

Operating Azure Kubernetes Service

Learnings from the field

Daniel Neumann

Staff Software Engineer for LeanIX

Microsoft MVP – Microsoft Azure

<https://www.danielstechblog.io>

@neumannndaniel

Session objectives

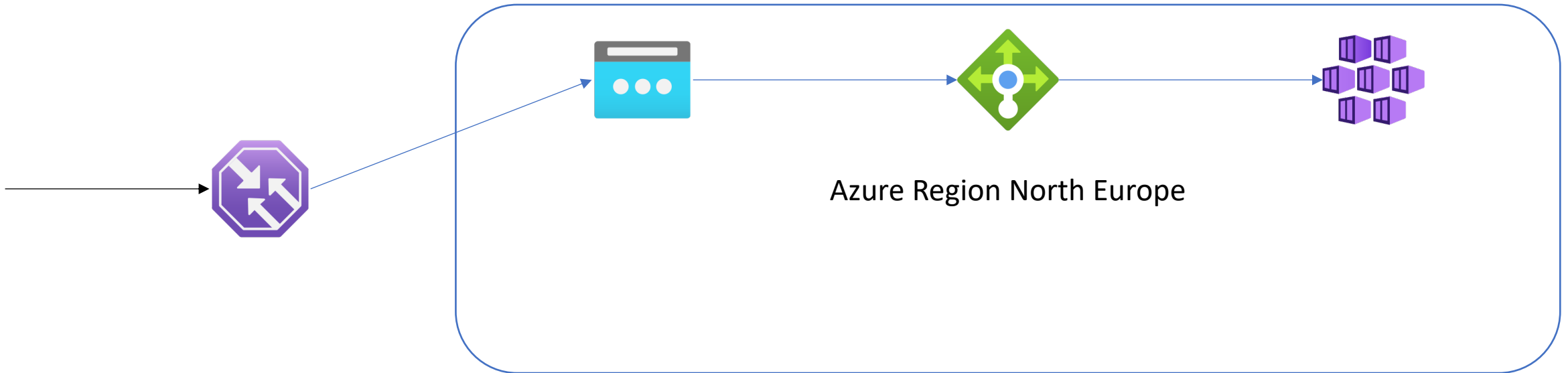
- Business Continuity and Disaster Recovery
- High availability
- Patch and upgrade management
- Governance, monitoring, and alerting

Business Continuity and Disaster Recovery

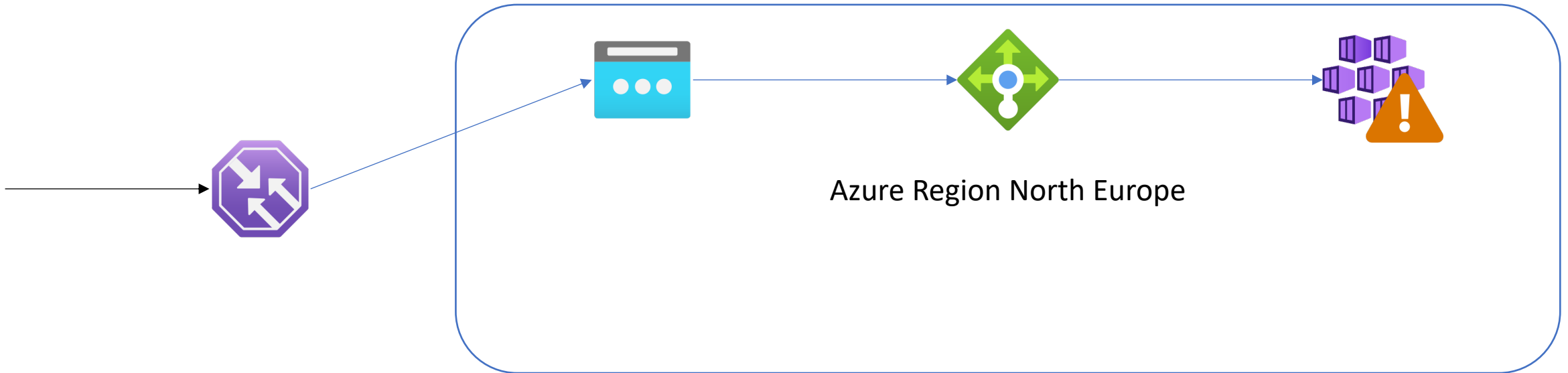
Business Continuity and Disaster Recovery

- Have a DR strategy in place
- Use Traffic Manager as DNS load balancer in front of AKS
- Stateless services – Use PaaS services for storing state
 - Azure Storage
 - Azure Database for PostgreSQL
 - Azure Cache for Redis
 - etc.

