



**Università degli Studi di Roma “Tor Vergata”**

---

**FACOLTÀ DI INGEGNERIA**

**Corso di Laurea Magistrale in Ingegneria Informatica**

**Performance Modeling of  
Computer Systems and Networks  
Project**

Studente:

**Andrea Graziani**

**Matricola 0273395**

Docente:

**Prof. Vittoria de Nitto Personè**

# Contents

<b>1</b>	<b>Introduction</b>	<b>2</b>
<b>2</b>	<b>Goals</b>	<b>2</b>
<b>3</b>	<b>Conceptual Model</b>	<b>3</b>
<b>4</b>	<b>Specification Model</b>	<b>6</b>
4.1	System's variables . . . . .	6
4.2	Statistics computing . . . . .	7
4.3	System's assumptions . . . . .	8
4.4	System's events and assumptions . . . . .	9
<b>5</b>	<b>Computational Model</b>	<b>11</b>
5.1	Events implementation . . . . .	11
5.2	Event list . . . . .	13
5.3	Simulation clock . . . . .	13
5.4	Next-event simulation logic . . . . .	14
5.5	System network implementation . . . . .	16
5.6	Output analysis . . . . .	17
<b>6</b>	<b>Analytical solution</b>	<b>18</b>
6.1	System based on access control Algorithm 1 . . . . .	18
6.1.1	Balance equation . . . . .	18
6.1.2	Probabilities computing . . . . .	20
6.1.3	Average arrival rates . . . . .	20
6.1.4	Average population . . . . .	21
6.1.5	Average response time . . . . .	22
6.1.6	Throughput . . . . .	23
6.2	System based on access control Algorithm 2 . . . . .	25
6.2.1	Balance equation . . . . .	25
6.2.2	Probabilities computing . . . . .	26
6.2.3	Average arrival rates . . . . .	27
6.2.4	Other metrics . . . . .	28
6.3	Summary of analytical results . . . . .	29

# 1 Introduction

In this report we will present how we have designed and implemented a discrete-event simulator for required system using the next-event approach.

In order to do that, we have strictly followed each step required by typical discrete-event simulation model<sup>1</sup> whose main results are reported in each sections of this document.

# 2 Goals

Our simulation model's goal is to capture some relevant global and local characteristics, including time-averaged response time, population and throughput for each class of jobs, through the design and implementation of two version of system based on two different routing policies, called *access control algorithms*.

Obviously we will compare results obtained by our simulations when based on different routing policies and we will care about to derive an analytical solution in order to validate that results not before to have properly define a queueing model.

---

<sup>1</sup>Lawrence M. Leemis, Stephen K. Park, *Discrete-Event Simulation: A First Course* (Pearson; 1 edition January 6, 2006), Algorithm 1.1.1, page 4

### 3 Conceptual Model

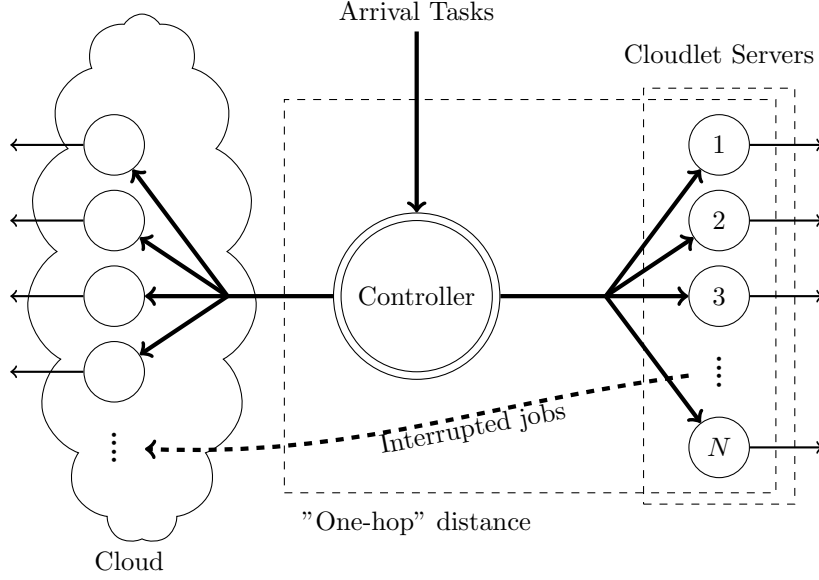


Figure 1: System diagram.

As known, building a system *conceptual model* means to describe it highlighting which state variables are important specifying how they are interrelated.<sup>2</sup>

Our system represents a network of nodes to which some application, running in a set of external mobile devices, send their tasks, called *jobs* or *task*, because of performance or energy saving reasons. Not all jobs sent to system are the same, in fact we distinguish two class type of jobs, that is *class 1* and *class 2* type jobs.

To minimize mean response time experienced by mobile device users, the system is been built as a *two-layer* network in such a way that any arriving job is sent, if possible, to the nearest node to be run.

To be more precise, as you can see from figure 1, our system is made up of a set of nodes including:

**Cloudlet node** Also called *edge cloud server*, it is a computer cluster located at an "*one-hop*" distance from mobile devices users, that is very close to them, therefore network delay experienced is very low. That cluster is able to guarantee absence of interferences among tasks allocated to it as long as their number does not exceed a given threshold  $N$ .

Consequently, from our conceptual model point of view, cloudlet node represents a *fixed-capacity multi-server service node with no queue* which capacity, that is the maximum possible number of jobs in it, is fixed to  $N$ .

<sup>2</sup>*Ibid.*

Since jobs inter-arrival times and service times are *exponentially* distributed according to our project specification, a queuing model able to properly represent that node is the M/M/N/N queueing model system<sup>3</sup> although there is a very important different with it: due to adopted routing policies, when cloudlet node is "full" and so not capable to run a further job, any arriving job is forwarded somewhere else rather than to cloudlet node, therefore no jobs are really lost.

**Cloud node** Also called *remote cloud server*, it is another computer cluster which, although it *suffers for greater network delay*, due of its high distance from mobile devices, being virtually made up of unlimited number of servers, has *unlimited resources*, therefore it can to process any number of tasks allocated to it.

Consequently we had modelled this system's node as an *unlimited-capacity multi-server service node with no queue*.

Since, as we previously said, jobs inter-arrival and service times are exponential, due to the presence of an infinite number of server, such a system's node can be modelled by a M/M/ $\infty$  queueing model system.<sup>4</sup>

**Controller node** It represents system's entry point being all arriving jobs sent firstly to it by mobile devices. Its duty is to decide about whether an arrived job should be sent to cloudlet or cloud node according to adopted routing policy.

It is located at the same distance of cloudlet node from mobile devices users, therefore, also in this case, network delay experienced is very low. Although it can run only one job at a time, it has the capability to process each arriving job *instantaneously*.

We have simply modelled it as *single-server service node with no queue*. It's correspondent queuing model system is M/M/1.

---

#### Algorithm 1

---

```

1: function ACCESSCONTROLFUNCTION(arrivalJob)
2:   if ( $n_1 + n_2 = N$ ) then
3:     Send arrivalJob on the cloud.
4:   else
5:     Send arrivalJob on the cloudlet.
```

---

At this point, in order to properly define our conceptual model, we need to specify which *system state variables* are truly important for our goals.

A complete characterization of the our system state is simply provided by the number of each class jobs running both on cloudlet and cloud node at a given instant of time; we will formalize this statement in the next section where we will provide relations and computation methods too.

For now we have to just focus on routing policies. As we previously said, there are two type of access control algorithms used by controller node to make

---

<sup>3</sup>Mor Harchol-Balter, *Performance Modeling and Design of Computer Systems* (Carnegie Mellon University, Pennsylvania 2013), 255-256

<sup>4</sup>*Ibid.*, 271

decisions; they are fully described by pseudo-code shown in Algorithm 1 and Algorithm 2. Note that you can use only one access control algorithm type during a simulation run.

Although they are both *work-conserving scheduling policy*, that is they always perform work on some job when there is a job in the system,<sup>5</sup> there is one difference between them: access control algorithms 1 does not provide job pre-emption once it starts service on cloudlet node unlike access control algorithms 2, which provides, instead, pre-emption when some condition are achieved requiring also a set-up time to restart interrupted job (note that Algorithm 2 use a further threshold variable to make its decision denoted by  $S$  which is less or equal than  $N$ ).

---

**Algorithm 2**

---

```

function ACCESSCONTROLFUNCTION(arrivalJob)
  if (arrivalJob.isClassOne) then
    if  $n_1 = N$  then
      Send arrivalJob on the cloud.
    else if  $n_1 + n_2 \leq S$  then
      Send arrivalJob on the cloudlet.
    else if  $n_2 \geq 0$  then
      Interrupt a class 2 job currently running on cloudlet.
      Send interrupted job to cloud.
      Send arrivalJob on the cloudlet.
    else
      Send arrivalJob on the cloudlet.
  else
    if ( $n_1 + n_2 \geq S$ ) then
      Send arrivalJob on the cloud.
    else
      Send arrivalJob on the cloudlet.

```

---



---

<sup>5</sup>*Ibid.*, Definition 28.1, page 474

## 4 Specification Model

In this section we will provide a *specification model* of our system in which we will turn all system's states into a collection of mathematical variables together with equations and logic describing how the state variables are interrelated including algorithms for computing their interaction and evolution in time.<sup>6</sup>

In other words, we will provide:

- A set of mathematical state variables that together provide a complete system description.
- A set of system event types.
- A collection of mathematical methods that will take place when each type of event occurs in order to update system state variables.

