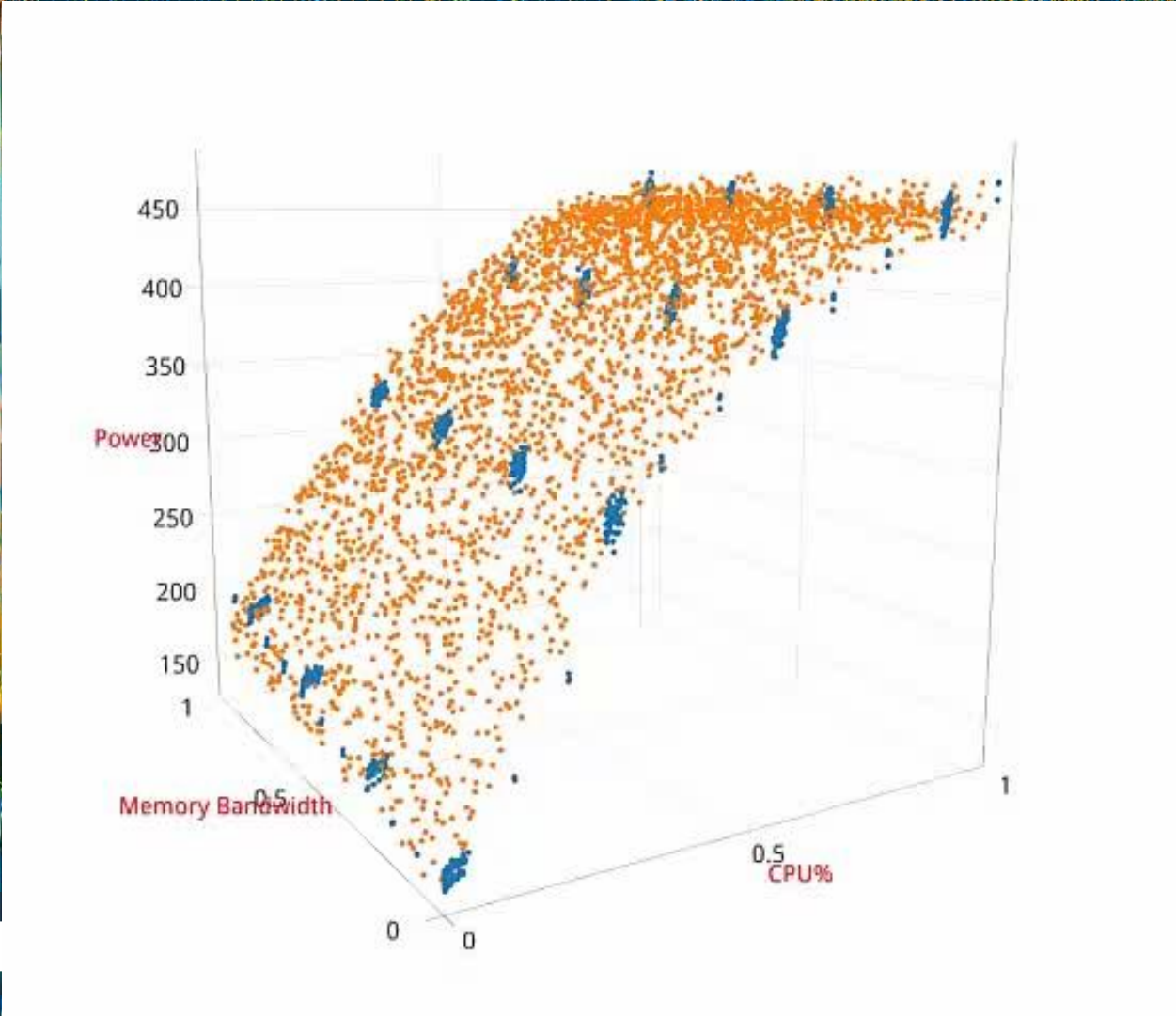DCD Best Data Center Initiative 2017

IEEE
