

UNIVERSITY OF TARTU
Faculty of Science and Technology
Institute of Computer Science
Computer Science Curriculum

Jeyhun Abbasov

AI-powered cloud usage controlling system in smart home environment

Master's Thesis (30 ECTS)

Supervisor: Chinmaya Kumar Dehury, Ph.D.

Tartu 2023

AI-powered cloud usage controlling system in smart home environment

Abstract: Cloud computing is utilized for handling and keeping the large amounts of data produced by IoT devices in a smart home environment. One of the major limitations of cloud computing is network latency. Due to these limitations, the Fog computing is introduced. Fog computing provides real-time services and saves network resources. A large portion of research are aimed at ensuring high-quality service delivery through the utilization of either fog or cloud environments. In our study, we present a Deep Reinforcement Learning service distribution solution using both fog and cloud computing environments in real-time manner. The demand for services is divided and distributed among the fog and cloud environments without compromising the Quality of Services (QoS). The results of implementation demonstrate substantial enhancement compared to state of the art studies.

Keywords: Cloud Computing, Fog Computing, Deep Reinforcement Learning, Service Delivery, Service Dispersal

CERCS: P170 Computer science, numerical analysis, systems, control

AI-jõuline pilvekasutuse juhtimissüsteem nutikas kodu keskkonnas

Lühikokkuvõte: Pilvandmetöötlust kasutatakse IoT-seadmete toodetud suurte andme-
mahtude haldamiseks ja hoidmiseks nutikas kodukeskkonnas. Üks pilvandmetöötluse
peamisi piiranguid on võrgu latentsus. Nende piirangute tõttu võetakse kasutusele uduar-
vutus. Uduarvutus pakub reaajas teenuseid ja säästab võrguressursse. Suur osa uuringu-
test on suunatud kvaliteetse teenuse osutamise tagamisele kas udu- või pilvekeskkonna
kasutamise kaudu. Oma uuringus tutvustame Deep Reinforcement Learning teenuse
levitamislahendust, mis kasutab reaajas nii udu- kui ka pilvandmetöötluskeskkondi.
Nõudlus teenuste järele on jagatud ja jaotatud udu- ja pilvekeskkondade vahel, ilma et see
kahjustaks teenuste kvaliteeti (QoS). Rakendamise tulemused näitavad olulist paranemist
võrreldes nüüdisaegsete uuringutega.

Võtmesõnad: Pilvandmetöötlus, Uduarvutus, Sügav Õpetamisreinforcement, Teenuse
Tarnimine, Teenuse Jaotus

CERCS:P170 Arvutiteadus, arvutusmeetodid, süsteemid, juhtimine

Acknowledgements

I would like to thank my family, friends who always motivated me during this research and my supervisor Prof. Dr. Chinmaya Dehury for giving me detailed guidance, providing direction to the correct path.

Contents

1	Introduction	7
1.1	Motivation	9
1.2	Goal and Contributions	11
1.3	Outline	11
2	State of the Art	12
2.1	Background	12
2.2	Related Work	13
2.3	Summary	17
3	Methodology	18
3.1	Modelling Services	18
3.1.1	Application 1: Deep learning based object classification	19
3.1.2	Application 2: Text-audio synchronisation or forced alignment	20
3.1.3	Application 3: Speech-to-text conversion	21
3.2	Main Parameters of the System	23
3.2.1	Environment Parameters	23
3.2.2	User Parameters	23
3.2.3	Service Parameters	23
3.2.4	Priority Parameters	23
3.3	DRL-based Service Placement	23
3.3.1	State Space	23
3.3.2	Action Space	23
3.3.3	Rewards	23
3.4	Summary	23
4	Experiment	24
4.1	Experimental Setup	24
4.2	Summary	24
5	Results	25
5.1	Summary	25
6	Conclusion	26
	References	32

Appendix	33
I. List of Figures	33
II. List of Tables	34
III. Source Code	35
IV. Licence	36

1 Introduction

With the rise of the Internet of Things (IoT), smart homes have become a reality, offering greater convenience and control over our daily lives [Ela19]. IoT devices are equipped with sensors and actuators that allow us to monitor and control various aspects of our home environment, such as temperature, lighting, and security, from anywhere at any time through dedicated mobile applications or web services. This has revolutionized the way we live, work, and interact with our homes. In addition, the increasing demand for computing and storage resources from these smart home services is driving the development of cloud computing infrastructure, making it possible to store and process large amounts of data in real-time, further enhancing the capabilities of our smart homes.

Cloud computing is comprised of a vast variety of powerful physical servers that provide an almost infinite number of computation, storage, and networking capabilities to enable IoT services. In cloud computing, there are typically three categories of services: infrastructure as a service (IaaS), platform as a service (PaaS), and software as a service (SaaS). IaaS is the most basic type of cloud computing. It provides access to a virtualized environment that includes computer resources such as servers and virtual machines, storage, networks, and operating systems. PaaS provides an on demand platform for developing, testing, deploying, and managing software without having to worry about establishing or maintaining the necessary IT infrastructure such as servers, storage, networking, and databases. SaaS is a vendor-provided cloud-based software that may be accessed via a web browser. Customers may subscribe to it and enjoy a variety of applications like as accounting, sales, invoicing, and more. The advantages of SaaS include low cost, great performance, and dependability. Despite its advantages, cloud computing has drawbacks, such as geographical constraints and sophisticated network architecture [SDV18]. This is where Fog Computing comes in, filling the gap between the traditional cloud computing and IoT devices. It extends cloud computing closer to the edge devices, improving mobility support, enabling location awareness, and reducing latency for better performance and user experience. [SWHL16]

Fog computing is proposed as a solution to address the previously mentioned issue. Fog computing solves the issues of low latency, location awareness and improves quality-of-service for real time and streaming applications by being deployed at the edge of the network. Examples of fog computing applications can be found in industrial automation, transportation and networks that involve sensors and actuators. In addition, the fog infrastructure accommodates different types of devices, such as user end devices, access points, edge routers, and switches. Fog computing is ideal for real-time big data analytics and provides benefits in various fields like entertainment, advertising, personal computing, etc. by enabling densely distributed data collection points. By utilizing local resources in fog computing, a user's service requirements can be met with less reliance on distant cloud resources, leading to improved performance. Implementing fog computing can help reduce service latency, improve overall Service Response Time (SRT), reduce

network congestion, decrease energy consumption, and enhance system security by utilizing resources at the edge organized as fog nodes [SMBMT⁺18, YHQL15].

Another crucial element in a smart home environment is the Smart Gateway. The Smart Gateway is responsible for managing various IoT operations, such as collecting and preprocessing data, filtering and reformatting it, uploading relevant data to the fog/cloud, monitoring IoT objects and sensors, monitoring energy consumption, ensuring data security and privacy, and overall service management. Additionally, The Smart Gateway can play a key role in efficient service dispersal process. Communication between IoT and Smart Gateway can either be direct or through base station(s). There are two types (Single-hop communication, Multi-hop communication) of Smart Gateway-based communication. In single-hop connectivity, the sensors and IoT devices are directly connected to the Smart Gateway. The Smart Gateway then collects the data and sends it to either the Fog or Cloud. This type of communication is typically used in smaller networks where the sensor nodes have limited roles, such as in smart health or healthcare applications. Quick monitoring and response is possible through direct communication with the Smart Gateway. M2M communication takes place in this scenario, and the Smart Gateway can perform data refinement, filtering, trimming, and security measures, service dispersal process, depending on the application needs and in conjunction with fog computing. The extent of this type of communication depends on the capabilities of the Smart Gateway device. In a scenario where multiple sensor networks and IoTs are connected, direct connections become unfeasible. These networks and IoTs have their own base stations and sink nodes, and the gateway collects data from these. This creates a multi-hop communication setup, where nodes are more diverse and spread out and data is more heterogeneous, requiring more processing and analysis from the gateway. Sink nodes provide an additional layer of communication and make underlying sensors and things a "black box," adding security. Security can be customized according to the IoT and WSN, and sink nodes can manage sensor networks based on their constraints. The gateway must handle heterogeneous data from various devices, IoTs, and WSNs, requiring transcoding and interoperability. This can either be achieved through an intelligent gateway or through Fog computing resources. This setup is suitable for large scale IoTs/WSNs and mobile objects, such as vehicle tracking IoTs and environmental monitoring. [AH14] For the purpose of keeping things simple, we have selected to use Single-hop communication in our Smart Gateway setup. As an illustrative example, Figure 1 illustrates an abstract view of Smart Gateway with Cloud-Fog architecture.

