



Amazon EC2 Auto Scaling

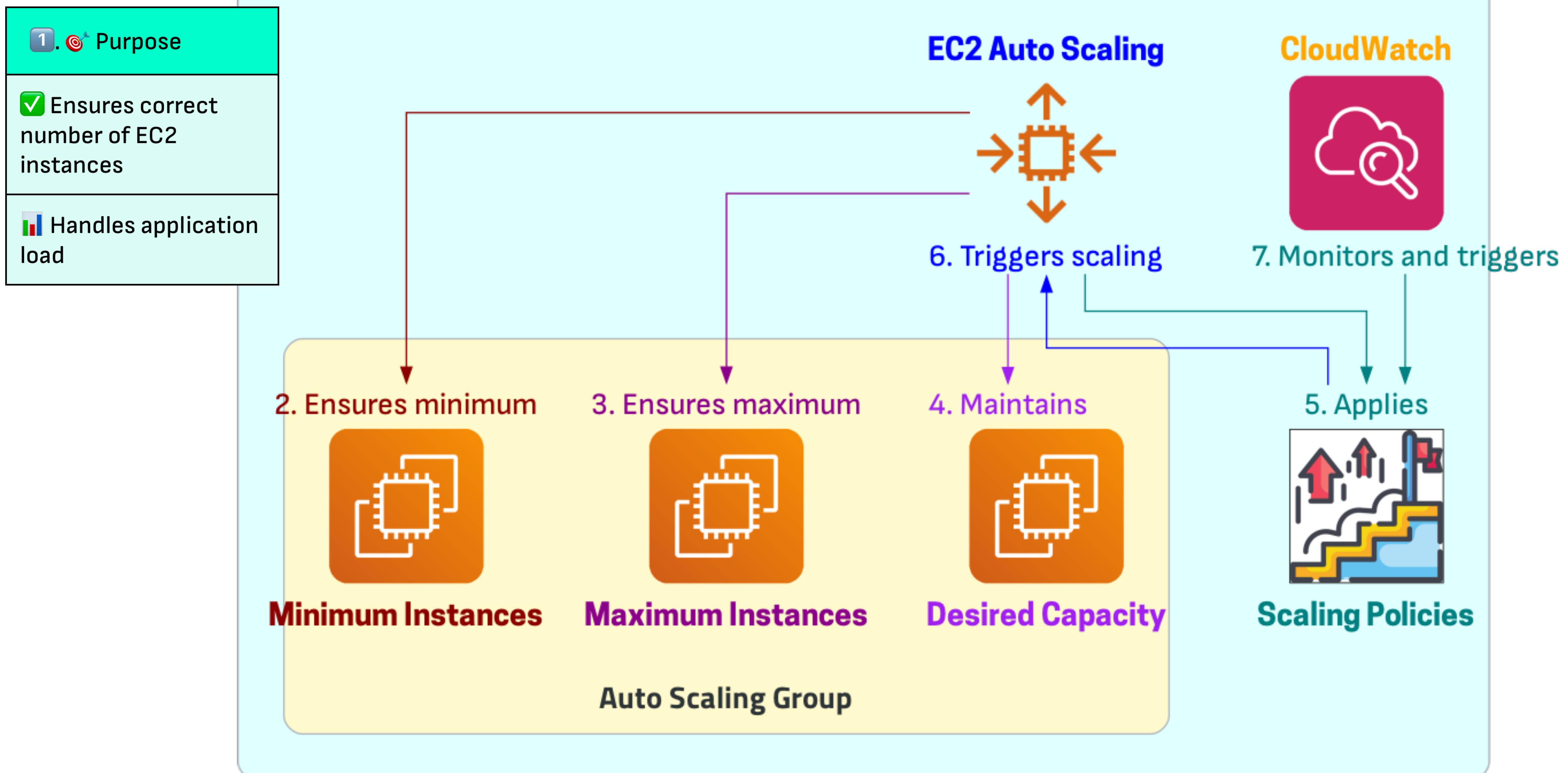
Table of Contents



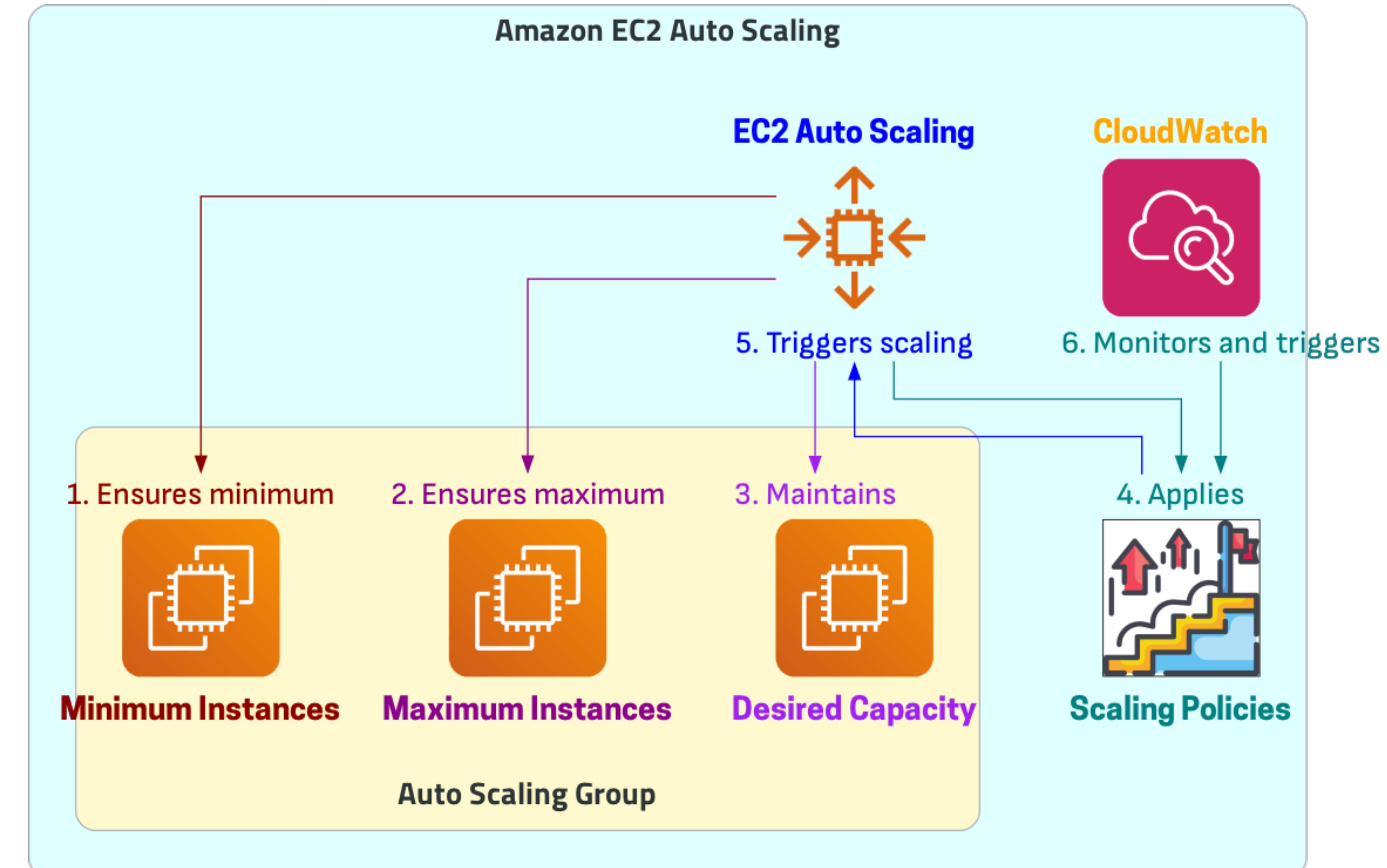
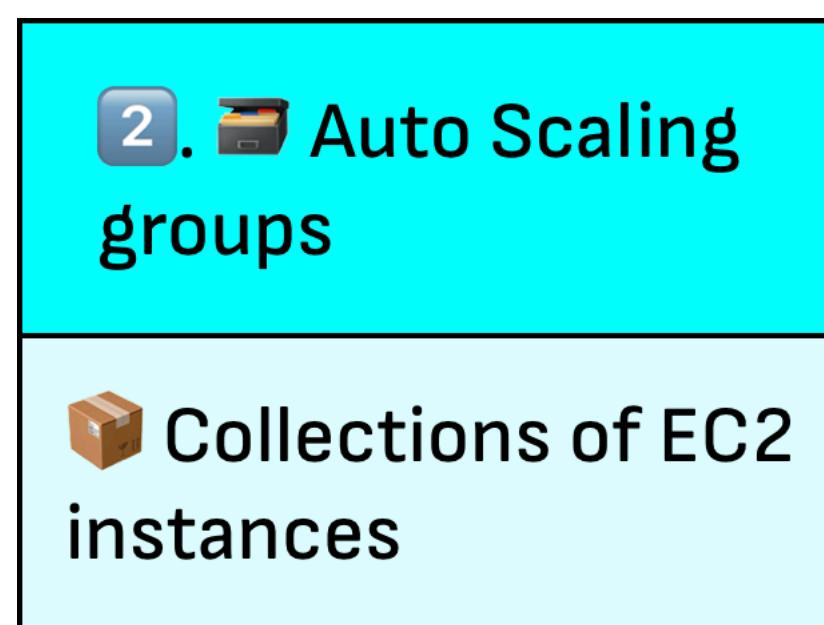
- 1. What is Amazon EC2 Auto Scaling?
- 2. Features
- 3. Auto Scaling Group Example
- 4. Amazon EC2 Auto Scaling Benefits
- 5. Cover Variable Demand
- 6. Typical Web app architecture
- 7. Distribute instances across Availability Zones
- 8. Auto Scaling launch templates
- 9. Auto Scaling Methods
- 10. Set Scaling Limits for Your Auto Scaling Group
- 11. Elastic Load Balancing with Auto Scaling
- 12. Use EventBridge for Auto Scaling Events

What is Amazon EC2 Auto Scaling?

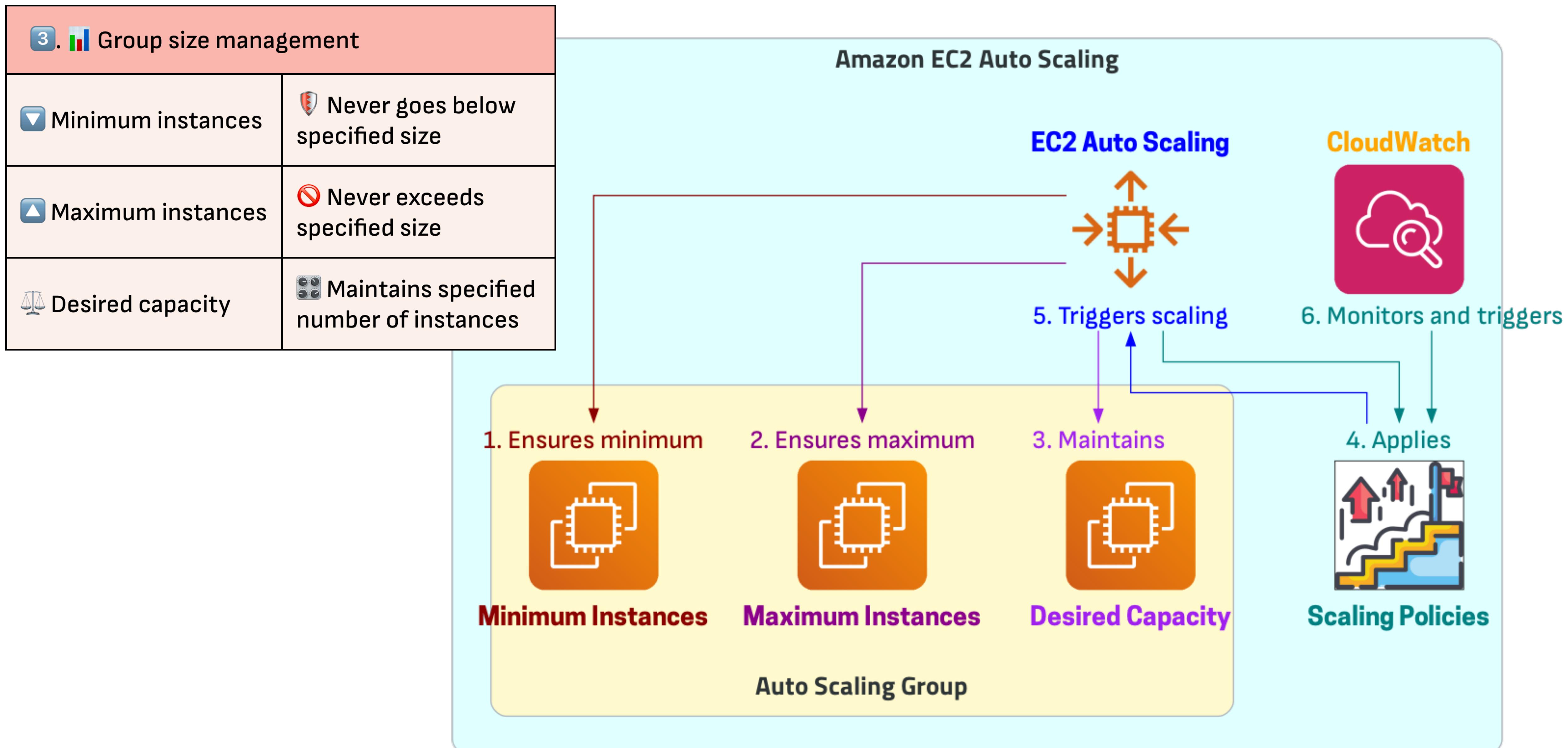
Amazon EC2 Auto Scaling



What is Amazon EC2 Auto Scaling?



What is Amazon EC2 Auto Scaling?



What is Amazon EC2 Auto Scaling?

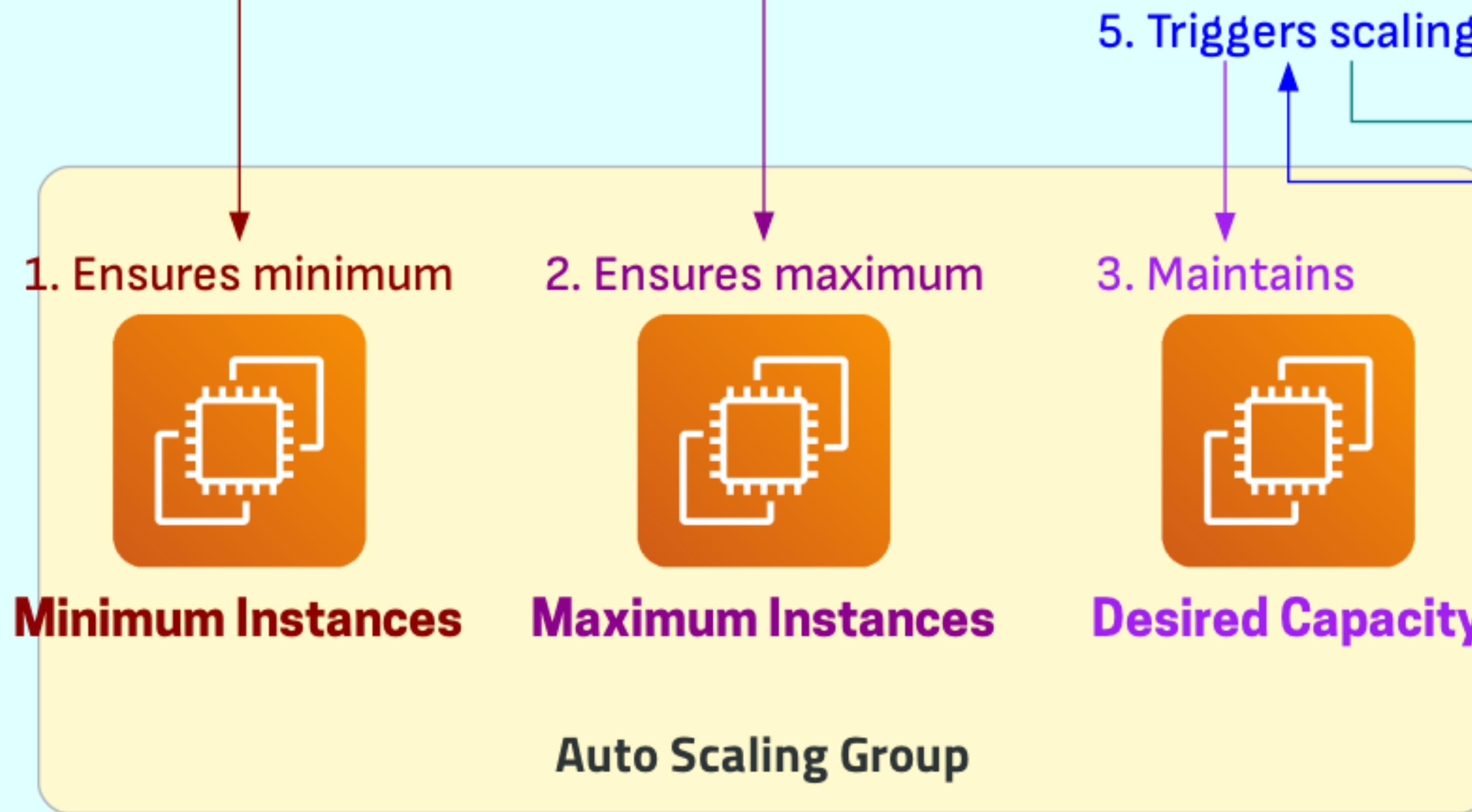
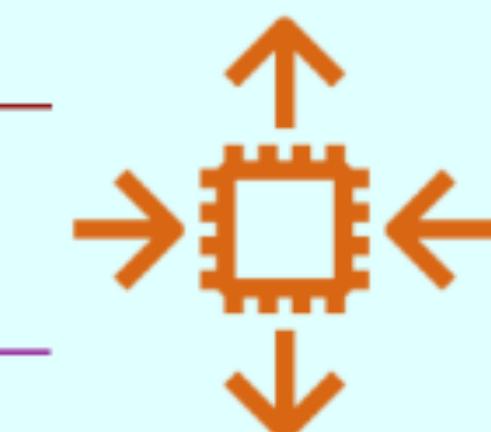
4. Scaling policies