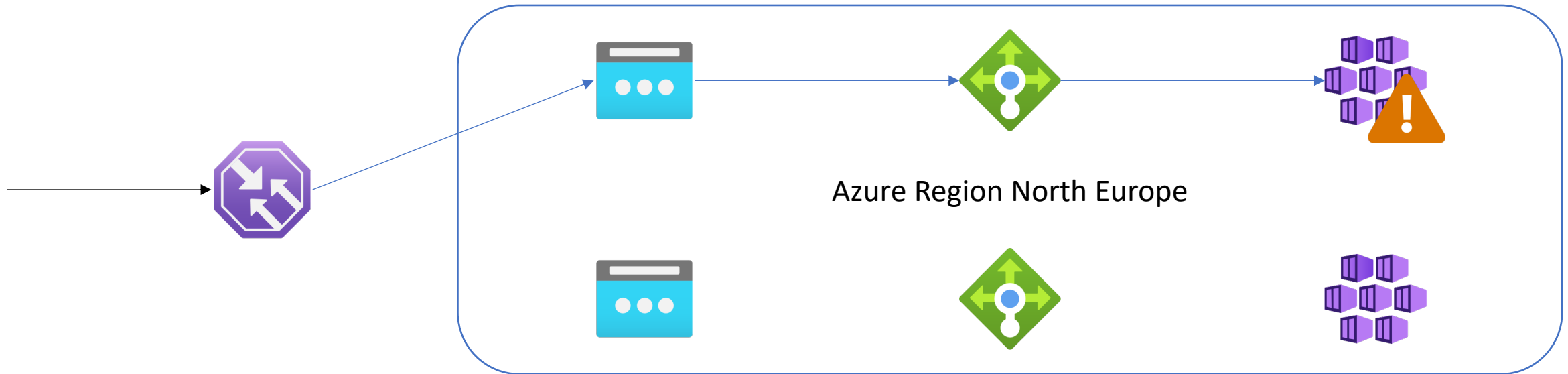
BCDR – Traffic Manager Architecture



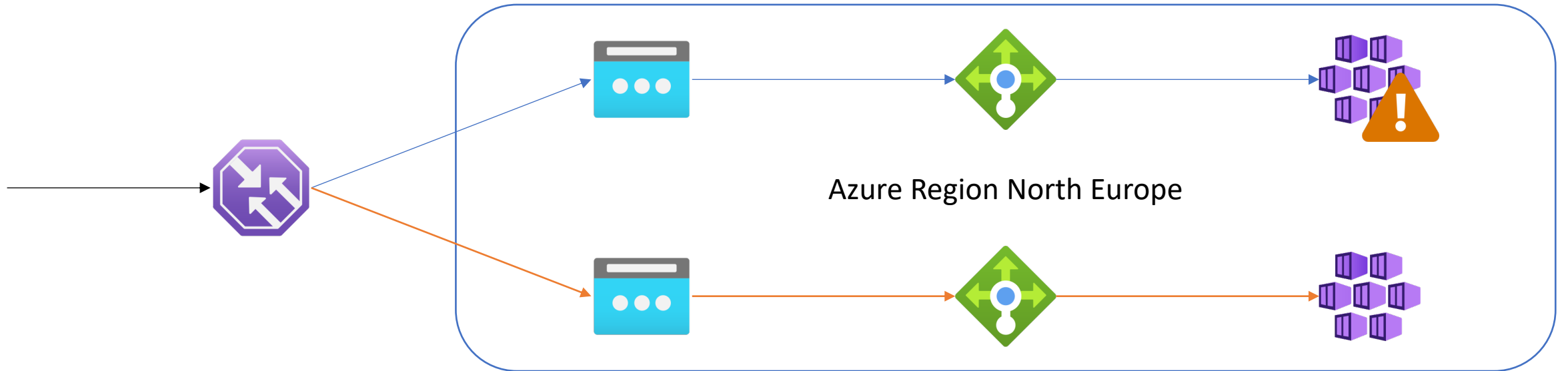
BCDR – Traffic Manager Architecture



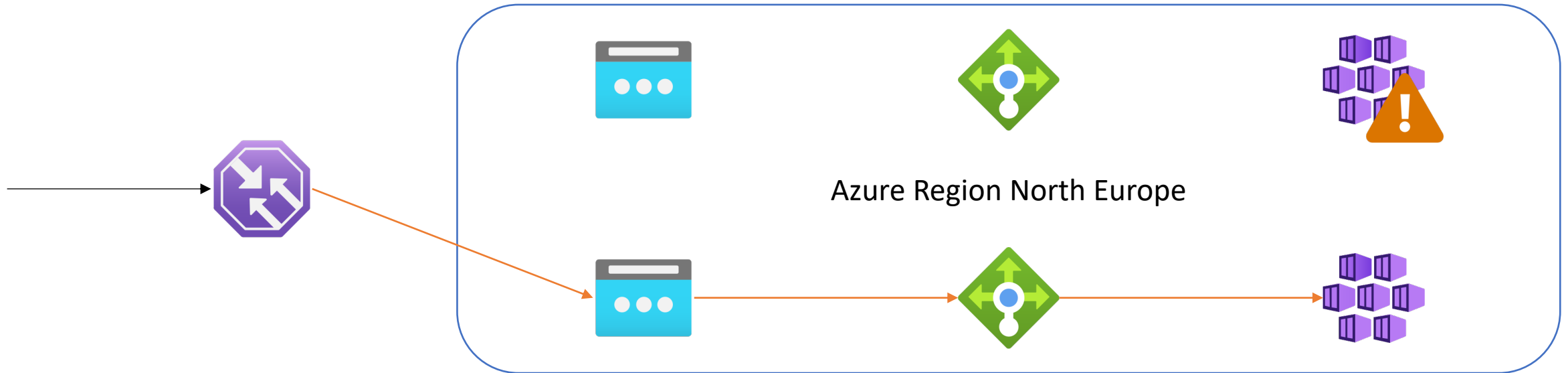
BCDR – Traffic Manager Architecture



BCDR – Traffic Manager Architecture



BCDR – Traffic Manager Architecture



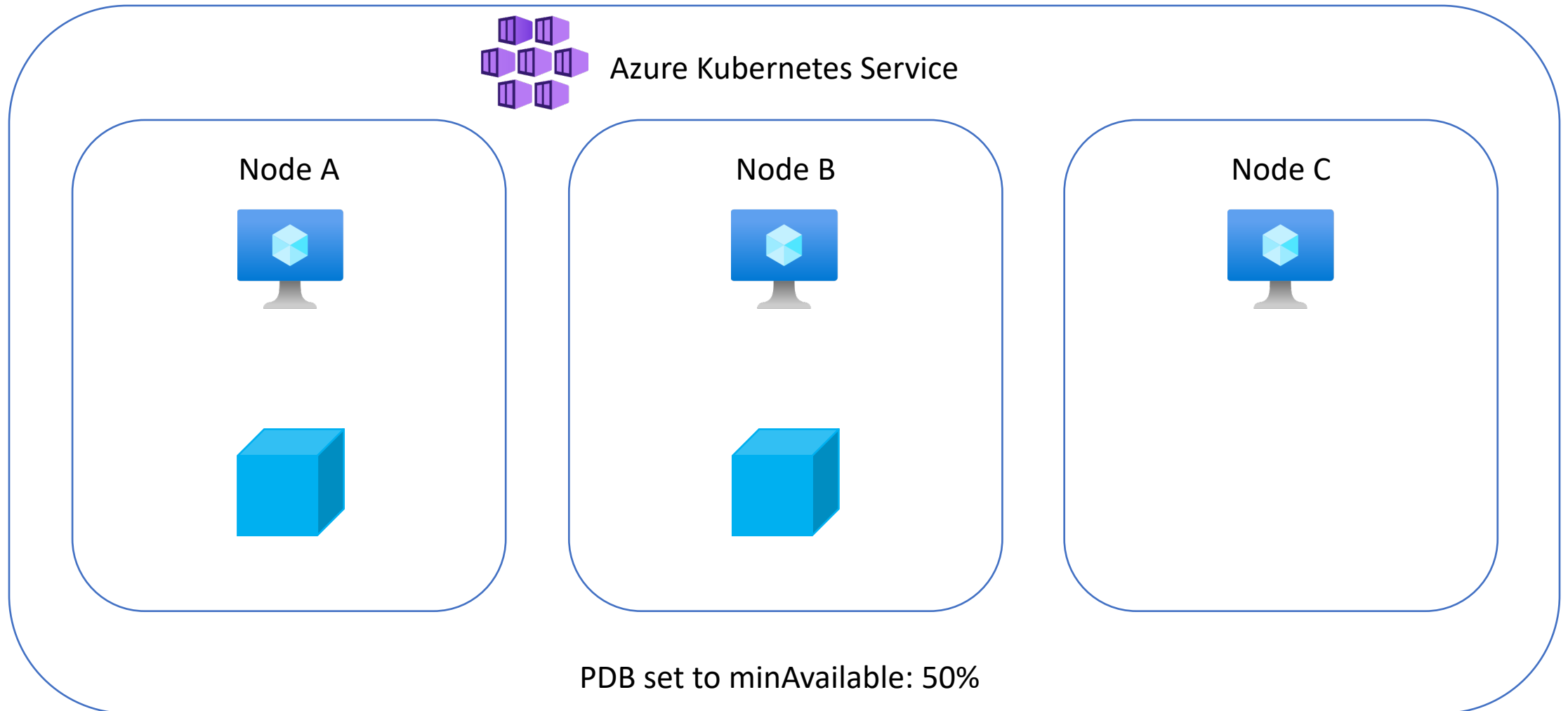
Demo – Traffic Manager

High availability

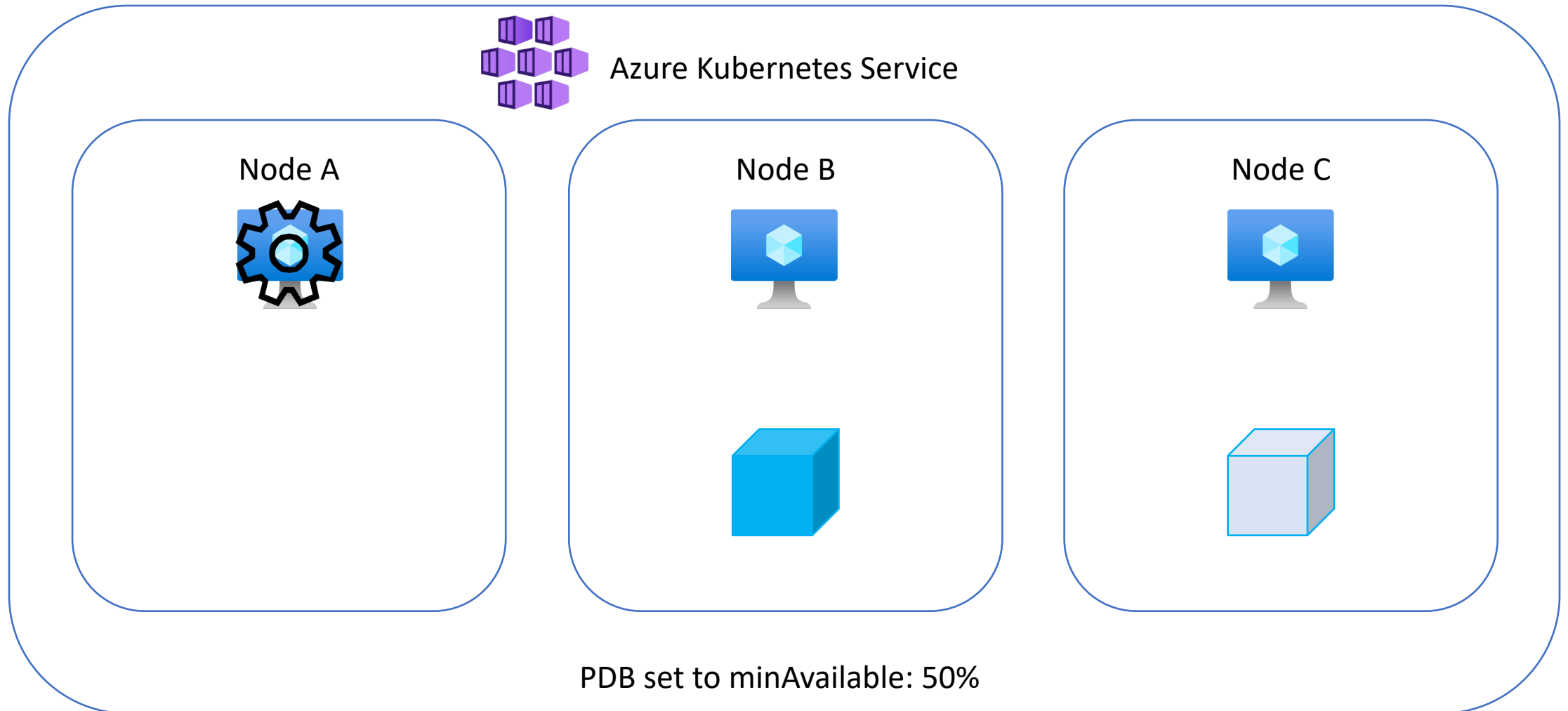
High availability

- Use Pod Disruption Budgets for your applications
 - Voluntary eviction
- Use Pod Anti Affinity settings to distribute the application across different nodes / availability zones
 - Nonvoluntary eviction
- Use Pod Topology Spread Constraints to distribute the application across different availability zones
 - Kubernetes 1.19 or higher
 - Nonvoluntary eviction