In our approach, the service workload is divided into smaller sub-services by Smart Gateway and is shared between fog and cloud servers. The fog node processes some slices of the service request and the rest is sending to the cloud for processing. The decision of which slices to be processed by the fog and cloud depends on various factors such as user proximity, user priority, service sensitivity, energy consumption, etc. Determining the correct distribution of the service between fog and cloud becomes increasingly complex

as the service complexity grows, leading to a new research challenge.

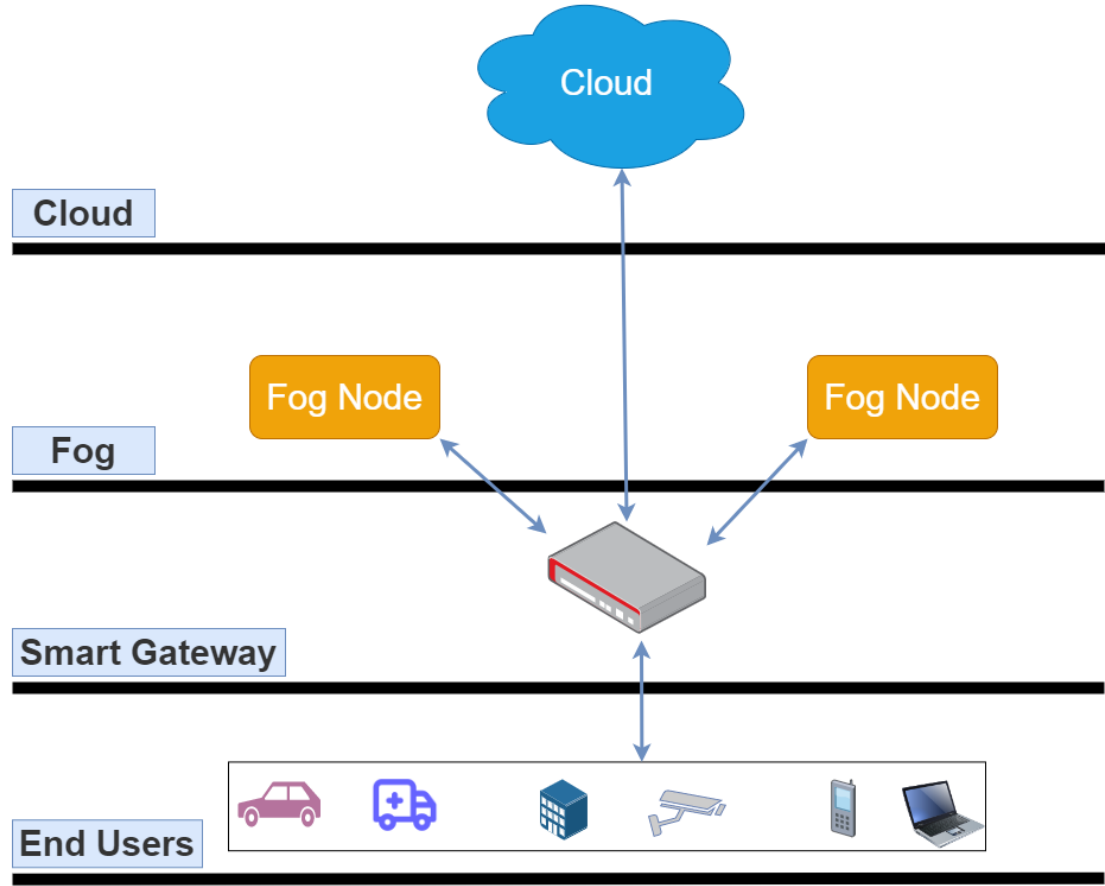


Figure 1. Abstract view of Smart Gateway with Cloud-Fog architecture.

1.1 Motivation

The allocation of services between the fog and cloud is managed by the Smart Gateway, as previously discussed, to provide services to the users through both fog and cloud environments. As seen in Figure 2, the division of services between fog and cloud for different services and users is not consistent, and an individual user may utilize a variety of services. For example, the User1 is using three services named as Service1, Service2, and Service3. The service request from User1 includes a request for 25% of Service1, a request for 35% of Service2, and a request for 40% of Service3. The 45% and 55% of Service3 that is requested by User1 is handled by Fog and Cloud, respectively. The distribution of service workload between fog and cloud servers can vary, depending on