### 4.1 System's variables

In order to properly describe all system's state variables, we need to introduce some mathematical notations.

$\tau \in (t_0, t)$  denotes a time's instant of our system simulation clock where  $t_0$  and  $t$  represent respectively *start moment* and *final moment* of our simulation; we will use also  $c \in \{1, 2\} = C$ , representing the class to which a job belongs, and  $x \in \{\text{cloudlet}, \text{cloud}, \text{global}\} = X$ , used, instead, to refer to a specific system's node or to the whole system.

At this point we can introduce all mathematical variables used in our model:

$n_x^{(c)}(\tau)$	=	Number of class $c$ jobs currently running at $x$ system's node at time $\tau$
$d_x^{(c)}(\tau)$	=	Number of class $c$ departed jobs from node $x$ at time $\tau$
$s_{x,i}^{(c)}$	=	Service time of class $c$ job $i$ served on $x$ node
$i_{cloudlet}^{(2)}(\tau)$	=	Number of class 2 interrupted jobs at time $\tau$ which were running on cloudlet

We define the *state* of our system at time  $\tau$  as follows:

$$\omega(\tau) = (\omega_{cloudlet}(\tau), \omega_{cloud}(\tau)) \quad (1)$$

Where  $\omega_{cloudlet}(\tau)$  and  $\omega_{cloud}(\tau)$  are, respectively, cloudlet and cloud node state at time  $\tau$ . Specifically:

$$\begin{aligned} \omega_{cloudlet}(\tau) &= (n_{cloudlet}^{(1)}(\tau), n_{cloudlet}^{(2)}(\tau)) \\ \omega_{cloud}(\tau) &= (n_{cloud}^{(1)}(\tau), n_{cloud}^{(2)}(\tau)) \end{aligned} \quad (2)$$

Thus:

$$\omega(\tau) = ((n_{cloudlet}^{(1)}(\tau), n_{cloudlet}^{(2)}(\tau)), (n_{cloud}^{(1)}(\tau), n_{cloud}^{(2)}(\tau))) \quad (3)$$

At this point we can show some equations describing constrains or relations among these variables.

Obliviously is true that:

---

<sup>6</sup>Lawrence M. Leemis, Stephen K. Park, *Discrete-Event Simulation: A First Course* (Pearson; 1 edition January 6, 2006), 186

$$n_{global}^{(c)}(\tau) = \sum_{x \in X \setminus \{global\}} n_x^{(c)}(\tau) \quad \forall c \in C, \forall \tau \in (t_0, t) \quad (4)$$

$$d_{global}^{(c)}(\tau) = \sum_{x \in X \setminus \{global\}} d_x^{(c)}(\tau) \quad \forall c \in C, \forall \tau \in (t_0, t) \quad (5)$$

If our system is based on Algorithm 1, is true that:

$$\sum_{c \in C} n_{cloudlet}^{(c)}(\tau) \leq N \quad \forall \tau \in (t_0, t) \quad (6)$$

When, instead, our system is using Algorithm 2, is verified that:

$$\begin{aligned} \sum_{c \in C} n_{cloudlet}^{(c)}(\tau) \leq S &\Leftrightarrow n_{cloudlet}^{(2)}(\tau) > 0 \quad \forall \tau \in (t_0, t) \\ n_{cloudlet}^{(1)}(\tau) \leq N &\Leftrightarrow n_{cloudlet}^{(2)}(\tau) = 0 \quad \forall \tau \in (t_0, t) \end{aligned} \quad (7)$$

## 4.2 Statistics computing

During this subsection we will describe how we can compute all system's statistics which we need using mathematical notation explained before.

**Time-averaged population** In order to compute  $E[N_x]^{(c)}$ , that is the *time-averaged number of  $c$  class jobs in  $x$  system's node*, we can do as follows:<sup>7</sup>

$$E[N_x]^{(c)} = \frac{1}{t - t_0} \int_{t_0}^t n_x^{(c)}(\tau) d\tau \quad \forall c \in C, \forall x \in X \quad (8)$$

To calculate instead  $E[N_x]$ , *not-class-based time-averaged number of jobs running in a  $x$  system's node*, we can proceed as shown below:

$$\begin{aligned} E[N_x] &= \frac{1}{t - t_0} \cdot \int_{t_0}^t \left( n_x^{(1)}(\tau) + n_x^{(2)}(\tau) \right) d\tau \\ &= \frac{1}{t - t_0} \cdot \int_{t_0}^t n_x^{(1)}(\tau) d\tau + \frac{1}{t - t_0} \cdot \int_{t_0}^t n_x^{(2)}(\tau) d\tau \\ &= E[N_x]^{(1)} + E[N_x]^{(2)} \\ &= \sum_{c \in C} E[N_x]^{(c)} \quad \forall x \in X \end{aligned} \quad (9)$$

**Time-average response time** Are  $E[S_x]^{(c)}$  and  $E[T_x]^{(c)}$ , respectively, the *time-averaged service time and response time of class  $c$  jobs in  $x$  system's node*.

Since our system hasn't queues, there is no waiting time or delay experienced by jobs in our system, therefore result that  $E[S_x]^{(c)} = E[T_x]^{(c)}$ . Said that, we just have to show how to compute  $E[S_x]^{(c)}$ :<sup>8</sup>

<sup>7</sup>*Ibid.*, 19

<sup>8</sup>*Ibid.*, 17



$$E[T_x]^{(c)} = E[S_x]^{(c)} = \frac{1}{d_x^{(c)}(t)} \cdot \sum_{i=0}^{d_x^{(c)}(t)} s_{x,i}^{(c)} \quad \forall c \in C, \forall x \in X \quad (10)$$

If we aren't interested about per-class metrics, we can compute  $E[S_x]$  and  $E[T_x]$ , respectively *time-averaged service time* and the *time-averaged response time experienced by any class jobs in a  $x$  system's node*, we can do as follows:

$$E[T_x] = E[S_x] = \frac{1}{d_x^{(1)}(t) + d_x^{(2)}(t)} \cdot \left( \sum_{i=0}^{d_x^{(1)}(t)} s_{x,i}^{(1)} + \sum_{i=0}^{d_x^{(2)}(t)} s_{x,i}^{(2)} \right) \quad \forall x \in X \quad (11)$$

**Throughput** *Per-class system's node throughput*, denoted with  $X_x(c)$ , is the rate of per-class  $c$  completions at  $x$  system's node, while *system's node throughput*, indicated with  $X_x$ , is the rate of any class job completions at  $x$  system's node<sup>9</sup>.

Generally the throughput is the ratio between the total number of jobs completed and observation time<sup>10</sup>, therefore aforesaid metrics are simply computable as follows:

$$X_x(c) = \frac{d_x^{(c)}(t)}{t} \quad \forall c \in C, \forall x \in X \quad (12)$$

$$X_x = \frac{d_x^{(1)} + d_x^{(2)}(t)}{t} \quad \forall x \in X \quad (13)$$

**Percentage class 2 job interruptions** Is  $p_{interruption}^{(2)}$  the percentage class 2 job interruptions previously running on cloudlet. Note that this metric is computable only when our simulation is using access control algorithm 2.

To compute them we can simply do:

$$p_{interruption}^{(2)} = \frac{i_{cloudlet}^{(2)}(t)}{d_{global}^{(2)}(t)} \quad (14)$$

### 4.3 System's assumptions

To complete to turn out our conceptual model into a proper specification model, we have to make additional assumptions which we have reported below:

- Although initial state variables can have any non-negative integer value<sup>11</sup>, we have chosen, as common, to set  $n_x^{(c)}(0) = 0$  and  $d_x^{(c)}(0) = 0$  setting, in this way, initial system status as *idle*.

<sup>9</sup>Mor Harchol-Balter, *Performance Modeling and Design of Computer Systems* (Carnegie Mellon University, Pennsylvania 2013), 18

<sup>10</sup>*Ibid.*

<sup>11</sup>*Ibid.* 190

- As a consequence of the previous point, first event must be either a class 1 or class 2 job arrival on controller node.
- Terminal state is also idle. Rather than specifying the number of jobs processed, our stopping criteria has been specified in terms of a time  $\tau^*$  beyond which no new jobs can arrive. This assumption effectively closes the door at time  $\tau^*$  but allows the system to continue operation until all jobs have been completely served.

Therefore, the last event must be either a class 1 or class 2 job completion.

#### 4.4 System's events and assumptions

Let's look to all possible *events* which can occurs during our simulations.

Since by definition any event is an occurrence that may change the state of the system<sup>12</sup>, obviously our system's state variables can change only when an event occurs.

All event capable to alter our system's state and their effect on state variables are reported in table 1 but pay attention to the fact that all events reported just focus on access control algorithm 1 based system. When you use access control algorithm 2, you have to remember that, when a class 1 job arrival on cloudlet node occurs, a class 2 job may be interrupted if algorithm 2 conditions would be achieved; in this case  $n_{\text{cloudlet}}^{(1)}(\tau)$  are then incremented by one while  $n_{\text{cloudlet}}^{(2)}(\tau)$  is decremented by one and a class 2 job arrival is scheduled.

