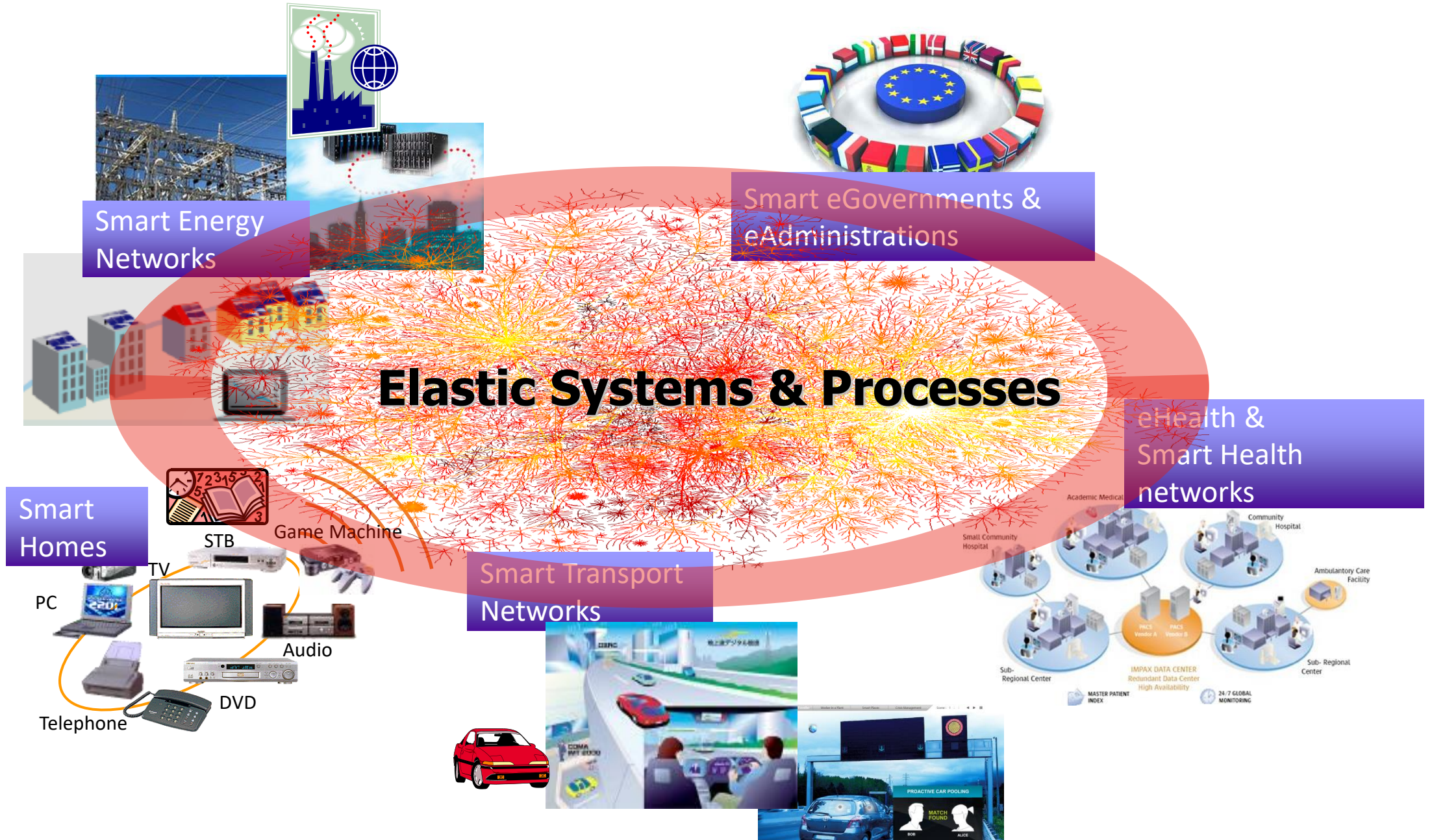# Edge Intelligence

## The Co-evolution of Humans, IoT, and AI
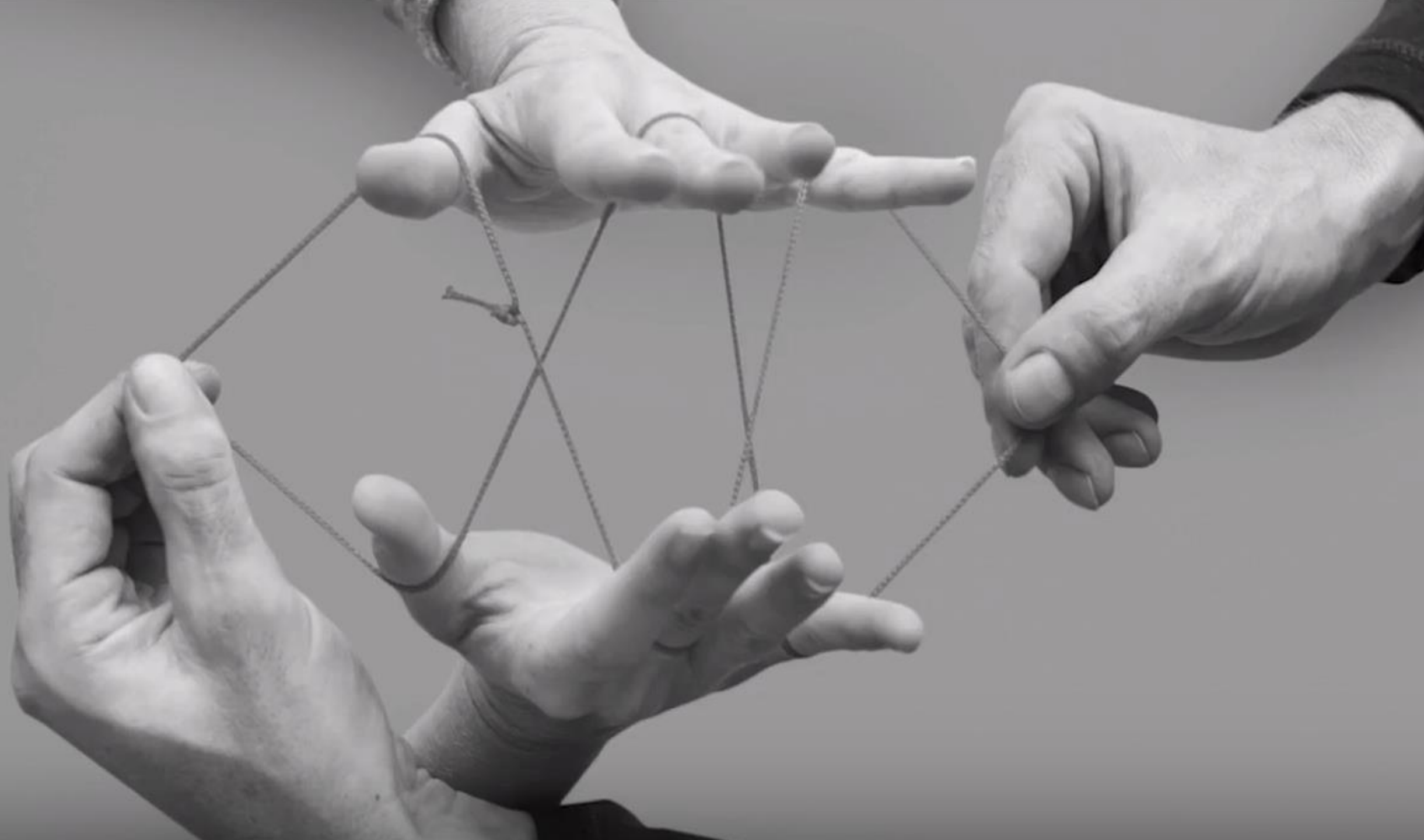
8 May 2020, IoTBDS 2020

Schahram Dustdar
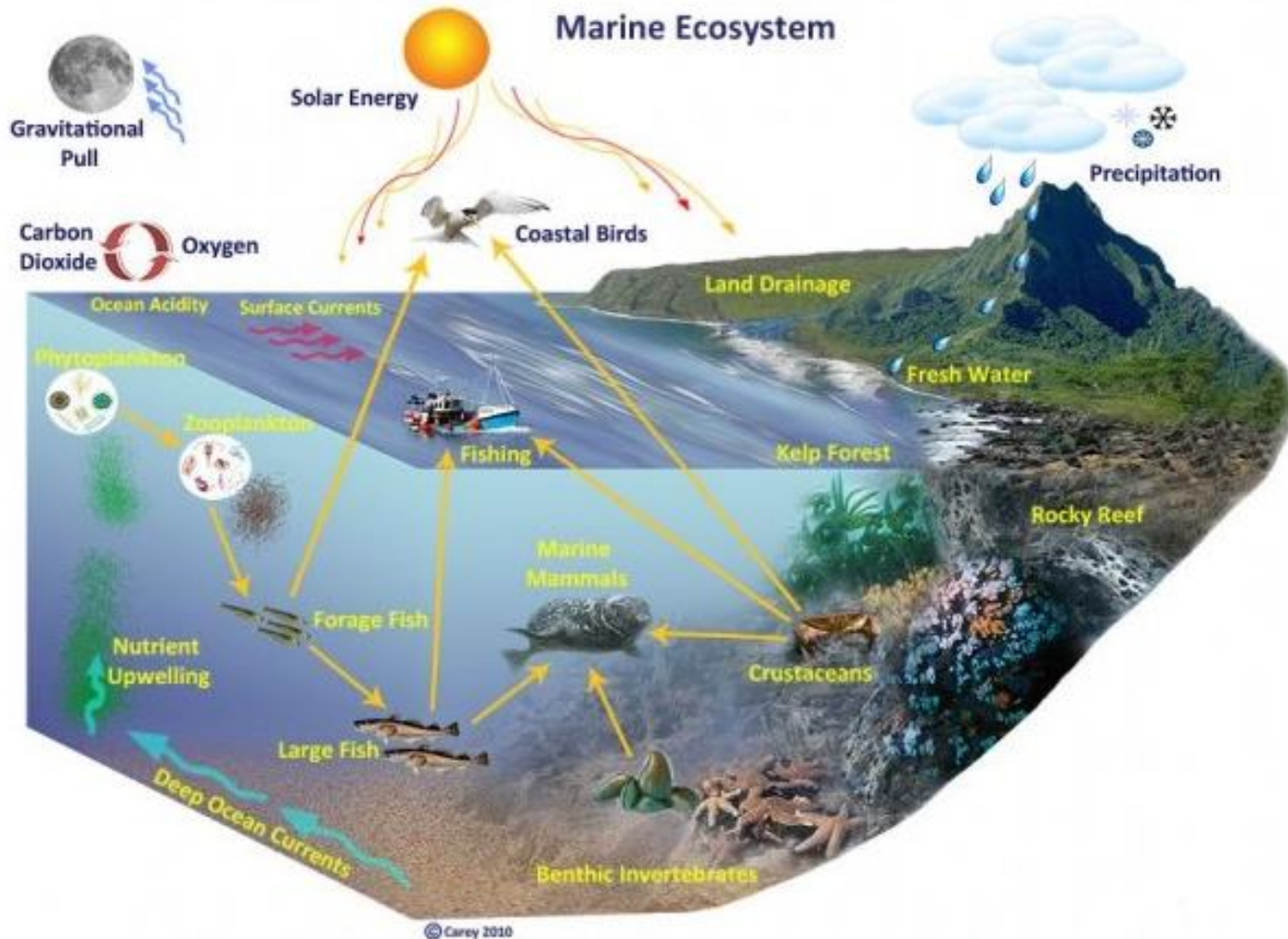
**dsg.tuwien.ac.at**

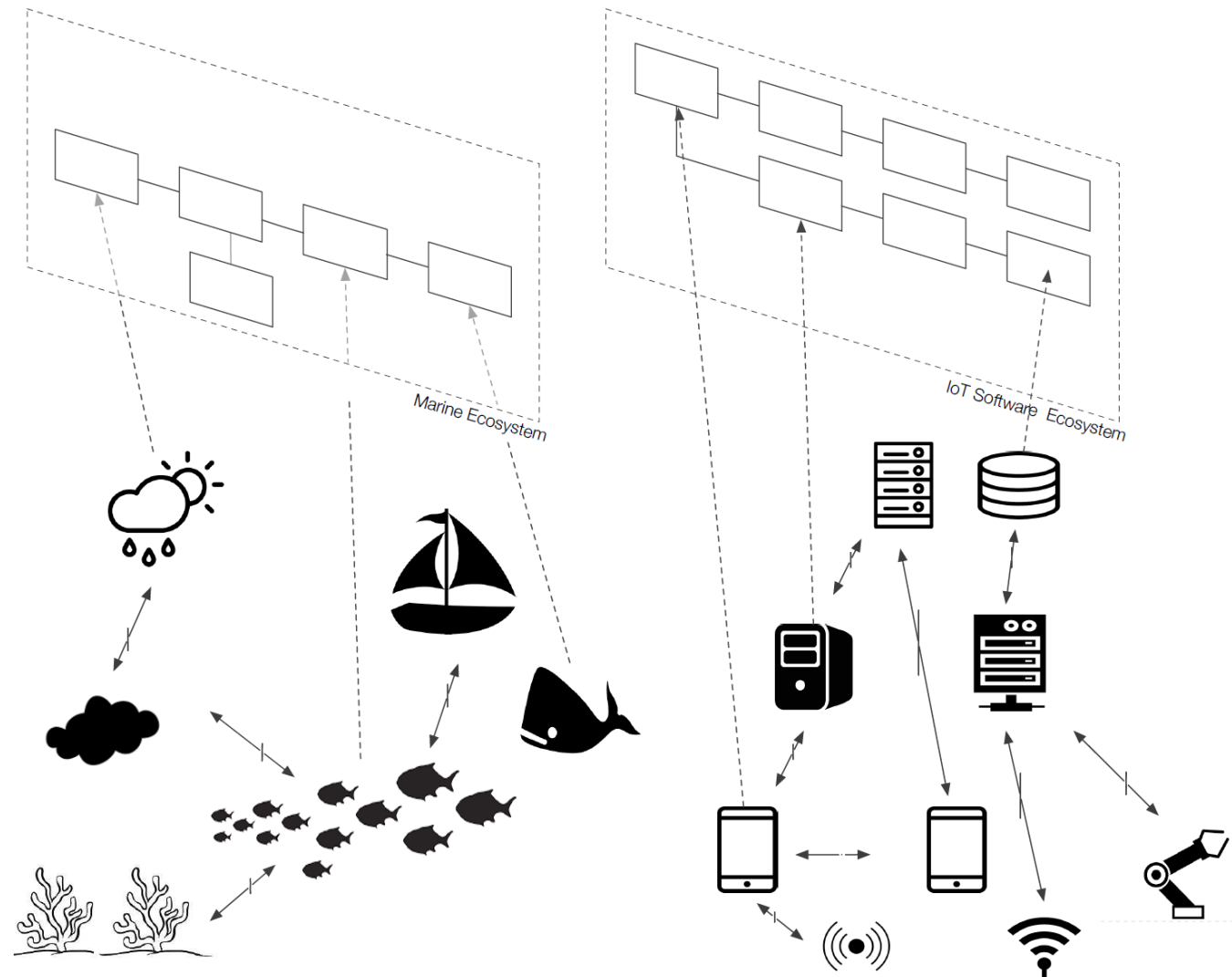# Smart Evolution – People, Services, and Things



Smart Energy Networks

Smart eGovernments & eAdministrations

Smart Homes

STB

