



ARIS

DISTRIBUTED INSTALLATION

COOKBOOK

ARIS 10.0 SR4

April 2018

Important disclaimer:

This document provides you with the knowledge for installing ARIS in distributed environments.

For acquiring that knowledge Software AG recommends that you attend the training course **ARIS Server Installation** available via <http://softwareag.com/education>, or ask for an ARIS installation service by Global Consulting Services.

TABLE OF CONTENTS

1	Introduction & Overview	7
1.1	Motivation	7
1.2	Structure of this document	7
1.3	Prerequisites	8
1.4	ARIS Architecture	8
1.4.1	ARIS provisioning & Setup	9
1.4.2	Advantages of distributed microservice architectures	10
2	High availability primer	11
2.1	Availability of software systems	11
2.2	Making ARIS highly available	13
2.2.1	Data redundancy	14
2.2.1.1	Application-managed replica consistency	14
2.2.1.2	Service-managed replica consistency	15
2.2.2	The CAP theorem	16
2.2.3	Maintaining replica consistency	18
2.2.3.1	Strict consistency	18
2.2.3.2	Eventual consistency	19
3	Preparation steps	23
3.1	Preparing a fresh installation	23
3.2	Install ARIS Cloud Controller on your admin machine	24
3.3	Create a project directory	25
3.4	Install the ARIS Agent on all nodes	25
3.5	Controlling multiple nodes with a single ACC instance	25
3.5.1	A look at the common ACC command line switches	25
3.5.2	A different way to start ACC	26
3.5.3	The "add node" command	27
3.5.4	Qualifying commands with the node name	28
3.5.4.1	The "on <nodename>" clause	28
3.5.4.2	Setting a "current node"	29
3.6	Adding our worker nodes to ACC	30
3.7	The "on all nodes" clause	32
3.8	Creating a node file	32

3.9	Using the node file	33
3.10	Preparing and using a remote repository	34
3.11	Making the agents on the worker nodes use the remote repository	37
3.12	Get and use a generated.apptypes.cfg file with your ACC	37
3.13	Securing access to your agents	38
3.14	Our distributed example scenario after performing the preparational steps.	39
3.15	Overview of the remaining steps	40
4	Zookeeper clustering and the Zookeeper manager	41
4.1	Motivation - The role of Zookeeper in the ARIS architecture	41
4.2	Zookeeper ensembles & failover	42
4.3	Failover scenarios	43
4.4	Setting up a Zookeeper cluster (a.k.a. "ensemble") manually	44
4.4.1	Example	45
4.5	Setting up a Zookeeper cluster with ACC's built-in Zookeeper manager	46
4.5.1	Configuration	46
4.5.2	Commands	46
4.5.2.1	Add zookeeper instance	46
4.5.2.2	Remove zookeeper instance	47
4.5.2.3	Commit zookeeper changes	47
4.5.2.4	Discard zookeeper changes	47
4.5.2.5	List zookeeper instances	47
4.5.2.6	Show zookeeper changes	48
4.5.2.7	validate zookeeper changes	48
4.5.2.8	validate zookeeper ensemble	48
4.5.3	Usage example	48
4.6	Bootstrapping the Zookeeper connection	49
4.6.1	How do runnables find Zookeeper?	49
4.6.2	The zookeeper.connect.string configure parameter	50
4.6.3	Automatic setting of zookeeper.connect.string (≥ 9.8.2)	51
4.7	Zookeeper in our example scenario	52
5	Elasticsearch clustering and the Elasticsearch manager	53
5.1	Motivation - The role of Elasticsearch	53
5.2	Elasticsearch terminology and clustering	53
5.3	Why clustering Elasticsearch became more difficult with ARIS 10	58

5.3.1	Manually setting up an Elasticsearch 2.x cluster	59
5.3.1.1	Configuring an Elasticsearch instance to find the other cluster members	59
5.3.1.2	Additional Elasticsearch cluster configure parameters	59
5.3.2	Example of a manually configured three node Elasticsearch cluster	61
5.4	Setting up an Elasticsearch cluster with the Elasticsearch manager	63
5.4.1	Configuration	63
5.4.2	Commands	63
5.4.2.1	Validate elasticsearch cluster	63
5.4.2.2	List elasticsearch instances	63
5.4.2.3	Add elasticsearch instance	64
5.4.2.4	Remove elasticsearch instance	64
5.4.2.5	Validate elasticsearch changes	64
5.4.2.6	Commit elasticsearch changes	65
5.4.2.7	Reset elasticsearch changes	65
5.4.3	Usage examples	65
5.5	Elasticsearch in our example scenario	66
6	Relational database backends	66
6.1	Motivation - The role of relational databases in the ARIS architecture	66
6.2	Scaling & high availability with relational databases	66
6.2.1	Distributing the load of multiple tenants across several DB instances (no high availability)	67
6.2.2	High availability for the relational database backend by using an external, highly available database cluster (Oracle or SQL Server)	67
6.2.2.1	Registering the external DBMS as an external service	68
	Registering an external Oracle DBMS	68
	Registering an external MS SQL DBMS	70
6.2.2.2	Creating the database schemas for each tenant	71
6.2.2.3	Assigning the tenant to the database service	72
6.2.2.4	Enhancing the runnables with the JDBC driver for the external database system	73
6.3	Relational databases in our example scenario	74
7	Cloudsearch clustering	75
7.1	Motivation - The role of Cloudsearch in the ARIS architecture	75
7.2	Cloudsearch clustering, variant 1: data load distribution of multiple tenants across several instances, no high availability, no query load distribution	77
7.3	About failure scopes	78

7.4	Cloudsearch clustering, variant 2: high availability and query load distribution, no data load distribution	78
7.5	Cloudsearch clustering, variant 3: high availability, query and data load distribution	79
7.6	Some general guidelines & recommendations	81
7.6.1	Multiple instances in one data center only makes sense when using more than one tenant	81
7.6.2	Size your Cloudsearch instances equally	82
7.6.3	Have the same number of Cloudsearch instances in each data center	82
7.6.4	Adding Cloudsearch instances retroactively will not help existing tenants	82
7.7	Cloudsearch in our example scenario	83
8	Clustering ARIS micro-service-based applications	84
8.1	Motivation - How distributing micro-service-based applications works	84
8.1.1	Load balancing and preferred routing	84
8.1.2	Failure handling	85
8.2	Distributing stateless applications	86
8.2.1	Configuring document storage for high availability	86
8.2.2	Restrictions of the tenant management application (inside UMCAAdmin) regarding high availability	90
8.2.3	Restrictions of the simulation component regarding high availability	91
8.2.4	Restrictions regarding the Dashboarding ("ARIS Aware") component regarding high availability	91
8.2.5	Restrictions regarding the ARCM component regarding high availability	93
8.3	Distributing applications in our example scenario	94
8.4	Finishing touches	96
8.4.1	Adding JDBC drivers	96
8.4.2	Adding help files	97
9	Highly available client access to a distributed installation	100
9.1	The roles of the ARIS load balancer (ARIS LB)	100
9.1.1	The same origin policy	100
9.1.2	TLS termination proxy	100
9.1.3	Load balancing and high availability	101
9.1.4	The X-Forwarded-* headers	102
9.2	Configuring load balancer in a distributed, highly available installation	106
9.2.1	Make the ARIS LBs correctly to play well with a HA loadbalancer	106
9.3	The load balancers in our example scenario	107

9.4	Optional - Simulating a highly available hardware load balancer with software (for educational purposes only)	108
9.4.1	Setting up Apache HTTPD	109
10	Administration of a distributed installation	113
10.1	Starting and stopping a distributed system	113
10.2	Collecting and deleting log files	115
10.2.1	The "collect log files" command	115
10.2.2	The "delete log files" command	116
10.3	Updating a distributed installation	117

1 Introduction & Overview

1.1 Motivation

The goal of this document is to provide the reader with the information necessary to manually set up a distributed (or "multi-node") ARIS installation. A distributed installation is an installation where the ARIS components (the "runnables") are not only located on a single machine, but where they are spread across several machines. There are several possible reasons for setting up a distributed ARIS installation instead of keeping it simple and just installing ARIS on one machine via the setup:

Available hardware resources & cost: For some more demanding usage scenarios, the total hardware requirements (in particular number of and performance of CPU cores, available RAM) of all ARIS components might be more than a single machine (or at least a single *affordable* machine) can handle. By installing the components across several nodes, it is possible to use several smaller (and thus cheaper) machines instead of one larger and perhaps prohibitively expensive machine for an ARIS installation.

Requirement for high availability: When installing ARIS in the standard way on a single machine, the availability of your ARIS system will depend on the availability of this single machine. ARIS will become unavailable for its users if that machine becomes unavailable (e.g., hardware defect, power or network outage of the machine, power or network outage of the entire data center etc.). In some usage scenarios, it is desirable to have ARIS available even in the case of such incidents. By spreading the ARIS components across multiple machines (which might even be located in different data centers, to also protect against the failure of an entire data center) and installing them redundantly, it is possible to keep ARIS up and running even under such circumstances.

Even if you do not actually have such requirements and thus do not plan to distribute the runnables across multiple nodes, this document will also give you the necessary information to manually (i.e., without the use of the setup) install a *single*-node installation - after all, a single node installation is just a special case of multi-node installation.

Throughout this document, as a running example we will develop a distributed ARIS Connect installation. We tried to keep this scenario as close as possible to what one would do in an actual customer environment, but in some places compromises have been made to keep the efforts and hardware requirements manageable. In those places, the compromises will be pointed out so that you know how to adjust when installing an actual production system.

1.2 Structure of this document

The document is structured in several chapters. After a brief overview of the ARIS 9/10 architecture in the remainder of this chapter, chapter 2 continues with an overview of the basic concepts of distributed computing and high availability. Chapter 3 explains all the necessary preparation steps to get started with a distributed installation. Chapters 4 to 7 explain the role of essential backend services of every ARIS installation, Zookeeper, Elasticsearch, the relational database and Cloudsearch, and detail how to make them highly available. Chapter 8 explains how the micro-service-based ARIS application components can be distributed for high availability. Chapter 9 describes how to make sure that also the access of clients to the installation does not become a single point of failure. Chapter 10 gives a few basic guidelines on how to operate a distributed ARIS installation.

1.3 Prerequisites

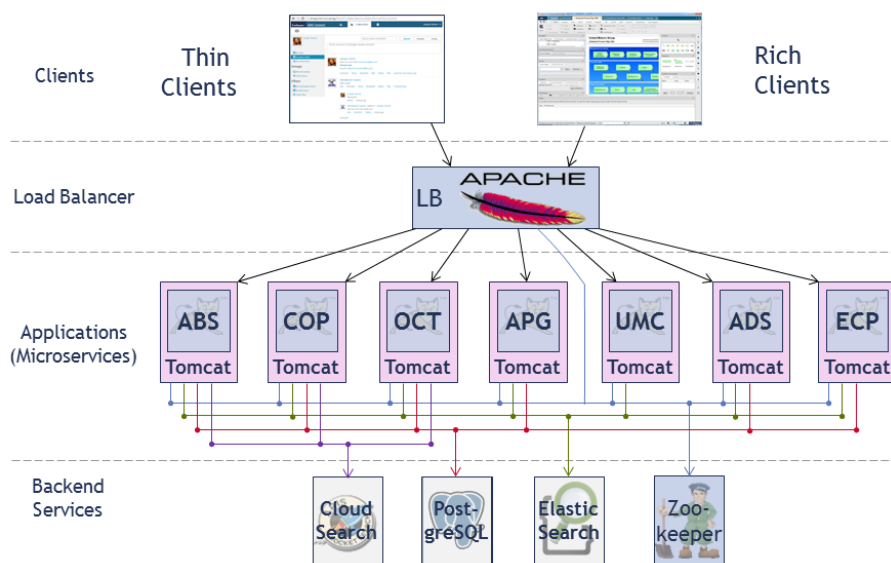
It is highly recommended that the reader already has a good understanding of and some familiarity with the ARIS provisioning framework, as this is the main toolset that will be used throughout this document.

To actually try out the sample scenario that is developed throughout this document, one needs access to three physical or virtual machines which will act as "worker nodes", i.e., they will be the machines on which the ARIS Cloud Agent and the ARIS runnables will be installed. For training purposes, the scenario can be kept small so that it should be possible to run it on a host with 4 CPU cores and 32GB RAM, running 3 virtual machines with 8GB RAM each. The example scenario can be easily adapted also for production use, but then a significantly higher amount of hardware resources will be needed.

1.4 ARIS Architecture

ARIS 9 and 10 are based on a microservice architecture, illustrated in figure 1 below. In a microservice architecture, the functionalities of the application are no longer contained inside a single, monolithic application server, but instead spread out over many smaller application servers, the **microservices**. The microservices are individual operating system processes that are stateless, i.e., they themselves do not store any persistent application state (i.e., the actual application data like ARIS models, users, documents, comments etc.). The persistent application state is instead stored in the **backend services**. A **load balancer (LB)** is used as a single point of entry for access of all kinds of **clients** to the application. ARIS offers both web-based thin clients running in a browser and more traditional rich clients that run as dedicated applications on the client machines.

FIGURE 1: ARIS 9/10 ARCHITECTURE

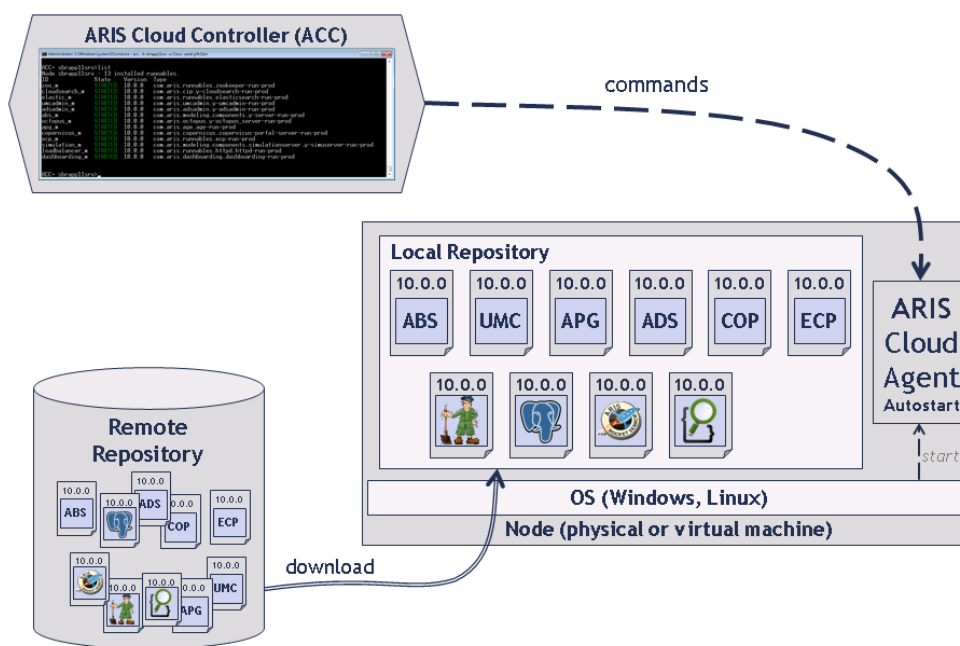


The individual microservices and backend services can be deployed independently. To make that possible, both microservices and backend services are packaged as so-called **runnables**, which can be installed, started, stopped and configured using the ARIS provisioning framework. Depending on the actual product being installed (ARIS Connect, ARCM,...), the actual types of runnables needed differs. A minimum installation of a specific product will contain each runnable exactly once, and all runnables will be located on the same node.

1.4.1 ARIS provisioning & Setup

The ARIS provisioning framework consists of two main components, illustrated in figure 2 below: ARIS Cloud Agent is a small service that needs to be installed on all machines that should serve as a worker node in an ARIS installation. It offers a REST API through which it can receive commands from the second component, ARIS Cloud Controller (ACC). The cloud controller is a console application which accepts commands entered via an interactive shell. It can control several cloud agents at once, which is useful in particular when working with distributed environments.