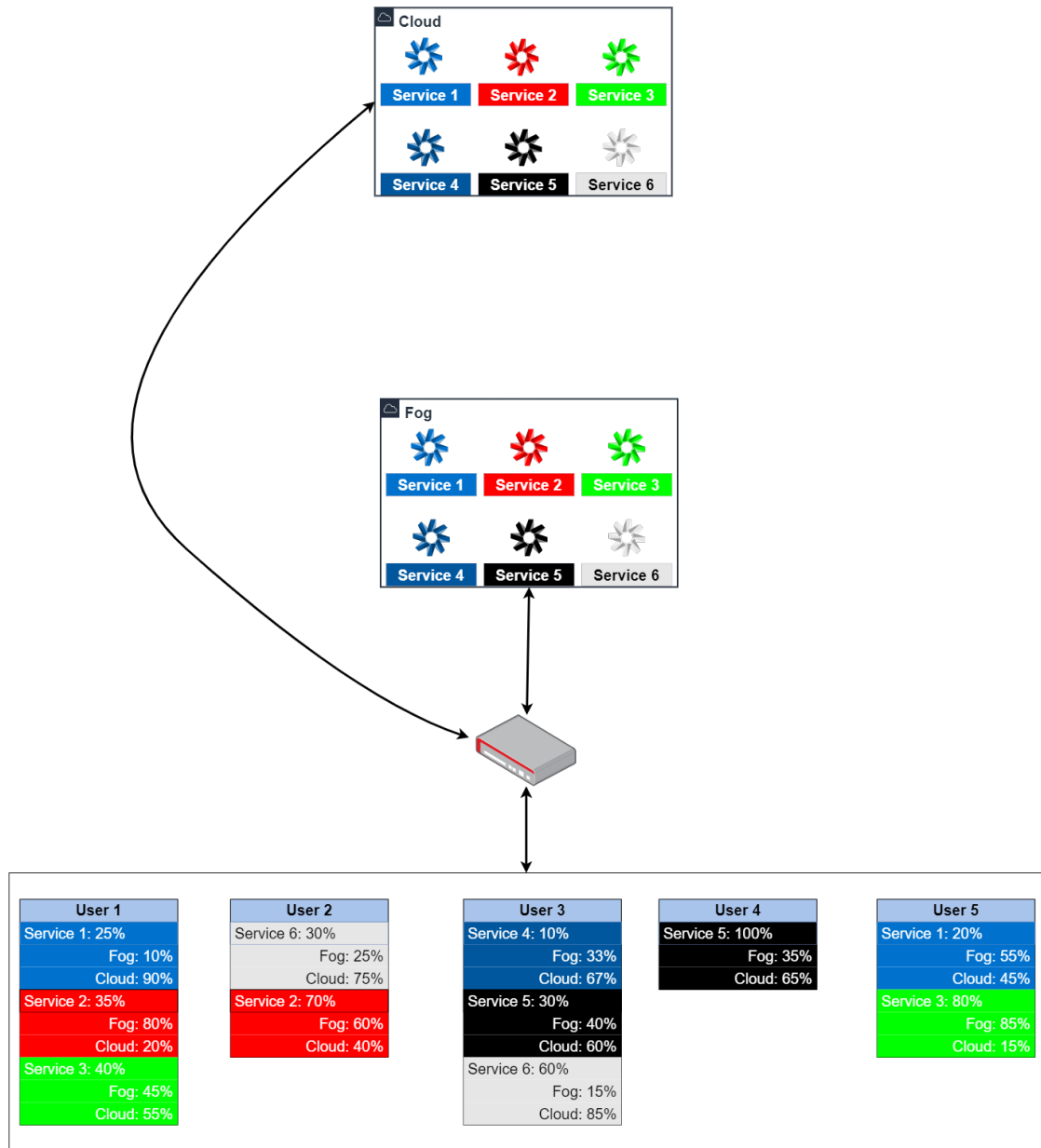


Figure 2. An example motivational scenario.

various factors such as the sensitivity of the service, user location, available resources, and others. For example, a critical service with low tolerance for latency cannot be assigned to the cloud, whereas a fog node experiencing resource constraints may need to delegate tasks to the cloud. Balancing the workload between fog and cloud environments becomes challenging when taking into account multiple parameters. This motivates us

to investigate a new solution for distributing service workload effectively between fog and cloud computing platforms while maintaining quality of service. In this thesis, to overcome the aforementioned issue, we are adopting deep reinforcement learning.

1.2 Goal and Contributions

The goals of this paper are to enhance QoS with reduced service latency by delegating a portion of the service to the cloud, serve the maximum number of users in real-time by letting the cloud handle part of the service request, and efficiently utilize the fog node's limited resources by incorporating the cloud in the real-time service delivery process. The main contribution of this work can be summarized as follows:

- This paper outlines the problem of distributing users' services among fog and cloud environments;
- Proposed a new deep reinforcement learning (DRL) based method to distribute service requests to fog and cloud computing environments;
- The proposed DRL approach performance is evaluated in terms of metrics such as success rate, distribution of service requests, and etc.
- The results are compared with existing state-of-the-art methods.

1.3 Outline

This chapter provides an overview of IoT, cloud computing, fog computing, smart gateway, the significance of service distribution, and how deep reinforcement learning can help resolve the issue. Chapter 2 delves into the background of service distribution and reviews the current state of the art. In Chapter 3, we present our solution to the service distribution problem by discussing the architecture, applications, methods, and algorithms used. Chapter 4 details the experiment process and results are discussed in Chapter 5. Finally, Chapter 6 concludes the work with a summary and future direction.

2 State of the Art

In this section, we delve into the current state of the art in the field of service delivery. Our discussion begins by providing an overview of the background of service delivery in both fog and cloud computing environments in Section 2.1. This section aims to give the reader a comprehensive understanding of the context in which the current work is being carried out. Following this, in Section 2.2, we delve into the related works that have been done in the area of service dispersal methods within the cloud domain. This section aims to give the reader an understanding of the prior efforts-made in this field and the gaps that the current work aims to fill. Finally, in Section 2.3, we provide a summary of the key points covered in this section, offering a clear and concise overview of the state of the art in the field of service delivery.

2.1 Background

Recently, the Internet of Things (IoT) has made a significant impact on our society by converting everyday objects like wearables, transportation, and augmented reality into communicating devices, bringing both new challenges and opportunities. It is evident that the existing infrastructure will not be capable of handling the vast amount of data generated as the number of devices connected to the network is expected to reach over 50 billion by 2020, according to Cisco [Kim16]. The current Cloud infrastructure is inadequate to support a substantial amount of the existing IoT applications, due to three main reasons. First, the large volume of generated data makes it difficult to transfer it from the end-devices where it is created to the Cloud servers where it is processed. This is due to limitations in bandwidth, processing overhead, and transmission costs. Second, applications that require real-time analysis, such as online gaming and video applications, can be negatively impacted by the significant end-to-end delay from end-devices to the faraway Cloud servers. Lastly, privacy and security concerns may dictate that certain sensitive data must not traverse the entire Internet, making it necessary to process it closer to its source [YQL15]. A solution has been identified to address the problems mentioned above by avoiding network congestions, reducing communication costs, and minimizing the delay in data transfer. This new concept that leverages the advantages of the Cloud and the decentralization of service processing on edge devices is called Fog Computing [BMNZ14].

Fog Computing extends Cloud Computing by locating processing, analysis, and storage closer to the source of requests, at the edge of the network. [DGC⁺16] The goal is to minimize the data sent to the Cloud and lower the time delay and computational expenses. One of the main challenges of Fog Computing, a new paradigm, is the management of services, specifically the problem of determining the appropriate placement for them. Efficiently deploying services on available Fog nodes is a significant challenge facing the widespread adoption of Fog Computing. The difficulty in deploying services on Fog

nodes arises from the fact that they are geographically dispersed, have limited resources, and can change dynamically, adding to the complexity of the issue. Therefore, it is crucial to determine an efficient, effective, and equitable way of allocating resources to IoT applications in order to provide end-to-end guaranteed services to the users. Additionally, the focus may vary depending on the situation and the priorities, including aspects such as resource utilization [TD17], Quality of Service (QoS) [SNSD17, BF17], and Quality of Experience (QoE) [MSRB19]. Over the past few years, multiple efforts have been made to address this challenge. An efficient service placement has been proposed by taking into account different characteristics, assumptions, and strategies.

The literature has used various models and definitions to describe the applications produced by IoT devices and how they are processed by Fog resources and Cloud servers. According to the literature reviewed, there are three primary categories of IoT applications: 1) a single, integrated service, 2) multiple interdependent components, and 3) a connected graph. An application submitted by end-users or IoT devices can be represented as a single, standalone component, which is known as a monolithic service. This can include things like image processing applications or data instances that need to be processed or stored on a single physical node. The second type of application consists of multiple components that work together to perform a specific task. Each component performs a specific function within the application, but the relationships between them are not taken into account. In connected graph scenario, the application is comprised of components that are connected and interdependent, represented as a graph. The graph vertices symbolize the processing components of the application, while the edges depict the communication requirements and inter-dependencies between nodes. [SDL20]

2.2 Related Work

The placement of services that are effective in fog and cloud environments is a complex task that is influenced by several factors. The first of these is the heterogeneous and resource-constrained nature of most Fog nodes, which have limited capacities. The second factor is the dynamic nature of the environment, where resources can change over time and workload can be vary. Thirdly, the issue is made even more complicated by the widespread distribution of fog devices across a large-scale infrastructure. The problem of SPP (service placement problem) in a fog and cloud environment is a difficult task due to various specificities and constraints. Within the literature, we can identify the primary approaches employed to solve this problem as heuristic approach, machine learning based algorithms, and deep reinforcement learning based algorithms.