Game Machine

TV

PC

Audio

DVD

Telephone

**Elastic Systems & Processes**

eHealth & Smart Health networks

Smart Transport Networks

# Ecosystems: People, Systems, and Things
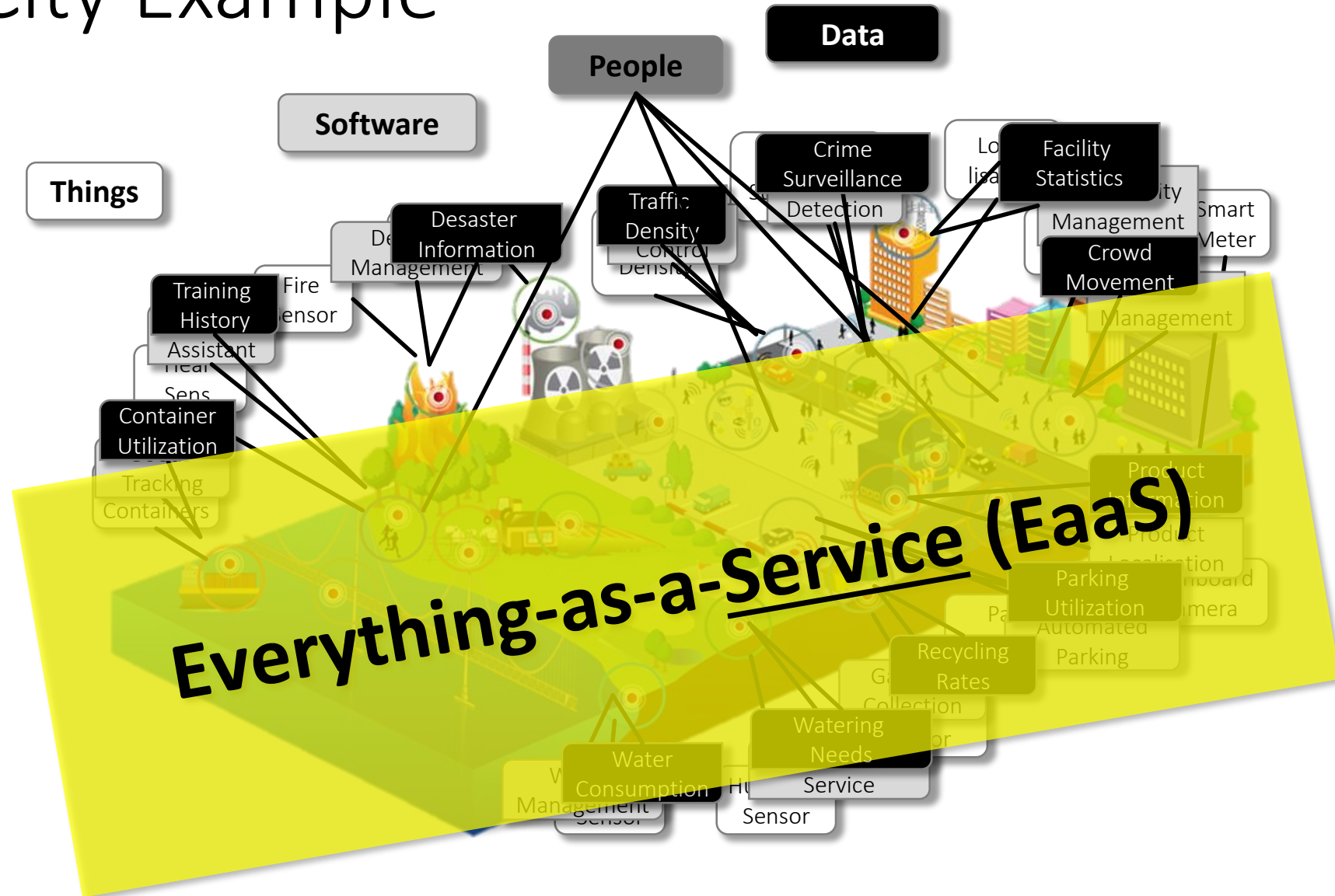


Marine Ecosystem

**Complex system** with networked dependencies and intrinsic adaptive behavior – has:

1. **Robustness & Resilience mechanisms**: achieving stability in the presence of disruption

