# Cloud Computing

1.  Why the Auto Scaling feature is significant in cloud computing?

    Auto Scaling is a crucial feature in cloud computing because it allows applications to automatically adjust their capacity based on traffic and workload demands. This capability ensures that resources are optimally utilized, costs are minimized, and performance remains consistent even during fluctuations in usage.

    Here's, why Auto Scaling is significant:

    ➢ Optimized Resource Usage: Auto Scaling ensures that you have the right amount of resources at any given time, scaling up when demand increases and scaling down during periods of low activity, thus optimizing resource allocation.

    ➢ Cost Efficiency: By scaling resources dynamically, Auto Scaling helps in cost management by only using what is needed at any moment, reducing unnecessary expenditure on idle resources.

    ➢ Improved Availability and Reliability: Applications can maintain consistent performance and availability by automatically scaling to meet varying workload demands, reducing the risk of downtime due to overload or insufficient capacity.

    ➢ Elasticity: Auto Scaling provides elasticity to applications, allowing them to quickly respond to changes in demand without manual intervention, ensuring that performance is maintained under varying conditions.

    ➢ Enhanced User Experience: Users experience consistent performance levels even during peak usage times, ensuring a smooth and reliable experience.

    ➢ Operational Simplicity: It simplifies operations by automating the scaling process based on predefined policies or metrics, reducing the need for manual intervention and allowing teams to focus on other aspects of application management.

2.  Describe the following

    1.  Metric-based autoscaling: Metric-based autoscaling, also known as dynamic autoscaling, adjusts the number of instances or resources allocated to an application based on realtime metrics and performance indicators. This approach allows systems to scale automatically in response to changes in demand, ensuring optimal performance and resource utilization at all times.

        Key points:
         Realtime Monitoring: Constantly monitors metrics such as CPU utilization, memory usage, network traffic, or custom application-specific metrics.

Thresholds and Policies: Defines thresholds (e.g., CPU > 70%) and scaling policies (e.g., add 2 instances when CPU exceeds threshold for 5 minutes).
Automated Scaling Actions: Automatically adds or removes instances, adjusts resource allocation (CPU, memory), or scales horizontally based on predefined rules.
Examples: Used in web applications, APIs, and services that experience fluctuating traffic patterns throughout the day or in response to external events (e.g., promotions, news events).

