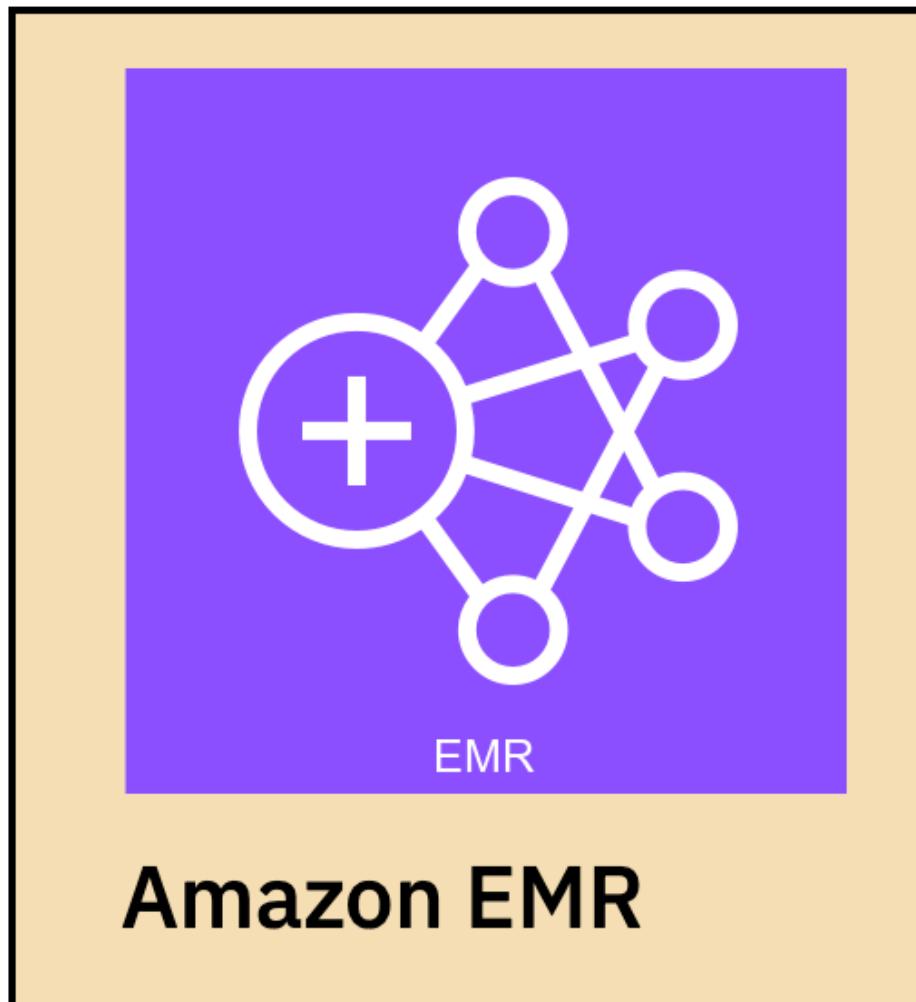




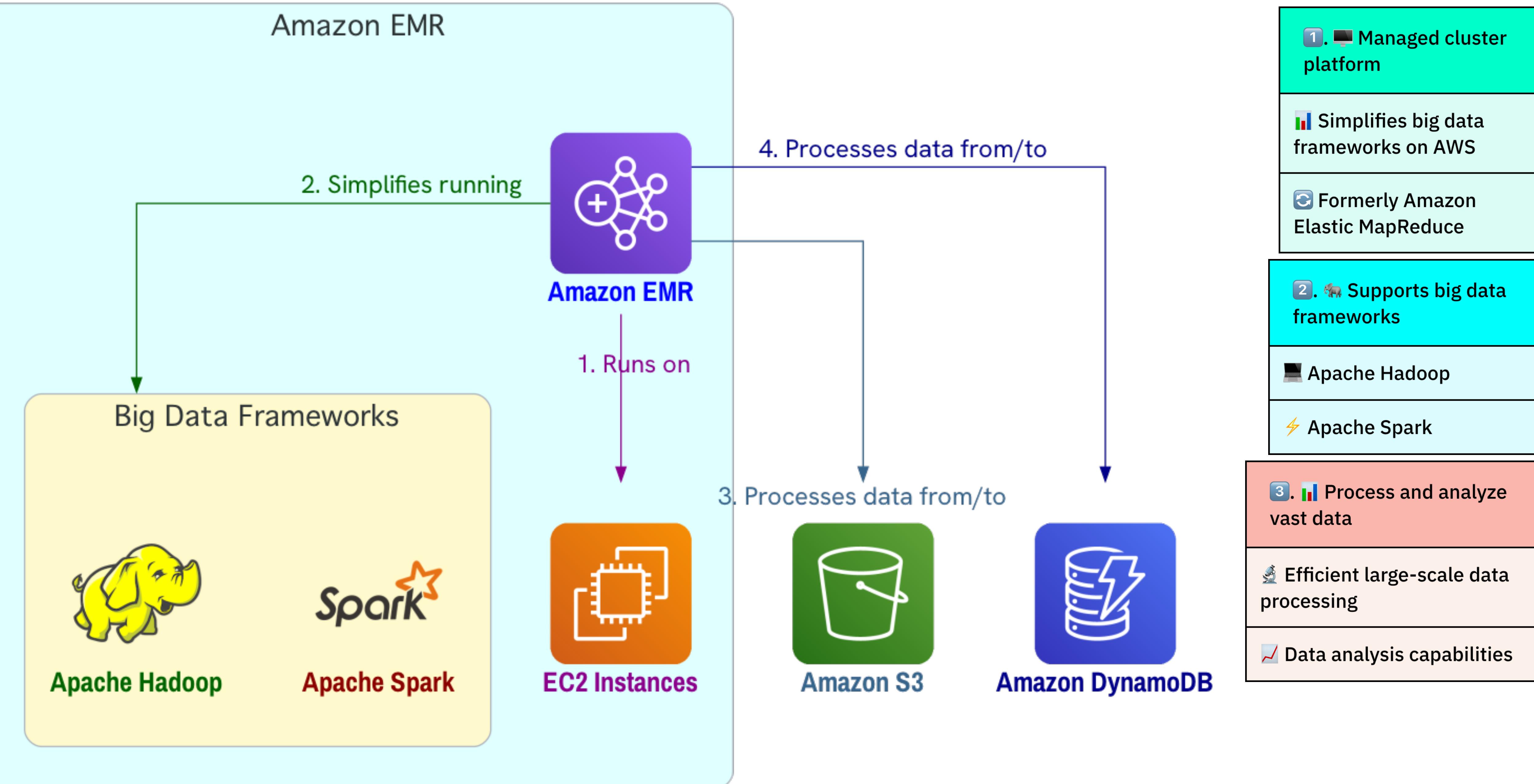
Amazon EMR

Table of Contents

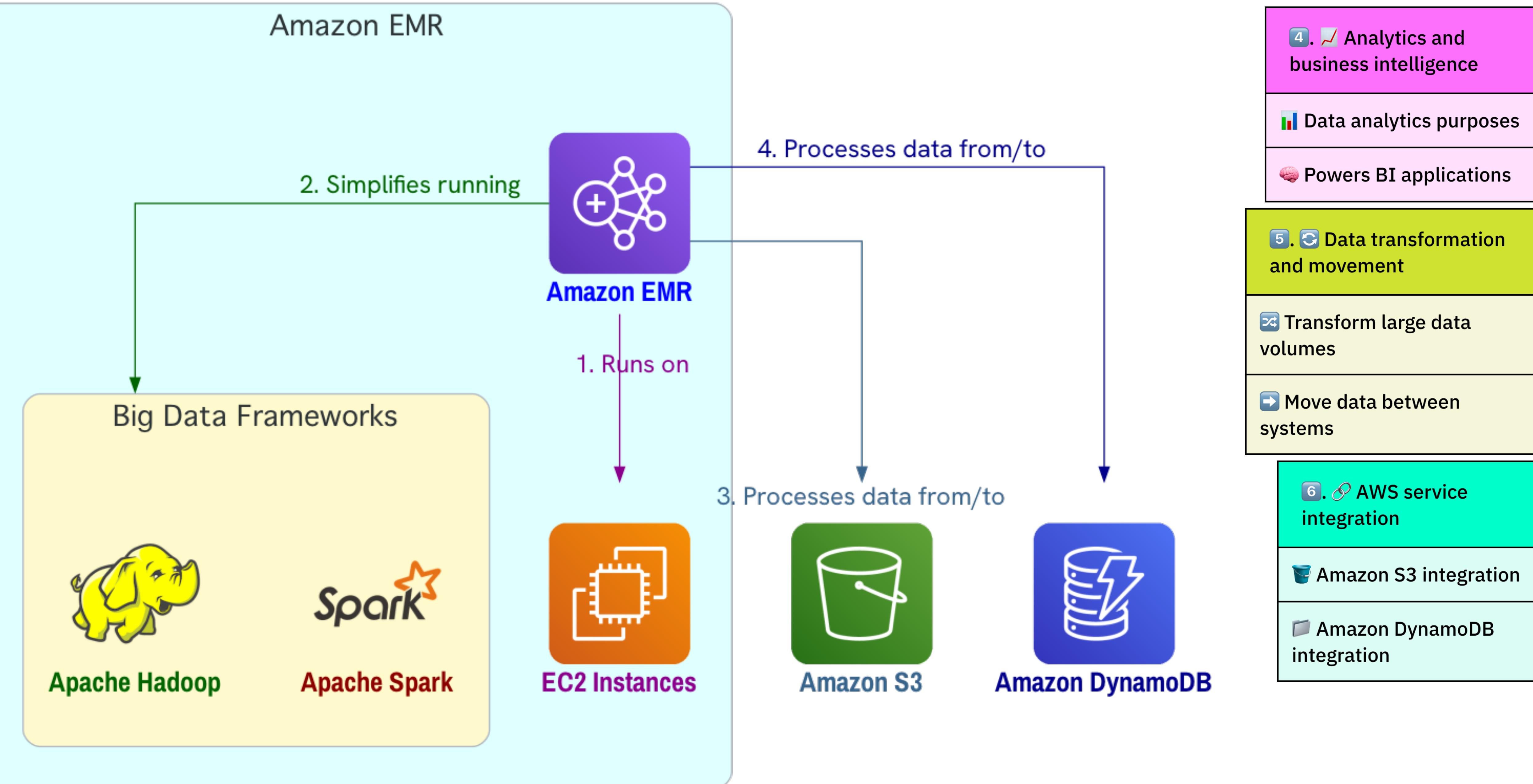


- 1. What is Amazon EMR?
- 2. Understanding clusters and nodes
- 3. Submitting work to a cluster
- 4. Running Steps to Process Data in Amazon EMR
- 5. Amazon EMR Cluster Lifecycle
- 6. Amazon EMR Cluster Failure Handling
- 7. AWS Integration
- 8. Architecture