FIGURE 2: ARIS PROVISIONING ARCHITECTURE



The core functionality of provisioning is to "install" runnables (with the "configure" command), start, stop, reconfigure or update them. When told to install a runnable, the agent will try to download it from a remote repository, whose URL needs to be specified before.

ARIS provisioning is also used under the hood by the ARIS setups. The setup will install the cloud agent as a service on the respective machine and then use an embedded ACC to execute the commands (in particular configure commands) needed to install all runnables required for the product being installed exactly once. While the setup is running, it also provides the remote repository from where the agent can download the runnables and other required artifacts (e.g., help files).

For many ARIS usage scenarios, it is sufficient to have a single instance of each type of runnable in the system. This is also what happens when you install ARIS with the setup: all runnables required for the respective product are installed exactly once inside ARIS Cloud Agent.

1.4.2 Advantages of distributed microservice architectures

An distributed microservice architecture has many advantages over traditional monolithic applications, where all functionality is contained in one big piece of software:

Scalability - If an application has to deal with increasing load (e.g., increasing number of concurrent users, increasing amounts of data etc.), it will sooner or later reach a point where performance degrades, i.e. response time increase, batch jobs take longer, and the overall user experience suffers. Aside from improving the performance of the software itself, which is of course a lengthy process and has its limits, the only other option is to make more hardware available to the application. In a traditional monolithic application, your only chance to remedy this situation is to **scale up** your hardware, i.e., buy a bigger server. While this will work well to some point, there are limitations of how big you can make a single server. And even if you are still far away from what is technically possible, once you leave the range of mainstream "run-of-the-mill" commodity servers, buying more powerful hardware becomes over proportionally more expensive. This restriction becomes even more severe when you are working with software virtualization in your data center or are even running your software on virtual machines provided by an Infrastructure-as-a-Service (IaaS) provider: usually there are rather strict limits of how big you can make an individual virtual machine. With a distributed application, you can instead simply add new instances of the individual application components, instead of making the single instance bigger, i.e., do a **scale-out**. If the capacity of a single server does not suffice anymore, you can simply add an additional server and put some of the components there, instead of buying a new bigger server.

Higher quality due to improved maintainability - In a monolithic application, there is no technical limitation about which application components can interact with each other - everything can directly talk with everything else. This makes it easy to create new dependencies among individual parts of the application (often without knowing). Unless very strict discipline is observed when making changes, this can over time lead to a code base that becomes harder and harder to change and maintain, as since everything can depend on everything else, the effects of what appears to be a small local change on the rest of the system become harder and harder to estimate. If instead you are working with isolated components with clearly assigned functionalities (i.e., microservices), you invariably have to interact with other components through dedicated, well-designed, public and usually quite stable interfaces, making the interdependencies clear and isolating the components from changes in others. This reduces the amount of coupling among components and makes it easier to change and test the individual components, making it much easier to guarantee good software quality even when adding new functionality.

High availability - If you have a monolithic application on a single server, if anything happens to that server (e.g., hardware failure, network or power outage etc.), your application will become at least temporary unavailable. If an application is mission-critical for an organization, such an outage can cause massive cost. While it is possible to make single servers highly available by using corresponding hardware with redundant components, this is again expensive and there are limits to what can be done. With a properly designed distributed application, it is instead possible to stick to commodity hardware (for the most part) and to achieve high availability by configuring the application in a redundant way, so that the outage of individual components allow it continue working (albeit perhaps with reduced capacity).

The remainder of this document will discuss how to leverage the advantages of the microservice-based ARIS architecture to set up distributed and highly available ARIS installations.

2 High availability primer

2.1 Availability of software systems

A software application that is vital for an organization to do its business should ideally be available to its user 24h/7 days each week. Availability is usually understood as a situation in which the system is able to fulfill its functions without restrictions. A system is highly available if it can fulfill its functions for a higher than normal period. Availability is usually given as a decimal number between 0 and 1, that is obtained by dividing uptime by total time, which can also be understood as a percentage. In high availability parlance, it is common to refer to the number of nines behind the decimal point to indicate the rough category of high availability of the system. For example, an availability of 0.999 would be a "three nines" or 99.9% availability. While this might seem like good availability at first glance, an outage of 0.1% of a year is already a whopping 525 minutes.

Terminology - High availability, fail-safe etc.

There is often quite some confusion when it comes to terminology. In particular, the term "fail-safe" is often wrongly used synonymously with "high availability". However, if a system is fail-safe, this does not mean that failure is impossible or improbable, but that it is built in a way such that in case of failure, it will fail in the safest possible way, i.e., with minimal harm to other equipment, the environment, or to people. High availability, on the other hand states that the change of failure of the system is lower than average, or, put differently, that it is available for a higher than normal period. Consequently, we will use the term "high availability" throughout this document.

Aside from the inevitable need to update the software (which in all but the most demanding environments will mean a maintenance window, i.e., a planned downtime of the system, and which is not in the scope of this document), there is a plethora of unplanned things that can possibly go wrong in a computing environment that reduce the availability of an application:

- **Software failure** - Probably the most obvious (and arguably most common) cause for a software system that can lead to it becoming unavailable is a crash of the software. The most obvious reasons for this are of course bugs in the software itself, and as we all know, all but the most trivial programs will inevitably have bugs. Preventing, detecting and fixing bugs (and that as early as possible, ideally before a customer encounters it) is perhaps the most essential task of any software R&D organization. Still, fact is, bugs *will* happen, and sometimes they will crash your entire software (i.e., the operating system process will be terminated) or the software will still be running, but cease to function properly. The operating system itself can also have bugs that can make individual processes or the entire OS crash.
- **Hardware failure** - The hardware the software is running on can have a defect. Quite obviously, mechanical components like hard disks can break down, but also solid state devices, including CPUs, chipsets, memory, SSDs, network devices or peripheral components can fail. Sometimes failures are not permanent. For example, even while quite rare, data held in RAM can become corrupted (e.g., a single bit can be inverted), which could cause the program to crash (and making it seem like a software failure).
- **Infrastructure failure** - The infrastructure required by an application to function and be available to its users can fail. For example, power supply or network connectivity can be interrupted, either due to human error (someone turning the wrong switch or pulling the wrong cable, an excavator used carelessly near vital power lines or network cables), or due to a large scale blackout affecting an entire region.

As it should have become clear, the **scope affected by** any such **failures** can vary:

- **Process failure** - An individual operating system process (i.e., in ARIS terminology, a single runnable) can crash, e.g., due to hitting a bug.
- **Node failure** - An entire machine can go down, e.g., due to problems with the power supply or a crash of the OS.
- **Datacenter failure** - A whole data center can be affected, in particular if vital components of the infrastructure are compromised.
- **Network failure** - Even if the software and the machines itself work fine, the network connectivity between them can fail. From the perspective of the individual node, it will often not be distinguishable from the failure of the other nodes, but a network failure can produce other kinds of issues compared to a process, node or data center outage, which is why we list it separately here.

The take away message here is: Everything can (and eventually will) fail, all the time!

The key to protect yourself and your vital software system against any of these failures is the introduction of redundancy.

The obvious approach to address a lot of the different kinds of failures above is to make hardware redundant. For example, you can buy servers with redundant power supplies (that might each be connected to an independent mains circuit or to an "uninterruptible power supply" device that uses battery power or even an auxiliary power generator to bridge blackouts for varying lengths of time) and with redundant network adapters that are connected to redundant network switches. You can store your vital application data redundantly on a RAID system to cover hard drive failures etc. The whole idea of redundancy is to reduce or ideally eliminate single points of failure in your system.

The amount of additional effort one is going to invest to protect your software system against outages will vary depending on how vital it is for the respective organization.

However, making hardware redundant has its limitations. While a lot can be done, at some point the cost will increase exponentially. Further, making your hardware redundant will not help you in the case of issues with the software itself.

In a classical monolithic application, there is little that can be done at this point. Fortunately, for a software built around a modern, microservice-based architecture like ARIS, the situation is much better. Not only the hardware the software runs on can be made redundant, but the individual software components themselves (i.e., the microservices and backend services) can be made redundant as well, simply by having more than one of each kind of service inside an application. If configured properly, an ARIS installation can be built so that it survives failures on various scopes:

- If you have properly configured each microservice and backend service redundantly on a single machine, you can make your ARIS installation **resilient to a process failure** (i.e., an individual crashed runnable).
- If you properly spread your microservices and backend services across multiple machines, you can make your ARIS installation **resilient against the failure of an entire node** (whether it is caused by a hardware or a software problem).
- If you have extremely demanding availability requirements for your ARIS installation, you can even configure the individual microservices and backend services redundantly across several data centers, **making your system resilient to a data center outage**.

At this point it shouldn't be forgotten to mention that while only a distributed software can potentially(!) reach a higher degree of availability than a monolithic one, a distributed software is also faced with new kinds of failures that simply cannot occur in a monolithic application: For example, in a monolithic application, the individual components of a software communicate locally within a single operating system process, so the software does not need to worry about a network outage that might temporarily make it impossible to communicate with another component. In a distributed system, where the individual components run as dedicated operating system processes that possibly reside on different machines, this communication has to happen over the network, and a network issue suddenly has to be taken into account by the communicating components. A distributed application that is to become highly available has to be built in a way that it handles such failures gracefully and in particular prevents cascading failures that ultimately bring down the whole system. Whenever writing code that communicates with other components that might or might not be located locally, developers have to build their code in a way such as to expect that communication might fail and devise ways to recover from such failures. Consequently, making software capable of being distributed introduces new levels of complexity, which can also introduce new bugs.

So for simple deployments, where high availability is not required, the additional failure modes and the potential for additional bugs will at first be a disadvantage of a distributed over a monolithic application. On the other hand, without being able to distribute an installation across nodes or even data centers (like it would be the case with a monolithic application), there would be no possibility to optionally make an application highly available at all.

2.2 Making ARIS highly available

Making an application like ARIS highly available is not solely the responsibility of the developers of the software, but also that of the administrators that set up an installation. Depending on which scope of errors you want to protect your ARIS installation against, there are ways to deploy it accordingly, but of course with increasing amounts of overhead and effort involved, the more severe the failure is you want the system to be able to tolerate.

If one doesn't care about high availability at all, one can just stick with a single-node installation. The automatic restart of crashed components will quickly resolve minor issues and in many environments a short outage of parts of an ARIS installation can be tolerated. However as we said above, the **key to high availability is having redundancy** built into the system: only if you have extra instances of components available, you can survive the outage of one (or more, if desired) instance(s).

If you want ARIS to survive the outage of a single runnable, you can simply configure the key components redundantly on a single node (or on multiple nodes).

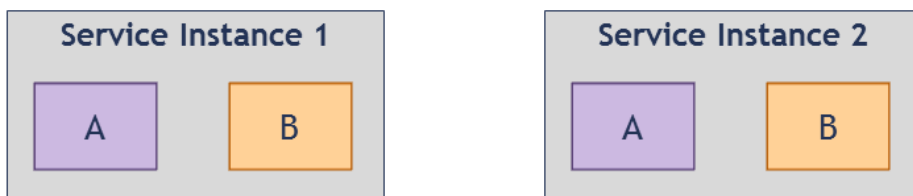
If you want to survive the outage of individual nodes or the network between them, you can deploy the components redundantly across multiple nodes. If it is desirable that an ARIS installation remains available even if an entire data center fails, you can deploy the components redundantly across data centers.

However, redundancy comes at a price: not only do you have to spend additional hardware resources that hold the redundant components, but you also need to **hold the actual data redundantly** (i.e., have several - ideally always identical - copies of the same data item), so that in case a backend service instance fails, you can still access all data. Here the problem is to keep these copies (also called "replica") consistent, not only during normal operation, but in particular in the face of any of the above mentioned failures.

2.2.1 Data redundancy

Let's consider that we have two instances of a backend service in which we store two copies of two data items A and B, as shown in figure 3 below. For the purposes of this chapter, here, a "data item" can stand in for data at any level of granularity, e.g., a single variable value, a single tuple in a relational database, an entire table, or an entire database.

FIGURE 3: DATA REPLICATION



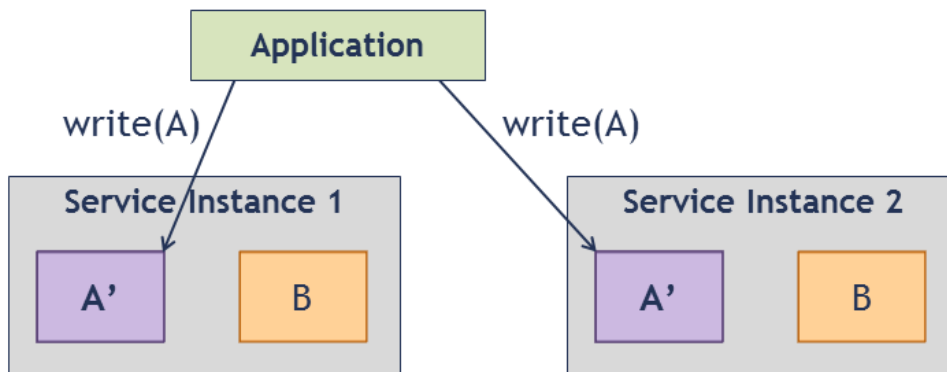
As long as the data itself doesn't change, the data from both instances can now be used to service (read) requests. The trouble starts when data is modified. There are basically two approaches:

- **Application-managed replica consistency:** the applications/microservices themselves (or some component inside them) are aware of the replication and make sure that all replica are updated.
- **Service-managed replica consistency:** the service itself (i.e., all instances together) is responsible for making sure that all replicas are updated.

2.2.1.1 Application-managed replica consistency

When the consistency of replica is managed by the applications, this means that all application instances need to be aware of the fact that data is replicated and need to know where (i.e., in which service instance) each replica is located. When a user request (or some background job etc.) requires that a data item is changed, the application is then responsible for updating all replica. Continuing from our example from figure 3 above, figure 4 below shows this approach. If a write happens on data item A, the application updates both copies, basically leading to a new, changed version of that data item, indicated by the prime symbol ('). Only after all copies are successfully updated, the write operation is really completed and the application can proceed with serving the user request (or executing the background job etc.).

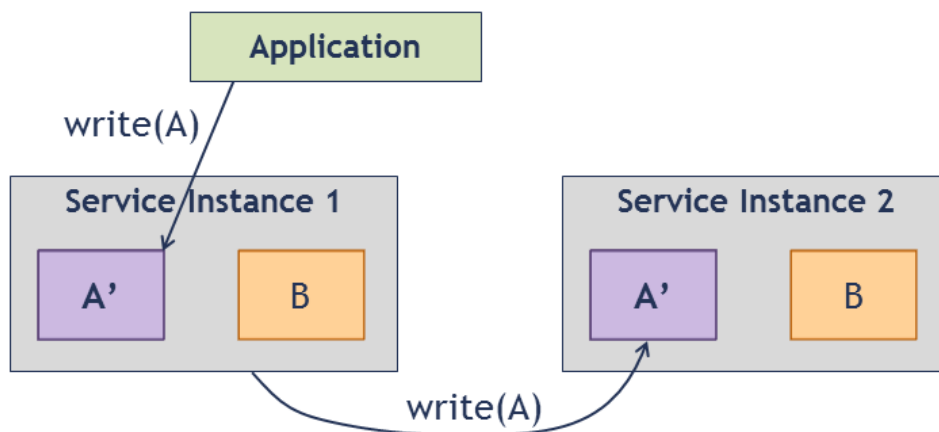
FIGURE 4: APPLICATION-MANAGED REPLIA CONSISTENCY



2.2.1.2 Service-managed replica consistency

The more common approach is to leave the management of replication, in particular that of updating redundant copies, to the service itself. So the application that wants to update the data item can do the write operation via *any* service instance (even via an instance that doesn't hold a replica of the particular data item), and leave it to the service instance itself to propagate the changes to the other instances and make them update their replica. This approach is illustrated in figure 5 below.

FIGURE 5: SERVICE-MANAGED REPLICA CONSISTENCY



Both application- and service-managed replica consistency have their pros and cons. But the actual difficulties that one encounters when trying to maintain replicas consistently are the same in both scenarios. In the next sections, we will discuss these problems, using service-managed replica consistency in our examples. But before, we will first have a look at the theoretical underpinnings of data consistency, in particular the so-called CAP theorem.

