# ADS Tuning Guide

—

Nicolas Peulvast
Performance Architect

IBM

# Disclaimer Official

# Contents

# Executive Summary

## Authentication

The Authentication mode chosen has an impact on the throughput of the ADS Runtime Pod and is managed by the Common Services Layer: be sure to have an adapted T-Shirt size of the Common Services Layer to optimize the Authentication performance.

## Network

The Network layer may have an important impact of the overall performance and especially on the throughput of the ADS Runtime Pod in case of multiple users in parallel : be sure to have an adapted T-Shirt size of the Zen Layer to optimize the network performance.

## BAI

Activating the BAI event emission has an overhead, that is dependent of the additional collected information.

You must increase the resource used for the system to get back to the expected throughput.

Note that you must tune the BAI stack to have a similar average event throughput ingestion from BAI in order to not overload your system.

## Tuning

A fine-tuning of the different layers is the key to a good throughput on the ADS Runtime, and especially the network layers.

See the following slides for the tuning guidelines.

# General Starter Tuning Guidelines

- In order to propose a Demo/Minimal sizing, the following configuration is suggested for the ADS Runtime service:

  - CPU Request == CPU Limit == 0.5

  - While startup time is longer, it does not reach probe limits

  - With this configuration, we still have a good level of performance at the ADS Runtime level

```yaml
decision_runtime_service:
  autoscaling:
    enabled: false
  replica_count: 1
  resources:
    requests:
      cpu: '500m'
      memory: '2Gi'
    limits:
      cpu: '500m'
      memory: '3Gi'
```

# General Prod Tuning Guidelines

- The network latency between the Cluster and the database has a huge impact on the Designer performance as the Git Service performs a lot of small requests on each repository access.

- Ephemeral Storage & impact on stability: if the Ephemeral storage is configured too small, then the ADS Pods will be evicted.

- A Horizontal Pod Autoscaler (HPA) is available out-of-the-box in the ADS Runtime delivery

  - Using this HPA increase the throughput of the ADS runtime but also increase the number of VPC billed to the customer

  - By default, this HPA is not set for the Small, Medium and Large T-Shirt size but only for X-Large T-Shirt size

# Tuning Foundational service

Align the profile (small, medium, large) of your
Common service/Foundational Service to the
targeted profile of your ADS product.

- In the OpenShift console:

  - Search (ibm-common-services ns) >
    Resources 'CommonService' > common-
    service

  - In the YAML, change the spec.size value

    - starterset

    - small

    - medium

    - large

# Tuning Shared Configuration

Align the profile (small, medium, large) of your Shared configuration to the targeted profile of your ADS product: it's particularly useful when the BAI events are used in ADS as it tunes the InsightEngine.

- Change it using the Shared Configuration

- In the OpenShift console:

    - Search (your ns) > Resources 'ICP4ACluster' > select your CR deployment

    - In the YAML, change the spec.shared_configuration.sc_deployment_profile_size value

        - small

        - medium

        - large

# Tuning IAF Layer

- This performance tuning is only applicable if you select the BAI event emitter

- You can add additional JVM parameter in the AutomationBase configuration

- If you do that, you must be sure to set shared_configuration.sc_install_automation_base to false

# Tuning Zen Layer

- You can tune the NGNIX layer

  - Edit the ConfigMap `<crname>-ads-designer-zen-configuration` & `<crname>-ads-runtime-zen-configuration`

  - Tune the `nginx.conf` part of the `ConfigMap`

- You can verify the zen resources via the `scaleConfig` parameter in Zen Service where the default profile is `small`, and we also support `medium` / `large` / `xlarge`

  - The Zen layer is adapted to the T-Shirt size that you selected in your Custom Resource

  - T-shirt size `medium` - 3 replicas , cpu limit 800M – Throughput of ADS Runtime Pod up to 11000 TPS

  - T-shirt size `large/xlarge` - 5 replicas, cpu limit 2 – Throughput of ADS Runtime Pod up to 11000 TPS

  - We recommend using the large configuration to avoid bottleneck in the Zen Layer

  - We can manually force the size of your Zen layer using `oc patch AutomationUIConfig iaf-system --type=merge -p '{"spec":{"zenService":{"scaleConfig":"large"}}}'`
    But note the operator will change it back to your Custom Resource value after one roundtrip of the Operator (usually 20 minutes)

- Put the `roundrobin` annotation in you `CPD` route

  - `haproxy.router.openshift.io/balance=roundrobin`

  - Since 21.0.3, This annotation is always override by the source value as some element behind the reverse-proxy have bug in using the roundrobin algorithm.

