**Name**: Tural Mehtiyev

**Project Title**: Scalability analysis of microservice architecture on a particular software application

**1. What are you going to do?**

I'm going to conduct a detailed analysis of a microservice architecture system to assess its scalability on a particular software application. Scalability in this context means the ability of the system to handle increased load without a detrimental impact on performance, or the capacity to expand to accommodate growth. This includes an examination of how well the system can handle increasing load, the strategies it uses to scale, and areas where scalability can be improved.

Below is the initial list of the scalability attributes to be used in the project:

➢ Performance: This refers to the ability of the architecture to maintain or improve its functioning as the load increases.

➢ Load Scalability: This measures how well the architecture can handle an increase in concurrent tasks or processes, including the number of user requests or transactions it can manage effectively.

➢ Resource Utilization: It's essential to understand how efficiently the architecture uses available resources (CPU, memory, storage, network bandwidth, etc.) when scaling up or out.

➢ Capacity: This refers to the maximum load that the system can handle under specific conditions. This could involve the number of microservices it can manage, the number of concurrent users it can support, the volume of data it can process, etc.

➢ Fault Tolerance and Resilience: Microservices architectures are designed to continue functioning even when individual services fail. This attribute refers to how well the architecture handles such failures and whether it can still scale under these conditions.

The following are key metrics to consider when assessing the scalability of a software architecture.

➢ Response Time: This measures how quickly the system responds to a user's request. In a scalability context, we want to observe how this metric changes as load increases. A scalable system should maintain relatively stable response times even under increased load.

➢ Throughput: This measures the number of requests the system can handle per unit of time, typically expressed as Requests Per Second (RPS) or Transactions Per Second (TPS). In terms of scalability, we want to see this number increase proportionately to increases in load or resources.

➢ Resource Usage: This metric monitors the consumption of resources like CPU, memory, disk I/O, and network I/O at different load levels. Efficient resource usage is a key element of scalability, as it means that adding resources will result in proportional increases in capacity.

➢ Error Rates: This represents the percentage of requests that result in errors. High or increasing error rates under load can indicate scalability problems, as it suggests that the system is struggling to handle the increased demand.

➢ Latency: This refers to the delay before a transfer of data begins following an instruction for its transfer. In distributed systems like microservices, minimizing latency is crucial for maintaining performance as the system scales.

➢ System Utilization: This measures how well the system's capacity is being used, often broken down into CPU utilization, memory utilization, etc. A system that can maintain high utilization rates as it scales is using its resources efficiently.

➢ Service Time: This is the time taken by the system to complete a request once it starts processing it. It's a measure of the system's internal efficiency and can provide insights into potential bottlenecks.

➢ Scalability Ratio: This is the ratio of the throughput achieved by adding more resources to the system to the throughput of a single resource. It can provide a direct measure of the system's scalability.

➢ Queue Length: This measures the number of incoming requests waiting to be processed. A consistently long queue could mean that the system is unable to keep up with the incoming load, indicating scalability issues.

➢ Availability: This refers to the ability of the system to remain available and operational as it scales. High availability is crucial for ensuring a consistent user experience, especially in a microservices architecture where the failure of a single service can impact the entire system.

## 2. How is it done today? Current Limitations?

Today, scalability analysis is usually done through a mix of load testing, code reviews, and architectural analysis. Current limitations include the time and expertise required to set up and interpret load tests, difficulty predicting future loads, and the challenge of identifying bottlenecks in complex systems. In addition, microservices, due to their distributed nature, introduce additional complexity in managing and coordinating different services, data consistency, network latency, and fault tolerance, which can also affect scalability.

## 3. What is your idea to do something better?

My plan is to develop a comprehensive methodology that combines empirical testing with theoretical analysis. This includes setting up automated load tests that simulate different types of load increase, performing a detailed review of the system's code and architecture to identify potential bottlenecks. The results of this analysis will be compiled into a detailed report that includes both specific, actionable recommendations for improving scalability and general insights that can be applied to other microservice architectures.

## 4. Who will benefit from your work? Why?

Both developers and organizations can benefit from this work. Developers can use the insights gained from this analysis to design more scalable microservice architectures, while organizations can benefit from the improved capacity planning and potentially lower infrastructure costs that result from better scalability. In addition, better scalability can improve user experience by ensuring that the system remains responsive even under high load.

## 5. What risks do you anticipate?

➢ Representativeness of the chosen system: There's a risk that the chosen system might not be representative of other microservice architectures, limiting the applicability or generalizability of the findings. Different systems can have unique characteristics that impact scalability in ways not observed in the chosen system.

➢ Incomplete identification of scalability issues: Load tests and analysis might not reveal all potential scalability issues, particularly those that only manifest under very specific or extreme conditions. This might lead to an overly optimistic assessment of the system's scalability.

➢ Quality and availability of documentation: Finding a system with comprehensive, accurate, and up-to-date documentation can be challenging. If the system's documentation is incomplete, outdated, or not well-structured, this could complicate the analysis process and potentially lead to misunderstandings about the system's design and behavior.

➤ Technical complexities and unknowns: Given the complex nature of microservices-based architectures, there might be technical challenges and unforeseen issues that arise during the analysis. This could impact the project timeline or the depth of analysis possible within the project scope.

**6. Out of pocket costs? Complete within 11 weeks?**

There may be minimal out-of-pocket costs for cloud computing resources if load tests are performed in a cloud environment. With good project planning and a clear focus, the project can likely be completed within 11 weeks.

**7. Midterm results?**

By the midterm point, I expect to have completed the setup of the load tests, run initial tests, and started the code and architectural analysis. I should have some preliminary findings and be able to identify potential areas for improvement.

**8. Final Demonstration?**

The final demonstration will include a presentation of the findings and recommendations from the scalability analysis. I will also demonstrate the load testing setup and discuss how the tests were used to assess scalability. If possible, I may also show before-and-after performance results to illustrate the impact of any recommended changes that were implemented.