# CLOUD COMPUTING

## Intro to Cloud: Background, Terminology, Basics, Data Centers

**PAUL TOWNEND**

ASSOCIATE PROFESSOR, UMEÅ

# A LITTLE BIT ABOUT ME

Originally from Leeds, UK. Joined Umeå in August 2020

PhD in Fault-Tolerant Grid Systems (2006)

General Chair of 13 IEEE conferences, TPC member of 28 IEEE conferences

Co-founded and worked for 3 years as CTO of spin-out company

1500+ citations, 70+ peer-reviewed publications in Distributed Systems

Member of the WASP Graduate School Management board

Examiner for this course and Learning Feature Representations

# INDUSTRY COLLABORATIONS

# EXTERNAL / INTERNATIONAL ACTIVITY

| | |
|---|---|
| **Leader** | University of Leeds Distributed Systems and Services group (2004-2017)<br>UK-China COLAB collaboration (Leeds, Beihang, Chongqing, NUDT) (2006 – present) |
| **PhD examiner** | University of Newcastle, UK (2017-2018) |
| **Editor** | IEEE Transactions on Services Computing 12(1) (*Guest editor, 2019*)    Philosophical Transactions of the Royal Society A: Cloud Computing (*January 2013*)<br>ACM Transactions on Embedded Systems (*Late 2015*)    IEEE Distributed Systems Online Journal, Dependability Topic (*2004-2007*) |
| **Steering committee** | IEEE ISORC    IEEE SOSE    IEEE JCC |
| **General Chair** | IEEE BigDataService 2021 (Online)    IEEE SOSE 2015 (San Francisco, USA)    IEEE SOSE 2016 (Oxford, UK)<br>IEEE ISORC 2018 (Singapore)    IEEE RTCPS 2014 (Reno, USA)    IEEE ISORC 2015 (Auckland, NZ)    IEEE iVCE 2013 (San Francisco, USA)<br>IEEE ISORC 2016 (York, UK)    WODSOG 2006 (Leeds, UK)    GCC 2014 (London, UK) |
| **PC Chair** | IEEE BigDataService 2020 (Oxford, UK)    IEEE SOSE 2014 (Oxford, UK)<br>IEEE JCC 2020 (Oxford, UK)    IEEE iVCE 2012 (Shenzhen, China)    UK e-Science AHM 2011 (York, UK) |
| **Workshop Chair** | IEEE ISORC 2014 (Reno, USA) |
| **Publicity Chair** | IEEE ISORC 2013 (Paderborn, Germany)    IEEE SRDS 2007 (Beijing, China) |
| **Local Chair** | RNEC 2008 (Leeds, UK)    IEEE SRDS 2006 (Leeds, UK) |
| **PC member** | IEEE SOSE 2019 (San Francisco, USA)    IEEE EUC 2014 (Milan, Italy)    IEEE SmartIoT 2018 (Xian, China)    IEEE SEUS 2014 (Reno, USA)<br>IEEE SOSE 2018 (Bamberg, Germany)    IEEE ICDCS 2013 (Philadelphia, USA)    IEEE SOSE 2017 (San Francisco, USA)    IEEE SOSE 2013 (San Francisco, USA)<br>IEEE BDCAT 2018 (Zurich, Switzerland)    IEEE CSE 2013 (Sydney, Australia)    IEEE iThings 2015 (Sydney, Australia)    IEEE ISORC 2012 (Shenzhen, China)<br>IEEE BDSEA 2016 (Shanghai, China)    IEEE DSN 2012 (Boston, USA)    IEEE CCBD 2015 (Taipei, Taiwan)    IEEE SOSE 2011 (Los Angeles, USA)<br>IEEE iVCE 2015 (San Francisco, USA)    IEEE HASE 2008 (Nanjing, China)    IMTIC 2015 (Jamshoro, Pakistan)    SDMCMM 2012 (Montreal, Canada)<br>IEEE DSAA 2015 (Paris, France)    IEEE ICEBE 2019 (Shanghai, China)    IEEE iVCE 2014 (Oxford, UK)    IEEE CCBD 2015 (Taipei, Taiwan) |

# ABOUT THIS MODULE

Fairly lecture heavy , but hopefully we can have discussion along the way

Light assignment – your major work should be your PhD

Invited keynotes from Ericsson Research and Google USA

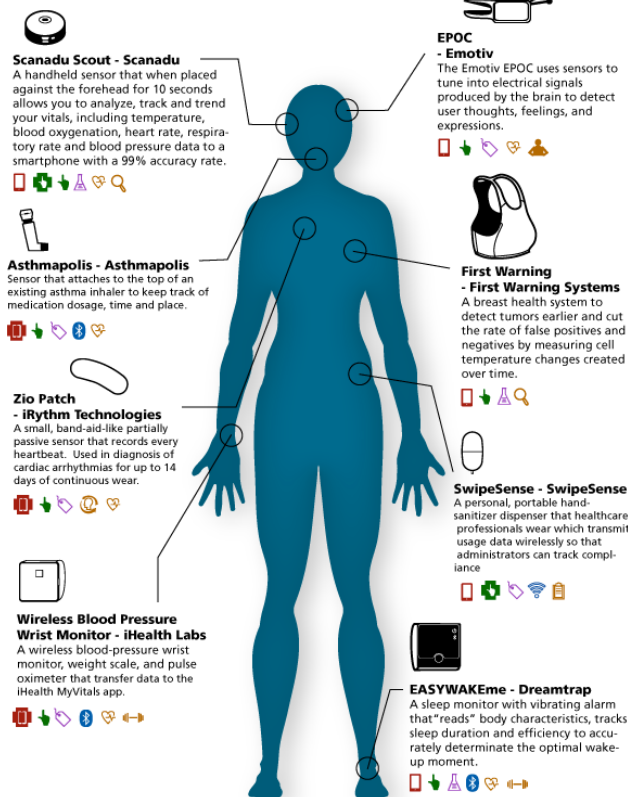A practical session to introduce the Ericsson Research Data Center and OpenStack