---

<sup>12</sup>Lawrence M. Leemis, Stephen K. Park, *Discrete-Event Simulation: A First Course* (Pearson; 1 edition January 6, 2006), 187

Table 1: Access control algorithm 1 system's events

Event name that occurred at time $\tau$	Event's place	Event's effects
Class 1 job arrival	Cloudlet	$n_{\text{cloudlet}}^{(1)}(\tau)++$
	Cloud	$n_{\text{cloud}}^{(1)}(\tau)++$
	Controller	$n_{\text{global}}^{(1)}(\tau)++$
Class 2 job arrival	Cloudlet	$n_{\text{cloudlet}}^{(2)}(\tau)++$
	Cloud	$n_{\text{cloud}}^{(2)}(\tau)++$
	Controller	$n_{\text{global}}^{(2)}(\tau)++$
Class 1 job departure	Cloudlet	$n_{\text{cloudlet}}^{(1)}(\tau)--$ $n_{\text{global}}^{(1)}(\tau)--$ $d_{\text{cloudlet}}^{(1)}(\tau)++$ $d_{\text{global}}^{(1)}(\tau)++$
	Cloud	$n_{\text{cloud}}^{(1)}(\tau)--$ $n_{\text{global}}^{(1)}(\tau)--$ $d_{\text{cloud}}^{(1)}(\tau)++$ $d_{\text{global}}^{(1)}(\tau)++$
Class 2 job departure	Cloudlet	$n_{\text{cloudlet}}^{(2)}(\tau)--$ $n_{\text{global}}^{(2)}(\tau)--$ $d_{\text{cloudlet}}^{(2)}(\tau)++$ $d_{\text{global}}^{(2)}(\tau)++$
	Cloud	$n_{\text{cloud}}^{(2)}(\tau)--$ $n_{\text{global}}^{(2)}(\tau)--$ $d_{\text{cloud}}^{(2)}(\tau)++$ $d_{\text{global}}^{(2)}(\tau)++$

## 5 Computational Model

Let's start now the description of our system *computational model* in which, as known, system's states exist as a collection of data structure, classes and variables that collectively characterize the system and are systematically updated as simulated time evolves.

Our model is been implemented using Java programming languages, which source code is fully available on GitHub<sup>13</sup>, famous web-based hosting service for version control using `git`; we remind that  $\text{\LaTeX}$  source code of this report is available too.<sup>14</sup>

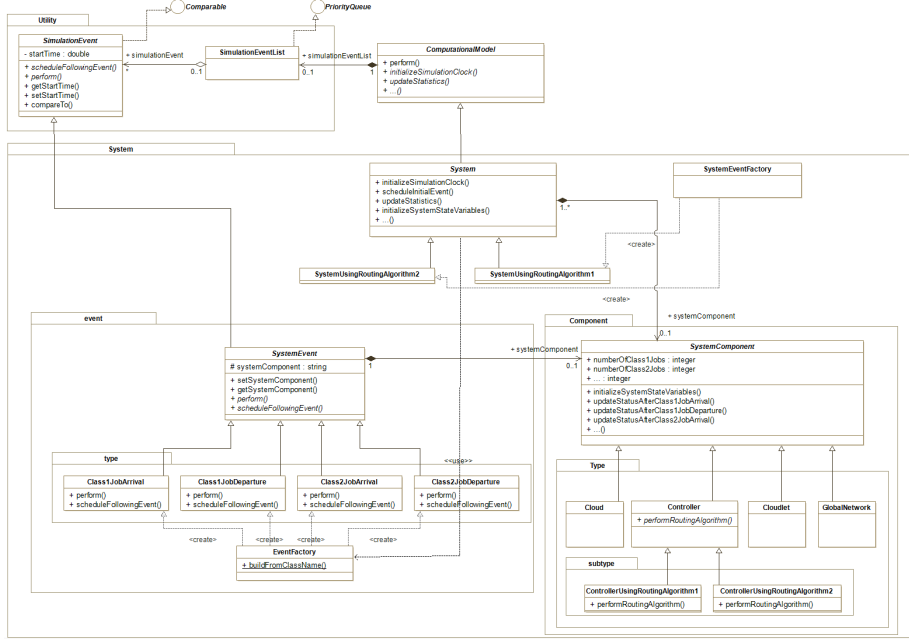


Figure 2: Class diagram showing main classes and methods only.

### 5.1 Events implementation

In order to properly describe how our implementation works we need to understand how all system's events, seen in previous section, are represented and managed.

As you can see in the class diagram 2, a *generic* system's event has been implemented and represented using an *abstract* Java class called **SimulationEvent** (located into `nexteventsimulation.utility` Java package) which has two main features:

- Has only one `double` type field, called **StartTime**, which is used to hold the instant of time according to which an event occurs. Setter and getter methods are provided in order to manage this field.

<sup>13</sup>Source code available on <https://github.com/AndreaG93/PMCSN-Project>

<sup>14</sup>See <https://github.com/AndreaG93/PMCSN-Project-Report>

- It declares two very important abstract methods which, as we will see later, are called by next-event simulation algorithm implemented in `ComputationalModel` class; these methods are called `perform()`, whose purpose is to update system's state variables when an event occurs, and `scheduleFollowingEvent()`, which is used, instead, to insert into an *event list* events that follows the current one.

Although it is obvious that exact methods behaviour depends strictly by which concrete class implements them, it is however very important to point out that `SimulationEvent` class does *not* represent system-specific events. In fact, to ensure low coupling and high cohesion among classes, neither `SimulationEvent` or `ComputationalModel` class need to know any specific system information.

System-specific events are instead represented by instances of `SystemEvent` class (located into `nexteventsimulation.computationalmodel.model.system.event` package) whose code is partially shown in Listing 1.

Listing 1: `SystemEvent` class implementation.

---

```

1 public abstract class SystemEvent extends SimulationEvent {
2
3     protected SystemComponent systemComponent;
4
5     public void setSystemComponent(SystemComponent
6         systemComponent) {
7         this.systemComponent = systemComponent;
8     }
9
10    public SystemComponent getSystemComponent() {
11        return systemComponent;
12    }
13    ...

```

---

As you can see from reported snippet, `SystemEvent` class has a `SystemComponent` type field, called `systemComponent`, which is used to hold a reference to that system's node where a generic event occurs; in this way any `SystemEvent` instance is able to update state variables of the component of which it has the reference, without altering others parts of system.

However, in order to do that, required abstract method implementation is needed. In our implementation there are four concrete `SystemEvent` type classes used to represent all system's events:

- `Class1JobArrival` class which represents a class 1 job arrival.
- `Class2JobArrival` class which models a class 2 job arrival.
- `Class1JobDeparture` class is used to model a class 1 job departure.
- `Class2JobDeparture` class for class 2 jobs departures.

Any instance of above classes, having knowledge about methods published by `SystemComponent` class, can finally update states of any component of our system. In this way we can achieve an high degree of cohesion among classes

and low coupling taking advantage of polymorphism pattern; in fact, as we can see also from shown class diagram, `ComputationalModel` class has to know only about `SimulationEvent` class and, thorough invocation of its methods, being able to update our system ignoring all specific information about our system.

## 5.2 Event list

As known, an *event list*, also called *calendar*, represents a data structure used to keep track of scheduled time of occurrence for the next possible event.

For reasons of efficiency, our event list is been implemented as a *min-heap binary tree*, that is a *priority queue* in which all events are ordered by time, in such a way that the most imminent event is found on the tree's root. This data structure is able to provide  $O(\log_2(n))$  time for the en-queuing and dequeuing operations and constant time for root element retrieval.<sup>15</sup>

Fortunately a Java implementation is already available in JDK through built-in class called `PriorityQueue<E>` which provides all method that we need, like `poll`, for root element retrieval, `add` and `remove`, used respectively for en-queuing and dequeuing operations.

However, in order to that data structure works properly, is required that a total ordering among all possible event exists. To provide it, we need that all `SimulationEvent` type instances implement an interface called `Comparable` providing an implementation, shown in listing 2, of a method called `compareTo`.<sup>16</sup> In this way, is very easy to order all `SimulationEvent` type instances based on their start event time, from the youngest to the oldest.

Listing 2: Snippet of `SimulationEvent` class implementation

---

```

1  @Override
2  public int compareTo(SimulationEvent o) {
3      return Double.compare(this.startTime, o.startTime);
4  }

```

---

## 5.3 Simulation clock

As we know a *simulation clock* is necessary to keep track of the current value of simulated time during simulation.

In our system's implementation, our simulation clock is represented by a `SimulationClock` type instance, whose code is partially reported in Listing 3.

Simulation clock state is represented by two `double` type fields: the first, called `currentEventTime`, is used to holding current simulated time value, while the second, called `nextEventTime`, is used instead to keep track of next event simulated time value; obviously, proper setter and getter methods are provided in order to manage them.

From software architecture point of view, believing that multiple `SimulationClock` instance are not necessary for our purposes, we have choosen to design this class using the *singleton* pattern.

<sup>15</sup>Camil Demetrescu, Irene Finocchi, Giuseppe F. Italiano - *Algoritmi e strutture dati* - second edition, McGraw-Hill, pag. 200

<sup>16</sup>See [https://docs.oracle.com/en/java/javase/12/docs/api/java.base/java/lang/Comparable.html#compareTo\(T\)](https://docs.oracle.com/en/java/javase/12/docs/api/java.base/java/lang/Comparable.html#compareTo(T))

Listing 3: SimulationClock class implementation.

