# AKS Hero Updates

Kunal Chandratre (Sr. Cloud Solution Architect)
@Microsoft

Microsoft

# Disclaimer*

- All views expressed in this session are personal, in no way it represents the company I work for.

Microsoft

# Agenda

- Orchestration in containers

- AKS Updates

- AKS and AI

- Demos

- Wrap up!

Microsoft

# Is container a new concept?

# Popular Orchestrators

Docker Swarm

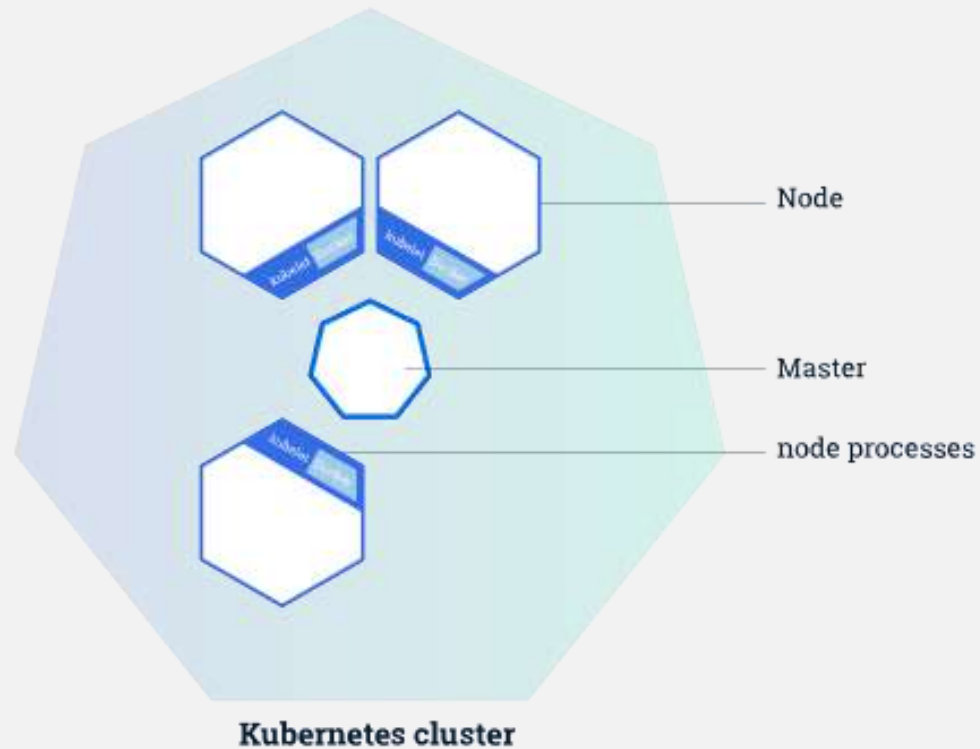Apache Mesos – Marathon

Kubernetes

Microsoft

# What is Kubernetes?

*"Kubernetes is an open-source system for automating deployment, scaling, and management of containerized applications."*

**Kubernetes** comes from the Greek word **κυβερνήτης:**, which means *helmsman* or *ship pilot*, ie: the captainer of a container ship.