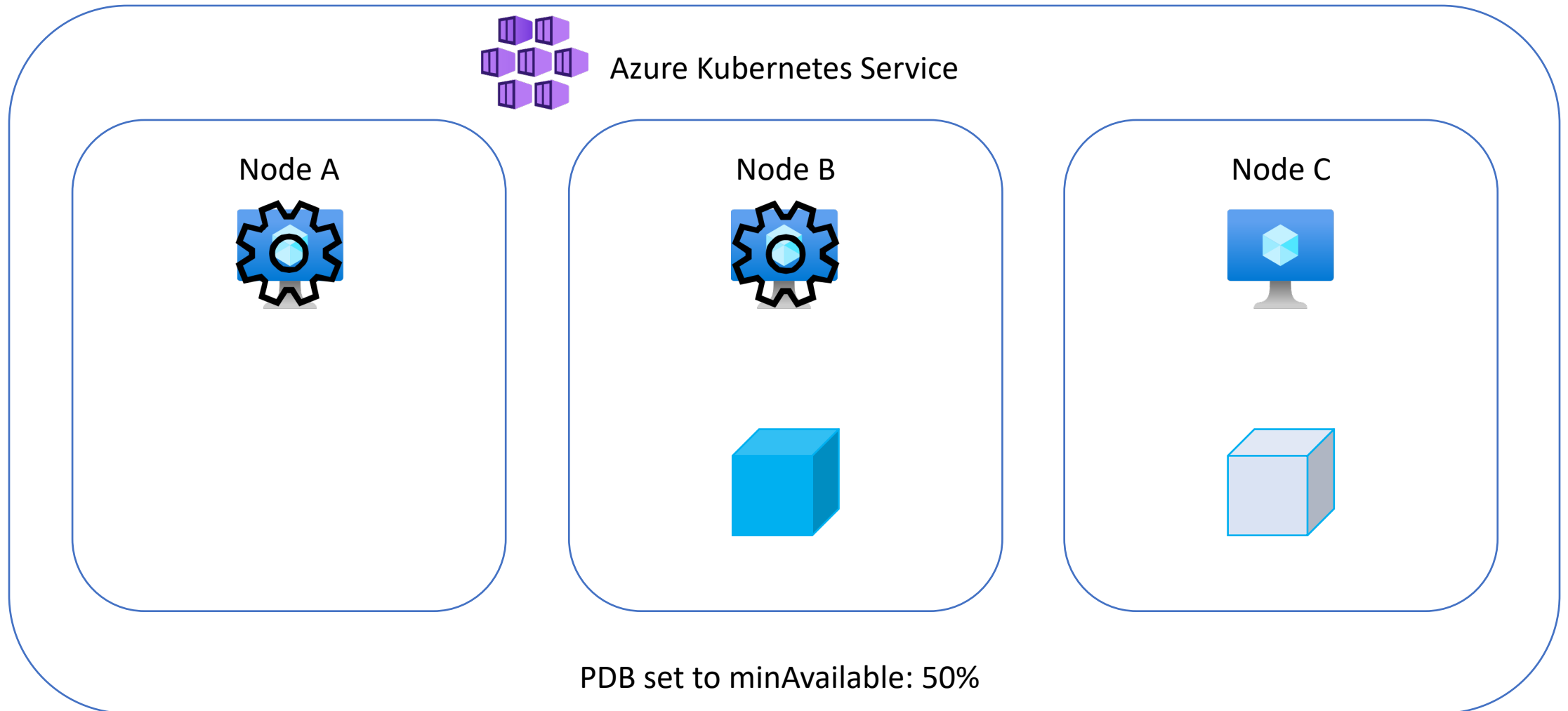
HA – Pod Disruption Budgets



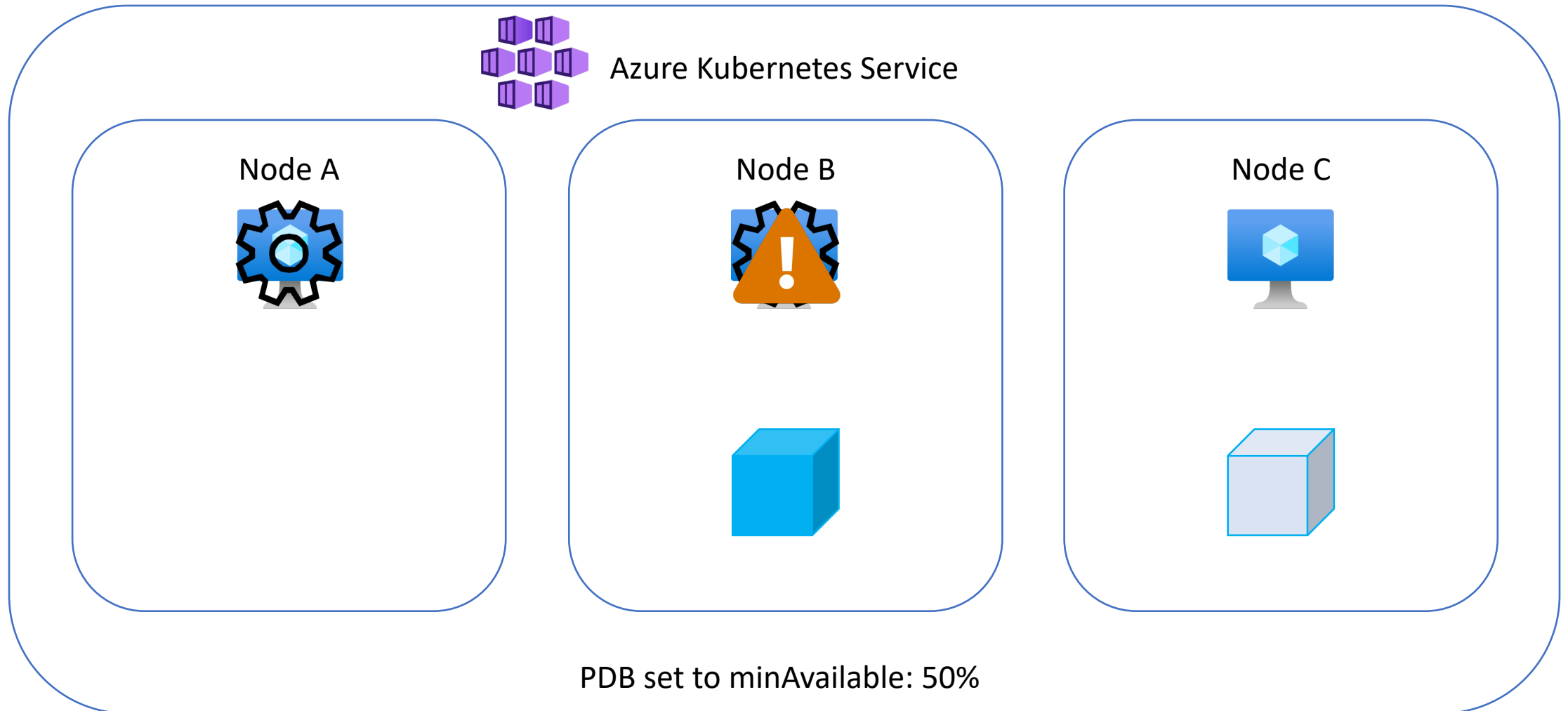
HA – Pod Disruption Budgets



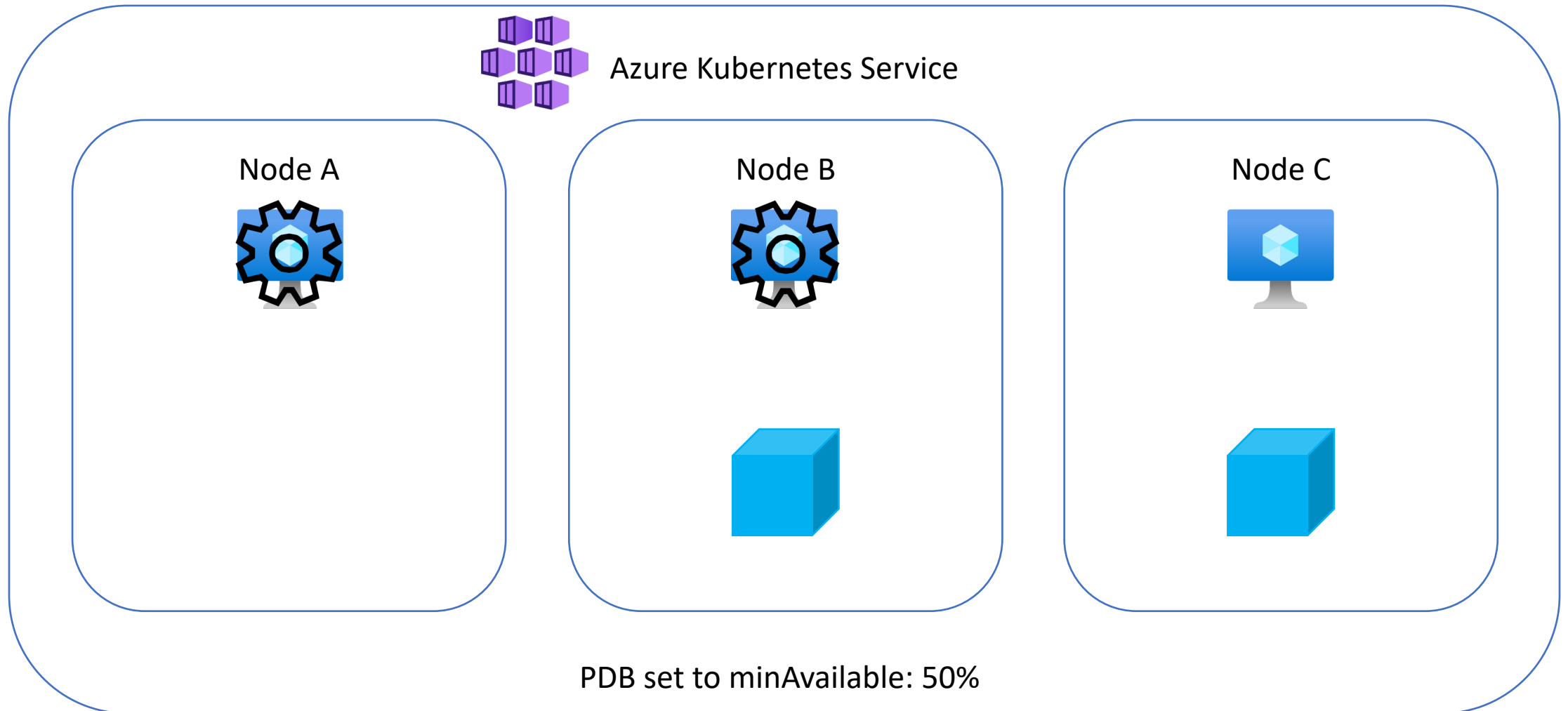
HA – Pod Disruption Budgets



HA – Pod Disruption Budgets



HA – Pod Disruption Budgets



HA – Pod Topology Spread Constraints



Azure Kubernetes Service

Node A – Zone A



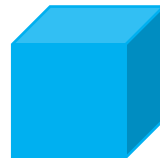
A-B: $2-1=1$
A-C: $2-0=2$



Node B – Zone B



B-A: $2-1=1$
B-C: $2-0=2$



Node C – Zone C

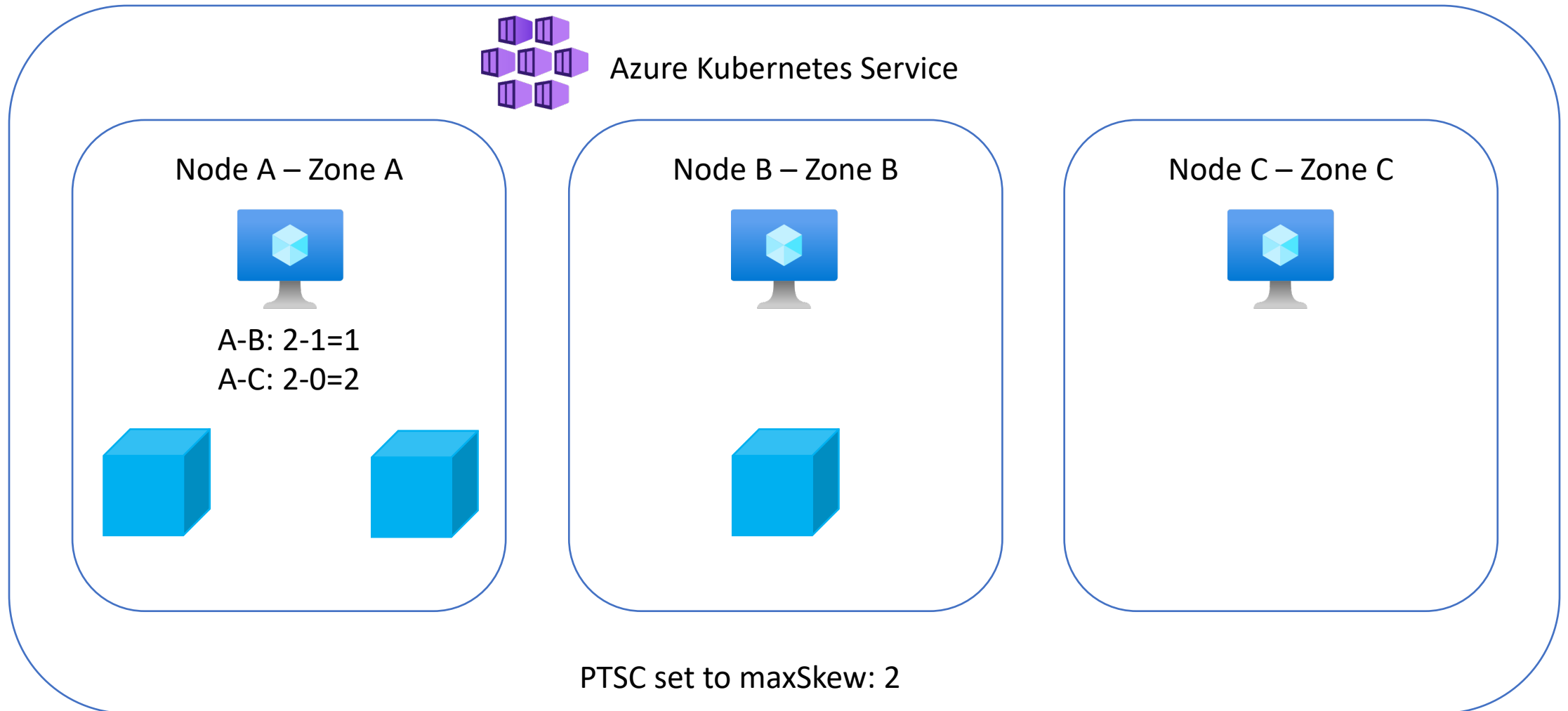