Ultimate goal: learn at least one thing that can inform/be applied to your PhD

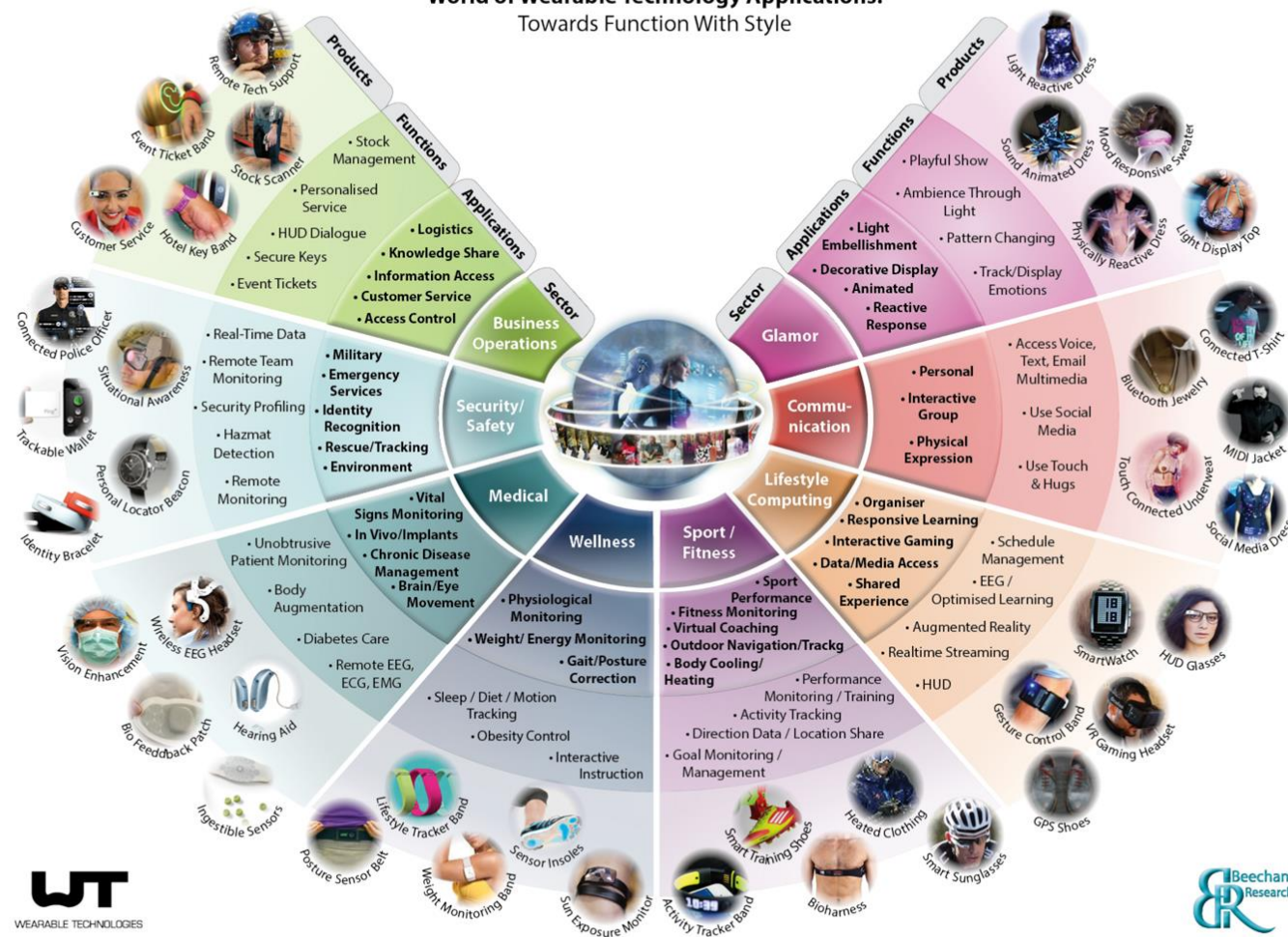WASP | WALLENBERG AI, AUTONOMOUS SYSTEMS AND SOFTWARE PROGRAM

# SCHEDULE

| DAY ONE | | May 4, 2021 |
|---|---|---|
| 09:30 – 09:45 | | Welcome to the Course |
| 09:45 – 11:15 | L1 | Introduction to Cloud: Background, terminology, basics, data centres |
| 11:15 – 11:45 | | BREAK |
| 11:45 – 12:45 | L2 | Distributed processing: Hadoop |
| 12:45 – 14:00 | | BREAK (LUNCH) |
| 14:00 – 15:00 | L3 | Distributed processing: Apache Spark and Apache Storm |
| 15:00 – 15:30 | | BREAK |
| 15:30 – 16:45 | L4 | Edge, Fog, and Serverless Computing |
| | | |