2. **Measures of health**: diversity, population trends, other key indicators
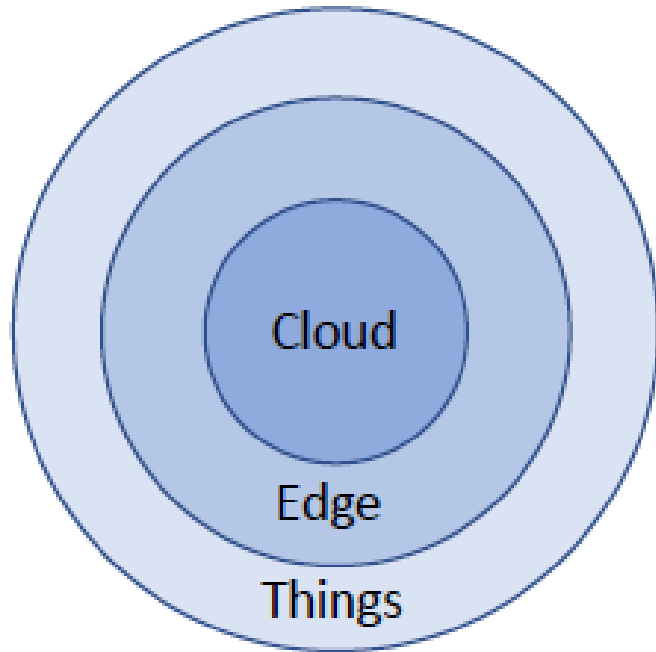
3. **Built-in coherence**

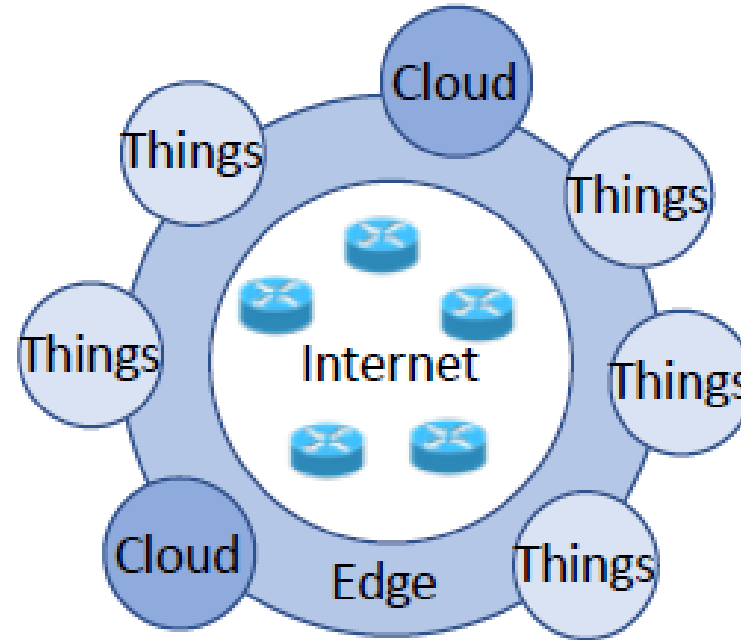4. **Entropy-resistence**

# Ecosystems for IoT Systems



Marine Ecosystem

IoT Software Ecosystem

# Smart City Example



Things

Software

People

Data

Training History
Assistant

Container Utilization

Tracking Containers

Fire Sensor

Desaster Information
Management

Traffic Density
Control Density

Crime Surveillance Detection

Facility Statistics
Management

Crowd Movement
Management

Smart Meter

Everything-as-a-Service (EaaS)

Product Information

Product Localization

Parking Utilization
Automated Parking

Recycling Rates

Watering Needs
Service

Water Consumption
Management Sensor

# Perspectives on the IoT: Edge, Cloud, Internet



**(a)** A cloud-centric perspective: Edge as "edge of the cloud"

**(b)** An Internet-centric perspective: Edge as "edge of the Internet"

Kim, H., Lee, E.A., Dustdar, S. (2019). Creating a Resilient IoT With Edge Computing, *IEEE Computer, 52/8, August 2019*

# Cloud-centric perspective

## Assumptions

- Cloud provides core services; Edge provides local proxies for the Cloud (offloading parts of the cloud's workload)

## Edge Computers

- play supportive role for the IoT services and applications

- Cloud computing-based IoT solutions use cloud servers for various purposes including massive computation, data storage,  communication between IoT systems, and security/privacy

## Missing

- In the network architecture, the cloud is also located at the network edge, not surrounded by the edge

- Computers at the edge do not always have to depend on the cloud; they can operate autonomously and collaborate with one another directly without the help of the cloud
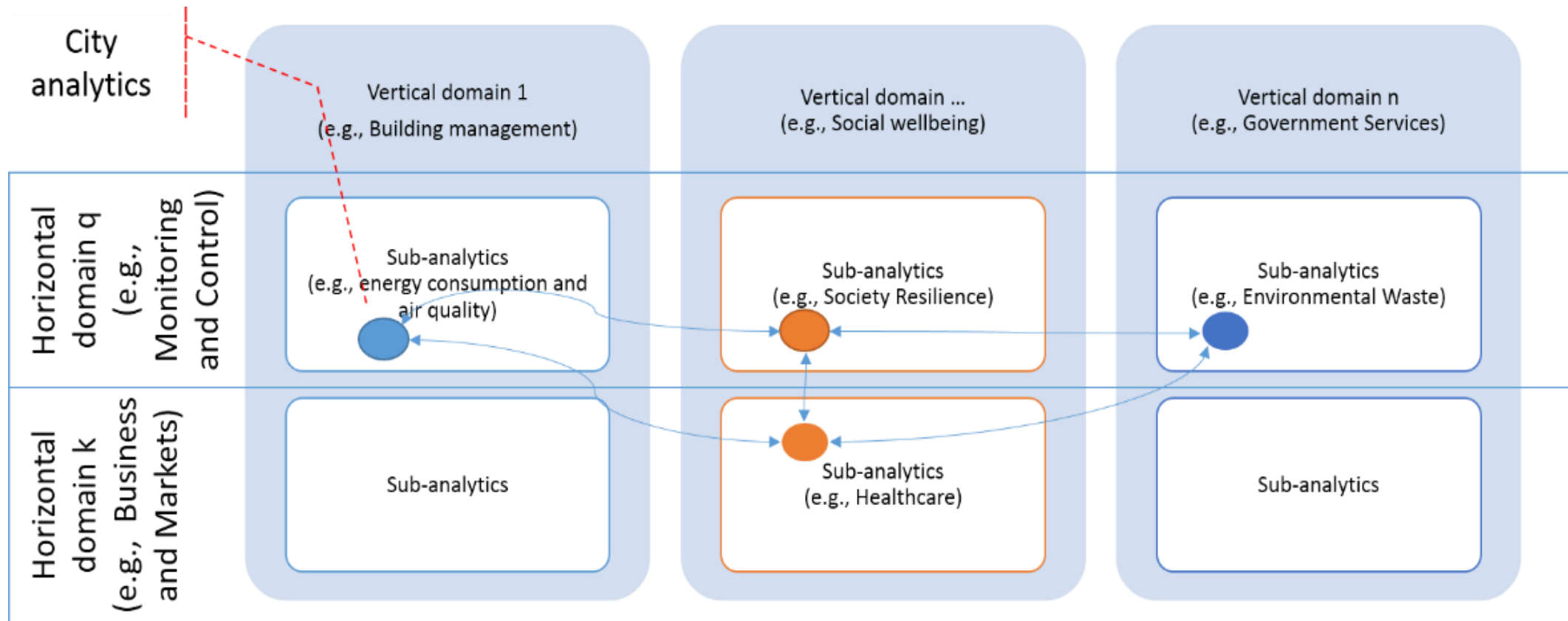
# Internet-centric perspective

## Assumptions

- Internet is center of IoT architecture; Edge devices are gateways to the Internet (not the Cloud)

