**Steps to Set Up an AWS EMR Cluster**
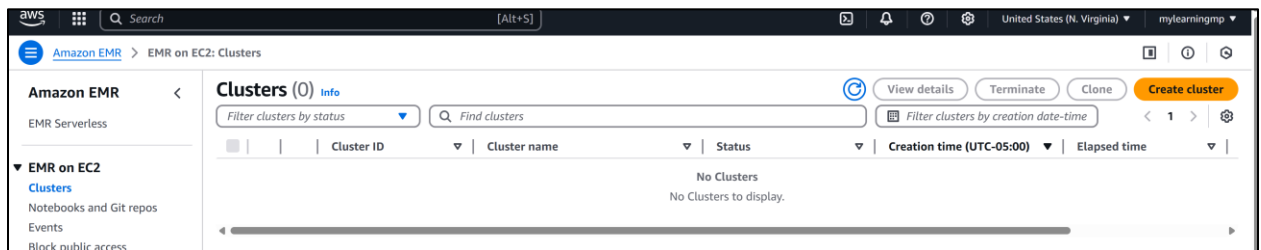
EMR Cluster can be setup from AWS Console, SDK or AWS CLI. For the blog we will cover setup via **AWS Console**.

**Step 1: Log in to AWS Console**

1. Login to AWS Console.

2. Once logged in, search for "EMR" in the search bar and navigate to **Amazon EMR.**
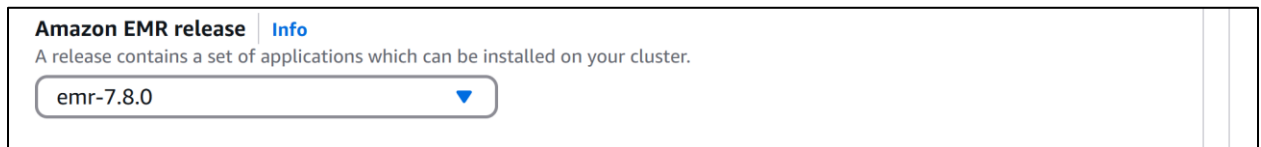
**Step 2: Create EMR Cluster.**

1. In the EMR Cluster, click on "Create Cluster"
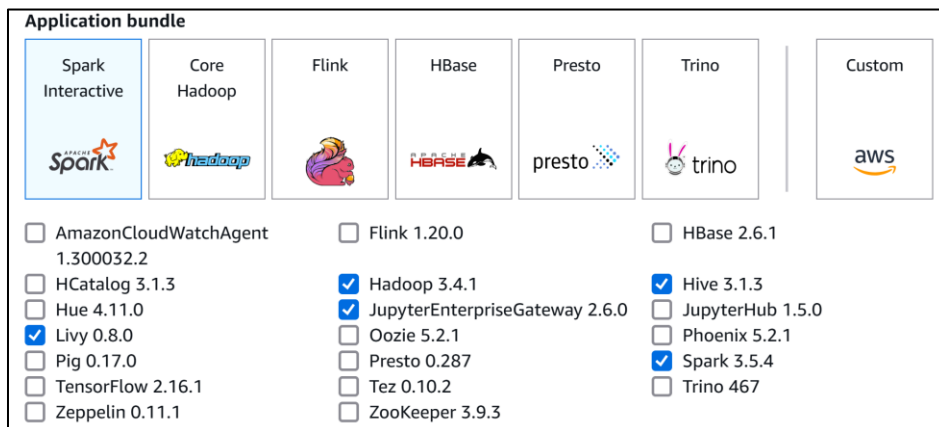


**Step 3: Configure the Software and Applications:**

1. Choose the EMR Release version: We are going to choose **emr-7.8.0**



2. Choose the Application bundle: From the available choices we choose "Spark". Checkmarks are chosen automatically based on the application choice you make.

## Step 4: Configure Instance Types and Cluster Size :

Choose a configuration method for the primary, core, and task node groups for your cluster.

Chose the Primary or Master Node:

1. **Choose EC2 Instance Types**:
    1. Master Node: Choose an instance type like m5-xlarge
    2. Core Nodes : Something like r5.xlarge for memory intensive applications like Spark.
    3. Task Nodes : If your workload needs extra processing power without additional storage , make use of Task node for optimizing Cluster performance.
2. Cluster Sizing: For optimal performance and cost-effectiveness, tailor your cluster size to your specific needs. Define the initial number of instances for each node type, considering both your data volume and processing intensity. EMR's built-in auto-scaling functionality provides the flexibility to automatically scale the cluster up or down based on real-time workload fluctuations, preventing unnecessary resource overhead.

**Step 5: Enabling EMR-Managed Scaling for Auto-Scaling**

Chosen below is the "Use EMR-managed scaling", since we want EMR to dynamically auto scale based on workload. Here you can also specify the Min and Max Cluster Size and also choose Max Core Nodes and On-Demand instances in cluster.
You can also specify number of instance that the cluster starts with under the "Provisioning Configuration" setting.

▼ **Cluster scaling and provisioning - *required*** Info
Choose how Amazon EMR should size your cluster.

Choose an option

| ○ Set cluster size manually | ● Use EMR-managed scaling | ○ Use custom automatic scaling |
|---|---|---|
| Use this option if you know your workload patterns in advance. | Monitor key workload metrics so that EMR can optimize the cluster size and resource utilization. | To programmatically scale core and task nodes, create custom automatic scaling policies. |

**Scaling configuration**

**Minimum cluster size**                           **Maximum cluster size**

| 2 | instance(s) | | 20 | instance(s) |

**Maximum core nodes in the cluster**
Limit the number of core nodes in your cluster.

| 20 | instance(s) |

**Maximum On-Demand instances in the cluster**
To provision the primary node to use On-Demand pricing and other nodes in the cluster to use Spot pricing, set this value to 1. To provision the entire cluster to use On-Demand pricing, use the same value as your maximum cluster size.

| 20 | instance(s) |

**Provisioning configuration**

Set the size of your core and task instance groups. Amazon EMR attempts to provision this capacity when you launch your cluster.

| Name | Instance type | Instance(s) size | Use Spot purchasing option |
|---|---|---|---|
| Core | m5.xlarge | 1 | ☐ |
| Task - 1 | m5.xlarge | 1 | ☐ |

**Step 6: Networking and Security Configuration**

1. **VPN and Network Configuration:** Choose the appropriate VPC and Firewall for your EMR Setup.
2. **IAM Role configuration :**
   Choose or create new IAM Role granting necessary permissions to the clusters to interact with other AWS services like S3, Dynamo DB etc.
      3. EMR Default Role: Allows cluster to interact with S3 and Dynamobo DB and other services.
      4. EMR Default EC2 Role : Grants EC2 instances in the cluster necessary permission

   In our case we choose Amazon to create the Service Roles and Atoscaling Role for EC2. Please see our preferences below.



**Step 7 : Logs:**

Configure AWS S3 location for logging.

**Step 8: Review and Create Cluster:**

1. Review all configuration.
2. After a through review hit the "Create Cluster"



**Step 9: Manage and Monitor the Cluster.**

1. Once your cluster is up and running, you can monitor its status via the **EMR Dashboard**.
2. Use **AWS CloudWatch** for real-time monitoring of cluster metrics, resource utilization, and job progress.
3. You can also SSH into the master node if you need to directly manage or troubleshoot the cluster.