---

```
1 public class SimulationClock {
2
3     private static SimulationClock instance = null;
4
5     private double currentTime;
6     private double nextEventTime;
7
8     private SimulationClock(){
9     }
10
11    public static SimulationClock getInstance() {
12        if (instance == null)
13            instance = new SimulationClock();
14        return instance;
15    }
16    ...
```

---

## 5.4 Next-event simulation logic

The most important code is included into `ComputationalModel` class because it contains main simulation logic.

To start a next-event simulation is very easy because we just need to invoke an its method called `perform` which contains the algorithm<sup>17</sup> used to perform a simulation based on next-event approach, whose Java implementation is reported in Listing 4. Our algorithm consists of the following steps:

**Simulation initialization** As phase name suggests, during this step system state variables, simulation clock, event list and everything is needed for simulation is initialized compatibly with what we have described in previous section, during specification model analysis.

Remember that `ComputationalModel` is an abstract class which knows anything about our specific system, therefore, obviously, system initialization phase differs according to which concrete class implements that class. Anyhow, as you can see from code snippet reported below, this task is performed invoking `initializeSimulation()` method; that invocation represents a short-cut used to call a large set of abstract methods for initialization purposes, like `initializeSystemStateVariables`, `initializeSimulationClock` and `scheduleInitialEvent`.

**Process current event** Although we believe the code is self explanatory, we sum up this phase as follows: until event list is not empty (line 6), the most imminent possible event is extracted and removed from event list (line 8); then simulation clock is advanced to this event's scheduled time of occurrence (line 12) and finally system state is updated (line 14).

**Schedule following event** Until a fixed simulation time  $\tau^*$  is not exceeded, a new event is generated by current event and consequently placed into

---

<sup>17</sup>Lawrence M. Leemis, Stephen K. Park, *Discrete-Event Simulation: A First Course* (Pearson; 1 edition January 6, 2006), Algorithm 5.1.1, page 189

event list. Subsequently simulation clock state is updated in order to compute some system metrics.

**Statistic update** While simulation advance from one event time to the next, through an invocation to an abstract method called `updateStatistics`, is possible to update some variables for system's statistic metrics computation.

**Simulation result generation** During this final step, all collected data are elaborated using various approaches, producing some statistic metrics including graphics like histograms or scatter plots. We will examine involved classes later.

Listing 4: Snippet of `perform` method

---

```
1 public void perform() {
2
3     initializeSimulation();
4
5     while (!this.simulationEventList.isEmpty()) {
6
7         SimulationEvent actualEvent = this.
            simulationEventList.poll();
8
9         if (actualEvent != null) {
10
11             SimulationClock.getInstance().
                setCurrentEventTime(actualEvent.
                    getStartTime());
12
13             actualEvent.perform();
14             actualEvent.scheduleFollowingEvent();
15
16             SimulationEvent nextEvent = this.
                simulationEventList.peek();
17
18             if (nextEvent != null)
19                 SimulationClock.getInstance().
                    setNextEventTime(nextEvent.
                        getStartTime());
20         }
21         updateStatistics();
22     }
23
24     BatchMeansRegister.getInstance().
        computeStatisticsAndWriteData();
25     ScatterPlotRegister.getInstance().writingOutputData
        ();
26     HistogramsRegister.getInstance().
        computeStatisticsAndWriteData();
27
28     manageCurrentSimulationResult();
29 }
```

---



## 5.5 System network implementation

As we have seen in previous subsection, `ComputationalModel` is used to represent a very *generic* form of system (or queues network) because it includes only the code, seen previously, used to perform a next-event simulation.

Our specific queues network is been modelled using an abstract class called `System`, which provide all required methods implementations in order to perform our simulation.

The state of any `System` type instance is made up of four `SystemComponent` type fields which are used to hold a reference to any system's component, that is cloud node, cloudlet node and controller; an extra `SystemComponent` type field is a short-cut used to represent the global system.

A `System` type instance has very important responsibilities including:

- Allocation all `SystemComponent` type objects in order to represent each system's nodes.
- Methods implementation providing for simulation clock and system nodes state variables initialization.
- To manage event scheduling through some methods invocations like `scheduleEventOnCloud`, `scheduleEventOnCloudlet` etc.
- To provide a method able to manage class 2 job interruption; the so called `removeCloudletClass2JobDeparture` method

As you can see from class diagram reported in Figure 2, there are two `System`'s concrete sub-class called `SystemUsingRoutingAlgorithm1` and `SystemUsingRoutingAlgorithm2`; as their names suggest, they are used respectively to represents a system based on access control algorithm 1 and 2.

A *generic* system's node, like cloud or cloudlet node, is modelled by a `SystemComponent` class.

class has key role during our simulation performing because:

- Has a set of fields used to keep track of component's variable states, like number of departed jobs or current jobs population, including those used to compute time-average statistics.
- Provides a very important set of functions for status updating after event occurrences; for example, when a class one job arrives on a certain system's component, `updateStatusAfterClass1JobArrival()` method can be invoked in order to update the number of class one jobs actually present on that component.
- Contains a set of method used to computed inter-arrival and service times, necessary for future event scheduling. Obviously each concrete `SystemComponent` instance has his own mean rates to compute those values.
- Includes another set used instead to schedule future events after event occurrence, like a departure after an arrival.

Is very important to precise that these methods are abstract and have to be implemented by each concrete implementations in order to reflect each component behaviour; In fact controller node doesn't need to schedule departure after an arrival because of his routing role.

All system's component, that is cloud, cloudlet and controller, are represented, obviously, by `Cloud`, `Cloudlet` and `Controller` classes each of which is a `SystemComponent` subclass; however we have used another subclass, called `GlobalNetwork`, in order to represent whole system keeping track of global state variables.

Observe that, `Controller` class is also an abstract class because required implementation of a method call `performRoutingAlgorithm()` used to perform the two required routing algorithms whenever an arrival occurs on controller node; those algorithms are provided by `ControllerUsingRoutingAlgorithm1` and `ControllerUsingRoutingAlgorithm2` classes.

## 5.6 Output analysis

each statistic is associated to an `TransientStatistics` type instance.

## 6 Analytical solution

In this last section we will develop an analytical solution to validate the results obtained previously through our simulations.

### 6.1 System based on access control Algorithm 1

Let's start with presentation of analytical solution of system based on access control algorithm 1.

In order to compute all parameters and metrics of interest associated with above-mentioned system, due to cloudlet's limited resources according to which it can accept jobs until their number does not exceed a given threshold  $N$ , is crucial compute first the **fraction of jobs that are forwarded to cloudlet and to cloud**; to do it, we must compute at first probability according to which the sum of job of each class in cloudlet system is equal to that threshold.

To determine this probability, we had modelled cloudlet with a **continuous-time Markov chain (CTMC)**, of which you can see a graphical representation in Figure 3, where each chain's state, denoted with  $(n_1, n_2)$ , is represented by the number of class 1 job,  $n_1$ , and class 2 job,  $n_2$ , present in system at a certain moment.

#### 6.1.1 Balance equation

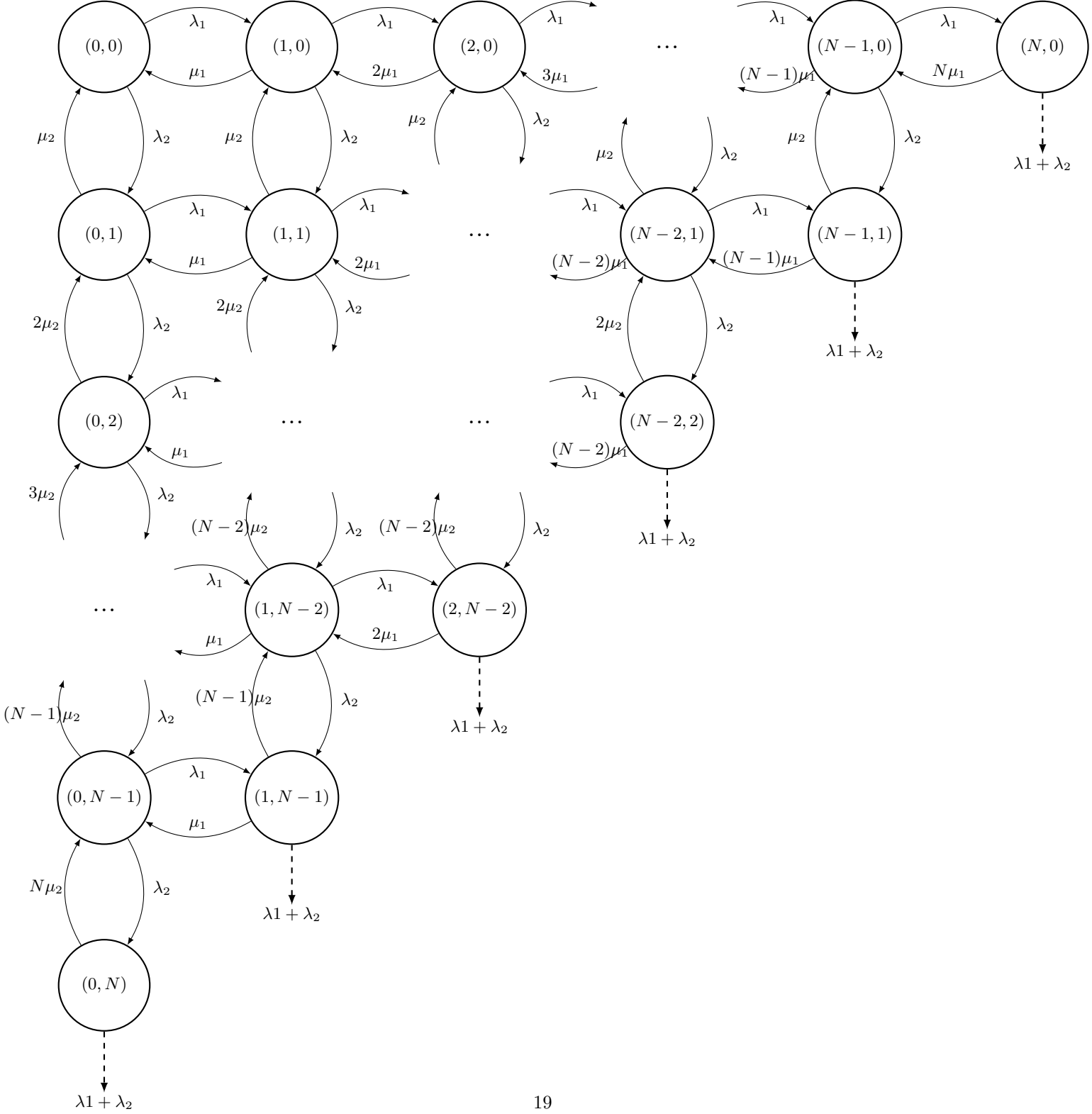
Obviously we can compute **limiting probabilities**  $\pi_{(n_1, n_2)}$ , namely the probability according to which the chain is in a certain state, say  $j$ , independently of the starting state, say  $i$ , via **balance** (or **stationary**) **equations**, in which we can equate the rate at which the system leaves state  $j$  with the rate at which the system enters state  $j$ <sup>18</sup>, **remembering that limiting probabilities sum to 1** (i.e.,  $\sum_{j=0}^{\infty} \pi_j = 1$ ).