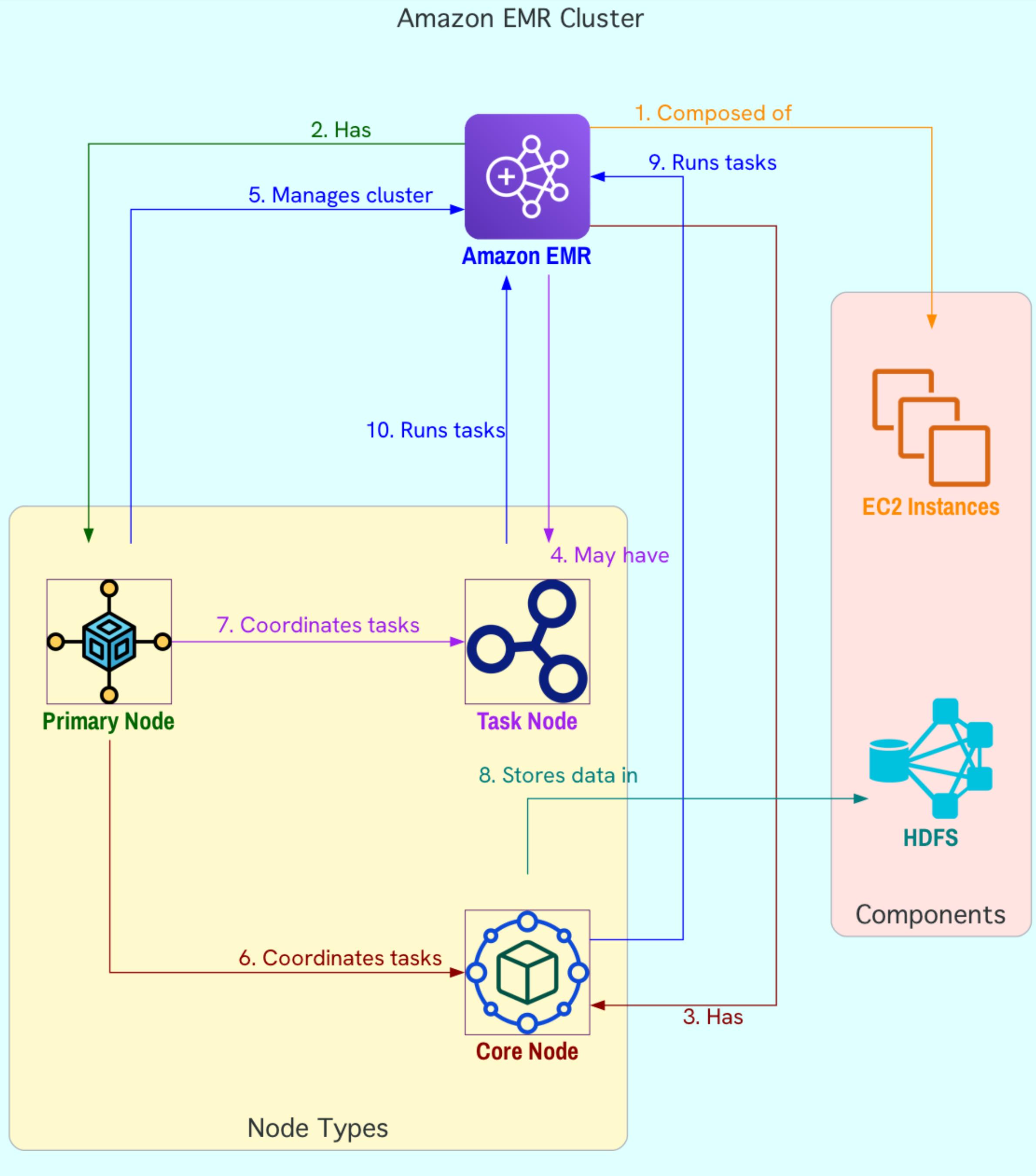
What is Amazon EMR?



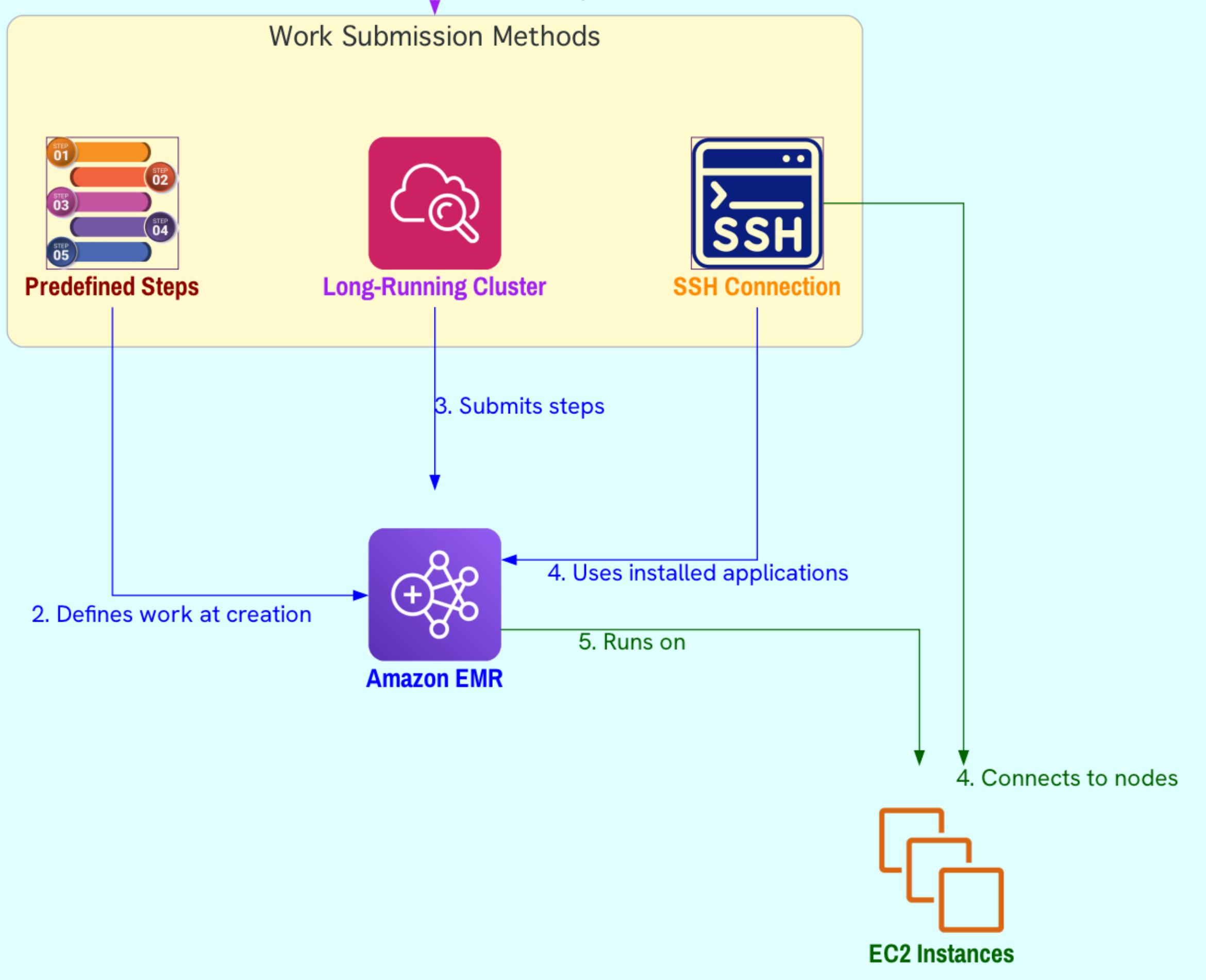
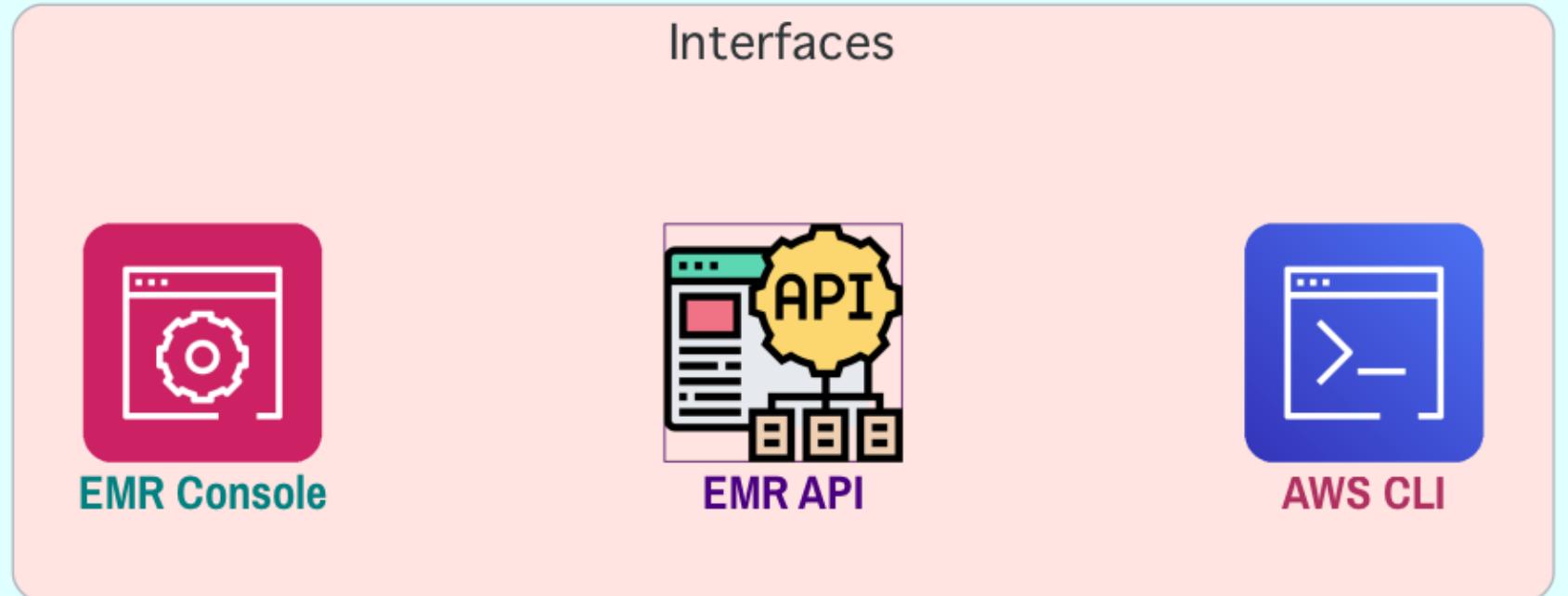
What is Amazon EMR?