# Kubernetes



Kubernetes cluster
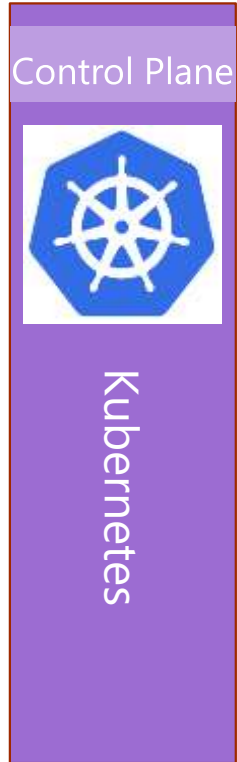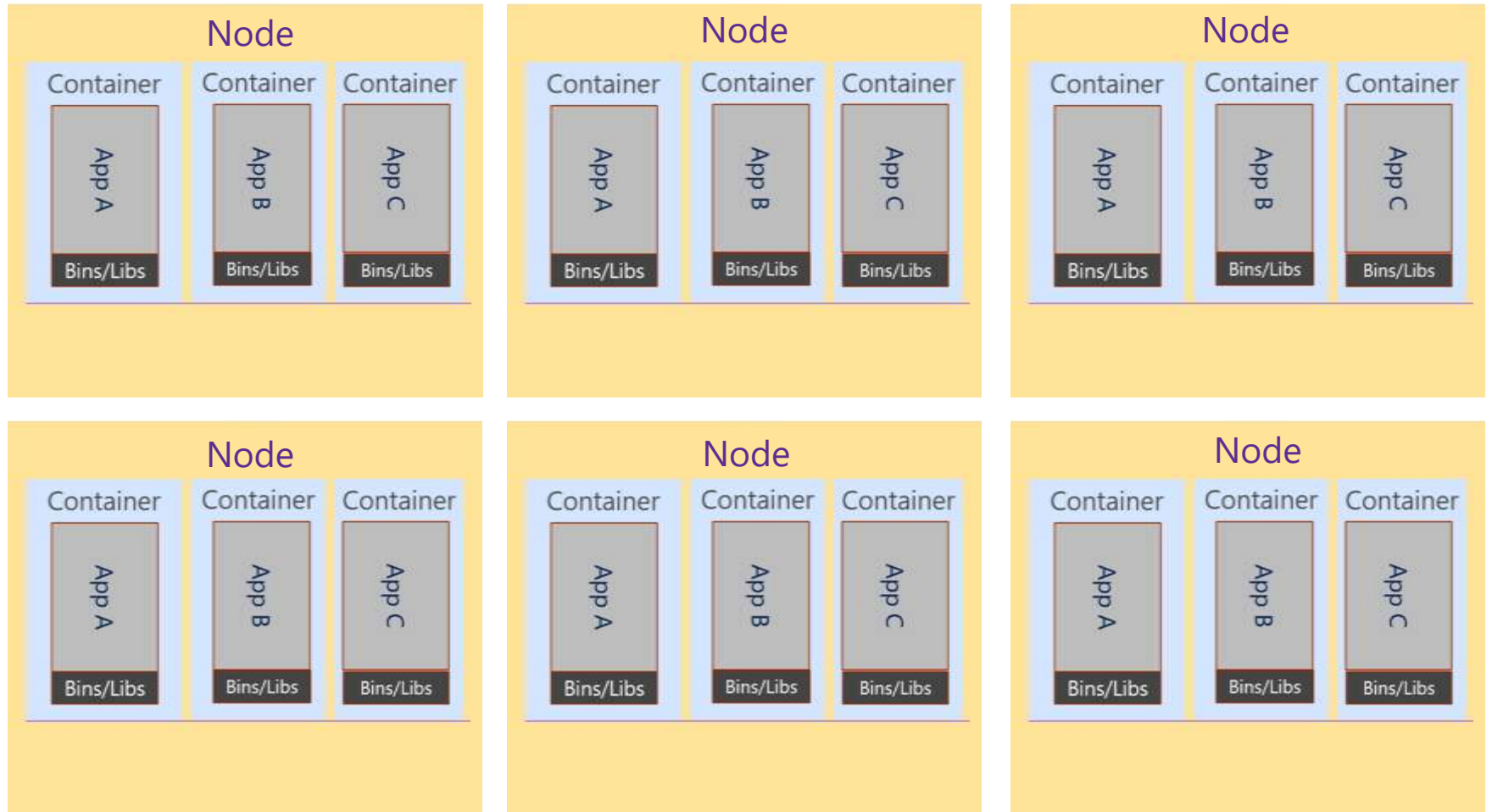
- **Master**: The server that runs the Kubernetes management processes, including the API service, replication controller and scheduler.

- **Node**: The host that runs the kubelet service and the Docker Engine. Minions receive commands from the master.

- **Kubelet**: The node-level manager in Kubernetes; it runs on a minion.

- **Pod**: The collection of containers deployed on the same minion.

- **Replication controller**: Defines the number of pods or containers that need to be running.

- **Service**: A definition that allows the discovery of services/ports published by each container, along with the external proxy used for communications.

- **Kubecfg**: The command line interface that talks to the master to manage a Kubernetes deployment.