Table 2: Balance equations.

$(\lambda_1 + \lambda_2)\pi_{(0,0)}$	$=$	$\mu_1\pi_{(1,0)} + \mu_2\pi_{(0,1)}$	
$(\lambda_1 + \lambda_2 + n_1\mu_1)\pi_{(n_1,0)}$	$=$	$\lambda_1\pi_{(n_1-1,0)} + \mu_1(n_1+1)\pi_{(n_1+1,0)} + \mu_2\pi_{(n_1,1)}$	$\forall n_1 \in \mathbb{N} \cap [1, N-1]$
$(\lambda_1 + \lambda_2 + n_2\mu_2)\pi_{(0,n_2)}$	$=$	$\lambda_2\pi_{(0,n_2-1)} + \mu_1\pi_{(1,n_2)} + \mu_2(n_2+1)\pi_{(0,n_2+1)}$	$\forall n_2 \in \mathbb{N} \cap [1, N-1]$
$\mu_1 N \pi_{(N,0)}$	$=$	$\lambda_1 \pi_{(N-1,0)}$	
$\mu_2 N \pi_{(0,N)}$	$=$	$\lambda_2 \pi_{(0,N-1)}$	
$(n_1\mu_1 + n_2\mu_2)\pi_{(n_1,n_2)}$	$=$	$\lambda_1\pi_{(n_1-1,n_2)} + \lambda_2\pi_{(n_1,n_2-1)}$	$\forall n_1, n_2 \in \mathbb{N} \cap [1, N-1] \mid n_1 + n_2 = N$
$(\lambda_1 + \lambda_2 + n_1\mu_1 + n_2\mu_2)\pi_{(n_1,n_2)}$	$=$	$\lambda_1\pi_{(n_1-1,n_2)} + \lambda_2\pi_{(n_1,n_2-1)} + \mu_1(n_1+1)\pi_{(n_1+1,n_2)} + \mu_2(n_2+1)\pi_{(n_1,n_2+1)}$	$\forall n_1, n_2 \in \mathbb{N} \cap [1, N-1] \mid n_1 + n_2 < N$

<sup>18</sup>Cfr. Mor Harchol-Balter - *Performance Modeling and Design of Computer Systems* - Carnegie Mellon University, Pennsylvania, pag. 237

Figure 3: Access control algorithm 1 based cloudlet's CTMC.



### 6.1.2 Probabilities computing

Having found the stationary probabilities, we can now find  $\Pi_{\text{SendToCloud}}$ , that is the **probability that an arriving job on controller has to be forwarded to cloud**. Observe that the class to which an arrival job belongs to is not important because, according to access control Algorithm 1, **jobs of both classes have same probability to be sent to cloud**. To be more precise,  $\Pi_{\text{SendToCloud}}$  is the probability that an arrival job find that the number of jobs present in cloudlet has exceeded threshold  $N$ , which occurs when  $n_1 + n_2 = N$ . Formally:

$$\begin{aligned}
 \Pi_{\text{SendToCloud}} &= P\{\text{An arrival job on controller sees } N \text{ jobs in cloudlet}\} \\
 &= \text{Limiting probability that there are } N \text{ jobs in system} \\
 &= \sum_{\substack{n_1, n_2 \in \mathbb{N} \cap [0, N] \\ n_1 + n_2 = N}} \pi_{(n_1, n_2)}
 \end{aligned} \tag{15}$$

At this point we can easily compute  $\Pi_{\text{SendToCloudlet}}$ , which instead represents the **probability according to which an arriving job on controller has to be accepted on cloudlet** and it is same for both job classes too.

$$\begin{aligned}
 \Pi_{\text{SendToCloudlet}} &= P\{\text{An arrival job on controller sees less than } N \text{ jobs in cloudlet}\} \\
 &= 1 - P\{\text{An arrival job on controller sees } N \text{ jobs in cloudlet}\} \\
 &= 1 - \Pi_{\text{SendToCloud}}
 \end{aligned} \tag{16}$$

### 6.1.3 Average arrival rates

Is  $\lambda_i(c)$  the **total arrival rate into a system's component  $i$  of class  $c$  job**. Applying previous results, using an appropriate equation<sup>19</sup>, we can now compute per-class average arrival rates as follow:

$$\begin{aligned}
 \lambda_{\text{cloud}}(1) &= \lambda_1 \cdot \Pi_{\text{SendToCloud}} \\
 \lambda_{\text{cloud}}(2) &= \lambda_2 \cdot \Pi_{\text{SendToCloud}} \\
 \lambda_{\text{cloudlet}}(1) &= \lambda_1 \cdot \Pi_{\text{SendToCloudlet}} \\
 \lambda_{\text{cloudlet}}(2) &= \lambda_2 \cdot \Pi_{\text{SendToCloudlet}}
 \end{aligned} \tag{17}$$

Then we get  $\lambda_i$ , that is **total arrival rate to system's component  $i$** , by summing the per-class rates as follows:

$$\begin{aligned}
 \lambda_{\text{cloud}} &= \lambda_{\text{cloud}}(1) + \lambda_{\text{cloud}}(2) \\
 \lambda_{\text{cloudlet}} &= \lambda_{\text{cloudlet}}(1) + \lambda_{\text{cloudlet}}(2)
 \end{aligned} \tag{18}$$

---

<sup>19</sup> Cfr. Ivi, pag. 315, equation (18.1)

#### 6.1.4 Average population

Is  $E[N_x]^{(c)}$  the time-average population of class  $c$  jobs into a system's node  $x$ .

We can use previously computed stationary probabilities to get average population for cloudlet; it's enough to sum each state's limiting probability multiplied by corresponding number of job as follows:

$$\begin{aligned}
E[N_{\text{cloudlet}}](1) &= \sum_{(n_1, n_2) \in M} n_1 \cdot \pi(n_1, n_2) \\
E[N_{\text{cloudlet}}](2) &= \sum_{(n_1, n_2) \in M} n_2 \cdot \pi(n_1, n_2) \\
E[N_{\text{cloudlet}}] &= E[N_{\text{cloudlet}}](1) + E[N_{\text{cloudlet}}](2) \\
&= \sum_{(n_1, n_2) \in M} (n_1 + n_2) \cdot \pi(n_1, n_2)
\end{aligned} \tag{19}$$

Since we have modelled cloud component as a  $M/M/\infty$  system, in which there is no job's waiting time due to presence of an infinite number of servers, we can simply to apply **Little's Law**<sup>20</sup> in order to get cloud's average population as shown in Equation ??.

$$\begin{aligned}
E[N_{\text{cloud}}](1) &= \lambda_{\text{cloud}}(1) \cdot E[T_{\text{cloud}}](1) \\
&= \lambda_{\text{cloud}}(1) \cdot (E[T_{Q_{\text{cloud}}}](1) + E[S_{\text{cloud}}](1)) \\
&= \lambda_{\text{cloud}}(1) \cdot \left(\frac{1}{\mu_{\text{cloud}}(1)}\right) \\
&= \frac{\lambda_{\text{cloud}}(1)}{\mu_{\text{cloud}}(1)}
\end{aligned} \tag{20}$$

Similarly:

$$E[N_{\text{cloud}}](2) = \frac{\lambda_{\text{cloud}}(2)}{\mu_{\text{cloud}}(2)} \tag{21}$$

$$E[N_{\text{cloud}}] = E[N_{\text{cloud}}](1) + E[N_{\text{cloud}}](2) \tag{22}$$

Finally we can get global average job populations as follows:

$$\begin{aligned}
E[N](1) &= E[N_{\text{cloudlet}}](1) + E[N_{\text{cloud}}](1) \\
E[N](2) &= E[N_{\text{cloudlet}}](2) + E[N_{\text{cloud}}](2)
\end{aligned} \tag{23}$$

$$E[N] = E[N](1) + E[N](2)$$

---

<sup>20</sup> Cfr. Ivi, pag. 95, theorem (6.1)

### 6.1.5 Average response time

Is  $E[T_i](c)$  the **mean response time experienced by a class  $c$  jobs into a component  $i$** .