- Each LAN can be organized around edge devices autonomously
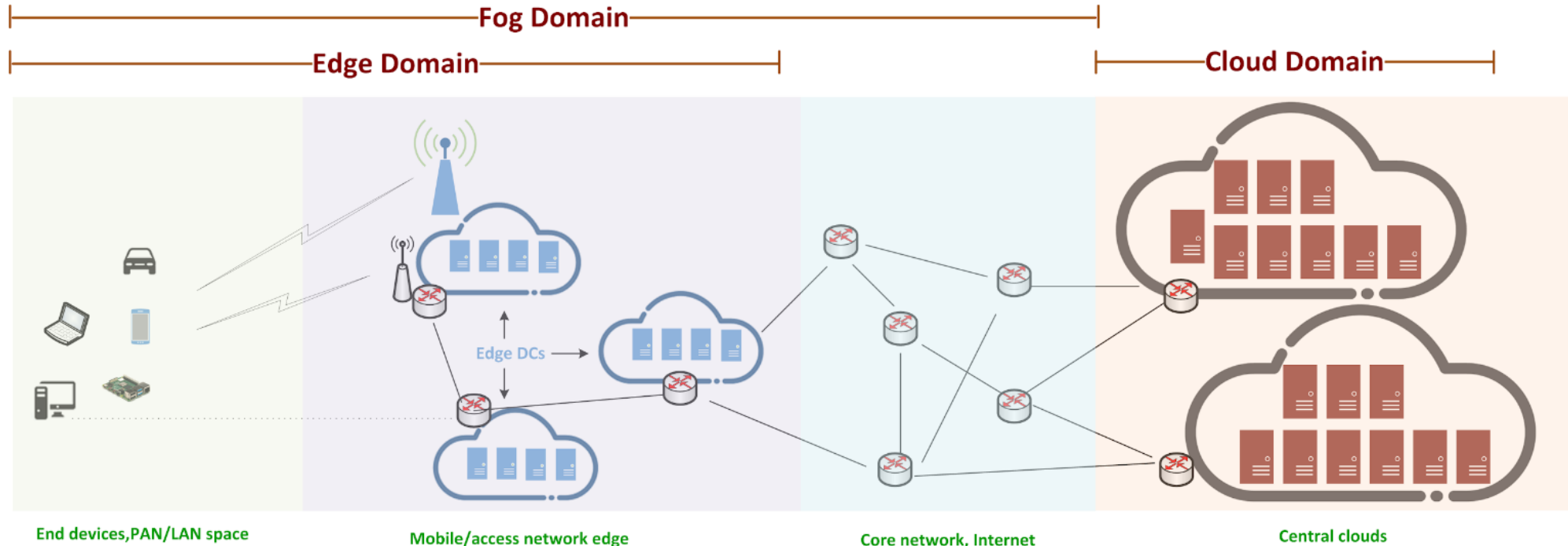
- Local devices do not depend on Cloud

## Therefore

- Things belong to partitioned subsystems and LANs rather than to a centralized system directly
- The Cloud is connected to the Internet via the edge of the network

- Remote IoT systems can be connected directly via the Internet. Communications does not have to go via the Cloud

- The Edge can connect things to the Internet and disconnect traffic outside the LAN to protect things -> IoT system must be able to act autonomously

# Dynamic Analytics (e.g., Smart City)

# IoT/Edge/Fog/Cloud Continuum: A high level view



**Fog Domain**

**Edge Domain** | **Cloud Domain**

Edge DCs

End devices, PAN/LAN space | Mobile/access network edge | Core network, Internet | Central clouds

Low reliability

Volatility

Mobility

(Mostly) Wireless connectivity

Small form factor

Battery constraints

Mobile, IoT, smart home, vehicles, …

**User/Service provider controlled**

Edge of the (mobile) network

Low latency to end device

Close to/collocated with 4G/5G base stations

General purpose compute infrastructure

Standards-based architectures & management/orchestration stacks
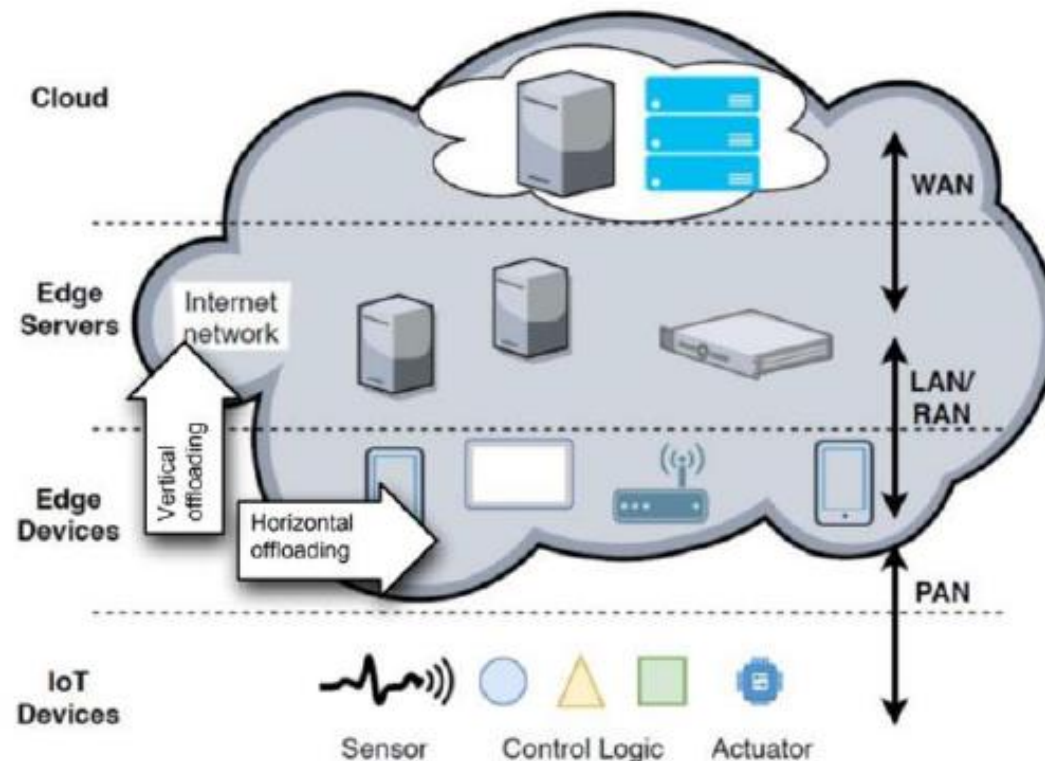
**Telecom operator controlled**

"Unlimited" compute/storage resources

Full spectrum of cloud services

High availability

Lower cost

Higher latency vs. edge/fog

**Cloud provider controlled**

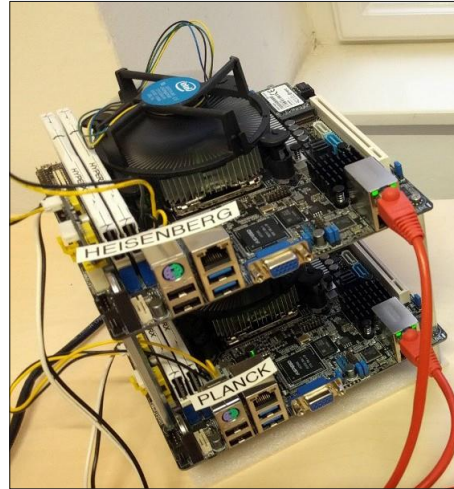# Vertical vs. Horizontal Edge/Fog/Cloud Architecture

Cloud Computing
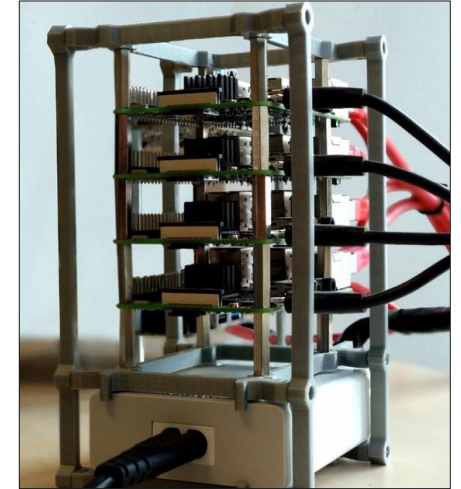
Fog Computing

Edge Computing

# Computing Continuum



Sun Modular Datacenter

Mini-ITX Servers [1]

Ubuntu Orange Box
(Intel NUC cluster)

"Micro Clouds" [2]

Server Computers

SOC & Single Board Computers

1. Rausch T., Avasalcai C., Dustdar S. (2018). Portable Energy-Aware Cluster-Based Edge Computers. 3rd ACM/IEEE Symposium on Edge Computing (SEC 2018), October 25-27, 2018, Bellevue, WA, USA

2. Elkhatib et al., 2017, "On Using Micro-Clouds to Deliver the Fog"

# Towards Edge Intelligence

## Computational Fabric

- dispersed resources allow training, monitoring, serving of models
- Heterogeneity of applications and models requires
  - (1) flexible and modular **infrastructure** and
  - (2) intelligent operations **mechanisms** (due to the <u>scale</u> of the infrastructure)

## Operationalization

- Automated AI application lifecylce management to the Edge

Rausch, T., Dustdar, S. (2019). Edge Intelligence: The Convergence of Humans, Things, and AI. In *IEEE International Conference on Cloud Engineering (IC2E) 24-27 June 2019*.