2.2.2 The CAP theorem

The CAP theorem¹, also known as Brewer's theorem, states that in a distributed system, it is impossible to provide more than two out of the following three properties²:

- **Consistency:** Every read operation is guaranteed to always see the most recently written value of a data item (or an error, if that is not possible)
- **Availability:** Every request receives a non-error response (but without guarantee that the response will contain the most recently written values of the involved data item)
- **Partition tolerance:** The system continues to operate despite an arbitrary number of messages being dropped (or delayed) by the network between nodes

What this boils down to is that in the case of a network partitioning, one has to choose either consistency, or availability, i.e., you can either get the data, with the risk of not seeing the most recent state of it (i.e., compromise consistency), or not getting it at all (but an error, i.e., compromising availability). In the absence of any network problems, you can have both consistency and availability.

A distributed, clustered installation of a traditional relational database systems will always choose consistency over availability (as only this will allow the system to provide the ACID guarantees³ that RDBMSs always provide, i.e., **A**tomicity, **C**onsistency, **I**solation and **D**urability), while most NoSQL systems will at least by default prefer availability, offering so-called eventual consistency (often called BASE, "**B**asically **A**vailable, **S**oft state, **E**ventual consistency", in contrast to ACID). In an eventually-consistent environment, it is not guaranteed that all redundant copies of a data item will always have the latest state, but they *eventually* will.

A few words about ACID transactions

ACID transactions⁴, which are the underlying paradigm of any "traditional" database system for updating data, offer the following four properties:

Atomicity: From the perspective of clients, all changes to the data done in a single transaction are either done all at once, or not all ("all or nothing"). So either the transaction completes with a successful "commit", in which case all its changes are applied to the database, or the transaction fails or is aborted, in which case no changes at all will be applied to the database. Internally, this is usually achieved by some form of undo mechanism. For example, the database system will write transaction log files that allow it to undo any change that has already been applied in case of an error or abort by the client or redo any changes of a committed transaction that have not yet been persisted.

Consistency: The changes done in a transaction will bring the content of the database from one consistent state to a new consistent state. Note that here, consistency means "being consistent from the perspective of the database system's data model", i.e., all constraints defined by the data model (in the case of the relational model, such constraints are in particular referential integrity, the uniqueness of primary keys etc.) are maintained. This is in general handled by checking the consistency constraints at the end of the transaction and aborting the transaction in case of a constraint violation (or by actively performing the necessary operations to make the database consistent again, e.g., by cascading deletes).

¹ Gilbert, Lynch: Brewer's conjecture and the feasibility of consistent, available, partition-tolerant web services, ACM SIGACT, Volume 33 Issue 2, June 2002, Pages 51-59 (PDF)

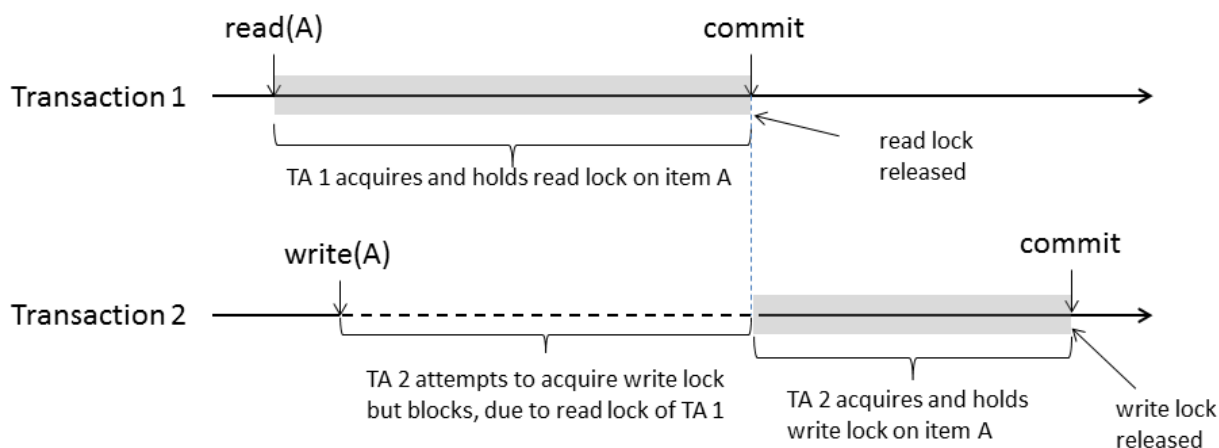
² Wikipedia: The CAP theorem https://en.wikipedia.org/wiki/CAP_theorem

³ Wikipedia: ACID <https://en.wikipedia.org/wiki/ACID>

⁴ Wikipedia: ACID <https://en.wikipedia.org/wiki/ACID>

Isolation: Different transactions working concurrently on the same set of data will see the database state as if they were executed strictly sequentially (in any possible order). A simple approach to achieve this guarantee is by using a locking approach with read and write locks: before a data item can be read by an application, a read lock has to be set, before a data item can be changed, a write lock has to be set. Locks, once they are set, are then held until the end of the transaction (i.e., until it is committed or aborted). When a transaction tries to set a lock on a data item and no other transaction is currently holding a lock on it, the lock operation will complete successfully, the transaction will obtain the lock and can continue working with the data item. Different read locks on the same data item acquired by different transactions are "compatible" and can coexist (since concurrent reading is possible without any conflicts). Write locks are not compatible with read locks or other write locks. So if a transaction attempts to obtain a write lock in order to change a data item, but other transactions have already set a read lock, the write lock cannot be obtained. This situation is shown in the figure below.

A simple example of locking-based concurrency control



Transaction 1 (TA 1) performs a read on a data item A. In order to be allowed to do this, it acquires and holds a read lock, which is held until it commits. Briefly after TA 1 acquires the read lock on A, a second transaction TA 2 attempts to change the same data item, for which it thus needs to acquire a write lock. Due to the read lock already held by TA 1, the attempt to get the write lock blocks, so TA 2 is suspended. Only after TA 1 commits and releases its read lock, the blocked attempt of TA 2 to get the write lock is de-blocked, TA 2 can obtain the write lock and modify A, itself releasing the write lock when it commits. Similarly, if a transaction tries to set a read lock but a write lock is already set, the reading transaction will block, until the writing transaction completes (or aborts) and releases the write lock.

In effect, any transaction will see the database in a state as if it was running sequentially with all other applications (i.e., one transaction after the other), even though of course the transactions can run concurrently to some degree (until they have to wait for getting a lock on a data item that is currently locked by another transaction). Much more sophisticated locking mechanisms, or isolation mechanisms that do not require locking (so-called *optimistic concurrency control* mechanisms, e.g., multi-version concurrency) exist, which allow higher degrees of concurrency, but all offer the same level of isolation, i.e., the database state will appear to each transaction as if it was running exclusively on the database.

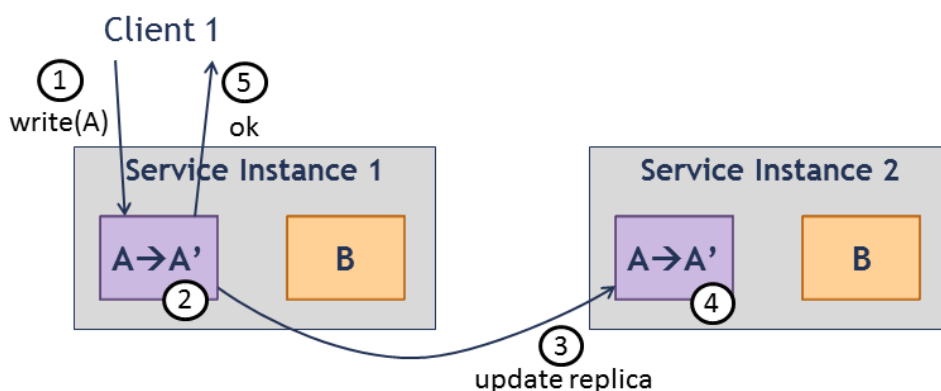
Durability: The changes done in a transaction that successfully committed stay that way, even in the case of errors, power failures etc. This property is guaranteed by making sure that changes are persisted to the data storage (which is of course slow) or by making sure that the change is written to the transaction log (which is written sequentially, which can be done faster), so that they can be reapplied in the case of a system crash.

2.2.3 Maintaining replica consistency

2.2.3.1 Strict consistency

If one wants to preserve strict consistency of data replica, this means that basically all replica have to be updated synchronously, i.e., within the write request duration, i.e., the duration of an application's write requests will increase. Even if the write operations are done in parallel, the **slowest write will be the limiting factor of performance**. Further, guaranteeing that all replica are updated consistently, incurs a significant overhead (e.g., a variant of the two-phase-commit protocol, 2PC⁵). This approach is sketched in figure 6 below.

FIGURE 6: SYNCHRONOUS, STRONGLY-CONSISTENT UPDATE OF REPLICA



The client write request reaches a service instance (1), which updates its local replica of the data item (2), then updates the remote replica (3 and 4), and only then returns a "successful" response to the client. If an error would occur while updating the remote replica (e.g., if service instance 2 became unavailable), the change to the data item would be rolled back and an error would be reported to the client instead. This makes sure that either both replica are updated or none of them. Further, concurrency control mechanisms (as the ones described when we discussed ACID transactions above) make sure that other clients that try to access the data replica in either instance during the write operation by client 1 will either see the old state of A and block the write operation of client 1 until their reads have been completed, or will themselves block until the write operation of client 1 has completed successfully (after which they will see the modified version A') or failed (in which case they will see the original version of the data).

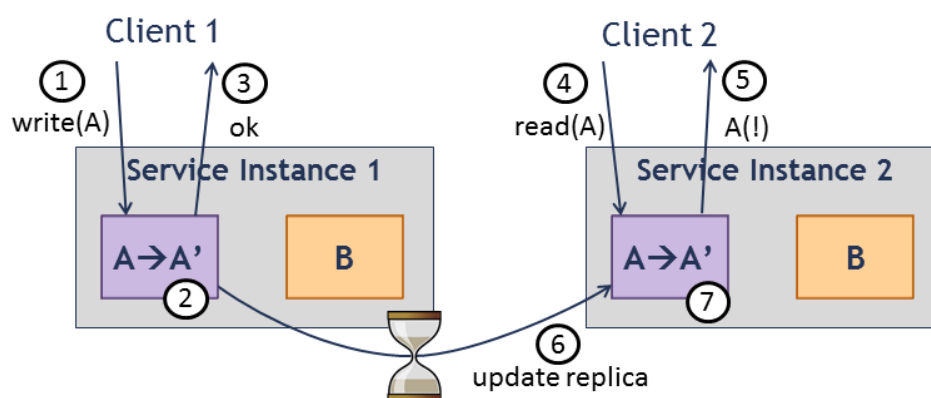
⁵ Two-phase commit protocol https://en.wikipedia.org/wiki/Two-phase_commit_protocol

Things start to get really difficult if you consider various error cases: In particular, what happens if a service instance holding a replica of a data item that is to be updated is currently unavailable? Should the application wait for it to become available? Even in the best case, that could take easily take minutes, or the instance might never come back, either causing the write operation to hang indefinitely, or run into a timeout and yield an error. So **if you favor consistency, you sacrifice availability**, just as the CAP theorem states.

2.2.3.2 Eventual consistency

Another option is to live with the situation that replica are not always in a consistent state, at least for some time. figure 7 below shows how a hypothetical service could handle this: The client write request reaches a service instance (1), which updates its local replica of the data item (2) and then immediately returns a "successful" response to the client. The replica of data item A located in service instance 2 is not yet updated. A second client could now read data item A. If it accessed service instance 2 which still has the replica that hasn't been updated yet, client 2 would still see the old value of A (4 and 5). Only after service instance 1 has *eventually* informed service instance 2 that it should update its replica (6 and 7), all replica are again in a consistent state and only then will clients that access the data item through any service instance all see the same value.

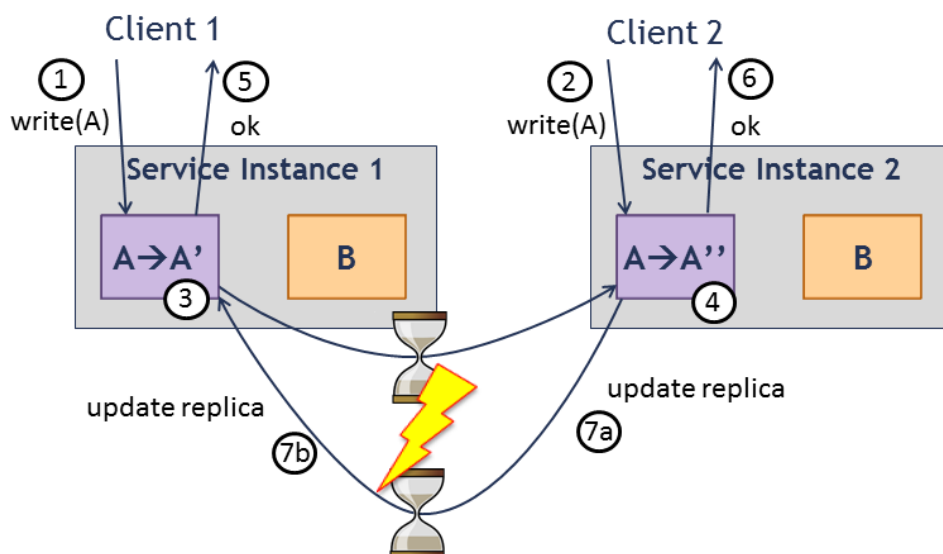
FIGURE 7: DEFERRED, EVENTUALLY-CONSISTENT UPDATE OF REPLICA



An obvious advantage of this approach - aside from replica updates happening asynchronously to the actual write operations and thus do not slow down client requests - is that it is now possible to still service write requests even if a service instance holding a replica is currently unavailable. However, now applications have to deal with the situation that they might see "stale" data. Depending on the actual application requirements, in particular for applications where a large part of the clients are only browsing the data, this might be perfectly acceptable. If for example some clients of a webpage see the old version of the page for a short time, this can often be tolerated.

However, in an eventual consistency model, keeping the replica consistent is much more challenging. Consider the situation depicted in figure 8 below:

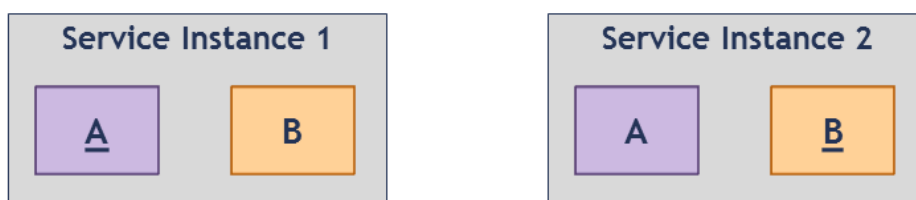
FIGURE 8: CHALLENGES WHEN UPDATING REPLICA IN AN EVENTUALLY CONSISTENT SERVICE



Here, two clients make a change to the same data item, but through different service instances. Client 1 changes the replica of data item A in service instance 1 (1, 3 and 5) to state A', while concurrently, client 2 changes the replica in service instance 2 (2, 4 and 6) from A to A''. When now the service instances try to update the respective other replica, a conflict occurs, since it is not clear which changed state of data item A, A' or A'', is the correct one. This situation has to be avoided at all cost, as it is in general not possible to resolve the resulting conflict automatically.

A common approach is a simple master-slave replication scheme. Here either an entire service instance or one explicitly selected replica of each data item is declared the "master", through which all write operations happen. In figure 9 below, we indicate which replica is the master replica by underlining it. Here, service instance 1 holds the master replica of data item A and a "slave" or secondary replica of item B, while service instance 2 holds the master replica of data item B and a secondary of data item A.

