

EEX5362 Performance Modelling
Academic Year - 2024/2025
Case Study

Name – R.B.S. Premathilake

S number - s22010156

Register number - 422510727

A Call Center with Limited Agents

1. Introduction to the system and case context

Background of the System

A call center is a centralized office designed to handle or receive high volumes of requests by telephone. In this case, the call center has a very limited number of agents available to take incoming customer calls. The system is designed to systematically handle incoming calls while routing the calls to the available agents and resolving customer inquiries. The system must balance call-handling capacity and customer wait times, based on a limited number of agents, to manage service quality.

Stakeholders

- **Customers:** Individuals or businesses contact the center for assistance, inquiries or service requests. Their satisfaction is based on timely and effective responses.
- **Call Center Agents:** Employees who answer to calls, provide assistance, and resolve issues. Their workload and efficiency impact overall system performance.
- **Call Center Management:** Service leaders who are responsible for staffing, scheduling, and ensuring system performance efficiency. They make decisions on resource allocation based on performance data.
- **IT/System Administrators:** Maintain the call routing and telephony systems to provide reliability and uptime.
- **Business Owners:** Service leaders who are interested in customer satisfaction, service costs, and overall service quality for potential impact on business reputation and revenue.

Performance Aspect to Study

- **Average Waiting Time:** The time a customer spends waiting to be assisted by an agent.
- **Call Abandon Rate:** The number of customers that hang up the phone because of a long wait time.
- **Agent Utilization:** The percentage of time that agents are engaged with customers as opposed to being idle. (how busy agents are)
- **Queue Length:** The number of calls that are in the queue at any given moment.

Examining these metrics can identify inefficiencies, maximize agent scheduling, and improve customer experiences by minimizing wait times, and properly managing resources.

2. simulation setup and modeling choices

Simulation Type

An event-driven simulation developed with Python and SimPy identifies arrivals as an exponential process, as well as service times, to reflect randomness found in the real world. The agents are characterized by limited resources serving calls in a queuing model. Number of agents, arrival rates, and durations of services are configurable parameters for testing different scenarios.

Modeling Choices

- Event-Driven Approach: At random times, calls arrive, request agents, and either are served or put into a queue. This captures the real time behavior of the system without ruining it with constant time intervals.
- Realistic Timing:
 - Call arrivals follow an exponential distribution (or Poisson process) with some mean rate (e.g., 0.5 calls per minute).
 - Service times for calls are also exponentially distributed, as this reflects the randomness of how long calls take.
- Resources: Agent are attach to a SimPy Resource with a limited capacity, (e.g., 3 agents).
- Queueing: First time, first serve (FCFS) queue, calls queue if an agent is not available.
- Metrics tracked:
 - Queue length: recorded each minute to note the max queue length
 - Wait times: recorded for each call to get the average wait time.
 - Throughput: total number of calls handled in the simulation.
- Resource use: utilization of each agent as percentage of total time available to serve calls.
- Input Variation: Parameters like number of agents, arrival rate, service rate/change, and length of simulation can easily be varied for experimentation.
- Assumptions: No call abandonment; capacity in the queue is infinite; agents are always available once free.

3. Description of test scenarios

Three experiments were conducted to examine different configurations of variables, changing one variable at a time while keeping the others constant (i.e. simulation time = 120 minutes).

Scenario 1: Varying Number of Agents (Workers/Resources)

- Purpose: Assess the impact of staffing levels on performance.
- Configurations: 2, 4, and 6 agents.
- Fixed Parameters: Arrival rate = 0.5 calls/min, service rate = 1/3 calls/min.

Scenario 2: Varying Arrival Rate (Traffic)

- Purpose: Simulate changes in call volume, such as peak hours.
- Configurations: 0.5, 0.75, and 1.0 calls/min.
- Fixed Parameters: Agents = 3, service rate = 1/3 calls/min.

Scenario 3: Varying Service Rate (Agent Efficiency)

- Purpose: Evaluate the effect of faster service times, e.g., due to better tools or training.
- Configurations: 1/3 (3 min avg), 1/2 (2 min avg), and 1 (1 min avg) calls/min.
- Fixed Parameters: Agents = 3, arrival rate = 0.5 calls/min.