Launches instances when demand increases

Terminates instances when demand decreases

Amazon EC2 Auto Scaling

EC2 Auto Scaling



CloudWatch

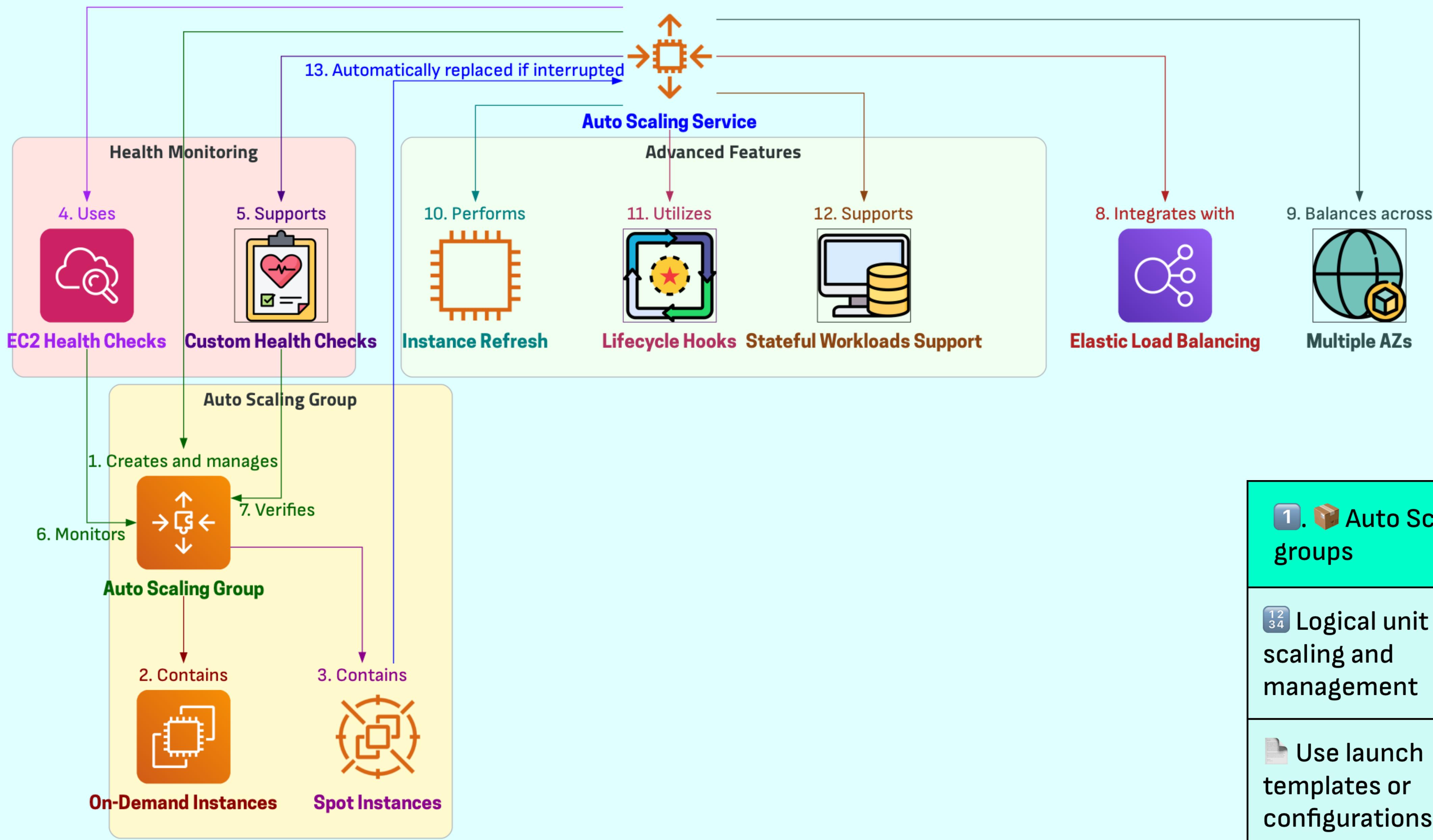


Monitors and triggers



Scaling Policies

Features

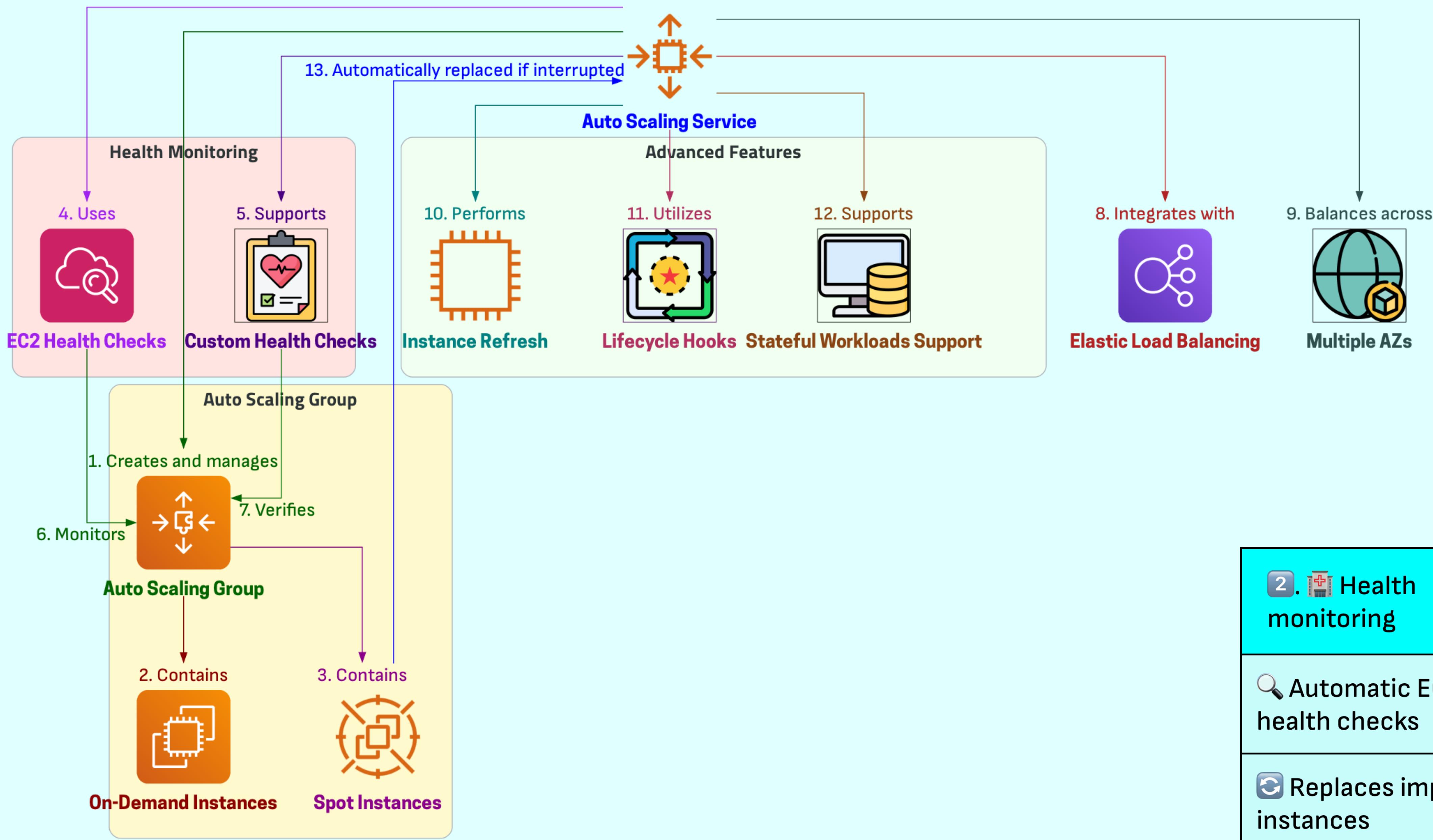


1. Auto Scaling groups

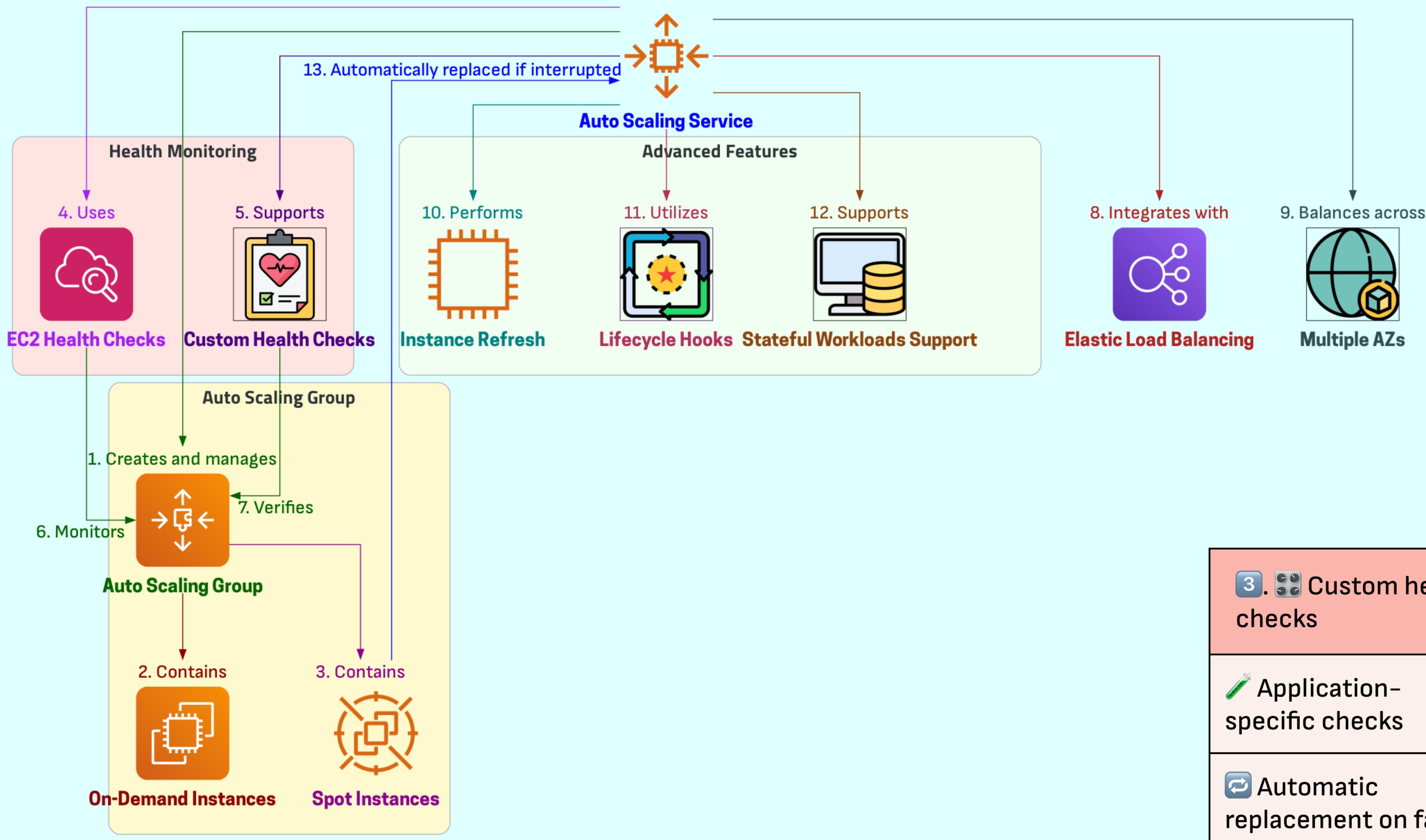
Logical unit for scaling and management

Use launch templates or configurations

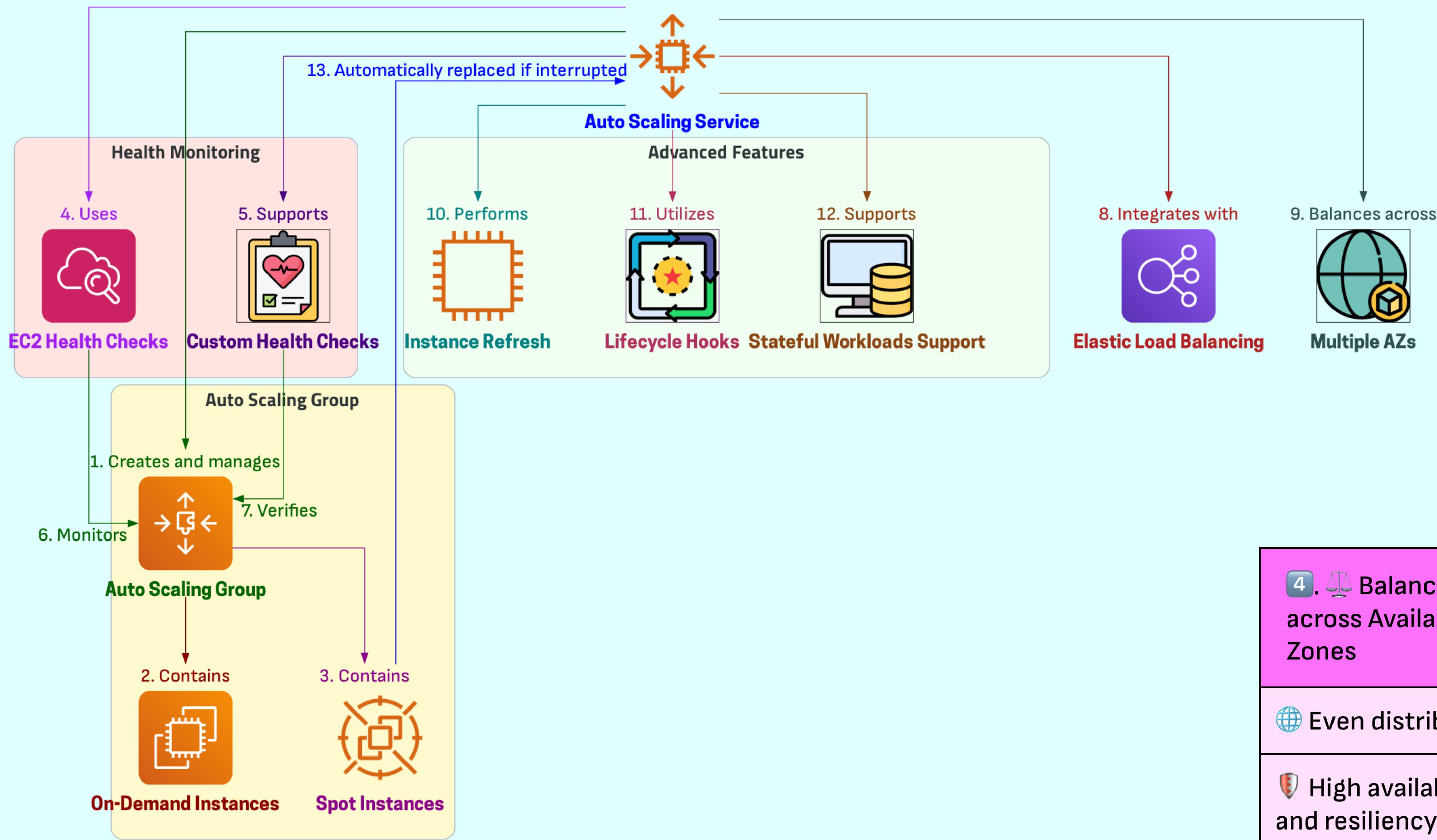
Features



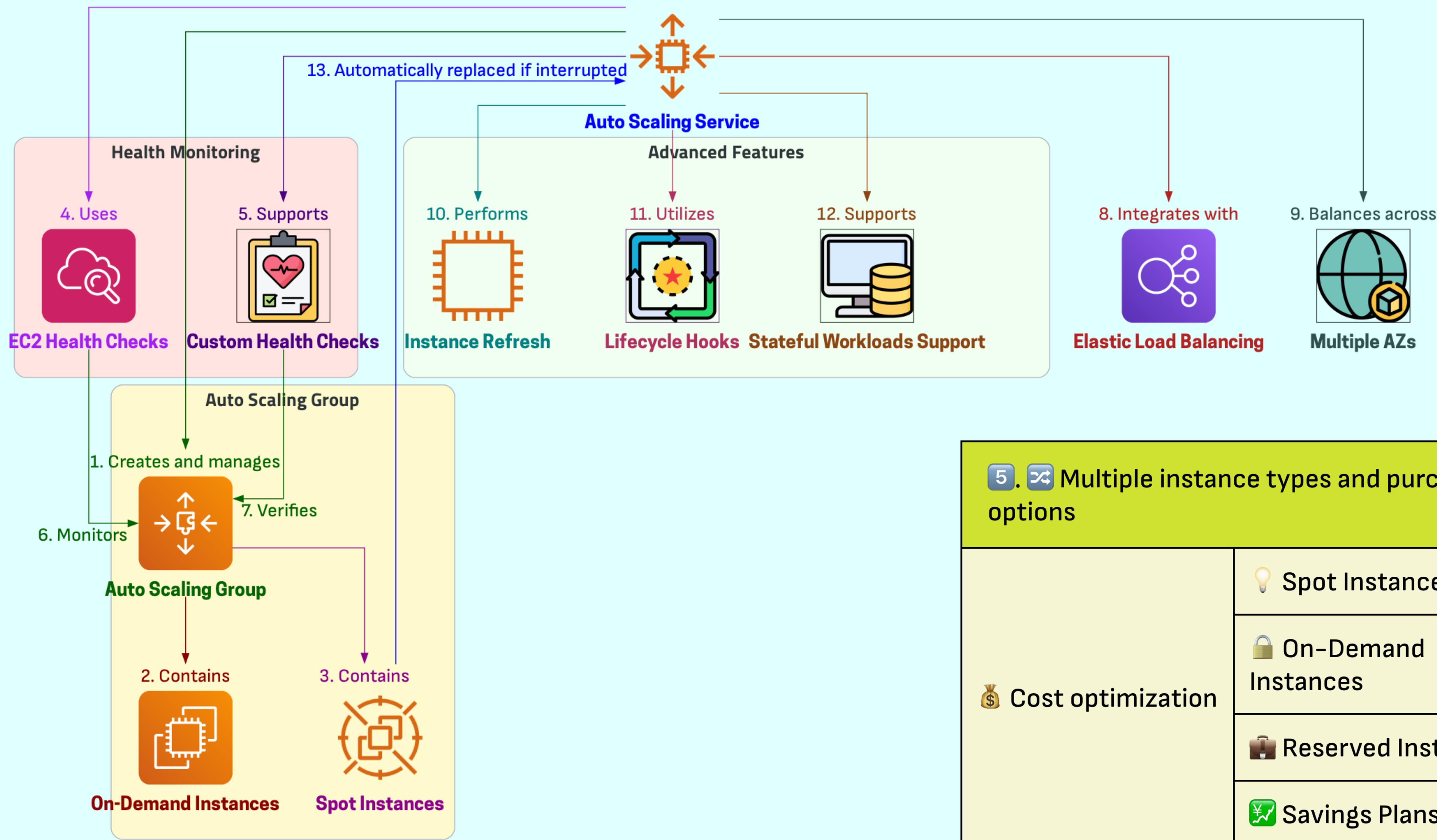
Features



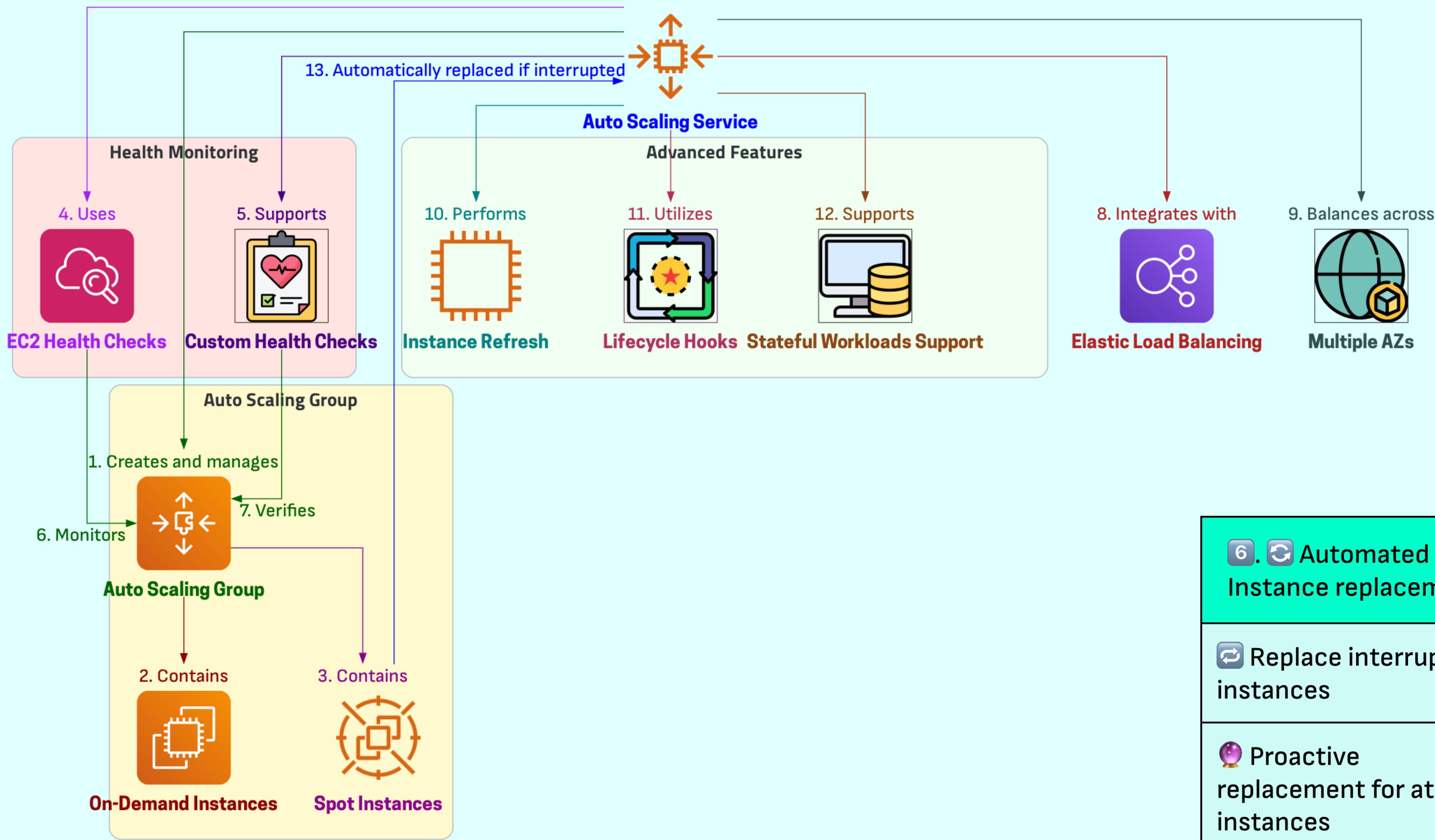
Features