**Note – In any orchestrator High Availability of Infra has to be managed by you.** ■■ Microsoft

# AKS is Kubernetes (Control Plane) as a Service

**Control Plane**

Kubernetes

## Cluster

| Node | Node | Node |
|---|---|---|
| Container App A Bins/Libs · Container App B Bins/Libs · Container App C Bins/Libs | Container App A Bins/Libs · Container App B Bins/Libs · Container App C Bins/Libs | Container App A Bins/Libs · Container App B Bins/Libs · Container App C Bins/Libs |
| Node | Node | Node |
| Container App A Bins/Libs · Container App B Bins/Libs · Container App C Bins/Libs | Container App A Bins/Libs · Container App B Bins/Libs · Container App C Bins/Libs | Container App A Bins/Libs · Container App B Bins/Libs · Container App C Bins/Libs |

# AKS Hero Features - Updates

# AKS Policies

- https://learn.microsoft.com/en-us/azure/aks/policy-reference

- Security baseline - https://learn.microsoft.com/en-us/security/benchmark/azure/baselines/aks-security-baseline

# Example - AKS Policy to save cost

1. Assign policy CPU Memory resource limits for PODs.

2. Deploying pods with higher cpu/ memory threshold than set limit, should fail.

3. Deploying pods with equal to or below cpu/ memory threshold than set limit, should succeed.

# Cluster Optimization workbook for Azure Monitor Container Insights

**Azure Monitor Container insights for collection of logs and events now has capability to help you optimize your cluster**
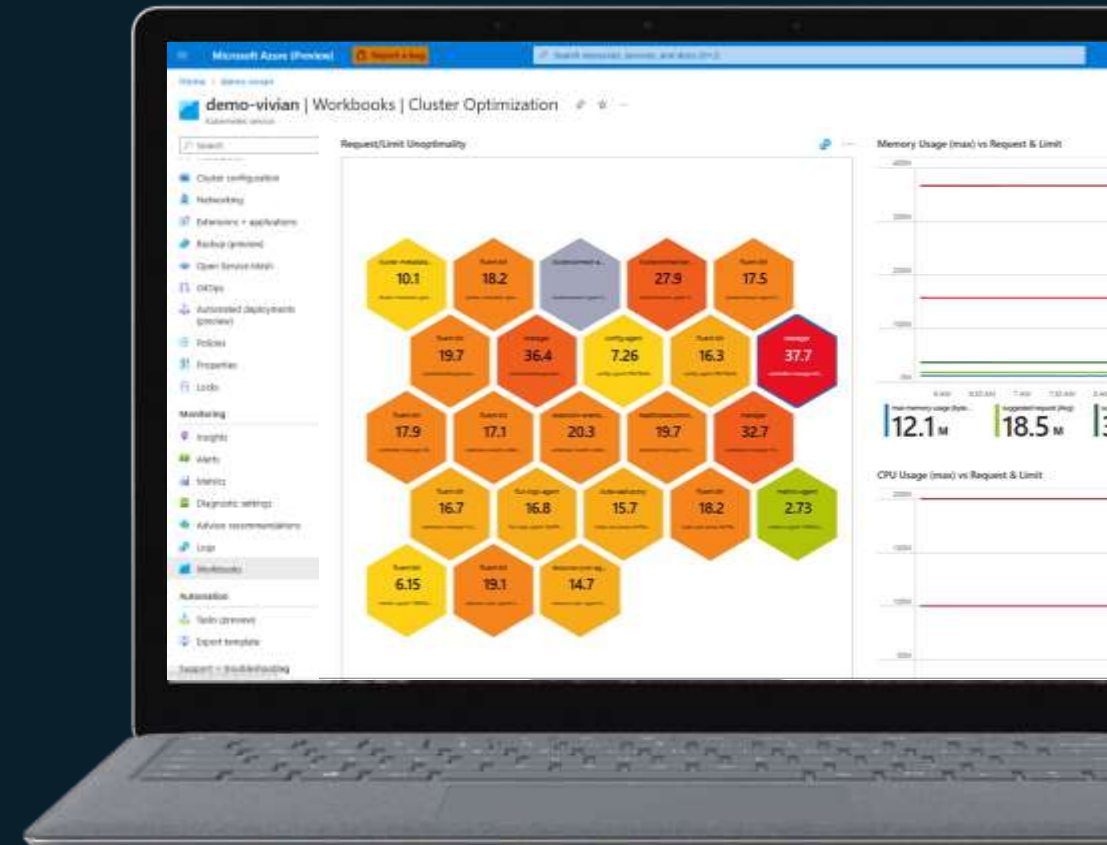
## Detect liveness probe failures

Detect liveness probe failures and their frequencies.