Scale Award 2017

Building 2000+ Node Container Facility
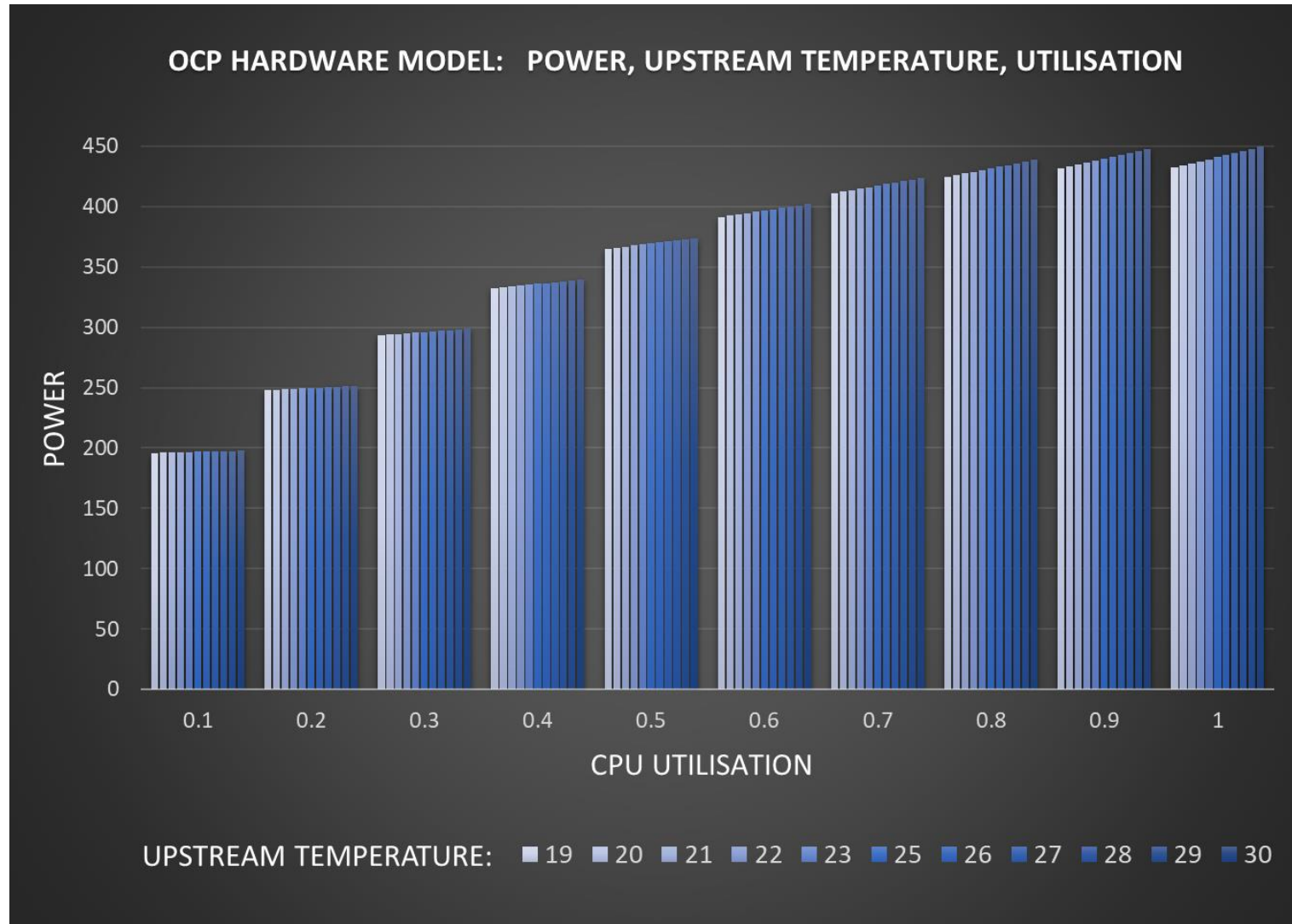