Understanding clusters and nodes



1. 🖥️ Cluster: Collection of EC2 instances	2. 🏠 Node: Individual instance in a cluster	3. 🎭 Node types: Different roles in cluster
Group of EC2 instances	Building blocks of cluster	Assigned roles within cluster
Central component of EMR	Specific role in distributed application	Specific software components installed
Enables distributed big data processing	Example: Apache Hadoop	
4. 🧠 Primary node: Manages cluster	5. 💪 Core node: Runs tasks & stores data	6. 💨 Task node: Runs tasks only (optional)
Coordinates data and task distribution	Runs tasks	Provides additional compute capacity
Tracks task status	Stores data in HDFS	🚫 Does not store data in HDFS
Monitors cluster health	At least one in multi-node clusters	⚖️ Increases processing without adding storage
1 One per cluster (minimum)		



Submitting work to a cluster

1. Define work as steps during cluster creation

Specify functions as steps

Process set amount of data

Terminate after completion

Ideal for batch processing

2. Submit steps to long-running clusters

Amazon EMR console

Amazon EMR API

AWS CLI

Steps contain one or more jobs

Flexible for ongoing workloads

3. Connect via SSH and use application interfaces

SSH to primary and other nodes

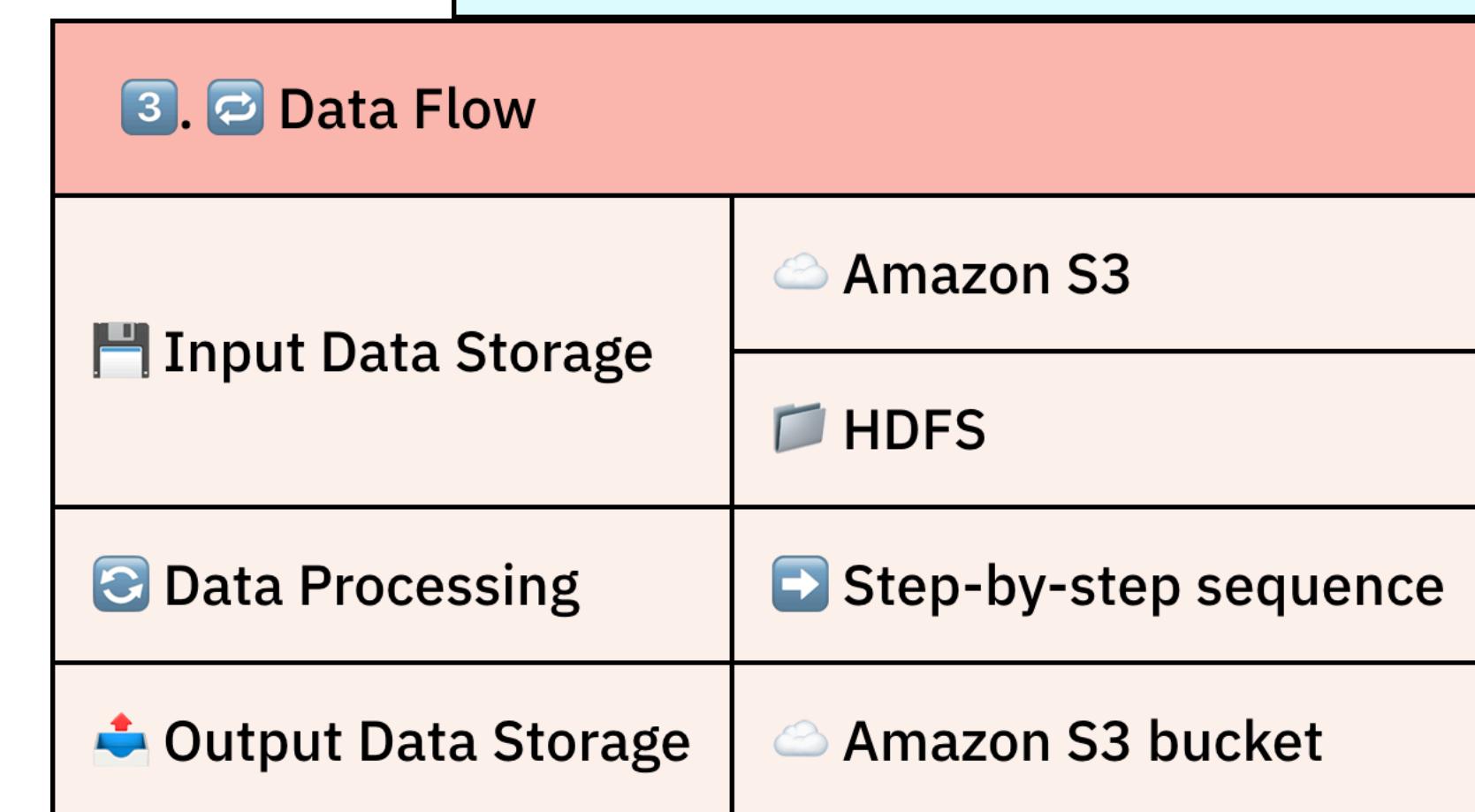
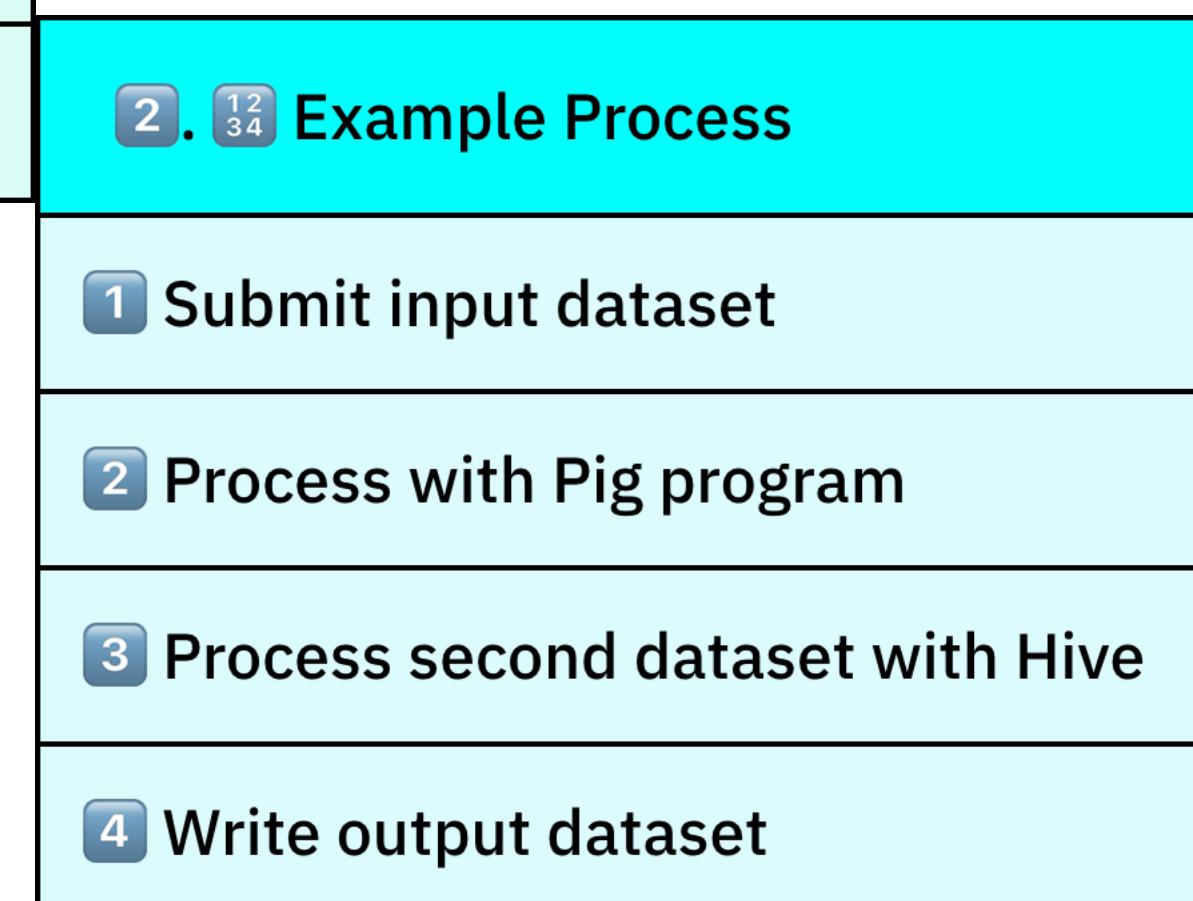
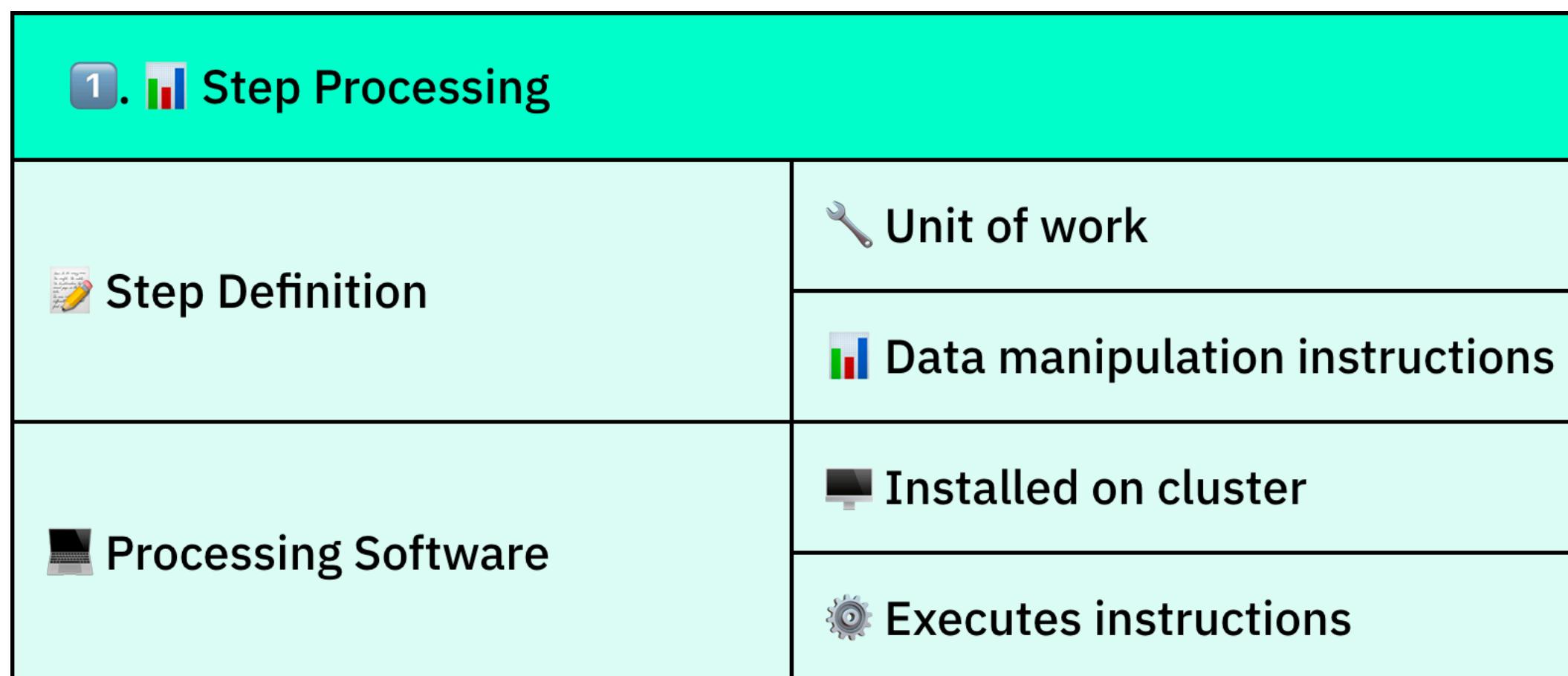
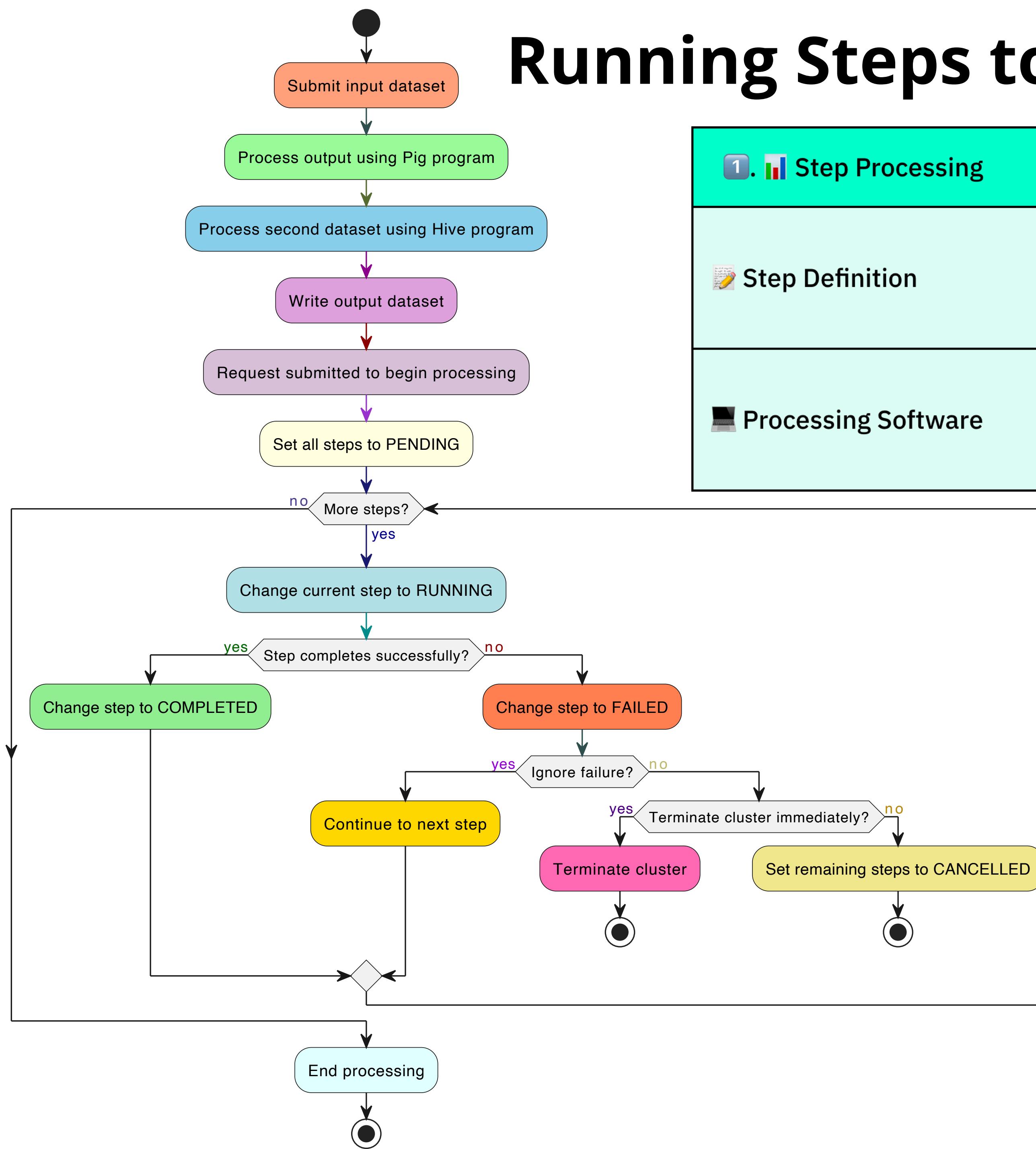
Use installed application interfaces