C-A: $0-1=-1$
C-B: $0-1=-1$

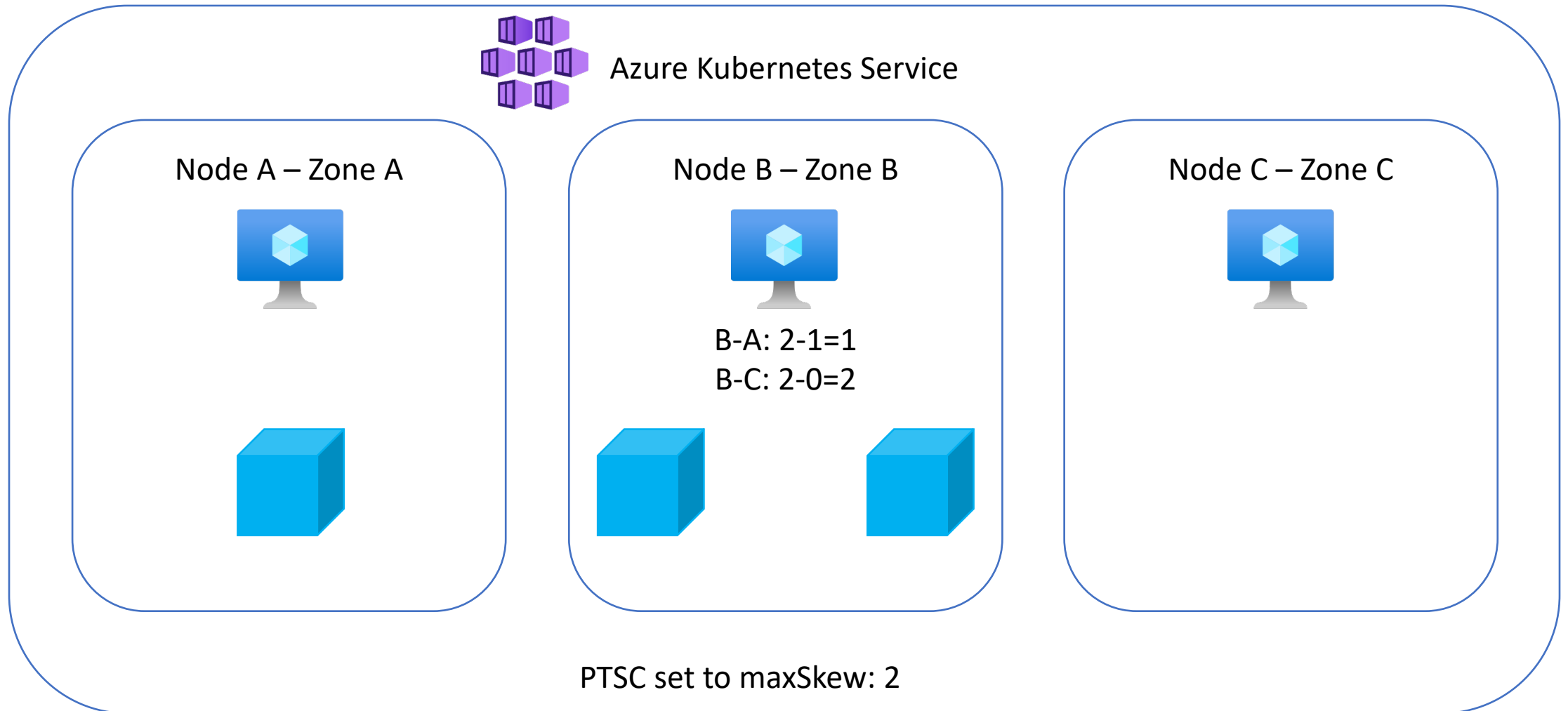


PTSC set to maxSkew: 1

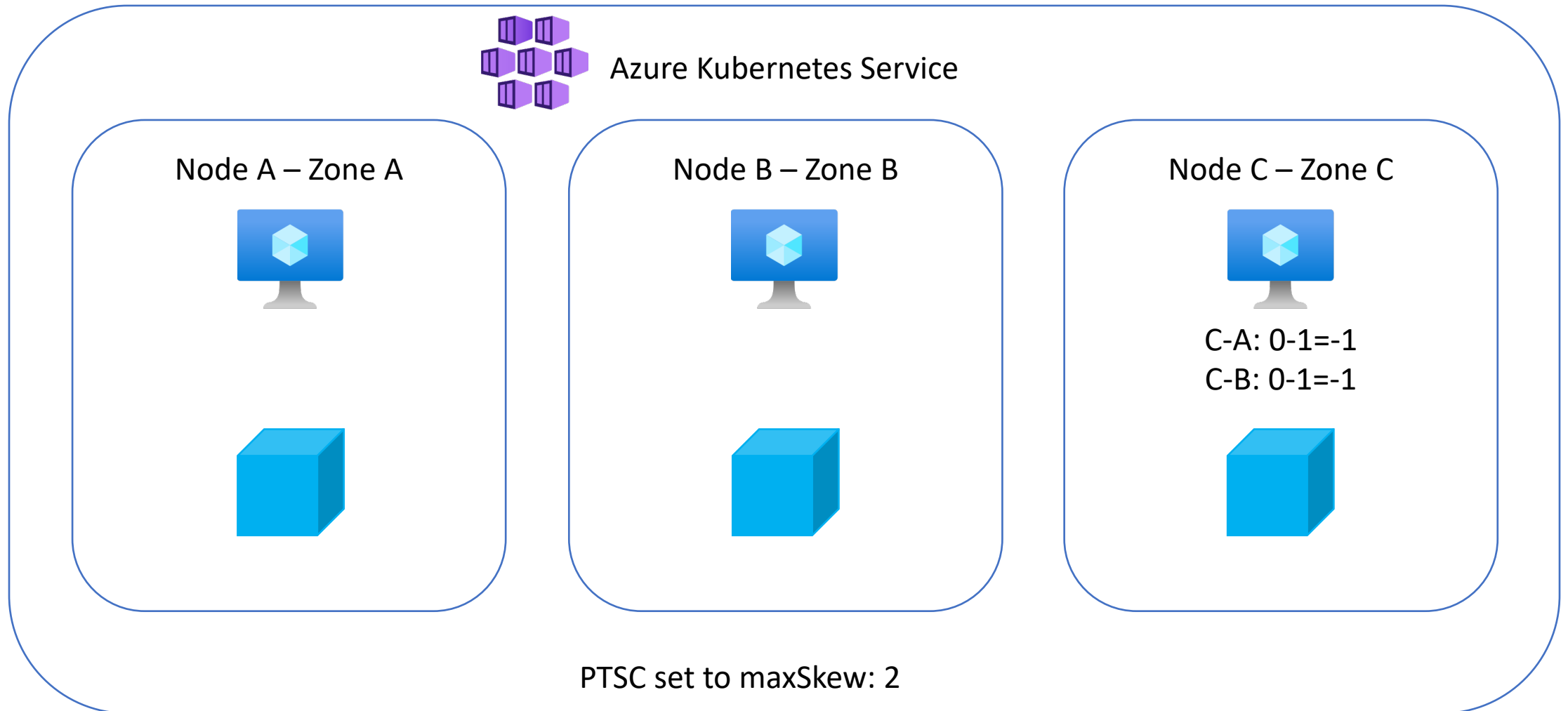
HA – Pod Topology Spread Constraints



HA – Pod Topology Spread Constraints



HA – Pod Topology Spread Constraints



Cluster autoscaling

- Use cluster autoscaler
 - Let you easily replace unhealthy nodes that have not been detected by the node auto-repair functionality
- Check your core quota
- AKS Uptime SLA
 - 99,95% SLA with AZ (Availability Zones)
 - 99,9% SLA
 - 99,5% **SLO** with free tier

Application autoscaling

- Use KEDA for advanced pod scaling scenarios
 - Standard HPA (Horizontal Pod Autoscaler) in Kubernetes only supports CPU and memory out of the box



Patch and upgrade management

Keep Azure Kubernetes Service up-to-date

- kured – Kubernetes Reboot Daemon
 - Apply security patches requiring a reboot
- Node OS image upgrade
 - New images released every 1-3 weeks
- Automatic Kubernetes version upgrade
 - Protects you from using an unsupported version