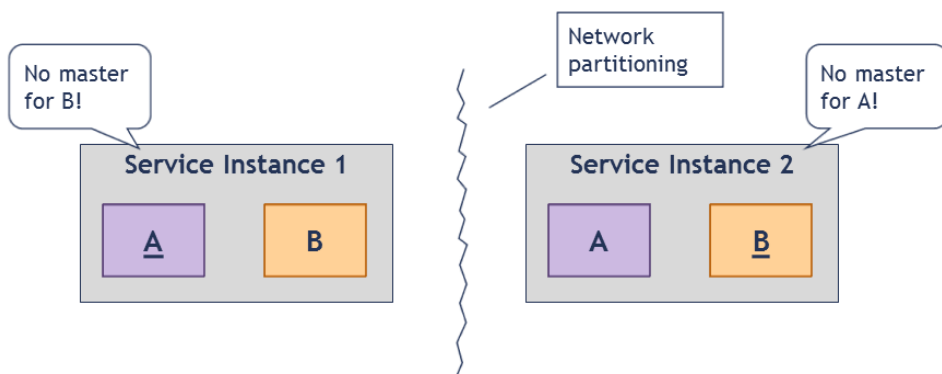
FIGURE 9: MASTER AND SLAVE REPLICA



Write operations to any data item are now managed by the instance holding the master replica, which is then responsible for propagating the change to all other instances. The instance holding the master replica of a data item is of course a possible bottleneck for writes to it. However, since the master replica of different data items can be distributed across all instances, the overall write load over all data items can also be spread out.

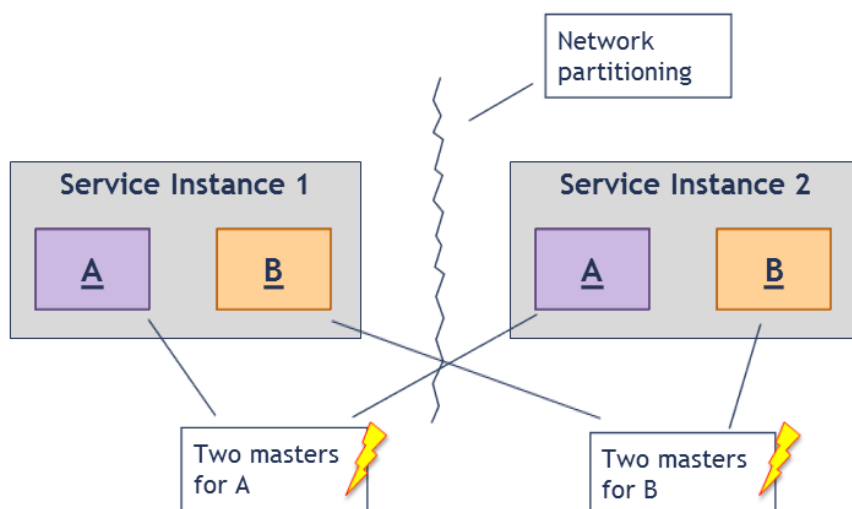
The next problem is now how to decide which replica is to be the master, and how to make sure that there is exactly one master at any given time - if there is no master, no writes can be performed, and if there are multiple masters, a problem like the one shown in figure 7 above could still occur. While the details of the actual mechanism with which this is achieved go beyond the scope of this document, the general idea is that an "election process" for the master is performed. However, even with this approach, problems can occur. Consider the situation shown in figure 10 below:

FIGURE 10: NETWORK PARTITIONING



Here, the communication between the two service instances is interrupted, a common situation in distributed computing environments. From the perspective of each instance, it might now appear that the other instance has failed, even though the problem is just the network between them. Instances 1 and 2 could now notice that they do not have access to a master of data item B or A, respectively. Without precautions, the instances might now start an election process for a new master, and since they are the only instance voting, each instance could conclude that it is now master for both data items. Observe that we now again end up with two masters for each data item, as shown in figure 11 below. This dreaded situation is called the **split brain problem**⁶.

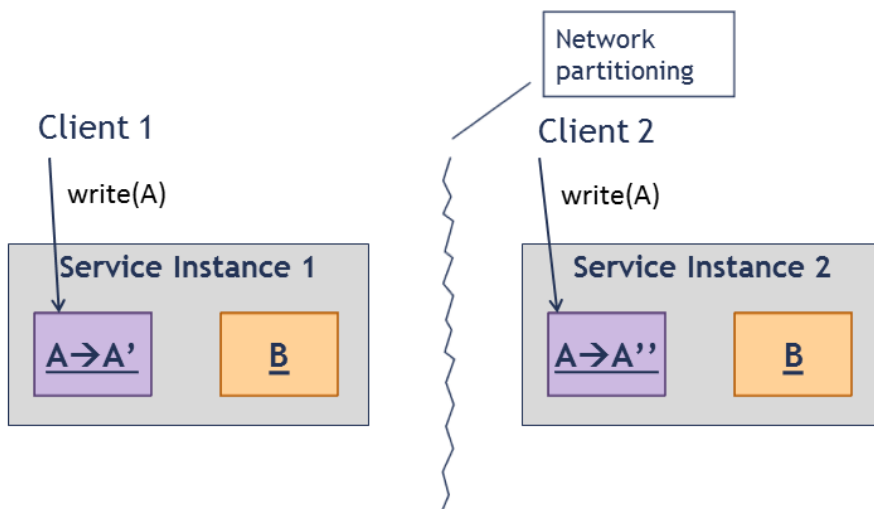
FIGURE 11: SPLIT-BRAIN PROBLEM - MULTIPLE MASTERS IN THE CASE OF A NETWORK PARTITIONING



⁶ Wikipedia: Split-brain (computing) [https://en.wikipedia.org/wiki/Split-brain_\(computing\)](https://en.wikipedia.org/wiki/Split-brain_(computing))

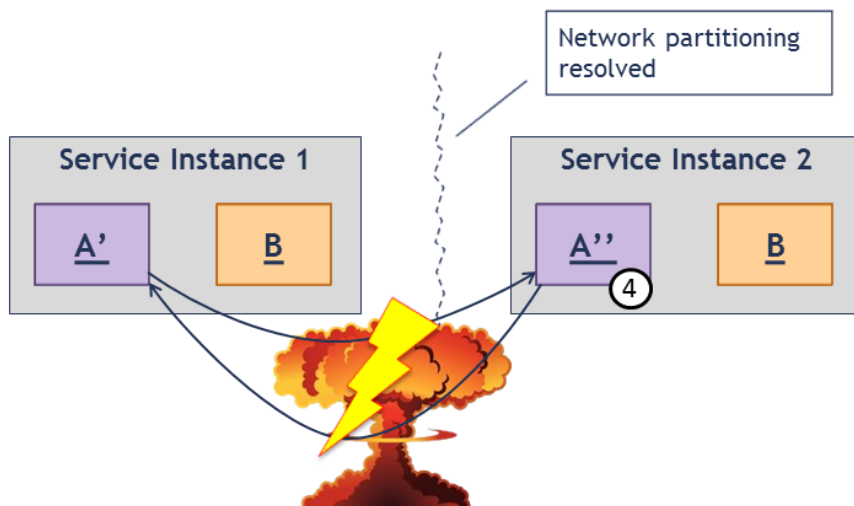
Once we are in this split-brain situation, clients on each side of the partitioning can modify the same data items independently, as shown in figure 12.

FIGURE 12: SPLIT-BRAIN IN CASE OF MULTIPLE MASTERS DURING A NETWORK PARTITIONING



The state of the data items on both sides of the network partitioning could now diverge from each other. Once the network partitioning is resolved, the issue would become apparent, once each master notices the competing master from the other side of the former partitioning. In general, the different states of the data items cannot be resolved automatically (at least not satisfactory), requiring human interaction to resolve the conflicts.

FIGURE 13: CONFLICT AFTER THE NETWORK PARTITIONING HAS RESOLVED



A common approach to avoid the election of multiple masters is to require the instances on each side of the network partitioning to obtain a quorum, i.e., require more than 50% of the voting instances (>50% **not** ≥50%) to be available before they can start the voting process for a new master. That way, only the bigger partition (i.e., the one having more instances in it) can elect a master and continue to operate, while the smaller partition will either not allow any operations at all or at least not allow any write operations. In our example scenario that we sketched so far, we only had two instances of our service.

Consequently, in a network partitioning situation like the one depicted above, neither partition would obtain the necessary number of voting instances for a quorum, and thus the service would become unavailable on both sides. There are a few approaches to avoid such ties: the most common one is to make sure that there is always an odd number of instances in the total system, so that no matter how the instances are partitioned, there is one side with a majority. Another "tie breaker" is the assignment of weights to instances, so that the side with at least 50% of the total weight can continue to operate. For example, in our two-instance scenario, we could have assigned a weight of 51% to service instance 1, and 49% to instance 2.

Notice how all the different problem scenarios discussed above are all a manifestation of what is stated by the CAP theorem: In the presence of a network partitioning, you can either live with reduced consistency (by using an eventual consistency model and/or risking a split brain situation) or with reduced availability (by using a strict mechanism to maintain replica consistency or by using a quorum approach, where the smaller partition will not be able to function).

3 Preparation steps

3.1 Preparing a fresh installation

Before you can do the actual installation of a distributed ARIS installation, you need to do a few preparation steps, which are described here. A fundamental requirement is that you have the desired number of (virtual or physical) worker machines available that are needed for the scenario you intend to set up and that you have administrative access to the worker machines.

You might wonder about how to figure out how many machines you need for a distributed scenario. Alas, there is no simple answer to that question, as the correct answer depends on the system environment, the expected load (e.g., number of current modelers and viewers), the amount of application data (e.g., number and size of ARIS databases, amount of documents, total number of users).

However, for a number of reasons that we are going to learn about over the course of this document, a minimum requirement for a distributed system that also offers some degree of high availability is to have at least three nodes. For the example scenario that we are going to use as a running example throughout this document, we assume that we have three nodes available. In our example commands, we further assume that the hostnames of these three nodes are `arisworker1.company.com`, `arisworker2.company.com`, and `arisworker3.company.com`, to which we will briefly refer to as "node 1" (or "n1"), "node 2" (or "n2"), and "node 3" (or "n3"), respectively, as illustrated in the figure below.

FIGURE 14



Over the course of this document, we will extend this basic diagram so that it contains all the components that we already installed.


In many less demanding or run-of-the-mill scenarios where customers are interested in a highly available ARIS installation, the example scenario we are going to develop over the course of this document will probably work just fine. However, for foreseeably more demanding requirements the best approach is a staged ramp-up of the introduction of ARIS: After starting with an initial system configuration and initial set of users, the system is then carefully monitored for overall performance. As more and more users are allowed access and more data is put into the system, the system is then reconfigured and/or expanded accordingly (i.e., scaled out, by adding more machines and runnables) depending on where performance bottlenecks show up.

3.2 Install ARIS Cloud Controller on your admin machine

First, you need to install ARIS Cloud Controller on your "admin machine", i.e., the machine from which you will perform the installation. You can do so with the installers provided on the ARIS DVD.

On Windows, you can install ACC with the ARIS Agent setup (located at <DVD_Root>\Setups\Windows\ARIS_Agent). The setup program is straightforward to use, you only need to choose an installation directory.

On Linux, you can use the installation package for your respective Linux distribution, which you can find at <DVD_Root>\Setups\Linux_RedHat\ARIS_Cloud_Controller or <DVD_Root>\Setups\Linux_SUSE\ARIS_Cloud_Controller, respectively. You can install these packages with the usual commands of your distributions package manager (in short: "rpm -i <package>" for RPM-based distributions, "dpkg -i <package>" for Debian-based distributions).

 Alas, there is currently no separate installer available to install only ACC on Windows; it is always installed together with the agent via the ARIS Agent setup. If you do not want to install the agent on your admin machine, you can just install the agent (incl. ACC) via setup, then stop the Windows service ("ARIS Agent") and set its start type to "manual".

You can later invoke ACC from a command window/shell using <pathToAccInstallationDirectory>\acc.bat (on Windows) or <pathToAccInstallationDirectory>\acc.sh (on Linux).

On Windows, the ACC batch file will be found at <installationDirectory>\server\acc (where <installationDirectory> is the path you chose in the agent setup directory. On Linux distributions, the ACC bash script file will usually be found at /usr/bin/acc (but that might differ depending on distribution).

On Windows you might want to add the directory containing acc.bat to your PATH environment variable, so that you can invoke ACC from a command prompt positioned in any directory. On Linux, acc.sh is usually already made available on the path by the package installer.

If you do not do that, you will later have to specify the full path to acc.bat when running ACC. In all the following steps we assume that you changed the PATH environment variable, so will call ACC without specifying the full path to it.

3.3 Create a project directory

On the admin machine from which you will be using ACC, **create a project directory** (e.g., `c:\aris_distributed`) to which you will copy all the various config files we create in the next sections.

i If you don't want to bother with changing the PATH of your Windows system, you can simply use the directory containing `acc.bat` as your project dir and put the files there.

3.4 Install the ARIS Agent on all nodes

Next, ARIS Agent needs to be installed on all nodes that you want to use for the distributed installation, by using the installers/packages for the respective OS.

On Windows, you can install the agent with the agent setup (located at `<DVD_Root>\Setups\Windows\ARIS_Agent`). The setup program is straightforward to use, you only need to choose an installation directory.

On Linux, you can use the installation package for your respective Linux distribution, which you can find at `<DVD_Root>\Setups\Linux_RedHat\ARIS_Agent` or `<DVD_Root>\Setups\Linux_SUSE\ARIS_Agent`, respectively.

You can install these packages with the usual commands of your distribution's package manager (In short, run `"rpm -i <package>"` for RPM-based distributions, or `"dpkg -i <package>"` for Debian-based distributions. For a more detailed description you can consult the "Server Installation Guide - Linux" on the ARIS DVD).

3.5 Controlling multiple nodes with a single ACC instance

While many users are familiar with using ACC to control a simple single-node installation, a single ACC instance is capable of working with multiple agents of a distributed installation. This section describes the basic commands used to register nodes with ACC.

3.5.1 A look at the common ACC command line switches

The well-known way with which ACC is started in a single node installation is to pass the hostname, username, password, the path to the generated `apptypes.cfg` file (more on that file later), and optionally the agent's port to the ACC script file (The ACC script file is named `acc.bat` and located at `<installationDirectory>\server\acc` on Windows and `acc.sh` located at `/usr/bin/acc` on Linux, respectively) via the command line, following this general structure:

```
acc.{bat|sh} -h <hostname> -p <agentPort> -u <agentUsername> -pwd
<agentPassword> -c <pathToGeneratedAppTypesCfgFile>
```

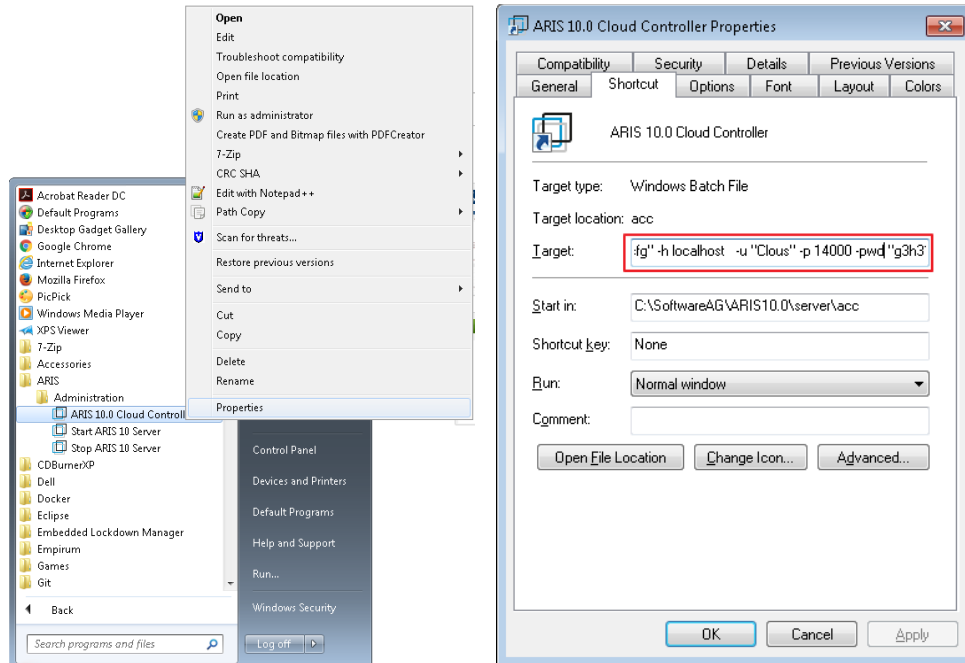
For example, in a single node Windows installation at the standard installation directory, where you didn't yet change agent username and password, you would start ACC with

```
acc.bat -h localhost -p 14000 -u Clous -pwd g3h31m -c
c:\SoftwareAG\ARIS10\server\acc\generated.apptypes.cfg
```

Note that you can leave out the `-p` command line switch if the agent is running on its default port, which is 14000 in ARIS 10 (and used to be 9001 in ARIS 9.x).