Submit tasks via scripts

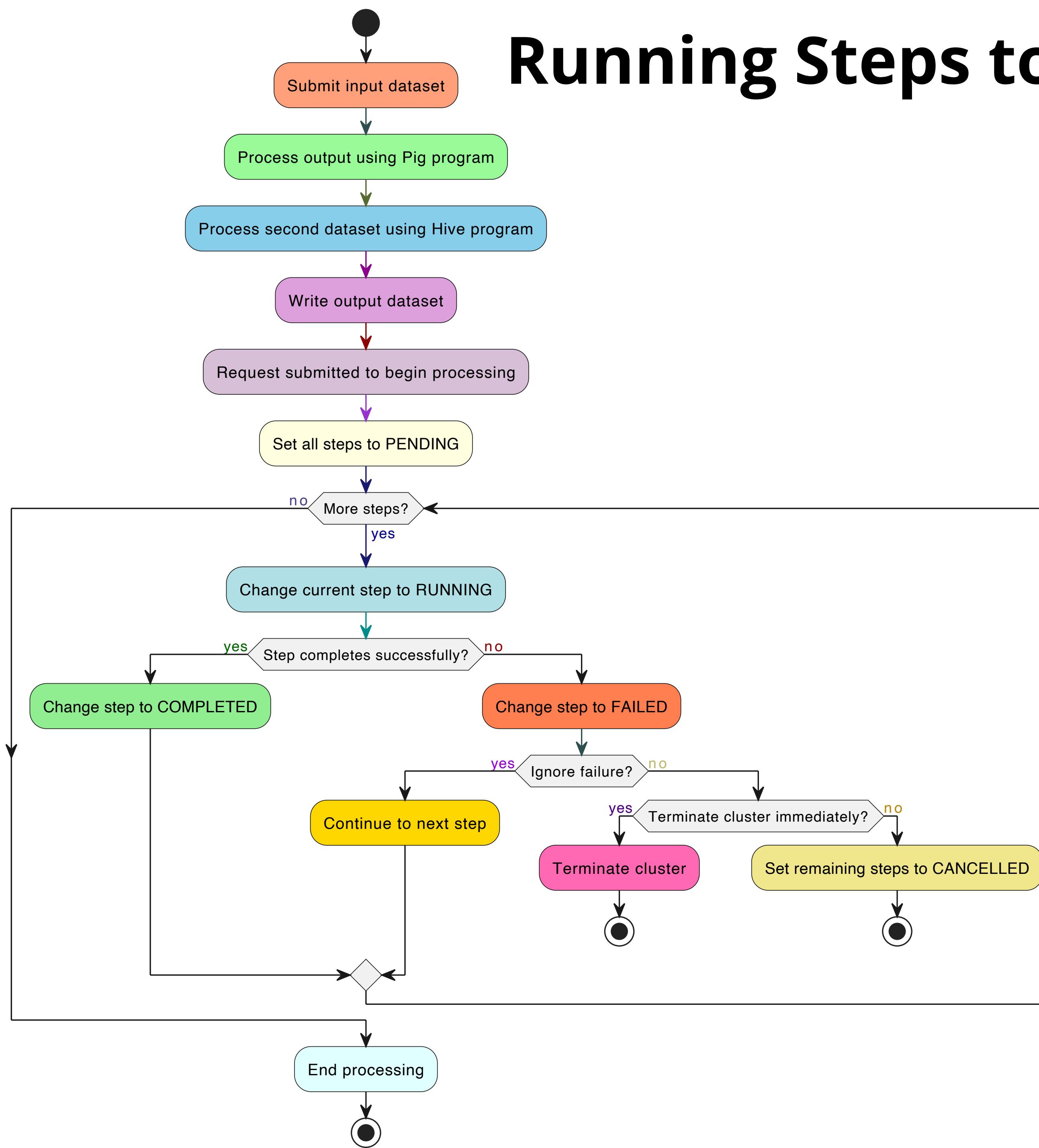
Interactive work possible

Ideal for complex tasks and debugging

Running Steps to Process Data in Amazon EMR



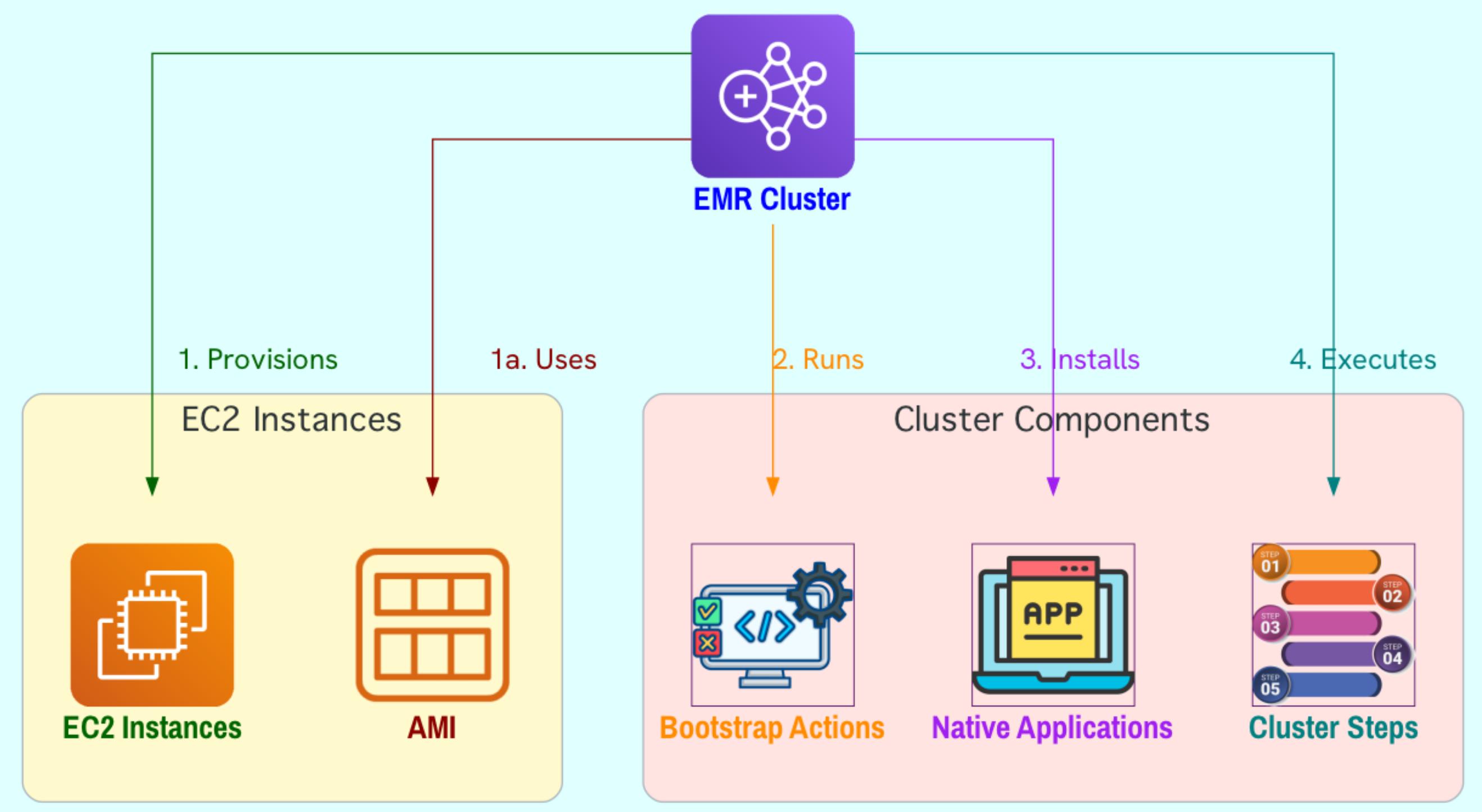
Running Steps to Process Data in Amazon EMR



4. ⏪ Step Execution Sequence	
🏁 Initiation	➡ Submit processing request
⏳ Initial State	🟡 All steps PENDING
▶ First Step Execution	🟢 First step RUNNING
	🟡 Other steps PENDING
✅ Step Completion	🔵 Completed step COMPLETED