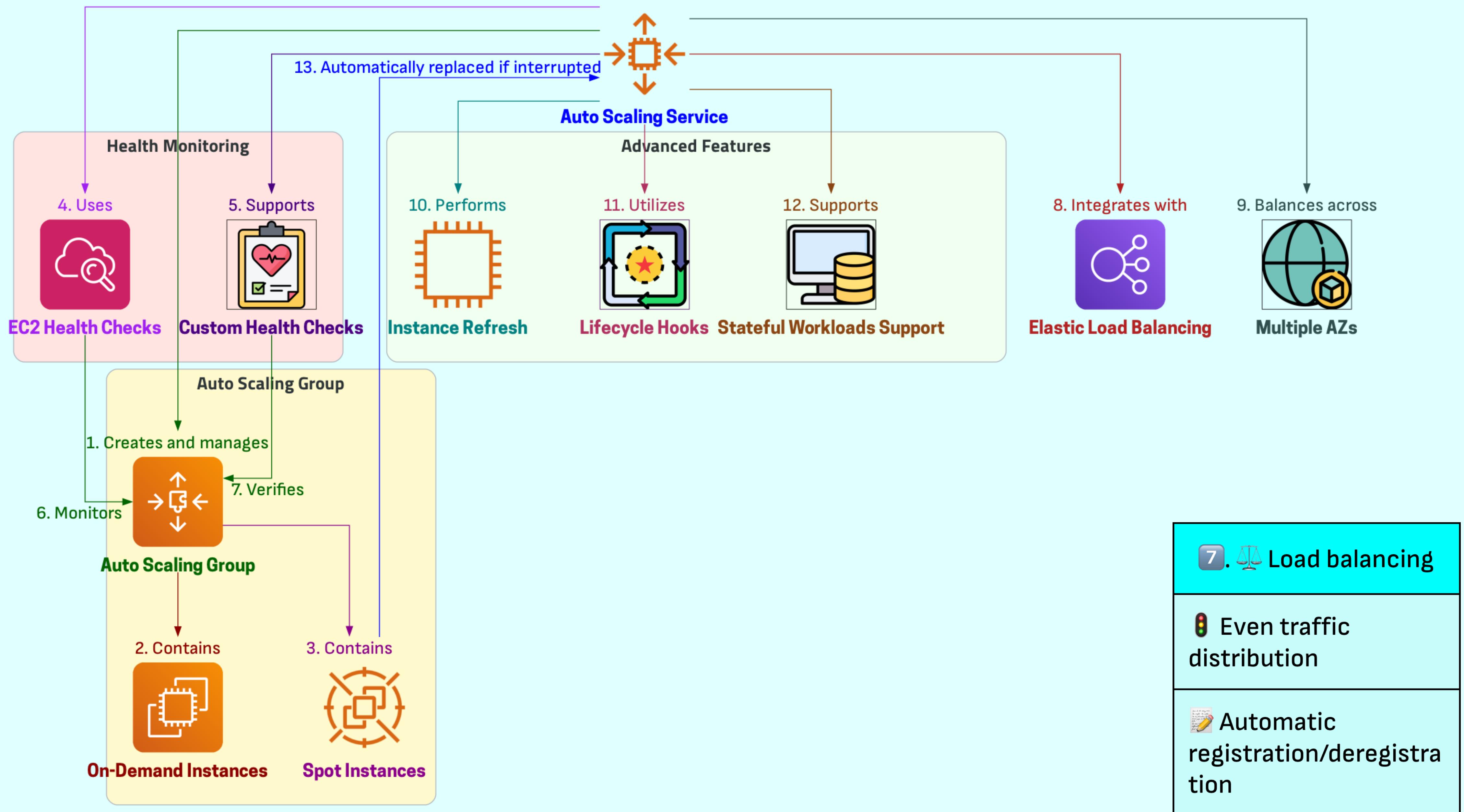
Features



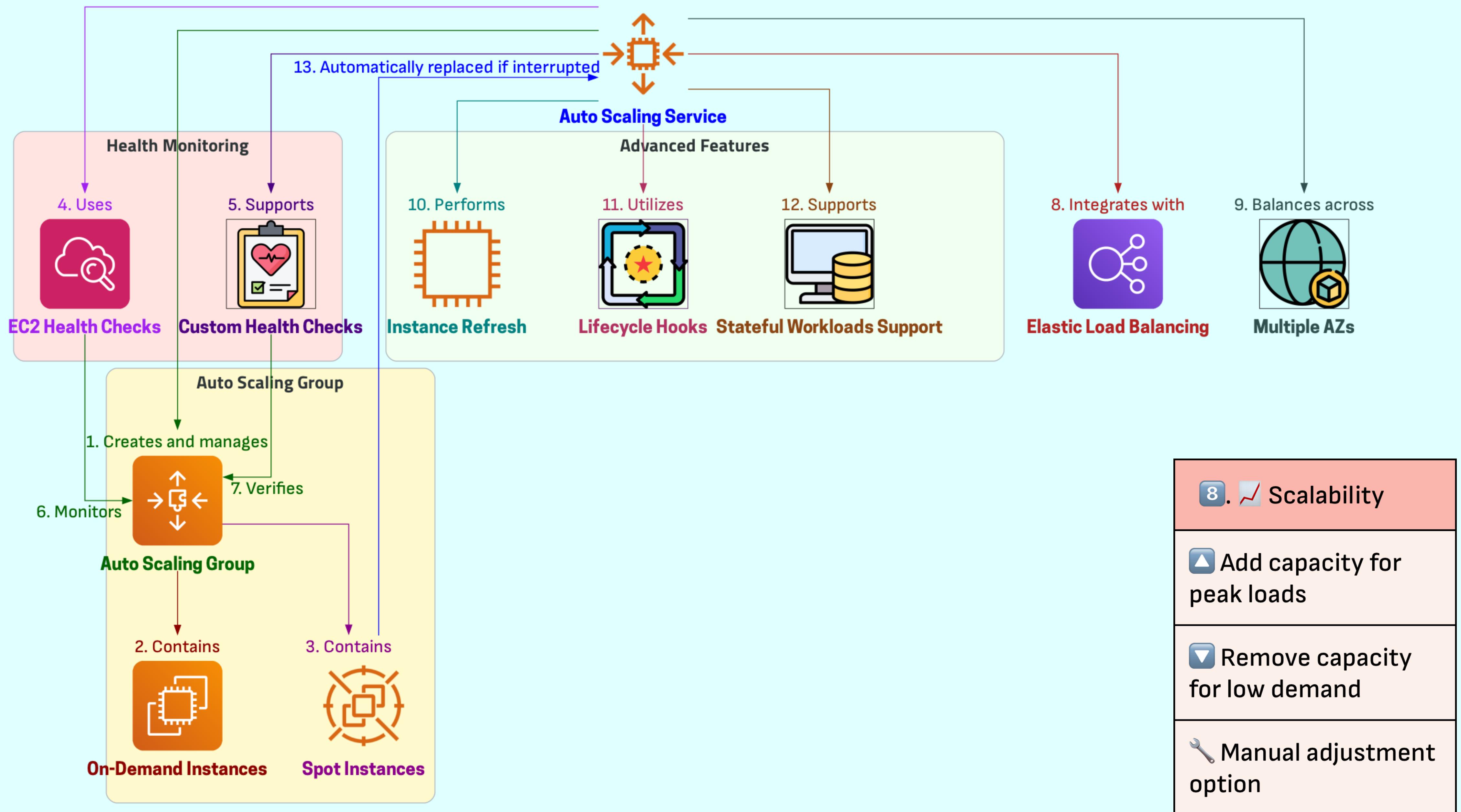
Features



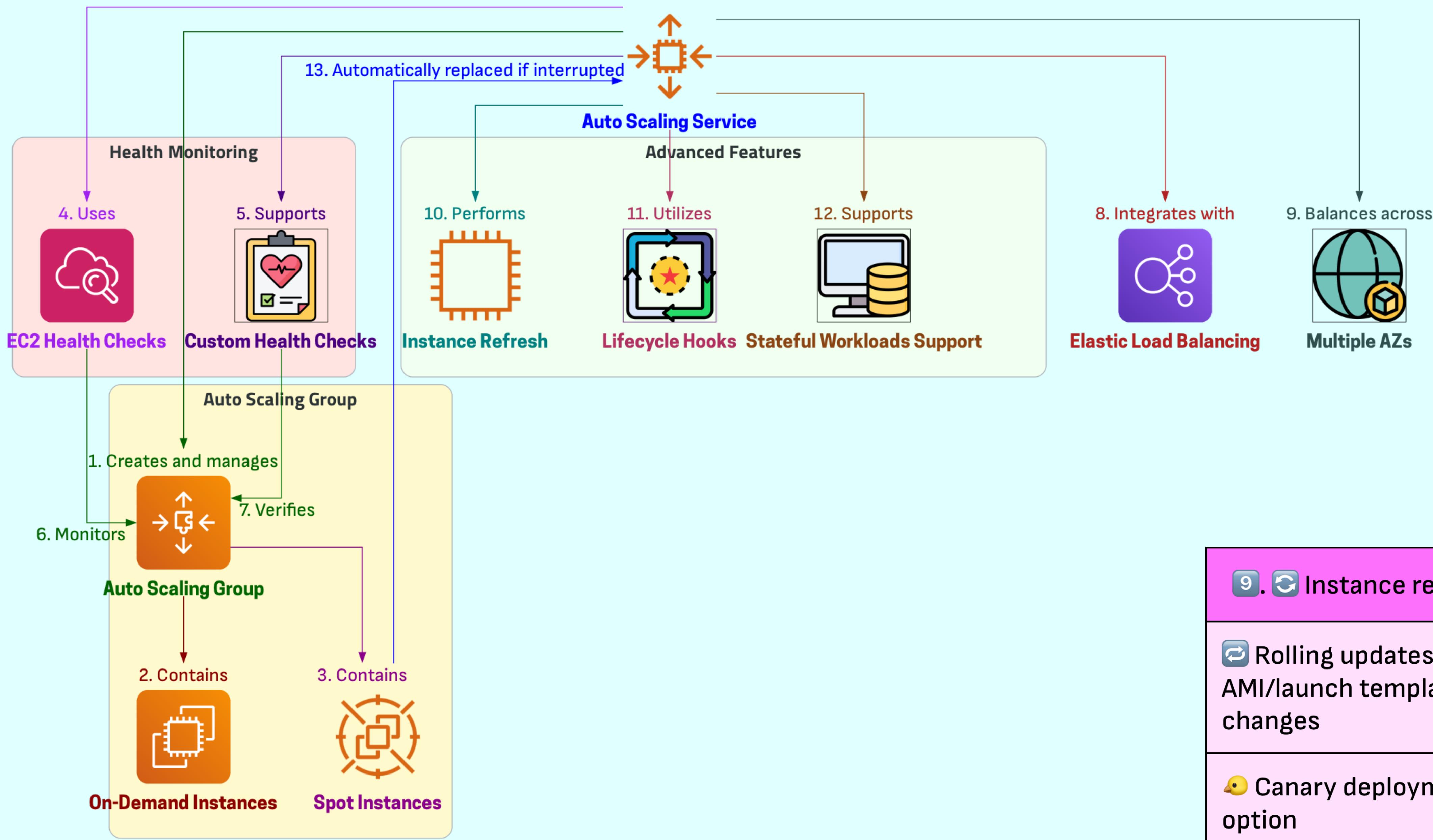
Features



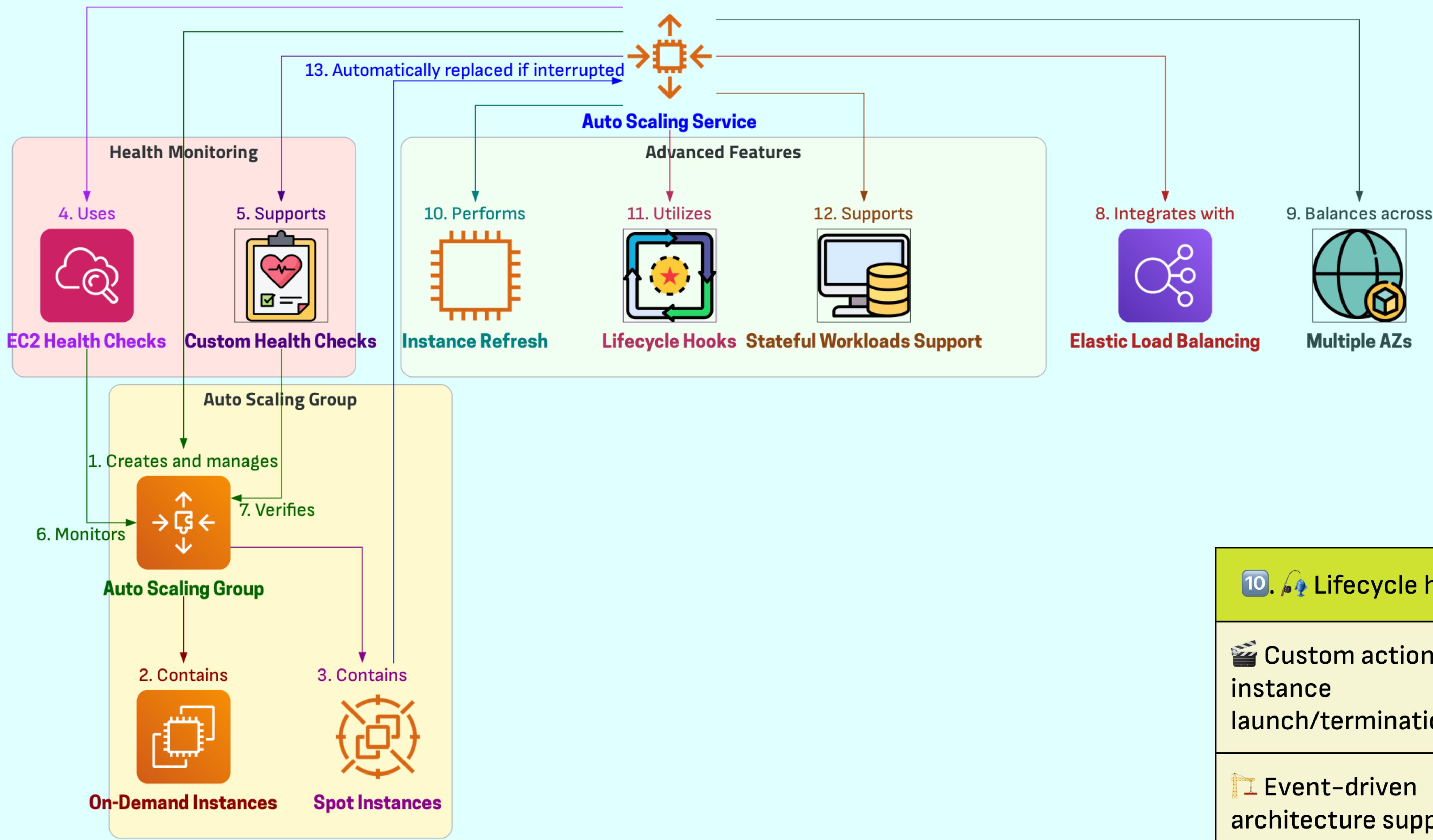
Features



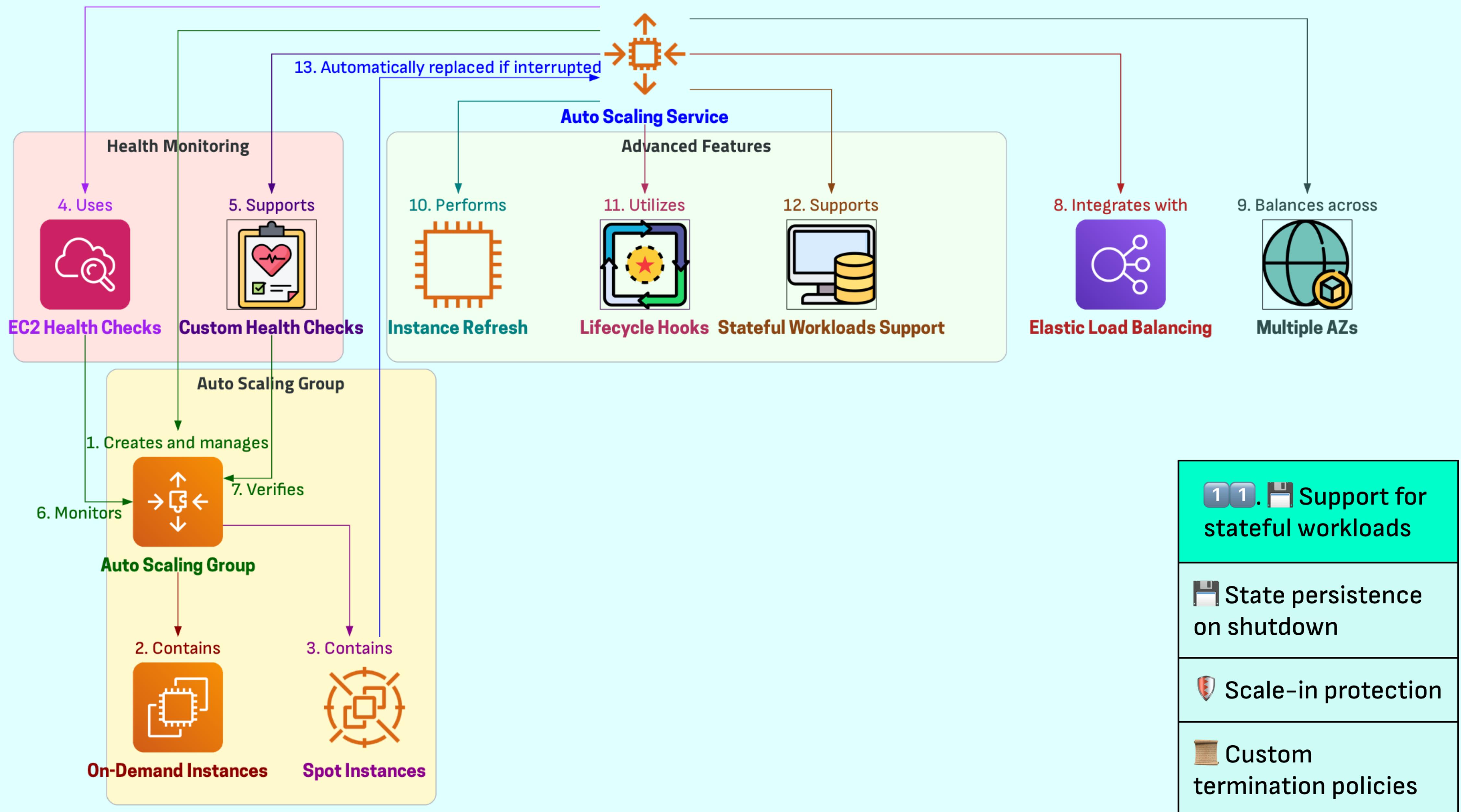
Features

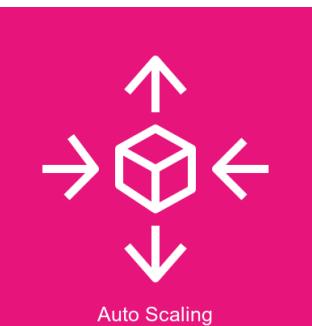


Features

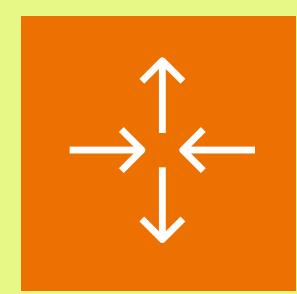


Features

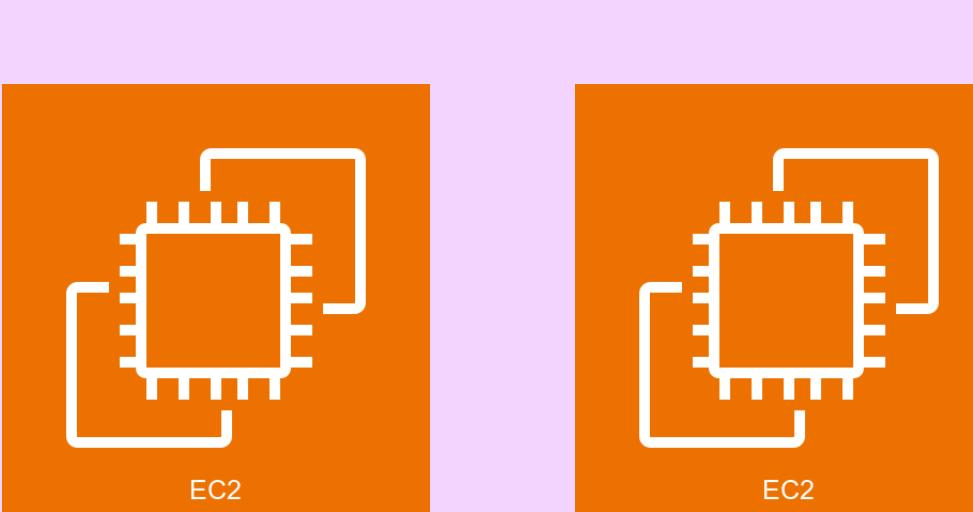




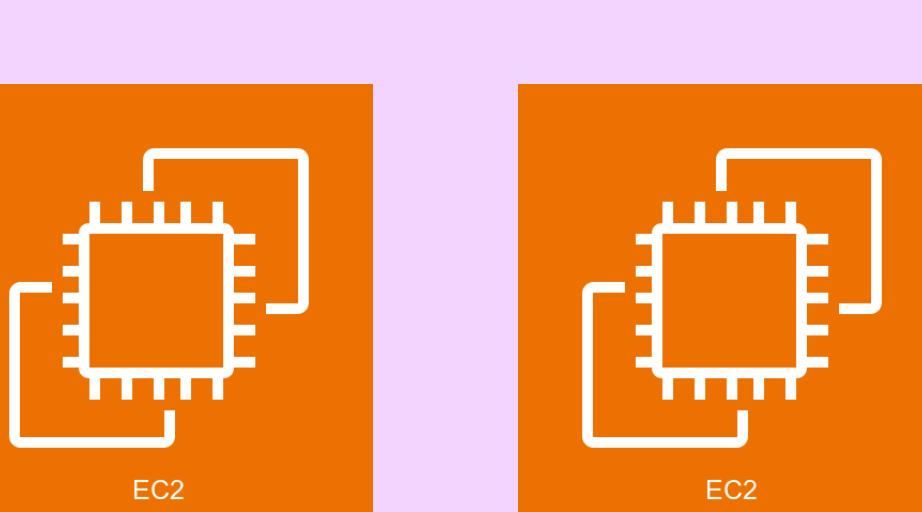
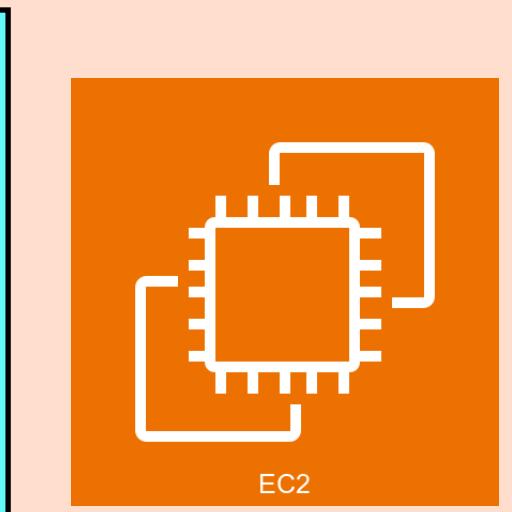
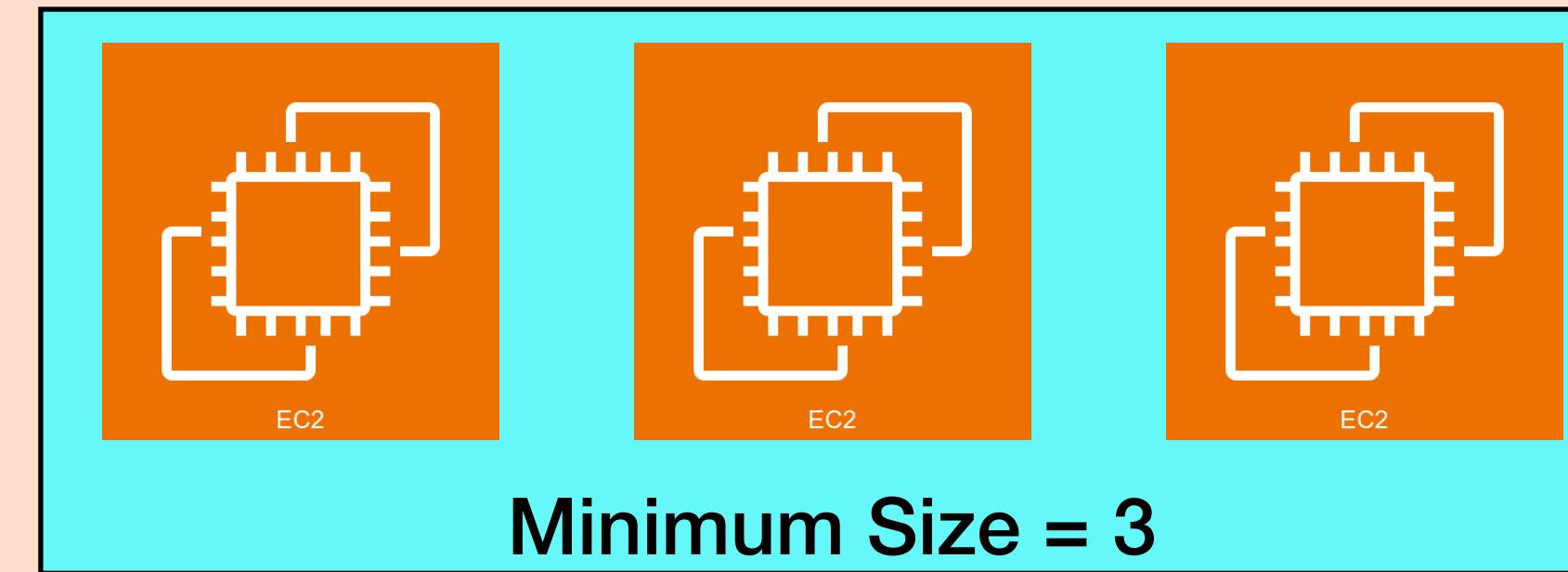
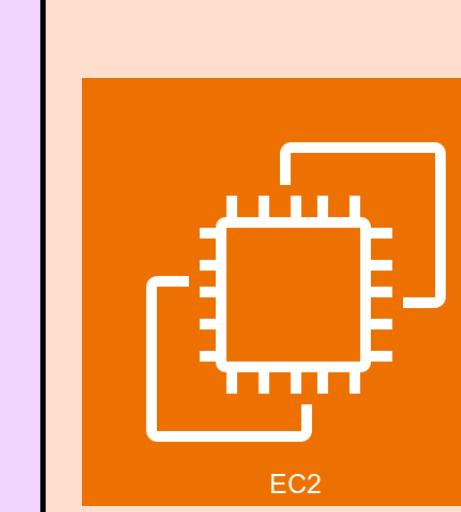
Auto Scaling Group Example