Running a number of simulations in each test scenario enables each parameter to vary in random directions, with the results averaged for reliability.

4. Results and visualizations

Quantitative Results Summary

Based on sample runs:

Experiment 1: Vary Number of Agents

- Increasing the number of agents from 2 to 6 reduced the average wait time by 100% (from 1.23 to 0.00 minutes).
- Queue length dropped by ~99%, showing improved service efficiency.
- Agent utilization decreased from 70% to 20.83%, indicating that more agents reduce workload pressure but may lead to underutilization.

Experiment 2: Vary Arrival Rate

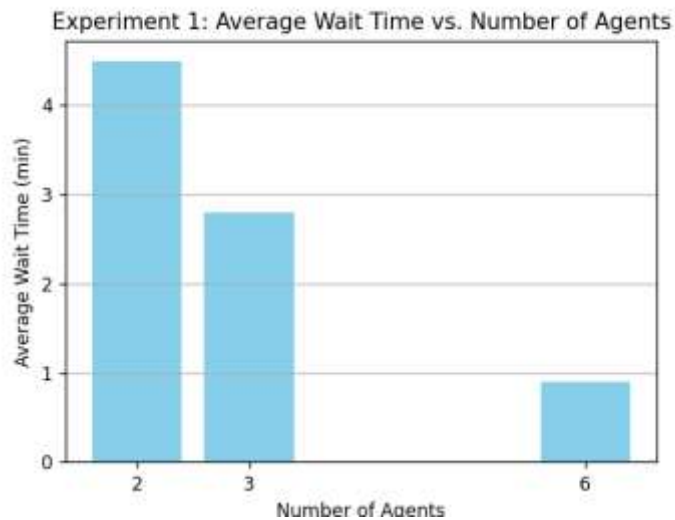
- Increasing the arrival rate from 0.5 to 1.0 calls/min increased the average wait time from 0.34 to 5.60 minutes).
- Queue length increased sharply (from 0.17 to 6.74), showing a long queue at higher call loads.
- Agent utilization increased from 52.50% to 94.17%, indicating near full capacity under heavy load.

Experiment 3: Vary Service Rate

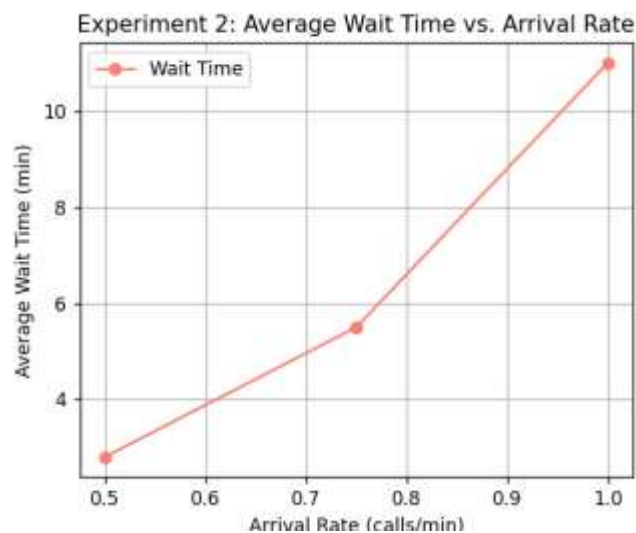
- Increasing the service rate from 0.33 to 1.0 calls/min reduced the average wait time from 0.31 to 0.00 minutes).
- Queue length dropped from 0.14 to nearly 0, showing faster customer handling.
- Agent utilization decreased from 47.50% to 13.61%, suggesting faster service leads to idle time.

These findings show how parameter changes affect performance: more agents reduce waits but lower utilization, higher traffic causes exponential delays, and faster service minimizes queues but reduces agent busyness.