This is also the command line used in the Windows start menu shortcut ("ARIS 10.0 Cloud Controller") that is created by the setup in a single-node installation, which you can inspect via the Properties context menu entry:

FIGURE 15



For the next few steps, best is that you use any single node installation of ARIS 10 at your disposal and start ACC with this command line. It will give you the familiar prompt:

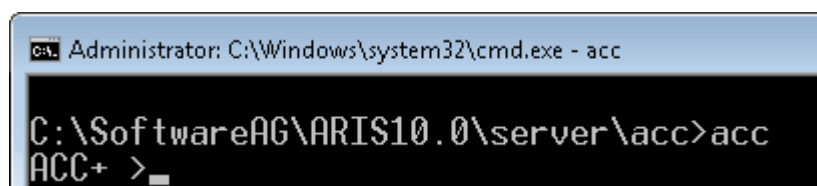


You can now start issuing the well-known commands (e.g., "list" to display all currently configured runnables).

3.5.2 A different way to start ACC

However, you can also start ACC using the script file without specifying any command line options. Open a command prompt in the directory containing the ACC script file (or in any directory when on Linux) and just enter "acc".

This will give you a prompt like this:



Notice that the prompt doesn't say localhost. If you now try running a command (try "list"), ACC will return an error:

```
ACC+ >list
No node specified and no current node set. You need to specify a node on which
this command should be executed, either by prefixing it with "on <nodeName>"
or by setting a current node using the "set current node <nodename>" command.
ACC+ >_
```

This is ACC's way of telling you that it does not know to which node to send a command.

But actually, at this point ACC doesn't even know about *any* nodes. You can check this with the "list nodes" command, which lists all nodes currently known to ACC:

```
ACC+ >list nodes
No nodes.

ACC+ >_
```

Without our usual command line options, ACC does not know to which agent to talk to. Instead of going back to using the command line switches, you can also use a command to make a new agent known to ACC, the **add node** command.

3.5.3 The "add node" command

The basic structure of the command is as follows (you can also get a help text displayed in ACC by typing "help add node"):

```
add node <logicalNodeName> <ipAddressOrHostname> [:<agentPort>]
<agentUsername> <agentPassword>
```

You will recognize the familiar parameters that we normally pass to ACC on the command line, i.e., the hostname or IP address of the machine (or often just "localhost"), the optional agent port, and the username and password of the agent user.

However you will also notice an additional parameter, <logicalNodeName>. Here you can specify a shorthand name with which you can later refer to the node in ACC commands. For starters, let's just add the agent on localhost with the command, and let's name it "localhost" as well (In general, it is not a good idea to use the hostname also as the logical node name... we do this in this section for didactic purposes).

```
add node localhost localhost Clous g3h31m
```

(Note how the optional port was omitted, as the agent installed as part of an ARIS 10 single-node installation will be listening for requests on the default port.)

If successful, the output will look like this and you can use the "list nodes" command to check that indeed your local agent has been registered with ACC:

```
ACC+ >add node localhost localhost Clous g3h31m
Agent at localhost:14000 added successfully with logical node name localhost.
ACC+ >list nodes
One node:
localhost : localhost (14000) OK
ACC+ >_
```

Of course you also again remove a node from the ACC's list of registered nodes, using the "remove node" command. The basic syntax is

```
remove node <logicalNodeName>
```

To remove our node localhost from ACC, type

```
remove node localhost
```

After trying out "remove node", re-add the node by running the previous "add node" command again (in case you didn't know: on Windows, you can scroll through the previous commands using the arrow up and arrow down keys).

3.5.4 Qualifying commands with the node name

So by now, we should be able to use our usual commands, right? Let's try "list" for starters:

```
ACC+ >list
No node specified and no current node set. You need to specify a node on which
this command should be executed, either by prefixing it with "on <nodeName>"
or by setting a current node using the "set current node <nodename>" command.
ACC+ >_
```

Dang! Same error as before! What went wrong?

Well, the error message already points out the problem (and the solution): While it now knows about a single agent, it still insists that you explicitly tell it to which agent to talk to. There are two ways to do this.

3.5.4.1 The "on <nodename>" clause

1. All commands that pertain to a single runnable, or to a group of runnables on a single node, can be prefixed with the "on <logicalNodeName>" clause, in which you specify the logical node name of the node to which to send the command, for example

```
on localhost list
```

Let's try this out now:

```
ACC+ >on localhost list
Node localhost - 14 installed runnables.
ID          State      Version  Type
zoo_m       STOPPED   10.0.0.2 com.aris.runnables.zookeeper-run-prod
postgres_m  STOPPED   10.0.0.2 com.aris.runnables.PostgreSQL-run-prod
cloudsearch_m STOPPED   10.0.0.2 com.aris.cip.y-cloudsearch-run-prod
elastic_m   STOPPED   10.0.0.2 com.aris.runnables.elasticsearch-run-prod
umcadmin_m  STOPPED   10.0.0.2 com.aris.umcadmin.y-umcadmin-run-prod
adsadmin_m  STOPPED   10.0.0.2 com.aris.adsadmin.y-adsadmin-run-prod
abs_m       STOPPED   10.0.0.2 com.aris.modeling.components.y-server-run-prod
octopus_m   STOPPED   10.0.0.2 com.aris.octopus.y-octopus_server-run-prod
app_m       STOPPED   10.0.0.2 com.aris.age.age-run-prod
copernicus_m STOPPED   10.0.0.2 com.aris.copernicus.copernicus-portal-server-run-prod
ecp_m       STOPPED   10.0.0.2 com.aris.runnables.ecp-run-prod
simulation_m STOPPED   10.0.0.2 com.aris.modeling.components.simulationserver.y-simuserver-run-prod
loadbalancer_m STOPPED   10.0.0.2 com.aris.runnables.httpd.httpd-run-prod
dashboarding_m DEACTIVATED 10.0.0.2 com.aris.dashboarding.dashboarding-run-prod
```

As you can see, ACC now knows that it should indeed do a list of all runnables on the node that we registered with the logical node name "localhost". However, having to prefix each command with "on <logicalNodeName>" adds significantly to the amount of typing one has to do, in particular when one is currently only working with a specific node, or even only has a single node registered in ACC; as in our case. For these cases, there is a better option, also mentioned in the error message above.

3.5.4.2 Setting a "current node"

As an alternative to the specifying the node to use with each command, you can tell ACC that all subsequent commands that pertain to a single runnable, or to a group of runnables on a single node, should by default be sent to a specific node, by using the "set current node" command. The "set current node" command has the following syntax:

```
set current node <logicalNodeName>
```

So in our example, where we registered the agent with the logical node name "localhost", the command would be

```
set current node localhost
```

Try this out now:

```
ACC+ >set current node localhost
Node localhost is set as current node.
ACC+ localhost>_
```

Observe how the ACC prompt changed and now shows the logical node name.

Now you can again run your "list" command without having to prefix it with "on <logicalNodeName>".

i After using the "set current node" command on an ACC that was started without any command line options, and where we added a single agent using the add node command and the set this node as the current node, the prompt will now look exactly as it does if the ACC is started by using the -h, -u, and -pwd command line switches on "localhost".

Basically, these command line switches are a shortcut for the commands

```
add node <hostNameGivenOnAccCommandLine> <hostNameGivenOnAccCommandLine>
<userNameGivenOnAccCommandLine> <passwordGivenOnAccCommandLine>
set current node <hostNameGivenOnAccCommandLine>
```


For example when starting ACC from a OS command prompt like this


```
acc.bat -h localhost -u Clous -pwd g3h31m [...]
```

as it is done in the ACC shortcut created by the setup in a single node installation, under the hood the following commands are executed:

```
add node localhost localhost Clous g3h31m
set current node localhost
```

The intention is to make the most common use case - that of using ACC to control a simple single-node installation - as convenient as possible - ideally without the casual ACC user even knowing that ACC can manage multiple nodes.

 While in most cases, the hostname given with "-h" is also used directly as logical node name, there is some conversion being done to turn the hostname into a valid node name. In particular, when an IP address is given with the "-h" switch instead of a hostname, an attempt is made to find a proper hostname for that IP (by using a reverse DNS lookup). If a hostname can't be found, the dots (".") (for IPv4 addresses) or colons (":") (for IPv6 addresses) are replaced with the underscore ("_") character, as neither dots nor colons are allowed in logical node names.

 If you have another ARIS installation at your disposal, you could now add its agent to your ACC by issuing a corresponding "add node" command. However, unless you really know what you are doing, **never add nodes belonging to a different installation to the same ACC**. While those commands for which you explicitly have specify the node and the runnable they apply to will work fine, other commands will implicitly be applied on all runnables (usually of a certain type) across all nodes, which will lead to unforeseeable behavior of these commands.

The ACC has certain indications to find out that nodes belonging to different installations have been added. If the ACC does notice such a situation, you will see a warning like the following:

```
ACC+ localhost>add node n2 [redacted] Clous g3h31m
Agent at [redacted]:14000 added successfully with logical node name n2.
There was a warning:
There are 2 Zookeeper ensembles across the node(s) n2 and localhost.
ACC+ localhost>_
```

Basically ACC noticed that across the nodes you registered there are multiple Zookeeper instances, which however are not connected to form a single Zookeeper cluster (we will learn about clustering Zookeeper in chapter 4) - this is a strong indication that either someone made a mistake when configuring Zookeeper or that nodes were added to ACC that belong to different installations.

3.6 Adding our worker nodes to ACC

With what we learned about the "add node" command and controlling multiple nodes with a single ACC, we could now use multiple "add node" commands to make a single ACC instance know about all our agents we just installed on our worker nodes.

Let's try this out now: start the ACC you installed on your admin machine, without specifying any command line options.

Now enter an add node command for each of your worker nodes. Remember the basic syntax of the "add node" command introduced above:

```
add node <logicalNodeName> <ipAddressOrHostname> [:<agentPort>]
<agentUsername> <agentPassword>
```

Replace <logicalNodeName> with that logical node name you want to assign with which you will later refer to it in ACC commands. You should choose a reasonably short but expressive name. If you do not have clearly defined roles for certain nodes, you can simply number them (I often just use "n1", "n2", ...). You only need to specify the agent port (separated from the hostname by a colon ":") if you want to run the agent on a different port than the default (14000 for ARIS 10, 9001 for ARIS 9).

i Why we introduce logical node names

So why do we actually introduce the logical node names and not just later refer to the nodes by their hostname or IP address?

Well, first hostnames or IP addresses are often long and easily mistyped. But that aside, assigning a logical node name and referencing it by this name in all commands allows you to easily exchange the actual nodes on which you run your commands - all you need is update the hostname or IP addresses in your add node commands. All commands you created (in particular the "template" file we are going to create which contains the configure commands that decide which runnables are put on which node) can then be easily be applied to a different environment (same number of nodes and roughly the same hardware sizing presupposed), i.e., your scripts become portable!

Assuming that the hostnames of our worker nodes are named arisworker1.company.com, arisworker2.company.com, and arisworker3.company.com and we just named those nodes n1, n2, and n3, the corresponding add node commands would be

```
add node n1 arisworker1.company.com Clous g3h31m
add node n2 arisworker2.company.com Clous g3h31m
add node n3 arisworker3.company.com Clous g3h31m
```

If everything works fine, you should now see an output similar to this:

```
C:\aris_distributed>acc
ACC+ >add node n1 arisworker1.company.com Clous g3h31m
Agent at arisworker1.company.com:14000 added successfully with logical node name n1.
ACC+ >add node n2 arisworker2.company.com Clous g3h31m
Agent at arisworker2.company.com:14000 added successfully with logical node name n2.
ACC+ >add node n3 arisworker3.company.com Clous g3h31m
Agent at arisworker3.company.com:14000 added successfully with logical node name n3.
ACC+ >_
```

If you have problems adding a node to ACC, this could be caused by the agent not having been installed properly, or because the agent is not running.

You can now run the "list nodes" command to see all nodes registered in ACC:

```
ACC+ >list nodes
3 nodes:
n1 : arisworker1.company.com (14000) OK
n2 : arisworker2.company.com (14000) OK
n3 : arisworker3.company.com (14000) OK
ACC+ >_
```

3.7 The "on all nodes" clause

In addition to explicitly specifying the node to send a command to with the "on <logicalNodeName>" clause or using the "set current node" command, a number of commands can be executed on all nodes currently known to the ACC, by prefixing them with the "on all nodes" clause.

For example, you can use "on all nodes" with the list command:

```
ACC+ >on all nodes list
Node n1 - NO installed runnables.
Node n2 - NO installed runnables.
Node n3 - NO installed runnables.
ACC+ >_
```

Or you can use the "on all nodes" clause with "get" and "set" commands that allow you to read resp. change the value of an agent configuration parameter. For example, type

```
on all nodes get rest.port
```

to get the value of the "rest.port" agent config parameter (which is the network port the agent's REST interface is listening on for requests and will be 14000 by default):

```
ACC+ >on all nodes get rest.port
On node n1 the value for "rest.port" is : "14000"
On node n2 the value for "rest.port" is : "14000"
On node n3 the value for "rest.port" is : "14000"
ACC+ >_
```

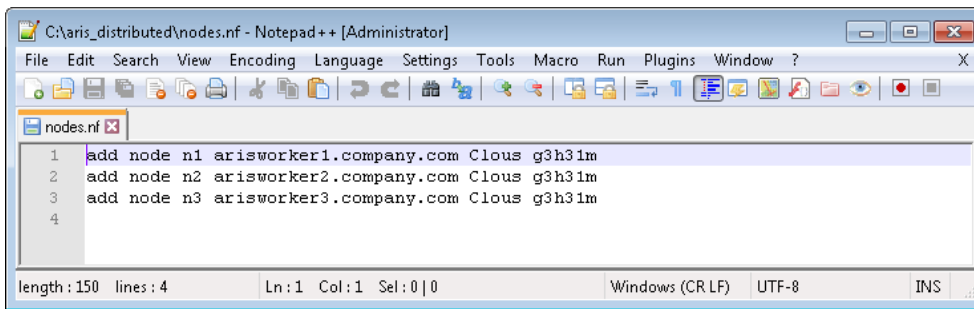
Other commands that can be used with "on all nodes" are the startall, stopall, killall, updateall, and deconfigureall. We will learn a bit about most of these commands in the chapter about administration of a distributed installation.

3.8 Creating a node file

So we can now control our three worker nodes with this single ACC instance. However, the regular ACC console doesn't persist nodes added to it in a previous session, an quite obviously, manually entering the "add node" commands every time we start ACC will be tedious. Fortunately, ACC offers a better way to tell it about the nodes that are part of an installation, a so-called "node file".