Auto Scaling Group

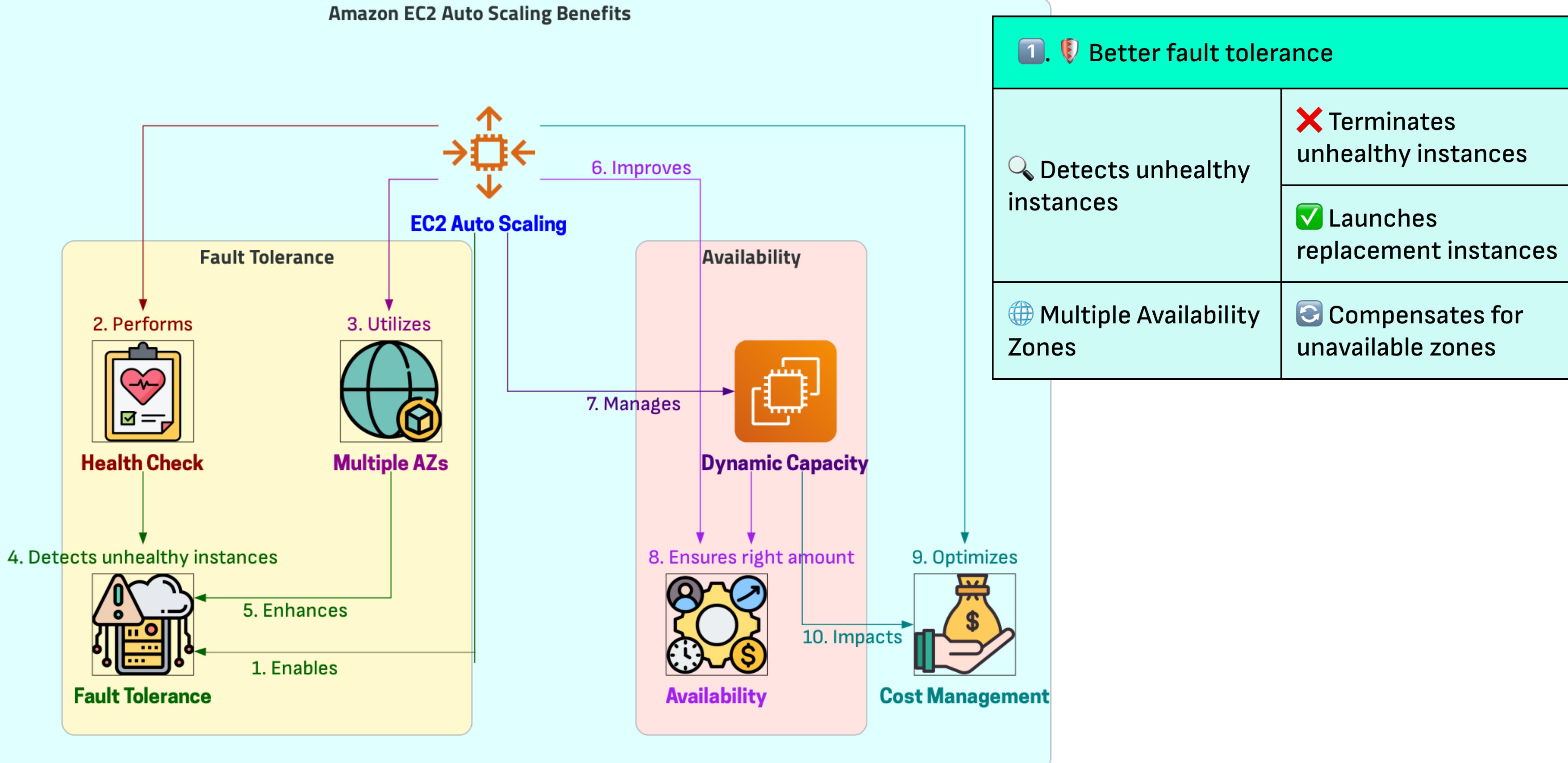


Scale Out/In

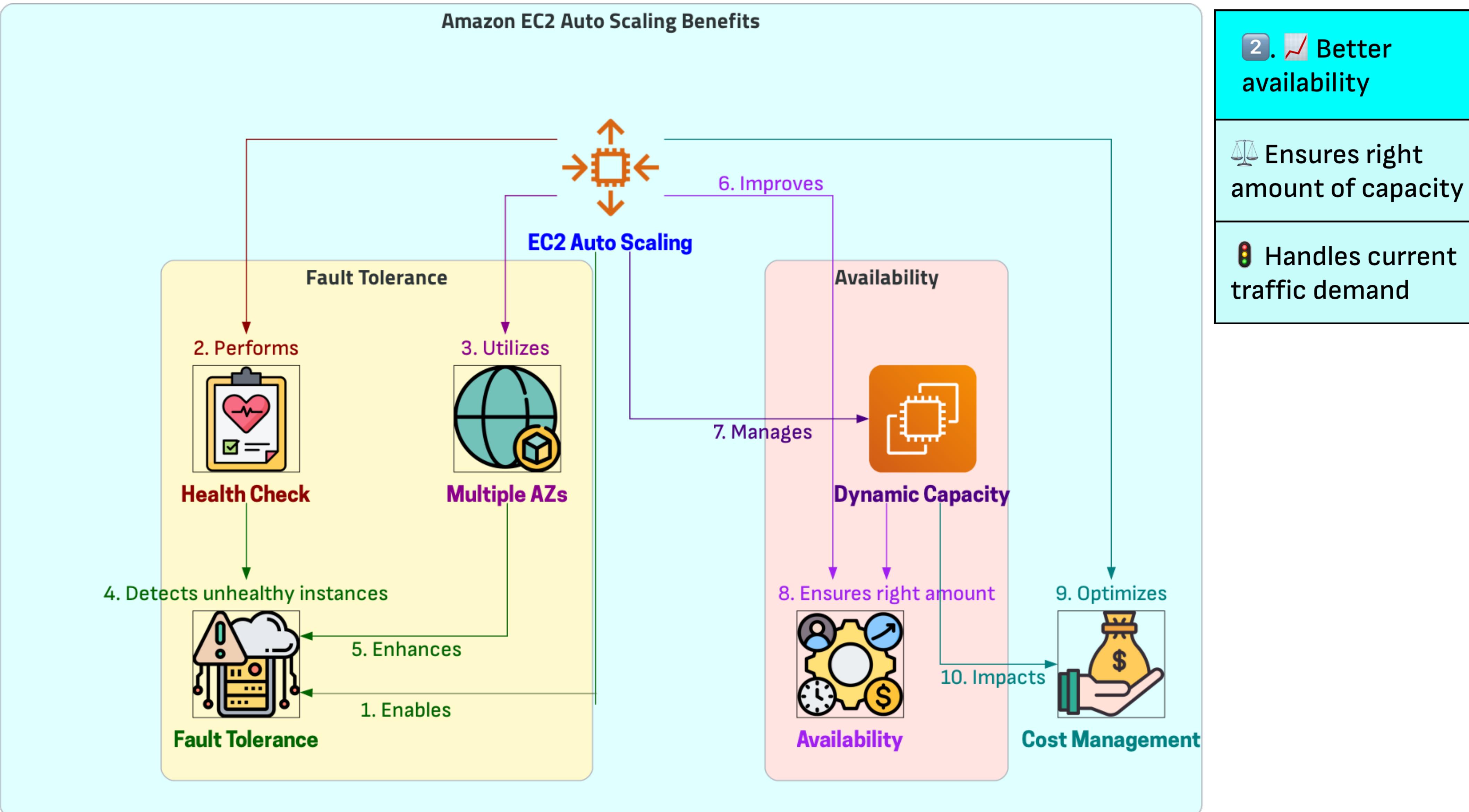


Scale Out/In

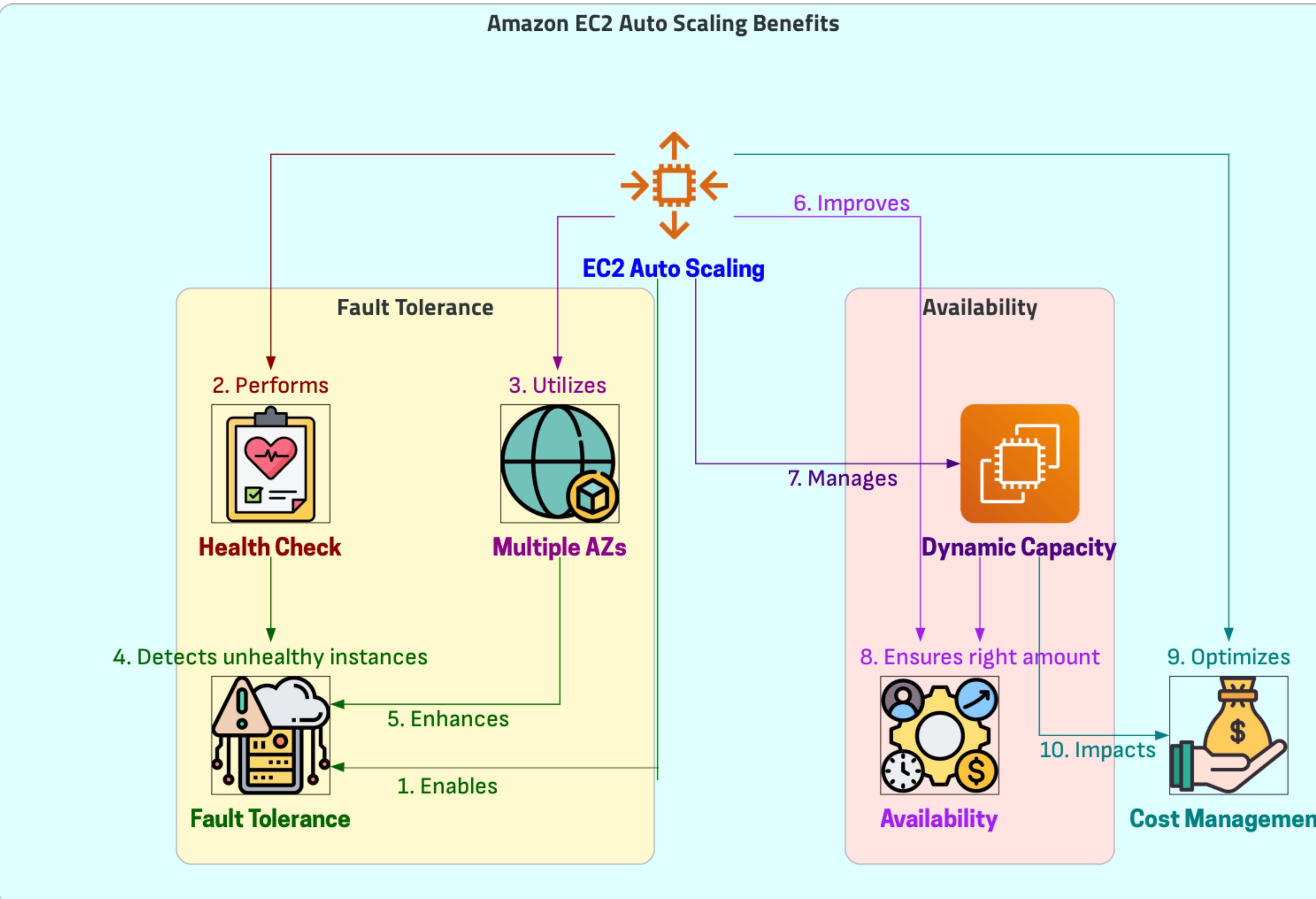
Amazon EC2 Auto Scaling Benefits



Amazon EC2 Auto Scaling Benefits

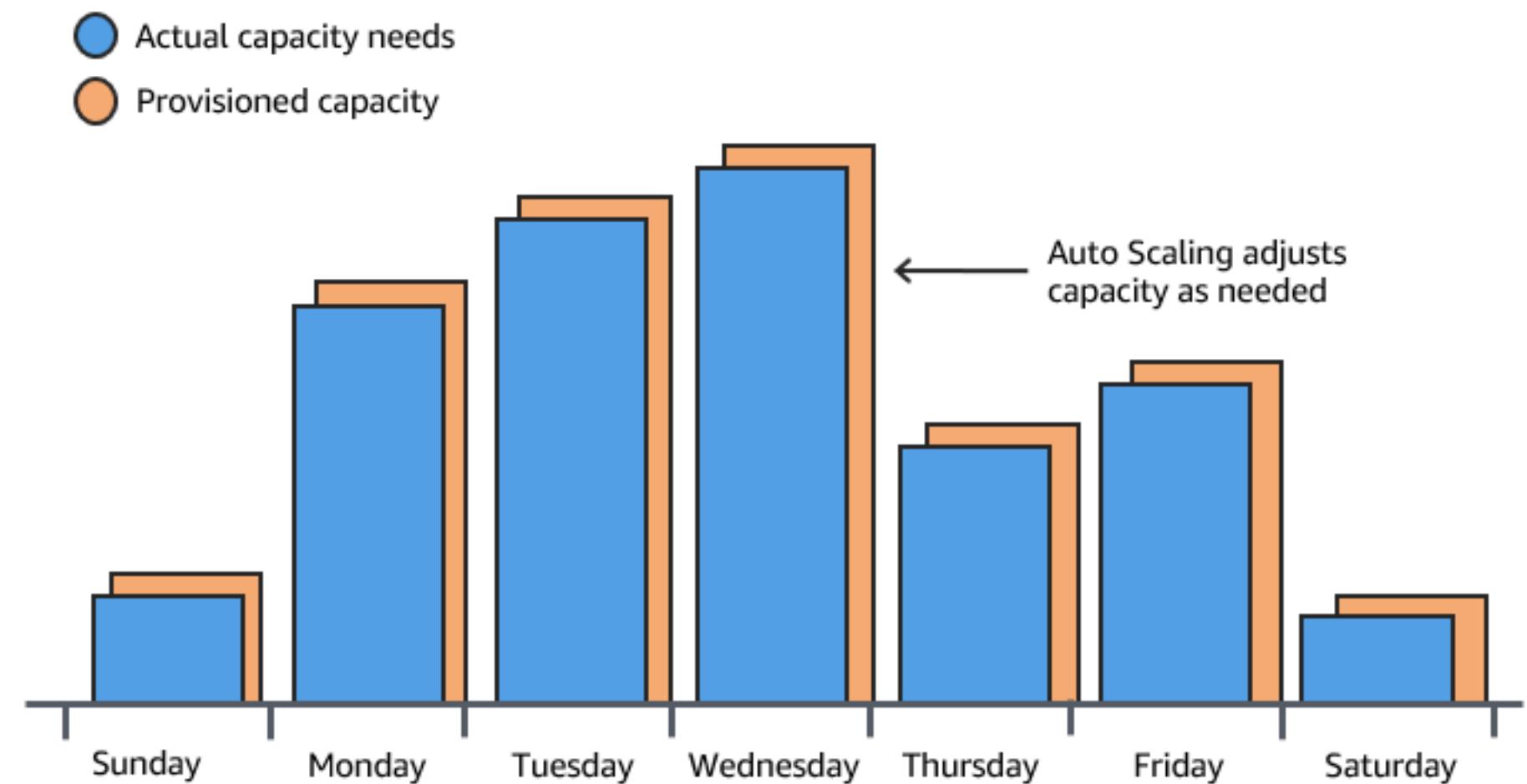
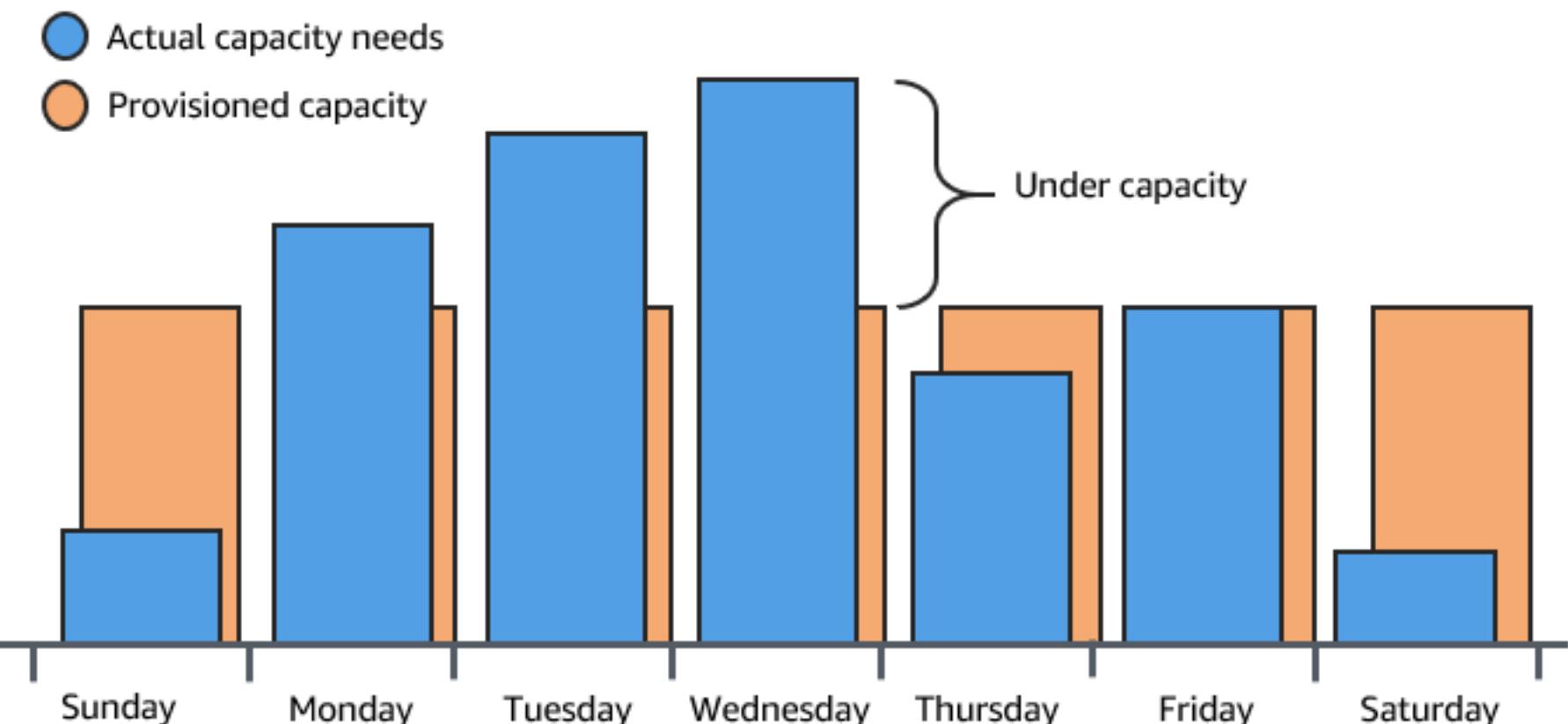
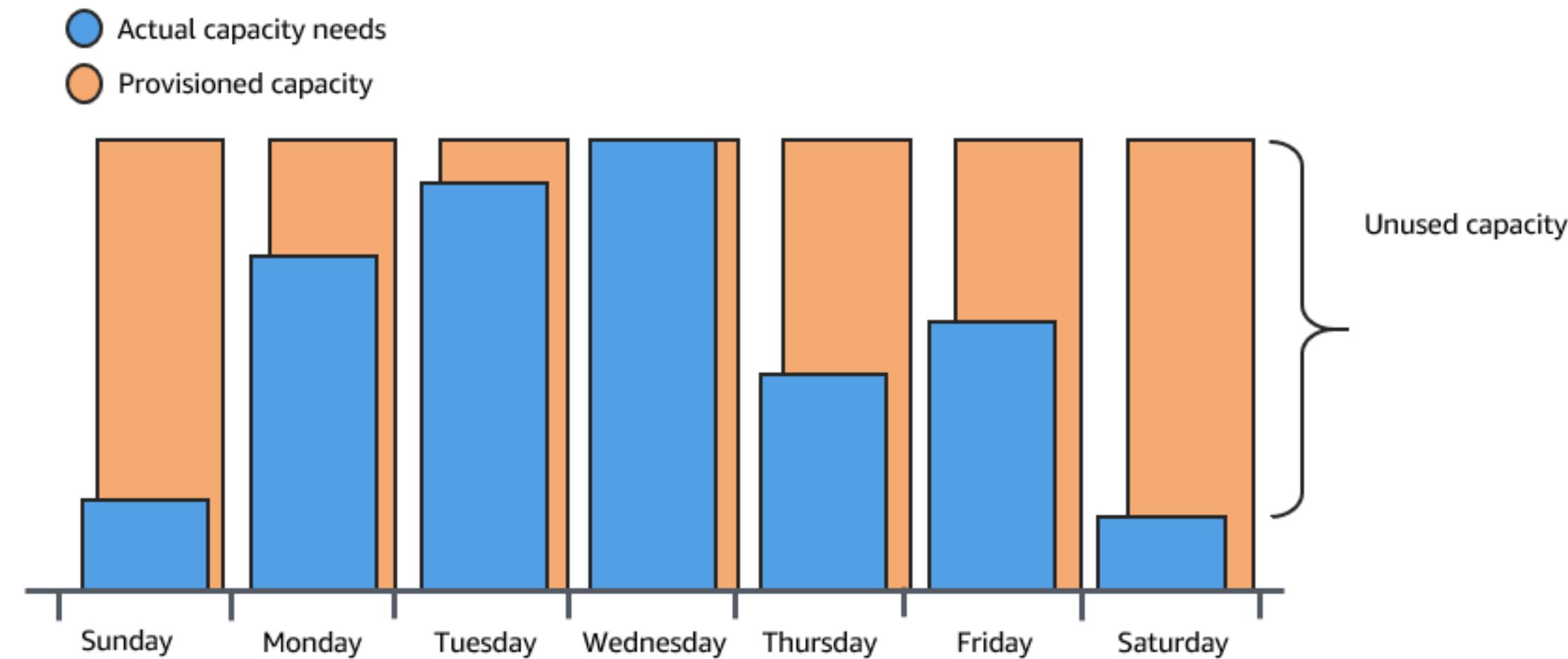
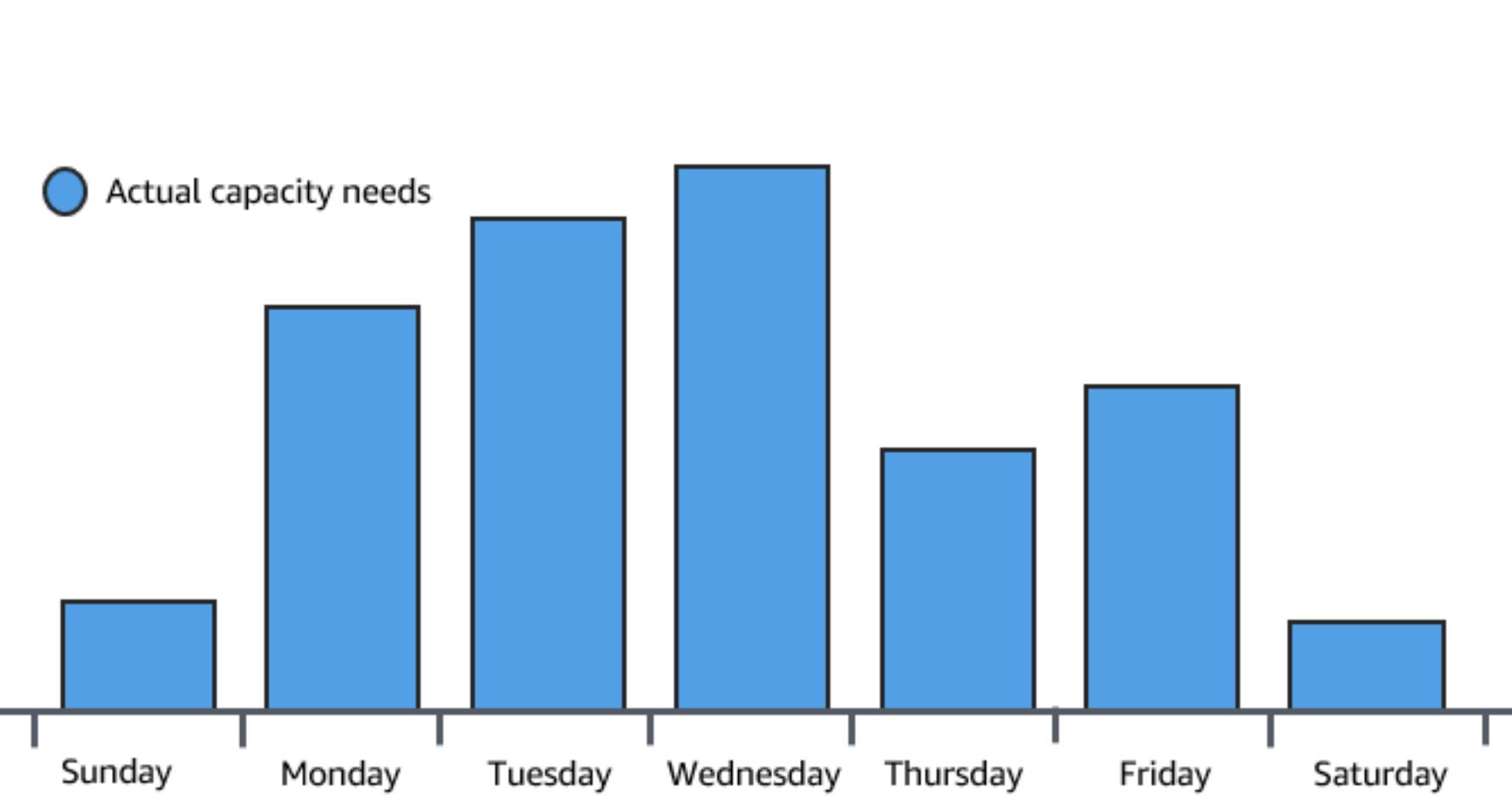


Amazon EC2 Auto Scaling Benefits



3. 💰 Better cost management	
Increases as needed	
Decreases as needed	