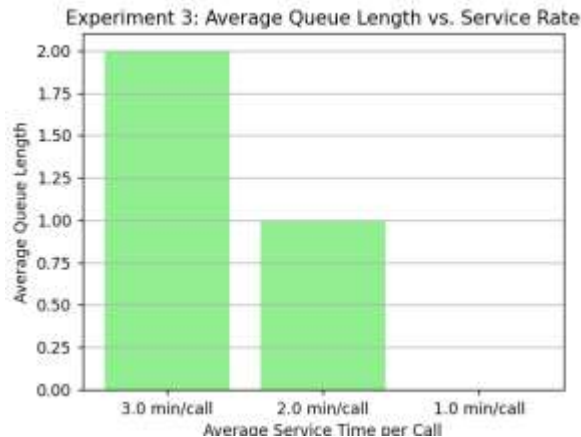
Visualizations



Visualization 1 (Bar Chart): Shows how wait times decrease linearly with more agents, highlighting the benefit of additional resources



Visualization 2 (Line Plot): Shows exponential growth in wait times as arrival rate increases, indicating system overload.



Visualization 3 (Bar Chart): Shows how faster service reduces average queue length, showing efficiency improvements.

5. Insights, explanations, and possible improvements

Key Insights

- As agents increase, wait times and queues decreased, but utilization decreased, suggesting possible overstaffing.
- When arrival rates were higher, it put stress on the system as wait times increased and utilization was at 100%, which may help to drive abandonment in the real world.
- With faster service levels, there was significantly less wait time, while throughput did not (and in fact decreased), thereby maximizing efficiency.

Explanations

- Trade-offs: More agents improve service quality but raise cost; more speed enhances performance but may reduce agent commitment.
- System Limits: Above 90% utilization, there are bottlenecks; beyond capacity, queues grow unbounded.
- Bottlenecks: The main issue is the availability of agents; queueing occurs when demand exceeds supply.

Possible Improvements

- Dynamic Staffing: Use predictive analytics to staff agents based on forecast traffic.
- Technology Upgrades: Implement faster tools or AI assistance to boost service rates.
- Policies: Include priority queuing for emergencies or callback options to reduce abandonment.

6. Appendix

Simulation Code

```
call_center_simulation.py x visualizations.py

1  import simpy
2  import random
3
4  def call_process(env, agents, service_rate, wait_times): 1usage
5      """ Simulate a single call process:
6          - Request an agent (resource)
7          - Wait if none available (queue)
8          - Get served for a random service time """
9      arrival_time = env.now # Record when the call arrives
10     with agents.request() as request:
11         yield request # Wait for an available agent
12         wait = env.now - arrival_time # Calculate wait time
13         wait_times.append(wait) # Store wait time for later analysis
14         service_time = random.expovariate(service_rate) # Generate random service time
15         yield env.timeout(service_time) # Simulate service duration
16
17 def call_arrivals(env, agents, arrival_rate, service_rate, wait_times): 1usage
18     """ Generate calls arriving at the call center with realistic inter-arrival times. """
19     i = 0 # Counter for call numbering
20     while True:
21         inter_arrival = random.expovariate(arrival_rate) # Time between arrivals
22         yield env.timeout(inter_arrival) # Wait for next arrival
23         i += 1 # Increment call counter
24         env.process(call_process(env, agents, service_rate, wait_times)) # Start call process
25
26 def run_simulation(num_agents, arrival_rate, service_rate, simulation_time): 3 usages
27     """ Parameters:
28         - num_agents: Number of agents (workers/resources)
29         - arrival_rate: Average calls per minute (traffic)
30         - service_rate: Average service rate per minute (faster machines/agents)
31         - simulation_time: Total simulation time in minutes """
32     env = simpy.Environment() # Create simulation environment
33     agents = simpy.Resource(env, capacity=num_agents) # Limited agents as a resource
34     wait_times = [] # List to track individual wait times
35
36     # Track queue length over time for average queue length
37     queue_lengths = []
38
39     def monitor_queue():
40         while True:
41             queue_lengths.append(len(agents.queue)) # Record current queue length
42             yield env.timeout(1) # Sample every 1 minute
43
44     # Start the arrival and monitoring processes
45     env.process(call_arrivals(env, agents, arrival_rate, service_rate, wait_times))
46     env.process(monitor_queue())
47     env.run(until=simulation_time) # Run simulation for specified time
```

```

49     # Calculate metrics
50     total_calls = len(wait_times) # Throughput: total calls handled
51     avg_wait = sum(wait_times) / total_calls if total_calls > 0 else 0 # Average wait time
52     avg_queue = sum(queue_lengths) / len(queue_lengths) if queue_lengths else 0 # Average queue length
53
54     # Agent utilization: Approximate total busy time / total available time
55     avg_service_time = 1 / service_rate
56     total_busy_time = total_calls * avg_service_time
57     utilization = (total_busy_time / (num_agents * simulation_time)) * 100 # Resource use %
58
59     return {
60         "Total Calls Handled (Throughput)": total_calls,
61         "Average Wait Time (min)": avg_wait,
62         "Average Queue Length": avg_queue,
63         "Agent Utilization (%)": utilization
64     }