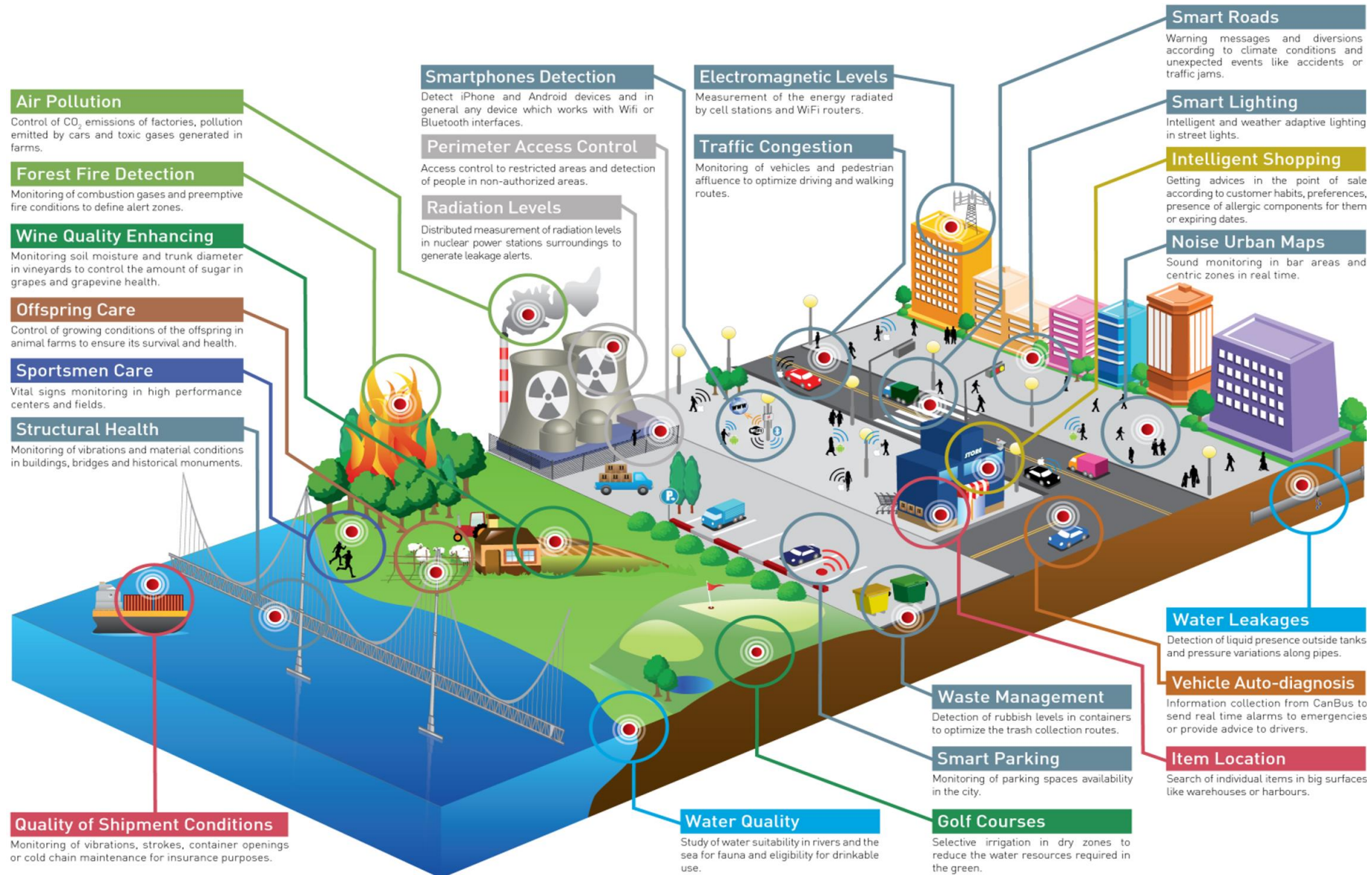| DAY TWO | | May 5, 2021 |
|---|---|---|
| 09:30 – 10:45 | L5 | Cloud Orchestration |
| 10:45 – 11:15 | | BREAK |
| 11:15 – 11:45 | K1 | Keynote: Johan Eker (Ericsson) |
| 11:45 – 12:30 | P1 | Cloud Practical Session (ER DC) |
| 12:30 – 14:00 | | BREAK (LUNCH) |
| 14:00 – 15:00 | P2 | Cloud Practical Session (ER DC) |
| 15:00 – 16:00 | L6 | Cloud Economics |
| 16:00 – 16:30 | | BREAK |
| 16:30 – 17:30 | K2 | Industry Keynote: "Challenges & Opportunities in Cloud" - Steve Webster, Google |
| 17:30 – 17:45 | | Closing remarks and assignment information |

# Why study Cloud?

## KEY

**Device Type:** Medical | Consumer

**User:** Healthcare Provider | Consumer

**Wireless Capability:** Bluetooth | Wifi

**Availability:** On the Market | In Development

**Used for:** Fitness/Wellness | Chronic Condition/Disease Management | Early Detection | Continuous Monitoring | Adherence | Rehab

**Scanadu Scout - Scanadu**
A handheld sensor that when placed against the forehead for 10 seconds allows you to analyze, track and trend your vitals, including temperature, blood oxygenation, heart rate, respiratory rate and blood pressure data to a smartphone with a 99% accuracy rate.

**Asthmapolis - Asthmapolis**
Sensor that attaches to the top of an existing asthma inhaler to keep track of medication dosage, time and place.

**Zio Patch - iRythm Technologies**
A small, band-aid-like partially passive sensor that records every heartbeat. Used in diagnosis of cardiac arrhythmias for up to 14 days of continuous wear.

**Wireless Blood Pressure Wrist Monitor - iHealth Labs**
A wireless blood-pressure wrist monitor, weight scale, and pulse oximeter that transfer data to the iHealth MyVitals app.

**EPOC - Emotiv**
The Emotiv EPOC uses sensors to tune into electrical signals produced by the brain to detect user thoughts, feelings, and expressions.