Launches when needed	
Pay for used instances only	
	Terminates when not needed

Cover Variable Demand



Typical Web app architecture

1.  Multiple app copies

Simultaneously running

Handles customer traffic volume

2.  EC2 instances

 Identical cloud servers

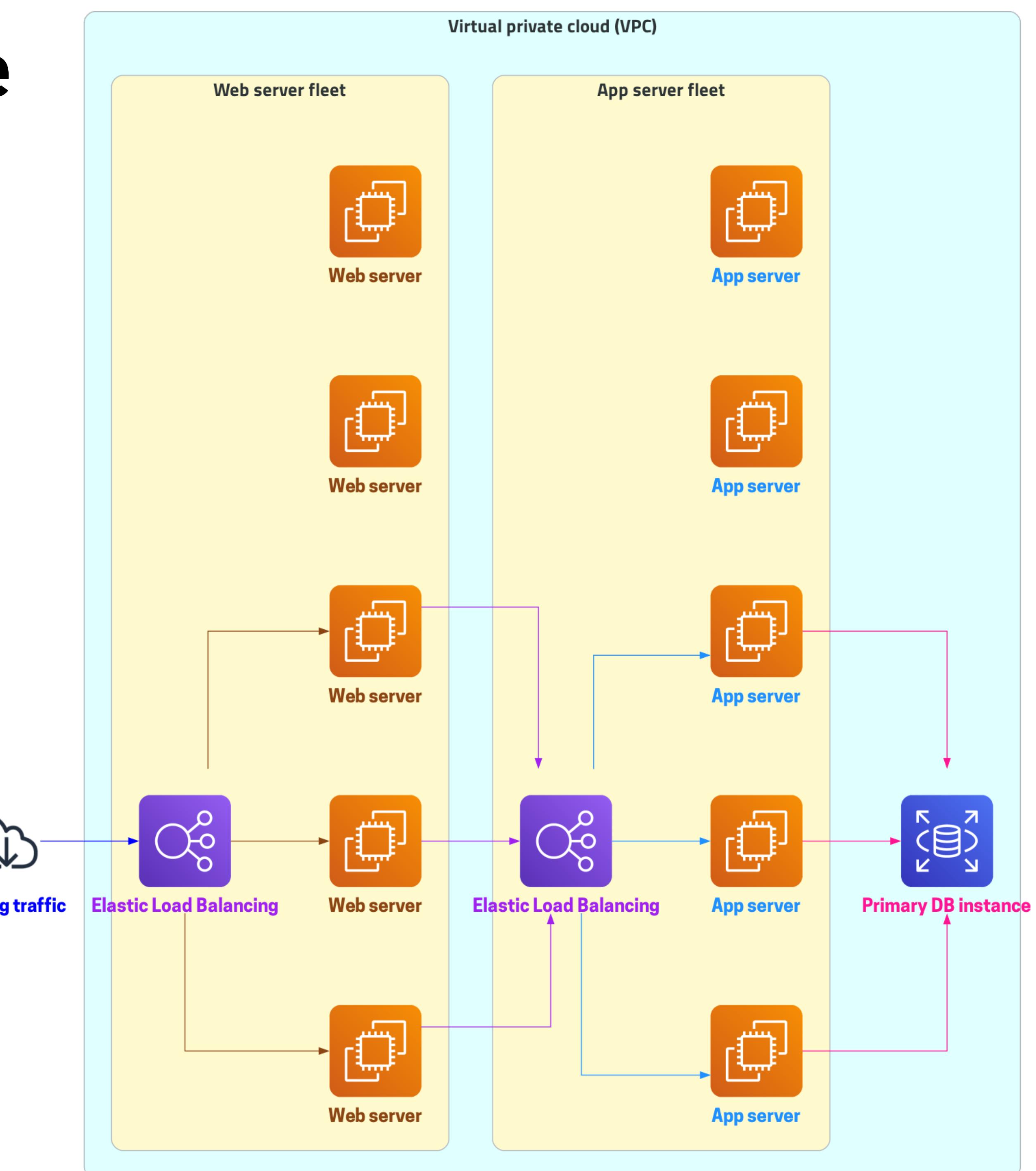
 Process customer requests

3.  Amazon EC2 Auto Scaling

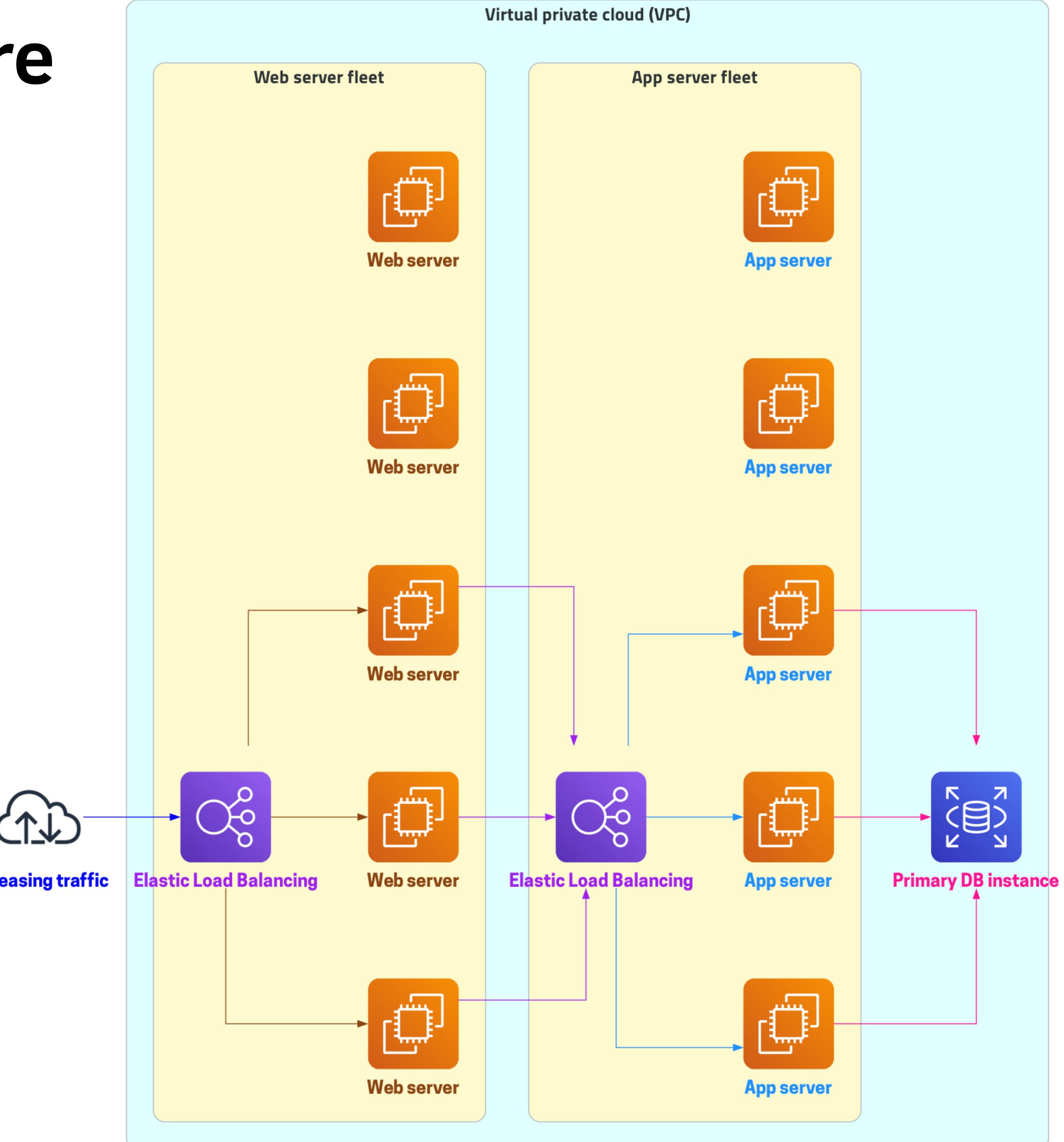
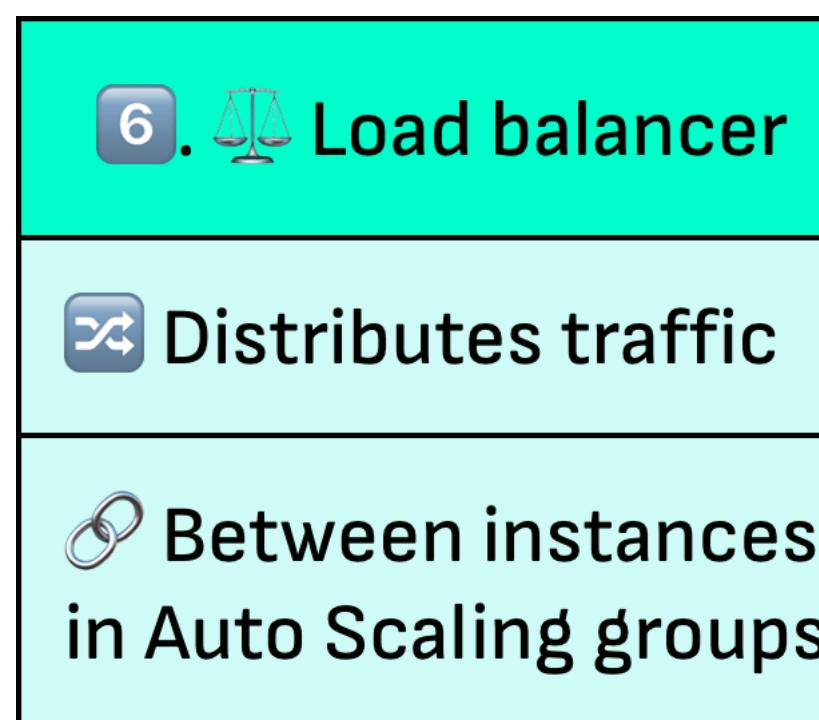
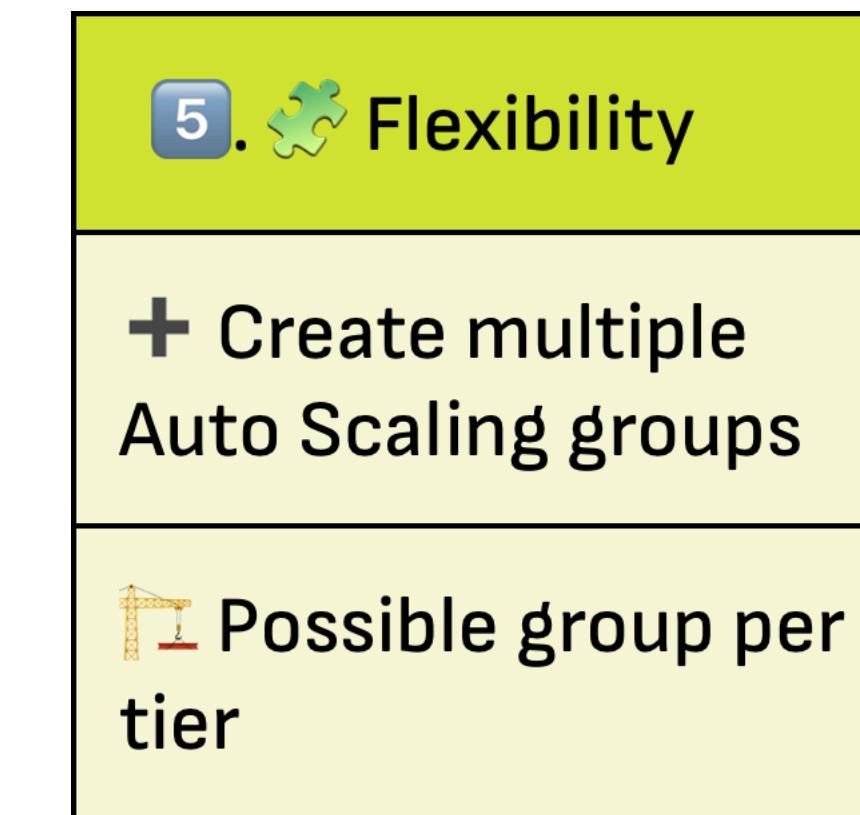
 Manages instance launch