## Identify event anomalies

Identify and group event anomalies that indicate recent increases in event volume for easier analysis.

## Optimize container limits and requests

Identify containers with high or low CPU and memory limits and requests, along with suggested limit and request values.

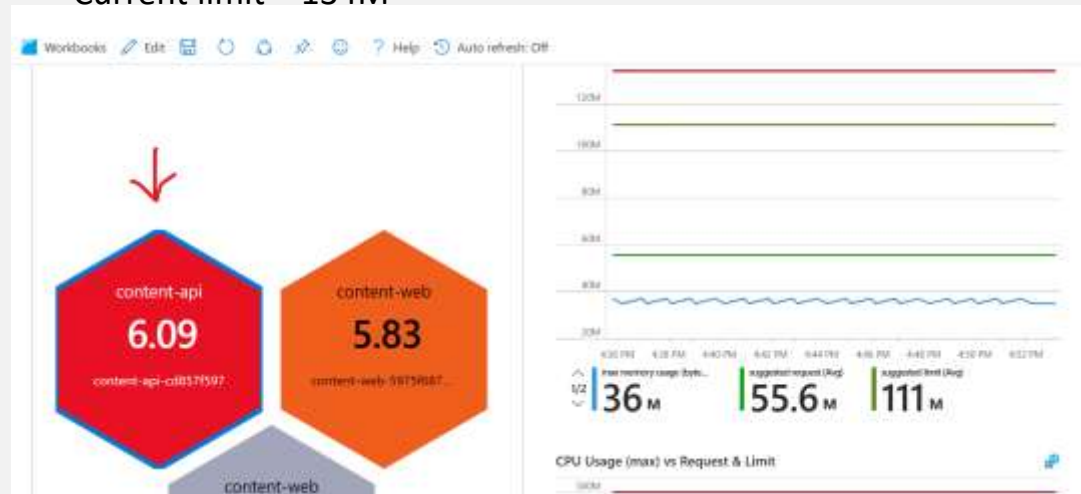# Container Optimizer Example

**For this container – Memory -**

Max used – 36MB
Suggested request – 55MB
Suggested limit – 111MB
Current request – 134MB
Current limit – 134M

# Demo – Cluster Optimization Workbook

1. Red – excessive request and limit is assigned to containers within a pod
2. Green – well set request and limit
3. Gray – no limit or request is set
4. Value closure to zero would be better
5. Show on portal

# AKS Cost Analysis

*aka.ms/aks/cost-analysis*

Azure native experience for cost visibility and allocation

Built on top of open source, vendor neutral CNCF sandbox project OpenCost

Available for Standard and Premium tier AKS clusters at no additional cost

Ensure costs are allocated to the right teams to drive accountability

Identify high spend areas and opportunities to optimize costs

Proactively identify cost anomalies to prevent unanticipated overspending

**Kubernetes specific views**

*Kubernetes cluster view*

# Demo – Aks Cost Analysis

1. Show on portal.

# Azure Arc Enabled Kubernetes

- Connect any Kubernetes deployment to Azure

- Traffic flows over https – private endpoint is in preview

- Enables below features –
  - Defender
  - Policies
  - AAD
  - Cluster Connect – Connect to cluster without opening inbound ports
  - Kubernetes Partners distribution supported - https://learn.microsoft.com/en-us/azure/azure-arc/kubernetes/validation-program