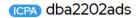  - If you want to switch to the `roundrobin` algorithm, you have to de-activate the CP4A operator.

# Tuning ADS configuration

Adapt profile (small, medium, large) of your ADS configuration to the targeted profile of your ADS product.

- Change it using the ADS Configuration

- In the OpenShift console:

  - Search (your ns) > Resources 'ICP4ACluster' > select your CR deployment

  - In the YAML, change the spec.ads_configuration.deployment_profile_size value

    - small

    - medium

    - large

    - extra-large

# Tuning Frontend Layer

- Check the number of thread allocated to your frontend/HAProxy and verify that the configuration is using the roundrobin algorithm

- As an example, edit the `/etc/haproxy/haproxy.cfg` file and change

  - the `ingress-https` backend from "`balance source`" to "`balance roundrobin`"

  - the `nbproc` should be at least set to 5

# Tuning Network Layer

- During our testing, we reached the network bandwidth capacity (1Gbps) hence the response curve flattens so switch to 10Gbps network should be considered.

- You can also tune the network layer as follow

  - Search [in ns `openshift-ingress-operator`] > Resources '`IngressController`' > select default

    - Edit the yaml and change the `spec.replicas` to at least 5

  - Search [in your cp4ba ns] > Deployment '`ibm-nginx`' > scale up the number of pod in order to reach 100% of CPU usage in your ODM Runtime

# Tuning HPAs

- HPAs are automatically created with the extra-large sizing configuration

- You have 3 additional HPA created in your namespace in that case



- You can change your HPAs using the following Custom Resource customization:

```yaml
spec:
  ads_configuration:
    decision_runtime_service:
      autoscaling:
        enabled: true
        max_replicas: 2
        min_replicas: 5
    parsing_service:
      autoscaling:
        enabled: true
        max_replicas: 2
        min_replicas: 5
    run_service:
      autoscaling:
        enabled: true
        max_replicas: 2
        min_replicas: 5
```

# Storage consideration

Official CloudPak documentation

https://www.ibm.com/docs/en/cloud-paks/cp-biz-automation/22.0.2?topic=deployment-storage-considerations

Align the profile (small, medium, large) of your Foundational service to the targeted profile of you ADS product.

# Tip for CR updating

- In order to be sure that a CR change is considered as soon as possible, you can delete the pod `ibm-cp4a-operator-XX-XX` in order to shutdown the current reconciliation loop (running on the old CR).

- It will result with a new resolution loop that will consider your new values.

# Additional resources

- https://access.redhat.com/documentation/en-us/openshift_container_platform/4.10/pdf/scalability_and_performance/openshift_container_platform-4.10-scalability_and_performance-en-us.pdf

# Abbreviations

| | Definition |
|---|---|
| **ADS** | IBM Automation Decision Services – Tested product from the CloudPak that provides a comprehensive environment for authoring, managing, and running decision services |
| **BAI** | IBM Business Automation Insights – Product from the CloudPak that processes event data so that you can derive insights into the performance of your business. You can use this data to drive automations and visualize the state of the indicators that matter most to you in near real-time |
| **CPU** | Central Processing Unit - In Kubernetes, 1 CPU unit is equivalent to 1 physical CPU core, or 1 virtual core, depending on whether the node is a physical host or a virtual machine running inside a physical machine. See: https://kubernetes.io/docs/concepts/configuration/manage-resources-containers/#meaning-of-cpu |
| **CS** | IBM Cloud Pak Common Services – a.k.a IBM Cloud Pak foundational services, this Cloud Pak  provides key foundational services See: https://www.ibm.com/docs/en/cpfs?topic=about |
| **CR** | Custom Resource: Description YAML file of object ICP4ACluster |
| **PSIT** | Performance and System Integration Testing team |
| **PVC** | A PersistentVolumeClaim (PVC) is a request for storage by a user See: https://kubernetes.io/docs/concepts/storage/persistent-volumes/ |