Wang and Li et. al [WL19] proposes a solution for improving task scheduling in fog computing using a hybrid heuristic algorithm. It performs better than other strategies in handling limited computing resources and high energy consumption in terminal devices. However, the study is limited because it doesn't address task clustering and fog node clustering. Zahoor et al. [ZJJ⁺18] proposes a model for managing resources in SGs using

cloud-fog computing, and uses several load balancing algorithms including HABACO to optimize response time. However, the proposed model needs to be extended to manage multiple load balancing applications and consider the scheduling of appliances and generation of microgrids. Rafique et al. [RSI⁺19] proposes a scheduling strategy to optimize resource utilization and reduce average response time in fog computing. The proposed approach outperforms other techniques in terms of energy consumption and execution time. However, it is required to apply reinforcement learning techniques to improve fog-IoT resource management. Yasmeen et al. [YJR⁺18] proposes a three-layered cloud and fog-based model to decrease consumer load and power generation system. Resource balancing is achieved through the use of three algorithms: Round Robin, throttled, and Particle Swarm Optimization with Simulated Annealing (PSOSA). However, the study is limited because it doesn't address examining regional fogs and investigating whether localized communication can reduce service latency. Yadav et al. [YNG19] proposes a hybrid algorithm that combines Genetic Algorithm and Particle Swarm Optimization to allocate services efficiently in a fog computing environment. The algorithm minimizes total makespan and energy consumption for IoT applications. However, the study's needs to be extended to investigate load balancing in multi-level fog infrastructure for IoT applications. Ren et al. [RZA21] suggests a model to decrease energy consumption and efficiently manage energy resources in fog systems. The proposed architecture aims to manage resources in a way that reduces losses and ensures quality of service. The study is lack of managing on cost improvement. Hosseinioun et al. [HKTG20] suggests a method for saving energy in the fog computing environment using the Dynamic Voltage and Frequency Scaling (DVFS) technique. A hybrid algorithm, IWO-CA, is used to create valid task sequences. Javanmardi et al. [JSPP21] The article describes the FPFTS fog task scheduler, which uses particle swarm optimization and fuzzy theory to optimize application loop delay and network utilization. The drawback of the study is that it requires developing a mobility-aware task scheduling algorithm. Xu et al. [XHZS19] proposes the LBP-ACS algorithm for task scheduling in cloud-fog environments, which considers priority constraints, energy consumption, and mixed deadlines. The algorithm reduces energy consumption and failure rates while ensuring reasonable scheduling length. However, the study doesn't address the scheduling of independent and associated tasks. Djemai et al. [DSMP19] proposes a strategy for IoT services placement in Fog architecture using a Discrete Particles Swarm Optimization algorithm. The placement strategy minimizes application delay violations and considers energy consumption. All previous related works use heuristic-based solutions, which lack performance guarantees. Also, our study mainly focuses on service delivery by fog and cloud computing using DRL-based approach.

Rahbari and Nickray et al. [RN20] presents a Module Placement method by Classification and regression tree Algorithm (MPCA) for selecting the best FDs for modules. MPCA selects the best FDs based on decision parameters such as authentication, con-

confidentiality, integrity, availability, capacity, speed, and cost. Bashir et al. [BLK22] proposes a dynamic resource allocation strategy for cloud, fog nodes, and users. The strategy uses TOPSIS (technique for order performance by similarity to ideal solution) to rank the fog nodes and logistic regression to calculate the load of individual fog nodes. Ferrández-Pastor et al. [FPMJMV18] proposes a method to design smart services based on the edge computing paradigm. The approach addresses interoperability and scalability issues of existing designs and implements energy management, security system, climate control, and information services subsystems on embedded devices. He et al. [HWC⁺17] proposes a new fog computing model with functional modules can reduce problems related to dedicated computing infrastructure and slow response times in cloud computing. Experiments were conducted to measure the performance of analytics tasks on fogs formed by Raspberry Pi computers using a distributed computing engine. Quality of services (QoS) aware admission control, offloading, and resource allocation schemes were designed to support data analytics services, and a scalable system-level simulator was developed to evaluate the fog-based analytics service and the QoS management schemes. Fei et al. [FSV⁺19] investigates machine learning techniques for CPS data stream analytics in a Cloud and Fog architecture. It provides guidance on how these methods should be deployed and presents the challenges associated with using them. This study is the first of its kind to systematically explore this topic. Priyabhashana et al. [PJ19] compares the use of various activation functions on a TensorFlow-based neural network model built using the Keras library. The proposed fog station was evaluated, and the experiment results showed the feasibility, efficiency, and applicability of the system. Alsaffar et al. [APH⁺16] proposes an architecture for IoT service delegation and resource allocation that combines fog and cloud computing. The authors introduce a decision algorithm based on linearized decision trees to manage and delegate user requests and balance workload based on services size, completion time, and VMs capacity. They also propose an algorithm to allocate resources to meet SLA and QoS requirements and optimize big data distribution. Selimi et al. [SCAF⁺19] proposes a fast heuristic algorithm that leverages state information about the network to inform service placement decisions, and demonstrates that it consistently outperforms random placement by 2x in bandwidth gain. The proposed method is evaluated using a real live video-streaming service and a popular Web 2.0 service, and results show significant improvements in packet loss rate and client response times, which ultimately translates to higher Quality of Experience (QoE) for the users. Yadav et al. [YMY20] presents a new technique for task allocation that minimizes response time and system cost using clustering. The proposed technique uses Fuzzy C-Means clustering and Hungarian method for task allocations, and the performance is evaluated and compared with other models.

Farhat et al. [FSM20] proposes R-learning model that aims to decrease the cloud's load by utilizing available fog resources in different locations. The model is demonstrating its ability to adapt to user demand and efficiently use fog resources. Orhean et al.

[OPR18] presents a reinforcement learning algorithm for solving the scheduling problem in distributed systems. The proposed algorithm considers the heterogeneity of nodes and the dependency graph of tasks for determining a scheduling policy that minimizes execution time. The authors also propose a platform that implements the algorithm and provides scheduling as a service to distributed systems. Tang et al. [TZZ⁺18] proposes a new architecture and algorithms for container migration to support mobility tasks with different application requirements. The proposed algorithms consider hosting mobile application tasks in containers of corresponding fog nodes, and modeling container migration strategies as multiple dimensional Markov Decision Process spaces. They used deep reinforcement learning algorithms in order to reduce the large MDP spaces. Gazori et al. [GRN20] proposes a Double Deep Q-Learning (DDQL)-based scheduling algorithm for task scheduling of fog-based IoT applications. The aim is to minimize long-term service delay and computation cost under resource and deadline constraints. The reinforcement learning approach, combined with the target network and experience replay techniques, is used to address this problem. Alam et al. [AHU⁺19] proposes a deep Q-learning based autonomic management framework for computation offloading in mobile edge/fog. The framework uses a distributed edge/fog network controller to allocate resources and minimize latency of service computing while maintaining energy efficiency. Zhang et al. [ZLY⁺18] proposes a double deep Q-learning model for energy-efficient edge scheduling (DDQ-EES). The model includes a generated network and a target network for producing Q-values, and rectified linear units (ReLU) as the activation function to avoid gradient vanishing. Lu et al. [LHD⁺20] proposes a QoE model for computation offloading that considers service latency, energy consumption, and task success rate. They also present an improved DDPG algorithm called double-dueling-deterministic policy gradients (D3PG) that utilizes Double Q-learning and Dueling networks. Wang et al. [WLLC19] proposes a Q-learning-based hierarchical service tree placement strategy to optimize the net utility in large networks with complex service structures. The proposed method recursively decouples a service tree into sub-trees each with a single computing sub-task and its associated data flows. Dehury and Srirama et al. [DS19, DS20] a reinforcement learning-based (RLPSD) and a deep reinforcement learning-based (DRLSD-FC) service delivery mechanism that distributes users' service requests between fog and cloud environments to minimize the workload on fog while maintaining service quality. The algorithms consider user constraints such as distance from fog, service sensitivity, and other metrics. However, the study is limited because it does not address service-placement in real-time manner as our study, and service dispersal mechanism inside of fog nodes, in our study we use smart gateway in order to do service dispersal jobs.