# AKS Construction Helper

- https://azure.github.io/AKS-Construction/?default=es

- https://github.com/Azure/Aks-Construction

# There is difference in infra scale and App based infra scale

KEDA - 
https://keda.sh/docs/2.8/scalers/

# Azure Linux
## Container Host OS for Azure Kubernetes Service

### Just enough OS

Smaller
OS footprint

### Reliability

Shift left in build pipeline

Prevent defective builds
from advancing

Stringent package tests

Performance tests

### Security

Secure defaults

Fast CVE patching

Secure supply chain

# Fleet Manager

- Why?
  - Managing 20+ AKS clusters in any org is always pain
  - Complexity increases in case of multi cloud, on premises Kubernetes
  - No Central governance exists as of today across all clusters across subscriptions, regions, resource groups etc.
- Join cluster across RG, Sub, Regions
- Member clusters should be under same AAD tenant.
- Support for multi cloud, hybrid AKS clusters
- Selective member configuration supported
- Max 20 member as of today
- Sample Use cases –
  - Upgrade AKS version of all clusters
  - Create same namespace across all clusters. Example, same namespace for ingress
  - Create common RBAC across all clusters
  - Centralised pod to pod communication policies

# Azure Kubernetes Fleet Manager

**Azure AD** **Azure Policy** **Microsoft Defender** **Azure Monitor**

Kubernetes configuration propagation

Networking

Identity & Security AuthN, RBAC

Infrastructure management

Flux GitOps Agent

Fleet capabilities

GitHub

Fleet – managed hub cluster

Join as member clusters

member-a

member-b

member-c

subscription-a

subscription-b
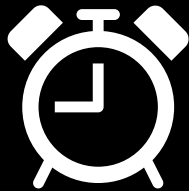
# AKS Kubernetes version Long Term Support (LTS)

Two years of Microsoft support, including CVEs and critical bugs

Upgrade available to the next AKS Kubernetes LTS

Ability to return to the upstream version train

Forked after upstream EOL, maintained by Microsoft in the open

# Azure powers OpenAI and ChatGPT

ChatGPT

**Runs on Azure Kubernetes Service (AKS)**

**Backed by Azure Cosmos DB**

**Developed on GitHub**

**Fastest**
growing consumer
product in history

**Every app** will be reinvented with AI | **New apps** will be built that weren't possible before

# Microsoft Copilot for Azure

An AI companion that simplifies how you design, operate, optimize, and troubleshoot both apps and infrastructure from cloud to edge.

Available initially in Azure portal. Expanding to Azure mobile app and CLI.

https://ms.portal.azure.com/?feature.drilldown=true&feature.apiversion=2023-04-01-preview&feature.aksview=true&feature.canmodifystamps=true&microsoft_azure_costm...

Microsoft Azure (Preview)     Search resources, services, and docs (G+/)     Copilot     MICROSOFT (MICROSOFT.ONMI...

Dashboard > petsupply-1 | Node pools >

# ws75043a084 | Overview
Node pool

Search

- Overview
- Nodes
- Configuration

▷ Start    ☐ Stop    ↑ Upgrade Kubernetes    ↑ Update image    ↗ Scale node pool    🗑 Delete    ↻ Refresh    Give feedback

∧ Essentials

| | | | |
|---|---|---|---|
| Provisioning state ⓘ | : Succeeded | Cluster | : petsupply-1 |
| Power state ⓘ | : Running (1/1 nodes ready) | Operating system | : Ubuntu Linux |
| Availability zones | : None | Kubernetes version | : 1.27.3 |
| Mode | : User | Node count | : 1 node |
| | | Node size | : Standard_NC12s_v3 |