A node file is nothing more but a text file containing an "add node" command for each of the nodes that are to be part of your distributed installation. For your node file, you can choose any name you want, but for the remainder of this document we will assume that you named it "nodes.nf" and created it in your project directory. Just add the three add node commands from above into that file and save it:

FIGURE 16



3.9 Using the node file

Now it is time to start an ACC and make it use our node file. You can pass in the path to the node file with the "-n" command line switch, like this:

```
acc -n <pathToNodeFile>
```

For example, assuming you created your node file in the directory "c:\aris_distributed" and named it "nodes.nf", you would run it with

```
acc -n c:\aris_distributed\nodes.nf
```

Try this out now. If everything works, the output should look like this:

```
C:\SoftwareAG\ARIS10_0\server>acc>acc -n c:\aris_distributed\nodes.nf
Executing command "add node n1 arisworker1.company.com Clous g3h31m" in line 1 of file c:\aris_distributed\nodes.nf
Agent at arisworker1.company.com:14000 added successfully with logical node name n1.
Executing command "add node n2 arisworker2.company.com Clous g3h31m" in line 2 of file c:\aris_distributed\nodes.nf
Agent at arisworker2.company.com:14000 added successfully with logical node name n2.
Executing command "add node n3 arisworker3.company.com Clous g3h31m" in line 3 of file c:\aris_distributed\nodes.nf
Agent at arisworker3.company.com:14000 added successfully with logical node name n3.
Validation of the nodes n1, n2 and n3 completed without errors.
ACC+ >_
```

If you added ACC to the PATH environment variable, best is to simply change to your project directory and start ACC from there, so you can leave out the path in the ACC command line, for example:

```
c:\aris_distributed>acc -n nodes.nf
```

If everything works as expected, ACC should start up and add the nodes listed in your node file without giving any errors, and leave you with a prompt without a current node set (i.e., "ACC+ >"):

```
C:\aris_distributed>acc -n nodes.nf
Executing command "add node n1 arisworker1.company.com Clous g3h31m" in line 1 of file c:\aris_distributed\nodes.nf
Agent at arisworker1.company.com:14000 added successfully with logical node name n1.
Executing command "add node n2 arisworker2.company.com Clous g3h31m" in line 2 of file c:\aris_distributed\nodes.nf
Agent at arisworker2.company.com:14000 added successfully with logical node name n2.
Executing command "add node n3 arisworker3.company.com Clous g3h31m" in line 3 of file c:\aris_distributed\nodes.nf
Agent at arisworker3.company.com:14000 added successfully with logical node name n3.
Validation of the nodes n1, n2 and n3 completed without errors.
ACC+ >_
```

From here on, we will assume that you added ACC to your PATH and are working on a shell whose current directory is your project directory.

You can now run the "list nodes" command once more. All nodes should be listed and show an "OK" state like this:

```
ACC+ >list nodes
3 nodes:
n1 : arisworker1.company.com (14000) OK
n2 : arisworker2.company.com (14000) OK
n3 : arisworker3.company.com (14000) OK
ACC+ >_
```

⚠ Important: From now on, whenever we say "start ACC", always pass in the node file with the -n switch, in addition to any other command line switches mentioned.

3.10 Preparing and using a remote repository

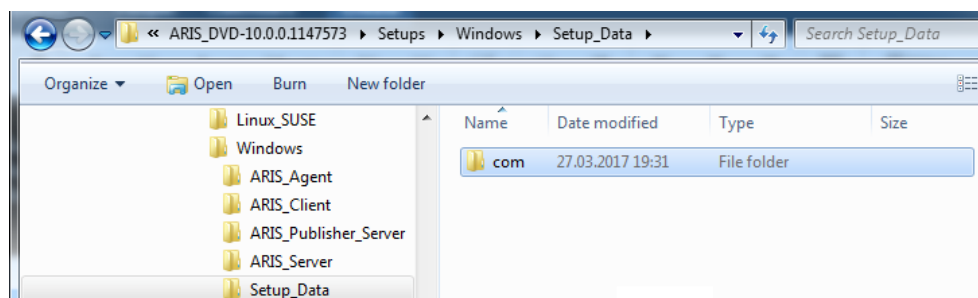
When installing ARIS on a single node with the setup, the setup program itself will start a repository server, i.e., a simple web server that allows the agent to download the runnables as they are found on the ARIS DVD image.

When installing a distributed environment with ACC, there is no setup program that will do this, so we need to provide a remote repository ourselves from which the agents can download the runnables to their local repositories.

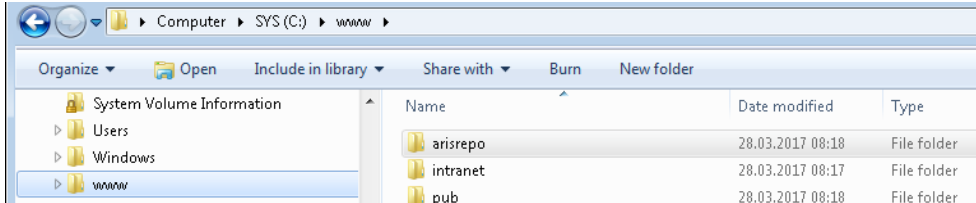
In many customer environments, you will have an FTP or web server available which can be used as remote repository. Assuming that such a server is available to you, first you need to create a subdirectory in that server's document root (i.e., the directory from which the server serves files) that will hold the runnables (e.g., "arisrepo").

Next, copy the contents of the Setup_Data directory, which can be found on the installation DVD at <DVD_Root>\Setups\Windows\Setup_Data (which usually only contains a folder named "com", as you can see in the screenshot below), to the directory you just created:

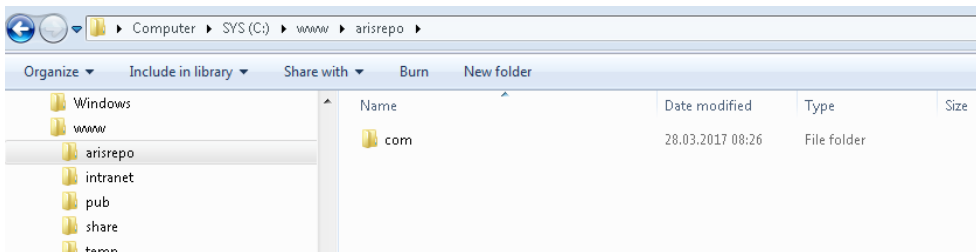
FIGURE 17



For example, assuming that you plan to use an HTTP server that can be reached via the URL "http://internal.company.com" and whose document root is the directory c:\www, and you created the "arisrepo" folder inside this location, you would copy the content of the Setup_Data folder to this folder:

FIGURE 18

So that afterwards the content of the folder c:\www\arisrepo looks like this:

FIGURE 19

If you now access this directory with a browser (in our example, the URL would be "http://internal.company.com/arisrepo") and assuming that directory browsing is enabled in your web server, you should now see a page similar to this:

FIGURE 20

i Installing your own remote repository server

In case you do not have a FTP or web server available that you can use as repository server, you can install and run your own. The simplest approach is to install the Apache HTTPD web server. On Linux machines, you usually can install it using the package manager of your distribution. On Windows, you can download a binary build of Apache from different places on the web.

The following is a quick guide about how to install Apache HTTPD on a Windows machine. It was done using a binary build from apachelounge.com (Apache 2.4.25 win64 VC11), downloaded from <https://www.apachelounge.com/download/VC11/> (Jan 28th, 2017).

⚠ Please note that the suggested source for Apache binaries is an external page not related to Software AG in any way. As always, be careful when downloading and running executables from 3rd party pages (i.e., run a virus scanner on the files first).

⚠ The Apache server set up with the following steps has not been secured against access in any way, so do not run it on a machine exposed to the public internet!

First unzip the ZIP file with the Apache binary distribution to a directory (e.g., somewhere under "Program Files", "c:\Program Files\Apache\httpd" will do nicely, we will refer to this directory as <apacheDir>).

Go to <apacheDir>\apache24\conf and open the file httpd.conf in an editor.

Search for a line like this

```
ServerRoot "c:/Apache24"
```

and change it so that it points to the "Apache24" subdirectory inside the directory to which you unzipped the ZIP file:

```
ServerRoot "C:\Program Files\Apache\httpd\Apache24"
```

Next, choose a directory which is to be your "document root", i.e., the directory where the files served by Apache will be located. As in the example above, we will use "c:\www". First create this directory if it doesn't exist, and the "arisrepo" subdirectory.

In httpd.conf, search for a line like this

```
DocumentRoot "c:/Apache24/htdocs"
```

and change it so that it points to your new document root directory, i.e.

```
DocumentRoot "c:/www"
```

(Note that you can use forward- and backward slashes in the Apache config files on Windows as path separators).

Then add a "Directory" section for the directory to which you copied the contents of the Setup_Data directory (in our example this was "c:\www\arisrepo"), like this

```
<Directory "c:/www/arisrepo">
  AllowOverride None
  Require all granted
  Options +Indexes
</Directory>
```

This will give unrestricted access to the contents of that directory and also enable directory listings, i.e., when accessing a page which corresponds to a directory in the document root and there is no index file, Apache will generate a clickable index, as it was shown in the screenshot above.

Don't forget to save the httpd.conf file.

You can now start Apache via the httpd.exe file located in <apacheDir>\apache24\bin. If everything works fine, you should see a windows cmd popping up (and staying open).

If there is a problem (usually you have a typo in the config file or the Apache's port - by default port 80 - is already occupied on your machine), the window will show briefly and exit immediately, giving you no time to see an error message. In that case, open a command prompt in the <apacheDir>\apache24\bin directory and start Apache by typing httpd.exe on the prompt. You will then be able to see the error message.

Once Apache is started, you should now be able to access it with any browser using a URL

```
http://<hostnameOfMachineRunningApache>:<apachePort>/<arisRepositoryDirectory>
```

Assuming that the machine on which Apache is running is named "internal.company.com" and you didn't change the port from the default 80, any you named the directory under your document root to which you copied the content of the Setup_Data directory "arisrepo", the URL would be

```
http://internal.company.com/arisrepo
```

(Note that you can leave out the port if you use the default port 80 of the http protocol).

This is also the URL that you will specify as repository URL in the next step.

3.11 Making the agents on the worker nodes use the remote repository

To now tell the agents on your worker nodes to use this server as the remote repository. Start the ACC and tell all agents about the repo server typing the command

```
on all nodes set remote.repository.url="<URL_of_Repo_Server>"
```

for example

```
on all nodes set
remote.repository.url="http://internal.company.com/arisrepo"
```

3.12 Get and use a generated.apptypes.cfg file with your ACC

To actually do anything useful with our ACC, we need to provide it with another file (in addition to our node file) - the generated.apptypes.cfg file we already saw when we looked at the "common" ACC command line switches as they are added to the start menu shortcut that is created by the single-server setup. This file contains the configuration of ACC, in particular it contains information that makes it "know" the technical names and versions of the runnables to install. The path to that file can be specified when starting ACC with the "-c" command line switch.

First we need to get this file fitting to the version of ARIS you want to install. You can find this file on the ARIS DVD image at <DVD_Root>\Setups\Windows\ARIS_Server.

Copy this file to your project directory (i.e., right next to your node file).

If you haven't done so already, exit the ACC console now and start it again, this time passing both the node file (with the "-n" command line switch) and the generated.apptypes.cfg file (with the "-c" command line switch), like this:

```
acc -n <pathToNodeFile> -c <pathToGeneratedAppTypesFile>
```

For example, if you opened your command prompt in your project directory (e.g., c:\aris_distributed) and both the node file and generated.apptypes.cfg file are located there as we suggested, you can run the following command to start ACC with both files:

```
acc -n nodes.nf -c generated.apptypes.cfg
```

If everything works OK, you should see the familiar output we already had above:

```
C:\aris_distributed>acc -n nodes.nf -c generated.apptypes.cfg
Executing command "add node n1 win10vm1 Clous g3h31m" in line 1 of file nodes.nf
Agent at win10vm1:14000 added successfully with logical node name n1.
Executing command "add node n2 win10vm2 Clous g3h31m" in line 2 of file nodes.nf
Agent at win10vm2:14000 added successfully with logical node name n2.
Executing command "add node n3 win10vm3 Clous g3h31m" in line 3 of file nodes.nf
Agent at win10vm3:14000 added successfully with logical node name n3.
Validation of the nodes n1, n2 and n3 completed without errors.
ACC+ >_
```

! Important: From now on, whenever we say "start ACC", always pass in **both** the **generated.apptypes.cfg** file with the "-c" command line switch and the node file with the "-n" command line switch.

3.13 Securing access to your agents

Unless you are in a private network environment, everyone knowing the default agent username and password can use his or her own ACC to take control of the freshly installed ARIS agents.

To secure your agents, it is recommended to change the password (and optionally also the username) of each agent to something more secure.

You can do this for all agents currently registered with ACC in one go by using the "on all nodes" prefix for the "set" command, like this:

```
on all nodes set username = <newAgentUserName>
on all nodes set password = <newAgentPassword>
```

Of course you could also have different usernames and passwords on each agent, in which case you would need to use the "on <nodename>" prefix to change the username and password of each node individually:

```
on <nodename> set username = <newAgentUserName>
on <nodename> set password = <newAgentPassword>
```

For example, to change the username to "admin" and the password to "sup3rT0PS3cr37", you would run the commands

```
on all nodes set username = "admin"
on all nodes set password = "sup3rT0PS3cr37"
```

While in ARIS 9, it was necessary to restart the agent for the change to take effect, for ARIS 10 this is no longer needed, the change becomes effective immediately. After making the username and/or password change, your current ACC session will also automatically know about the new username and password and will thus still be able to talk to the agents. However, you will also need to update your nodefile accordingly so that next time you start ACC, it is directly given the right credentials. In our example, you would now change your node file to

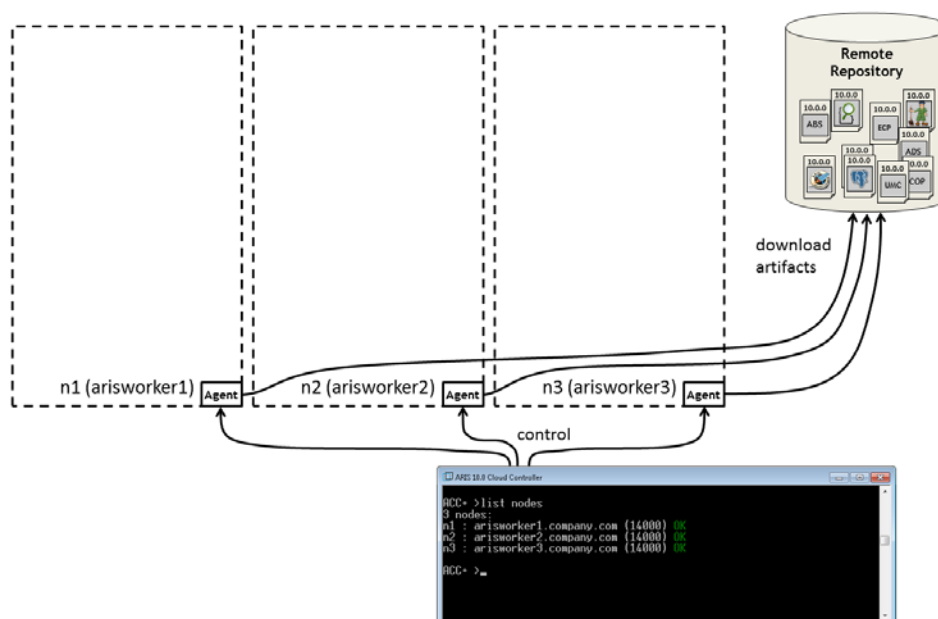
```
add node n1 arisworker1.company.com "admin" "sup3rT0PS3cr37"
add node n2 arisworker2.company.com "admin" "sup3rT0PS3cr37"
add node n3 arisworker3.company.com "admin" "sup3rT0PS3cr37"
```

⚠ Of course the new passwords are now contained in plaintext in your node file - it is therefore important that you restrict access to this file using OS-level file system permissions.

3.14 Our distributed example scenario after performing the preparational steps.

After successfully installing the agents on the three worker nodes used for our example scenario using the installer appropriate for the respective OS, setting up the remote repository and configuring ACC to control all three agents, the current state of our distributed example installation looks like this:

FIGURE 21



3.15 Overview of the remaining steps

Now that we have successfully prepared our agents, our remote repository and our working environment, we can finally get started with the actual installation procedure. The next few chapters will explain each major component in the ARIS architecture in more detail, in particular its function within ARIS and how to install it in a distributed environment.

Before you head on to next chapter, let's first get a brief overview of all the remaining steps:

- 1) You need obtain or create a template file (i.e. a file containing the configure (and similar) commands that will put the runnables on your worker nodes in a reasonable configuration). There are two ways to do this:
 - a) You can obtain an existing template that fits your requirements and the number and hardware sizing of your worker nodes. Alas, there are no official, pre-defined templates provided directly by R&D at the time of this writing.
However, Software AG's consulting unit should be able to provide you with templates for some common scenarios.
 - b) You can learn about how ARIS works and then use that knowledge to create your own template, where you decide on which worker nodes you want to put which runnables, if and how you make some or all runnables highly available etc. This step requires some in-depth knowledge about the ARIS architecture and is covered in detail in the remainder of this guide.