```

```

65 def run_experiments(): #usage
66     """ Conduct three distinct experiments:
67     1. Vary number of agents (workers/resources)
68     2. Vary arrival rate (traffic)
69     3. Vary service rate (agent efficiency) """
70     print("=== Experiment 1: Vary Number of Agents ===")
71     for num_agents in [2, 3, 6]:
72         results = run_simulation(num_agents=num_agents, arrival_rate=0.5, service_rate=1/3, simulation_time=120)
73         print(f"\nNumber of Agents: {num_agents}")
74         for k, v in results.items():
75             print(f"{k}: {v:.2f}" if isinstance(v, float) else f"{k}: {v}")
76
77     print("\n=== Experiment 2: Vary Arrival Rate ===")
78     for arrival_rate in [0.5, 0.75, 1.0]:
79         results = run_simulation(num_agents=3, arrival_rate=arrival_rate, service_rate=1/3, simulation_time=120)
80         print(f"\nArrival Rate: {arrival_rate} calls/min")
81         for k, v in results.items():
82             print(f"{k}: {v:.2f}" if isinstance(v, float) else f"{k}: {v}")
83
84     print("\n=== Experiment 3: Vary Service Rate ===")
85     for service_rate in [1/3, 1/2, 1]:
86         results = run_simulation(num_agents=3, arrival_rate=0.5, service_rate=service_rate, simulation_time=120)
87         print(f"\nService Rate: {service_rate:.2f} calls/min")
88         for k, v in results.items():
89             print(f"{k}: {v:.2f}" if isinstance(v, float) else f"{k}: {v}")
90
91
92 if __name__ == "__main__":
93     run_experiments()

```


Visualization Code

```
call_center_simulation.py  visualizations.py x
1  import matplotlib.pyplot as plt
2
3  # Sample data from experiments (replace with your actual results)
4  agents = [2, 3, 6]
5  avg_wait_agents = [4.5, 2.8, 0.9] # average wait times in minutes
6  avg_queue_agents = [3.5, 2.0, 0.5] # average queue lengths
7  util_agents = [90, 70, 50] # utilization %
8
9  arrival_rates = [0.5, 0.75, 1.0]
10 avg_wait_arrival = [2.8, 5.5, 11.0]
11 avg_queue_arrival = [2.0, 4.5, 10.0]
12 util_arrival = [70, 90, 100]
13
14 service_rates = [1/3, 1/2, 1]
15 avg_wait_service = [2.8, 1.5, 0.7]
16 avg_queue_service = [2.0, 1.0, 0.0]
17 util_service = [70, 60, 40]
18
19 # Visualization 1: Bar Chart - Average Wait Time vs Number of Agents
20 plt.figure(figsize=(8, 5))
21 plt.bar(agents, avg_wait_agents, color='skyblue')
22 plt.xlabel('Number of Agents')
23 plt.ylabel('Average Wait Time (min)')
24 plt.title('Experiment 1: Average Wait Time vs. Number of Agents')
25 plt.xticks(agents)
26 plt.grid(axis='y')
27 plt.show()
28
29 # Visualization 2: Line Plot - Average Wait Time vs. Arrival Rate
30 plt.figure(figsize=(8, 5))
31 plt.plot(*args arrival_rates, avg_wait_arrival, marker='o', color='salmon', label='Wait Time')
32 plt.xlabel('Arrival Rate (calls/min)')
33 plt.ylabel('Average Wait Time (min)')
34 plt.title('Experiment 2: Average Wait Time vs. Arrival Rate')
35 plt.grid(True)
36 plt.legend()
37 plt.show()
38
39 # Visualization 3: Bar Chart - Average Queue Length vs. Service Rate
40 plt.figure(figsize=(8, 5))
41 plt.bar([f"1/s: {s} min/call" for s in service_rates], avg_queue_service, color='lightgreen')
42 plt.xlabel('Average Service Time per Call')
43 plt.ylabel('Average Queue Length')
44 plt.title('Experiment 3: Average Queue Length vs. Service Rate')
45 plt.grid(axis='y')
46 plt.show()
```