2.3 Summary

In this chapter, we first introduced with background information on service placement in fog and cloud environments. Then we discussed the current state-of-the-art in service dispersial methods for both environments. We also saw the different approaches used to answer the SP problem in related works section. Finally, we saw what is new in our studies compared to the current state-of-the-art. In the next chapter, we will see our approach to the problem of service placement in fog and cloud environments.

3 Methodology

Multiple methods can be used to solve the same issue. In this section, we outline our approach to addressing the challenge of distributing services between fog and cloud environment. We start our discussion by modeling the services and mathematical representation and analysis of key parameters in Section 3.1, followed by a deep reinforcement learning, covering topics like state space, action space, and reward calculation in Section 3.2. Section 3.3 delves into the selected applications and their architectures for our study. Finally, we wrap up this section with a summary in Section 3.4.

3.1 Modelling Services

In Figure 2, the fog and cloud model offers multiple services to users. Users can obtain services by submitting a Service Request (SR) to the smart gateway. The Service Request can be partitioned, which is an assumption that can be demonstrated with various real world examples involving data intensive applications generated by IoT devices, distributed among the fog and cloud environments. The smart gateway collaborates with both fog and cloud environments to meet service resource demands. The original versions of applications were modified to satisfy the problem formulation, and the revised editions used in this section are available in the project repository. This section explains the services available in the fog and cloud environment.

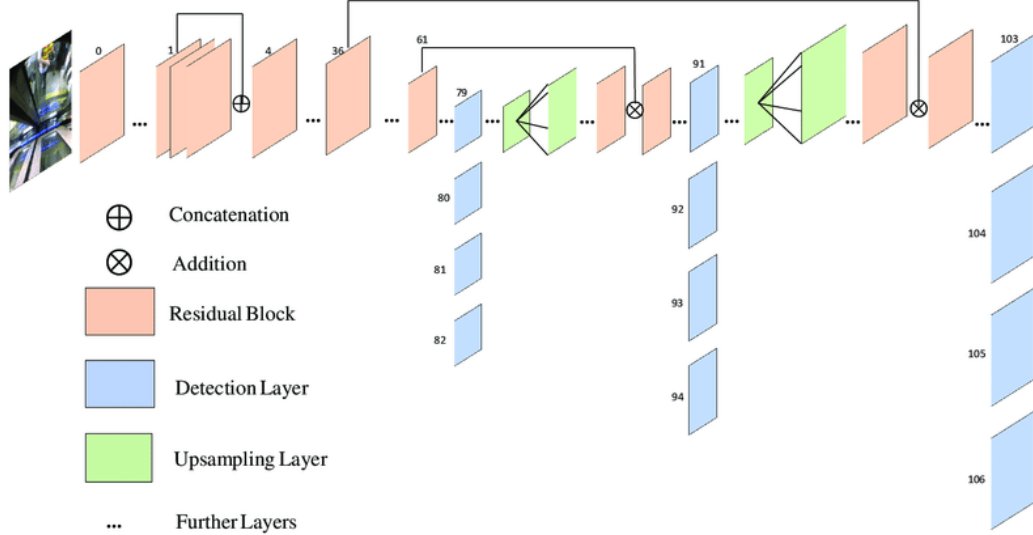


Figure 3. Network architecture of YOLOv3 [DLLL20].

3.1.1 Application 1: Deep learning based object classification

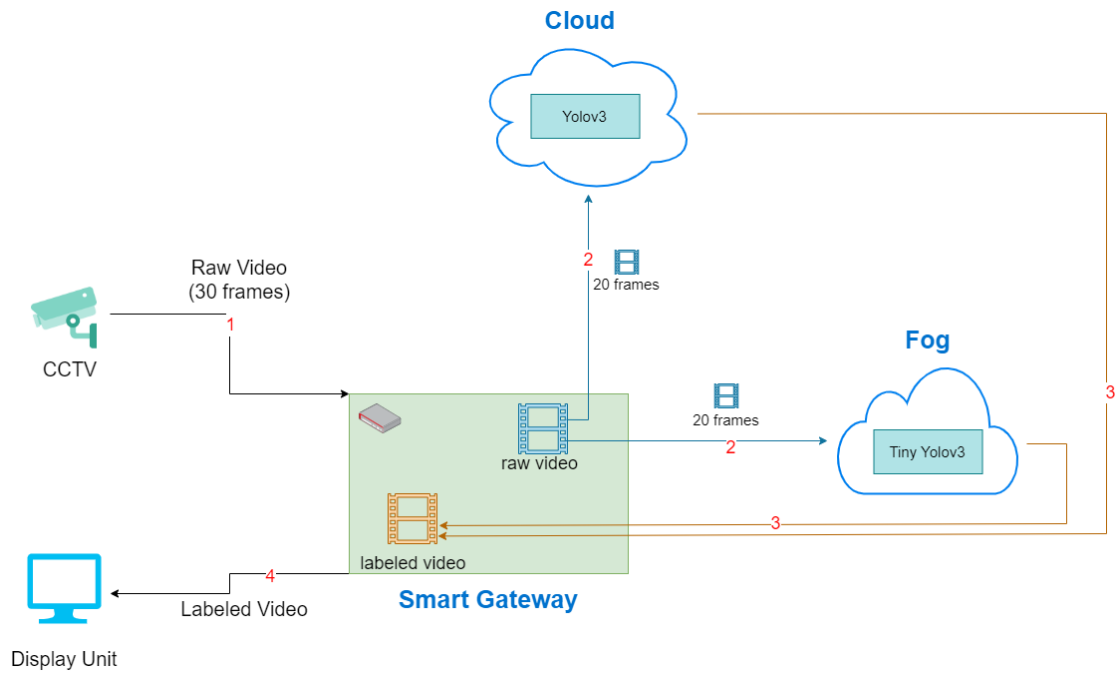


Figure 4. Architecture of Application 1: Deep learning based object classification

YOLOv3 (You Only Look Once) is real-time object detection algorithm that uses a single convolutional neural network (CNN) to detect objects in an image. It works by dividing an input image into a grid of cells, and predicting bounding boxes and class probabilities for each cell. Each bounding box consists of four coordinates (x, y, width, height) and a confidence score, which represents the likelihood that the box contains an object. Class probabilities are predicted for each box, indicating the likelihood that the object in the box belongs to a certain class. To improve the accuracy of object detection, YOLOv3 uses a few key techniques. One is feature pyramid networks, which use a pyramid of features extracted from multiple layers of the network to improve the detection of objects at different scales as shown in Figure 3. Another technique is residual connections, which allow gradients to flow more easily through the network during training, leading to better accuracy.

This application is a suitable choice for Fog environment, it can be deployed to detect objects in order to reduce communication delays in fog computing. The Algorithm 1 represents how Yolov3 used in this application. Additionally, The Figure 4 illustrates an architecture of deep learning based object classification application.