- 2) You can then start ACC and pass in the template file with the "-f" switch (in addition to generated.apptypes.cfg and nodefile!)

If everything was done correctly, the configure commands should now be executed one after the other.

- 3) Once all configure commands have completed successfully, again start ACC (this time passing only the nodefile and the generated.apptypes.cfg file) and type "on all nodes startall"
- 4) After everything has started up, you can access ARIS via any of the worker nodes on which an ARIS **loadbalancer** runnable was configured. If your workers are Windows machines, the ARIS load balancer (ARIS LB) will use the default ports for HTTP (port 80) and HTTPS (port 443) (unless of course the template you used chose different ports).

For example, if you have three worker nodes whose hostnames are arisworker1.company.com, arisworker2.company.com, and arisworker3.company.com, and you have a **loadbalancer** runnable configured and running on the first and third machine, you could access ARIS with a browser using the URL <http://arisworker1.company.com> or <http://arisworker3.company.com>.

If you decide to follow through this document and learn how to create a template file, you should now prepare a new empty text file in your project directory and name it accordingly (e.g., aris_3_nodes.tp).

All the commands we use in the next chapters to set up the distributed example environment can of course be run interactively, i.e., by directly typing them or copying and pasting them into the ACC shell as we go. As an alternative (or in in addition), you can also put these commands into your template file. This allows you to later reuse this template in other environments as a kind of "canned configuration".

4 Zookeeper clustering and the Zookeeper manager

4.1 Motivation - The role of Zookeeper in the ARIS architecture

This chapter explains how to properly configure the Zookeeper runnable in a distributed scenario.

In the ARIS 9/10 architecture, Zookeeper has several essential roles:

Service Discovery - Any ARIS component (the "runnables") registers itself in Zookeeper with essential information that allows other components to find and access it. The detailed information depends on the type of component, but the minimum information is the hostname (or IP address) and port on which the component can be accessed. Any services used by ARIS that are not provided by a runnable (so-called "external services") also have to be registered in Zookeeper using the ACC service registration functionality ("register external service", "update external service", "list (external) services" and "deregister external service" commands).

Coordination - Aside from the service discovery functionality, Zookeeper is used to solve a number of common problems in distributed software environments, in particular, but not limited to, the following:

- **Leader election:** Often, certain operations (e.g., background jobs, schema creation, etc.) should be done only exactly once in the system (or exactly once for each tenant). If you have multiple instances of an application in your installation (i.e., "horizontal scale out"), some way is needed for the instances to decide which of them is responsible for such an operation, i.e., which one becomes the so-called "leader". Further, it has to be assured that if the leader instance fails, a new leader is elected among the remaining instances (i.e., to make sure that the leader "seat" is never vacant for long). For, this the leader election⁷ pattern is used. Refer to the Zookeeper manual⁸ for details on how leader election can be handled with Zookeeper.
- **Distributed queue:** For some coordination tasks, messages need to be sent reliably between instances. One option would be to use a dedicated, persistent, transactional message queue in our infrastructure. However, to keep the footprint small we don't have such a component as a mandatory element of the ARIS infrastructure so far. For basic coordination tasks of the ARIS components with a rather low message volume (i.e., not for messages that transport any user generated events or content!), a distributed queue can be implemented on top of Zookeeper⁹.
- **Locking:** Some (usually long-running) operations need the guarantee that no other operation interferes with them. For example, a UMC LDAP import should at any given time only be running once in the system for each tenant. To prevent such operations from occurring concurrently, locking is used, which can also be handled through Zookeeper¹⁰.

Due to these essential functionalities handled through Zookeeper, if an ARIS installation is to be made highly available, the Zookeeper service itself has to be made highly available, as any outage of this service will lead to the system becoming unresponsive and/or crashing (almost) immediately

⁷ Wikipedia: Leader Election https://en.wikipedia.org/wiki/Leader_election

⁸ Zookeeper Recipes and Solutions: Leader Election https://zookeeper.apache.org/doc/trunk/recipes.html#sc_leaderElection

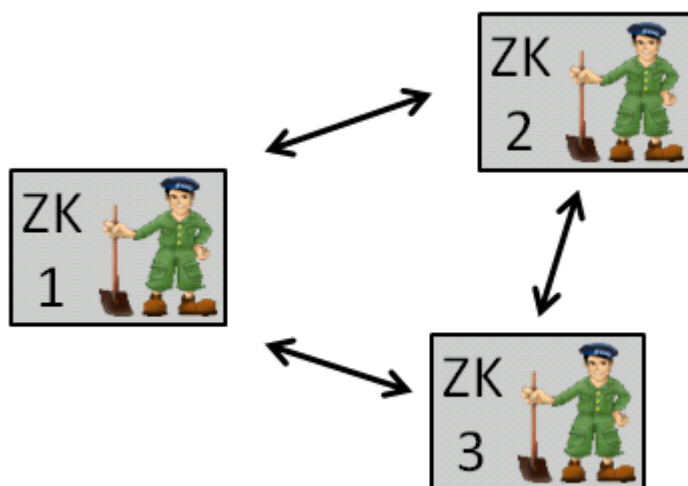
⁹ Zookeeper Recipes and Solutions: Queues https://zookeeper.apache.org/doc/trunk/recipes.html#sc_recipes_Queues

¹⁰ Zookeeper Recipes and Solutions: Locks https://zookeeper.apache.org/doc/trunk/recipes.html#sc_recipes_Locks

4.2 Zookeeper ensembles & failover

To make Zookeeper highly available, several Zookeeper instances can be clustered to form a so-called ensemble. In an ensemble, *all member instances will hold the full set of data* that is in the repository and all members of a Zookeeper ensemble keep the content of their repositories in sync with each other. For this to work, each Zookeeper instance needs to know all other ensemble members, as indicated in figure 13 below.

FIGURE 22: A THREE-NODE ZOOKEEPER ENSEMBLE



Details on how this knowledge is established are given in the next section.

So whenever a write operation occurs, the ensemble members communicate with each other to assure that all agree on the outcome of the operation. The coordination of this process is in the responsibility of the Zookeeper master instance. At any given time there must be exactly one instance appointed as master. So when you compare this to our discussion of replica consistency in chapter 2, Zookeeper works with an entire instance being named the master, instead of having masters for individual data items. While this means that all writes have to be coordinated by this instance, this is not a problem for the usual scenarios Zookeeper is used for. The master is chosen in a voting process, which can only occur if a true majority of the Zookeeper instances is available and can achieve a so-called quorum (>50%). Without a quorum, there will be no master, and without a master, no Zookeeper operations can happen. In other words, any Zookeeper operation (in particular writes, but also reads) can only happen if more than half (>50%, not $\geq 50\%$!) of the ensemble members are working OK and can communicate with each other and thus can "cast their votes" in the master election process. This means that a Zookeeper ensemble of size n can tolerate $\text{roundDown}((n-1)/2)$ instances failing before the entire ensemble becomes unavailable, as indicated by the following table:

Total instances	1	2	3	4	5	6	7	8	9
Tolerated failed instances	0	0	1	1	2	2	3	3	4

Note that in particular adding *a second instance* to your ensemble *will not improve the failure tolerance of your Zookeeper ensemble* (on the contrary - the likelihood that one instance out of two fails is higher than that one out of one fails!). Furthermore, note that having an even number of instances will not increase your level of failure tolerance over having an ensemble with one instance less.

So a first thought might be that having as many Zookeeper instances as possible is the way to go?

Well, *of course not!*

Reason is that the communication overhead for the master election procedure and in particular for the synchronization of the repository content during write operations increases with each additional instance (as it has to be made sure that all active instances have the same content in their repo). In general, for all but the most ambitious ARIS deployments, a Zookeeper ensemble consisting of three instances is perfectly fine. As a general rule of thumb the Zookeeper documentation states that when one puts more than seven instances into an ensemble, the performance of the ensemble will actually decrease compared to a smaller one.

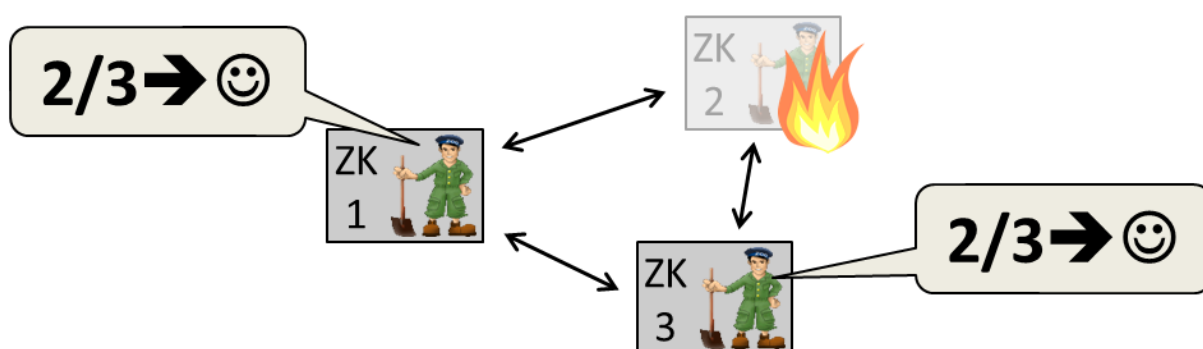
i Zookeeper observer instances

Zookeeper also allows to add so-called observer instances that do not participate in the master election process, but can otherwise be used by clients just as regular instances. This allows having a larger number of Zookeeper instances without increasing coordination overhead in particular to give clients that are placed remotely from the main data center (i.e., where the voting ensemble members are located) a local instance with which to communicate. While using observer instances can be beneficial in certain deployment scenarios, for ARIS deployments we never found a good reason to use them.

4.3 Failover scenarios

Let's now have a look at a few failure scenarios. Assume that in our three node ensemble from the figure above, instance 2 suddenly fails (e.g., the Zookeeper runnable crashes, the entire machine goes down, the data center it is located in goes up in flames etc.), as indicated in figure 23 below:

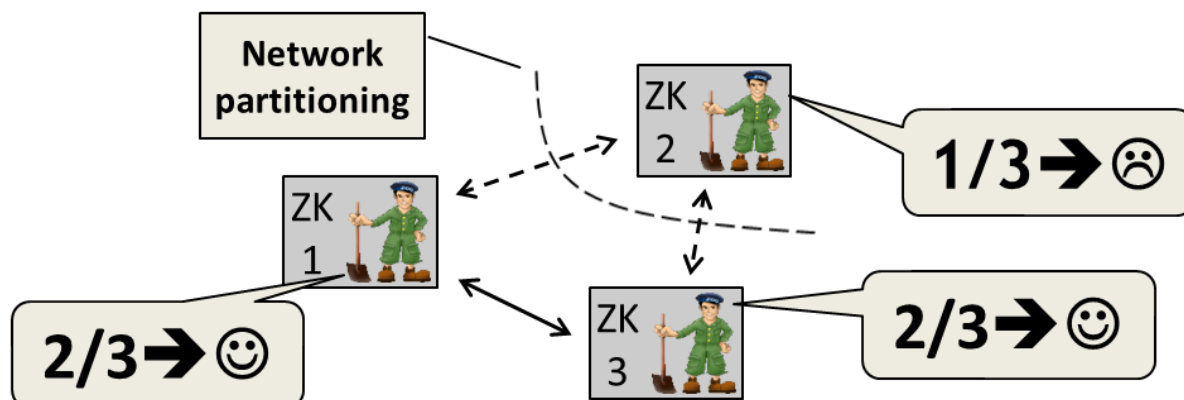
FIGURE 23: OUTAGE OF ONE INSTANCE IN A THREE-NODE ZOOKEEPER ENSEMBLE



Note that the remaining instances 1 and 3 together still have 2/3 of the votes, can therefore successfully elect a master instance among them and continue working normally (i.e. serve read *and* write requests). As soon as instance 2 becomes available again, it will rejoin the ensemble and sync its content to that of instances 1 and 3.

Let' now consider the apparently similar scenario shown in figure 24 below.

FIGURE 24: NETWORK OUTAGE (PARTITIONING) IN A THREE-NODE ZOOKEEPER ENSEMBLE



Here, instance 2 itself continues working fine, however, due to a network partitioning, it can no longer communicate with the other instances. The basic behavior is the same as before, the larger partition with instances 1 and 3 still has a quorum, so the instances can elect a master among them and can continue serving requests. Instance 2 in the smaller partition, while itself fully functional, cannot collect enough votes to elect a master, and does not serve any requests. That way, the dreaded split brain problem is avoided: assume for a moment that Zookeeper did not require a quorum to allow read and write operations. In that case, the smaller partition (consisting of only instance 2) could still continue working and could serve clients that can still reach it (because they are on the same side of the network partition). That way, the content of the repositories of the instances in the two partitions could be modified independently and thus diverge. While this wouldn't be noticed as long as the network partitioning persists, if the network partition were later resolved, it would not be clear which instances hold the "true" state of the repository (instances 1 and 3, or instance 2) and the repositories could not be consolidated into one consistent state, a situation which has to be avoided at all cost, as it could not be resolved automatically.

4.4 Setting up a Zookeeper cluster (a.k.a. "ensemble") manually


To make each ensemble member know all other instances, each instance has to be given a list of the hostname or IP address of *all* ensemble participants and their respective two server ports (named server port A and B), using the indexed, multi-valued "server" config parameter. The values of this parameter have the format

```
<hostname> : <serverPortA> : <serverPortB>
```

Being "indexed" means that instead of supplying multiple values for that parameter simply as a comma-separated list (e.g., $x = a,b,c$) you specify the *individual values as individual configure parameters*, by appending ".<index>" to the parameter key (e.g.: $x.1 = a$ $x.2 = b$ $x.3 = c$).

For example, in a three-node ensemble, where Zookeeper instances are on the nodes arisworker1.company.com, arisworker2.company.com, and arisworker3.company.com and are using the standard ports for server port A and B (for ARIS 10 these are the ports 14285 and 14290, respectively), you would need to set the configure parameters for each instance:

```
server.1=arisworker1.company.com:14285:14290
server.2=arisworker2.company.com:14285:14290
server.3=arisworker3.company.com:14285:14290
```

 In previous versions (<10), the server parameter could indeed be specified using a single multi-valued parameter, like this:

```
server="arisworker1.company.com:14285:14290",
"arisworker2.company.com:14285:14290", "arisworker3.company.com:14285:14290"
```

While multi-valued parameters still exist in ARIS 10, for technical reasons they cannot be used anymore for setting the server parameter(s) of a Zookeeper ensemble.

In addition to providing the proper values of the server.<index> parameters, each instance needs to be told which member in that list the instance itself is, using the "myid" parameter.

For example, for the instance located on node aris1, you would set myid to 1. Note that this way, you also implicitly tell each instance which server ports it should open.

4.4.1 Example

Assume you want to have a Zookeeper cluster spread across three nodes whose hostnames are arisworker1.company.com, arisworker2.company.com, and arisworker3.company.com and which you registered in ACC as n1, n2, and n3, respectively.

You would configure your three Zookeeper instances like this (note the use of a "\" as the last character in a line to make the command continue on the next line).

```
on n1 configure zoo_m zoo \
  myid=1 \
  server.1=arisworker1.company.com:14285:14290 \
  server.2=arisworker2.company.com:14285:14290 \
  server.3=arisworker3.company.com:14285:14290

on n2 configure zoo_m zoo \
  myid=2 \
  server.1=arisworker1.company.com:14285:14290 \
  server.2=arisworker2.company.com:14285:14290 \
  server.3=arisworker3.company.com:14285:14290

on n3 configure zoo_m zoo \
  myid=3 \
  server.1=arisworker1.company.com:14285:14290 \
  server.2=arisworker2.company.com:14285:14290 \
  server.3=arisworker3.company.com:14285:14290
```

Observe that the values for the three `server.<index>` configuration parameters are the same across all instances. Also observe that only the value of the `myid` parameter differs, telling each instance which "position" among the list of servers it holds.