**First Warning - First Warning Systems**
A breast health system to detect tumors earlier and cut the rate of false positives and negatives by measuring cell temperature changes created over time.

**SwipeSense - SwipeSense**
A personal, portable hand-sanitizer dispenser that healthcare professionals wear which transmits usage data wirelessly so that administrators can track compliance

**EASYWAKEme - Dreamtrap**
A sleep monitor with vibrating alarm that "reads" body characteristics, tracks sleep duration and efficiency to accurately determinate the optimal wake-up moment.

## World of Wearable Technology Applications: Towards Function With Style



© 2014 Beecham Research Ltd. & Wearable Technologies AG

**Air Pollution**
Control of $CO_2$ emissions of factories, pollution emitted by cars and toxic gases generated in farms.

**Forest Fire Detection**
Monitoring of combustion gases and preemptive fire conditions to define alert zones.

**Wine Quality Enhancing**
Monitoring soil moisture and trunk diameter in vineyards to control the amount of sugar in grapes and grapevine health.

**Offspring Care**
Control of growing conditions of the offspring in animal farms to ensure its survival and health.

**Sportsmen Care**
Vital signs monitoring in high performance centers and fields.

**Structural Health**
Monitoring of vibrations and material conditions in buildings, bridges and historical monuments.

**Quality of Shipment Conditions**
Monitoring of vibrations, strokes, container openings or cold chain maintenance for insurance purposes.

**Smartphones Detection**
Detect iPhone and Android devices and in general any device which works with Wifi or Bluetooth interfaces.

**Perimeter Access Control**
Access control to restricted areas and detection of people in non-authorized areas.

**Radiation Levels**
Distributed measurement of radiation levels in nuclear power stations surroundings to generate leakage alerts.

**Electromagnetic Levels**
Measurement of the energy radiated by cell stations and WiFi routers.

**Traffic Congestion**
Monitoring of vehicles and pedestrian affluence to optimize driving and walking routes.

**Smart Roads**
Warning messages and diversions according to climate conditions and unexpected events like accidents or traffic jams.

**Smart Lighting**
Intelligent and weather adaptive lighting in street lights.

**Intelligent Shopping**
Getting advices in the point of sale according to customer habits, preferences, presence of allergic components for them or expiring dates.

**Noise Urban Maps**
Sound monitoring in bar areas and centric zones in real time.

**Water Leakages**
Detection of liquid presence outside tanks and pressure variations along pipes.

**Vehicle Auto-diagnosis**
Information collection from CanBus to send real time alarms to emergencies or provide advice to drivers.

**Item Location**
Search of individual items in big surfaces like warehouses or harbours.

**Waste Management**
Detection of rubbish levels in containers to optimize the trash collection routes.

**Smart Parking**
Monitoring of parking spaces availability in the city.

**Water Quality**
Study of water suitability in rivers and the sea for fauna and eligibility for drinkable use.

**Golf Courses**
Selective irrigation in dry zones to reduce the water resources required in the green.

**WASP** | WALLENBERG AI, AUTONOMOUS SYSTEMS AND SOFTWARE PROGRAM

# TRENDS IN MODERN DISTRIBUTED SYSTEMS

Need good
security

Highly distributed
(often mobile)

Generate large amounts of
data

Latency can
be critical

Many small / low powered
devices

Need support for analytics
and decision making

# 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

# What is Cloud?

# NIST DEFINITION

A model for enabling ubiquitous, convenient, on-demand network access to a shared pool of configurable computing resources

(e.g., networks, servers, storage, applications, and services)

that can be rapidly provisioned and released with minimal management effort or service provider interaction.

You do not own computing infrastructure – you rent it

WASP | WALLENBERG AI, AUTONOMOUS SYSTEMS AND SOFTWARE PROGRAM

# MAJOR CLOUD CHARACTERISTICS

**On-demand self-service**

Can provision resources without human intervention

**Broad network access**

Service is accessible through the network

**Resource pooling**

Resources are pooled to several customers

**Rapid elasticity**

Can rapidly increase and decrease capacity

**Measured service**

Resource usage can be monitored and reported



WASP | WALLENBERG AI, AUTONOMOUS SYSTEMS AND SOFTWARE PROGRAM

# MAJOR CLOUD CHARACTERISTICS

**Virtualization**

Physical machines are split into "virtual" components

**Multi-tenancy**

Multiple jobs / users use the same physical machine

**Service model**

Multiple different ways Clouds can be used

**Pay-per-use**

Related to on-demand resources: users only pay for the resources that they actually use

# VIRTUALISATION



Cloud

Partition physical machines to maximise resource utilisation and elasticity

Subject to SLA

Resource management becomes critical

# VIRTUAL MACHINES AND CONTAINERS (1)

**VM 1**

App

Guest OS

**VM 2**

App

Guest OS

**VM 3**

App

Guest OS

Hypervisor

Host operating system (if Type 2 Hypervisor)

Physical server

**VMs are like neighbours**

**Container 1**

App

**Container 2**

App

**Container 3**

App

Container engine (Docker, etc.)

Host operating system

Physical server

**Containers are like housemates**

WASP | WALLENBERG AI, AUTONOMOUS SYSTEMS AND SOFTWARE PROGRAM

# VIRTUAL MACHINES AND CONTAINERS (2)

**Container 1**     **Container 2**     **Container 3**

| App | App | App |

Container engine (Docker, etc.)