Algorithm 1: Deep learning based object classification

Input: Video to be analyzed

Result: Detected objects in the video

- 1 Load the pre-trained YOLOv3 object detection model;
 - 2 **foreach** *frame in video* **do**
 - 3 Extract image from frame;
 - 4 Perform object detection using YOLOv3;
 - 5 Draw bounding boxes around detected objects;
 - 6 Add labels to bounding boxes;
 - 7 Insert the annotated image back into the frame;
 - 8 Add the frame to the output video;
-

3.1.2 Application 2: Text-audio synchronisation or forced alignment

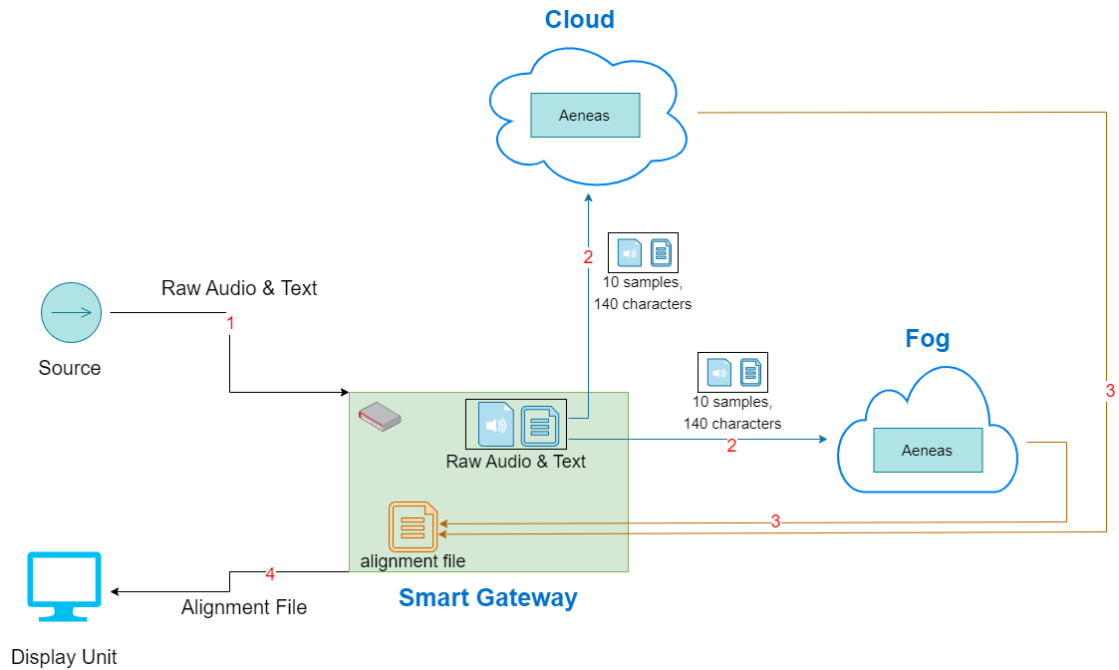


Figure 5. Architecture of Application 2: Text-audio synchronisation or forced alignment

Aeneas is a text-alignment library that aligns text with audio or video content. As shown in Algorithm 2, the library create sync map, allowing users to jump to the corresponding point in the audio or video when clicking on a specific word or phrase in the text. Aeneas works by first converting the audio or video into a time-aligned phonetic

transcript, known as an "alignment graph." It then generates a time-aligned version of the text, known as a "fragmented transcript," based on the words and timing information in the alignment graph. Next, Aeneas performs a process called "forced alignment," which uses statistical models to align the transcript to the audio or video. This process allows for a precise time-based synchronization of the text and media content. The resulting output of Aeneas can be an XML, JSON, SRT files that includes a detailed mapping of the text and media content, including time codes and other metadata. The Figure 5 illustrates an architecture of text-audio synchronisation or forced alignment application.

Algorithm 2: Text-audio synchronisation or forced alignment

Input: Text and Audio Files

Output: Alignment File (*.srt)

```

1 Procedure align(TextFile, AudioFile)
    /* Initialize Aeneas object */
2   aeneas = Aeneas(language="en")
    /* Create the sync map */
3   syncmap = SyncMap()
    /* Perform the alignment */
4   aeneas.execute(text=TextFile, audio=AudioFile, sync=syncmap)
    /* Save the sync map as an alignment file */
5   syncmap.writetofile("alignment.srt")
6   Return the saved alignment file as an output;
```

3.1.3 Application 3: Speech-to-text conversion

Algorithm 3: Speech-to-text conversion

Input: Audio file to be transcribed

Result: Transcription of the audio file

```

1 Load the audio file;
2 Initialize the PocketSphinx recognizer;
3 Set the language model and dictionary;
4 Process the audio file in chunks;
5 Feed the chunks to the recognizer;
6 Get the transcription from the recognizer;
7 Output the transcription;
```

PocketSphinx is a lightweight speech recognition engine. It is based on the Sphinx-4 open source speech recognition system, which uses Hidden Markov Models (HMMs) to model the acoustic features of speech. PocketSphinx is capable of recognizing speech in a variety of languages, and can be used for tasks such as voice search, dictation, and command recognition. The PocketSphinx library provides a set of Python APIs for developers to use in their applications. As shown in Algorithm 3, to use PocketSphinx, firstly start by configuring a decoder, which specifies the acoustic and language models to be used for speech recognition. The audio input is then passed to the decoder, which performs a series of signal processing and feature extraction operations to produce a stream of speech recognition hypotheses. These hypotheses are then scored and ranked according to their likelihood of being correct, and the most likely hypothesis is returned as the recognized text. The Figure 6 illustrates an architecture of speech-to-text conversion application.