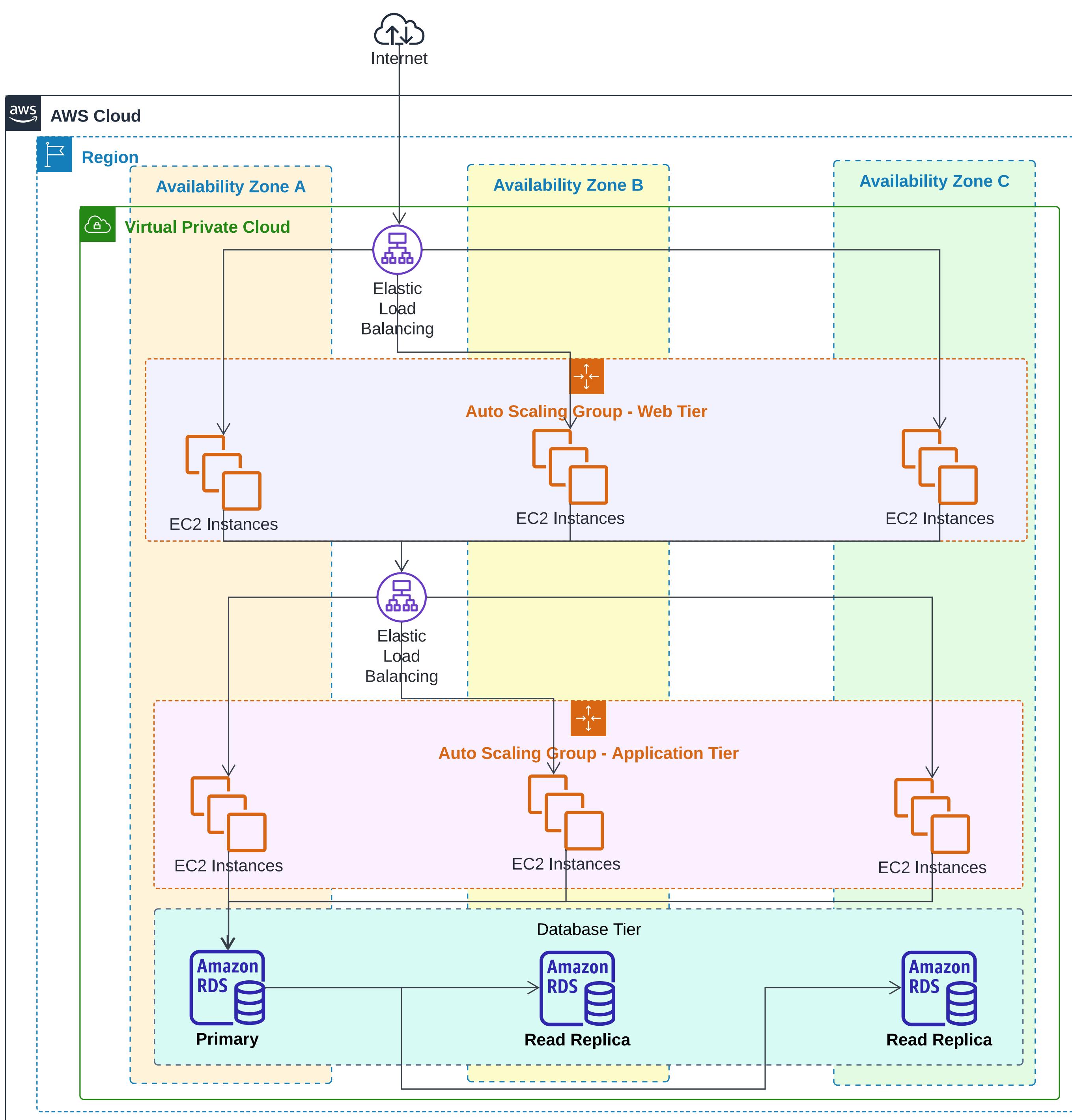
 Manages instance termination

 Based on defined criteria  Amazon CloudWatch alarms



Typical Web app architecture





Distribute instances across Availability Zones

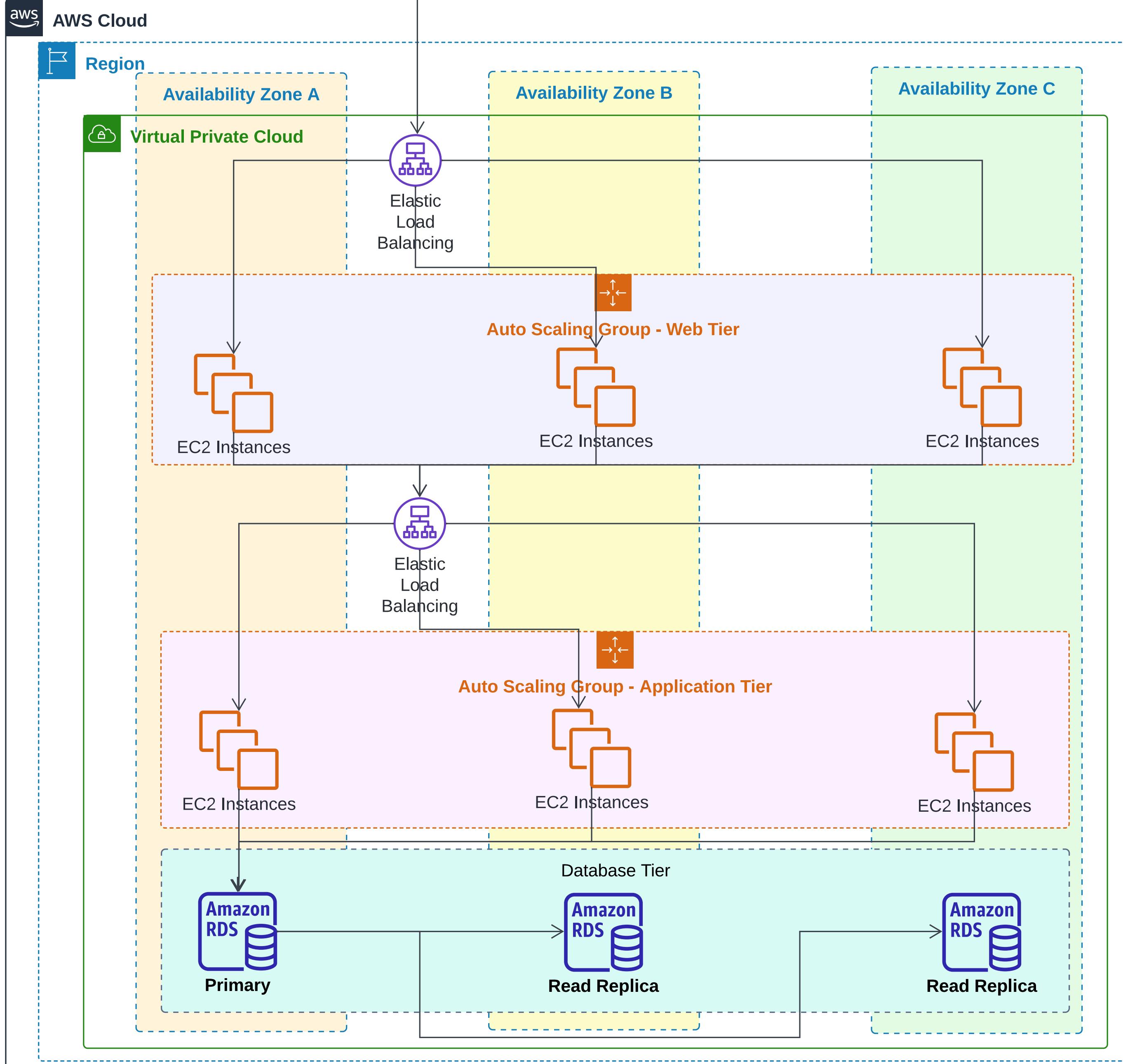
1. Availability Zones (AZs)

- Isolated locations in AWS Region
- Multiple AZs per Region
- Designed for high availability
- Identified by Region code + letter (e.g., us-east-1a)

2. Benefits of using multiple AZs

- Increases application availability
- Provides independence between zones

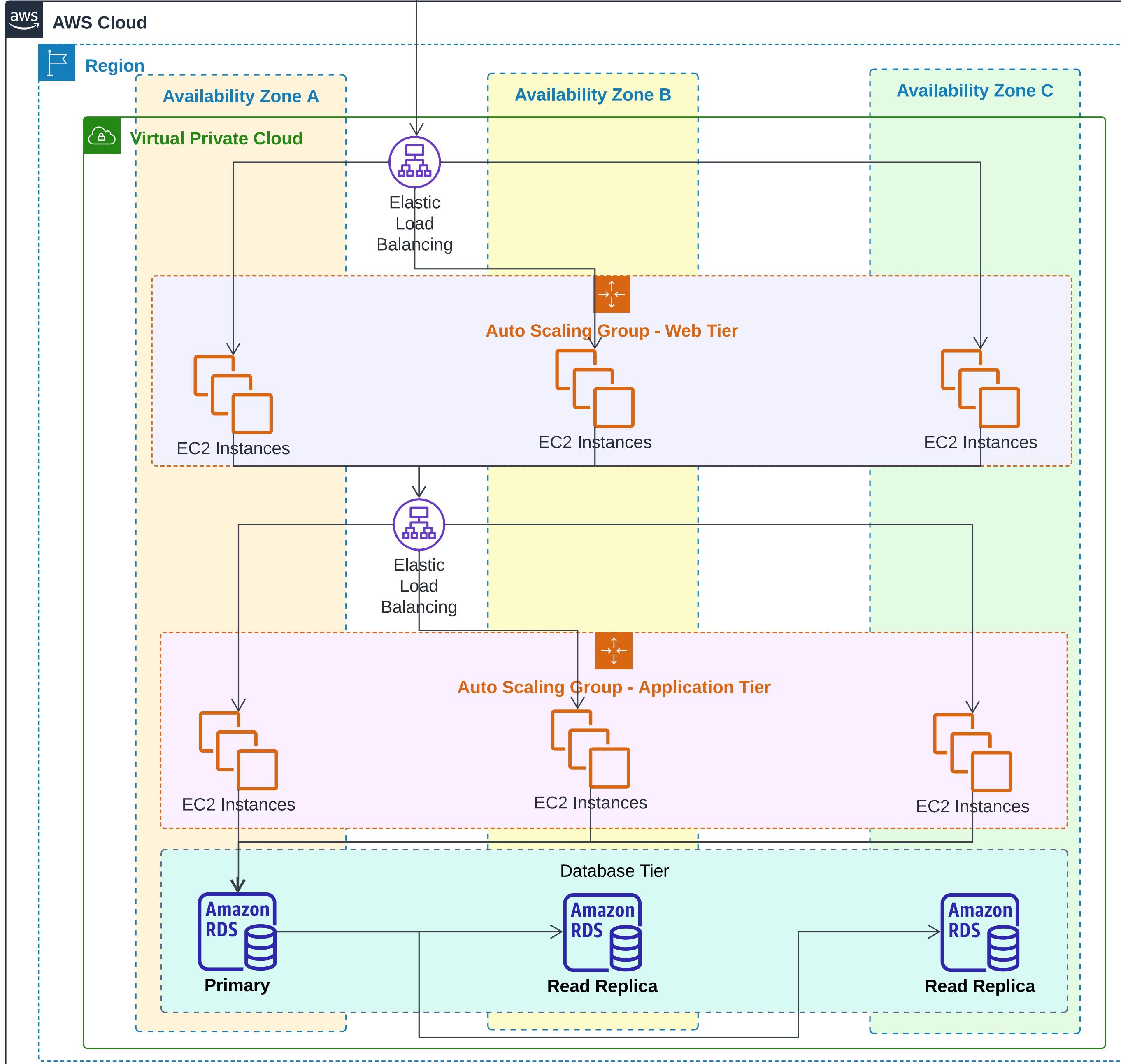
Distribute instances across Availability Zones



3. 🌐 VPC and Subnet Configuration

Create custom VPC and subnets	Define subnets in each AZ
Subnet restrictions	Must reside in one AZ
	Cannot span zones

Distribute instances across Availability Zones



4. 🚀 Auto Scaling Group Deployment

Choose VPC and subnets for deployment

EC2 instances created in chosen subnets

Instances associated with specific AZ

5. ⚖️ Instance Distribution

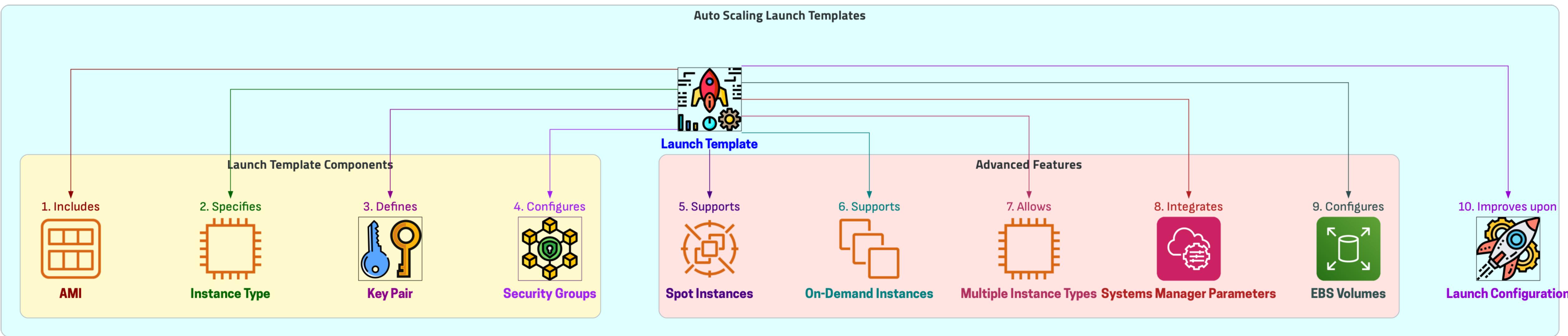
Amazon EC2 Auto Scaling evenly distributes instances

High availability

Goals

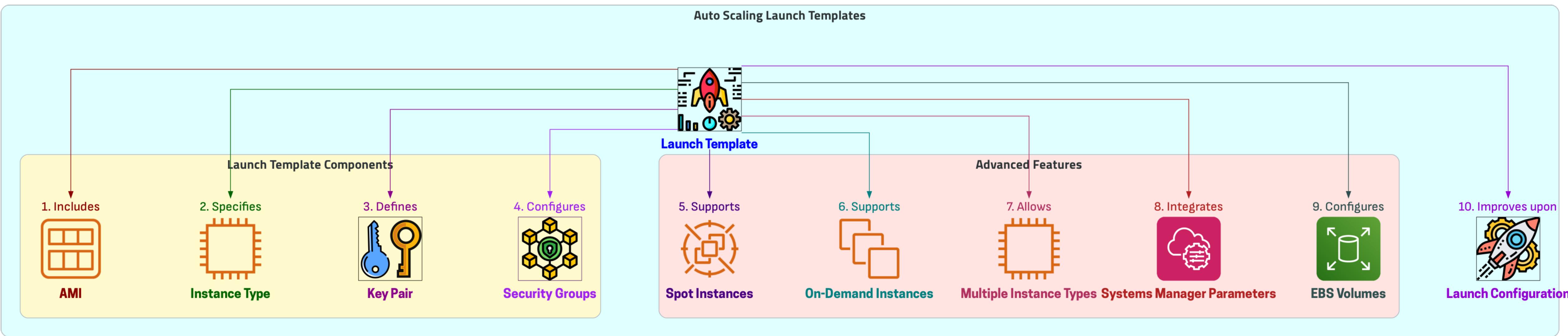
Reliability

Auto Scaling launch templates



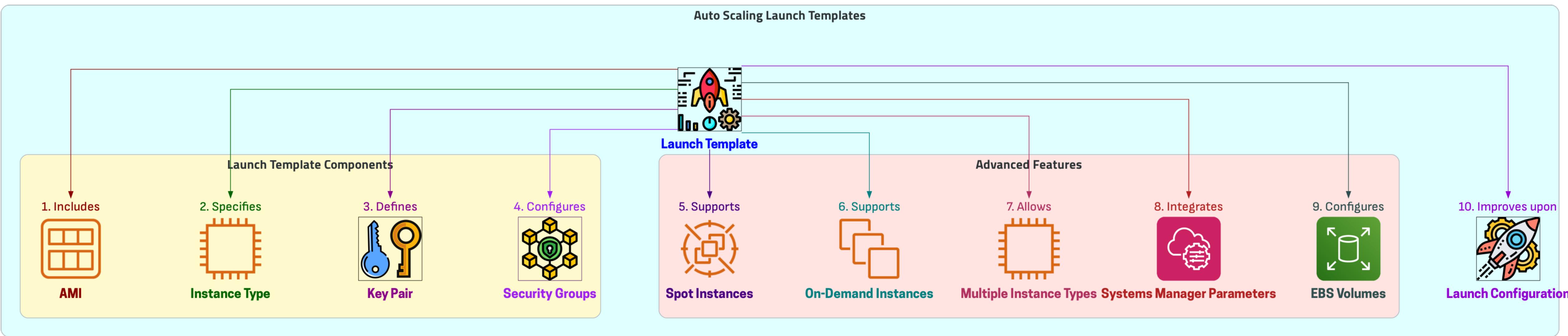
1. Launch Template Overview
Specifies instance configuration
Includes AMI ID, instance type, key pair, security groups
Supports multiple versions