Host operating system

Physical server

---

**VM 1**     **VM 2**     **VM 3**

| App | App | App |
| Guest OS | Guest OS | Guest OS |

Hypervisor

Host operating system (if Type 2 Hypervisor)

Physical server

---

"Light"

Fast to create

Do not have migration

Typically short lived and single function

---

"Heavy"

Slow to create

Slow to migrate

Typically long lived and multi-function

Global Search: Kubernetes vs OpenStack 2013-2020

— kubernetes: (Worldwide)    — openstack: (Worldwide)

# Cloud and Industry

# BENEFITS FROM AN INDUSTRY PERSPECTIVE

**Move from capital expenditure to operational expenditure**

Pay only for what is used (in some cases)

The **illusion** of having infinite resources

Offloads responsibility for security (physical and virtual) plus uptime, efficiency, maintenance, etc.

# CLOUD INDUSTRY DRIVERS

## Business Drivers for Implementing Cloud Solutions, Global 2019

| Driver | % |
|---|---|
| Increase app availability/uptime | 80% |
| Free up IT staff to focus on innovative solutions to business challenges | 80% |
| Deliver services and applications faster | 79% |
| Improve business continuity/disaster recovery | 78% |
| Eliminate hassle of integrating multi-vendor solutions | 78% |
| Reduce Costs | 77% |
| Support my company's Digital Transformation initiatives | 77% |
| "Go Green" – better manage environmental costs and impact | 77% |
| Reduce hardware/software maintenance | 76% |
| Position my company to take advantage of new technologies | 76% |
| Support business agility/market responsiveness | 75% |
| Shift costs from capital to operating budget | 74% |

*Reference: Succeed in the Digital Era with Cloud Communications, Frost & Sulllivan, 2019*

## How many of your applications come from the cloud?

Legend: ■ 2018 ■ 2020

| Range | 2018 | 2020 |
|---|---|---|
| <10 % | 74% | 17% |
| 10–29 % | 21% | 41% |
| 30–59 % | 3% | 26% |
| 60–89 % | 1% | 10% |
| >90 % | 1% | 6% |

Maturity Survey: 2018 Cloud Computing

Pie chart – Which of the following best describes your organization's use of container technology today?

- Not using containers and not interested at this time — 24%
- Researching container technology — 18%
- Experimenting with container technology — 17%
- Running containers in production — 16%
- Using containers for development & testing — 13%
- Don't know — 12%

**Q.** Which of the following best describes your organization's use of container technology today?

WASP | WALLENBERG AI, AUTONOMOUS SYSTEMS AND SOFTWARE PROGRAM

Big Data vs Cloud Computing

# EXAMPLES OF CLOUD PROVIDERS



| Hyperscalers | Smaller, often co-locational |

# PROBLEMS WITH USING THE CLOUD

| Customer perspective |
|---|
| Privacy and security: less control |
| Legal problems: who does the data belong to? etc |
| Designing scalable applications |
| Communication latency |
| Risk of lock-in: difficult to migrate to other providers |
| Ability to negotiate and manage the SLA |

| Provider perspective |
|---|
| SLAs: how much does it cost to respect them? |
| Cloud interoperability |
| Multi-tenancy management: interference |
| Energy efficiency management |
| Uncertainty and variability of service requests |

# Deployment and Service Models

# DEPLOYMENT MODELS (WHERE IS IT, WHO OWNS IT?)

| Private Cloud | Public Cloud | Hybrid Cloud |
|---|---|---|
| • Internal resources managed in a "Cloud-like" fashion<br><br>• Greater level of security and personalization<br><br>• Less scalability and higher costs<br><br><br>**(owned internally)** | • Resources rented to anyone who will pay<br><br><br><br><br><br>**(owned by someone else)** | • A combination of the above<br><br>• Typically links two or more cloud infrastructures (public or private) via a standard technology<br><br><br>**(owned by multiple orgs)** |

WASP | WALLENBERG AI, AUTONOMOUS SYSTEMS AND SOFTWARE PROGRAM

# SERVICE MODELS (WHAT DOES IT PROVIDE?)

**Software-as-a-Service (SaaS)**

Ready-to-use applications
(O365, Spotify, Dropbox, etc)

**Platform-as-a-Service (PaaS)**

Ready-to-use platform
(Windows, Google Apps Engine, etc)

**Infrastructure-as-a-Service (IaaS)**

Full access to a hosted machine / VM
(Amazon EC2, Windows Azure, etc)

| On-Premises | IaaS<br>Infrastructure as a Service | PaaS<br>Platform as a Service | SaaS<br>Software as a Service |
|---|---|---|---|
| Applications | Applications | Applications | Applications |
| Data | Data | Data | Data |
| Runtime | Runtime | Runtime | Runtime |
| Middleware | Middleware | Middleware | Middleware |
| O/S | O/S | O/S | O/S |
| Virtualization | Virtualization | Virtualization | Virtualization |
| Servers | Servers | Servers | Servers |
| Storage | Storage | Storage | Storage |
| Networking | Networking | Networking | Networking |

You Manage    Other Manages

WASP | WALLENBERG AI, AUTONOMOUS SYSTEMS AND SOFTWARE PROGRAM

# SERVICE MODELS: ON PREMISES ("ON PREM")

**Traditional approach - organisation owns and manages its own infrastructure**

(buy hardware, install / configure / run applications, etc.)

| **What is good about this?** | Total control of the infrastructure |
|---|---|

On-Premises

| Applications |
|---|
| Data |
| Runtime |
| Middleware |
| O/S |
| Virtualization |
| Servers |
| Storage |
| Networking |

**What is bad about this?**

**Capacity management**: how many machines?

**Hardware management**: buy, install, manage, upgrade

**Security**: updates, hardware decommission (hard drives)

**Networking management**: how to connect everything?

**Storage management**: redundancy, disaster recovery

# SERVICE MODELS: INFRASTRUCTURE-AS-A-SERVICE

**Service as a commodity (e.g. electricity)**

**Cloud providers own and manage the infrastructure, customers rent resources**

Provide hardware, VMs, storage, servers, monitoring, alerts, etc.