# Fabric for Edge Intelligence
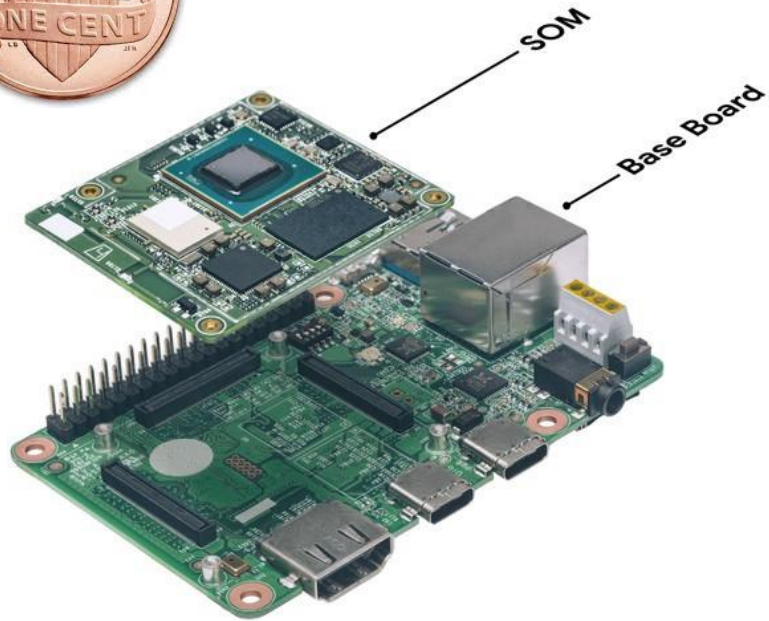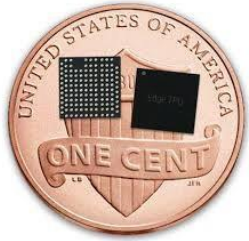
1. **Sensing (Sensor Data as a Service)**
   - Large number, dynamic and mobile nature of publishers/subscribers of sensor data + QoS requirements of edge intelligence
   ->> rethink centralized messaging services such as AWS IoT or MS Azure IoT Hub

   - Management and governance of such a distributed/decentralized sensing infrastructure

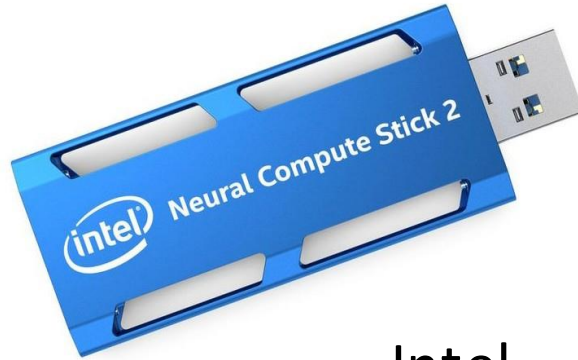2. **Edge computer network with modular AI capabilities**
   - New AI accelarators for edge devices (e.g., Google Edge TPU with an aplication specific integrated circuit; MS BrainWave with field-programmable gate arrays (FPGAs); Intel Neural Compute Stick; Baidu Kunlun, Huawei Atlas AI Platform

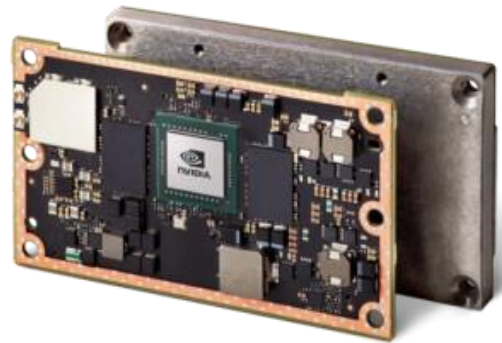3. **Intelligent orchestration mechanisms for decentralized and distributed infrastructure**

# Edge AI Accelerators



SOM

Base Board

Intel
Neural Compute Stick

Baidu Kunlun

Microsoft
Project BrainWave

Google Edge TPU
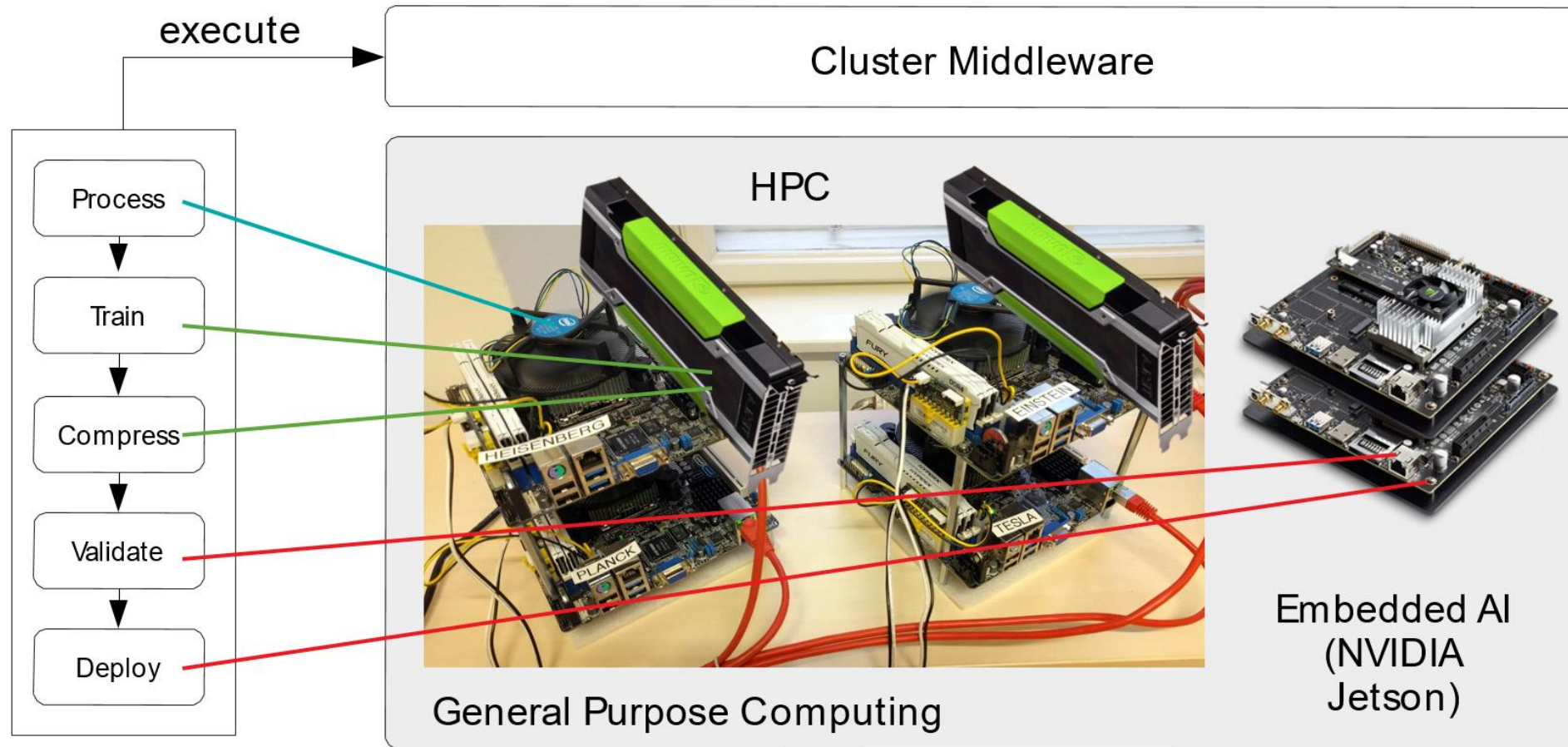
NVIDIA Jetson

Atlas 200

Atlas 300

Atlas 500

Huawei Atlas

# Edge Intelligence Fabric

Rausch T., Avasalcai C., Dustdar S. (2018). Portable Energy-Aware Cluster-Based Edge Computers. 3rd ACM/IEEE Symposium on Edge Computing (SEC 2018), October 25-27, 2018, Bellevue, WA, USA

Nastic S., Rausch T., Scekic O., Dustdar S., Gusev M., Koteska B., Kostoska M., Jakimovski B., Ristov S., Prodan R. (2017). A Serverless Real-Time Data Analytics Platform for Edge Computing. IEEE Internet Computing, Volume 21, Issue 4, pp. 64-71

Rausch T., Dustdar S., Ranjan R. (2018). Osmotic Message-Oriented Middleware for the Internet of Things.IEEE Cloud Computing, Volume 5, Issue 2, pp. 17-25

# Elasticity (Resilience)

(Physics) The property of returning to an initial form or state following deformation

**stretch** when a force stresses them

   e.g., **acquire** *new resources,* **reduce** *quality*

**shrink** when the stress is removed

   e.g., **release** *resources,* **increase** *quality*

# Elastic Computing > Scalability



**Resource** elasticity

Software / human-based computing elements, multiple clouds

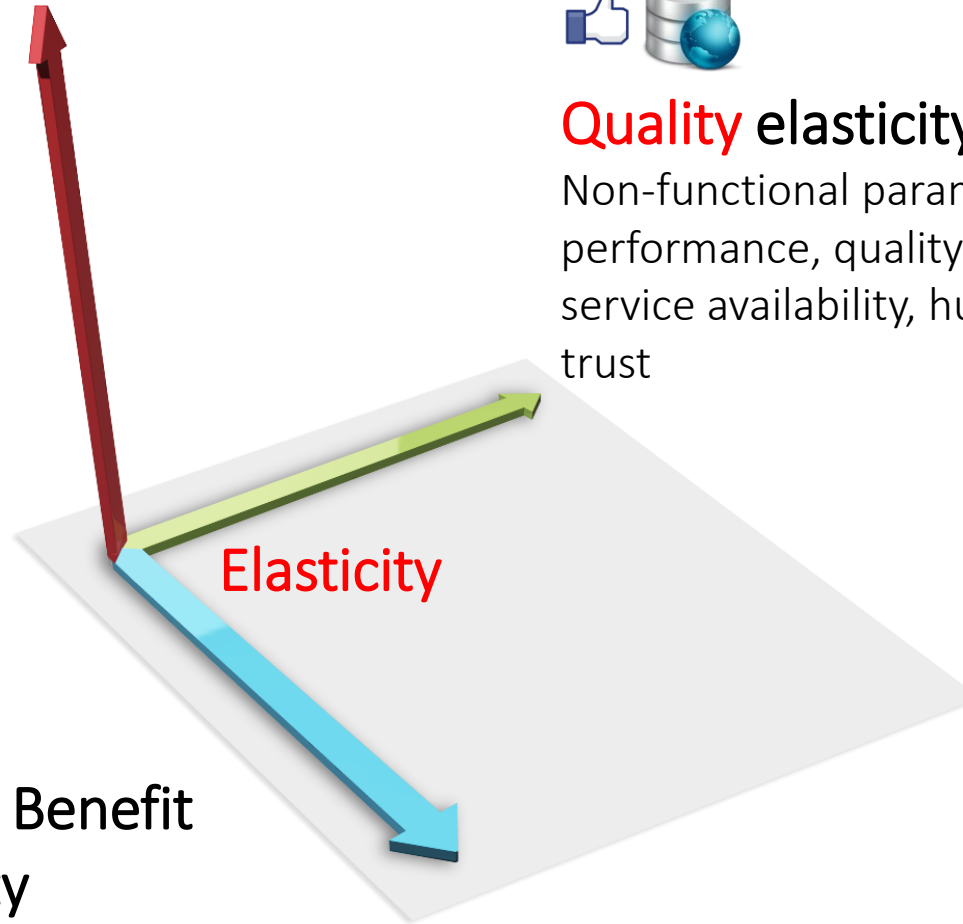**Quality** elasticity

Non-functional parameters e.g., performance, quality of data, service availability, human trust