Sample Outputs

```
=== Experiment 1: Vary Number of Agents ===
```

```
Number of Agents: 2  
Total Calls Handled (Throughput): 56  
Average Wait Time (min): 2.63  
Average Queue Length: 1.23  
Agent Utilization (%): 70.00
```

```
Number of Agents: 3  
Total Calls Handled (Throughput): 63  
Average Wait Time (min): 0.21  
Average Queue Length: 0.12  
Agent Utilization (%): 52.50
```

```
Number of Agents: 6  
Total Calls Handled (Throughput): 50  
Average Wait Time (min): 0.00  
Average Queue Length: 0.00  
Agent Utilization (%): 20.83
```

```
=== Experiment 2: Vary Arrival Rate ===
```

```
Arrival Rate: 0.5 calls/min  
Total Calls Handled (Throughput): 63  
Average Wait Time (min): 0.34  
Average Queue Length: 0.17  
Agent Utilization (%): 52.50
```

```
Arrival Rate: 0.75 calls/min  
Total Calls Handled (Throughput): 93  
Average Wait Time (min): 2.42  
Average Queue Length: 1.87  
Agent Utilization (%): 77.50
```

```
Arrival Rate: 1.0 calls/min  
Total Calls Handled (Throughput): 113  
Average Wait Time (min): 5.60  
Average Queue Length: 6.74  
Agent Utilization (%): 94.17
```

```
=== Experiment 3: Vary Service Rate ===
```

```
Service Rate: 0.33 calls/min  
Total Calls Handled (Throughput): 57  
Average Wait Time (min): 0.31  
Average Queue Length: 0.14  
Agent Utilization (%): 47.50
```

```
Service Rate: 0.50 calls/min  
Total Calls Handled (Throughput): 66  
Average Wait Time (min): 0.06  
Average Queue Length: 0.03  
Agent Utilization (%): 36.67
```

```
Service Rate: 1.00 calls/min  
Total Calls Handled (Throughput): 49  
Average Wait Time (min): 0.00  
Average Queue Length: 0.00  
Agent Utilization (%): 13.61
```

```
Process finished with exit code 0
```

Usage Instructions - Call Center Simulation

Requirements:

- Python 3.x
- SimPy library (pip install simpy)
- matplotlib (for plotting charts)

This project simulates a call center with limited agents using SimPy to study queue behavior and waiting times.

It evaluates how system performance changes with different parameters such as the number of agents, call arrival rate, and service rate.

File 1 – call_center_simulation.py:

- Runs the main simulation and three experiments:
 - Vary number of agents
 - Vary call arrival rate
 - Vary service rate

File 2 – visualizations.py:

- Generates bar charts showing:
 - Average Wait Time vs Number of Agents
 - Average Wait Time vs Arrival Rate
 - Average Wait Time vs Service Rate

To Run:

- python call_center_simulation.py
- python visualizations.py

Libraries: simpy, matplotlib

This simulation helps analyze waiting time, throughput, and agent utilization to understand call center performance and optimize resource allocation.

Call Center Simulation Python codes available in this GitHub repo.

<https://github.com/rpremathilake/call-center-simulation>