**What is good about this?**

No need to buy hardware (as a customer)

Elastic and pay-per-use

**What is bad about this?**

Lost control on the infrastructure

Possible higher cost in the long run (see: cloud economics)

Legal issues (how to manage data?)

Applications may require re-design for cloud deployment

## IaaS
Infrastructure as a Service

- Applications
- Data
- Runtime
- Middleware
- O/S
- Virtualization
- Servers
- Storage
- Networking

# SERVICE MODELS: INFRASTRUCTURE-AS-A-SERVICE

| Challenges |
| --- |
| How to deploy an application? |
| Manage availability and scalability, load balancing |
| Manage the operating system |
| Selection and configuration of the hardware |

**IaaS**
Infrastructure as a Service

- Applications
- Data
- Runtime
- Middleware
- O/S
- Virtualization
- Servers
- Storage
- Networking

# SERVICE MODELS: PLATFORM-AS-A-SERVICE

**"Rent the operating system"**

Customer uploads + controls applications, and can configure some of the environment

| | |
|---|---|
| **What is good about this?** | Infrastructure is managed by the cloud provider |
| **What is bad about this?** | Limitation in the deployment environment |
| **Challenges** | The customer is in charge of managing load balancing and networking |

PaaS
Platform as a Service

Applications

Data

Runtime

Middleware

O/S

Virtualization

Servers

Storage

Networking

# SERVICE MODELS: SOFTWARE-AS-A-SERVICE

**SaaS**
Software as a Service

Applications

Data

Runtime

Middleware

O/S

Virtualization

Servers

Storage

Networking

**Use an existing application that is provided on a Cloud infrastructure**

Same application is shared on different customer devices by a thin client

**What is good about this?**

Can use with no deployment or configuration decisions

**What is bad about this?**

Limited configuration decisions and control

WASP | WALLENBERG AI, AUTONOMOUS SYSTEMS AND SOFTWARE PROGRAM

# NEWER SERVICE MODELS

# SERVICE MODELS: CONTAINER-AS-A-SERVICE (CAAS)

Containers pack applications and dependencies into an image

Same isolation as virtualisation, lighter than VMs

Requires a container orchestrator
Manages the infrastructure and runs containerised applications
Manages the interactions between applications

**What is good about this?**

Declarative deployment (desired final state of infrastructure)
Avoids vendor lock-in (multi-cloud environment)
Multiple containers on a single machine

**Containers as a Service (CaaS)**

Function
Application
Runtime
Container
OS
Virtualisation
Hardware

# SERVICE MODELS: FUNCTION-AS-A-SERVICE (FAAS)

Serverless – a new trend.  Pay per request (no idle time)

Auto-scaling and availability provided out-of-the-box

**Computation is implemented as functions and execution is event driven**

Customers define functions

Users select functions and specify the events triggering them

Examples include AWS Lambda, Google Cloud Functions, etc.

**Function as a Service (FaaS)**

| Function |
| Application |
| Runtime |
| Container |
| OS |
| Virtualisation |
| Hardware |

WASP | WALLENBERG AI, AUTONOMOUS SYSTEMS AND SOFTWARE PROGRAM

# FaaS and Furious by Forrest Brazeal



The two tribes regarded each other suspiciously
in the glow of their brightly blazing production environments.

# Cloud Native Applications

# MAKING AN APPLICATION "CLOUD READY"

Almost any application can run in the Cloud, but they are **not** "cloud-ready"

A key advantage to cloud is the ability to **auto-scale**: an "elastic" application can scale with load

Traditional application: 1 VM per tier



Client

↓

Application

↓

Database

# AUTO-SCALING CLOUD APPLICATION ARCHITECTURES

# Physical Cloud Infrastructure

# (Data Centers)

# DATA CENTERS

A **data center** is a facility composed of networked computers and storage that businesses or other organizations use to organize, process, store and disseminate large amounts of data