2. Schedule-based autoscaling: Schedule-based autoscaling involves scaling resources up or down based on a predefined schedule or timebased rules rather than realtime metrics. This approach is useful for applications or workloads with predictable usage patterns or when demand varies predictably at certain times.
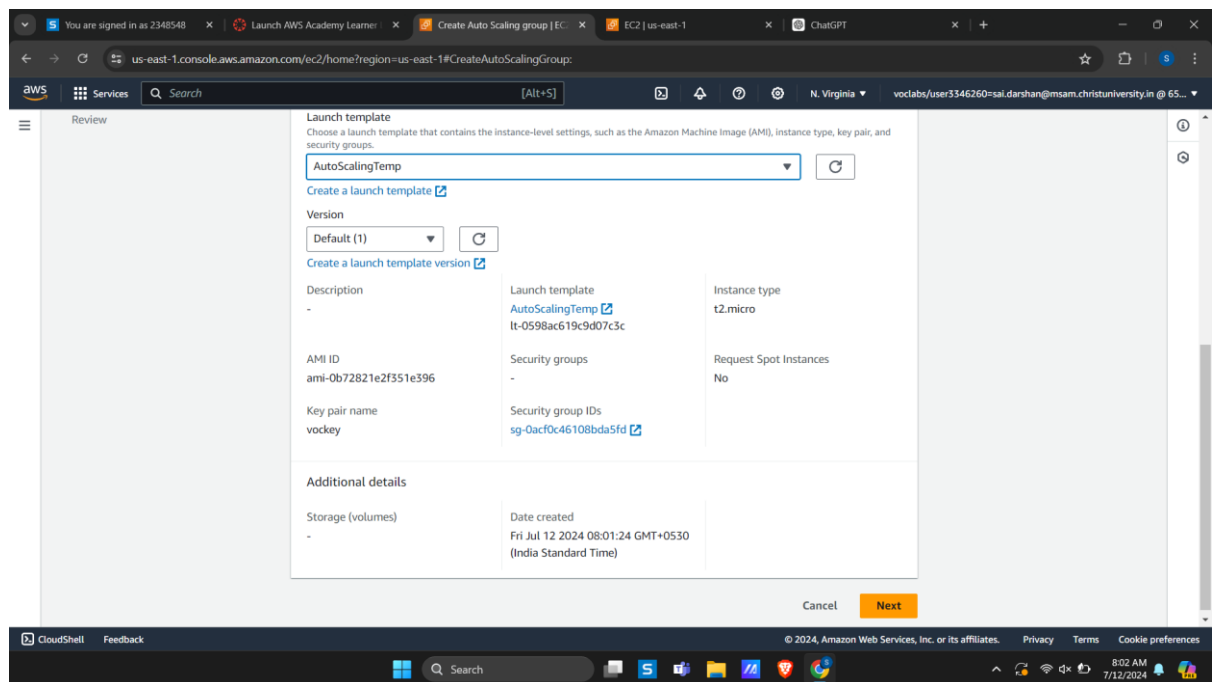
   Key points:
   Predefined Scaling Events: Scales resources at specific times or intervals (e.g., weekdays at 9 AM, weekends after midnight).
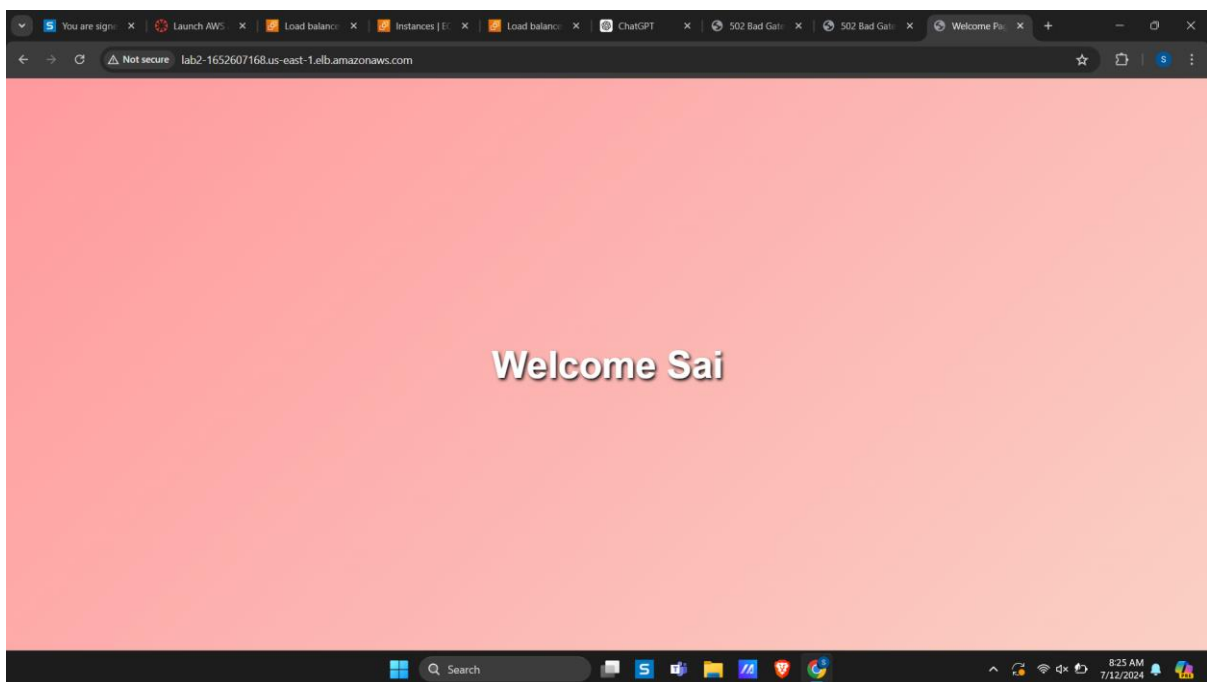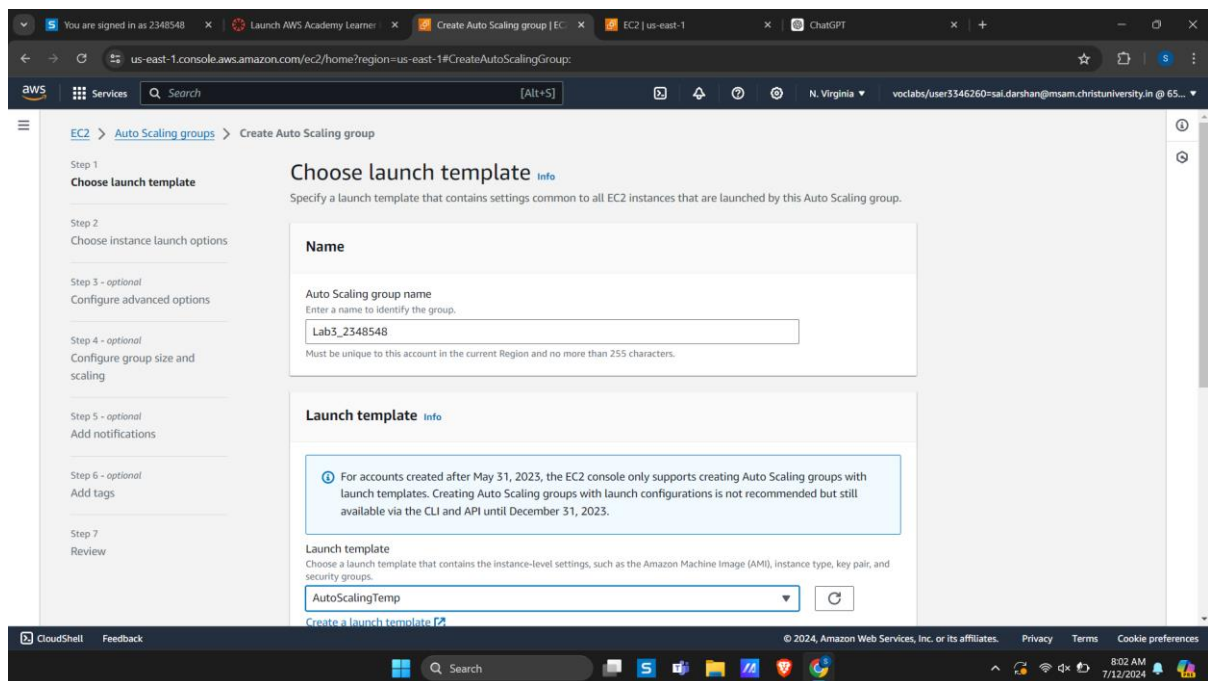   Predictable Workloads: Suitable for applications with known peak times or recurring periods of high and low demand.
   Cost Efficiency: Ensures resources are available when needed without overprovisioning during periods of low demand, optimizing cost.
   Examples: Used in batch processing jobs, financial reporting systems, or applications with scheduled data processing tasks.

3. Demonstrate Metric based autoscaling or Schedule based auto scaling to cater your organizations business requirements. (Specify the requirements)

**Financial Services and Trading Platforms:**

Financial services rely on real-time data processing, trading activities, and customer transactions that can vary significantly throughout the trading day. Auto-scaling enables these platforms to scale compute resources dynamically based on market conditions, transaction volumes, and data analytics needs. It ensures responsiveness, data accuracy, and operational resilience.