🔄 Process Repetition	▶ Next step starts RUNNING
🏁 Processing End	✅ All steps complete

5. ✗ Step Failure Handling	
🔴 Failure State	✗ Failed step FAILED
🚫 Default Behavior	🔴 Remaining steps CANCELLED
☒ Alternative Options	➡ Ignore failure and proceed
	✗ Terminate cluster immediately

Amazon EMR Cluster Lifecycle

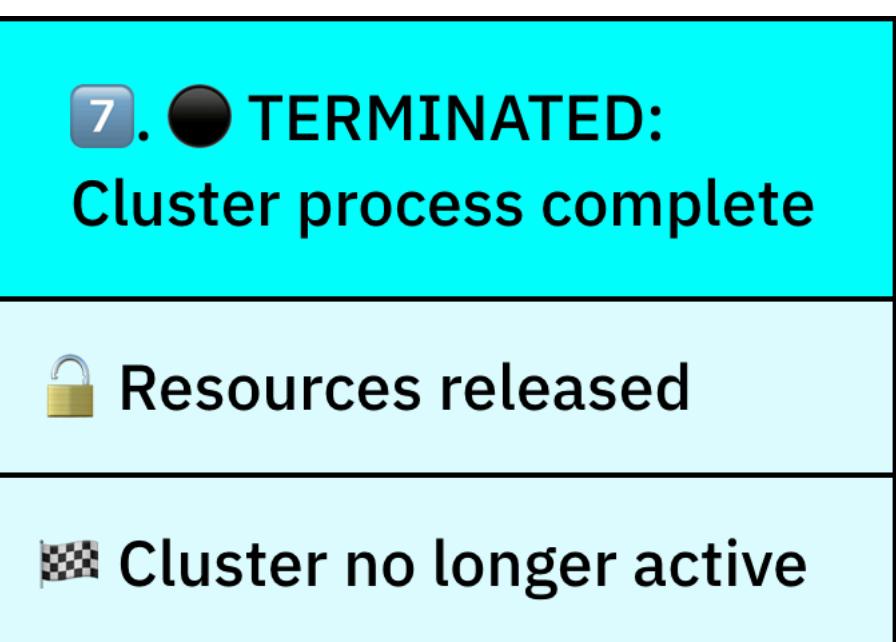
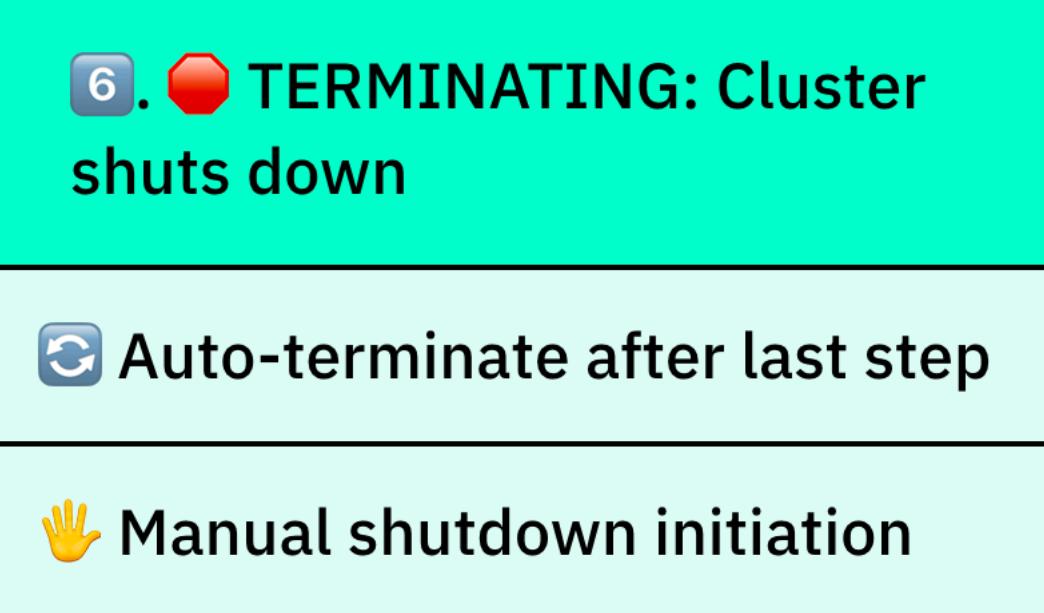
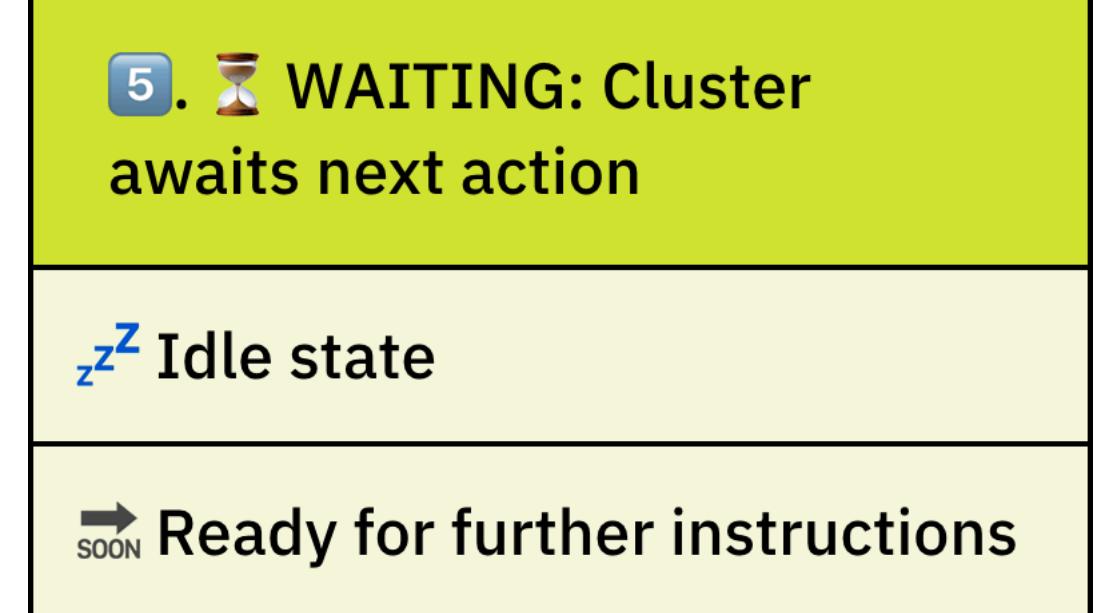
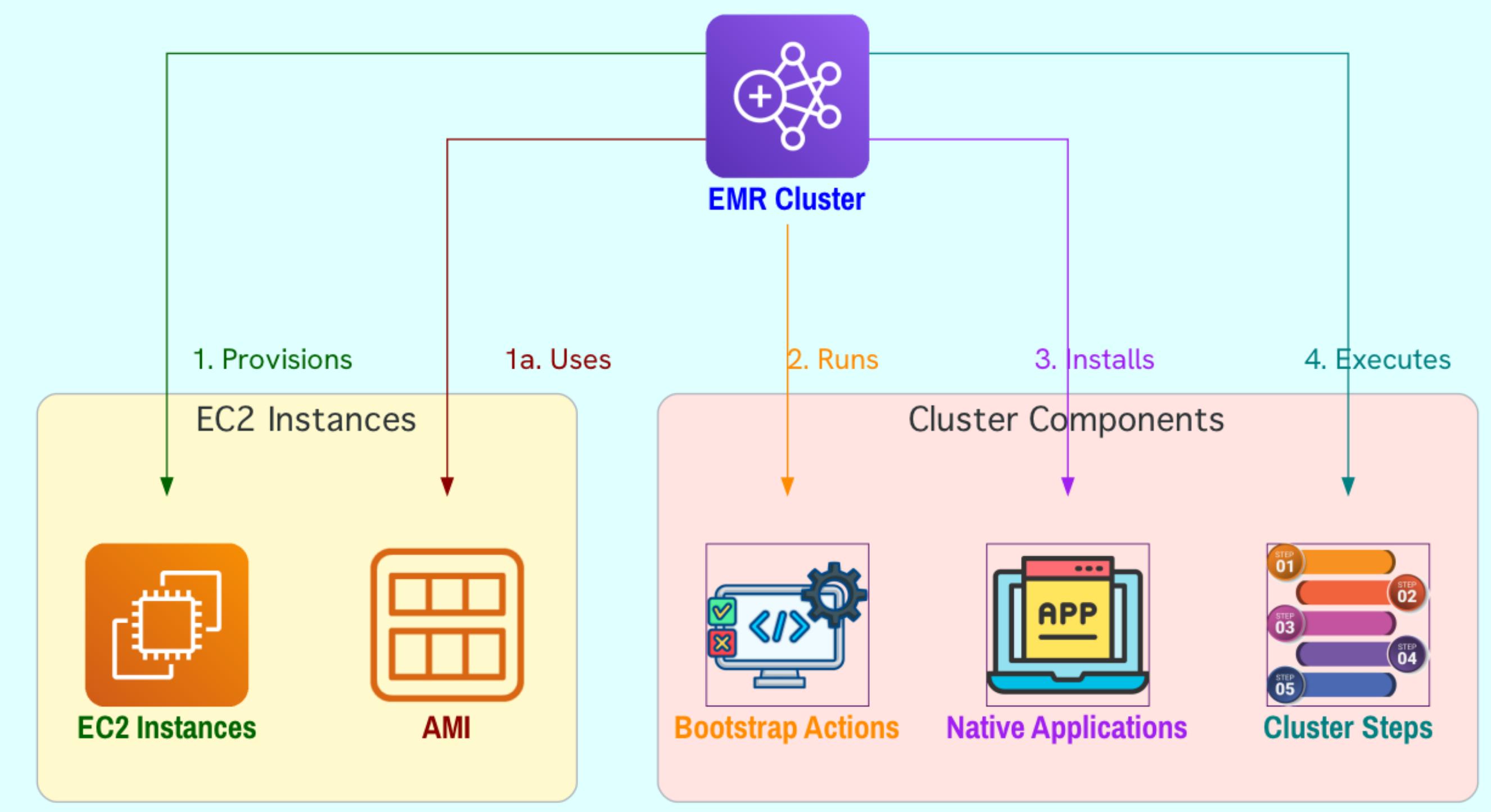