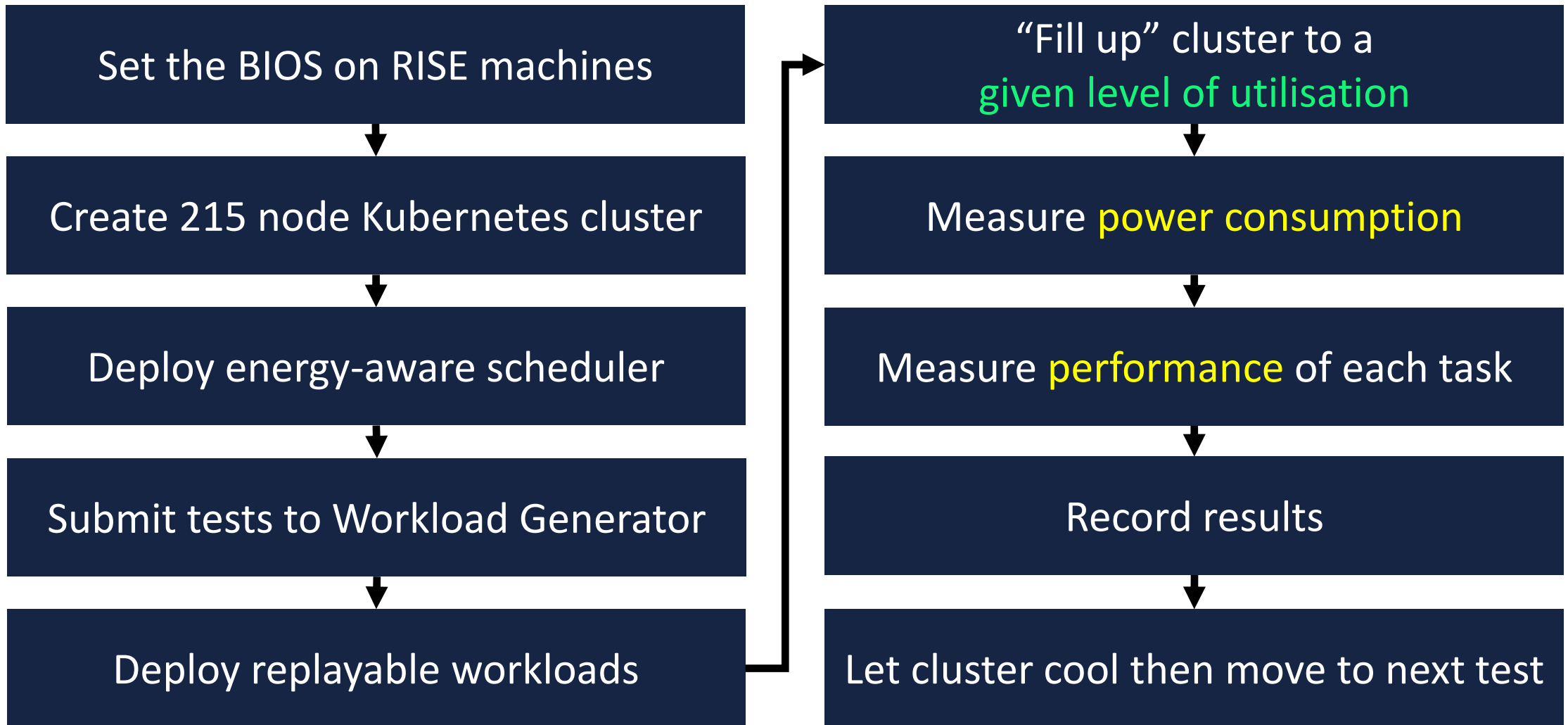
## RISE SICS, NORTHERN SWEDEN