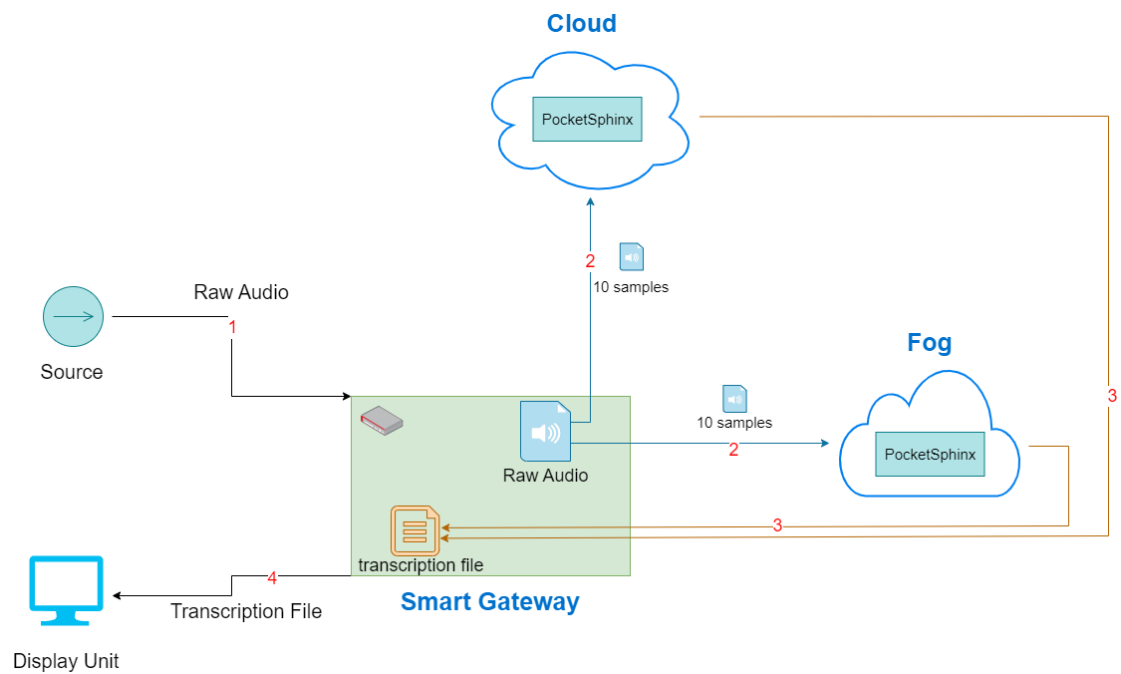


Figure 6. Architecture of Application 3: Speech-to-text conversion

3.2 Main Parameters of the System

3.2.1 Environment Parameters

3.2.2 User Parameters

3.2.3 Service Parameters

3.2.4 Priority Parameters

3.3 DRL-based Service Placement

3.3.1 State Space

3.3.2 Action Space

3.3.3 Rewards

3.4 Summary

4 Experiment

4.1 Experimental Setup

4.2 Summary

5 Results

5.1 Summary

6 Conclusion

References

- [AH14] Mohammad Aazam and Eui-Nam Huh. Fog computing and smart gateway based communication for cloud of things. In *2014 International conference on future internet of things and cloud*, pages 464–470. IEEE, 2014.
- [AHU⁺19] Md Golam Rabiul Alam, Mohammad Mehedi Hassan, Md Zia Uddin, Ahmad Almogren, and Giancarlo Fortino. Autonomic computation of-flooding in mobile edge for iot applications. *Future Generation Computer Systems*, 90:149–157, 2019.
- [APH⁺16] Aymen Abdullah Alsaffar, Hung Phuoc Pham, Choong-Seon Hong, Eui-Nam Huh, and Mohammad Aazam. An architecture of iot service delegation and resource allocation based on collaboration between fog and cloud computing. *Mobile Information Systems*, 2016, 2016.
- [BF17] Antonio Brogi and Stefano Forti. Qos-aware deployment of iot applications through the fog. *IEEE internet of Things Journal*, 4(5):1185–1192, 2017.
- [BLK22] Hayat Bashir, Seonah Lee, and Kyong Hoon Kim. Resource allocation through logistic regression and multicriteria decision making method in iot fog computing. *Transactions on Emerging Telecommunications Technologies*, 33(2):e3824, 2022.
- [BMNZ14] Flavio Bonomi, Rodolfo Milito, Preethi Natarajan, and Jiang Zhu. Fog computing: A platform for internet of things and analytics. *Big data and internet of things: A roadmap for smart environments*, pages 169–186, 2014.
- [DGC⁺16] Amir Vahid Dastjerdi, Harshit Gupta, Rodrigo N Calheiros, Soumya K Ghosh, and Rajkumar Buyya. Fog computing: Principles, architectures, and applications. In *Internet of things*, pages 61–75. Elsevier, 2016.
- [DLL⁺16] Ruilong Deng, Rongxing Lu, Chengzhe Lai, Tom H. Luan, and Hao Liang. Optimal workload allocation in fog-cloud computing toward balanced delay and power consumption. *IEEE Internet of Things Journal*, 3(6):1171–1181, 2016.
- [DLLL15] Ruilong Deng, Rongxing Lu, Chengzhe Lai, and Tom H. Luan. Towards power consumption-delay tradeoff by workload allocation in cloud-fog computing. In *2015 IEEE International Conference on Communications (ICC)*, pages 3909–3914, 2015.

- [DLLL20] Yuan Dai, Weiming Liu, Haiyu Li, and Lan Liu. Efficient foreign object detection between psds and metro doors via deep neural networks. *IEEE Access*, 8:46723–46734, 2020.
- [DS19] Chinmaya Kumar Dehury and Satish Narayana Srirama. Personalized service delivery using reinforcement learning in fog and cloud environment. In *Proceedings of the 21st International Conference on Information Integration and Web-Based Applications & Services*, pages 522–529, 2019.
- [DS20] Chinmaya Kumar Dehury and Satish Narayana Srirama. An efficient service dispersal mechanism for fog and cloud computing using deep reinforcement learning, 2020.
- [DSMP19] Tanissia Djemai, Patricia Stolf, Thierry Monteil, and Jean-Marc Pierson. A discrete particle swarm optimization approach for energy-efficient iot services placement over fog infrastructures. In *2019 18th International Symposium on Parallel and Distributed Computing (ISPDC)*, pages 32–40. IEEE, 2019.
- [Ela19] Hanan Elazhary. Internet of things (iot), mobile cloud, cloudlet, mobile iot, iot cloud, fog, mobile edge, and edge emerging computing paradigms: Disambiguation and research directions. *Journal of Network and Computer Applications*, 128:105–140, 2019.
- [FPMJMV18] Francisco-Javier Ferrández-Pastor, Higinio Mora, Antonio Jimeno-Morenilla, and Bruno Volckaert. Deployment of iot edge and fog computing technologies to develop smart building services. *Sustainability*, 10(11):3832, 2018.
- [FSM20] Peter Farhat, Hani Sami, and Azzam Mourad. Reinforcement r-learning model for time scheduling of on-demand fog placement. *the Journal of Supercomputing*, 76:388–410, 2020.
- [FSV⁺19] Xiang Fei, Nazaraf Shah, Nandor Verba, Kuo-Ming Chao, Victor Sanchez-Anguix, Jacek Lewandowski, Anne James, and Zahid Usman. Cps data streams analytics based on machine learning for cloud and fog computing: A survey. *Future generation computer systems*, 90:435–450, 2019.
- [GRN20] Pegah Gazori, Dadmehr Rahbari, and Mohsen Nickray. Saving time and cost on the scheduling of fog-based iot applications using deep reinforcement learning approach. *Future Generation Computer Systems*, 110:1098–1115, 2020.

- [HKTG20] Pejman Hosseinioun, Maryam Kheirabadi, Seyed Reza Kamel Tabbakh, and Reza Ghaemi. A new energy-aware tasks scheduling approach in fog computing using hybrid meta-heuristic algorithm. *Journal of Parallel and Distributed Computing*, 143:88–96, 2020.
- [HWC⁺17] Jianhua He, Jian Wei, Kai Chen, Zuoyin Tang, Yi Zhou, and Yan Zhang. Multitier fog computing with large-scale iot data analytics for smart cities. *IEEE Internet of Things Journal*, 5(2):677–686, 2017.
- [JSPP21] Saeed Javanmardi, Mohammad Shojafar, Valerio Persico, and Antonio Pescapè. Fpfts: A joint fuzzy particle swarm optimization mobility-aware approach to fog task scheduling algorithm for internet of things devices. *Software: practice and experience*, 51(12):2519–2539, 2021.
- [Kim16] HS Kim. Fog computing and the internet of things: extend the cloud to where the things are. *Int. J. Cisco*, 2016.
- [LGL⁺15] Tom H Luan, Longxiang Gao, Zhi Li, Yang Xiang, Guiyi Wei, and Limin Sun. Fog computing: Focusing on mobile users at the edge. *arXiv preprint arXiv:1502.01815*, 2015.
- [LHD⁺20] Haodong Lu, Xiaoming He, Miao Du, Xiukai Ruan, Yanfei Sun, and Kun Wang. Edge qoe: Computation offloading with deep reinforcement learning for internet of things. *IEEE Internet of Things Journal*, 7(10):9255–9265, 2020.
- [MSRB19] Redowan Mahmud, Satish Narayana Srirama, Kotagiri Ramamohanarao, and Rajkumar Buyya. Quality of experience (qoe)-aware placement of applications in fog computing environments. *Journal of Parallel and Distributed Computing*, 132:190–203, 2019.
- [MWT⁺19] Jonathan McChesney, Nan Wang, Ashish Tanwer, Eyal De Lara, and Blesson Varghese. Defog: fog computing benchmarks. In *Proceedings of the 4th ACM/IEEE Symposium on Edge Computing*, pages 47–58, 2019.
- [OPR18] Alexandru Iulian Orhean, Florin Pop, and Ioan Raicu. New scheduling approach using reinforcement learning for heterogeneous distributed systems. *Journal of Parallel and Distributed Computing*, 117:292–302, 2018.
- [PJ19] HMB Priyabhashana and KPN Jayasena. Data analytics with deep neural networks in fog computing using tensorflow and google cloud platform.

In *2019 14th Conference on Industrial and Information Systems (ICIIS)*, pages 34–39. IEEE, 2019.