1. 🚀 STARTING: Provision EC2 instances	2. 🔧 BOOTSTRAPPING: Run bootstrap actions
Provisions EC2 instances	Executes specified actions
Installs custom applications	Performs customizations
Uses default or custom AMI	
3. ⚒ Install native applications	4. 🎓 RUNNING: Execute cluster steps
Hadoop	Connect to cluster instances
Hive	Run specified steps
Spark	Submit additional steps

Cluster States

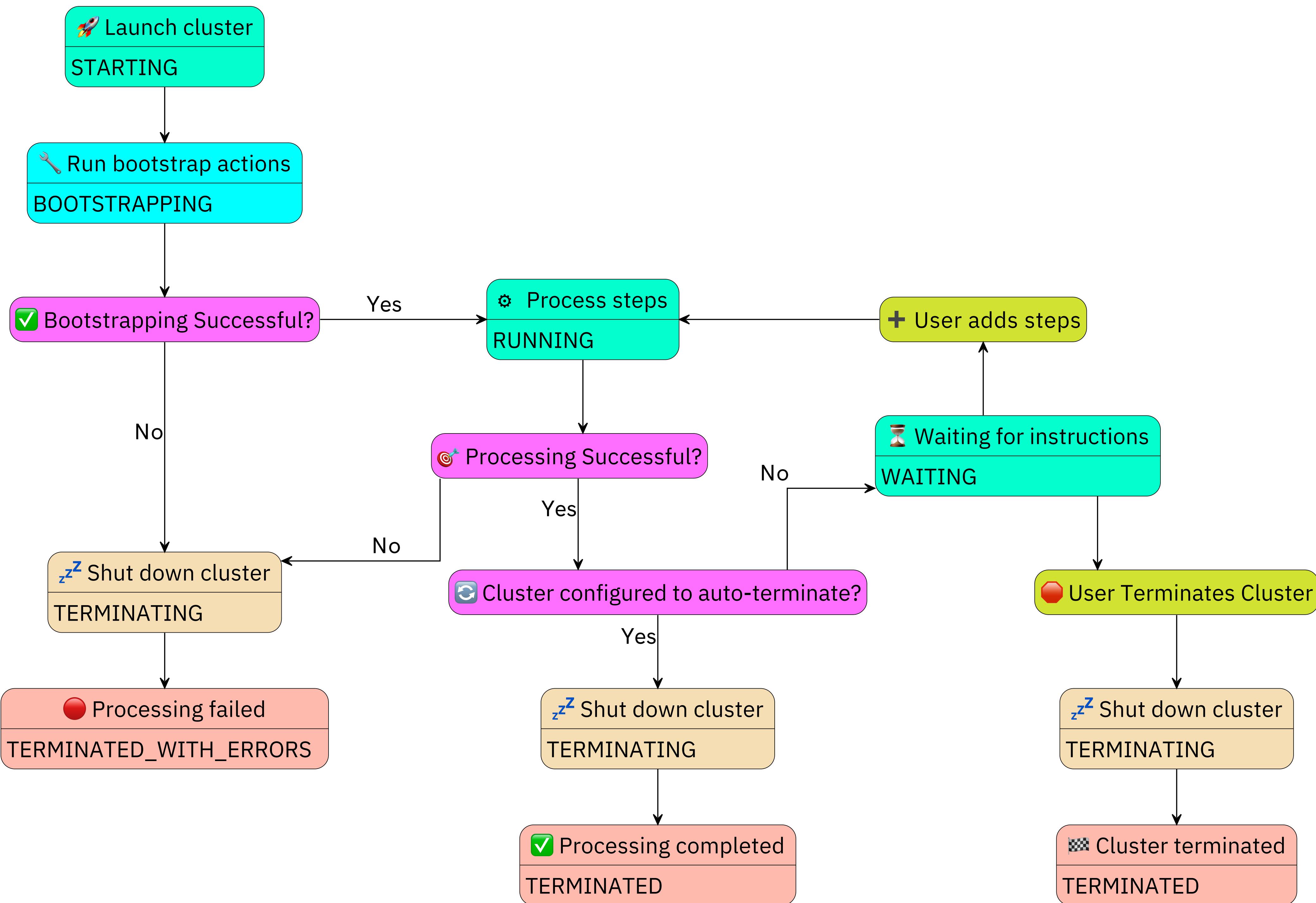


Amazon EMR Cluster Lifecycle



Cluster States

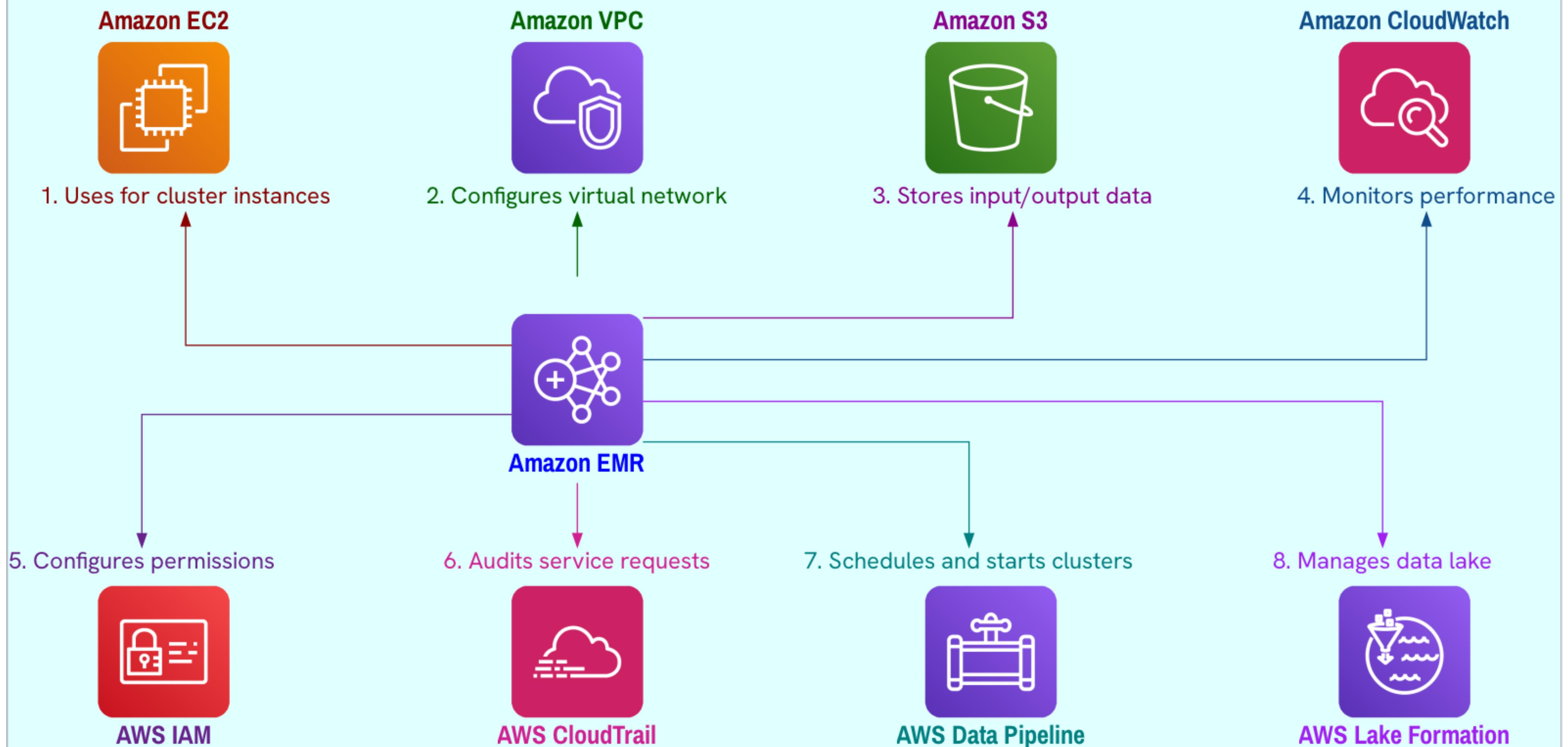




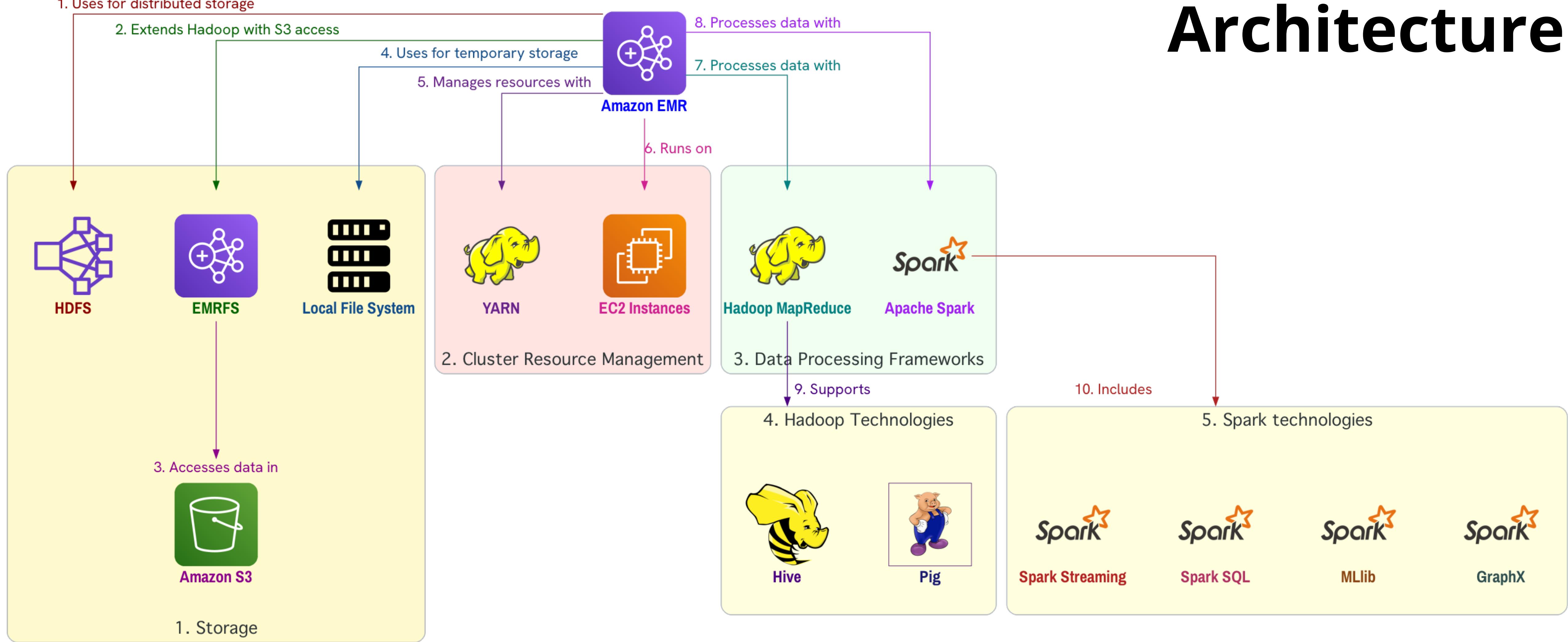
Amazon EMR Cluster Failure Handling

Amazon EMR Cluster Failure Handling	