**Properties**     Monitoring

### Node pool

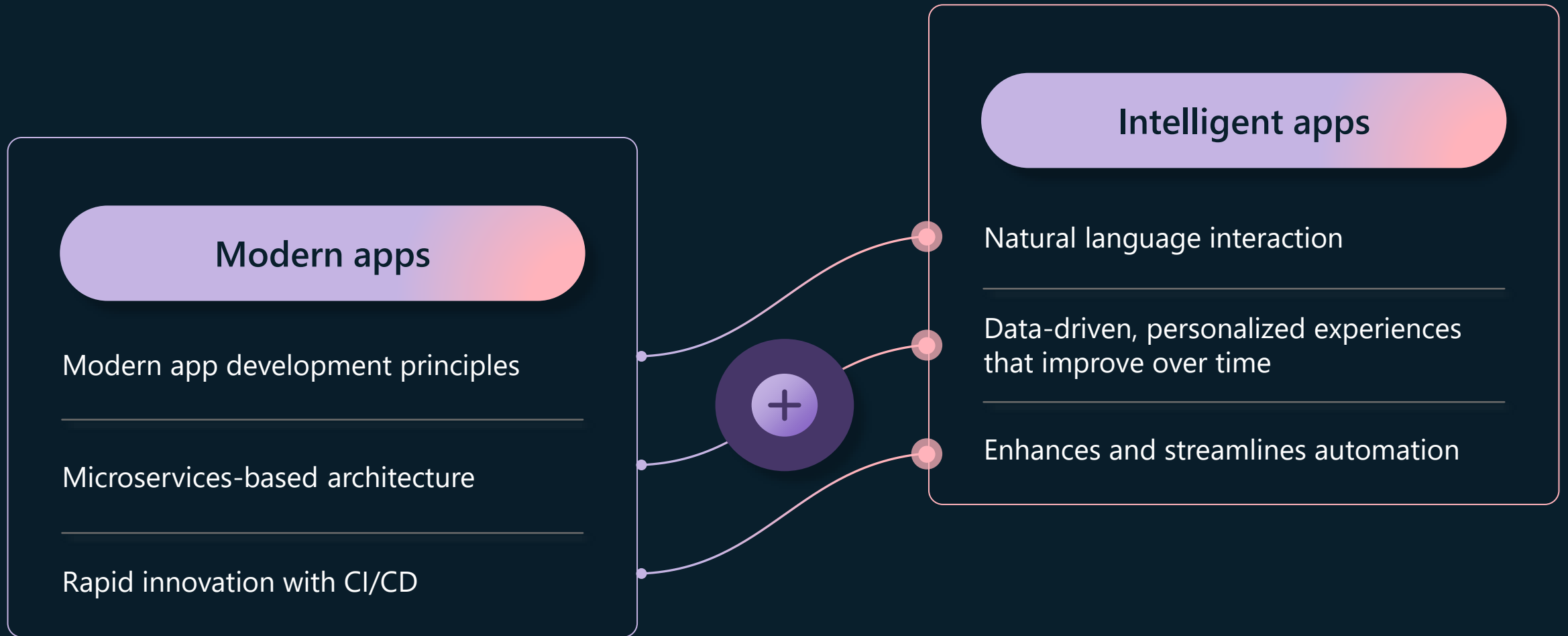| | |
|---|---|
| Max pods per node | 250 |
| Public IPs per node | Disabled |
| Autoscaling | Disabled |
| Azure Spot Instance | Disabled |
| Maximum price | N/A |
| Scale eviction policy | N/A |
| Node image version | AKSUbuntu-2204gen2containerd-202310.31.0 |
| Proximity placement group | N/A |

### Configuration

| | |
|---|---|
| Mode | User |
| Maximum surge | Default |
| Node drain timeout | 30 |
| OS disk size | 128 GB |
| OS disk type | Ephemeral |

### Taints and labels

Taints     sku=gpu:NoSchedule

Labels

apps : falcon-7b    kaito.sh/machine-type : gpu

kaito.sh/workspace : workspace-falcon-7b    kaito.sh/workspacenamespace : store

karpenter.sh/provisioner-name : default
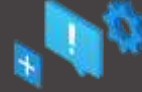
# Generative AI makes apps truly intelligent

**Modern apps**

Modern app development principles

Microservices-based architecture

Rapid innovation with CI/CD

**Intelligent apps**

Natural language interaction

Data-driven, personalized experiences that improve over time

Enhances and streamlines automation

# AI toolchain operator add-on for AKS

**Announcing**

Deploy

Inference

Innovate

**AI toolchain operator add-on for AKS**

## Workspace and infra setup to model inferencing
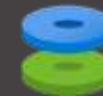
### in a matter of minutes

Load model weights

Model Containerization

Host image

Provision GPU infra

# Thank you...

Stay connected...

http://sanganakauthority.blogspot.com/

Twitter - @KunalChandratre

Active on LinkedIn!!