**Elasticity**

**Costs** & Benefit elasticity

rewards, incentives

Dustdar S., Guo Y., Satzger B., Truong H. (2012) Principles of Elastic Processes, IEEE Internet Computing, Volume: 16, Issue: 6, Nov.-Dec. 2012

# High level elasticity control

**#SYBL.CloudServiceLevel**
**Cons1: CONSTRAINT responseTime < 5 ms**
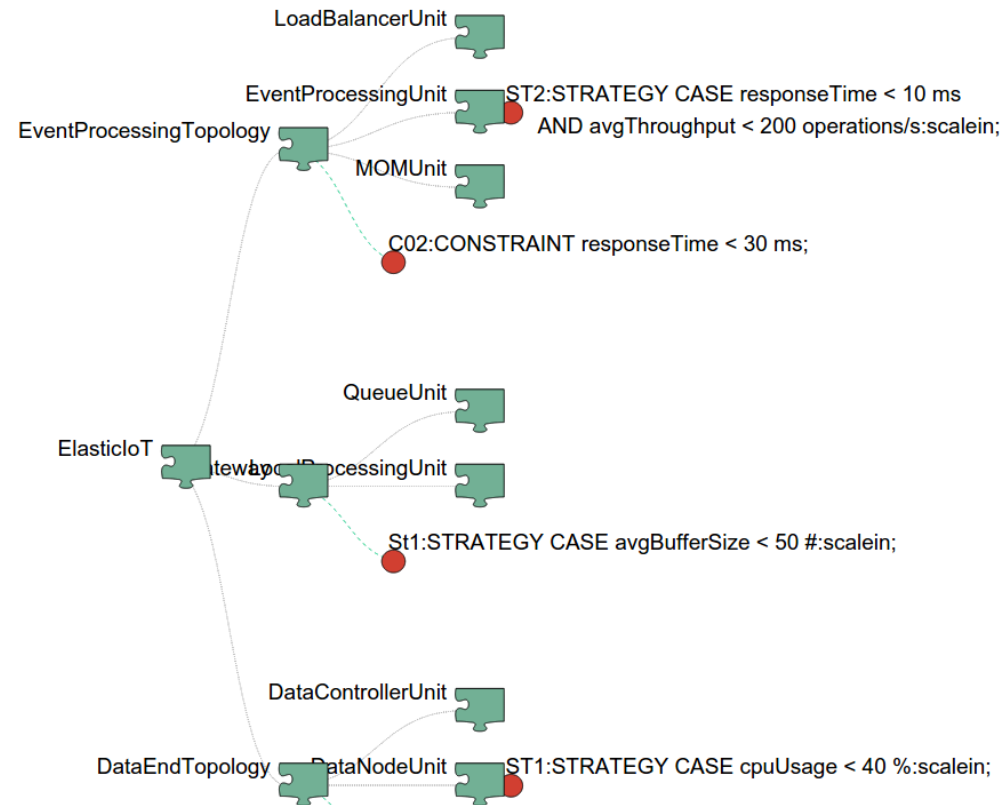**Cons2: CONSTRAINT responseTime < 10 ms**
**WHEN nbOfUsers > 10000**
**Str1: STRATEGY CASE fulfilled(Cons1) OR**
**fulfilled(Cons2): minimize(cost)**

**#SYBL.ServiceUnitLevel**
**Str2: STRATEGY CASE ioCost < 3 Euro :**
**maximize( dataFreshness )**

**#SYBL.CodeRegionLevel**
**Cons4: CONSTRAINT dataAccuracy>90% AND**
**cost<4 Euro**



Georgiana Copil, Daniel Moldovan, Hong-Linh Truong, Schahram Dustdar, **"SYBL: an Extensible Language for Controlling Elasticity in Cloud Applications"**, 13th IEEE/ACM International Symposium on Cluster, Cloud and Grid Computing (CCGrid), May 14-16, 2013, Delft, Netherlands

Copil G., Moldovan D., Truong H.-L., Dustdar S. (2016). **rSYBL: a Framework for Specifying and Controlling Cloud Services Elasticity**. *ACM Transactions on Internet Technology*

# Elasticity Model for Edge & Cloud Services

**Elasticity Pathway functions**: to characterize the elasticity behavior from a general/particular view

**Elasticity space functions**: to determine if a service unit/service is in the "elasticity behavior"

# Growing interest in federated learning

- Training on data **directly on remote devices**...

- ...**without revealing** the data themselves

- Sending the outcome of local training to server (**local updates**)

- Server aggregates these updates into a **global model**
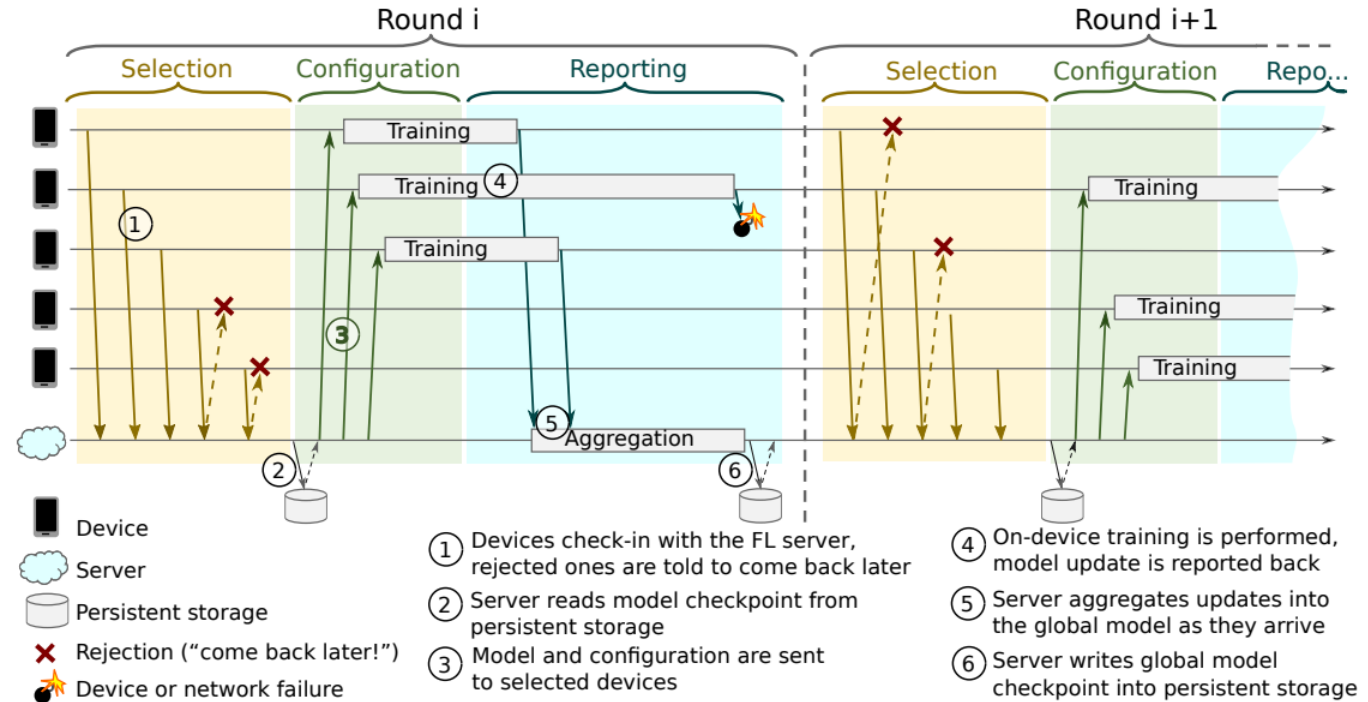
- Makes the model available to devices



Figure source & further reading: K. Bonawitz et al., "Towards Federated Learning at Scale: System Design," arXiV:1902.01046, March 2019. Available: https://arxiv.org/pdf/1902.01046.pdf

# Applications

- For mobile devices
  - Next-word prediction, face detection, voice recognition
  - Train on data from smartphone text editors, cameras, mics
  - Users do not wish to reveal their messages, photos, and videos
  - Also, they don't want to waste bandwidth and MBs from their data plan


- For organizations
  - Organizations such as hospitals have data, but should not expose them
  - Federating such data in a private way to apply ML for medical and other research


- For environmental, transportation, smart home, and other applications
  - Measurement devices with sensors (e.g., for air pollution) mounted on cars
  - Sensors in a smart home
  - Pushing data to servers for centralized training might leak driver patterns, daily habits, etc.