Governance, monitoring, and
alerting

Governance with Azure Policy

- Azure Policy for Kubernetes (Gatekeeper)
 - Enforce governance settings -> Labels
 - Enforce security settings -> No privileged containers
 - Enforce compliance settings -> Only allowed images
 - Mitigating CVEs -> CVE-2020-8554



Monitoring and alerting

- AKS – Diagnostic checks
- Azure Monitor – Container insights
 - Recommended alerts
 - Log query alerts
- Enable AKS control plane logs
 - Storage account recommended
- Kubernetes resources in the Azure portal

Demo

SNAT port exhaustion

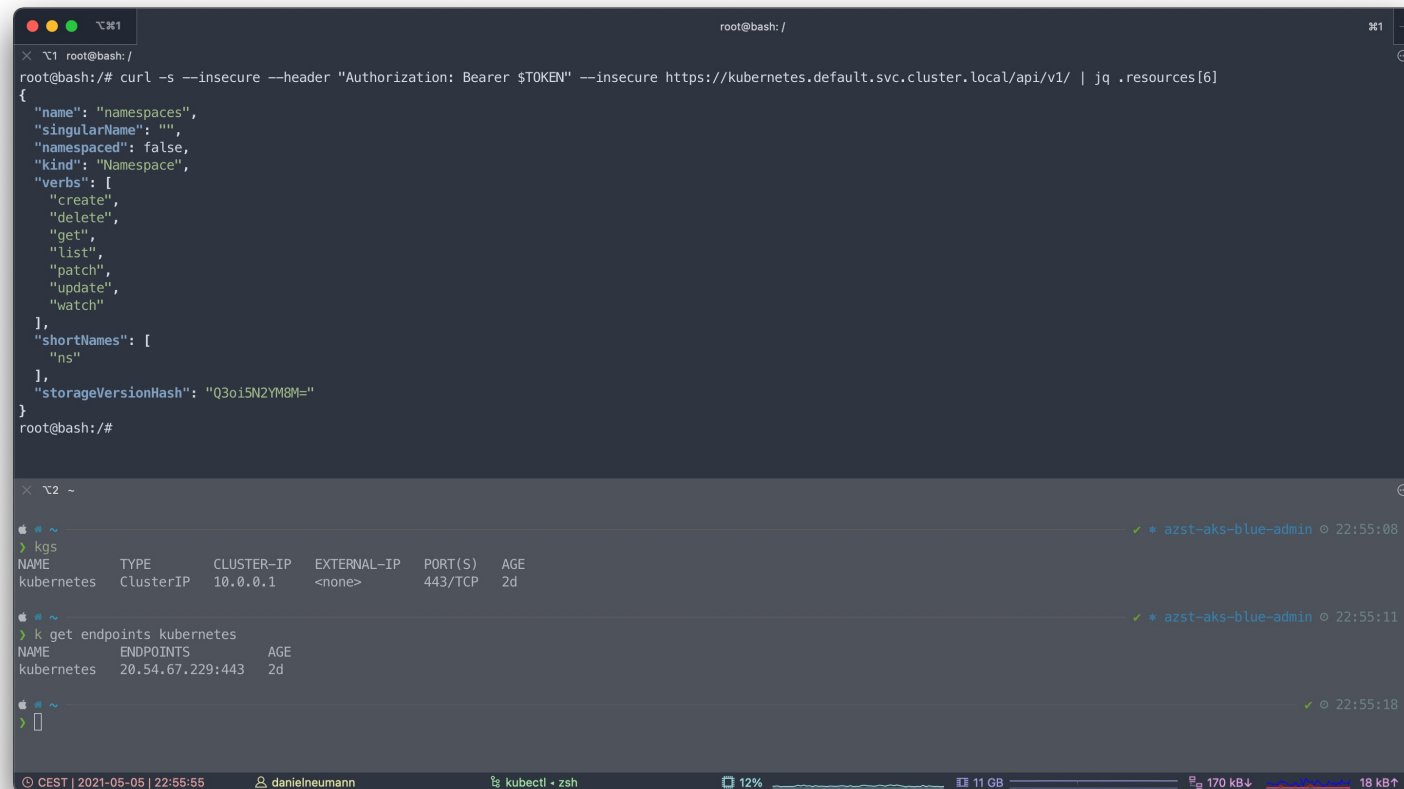
Governance, monitoring, and alerting

SNAT connections – What counts in?

- Access to Azure PaaS without Private Link
 - Azure Storage
 - Azure Database for PostgreSQL
 - Azure Cache for Redis
 - etc.
- Access to external services

SNAT connections – What counts in?

- AKS API Server access from within the cluster
 - `kubernetes.default.svc.cluster.local`



The image shows a terminal window with two tabs. The first tab, titled 'root@bash: /', contains a `curl` command that successfully retrieves the Kubernetes namespace configuration for 'kubernetes.default.svc.cluster.local'. The output is a JSON object describing the namespace. The second tab, titled '~', shows a series of `kubectl` commands and their outputs. First, `kubectl get kubernetes` is run, resulting in a table with one row for the 'kubernetes' ClusterIP. Then, `kubectl get endpoints kubernetes` is run, resulting in a table with one row for the 'kubernetes' endpoints. The terminal status bar at the bottom indicates the user is 'danielneumann' using 'kubectl - zsh' on '2021-05-05' at '22:55:55', with system metrics like 12% CPU, 11 GB memory, and network activity.

```
root@bash: /
root@bash:/# curl -s --insecure --header "Authorization: Bearer $TOKEN" --insecure https://kubernetes.default.svc.cluster.local/api/v1/ | jq .resources[6]
{
  "name": "namespaces",
  "singularName": "",
  "namespaced": false,
  "kind": "Namespace",
  "verbs": [
    "create",
    "delete",
    "get",
    "list",
    "patch",
    "update",
    "watch"
  ],
  "shortNames": [
    "ns"
  ],
  "storageVersionHash": "Q3oi5N2YM8M="
}
root@bash:/#

~
> kgs
NAME          TYPE          CLUSTER-IP   EXTERNAL-IP   PORT(S)    AGE
kubernetes    ClusterIP     10.0.0.1     <none>        443/TCP    2d

> k get endpoints kubernetes
NAME          ENDPOINTS          AGE
kubernetes    20.54.67.229:443   2d

>