In order to properly compute said metric for each system's component observe that:

- Knowing per-class average job population and per-class mean arrival rates, we can easily compute  $E[T_i](c)$  using Little's Law.
- Since our system haven't queues, because of there is no waiting time experienced by jobs, is true that  $E[T_i](c)$  is also equal to  $E[S_i](c)$ , that is **mean service time experienced by a class  $c$  jobs into a component  $i$** , which is equal to  $1/\mu_i(c)$ , where  $\mu_i(c)$  **means average service rate at which a class  $c$  jobs into a component  $i$  is served**.

Therefore we can get these metric as follows.

$$\begin{aligned}
 E[T_{\text{cloudlet}}](1) &= \frac{E[N_{\text{cloudlet}}](1)}{\lambda_{\text{cloudlet}}(1)} \\
 &= E[S_{\text{cloudlet}}](1) \\
 &= \frac{1}{\mu_{\text{cloudlet}}(1)}
 \end{aligned} \tag{24}$$

$$\begin{aligned}
 E[T_{\text{cloudlet}}](2) &= \frac{E[N_{\text{cloudlet}}](2)}{\lambda_{\text{cloudlet}}(2)} \\
 &= E[S_{\text{cloudlet}}](2) \\
 &= \frac{1}{\mu_{\text{cloudlet}}(2)}
 \end{aligned} \tag{25}$$

$$\begin{aligned}
 E[T_{\text{cloud}}](1) &= \frac{E[N_{\text{cloud}}](1)}{\lambda_{\text{cloud}}(1)} \\
 &= E[S_{\text{cloud}}](1) \\
 &= \frac{1}{\mu_{\text{cloud}}(1)}
 \end{aligned} \tag{26}$$

$$\begin{aligned}
 E[T_{\text{cloud}}](2) &= \frac{E[N_{\text{cloud}}](2)}{\lambda_{\text{cloud}}(2)} \\
 &= E[S_{\text{cloud}}](2) \\
 &= \frac{1}{\mu_{\text{cloud}}(2)}
 \end{aligned} \tag{27}$$

At this point we can get global per-class mean response times as follows:

$$\begin{aligned}
E[T](1) &= E[T_{\text{cloudlet}}](1) \cdot P\{\text{An arrival class 1 job is sent to cloudlet}\} \\
&\quad + E[T_{\text{cloud}}](1) \cdot P\{\text{An arrival class 1 job is sent to cloud}\} \\
&= E[T_{\text{cloudlet}}](1) \cdot \Pi_{\text{SendToCloudlet}} + E[T_{\text{cloud}}](1) \cdot \Pi_{\text{SendToCloud}}
\end{aligned} \tag{28}$$

Similarly:

$$\begin{aligned}
E[T](2) &= E[T_{\text{cloudlet}}](2) \cdot P\{\text{An arrival class 2 job is sent to cloudlet}\} \\
&\quad + E[T_{\text{cloud}}](2) \cdot P\{\text{An arrival class 2 job is sent to cloud}\} \\
&= E[T_{\text{cloudlet}}](2) \cdot \Pi_{\text{SendToCloudlet}} + E[T_{\text{cloud}}](2) \cdot \Pi_{\text{SendToCloud}}
\end{aligned} \tag{29}$$

Finally, using obtained results, we can get global mean response times as shown below:

$$E[T] = E[T](1) \cdot \frac{\lambda_1}{\lambda_1 + \lambda_2} + E[T](2) \cdot \frac{\lambda_2}{\lambda_1 + \lambda_2} \tag{30}$$

#### 6.1.6 Throughput

??

To determine system's throughput let's prove if our system is stable. As we know, every queueing system in which its mean arrival rate is less than its mean service rate is known as a stable system; formally, are  $\lambda$  and  $\mu$  our system's mean arrival rate and mean service rate respectively, if  $\lambda \leq \mu$ , given system is stable. Cloudlet subsystem is clearly *not* stable due to its limited resources compared to its workload; we have already seen that exists a not null probability according to which an arriving job on controller sees  $N$  jobs in cloudlet which implying its forward to cloud subsystem. Since it haven't a queue, cloudlet stability condition is achieved when  $\Pi_{\text{SendToCloud}}$  is null but, based on current system parameters, is not the case. Regarding cloud subsystem, no matter how high we make  $\lambda_{\text{cloud}}$  because it is made up of an infinite number of server for which completion rate is still bounded by the arrival rate. Accordingly to previous considerations we can conclude that **our system is stable** due to stability or its cloud subsystem therefore we can get system throughput as follows:

$$X = \lambda = \lambda_1 + \lambda_2 \tag{31}$$

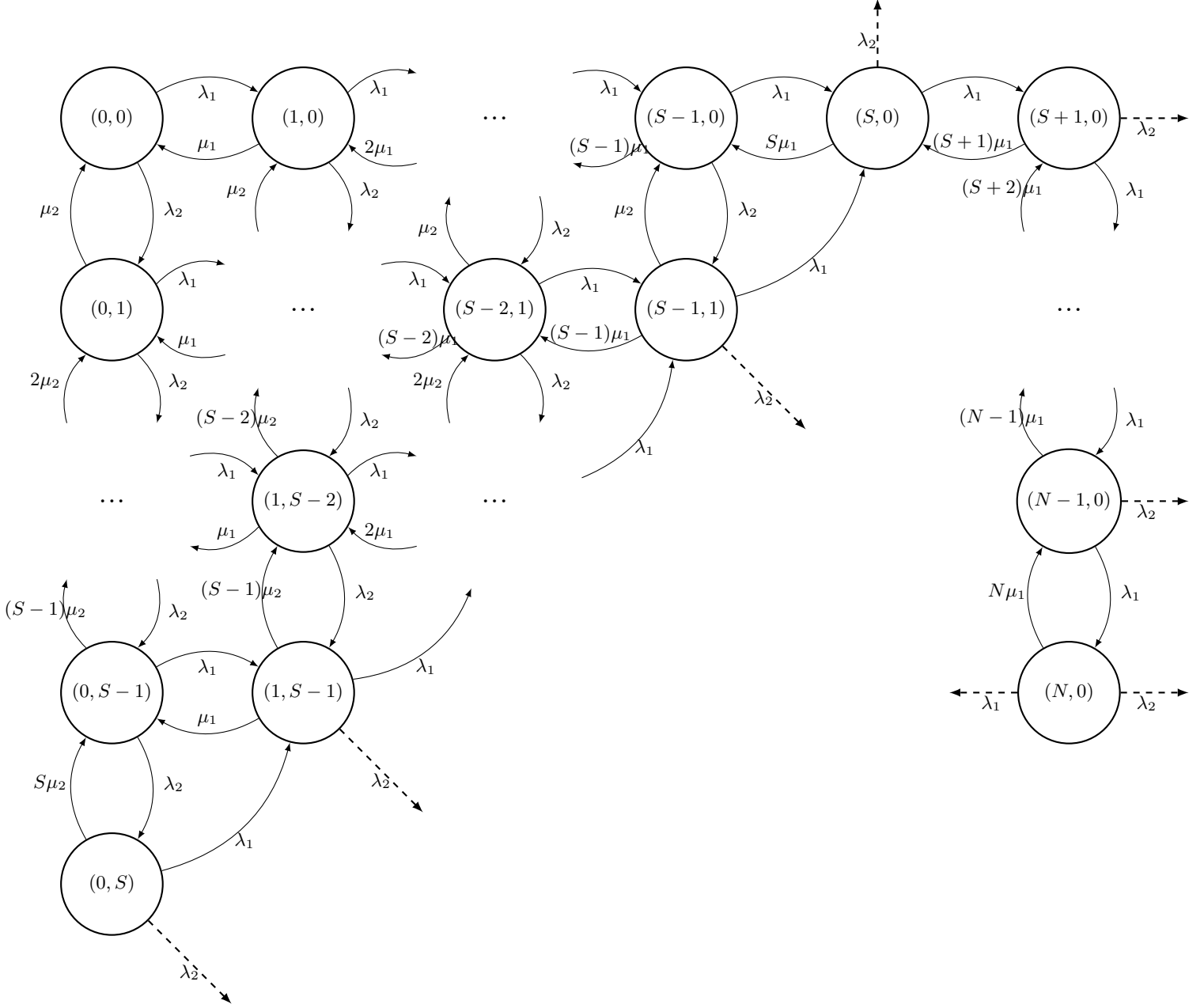
To get per-subsystems throughput we proceed as shown below:

$$X_{\text{cloud}} = \lambda_{\text{cloud}} \tag{32}$$

$$X_{\text{cloudlet}} = X - X_{\text{cloud}} \tag{33}$$



Figure 4: Cloudlet system component modeled using a CTMC



## 6.2 System based on access control Algorithm 2

Based on our analysis, performance achieved by our system when its routing policy is based on second algorithm, reported in Algorithm 2, are very different from previous case.