1. ⚡ Cluster Failure Consequences	⬅ Cluster termination
	💻 Instance termination
2. 🛡 Termination Protection	🚫 Prevents automatic termination
	💾 Allows data retrieval
3. 📊 Cluster State	X TERMINATED_WITH_ERRORS
4. 🔄 Recovery Process	⬇ Retrieve data
	🔓 Remove termination protection
	🔴 Manually terminate cluster

AWS Integration

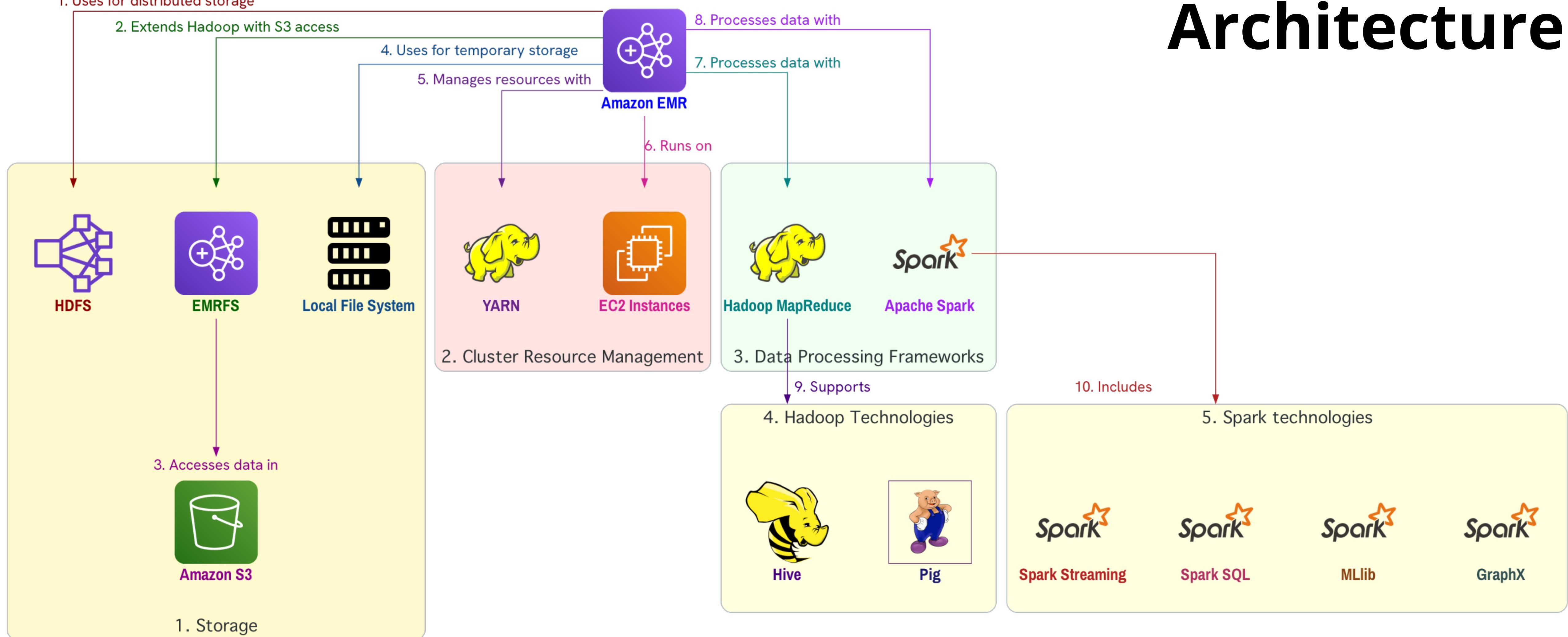


Architecture



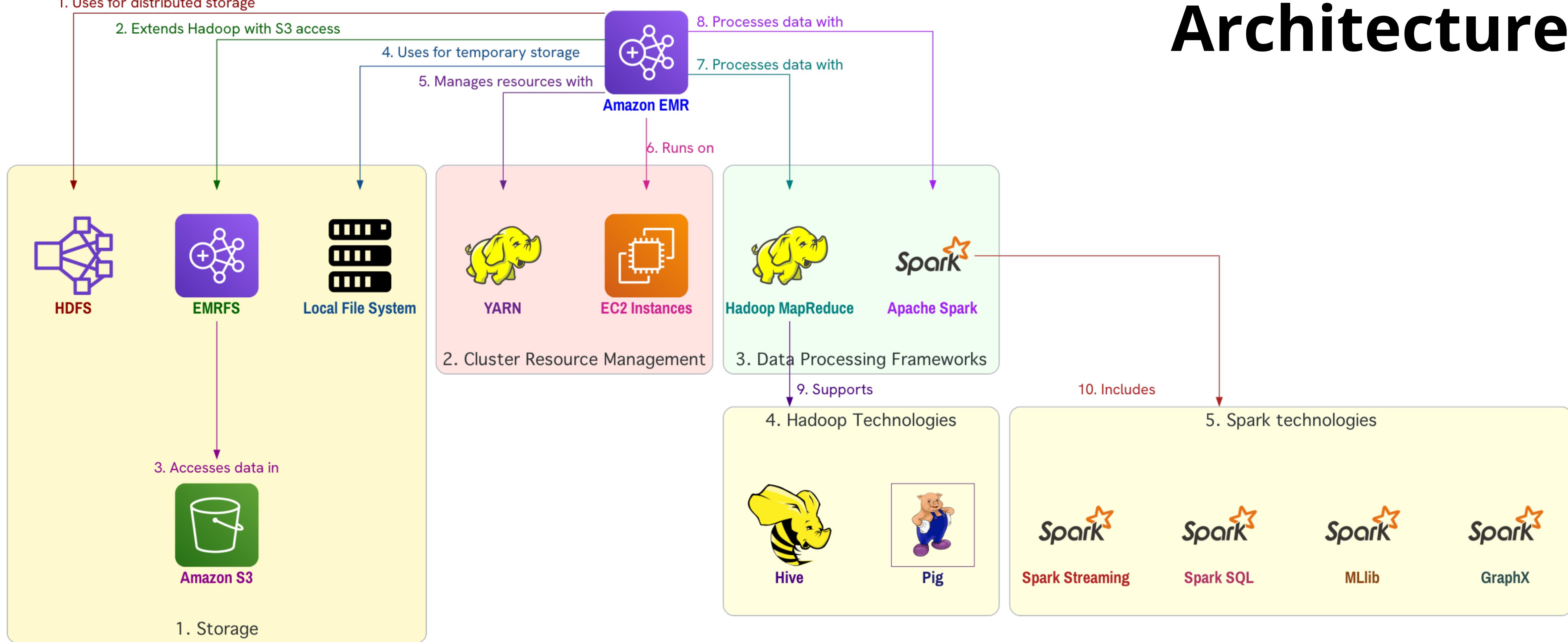
1. **Storage Layer:** HDFS: Distributed, scalable, Fault tolerance, Intermediate results; EMRFS: Direct S3 access, Input/output data; Local File System: EC2 instance store, Temporary storage