```

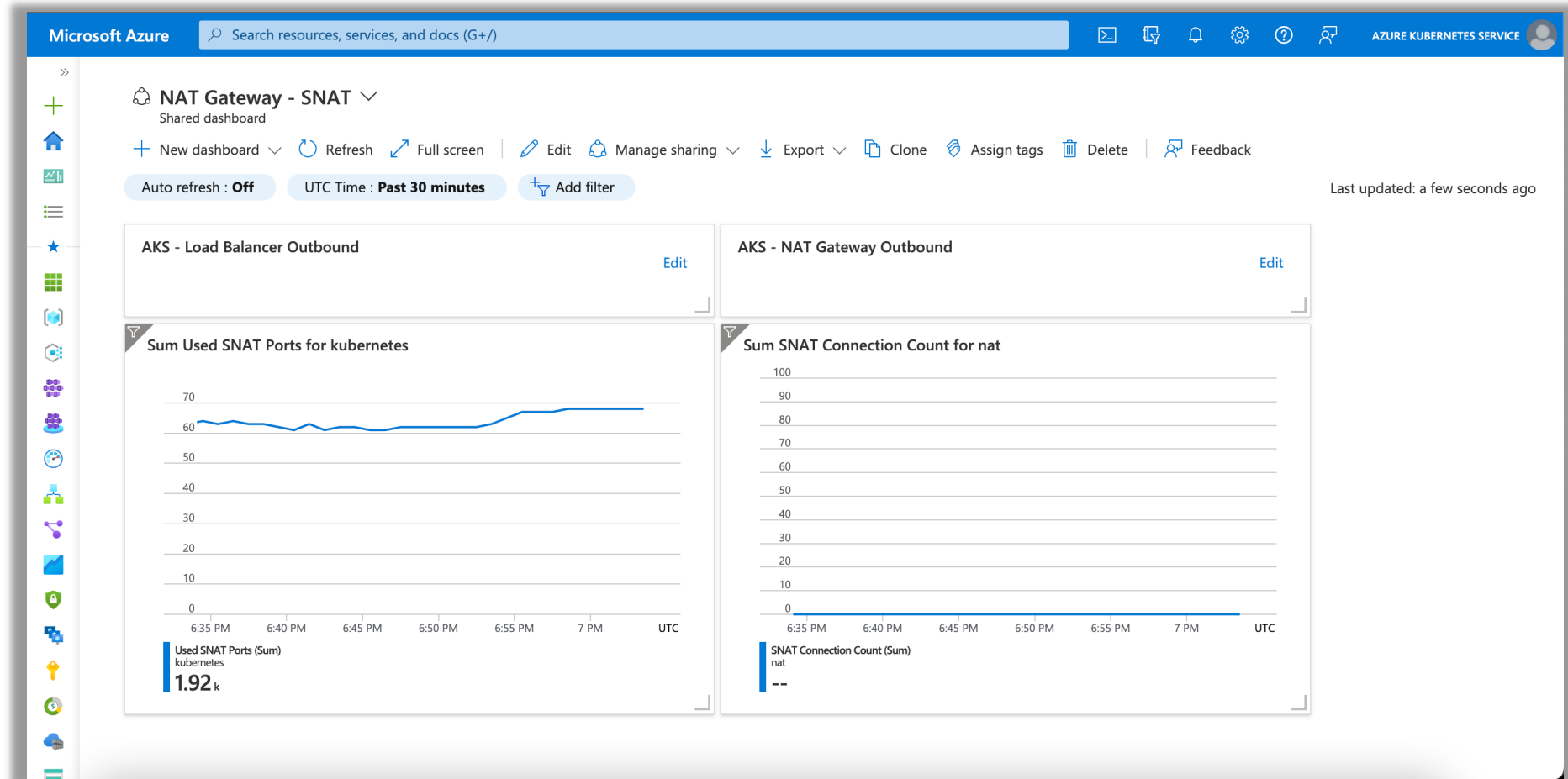

Mitigation

1. Assign enough public IPs to the load balancer, set a custom value for the allocated SNAT ports per node, and set TCP idle reset to 4 minutes
 - Automatic default depends on the cluster size and starts with 1024 SNAT ports and ends at 32 SNAT ports per node
 - Still the risk running into a SNAT port exhaustion
2. Use Azure Virtual Network NAT
 - Do not use *outboundType* managedNATGateway or userAssignedNATGateway in the AKS configuration

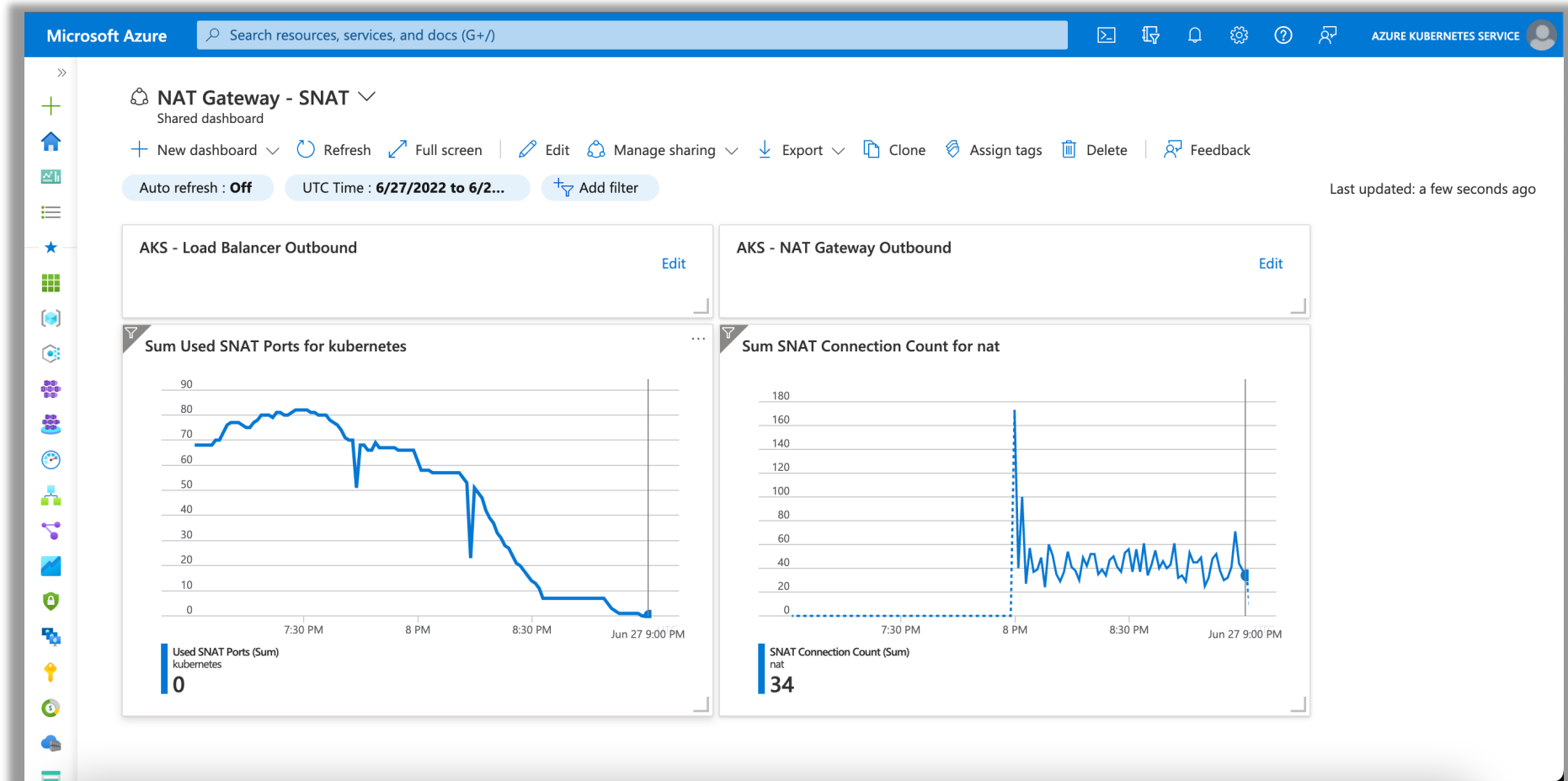
Why should I stick to *outboundType* loadBalancer?

1. Azure Virtual Network NAT takes precedence
 - *“A NAT gateway takes precedence over other outbound scenarios and replaces the default Internet destination of a subnet.”*
 - All outbound traffic will go through the Virtual Network NAT and not through the load balancer via its outbound rules.
2. When using managedNATGateway or userAssignedNATGateway you cannot recover yourself from a Virtual Network NAT outage without redeploying the AKS cluster
 - Using the *outboundType* loadBalancer lets you disassociate the Virtual Network NAT from the subnet and AKS will leverage the outbound rules from the load balancer for outbound traffic