# GOOGLE HAMINA DATA CENTRE, FINLAND

# DATA CENTERS AS SYSTEMS OF SYSTEMS

# HOT AISLE / COLD AISLE CONTAINMENT

Network switch

Server

Server

Server

Server

Server

Server

Server

Server

Server

Server

RACK

Hot Aisle

21 Billion SEK facility
(twenty one)

Server buildings

Cooling

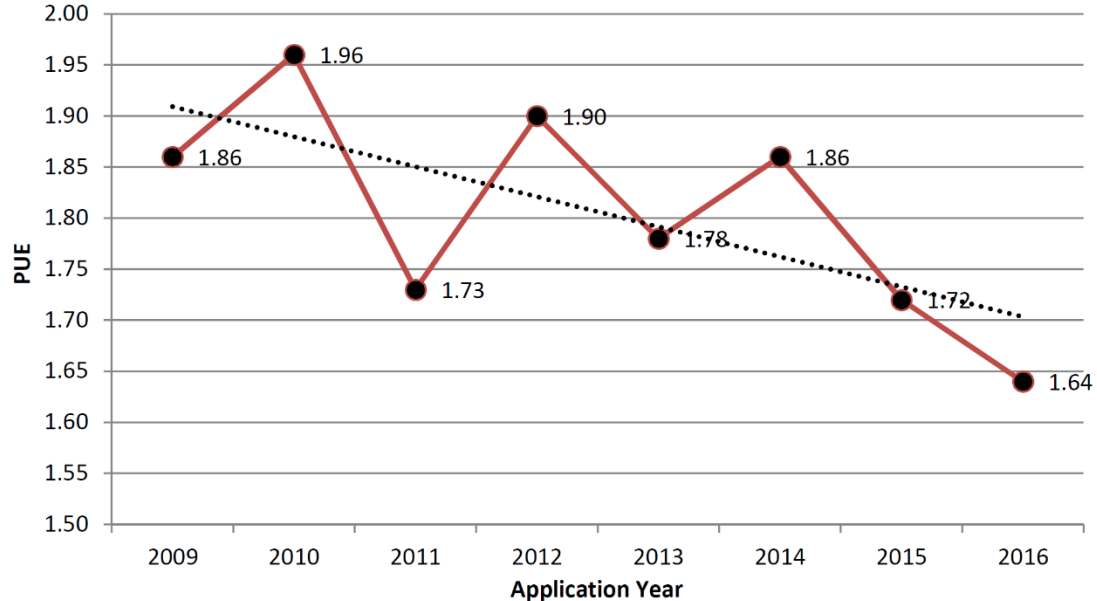Power facility

GOOGLE DATA CENTRE, CHARLESTON, USA

WASP | WALLENBERG AI, AUTONOMOUS SYSTEMS AND SOFTWARE PROGRAM

# PREDICTED GLOBAL POWER CONSUMPTION



**Typical DC utilisation is between 10-20%**

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

# POWER USAGE EFFECTIVENESS

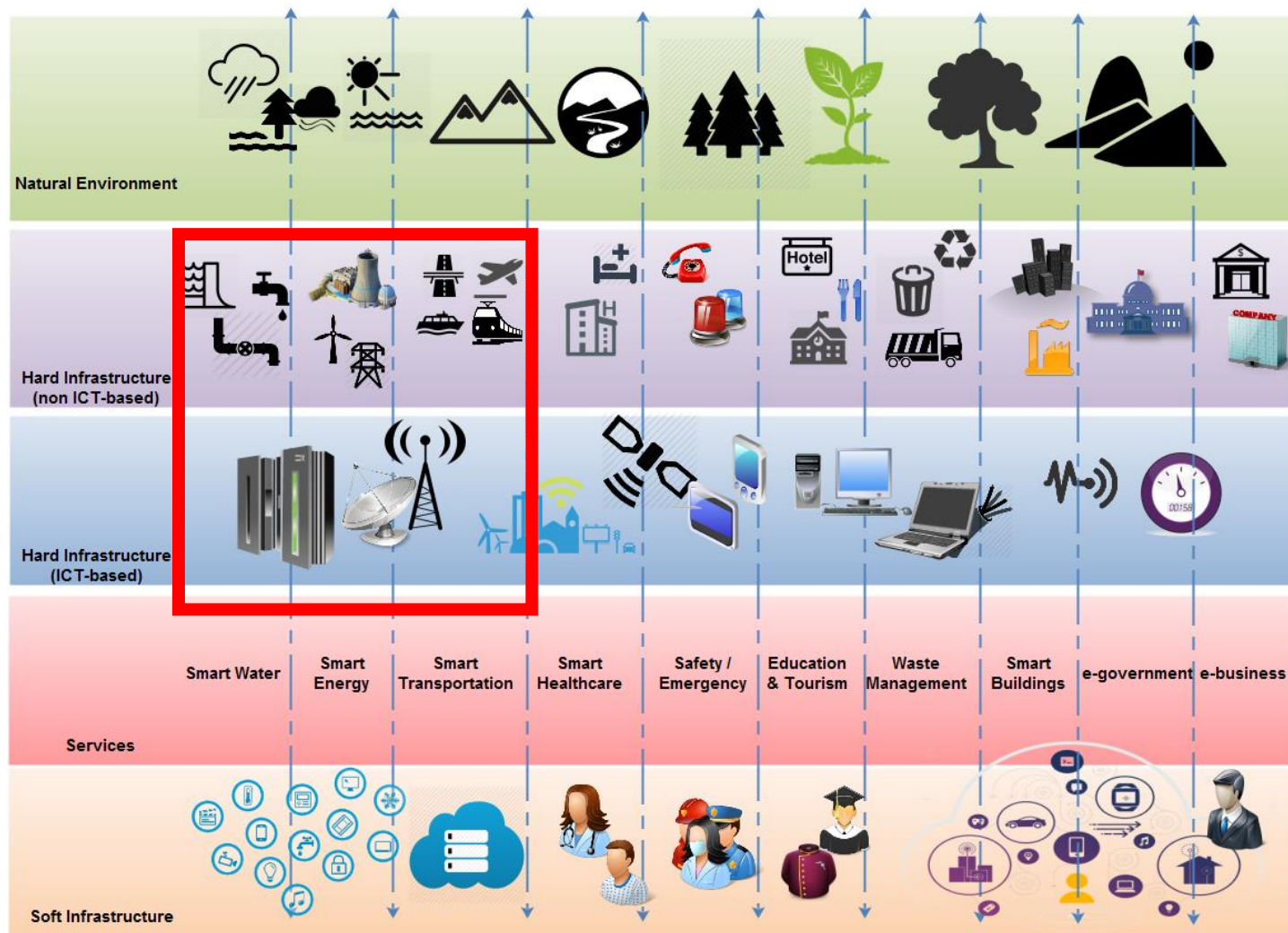$$PUE = \frac{Total\ Facility\ Power}{Total\ IT\ Power}$$



PUE is the data centre industry standard metric for efficiency

What are some of the problems with PUE?

As PUE approaches 1, why does IT power become increasingly important?