- [RN20] Dadmehr Rahbari and Mohsen Nickray. Task offloading in mobile fog computing by classification and regression tree. *Peer-to-Peer Networking and Applications*, 13:104–122, 2020.
- [RSI⁺19] Hina Rafique, Munam Ali Shah, Saif Ul Islam, Tahir Maqsood, Suleman Khan, and Carsten Maple. A novel bio-inspired hybrid algorithm (nbiha) for efficient resource management in fog computing. *IEEE Access*, 7:115760–115773, 2019.
- [RZA21] Xiaojun Ren, Zhijun Zhang, and Seyedeh Maryam Arefzadeh. An energy-aware approach for resource managing in the fog-based internet of things using a hybrid algorithm. *International Journal of Communication Systems*, 34(1):e4652, 2021.
- [SCAF⁺19] Mennan Selimi, Llorenç Cerdà-Alabern, Felix Freitag, Luís Veiga, Arjuna Sathiaselalan, and Jon Crowcroft. A lightweight service placement approach for community network micro-clouds. *Journal of Grid Computing*, 17:169–189, 2019.
- [SDL20] Farah Ait Salaht, Frédéric Desprez, and Adrien Lebre. An overview of service placement problem in fog and edge computing. *ACM Computing Surveys (CSUR)*, 53(3):1–35, 2020.
- [SDV18] Prasan Kumar Sahoo, Chinmaya Kumar Dehury, and Bharadwaj Veeravalli. Lvrn: On the design of efficient link based virtual resource management algorithm for cloud platforms. *IEEE Transactions on Parallel and Distributed Systems*, 29(4):887–900, 2018.
- [SMBMT⁺18] V.B. Souza, X. Masip-Bruin, E. Marín-Tordera, S. Sánchez-López, J. Garcia, G.J. Ren, A. Jukan, and A. Juan Ferrer. Towards a proper service placement in combined fog-to-cloud (f2c) architectures. *Future Generation Computer Systems*, 87:1–15, 2018.
- [SMW18] Hamed Shah-Mansouri and Vincent W. S. Wong. Hierarchical fog-cloud computing for iot systems: A computation offloading game. *IEEE Internet of Things Journal*, 5(4):3246–3257, 2018.
- [SNSD17] Olena Skarlat, Matteo Nardelli, Stefan Schulte, and Schahram Dustdar. Towards qos-aware fog service placement. In *2017 IEEE 1st international conference on Fog and Edge Computing (ICFEC)*, pages 89–96. IEEE, 2017.

- [SWHL16] Ivan Stojmenovic, Sheng Wen, Xinyi Huang, and Hao Luan. An overview of fog computing and its security issues. *Concurrency and Computation: Practice and Experience*, 28(10):2991–3005, 2016.
- [TD17] Mohit Taneja and Alan Davy. Resource aware placement of iot application modules in fog-cloud computing paradigm. In *2017 IFIP/IEEE Symposium on Integrated Network and Service Management (IM)*, pages 1222–1228. IEEE, 2017.
- [TZZ⁺18] Zhiqing Tang, Xiaojie Zhou, Fuming Zhang, Weijia Jia, and Wei Zhao. Migration modeling and learning algorithms for containers in fog computing. *IEEE Transactions on Services Computing*, 12(5):712–725, 2018.
- [WL19] Juan Wang and Di Li. Task scheduling based on a hybrid heuristic algorithm for smart production line with fog computing. *Sensors*, 19(5):1023, 2019.
- [WLLC19] Yimeng Wang, Yongbo Li, Tian Lan, and Nakjung Choi. A reinforcement learning approach for online service tree placement in edge computing. In *2019 IEEE 27th International Conference on Network Protocols (ICNP)*, pages 1–6. IEEE, 2019.
- [XHZS19] Jiuyun Xu, Zhuangyuan Hao, Ruru Zhang, and Xiaoting Sun. A method based on the combination of laxity and ant colony system for cloud-fog task scheduling. *IEEE Access*, 7:116218–116226, 2019.
- [YHQL15] Shanhe Yi, Zijiang Hao, Zhengrui Qin, and Qun Li. Fog computing: Platform and applications. In *2015 Third IEEE workshop on hot topics in web systems and technologies (HotWeb)*, pages 73–78. IEEE, 2015.
- [YJR⁺18] Anila Yasmeen, Nadeem Javaid, Obaid Ur Rehman, Hina Iftikhar, Muhammad Faizan Malik, and Fatima J Muhammad. Efficient resource provisioning for smart buildings utilizing fog and cloud based environment. In *2018 14th International Wireless Communications & Mobile Computing Conference (IWCMC)*, pages 811–816. IEEE, 2018.
- [YMY20] Seema Yadav, Rakesh Mohan, and Pradeep Kumar Yadav. Task allocation model for optimal system cost using fuzzy c-means clustering technique in distributed system. *Ingénierie des Systèmes d’Inf.*, 25(1):59–68, 2020.

- [YNG19] Vinita Yadav, BV Natesha, and Ram Mohana Reddy Guddeti. Gaps: Service allocation in fog computing environment using hybrid bio-inspired algorithm. In *TENCON 2019-2019 IEEE Region 10 Conference (TENCON)*, pages 1280–1285. IEEE, 2019.
- [YQL15] Shanhe Yi, Zhengrui Qin, and Qun Li. Security and privacy issues of fog computing: A survey. In *Wireless Algorithms, Systems, and Applications: 10th International Conference, WASA 2015, Qufu, China, August 10-12, 2015, Proceedings 10*, pages 685–695. Springer, 2015.
- [YXZ18] Ruozhou Yu, Guoliang Xue, and Xiang Zhang. Application provisioning in fog computing-enabled internet-of-things: A network perspective. In *IEEE INFOCOM 2018 - IEEE Conference on Computer Communications*, pages 783–791, 2018.
- [ZJJ⁺18] Saman Zahoor, Sakeena Javaid, Nadeem Javaid, Mahmood Ashraf, Faruh Ishmanov, and Muhammad Khalil Afzal. Cloud–fog–based smart grid model for efficient resource management. *Sustainability*, 10(6):2079, 2018.
- [ZLY⁺18] Qingchen Zhang, Man Lin, Laurence T Yang, Zhikui Chen, Samee U Khan, and Peng Li. A double deep q-learning model for energy-efficient edge scheduling. *IEEE Transactions on Services Computing*, 12(5):739–749, 2018.

Appendix

I. List of Figures

1	Abstract view of Smart Gateway with Cloud-Fog architecture.	9
2	An example motivational scenario.	10
3	Network architecture of YOLOv3 [DLLL20].	18
4	Architecture of Application 1: Deep learning based object classification	19
5	Architecture of Application 2: Text-audio synchronisation or forced alignment	20
6	Architecture of Application 3: Speech-to-text conversion	22

II. List of Tables

III. Source Code

The source code of applications and developed algorithms for this study is located in repository <https://github.com/chinmaya-dehury/AI4FogCloudServiceDisperse>. The access to the repository could be granted by sending an email to chinmaya.dehury@ut.ee

IV. Licence

Non-exclusive licence to reproduce thesis and make thesis public

I, **Jeyhun Abbasov**,

1. herewith grant the University of Tartu a free permit (non-exclusive licence) to reproduce, for the purpose of preservation, including for adding to the DSpace digital archives until the expiry of the term of copyright,

AI-powered cloud usage controlling system in smart home environment,
supervised by Prof. Dr. Chinmaya Dehury.

2. I grant the University of Tartu a permit to make the work specified in p. 1 available to the public via the web environment of the University of Tartu, including via the DSpace digital archives, under the Creative Commons licence CC BY NC ND 3.0, which allows, by giving appropriate credit to the author, to reproduce, distribute the work and communicate it to the public, and prohibits the creation of derivative works and any commercial use of the work until the expiry of the term of copyright.
3. I am aware of the fact that the author retains the rights specified in p. 1 and 2.
4. I certify that granting the non-exclusive licence does not infringe other persons' intellectual property rights or rights arising from the personal data protection legislation.

Jeyhun Abbasov

dd/mm/yyyy