The main feature of this kind of system is that a class 2 job, running on cloudlet, may be stopped through its execution on system's cloudlet node and then subsequently resumed in system's cloud node; is important to specify that adopted routing policy is *not* work-conserving because it re-runs interrupted jobs losing work previously done. Moreover, when a class 2 job is interrupted and sent on the system's cloud node, a setup time, denoted whit  $s_{setup}$ , has to be considered to restart the task on the cloud.

Most equations described and used in section 6.1 are still valid also for this system, although there are differences about numerical values of each metrics; therefore in this section we focus only to differences compared to previous case.

### 6.2.1 Balance equation

Regarding balance equation necessary for stationary probabilities computation, since there are no differences about mathematical notation used or significant conceptual differences compared to previously case, we limit ourselves to report them in table 3.

Table 3: Balance equations.

$(\lambda_1 + \lambda_2)\pi_{(0,0)}$	$=$	$\mu_1\pi_{(1,0)} + \mu_2\pi_{(0,1)}$	
$(\lambda_1 + \lambda_2 + n_1\mu_1)\pi_{(n_1,0)}$	$=$	$\lambda_1\pi_{(n_1-1,0)} + \mu_1(n_1 + 1)\pi_{(n_1+1,0)} + \mu_2\pi_{(n_1,1)}$	$\forall n_1 \in \mathbb{N} \cap [1, S - 1]$
$(\lambda_1 + \lambda_2 + n_2\mu_2)\pi_{(0,n_2)}$	$=$	$\lambda_2\pi_{(0,n_2-1)} + \mu_1\pi_{(1,n_2)} + \mu_2(n_2 + 1)\pi_{(0,n_2+1)}$	$\forall n_2 \in \mathbb{N} \cap [1, S - 1]$
$(S\mu_1 + \lambda_1)\pi_{(S,0)}$	$=$	$\lambda_1\pi_{(S-1,0)} + \lambda_1\pi_{(S-1,1)} + (S + 1)\mu_1\pi_{(S+1,0)}$	
$(\lambda_1 + S\mu_2)\pi_{(0,S)}$	$=$	$\lambda_2\pi_{(0,S-1)}$	
$(n_1\mu_1 + n_2\mu_2 + \lambda_1)\pi_{(n_1,n_2)}$	$=$	$\lambda_1\pi_{(n_1-1,n_2+1)} + \lambda_1\pi_{(n_1-1,n_2)} + \lambda_2\pi_{(n_1,n_2-1)}$	$\forall n_1, n_2 \in \mathbb{N} \cap [1, S - 1] \mid n_1 + n_2 = S$
$(\lambda_1 + \lambda_2 + n_1\mu_1 + n_2\mu_2)\pi_{(n_1,n_2)}$	$=$	$\lambda_1\pi_{(n_1-1,n_2)} + \lambda_2\pi_{(n_1,n_2-1)} + \mu_1(n_1 + 1)\pi_{(n_1+1,n_2)} + \mu_2(n_2 + 1)\pi_{(n_1,n_2+1)}$	$\forall n_1, n_2 \in \mathbb{N} \cap [1, S - 1] \mid n_1 + n_2 < S$
$(n_1\mu_1 + \lambda_1)\pi_{(n_1,0)}$	$=$	$\mu_1(n_1 + 1)\pi_{(n_1+1,0)} + \lambda_1\pi_{(n_1-1,0)}$	$\forall n_1 \in \mathbb{N} \mid S + 1 \leq n_1 \leq N - 1$
$N\mu_1\pi_{(N,0)}$	$=$	$\lambda_1\pi_{(N-1,0)}$	

### 6.2.2 Probabilities computing

Let's start clarifying all differences, compared to previous case, about routing probabilities necessary for analytical results computations.

Is  $\Pi_{\text{SendToCloud}}^{(c)}$  the probability according to which an arriving class  $c$  job on controller has to be forwarded to system's cloud node. Observing both CTMC and routing policy algorithm is very easy to understand that:

$$\begin{aligned}
\Pi_{\text{SendToCloud}}^{(1)} &= P\{\text{An arrival class 1 job on controller sees } N \text{ class 1 jobs in cloudlet.}\} \\
&= P\{n_1 = N\} \\
&= \pi_{(N,0)}
\end{aligned} \tag{34}$$

$$\begin{aligned}
\Pi_{\text{SendToCloud}}^{(2)} &= P\left\{ \begin{array}{l} \text{An arrival class 2 job on controller sees that number of} \\ \text{jobs in cloudlet exceed or is equal to a given threshold } S. \end{array} \right\} \\
&= P\{n_1 + n_2 \geq S\} \\
&= \sum_{\substack{0 \leq n_1 \leq S \\ 0 \leq n_2 \leq S \\ n_1 + n_2 = S}} \pi_{(n_1, n_2)} + \sum_{n_1=S+1}^N \pi_{(n_1, 0)} \\
&= \sum_{\substack{0 \leq n_1 \leq N \\ 0 \leq n_2 \leq S \\ n_1 + n_2 \geq S}} \pi_{(n_1, n_2)}
\end{aligned} \tag{35}$$

Is  $\Pi_{\text{SendToCloudlet}}^{(c)}$  the probability that an arriving class  $c$  job on controller has to be forwarded to system's cloudlet node.

In order to compute these probabilities we can proceed as shown below:

$$\begin{aligned}
\Pi_{\text{SendToCloudlet}}^{(1)} &= 1 - \Pi_{\text{SendToCloud}}^{(1)} \\
\Pi_{\text{SendToCloudlet}}^{(2)} &= 1 - \Pi_{\text{SendToCloud}}^{(2)}
\end{aligned} \tag{36}$$

To compute  $\Pi_{\text{interruption}}^{(2)}$ , that is the probability according to which a class 2 job, running on cloudlet, have to be stopped and forwarded to system's cloud node, is necessary know firstly  $\Pi_{\text{interruptible}}$ , the probability that system's cloudlet node is in a state in which any future class 1 job arrival is able to cause a class 2 job interruption. Formally:

$$\begin{aligned}
\Pi_{\text{interruptible}} &= P \left\{ \begin{array}{l} \text{System's cloudlet node is in a state in which any future class} \\ \text{1 job arrival is able to cause a class 2 job interruption.} \end{array} \right\} \\
&= P\{n_1 + n_2 = S, n_2 > 0\} \\
&= \sum_{\substack{0 \leq n_1 \leq S-1 \\ 1 \leq n_2 \leq S \\ n_1 + n_2 = S}} \pi(n_1, n_2)
\end{aligned} \tag{37}$$

At this point we can formally define  $\Pi_{\text{interruption}}^{(2)}$  as follows:

$$\begin{aligned}
\Pi_{\text{interruption}}^{(2)} &= P \left\{ \begin{array}{l} \text{A class 1 job is forwarded to system's cloudlet node when} \\ \text{it is in a state in which any class 1 job arrival is able to} \\ \text{cause a class 2 job interruption.} \end{array} \right\} \\
&= P \left\{ \left( \begin{array}{l} \text{System's cloudlet node is in a} \\ \text{state in which any future class} \\ \text{1 job arrival is able to cause a} \\ \text{class 2 job interruption.} \end{array} \right) \cap \left( \begin{array}{l} \text{A class 1 job is forwarded} \\ \text{to system's cloudlet node.} \end{array} \right) \right\} \\
&= \Pi_{\text{interruptible}} \cdot \Pi_{\text{SendToCloudlet}}^{(1)}
\end{aligned} \tag{38}$$

### 6.2.3 Average arrival rates

Finally, applying previous results, we can now compute per-class average arrival rates as follow of our system:

$$\begin{aligned}
\lambda_{\text{cloud}}^{(1)} &= \lambda_1 \cdot \Pi_{\text{SendToCloud}}^{(1)} \\
\lambda_{\text{cloud}}^{(2)} &= \lambda_2 \cdot \Pi_{\text{SendToCloud}}^{(2)} \\
\lambda_{\text{cloudlet}}^{(1)} &= \lambda_1 \cdot \Pi_{\text{SendToCloudlet}}^{(1)} \\
\lambda_{\text{cloudlet}}^{(2)} &= \lambda_2 \cdot \Pi_{\text{SendToCloudlet}}^{(2)}
\end{aligned} \tag{39}$$

Then we can get  $\lambda_x$ , that is total arrival rate to system's component  $x$ , proceeding as follows:

$$\begin{aligned}
\lambda_{\text{cloud}} &= \lambda_{\text{cloud}}^{(1)} + \lambda_{\text{cloud}}^{(2)} + \lambda_{\text{cloudlet}}^{(2)} \cdot \Pi_{\text{interruption}}^{(2)} \\
\lambda_{\text{cloudlet}} &= \lambda_{\text{cloudlet}}^{(1)} + \lambda_{\text{cloudlet}}^{(2)}
\end{aligned} \tag{40}$$

#### 6.2.4 Other metrics

In order to compute  $E[N_x]^{(c)}$  and  $E[T_x]^{(c)}$ , respectively time-averaged population and response time of class  $c$  jobs running in a  $x$  system's node, is possible to use same equations described in previous section because, although stationary probabilities and arrival rates have different numeric values due to different routing policy, there are no conceptual differences between the two kind of system. Analytical methods used to computed system throughput are still valid.