# Current research challenges

**Device recruitment strategies:** Which subset of the devices to assign a learning task at any given round? Processing, storage, battery, trustworthiness, data quality and other criteria to consider

**Volatility:** Devices can "disappear" or join at any time

**Asynchrony:** Algorithms face challenges when end devices do not submit their data in a timely manner

**Non independent and identically distributed data:** inaccuracies, personalization lost

**Heterogeneity in the volume of training data per device:** A device that contributes a lot may lead to a biased model
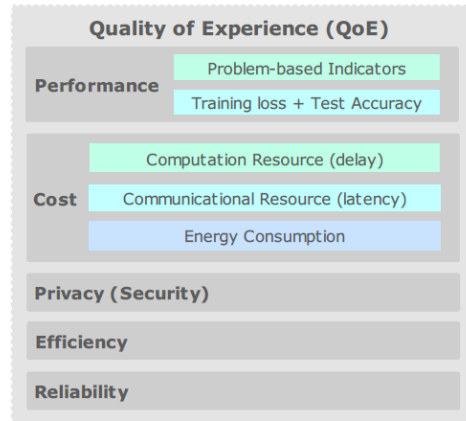
**Preventing privacy leaks:** Some private information may be inferred even if devices do not transmit the actual data

**Incentives to misbehave:** Why waste battery when I can let the others do all the work?

Further reading: T. Li at al., "Federated Learning: Challenges, Methods, and Future Directions," arXiv:1908.07873, August 2019. Available: https://arxiv.org/pdf/1908.07873.pdf

# Research Roadmap – Quality of Experience

Edge Intelligence: The Confluence of Edge Computing and Artificial Intelligence, *IEEE IoT Journal 2020, forthcoming*



1. **Performance**
E.g., the ratio of computation offloading

**2. Cost**
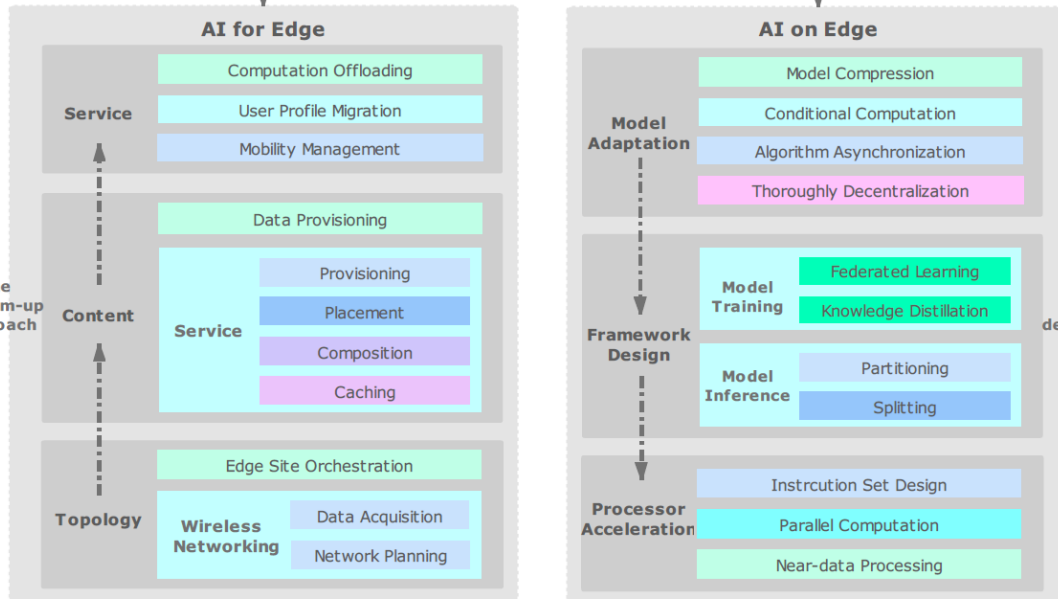Computation|Communication|Energy consumption costs

**3. Privacy & Security**
Federated learning, i.e., aggregating local machines models from distributed edge devices

**4. Efficiency**
Excellent performance with low overhead, e.g., model compression, conditional computation

**5. Reliability**
Relates to model upload and download and wireless network congestion

# AI for Edge



1. **Topology**
- Edge orchestration and coordination with small base stations
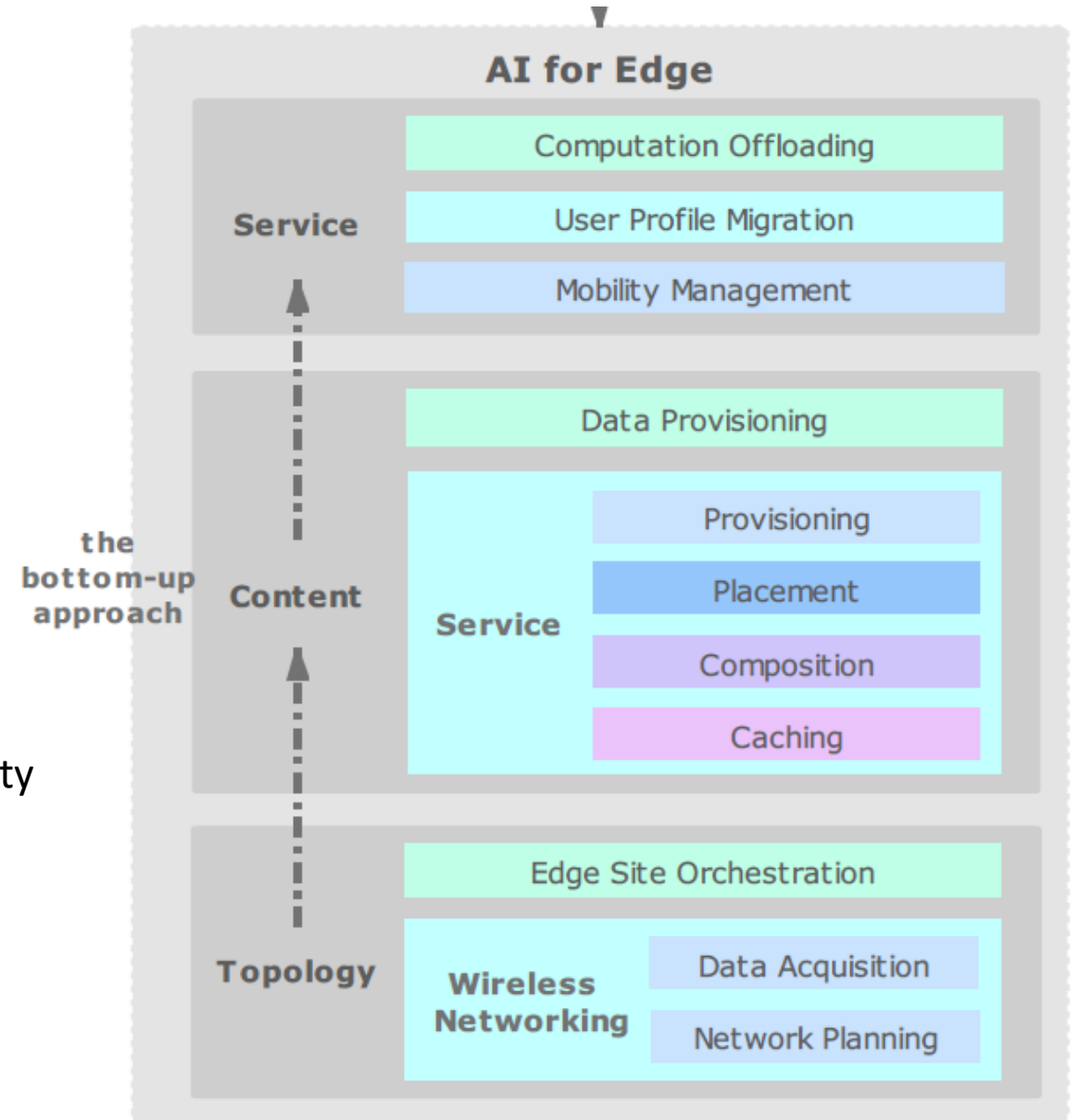- Unmanned Aerial Vehicles (UAVs) and access points

**2. Content**
Lightweight service frameworks for QoS-aware services, e.g., on mobile devices

**3. Service**
Computation offloading, User profile migration and mobility management

Edge Intelligence: The Confluence of Edge Computing and Artificial Intelligence, *IEEE IoT Journal 2020, forthcoming*

# Grand Challenges – AI for Edge

- **Model Establishment – restraining the optimization model**
  - Stochastic Gradient Descent (SGD)
  - MBGD (Mini-Batch Gradient Descent)

- **Algorithm Development**
  - Selection of which edge device should be responsible for deployment and execution in an online manner
  - SOTA formulates combinatorial and NP-hard optimization problems with high computational complexity

- **Trade-off between optimality and efficiency**
  - Consider resource constraint devices