# SMART INFRASTRUCTURE STACK: IMPLICATIONS

# SOFTWARE BASED EFFICIENCY IMPROVEMENTS

**Schedule virtual workloads** in a more effective manner

# WHAT CAN WE DO WITH 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 |

# TYPES OF DATA CENTER

**Co-locational**
Require customers to supply their own hardware

**Telecoms**
Typically feature high connectivity requirements and run specialised software services

**Dedicated hosting**
Provide server capacity to single customers with no sharing of machines (bare metal hosting)

**Managed hosting**
Provide servers and storage systems for customers, often as PaaS, IaaS or SaaS

**Shared hosting**
Cloud data centres are an example. Provide virtualised multi-tenant resources.

**Edge**
Typically smaller than traditional data centres, and located closer to where data is generated

**Hybrid**
A data centre facility with more than one of the above service models

# DATA CENTER TIERS

**Tiers classify the structure of a data centre**

### Tier I

- Single power supply system and single cooling system
- No backup policies or redundant components

### Tier II

- Single power supply system and single cooling system
- Some components are redundant and backup policies
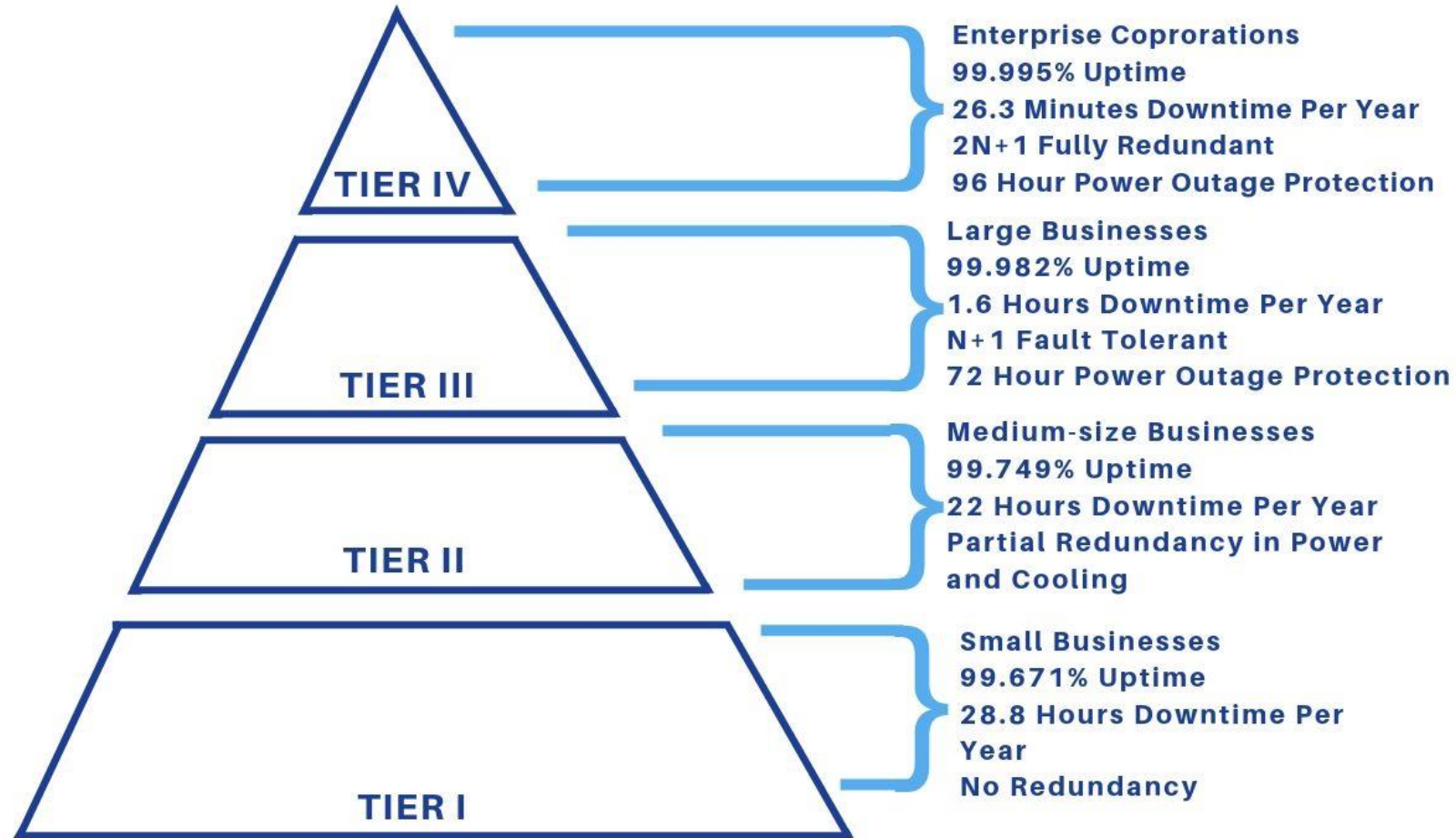- Minimum uptime must be 99.749% annually

### Tier III

- Multiple power supply systems + multiple cooling systems
- All components are redundant
- Most upgrades do not take centre offline
- Minimum uptime must be 99.982% annually

### Tier IV

- Designed and created to be totally fault-tolerant
- All components redundant
- Has more power and cooling systems
- Guaranteed uptime is 99.995%

FACEBOOK DATA CENTRE, LULEÅ, SWEDEN

WASP | WALLENBERG AI, AUTONOMOUS SYSTEMS AND SOFTWARE PROGRAM

# IN SUMMARY

NIST definition and characteristics

Virtual machines and containers

Industry uptake of Cloud and economics

Deployment and service models

What are data centres?

Data centre types and tiers

WASP | WALLENBERG AI, AUTONOMOUS SYSTEMS AND SOFTWARE PROGRAM