Regarding  $E[T]^{(2)}$ , that is time-averaged response time experienced by class 2 jobs, there are, instead, some conceptual differences due to jobs interruptions. Formally:

$$\begin{aligned}
E[T]^{(2)} &= E[T_{\text{cloudlet}}]^{(2)} \cdot P\{\text{An arrival class 2 job is sent to cloudlet}\} \\
&\quad + E[T_{\text{cloud}}]^{(2)} \cdot P\{\text{An arrival class 2 job is sent to cloud}\} \\
&\quad + E[S_{\text{setup}}] \\
&= E[T_{\text{cloudlet}}]^{(2)} \cdot \Pi_{\text{SendToCloudlet}}^{(2)} + E[T_{\text{cloud}}]^{(2)} \cdot \Pi_{\text{SendToCloud}}^{(2)} + E[S_{\text{setup}}] \\
&\hspace{15em} (41)
\end{aligned}$$

### 6.3 Summary of analytical results

All analytical results are been computed using a combination of Java and MATLAB<sup>21</sup> scripts.

In fact, all balance equations, previously shown in tables 2 and 3, are firstly automatically generated using simple Java scripts<sup>22</sup> able to built some MATLAB scripts as output<sup>23</sup> subsequently executed using MATLAB environment.

Is very important to precise that all equation systems are been resolved using a MATLAB function called `solve`,<sup>24</sup> which returns all systems solutions as *symbolic variables* or `sym` objects.<sup>25</sup>

A symbolic variable can hold an expression instead of just a numeric value allowing an higher precision than hardware numbers so is possible to compute more decimal places; obviously the trade off is that symbolic computations need more memory and are slower than hardware computation.

Practically all analytical solutions has been represented by MATLAB as fractions some of which are made up of very big numbers because which has been impossible to fit in this page; just think to symbolic variable holding  $\Pi_{\text{SendToCloud}}^{(2)}$  solution which has a 658 digits number as numerator!

In order to represent analytical solutions we have used two MATLAB functions called `digits` and `digits`<sup>26</sup> and `vpa`<sup>27</sup>; the first function sets the precision used by `vpa` to  $d$  significant decimal digits while the second is used to evaluate symbolic input to at least  $d$  significant digits. Using these functions was possible to list analytical results below.

---

<sup>21</sup>See. <https://www.mathworks.com/products/matlab.html>

<sup>22</sup>See `CTMCResolverScriptGenerator.java`, `ResolverUsingRoutingAlgorithm1.java` and `ResolverUsingRoutingAlgorithm2.java` files.

<sup>23</sup>See `MATLAB_ALG1-CTMC_ResolverScript.m` and `MATLAB_ALG1-CTMC_ResolverScript.m` files

<sup>24</sup>See <https://www.mathworks.com/help/symbolic/solve.html>

<sup>25</sup>See <https://it.mathworks.com/help/symbolic/sym.html>

<sup>26</sup>See <https://it.mathworks.com/help/symbolic/digits.html>

<sup>27</sup>See <https://it.mathworks.com/help/symbolic/vpa.html>

Table 4: Analytical results of Algorithm 1 based system. (Part 1)

Variable	Approximate value	Symbolic variable representation
$\Pi_{SendToCloud}$	0.413925392	$\frac{87859436926256646517251117638424741393475341796875000000}{212259113498085857520971393983611845585800622137044274729}$
$\Pi_{SendToCloudlet}$	0.586074607	$\frac{124399676571829211003720276345187104192325280340169274729}{212259113498085857520971393983611845585800622137044274729}$
$\lambda^{(1)}, X^{(1)}$	4	
$\lambda^{(2)}, X^{(2)}$	6.25	
$\lambda, X$	10.25	
$\lambda_{cloudlet}^{(1)}, X_{cloudlet}^{(1)}$	2.344298428	$\frac{2639445382678245}{1125899906842624}$
$\lambda_{cloudlet}^{(2)}, X_{cloudlet}^{(2)}$	3.662966295	$\frac{65986134566956125}{18014398509481984}$
$\lambda_{cloudlet}, X_{cloudlet}$	6.007264723	$\frac{108217260689808045}{18014398509481984}$
$\lambda_{cloud}^{(1)}, X_{cloud}^{(1)}$	1.655701571	$\frac{1864154244692251}{1125899906842624}$
$\lambda_{cloud}^{(2)}, X_{cloud}^{(2)}$	2.587033704	$\frac{46603856117306275}{18014398509481984}$
$\lambda_{cloud}, X_{cloud}$	4.242735276	$\frac{76430324032382291}{18014398509481984}$
$E[N_{cloudlet}^{(1)}]$	5.204009618	$\frac{1104598468206022474496767594370330477471953336910445636480}{212259113498085857520971393983611845585800622137044274729}$
$E[N_{cloudlet}^{(2)}]$	13.55210838	$\frac{2876558510953183527335332277006068951749878481537618845000}{212259113498085857520971393983611845585800622137044274729}$

Table 5: Analytical results of Algorithm 1 based system. (Part 2)

Variable	Approximate value	Symbolic variable representation
$E[N_{cloudlet}]$	18.75611799	$\frac{3981156979159206001832099871376399429221831818448064481480}{212259113498085857520971393983611845585800622137044274729}$
$E[N_{cloud}^{(1)}]$	6.622806284	$\frac{1864154244692251}{281474976710656}$
$E[N_{cloud}^{(2)}]$	11.759244113	$\frac{1165096402932656875}{99079191802150912}$
$E[N_{cloud}]$	18.382050397	$\frac{1821278697064329227}{99079191802150912}$
$E[N^{(1)}]$	11.826815902	$\frac{13315810923168857}{1125899906842624}$
$E[N^{(2)}]$	25.311352493	$\frac{2507828348524290043}{99079191802150912}$
$E[N]$	37.138168396	$\frac{3679619709763149459}{99079191802150912}$
$E[T_{cloudlet}^{(1)}], E[S_{cloudlet}^{(1)}]$	$2.\bar{2}$	$\frac{20}{9}$
$E[T_{cloudlet}^{(2)}], E[S_{cloudlet}^{(2)}]$	$3.\overline{7037}$	$\frac{100}{27}$
$E[T_{cloudlet}], E[S_{cloudlet}]$	3.125564588	$\frac{3460}{1107}$
$E[T_{cloud}^{(1)}], E[S_{cloud}^{(1)}]$	4	
$E[T_{cloud}^{(2)}], E[S_{cloud}^{(2)}]$	$4.\overline{54}$	$\frac{50}{11}$



Table 6: Analytical results of Algorithm 1 based system. (Part 3)

Variable	Approximate value	Symbolic variable representation
$E[T_{cloud}], E[S_{cloud}]$	4.332594235	$\frac{1954}{451}$
$E[T^{(1)}], E[S^{(1)}]$	2.958089587	$\frac{2497884592968457}{844424930131968}$
$E[T^{(2)}], E[S^{(2)}]$	4.052125751	$\frac{903333025213434725}{222928181554839552}$
$E[T], E[S]$	3.625184809	$\frac{33134390151034630493}{9140055443748421632}$

Table 7: Analytical results of Algorithm 2 based system. (Part 1)

Variable	Approximate value	Symbolic variable representation
$\Pi_{SendToCloud}^{(1)}$	0.0005377649	$\frac{703687441776640000000000000000}{1308540831314824230819390435296941}$
$\Pi_{SendToCloud}^{(2)}$	0.356859509	To big to be represented!
$\Pi_{SendToCloudlet}^{(1)}$	0.999462235	$\frac{1307837143873047590819390435296941}{1308540831314824230819390435296941}$
$\Pi_{SendToCloudlet}^{(2)}$	0.643140490	To big to be represented!
DA VEDERE $\Pi_{Interrupted}^{(2)}$	0.35705151957534072695649962275511	To big to be represented!
$\lambda^{(1)}, X^{(1)}$	4	
$\lambda^{(2)}, X^{(2)}$	6.25	
$\lambda, X$	10.25	
$\lambda_{cloudlet}^{(1)}, X_{cloudlet}^{(1)}$	3.997848940	$\frac{9002355498360251}{2251799813685248}$
$\lambda_{cloud}^{(1)}, X_{cloud}^{(1)}$	0.002151059	$\frac{4960006533878655}{2305843009213693952}$
$E[N_{cloudlet}^{(1)}]$	8.884108755	$\frac{11625219056649311918394581647083920}{1308540831314824230819390435296941}$
$E[N_{cloudlet}^{(2)}]$	9.608670684	To big to be represented!
$E[N_{cloudlet}]$	18.492779440	To big to be represented!
$E[N_{cloud}^{(1)}]$	0.008604239	$\frac{4960006533878655}{576460752303423488}$

Table 8: Analytical results of Algorithm 2 based system. (Part 1)

Variable	Approximate value	Symbolic variable representation
$E[T_{cloud}^{(1)}], E[S_{cloud}^{(1)}]$	4	
$E[T_{cloudlet}^{(1)}], E[S_{cloudlet}^{(1)}]$	$2.\bar{2}$	$\frac{20}{9}$
$E[T^{(1)}], E[S^{(1)}]$	2.223178248	$\frac{46136700210409393015}{20752587082923245568}$