# WIND TUNNEL BASED SERVER MODELLING

Actual readings (in blue)
Predicted values (orange)

OCP HARDWARE MODEL: POWER, UPSTREAM TEMPERATURE, UTILISATION

UMEÅ UNIVERSITY

Set the BIOS on RISE machines

↓

Create 215 node Kubernetes cluster

↓

Deploy energy-aware scheduler

↓

Submit tests to Workload Generator

↓

Deploy replayable workloads

"Fill up" cluster to a given level of utilisation

↓

Measure power consumption

↓

Measure performance of each task

↓

Record results

↓

Let cluster cool then move to next test

Benchmark Job Placement

UMEÅ UNIVERSITY



Power consumption during identical workloads
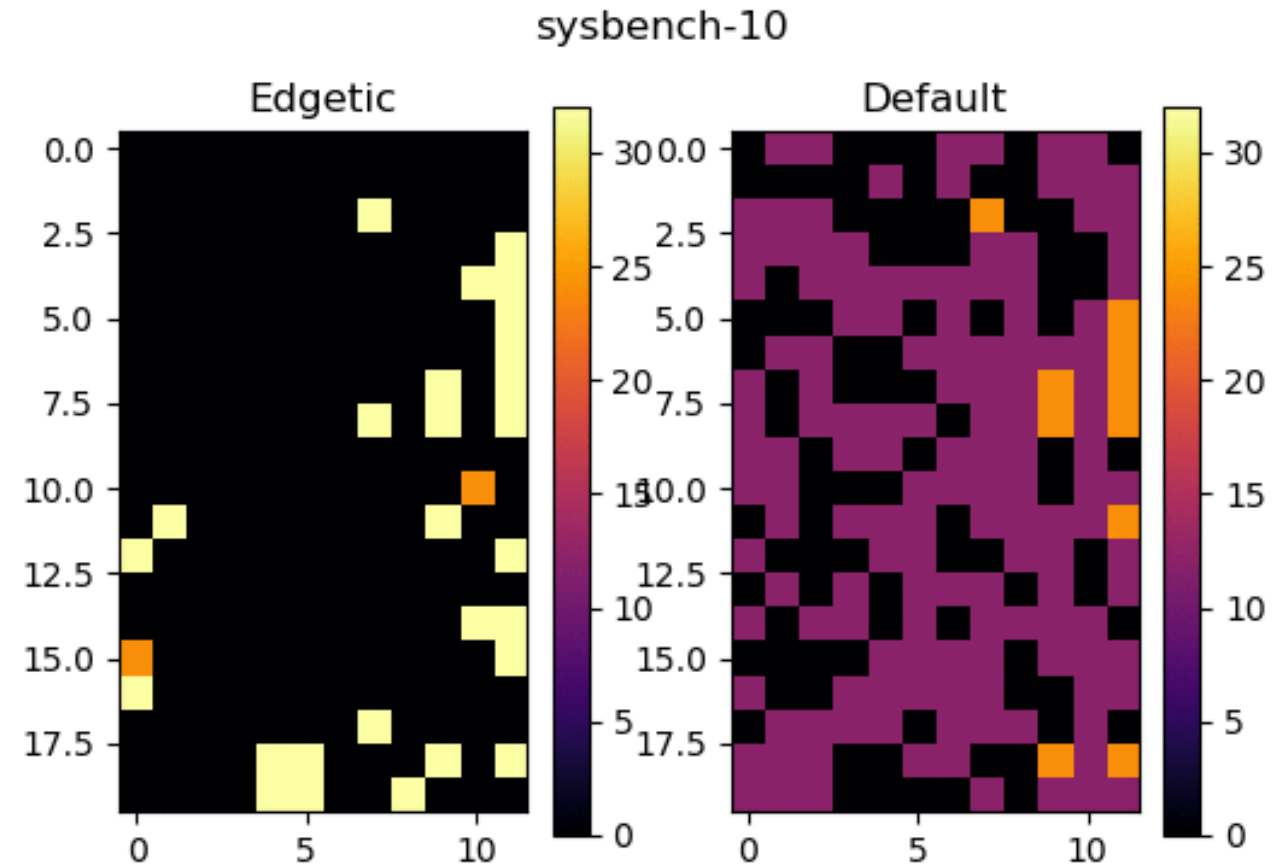
**215 nodes**

**OCP Hardware**

**Variety of workloads**

**Kubernetes containers**

**No prior workload knowledge**

**10-20% power savings**

**12 terabytes of telemetry data**



sysbench-10

**Overhead of Scheduler:  10ms per incoming workload**

cpi/watt v.s app. performance

CPU 1

CPU 2

Energy / perf
sweet spot

LATENCY!

norm_eps
norm_latency
norm_cpi_per_watt

Time(timestamp)

+1.5669e9

# paul.townend @ umu.se