Edge Intelligence: The Confluence of Edge Computing and Artificial Intelligence, *IEEE IoT Journal 2020, forthcoming*
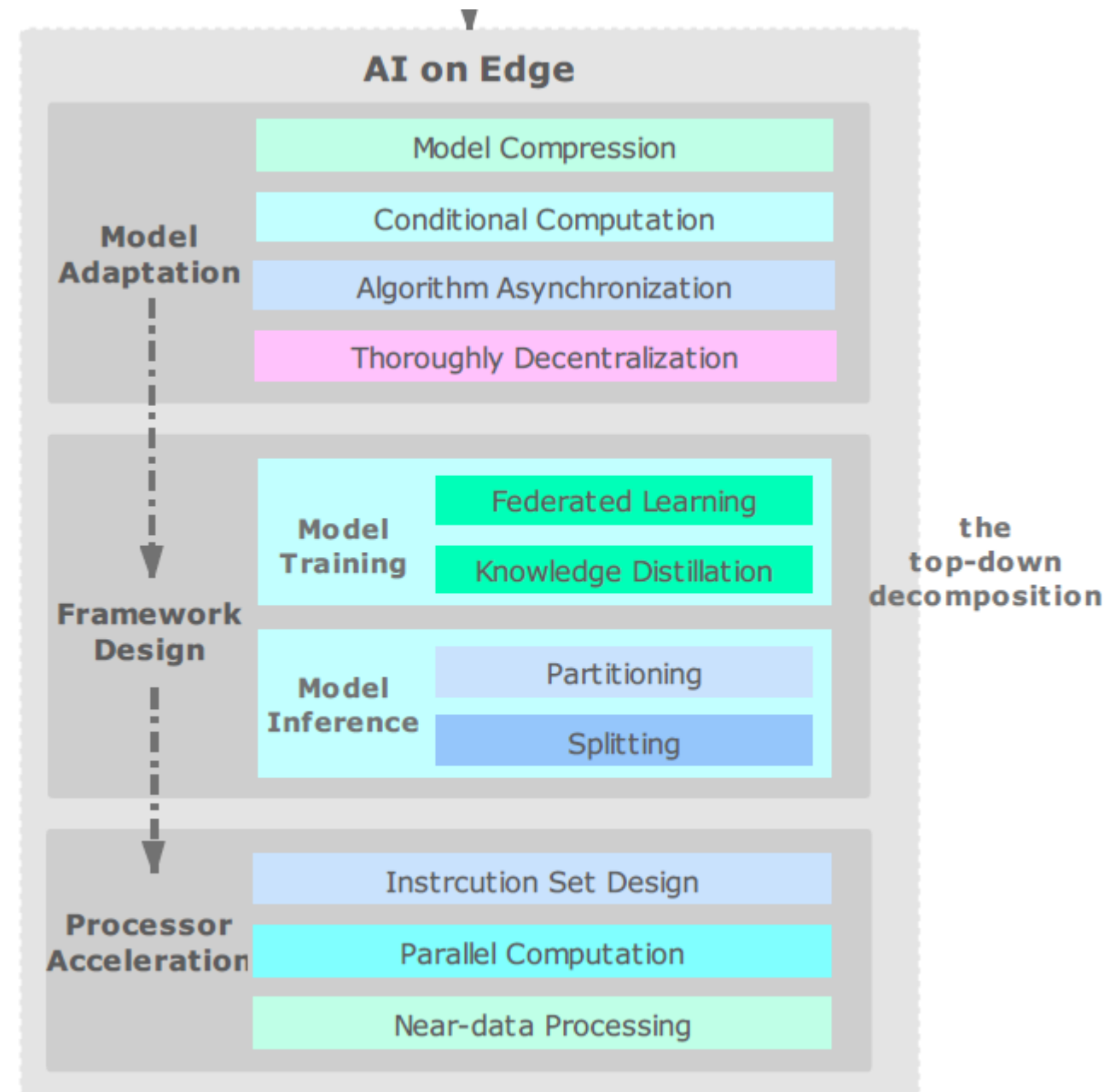
# AI on Edge

- **Data Availability**
  - Chellenge of lack of availability and usability of raw training data for model training and inference
  - Bias of raw data from various end user/mobile devices

- **Model Selection**
  - SOTA requires selection of need-to-be trained AI models has challenges
  - Threshold of learning accuracy and scale of AI models for quick deployment and delivery
  - Selection of probe training frameworks and accelerator architectures under limited resources
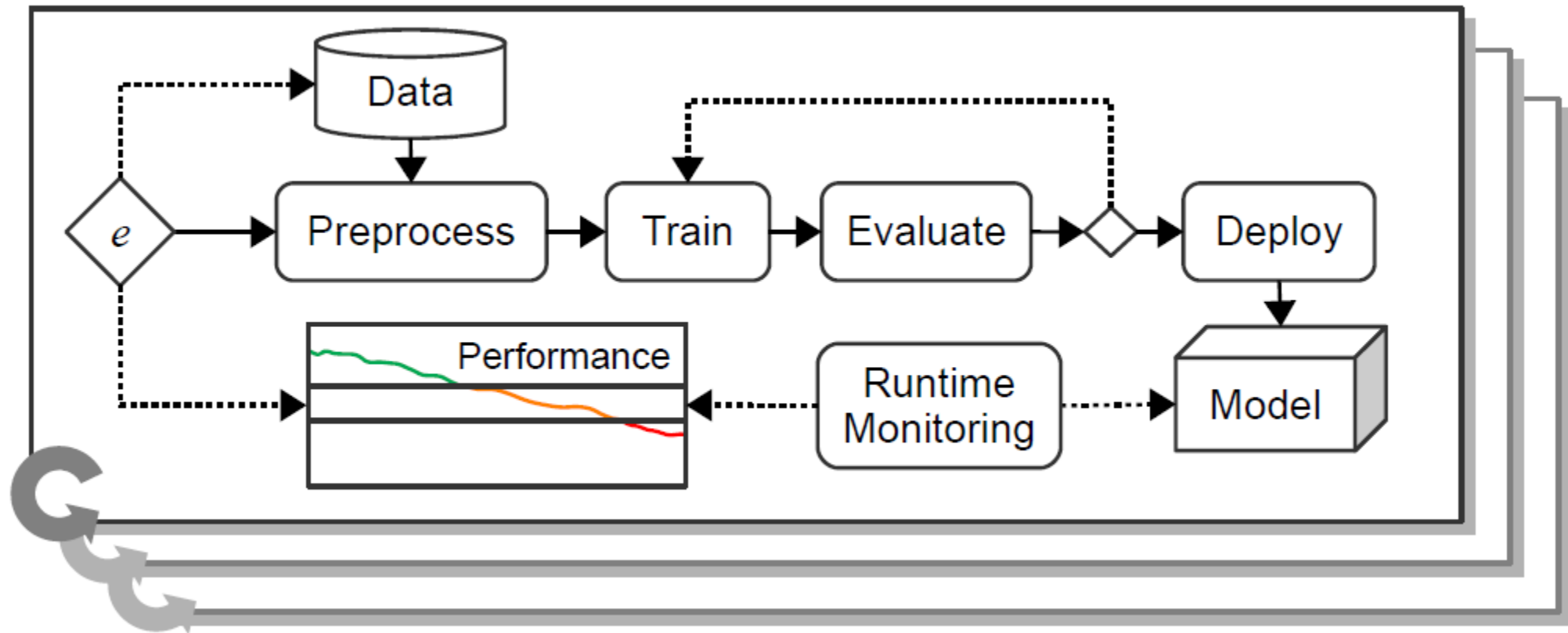
- **Coordination Mechanisms**
  - Cordination between heterogeneous edge devices, cloud, and various middlewares and APIs



Edge Intelligence: The Confluence of Edge Computing and Artificial Intelligence, *IEEE IoT Journal 2020, forthcoming*

# Managing the AI Lifecycle

AI lifecycle pipeline with a rule-based trigger *e* that monitors available data and runtime performance data to form an automated retraining loop

# AI Operations Workflows – Edge to Cloud

|  | Data characteristics | Model characteristics | Enabling technologies | Example use cases |
|---|---|---|---|---|
| **C2C** | - Training data is centralized<br>- Massive data sets | - Models are large<br>- Huge number of inferencing requests need to be load balanced | - Scalable learning infrastructure [39]<br>- Data warehousing | - Image search<br>- Recommender systems |
| **C2E** | - Training data is centralized<br>- Inferencing data may be sensitive | - Inferencing may need to happen in near-real time<br>- Large number of model deployments<br>- Models run on specialized hardware | - Model compression [42]<br>- Latency/accuracy tradeoff [43]<br>- Distributed inferencing [44]<br>- Transfer learning [45] | - Surveillance systems<br>- Self driving cars<br>- Fieldwork assistants |
| **E2C** | - Training data is distributed<br>- Training data may be sensitive | - Models can be centralized<br>- Huge number of inferencing requests need to be load balanced | - Decentralized/federated learning [41] | - Volunteer computing<br>- Novel Smart City use cases |
| **E2E** | - Training data is distributed<br>- Training and inferencing data may be sensitive | - Inferencing may need to be near-real time | - Decentralized/federated learning<br>- Distributed inferencing | - Industrial IoT (e.g., predictive maintenance)<br>- Privacy-aware personal assistants<br>- Novel IoT use cases |

Rausch, T., Dustdar, S. (2019). Edge Intelligence: The Convergence of Humans, Things, and AI. In *IEEE International Conference on Cloud Engineering (IC2E) 24-27 June 2019*.

# Conclusions

- Leverage the Computing "Continuum" from IoT->Edge->Fog->Cloud

- Differentiate between AI <u>for</u> Edge and AI <u>on</u> Edge. Both bring their distinct research challenges

- Need for an Edge Intelligence AI Fabric and a "clear" ecosystems understanding

# Thanks for your attention

Prof. Schahram Dustdar

IEEE TCSVC Outstanding Leadership
Award in Services Computing

Member and Chairman of Informatics
at *Academia Europaea*

IBM Faculty award

ACM Distinguished Scientist

IEEE Fellow

Distributed Systems Group
TU Wien, Austria
**dsg.tuwien.ac.at**