Auto Scaling launch templates



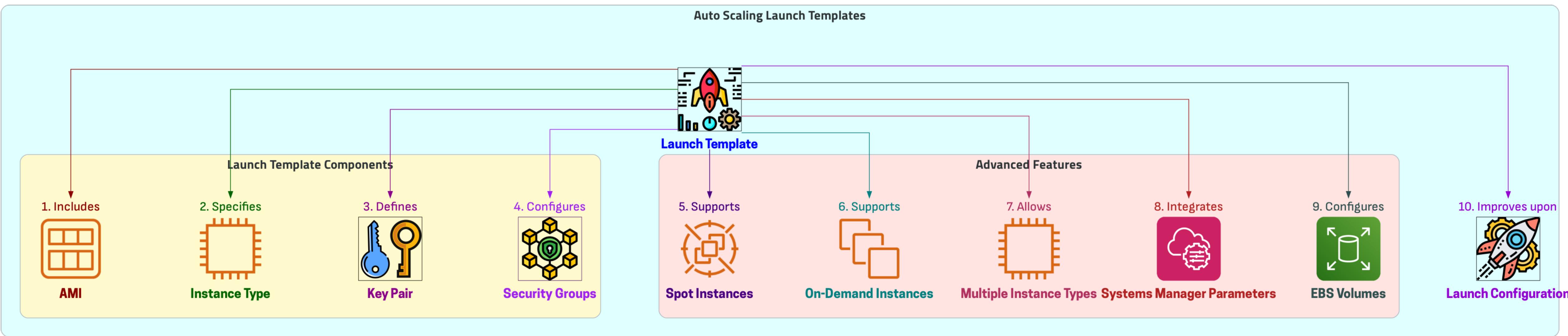
2. <small>12 / 34</small> Versioning Benefits	
	Create subset of parameters
	Reuse for other versions
	Base configuration without AMI/user data
	New version with AMI/user data for testing
	Delete test versions when not needed

Auto Scaling launch templates



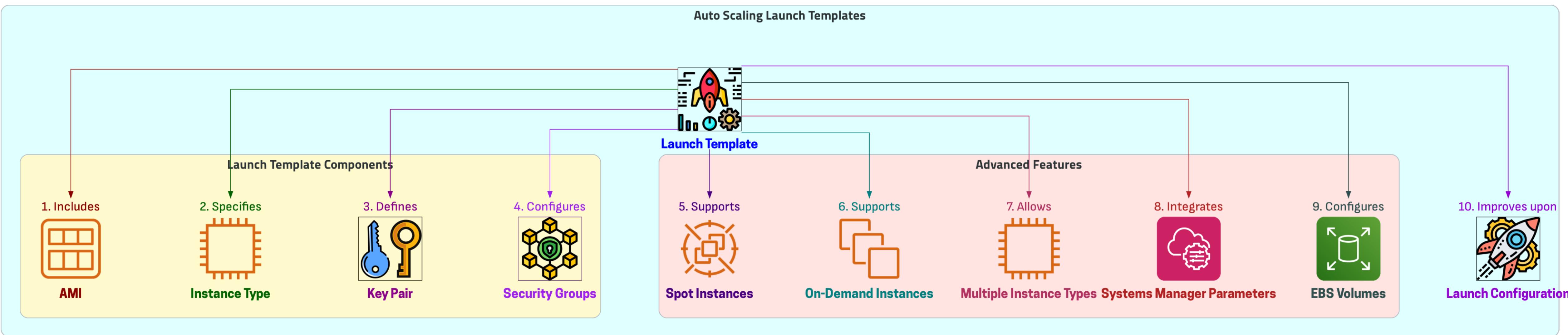
3. vs Comparison with Launch Configurations	
⚠ Some Auto Scaling features unavailable with launch configurations	🚫 Cannot launch both Spot and On-Demand Instances 🚫 Cannot specify multiple instance types
✓ Launch templates required for these features	

Auto Scaling launch templates



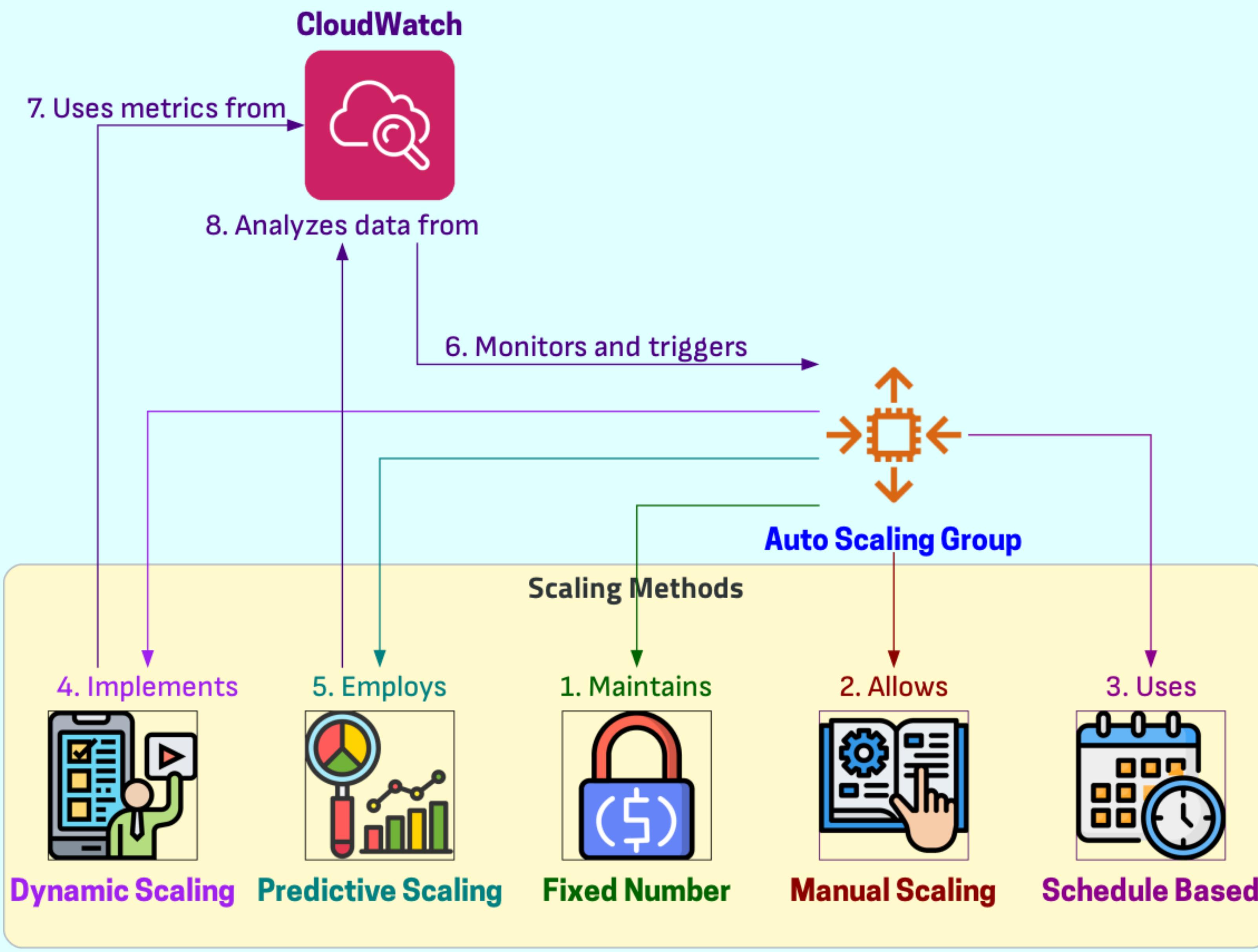
4. NEW	Access to Newer EC2 Features
	Systems Manager parameters (AMI ID)
	EBS Provisioned IOPS volumes (io2)
	EBS volume tagging
	T2 Unlimited instances
	Capacity Reservations and Blocks
	Dedicated Hosts

Auto Scaling launch templates



5. 🚀 Launch Template Creation	
🔒 All parameters optional	⚠️ Cannot add AMI later if not specified
⚠️ Limitations	💻 Can add instance types later if AMI specified