Architecture



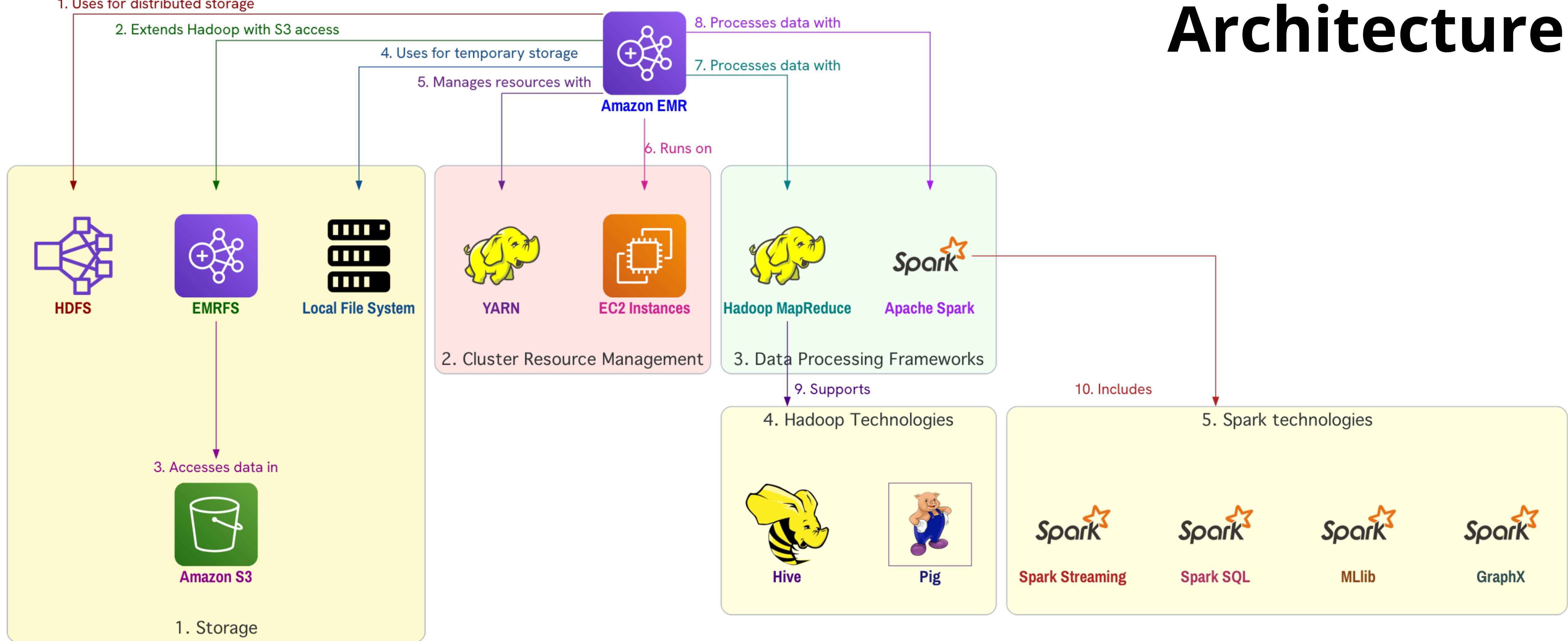
2. 🧠 Resource Management: 🧶 YARN: ⏳ Central resource management, ⚡ Multiple frameworks; 🤖 EMR Agent: 🔧 Administers YARN, 🏢 Maintains cluster health, 🛫 Communicates with EMR; 🎯 Spot Instance optimization: 🏠 Core nodes for app masters, ⏱ Job continuity

Architecture



3. Processing Frameworks: MapReduce: Open-source distributed computing, Simplifies parallel processing; Spark: Directed acyclic graphs, In-memory caching, Faster big data processing; Support: High-level interfaces, Multiple languages

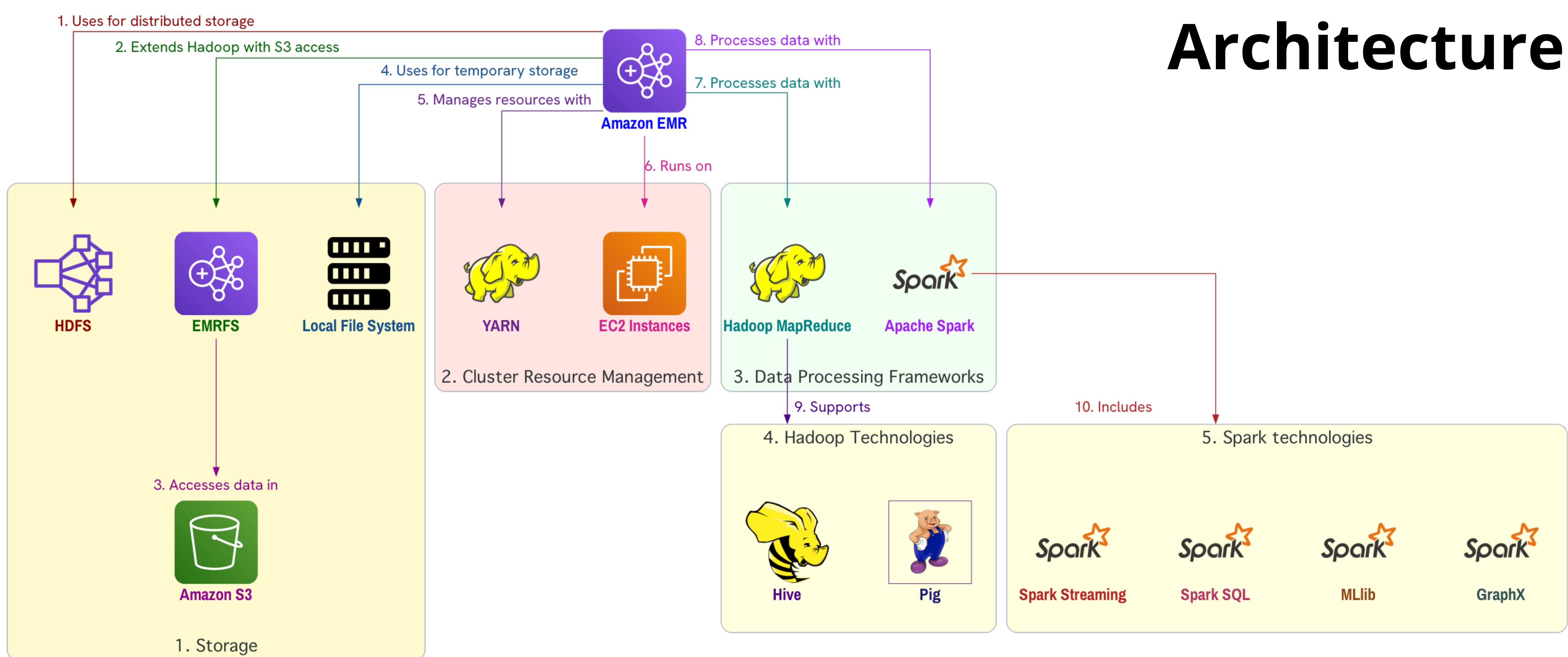
Architecture



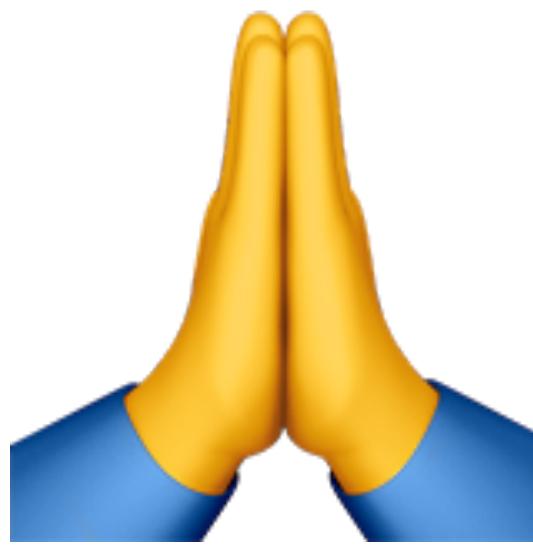
4. Applications & Programs:

- 🐝 Hive: 🏫 Data warehouse software, 🔎 SQL-like queries, 📊 Data summarization;
- 🐷 Pig: 🚧 High-level platform, 🐘 Creates MapReduce programs, 📋 Pig Latin language;
- ⭐ Spark Ecosystem: 🔎 Spark SQL, 💧 Streaming, 🧠 MLlib (machine learning), 🕸️ GraphX (graph processing);
- 📚 Multiple libraries and languages: ⚙️ Flexible data interaction

Architecture



5. 🔗 Integration: ☁ S3: 💰 Cost-effective storage, 🛡 Durable storage; 🖥 EC2: 📈 Scalable compute resources;💡 Spot Instances: 💰 Cost reduction, ⚡ Processing reliability



**Thanks
for
Watching**