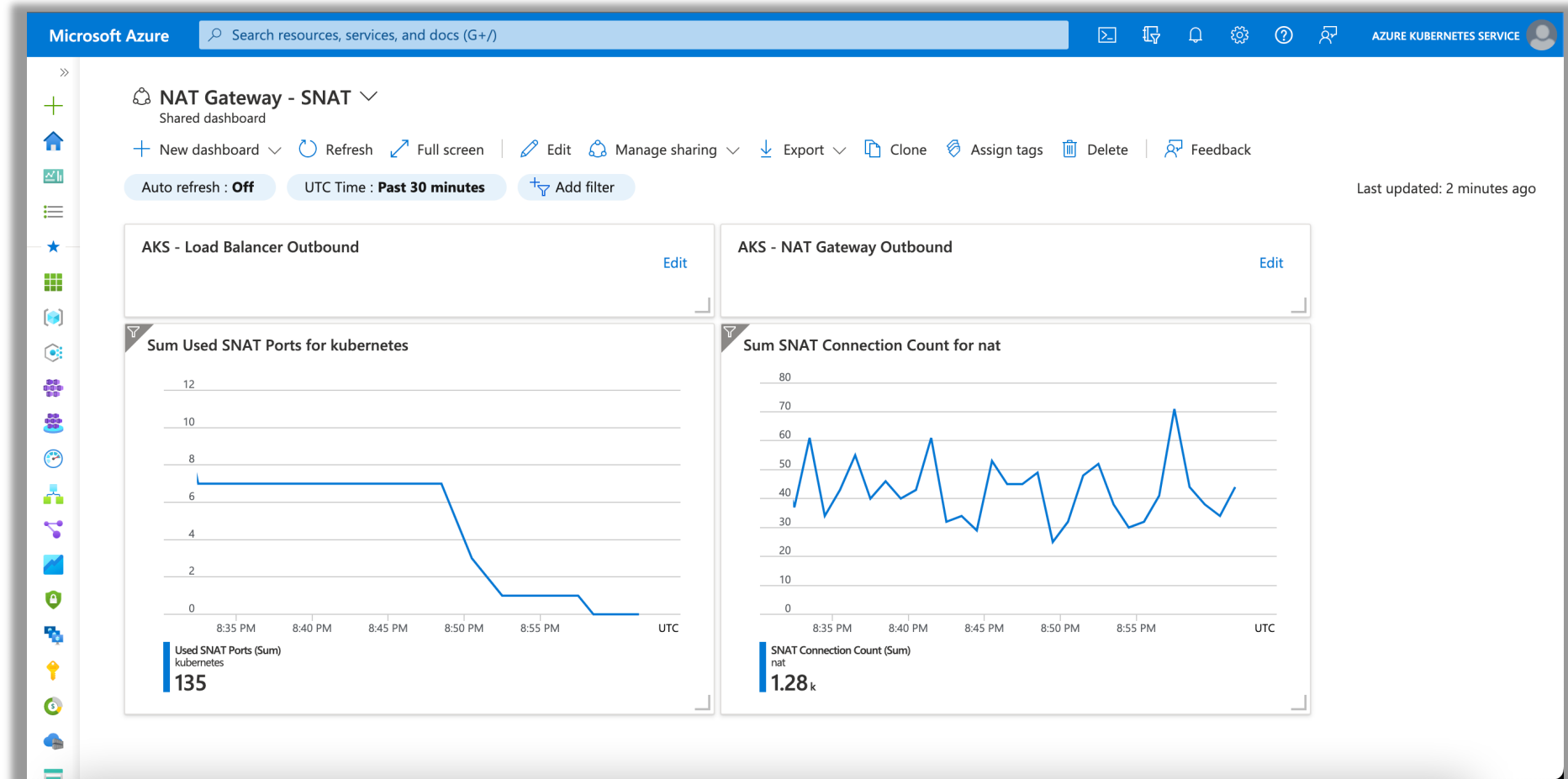
Why should I stick to *outboundType* loadBalancer?



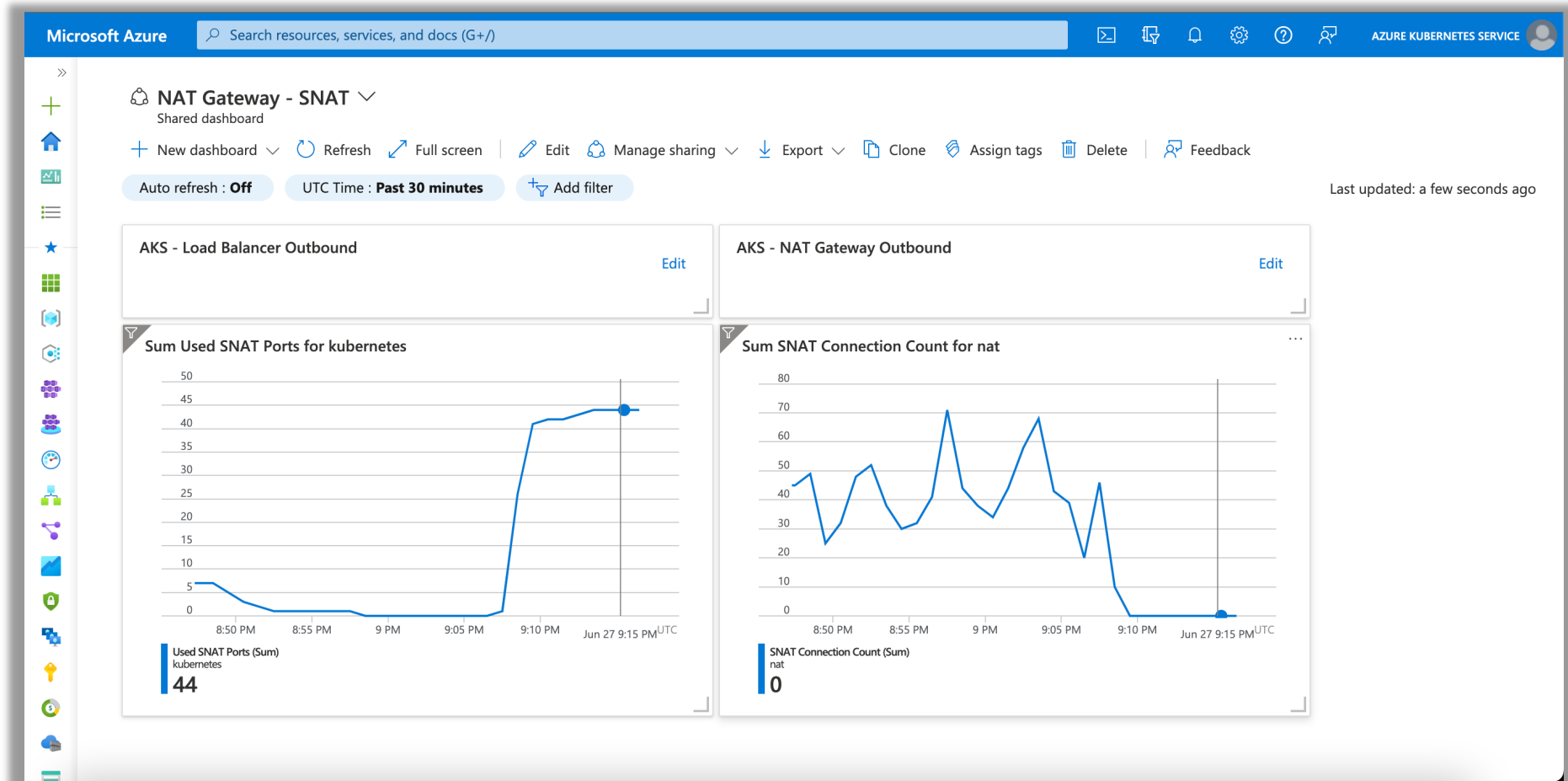
Why should I stick to *outboundType* loadBalancer?



Why should I stick to *outboundType* loadBalancer?



Why should I stick to *outboundType* loadBalancer?



Summary

- Have a DR plan in place and use PaaS services
- Educate your developers/engineers to use Kubernetes principles
- Keep AKS up-to-date
- Implement proper governance, monitoring, and alerting

Thank you!

Appendix

- Pod Disruption Budget
 - <https://kubernetes.io/docs/concepts/workloads/pods/disruptions/>
 - <https://www.danielstechblog.io/increase-your-application-availability-with-a-poddisruptionbudget-on-azure-kubernetes-service/>
- Pod Topology Spread Constraints
 - <https://kubernetes.io/docs/concepts/workloads/pods/pod-topology-spread-constraints/>
 - <https://www.danielstechblog.io/distribute-your-application-across-different-availability-zones-in-aks-using-pod-topology-spread-constraints/>
- Kured
 - <https://docs.microsoft.com/en-us/azure/aks/node-updates-kured>
- Node image upgrade
 - <https://docs.microsoft.com/en-us/azure/aks/node-image-upgrade>

Appendix

- Auto-upgrade channel
 - <https://docs.microsoft.com/en-us/azure/aks/upgrade-cluster#set-auto-upgrade-channel>
- Recommended Alerts
 - <https://docs.microsoft.com/en-us/azure/azure-monitor/containers/container-insights-metric-alerts>
- SNAT Port Exhaustion
 - <https://docs.microsoft.com/en-us/azure/load-balancer/load-balancer-outbound-connections>
 - <https://docs.microsoft.com/en-us/azure/load-balancer/load-balancer-standard-diagnostics#how-do-i-check-my-snat-port-usage-and-allocation>
 - <https://www.danielstechblog.io/detecting-snat-port-exhaustion-on-azure-kubernetes-service/>
- Uptime SLA
 - <https://docs.microsoft.com/en-us/azure/aks/uptime-sla>