Auto Scaling Methods



Auto Scaling Methods

1. Maintain a Fixed Number of Instances

Keep constant group size

Replace unhealthy instances

Maintain desired capacity

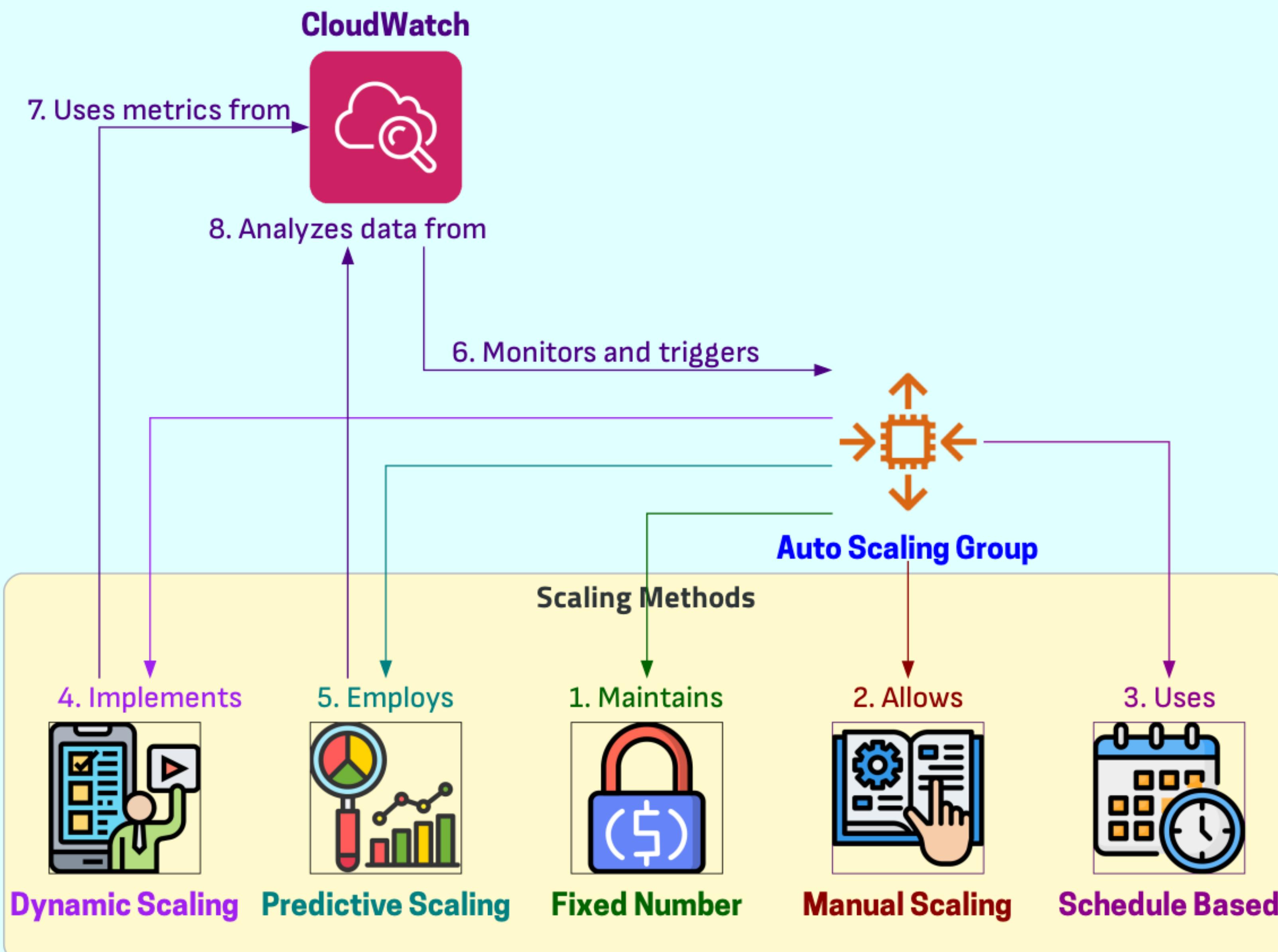
2. Manual Scaling

Adjust desired capacity

Terminate specific instances

Immediate control

Auto Scaling Methods



Auto Scaling Methods

3. Scale Based on a Schedule

Automate scaling actions

Based on predictable demand

Align with expected traffic

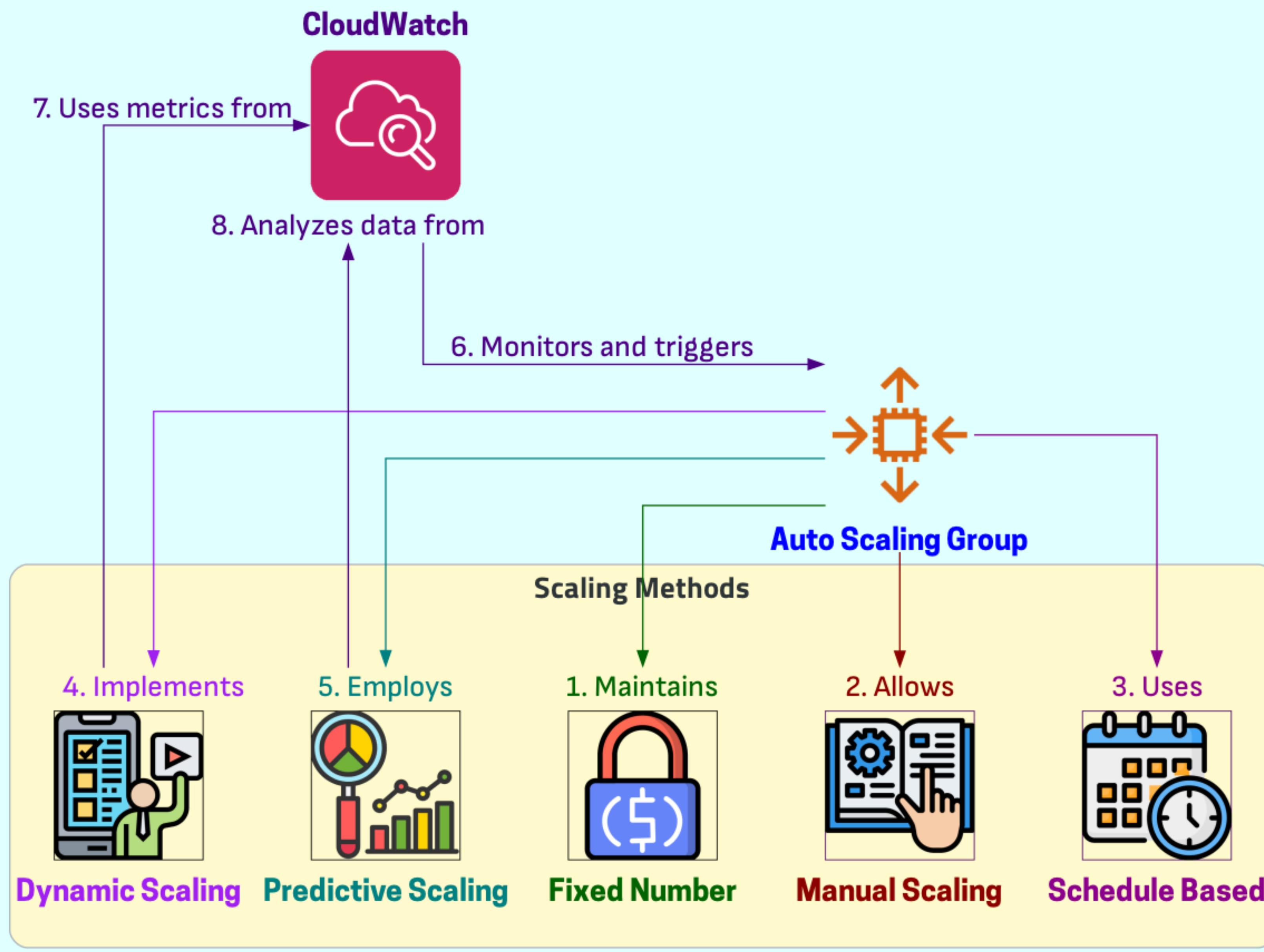
4. Scale Dynamically Based on Demand

Use scaling policies

Adjust capacity in real-time

Ensure optimal performance

Auto Scaling Methods



Auto Scaling Methods

5. Scale

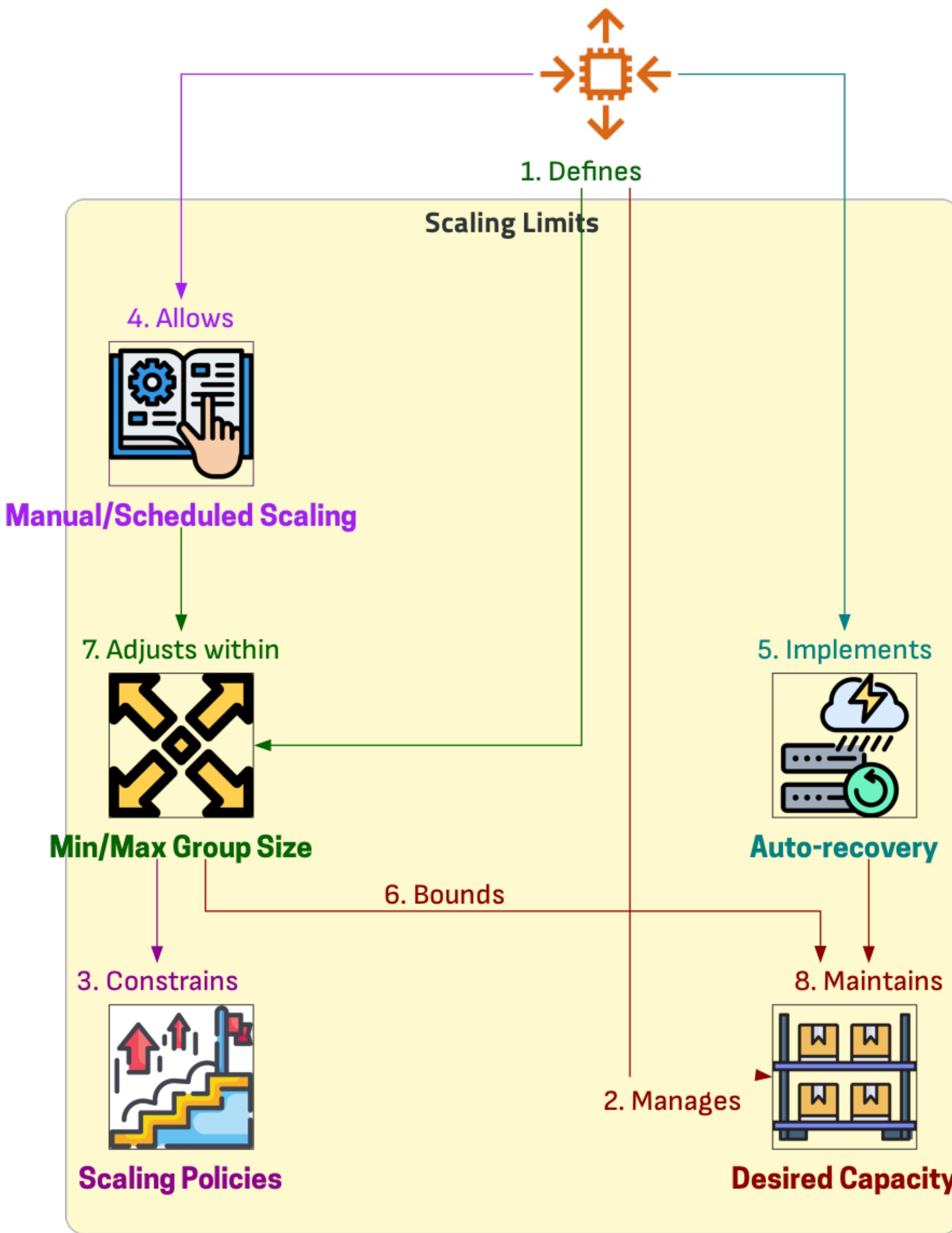
Proactively with
Predictive Scaling

Combine predictive
and dynamic scaling

Proactively adjust
resources

Ready for daily and
weekly fluctuations

Auto Scaling Group



Set Scaling Limits for Your Auto Scaling Group

1. 🚀 Minimum and Maximum Group Size

🛡️ Defines boundaries for group size

✖️ Ensures operational flexibility

↙ Within predefined limits

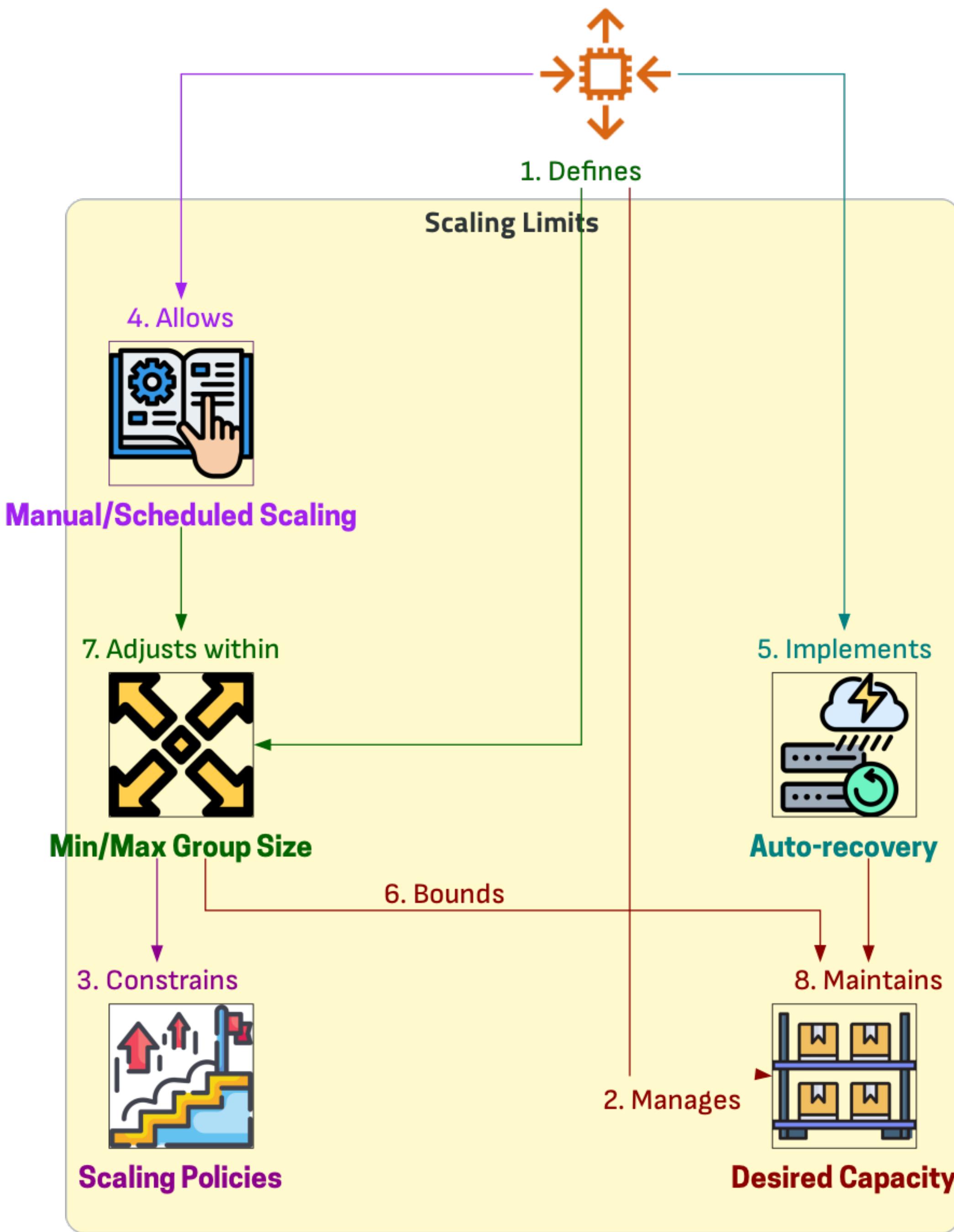
2. ⏹ Desired Capacity Management

🏁 Initial group capacity at creation

⚖️ Maintains level unless altered

⟳ Affected by scaling policies or scheduled actions

Auto Scaling Group



Set Scaling Limits for Your Auto Scaling Group

3. ▼ Minimum Capacity Constraints

🚫 Prevents reduction below minimum

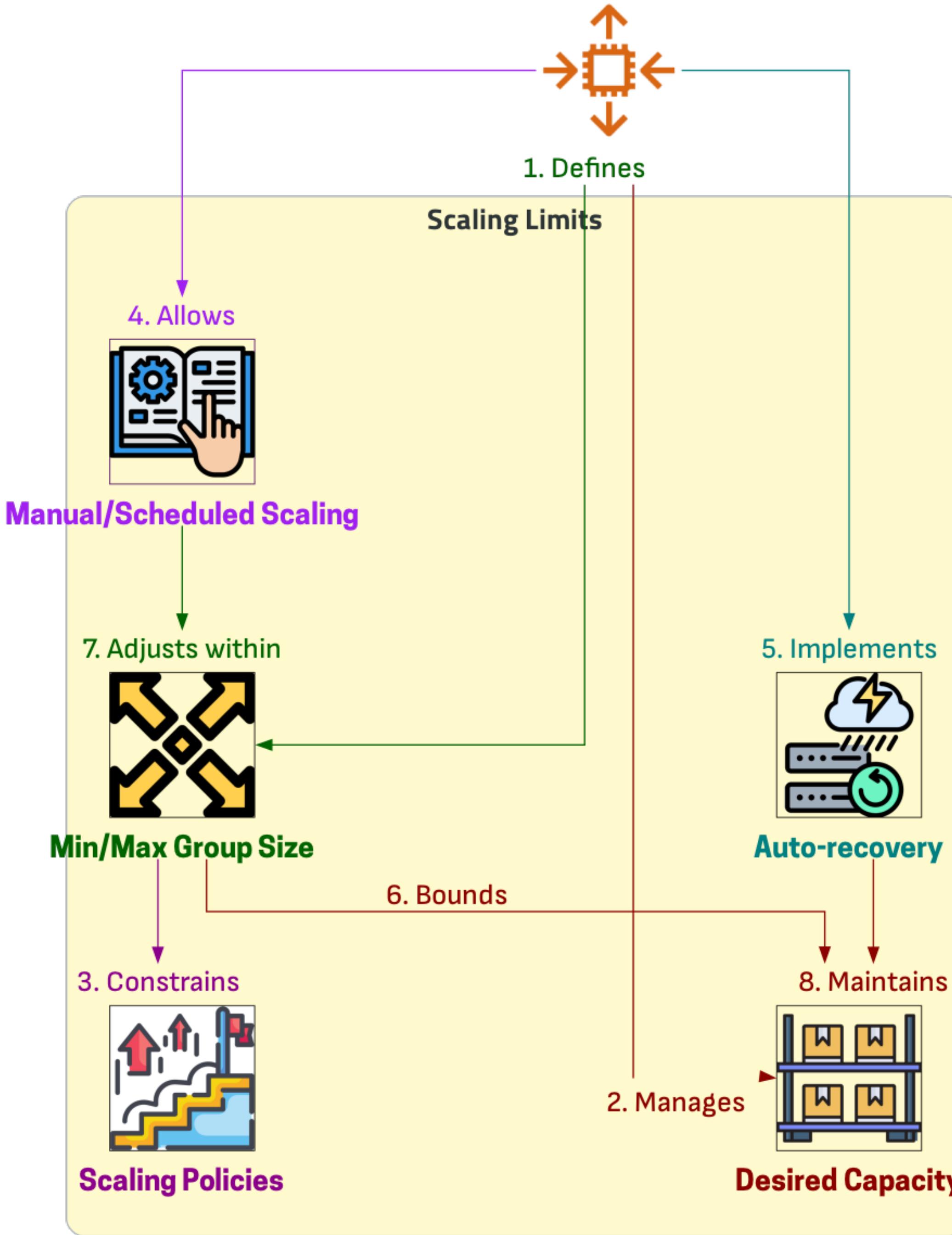
🛡 Safeguards against resource inadequacy

4. ▲ Maximum Capacity Constraints

🚫 Prevents exceeding maximum size

🎛 Controls resource overutilization

Auto Scaling Group



Set Scaling Limits for Your Auto Scaling Group

5. Manual and Scheduled Scaling Impact

⌚ Manual adjustment capability

⏰ Scheduled actions for capacity changes

🔧 Within set minimum and maximum limits

6. Auto-recovery of Desired Capacity

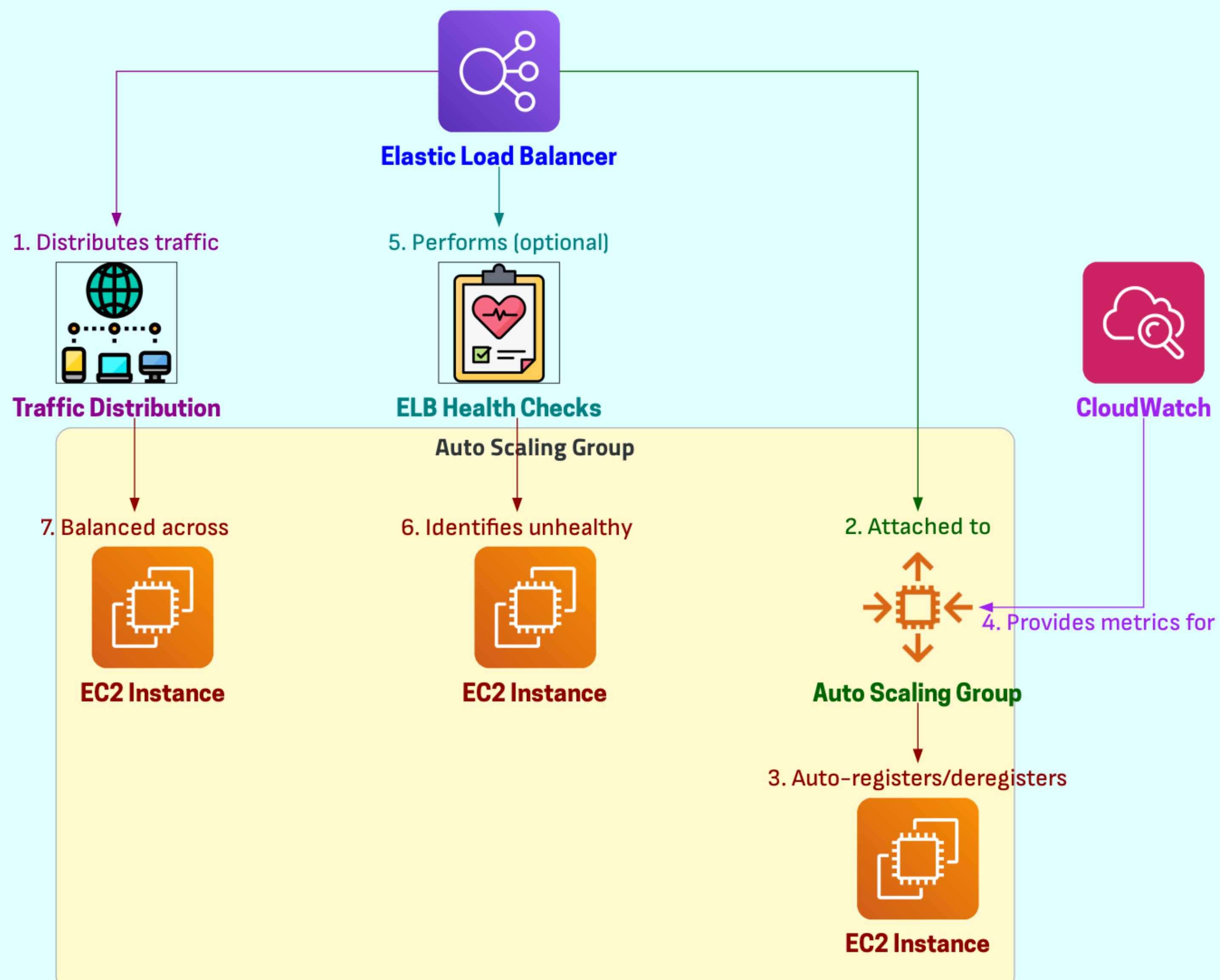
🚀 Automatically launches new instances

🎯 Maintains desired capacity

🛡 Handles unexpected terminations

💪 Ensures resilience and continuity

Elastic Load Balancing with Auto Scaling



Elastic Load Balancing with Auto Scaling

1. ⚡ Automatic Traffic Distribution

Evenly distributes incoming traffic

Prevents instance overload

2. 🔧 Attach Load Balancer to Auto Scaling Group

Single point of contact for traffic

Simplifies setup

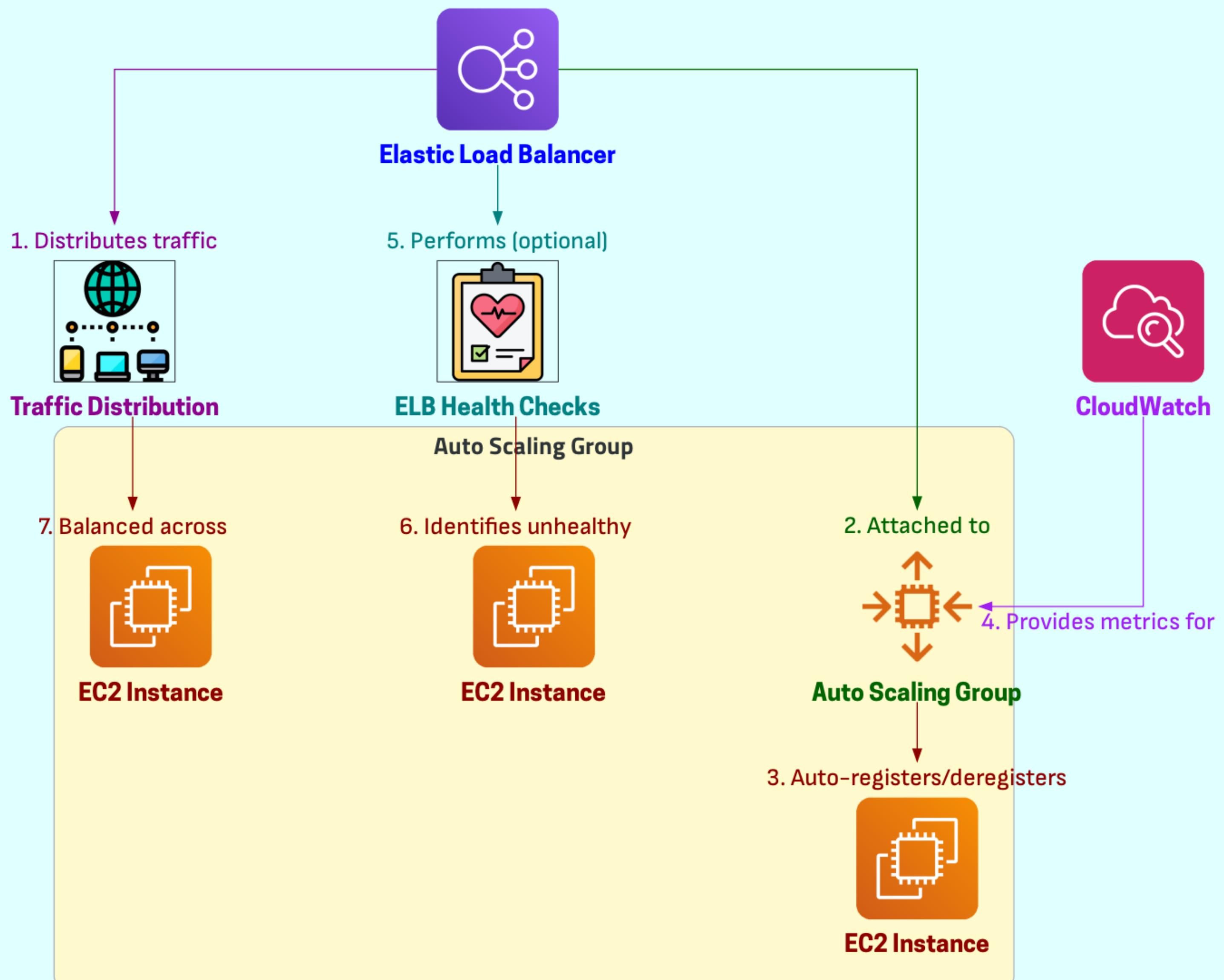
Enhances efficiency

3. 🚫 No Manual Instance Registration Needed

Automatic instance registration

Automatic instance deregistration

Elastic Load Balancing with Auto Scaling



Elastic Load Balancing with Auto Scaling

4. Use ELB Metrics for Scaling

Dynamic scaling based on demand

Utilizes ELB metrics

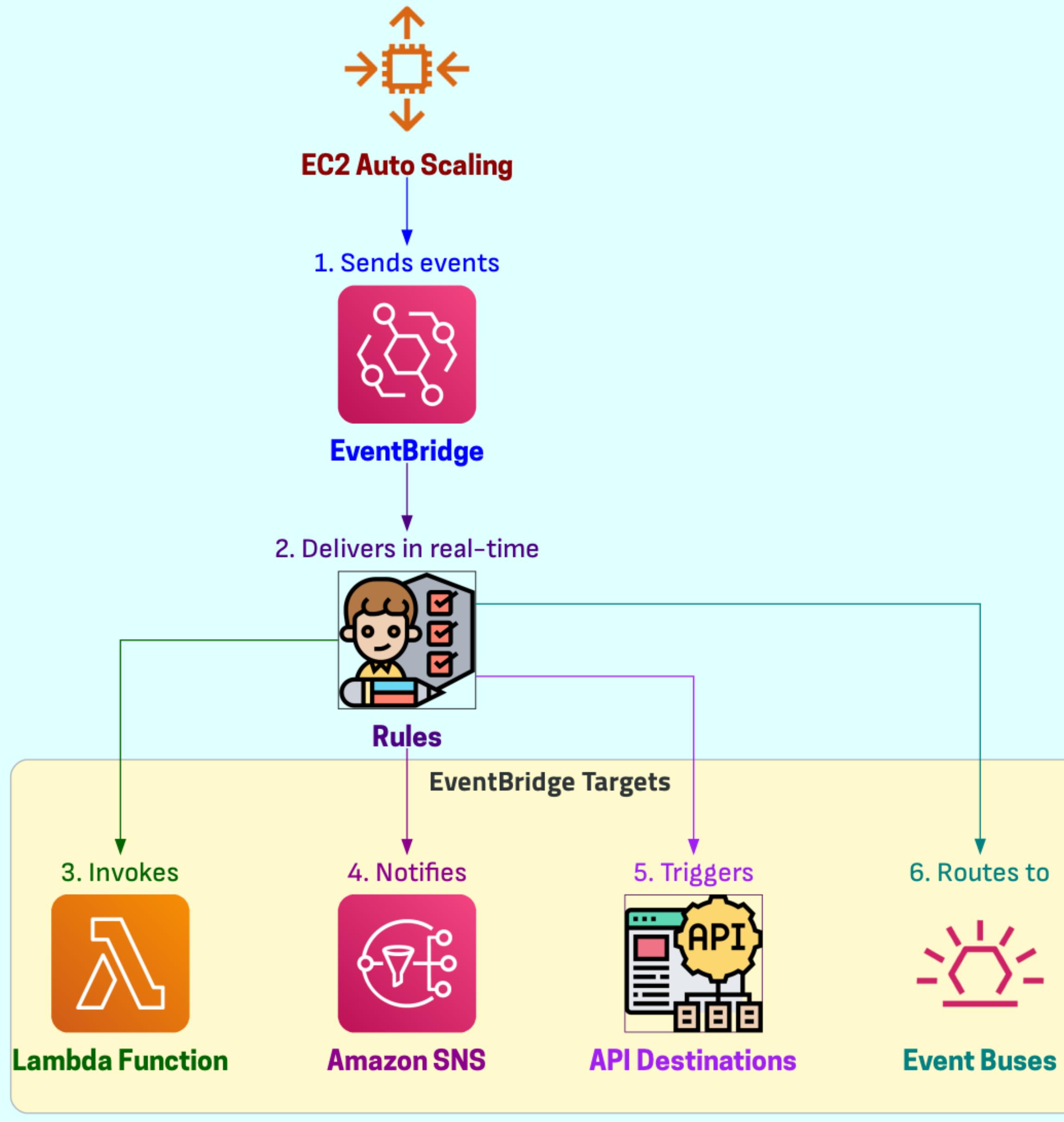
Application Load Balancer request count

5. Optional: ELB Health Checks

Identifies unhealthy instances

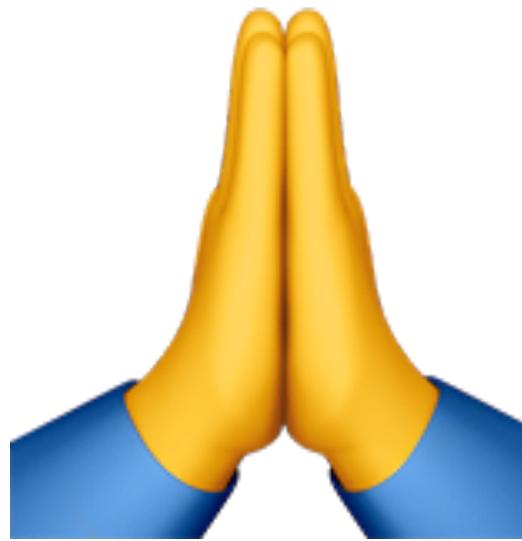
Replaces unhealthy instances

Enhances reliability



Use EventBridge for Auto Scaling Events

1. EventBridge Introduction	2. Real-time Event Delivery
Monitors resources	Nearly instantaneous delivery
Initiates actions	Allows timely responses
Based on rules	
3. Rules for Programmatic Actions	4. Wide Range of Targets
Create rules in EventBridge	AWS Lambda functions
Invoke actions or notifications	Amazon SNS topics
Launching instances	API destinations
Respond to EC2 Auto Scaling events	Event buses
Terminating instances	Diverse set of responses



**Thanks
for
Watching**