The obvious disadvantage of this manual configuration is that a) you can get it wrong and b) you cannot simply use the same provisioning template again in a different environment, since the `server.<index>` configuration parameters now contain values (the hostnames or IP addresses of the machines!) that depend on the actual deployment environment.

4.5 Setting up a Zookeeper cluster with ACC's built-in Zookeeper manager

To make setting up a Zookeeper ensemble easier and keep templates portable, the ACC contains the "Zookeeper manager" component, which offers a set of high-level commands to simply initial set up and later changes of an ensemble.

4.5.1 Configuration

The only ACC configuration parameter required by the Zookeeper manager is `"zkmgr.zookeeper.app.type"`, which has to be set to the name of the app type that will be used for configuring Zookeeper instances. Usually, this will be the app type `"zoo_m"`.

Note that if an app type with the name given for this parameter is not registered in ACC, Zookeeper manager will not work.

4.5.2 Commands

The following section describes all ACC commands that relate to ACC's Zookeeper manager component. Note that in general the command syntax allows to use the keywords `"zookeeper"` or `"zk"` interchangeably.

4.5.2.1 Add zookeeper instance

The `"add zookeeper instance"` command stages the addition of a new Zookeeper instance to the installation's Zookeeper ensemble on the specified node (or to the current node, if no node name is specified). It allows to optionally specify the `instanceld` to be used for the new Zookeeper instance. If none is specified, the cloud controller will automatically choose an `instanceld`. Note that the change is not submitted directly when executing the command, but it is only "staged". To apply the change, one needs to use the separate `"commit zookeeper changes"` command.

Syntax:

```
<nodeSpec>? "add" ("zk" | "zookeeper") ("master" | "observer")? "instance"? ("using" "instanceld"?  
<instanceld>)?
```

Examples:

```
add zk
```

This stages the addition of a new Zookeeper instance on the current node.

```
on n1 add zookeeper
```

And that stages the addition of a new Zookeeper instance on node `n1`.

4.5.2.2 Remove zookeeper instance

The "remove zookeeper instance" command stages the removal of a Zookeeper instance from the installation's Zookeeper ensemble. The Zookeeper instance to be removed can be identified either by specifying the node name and the instancelid of the Zookeeper instance or by specifying the myid of the Zookeeper instance (i.e., its sequential number among all instances of the ensemble).

Syntax:

```
<nodeSpec>? "remove" ("zk" | "zookeeper") "instance"? ("instancelid"? <instancelid>)?  
| "remove" ("zk" | "zookeeper") "instance"? ("myid"? <integer>)
```

Examples:

```
remove zk myid 3
```

This stages the removal of the Zookeeper with myid 3 from the ensemble.

```
on n2 remove zookeeper zoo2
```

And this removes the Zookeeper instance with instance id zoo2 on node n2.

4.5.2.3 Commit zookeeper changes

The "commit zookeeper changes" command applies all currently "staged" or "pending" changes of the Zookeeper ensemble.

Syntax:

```
"commit" "pending"? ("zk" | "zookeeper") ("change" | "changes") "force"?
```

Example:

```
commit zookeeper changes
```

4.5.2.4 Discard zookeeper changes

The "discard zookeeper changes" command discards all currently pending changes of the Zookeeper ensemble.

Syntax:

```
("discard" | "reset") "pending"? ("zk" | "zookeeper") ("change" | "changes")
```

Example:

```
discard pending zk changes
```

4.5.2.5 List zookeeper instances

The "list zookeeper instances" command displays all Zookeeper instances of the ensemble with their configuration.

Syntax:

```
"list" ("zk" | "zookeeper") "instances"
```

Example:

```
list zk instances
```

4.5.2.6 Show zookeeper changes

The "show zookeeper changes" command displays all currently pending changes of the Zookeeper ensemble that will be applied when using the "commit zookeeper changes" command.

Syntax:

```
"show" "pending"? ("zk" | "zookeeper") ("change" | "changes")
```

Example:

```
show pending zookeeper changes
```

4.5.2.7 validate zookeeper changes

The "validate zookeeper changes" command performs a validation of all currently pending changes of the Zookeeper ensemble, i.e., it checks whether the changes will result in a working ensemble.

Syntax:

```
"validate" "pending"? ("zk" | "zookeeper") ("change" | "changes")
```

Example:

```
validate zk changes
```

4.5.2.8 validate zookeeper ensemble

The "validate zookeeper ensemble" command performs a validation of the current state of all already configured Zookeeper instances. In particular, it checks if all instances form a single consistent ensemble. Any currently "staged" or "pending" changes are not considered by this command.

Syntax:

```
"validate" "pending"? ("zk" | "zookeeper") ("change" | "changes")
```

Example:

```
validate zookeeper ensemble
```

4.5.3 Usage example

If you now want to set up the three-node ensemble from our example above using Zookeeper manager instead of doing it manually, you would now issue an "add zookeeper instance" command for each instance, specifying the node on which it should be placed. Remember that this will not yet make any changes to the ensemble, but only "stages" the intended changes. To actually apply the changes, you then have to issue the "commit zookeeper changes" command. So, assuming that your three worker nodes were registered as n1, n2, and n3, respectively, the following four commands would automatically set up a working Zookeeper ensemble (and would also start all three instances).

```
on n1 add zk
on n2 add zk
on n3 add zk
commit zk changes
```


4.6 Bootstrapping the Zookeeper connection

In this section, we learned so far that Zookeeper is an essential element of the ARIS architecture and that it is - among other things - used to allow runnables to find each other.

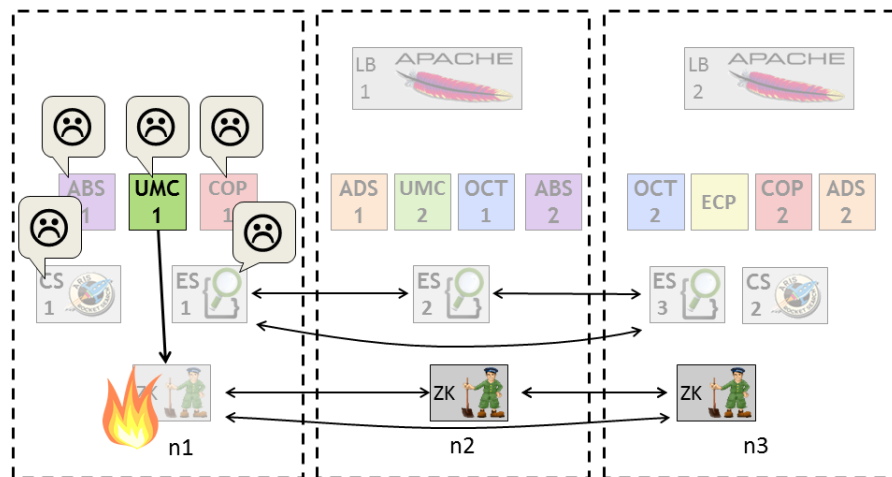
The obvious question you should have been asking yourself a few times up to now is: "How do runnables find Zookeeper?"

4.6.1 How do runnables find Zookeeper?

To keep things simple for single-node installations, by default, a runnable will try to access Zookeeper on the local machine on the default port of the Zookeeper runnable (which is 2181 for ARIS 9.x and 14281 for ARIS 10.x, respectively).

This default mechanism will also be enough for small distributed scenarios, where each node will hold an instance of the Zookeeper ensemble: Since the first Zookeeper instance on a node will be listening for client requests on the default port, all runnables on that node will find their local instance. However, if we leave things like that, there is already a bit of a restriction here: if a Zookeeper instance fails, all runnables on the same node will also lose their Zookeeper connection (even if the other Zookeeper instances on other nodes still have a quorum and work fine) and will eventually cease to work - in other words, the local Zookeeper instance becomes a single point of failure for all runnables on that node! figure 25 below illustrates this situation. Here we have a simple distributed scenario consisting of three nodes, each node holding a single Zookeeper instance of a three node ensemble plus a selection of our other runnables.

FIGURE 25

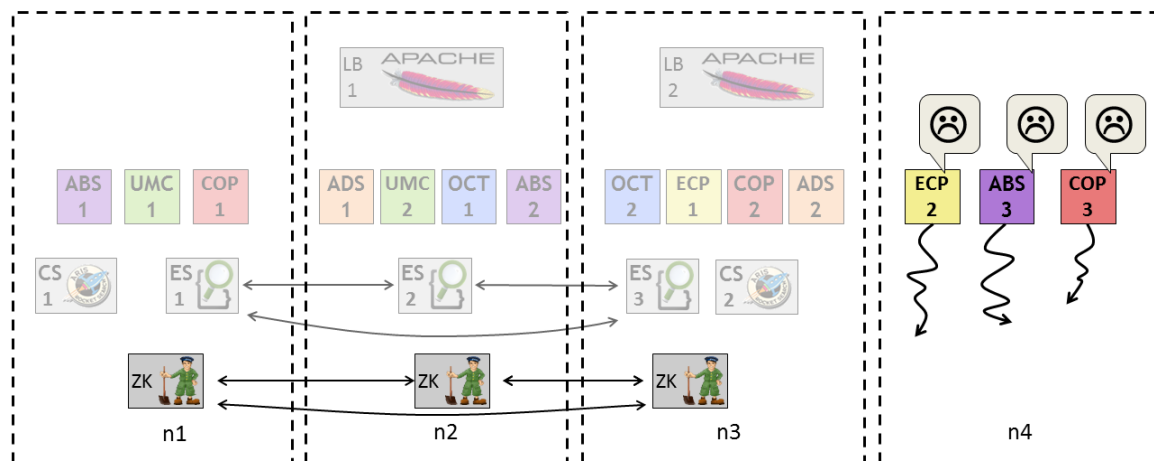


In the distributed scenario sketched here, the Zookeeper instance on node n1 fails, leaving all runnables on the same node without a Zookeeper connection, despite the ensemble as a whole still being available, as the Zookeeper instances on nodes n2 and n3 still have a quorum.

For small distributed scenarios with not the highest requirements regarding high availability, this is in practice not a strong limitation, as Zookeeper is usually quite stable and even if it does crash, the agent will restart it within a few seconds. Most of our applications tolerate a short Zookeeper outage quite well (they will usually just block until it is back, i.e., a request to the app will “hang” for a moment). Further, when aiming for a highly available ARIS installation, you anyway have to tolerate failure of an entire node (or loss of network connectivity etc.).

However, if this problem is not addressed, the installation as a whole will be less robust than it could be. It is obviously preferable to make each runnable know about all Zookeeper instances, so that they can switch to another one if their local instance fails, and, of course, for more complex deployments where we do not have a Zookeeper instance on each node. Consider for example the scenario sketched below:

FIGURE 26: A NODE WITHOUT A LOCAL ZOOKEEPER INSTANCE



Here, we added a fourth node to our previous example scenario to have additional capacity for scaling out some of our runnables so that we can handle a higher request load. However, since we learned that having an even number of Zookeeper instances is not useful, we decided to not add a Zookeeper instance to this node as well. By default, however, this would leave our new runnables without a Zookeeper connection, and they would not be able to work and actually handle user requests.

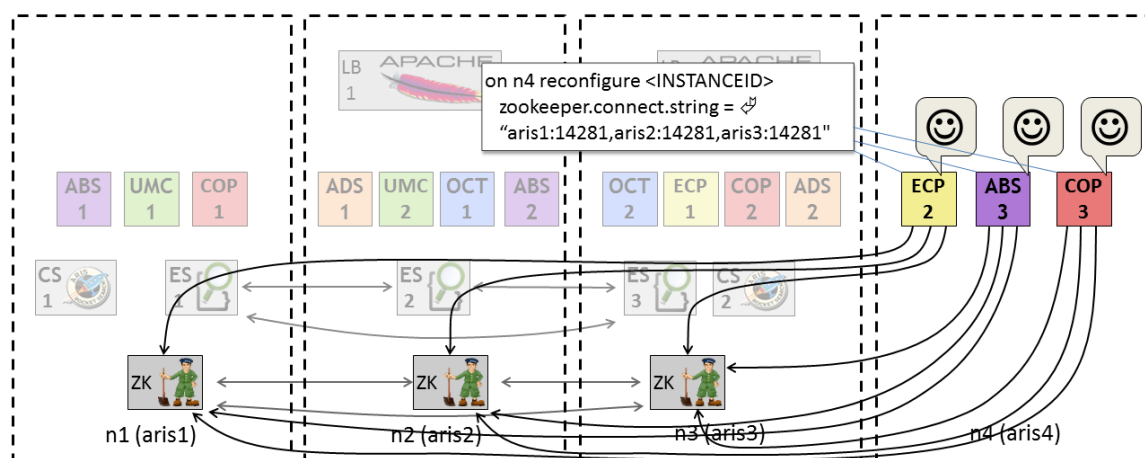
4.6.2 The zookeeper.connect.string configure parameter

To resolve this problem, each runnable can be told where to find Zookeeper instances with the `zookeeper.connect.string` parameter. The `zookeeper.connect.string` contains a comma-separated list of the hostnames and client(!) ports of Zookeeper instances. Assuming that our four nodes n1 to n4 above had the hostnames `aris1` to `aris4`, and assuming that our three Zookeeper instances on nodes n1 to n3 were left at the default clientPort (14281), the proper value of the Zookeeper connect string in this installation would be

```
aris1:14281,aris2:14281,aris3:14281
```

We can now set this value with a reconfigure command on our three additional runnables on node n4 as show here:

FIGURE 27: SETTING THE ZOOKEEPER.CONNECT.STRING PARAMETER



(and to improve failure tolerance, i.e., to not have the local Zookeeper instances on nodes n1 to n3 as a single point of failure for all runnables on that node, we should do the same for all runnables on all nodes, as discussed above).

While using the `zookeeper.connect.string` parameter allows to remedy the problem of Zookeeper connectivity, this is quite obviously not very convenient and also error-prone, as one has to maintain this parameter on all runnables. And again, the value of this parameter - while the same for all runnables of an installation - depends on the actual deployment environment, as the `zookeeper.connect.string` parameter contains the hostnames (and not the logical node names) of the nodes involved and thus would need to be changed for each deployment environment, again making templates not portable.

4.6.3 Automatic setting of `zookeeper.connect.string` (≥ 9.8.2)

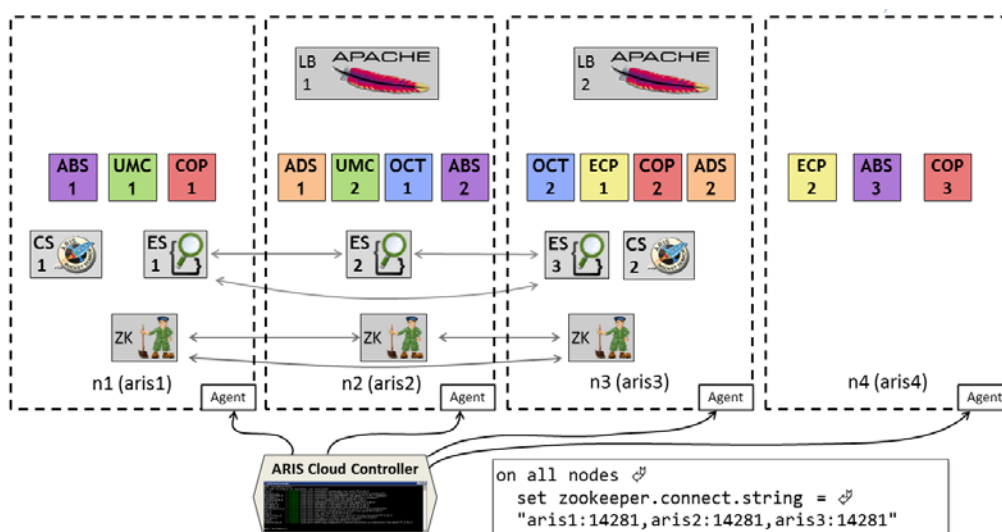
To resolve this problem, with version 9.8.2, a feature was introduced to make this more convenient: An agent configuration parameter "`zookeeper.connect.string`" was added. First, the value of this parameter is used by the agent itself to create its own Zookeeper connection (needed for the agent to be able to perform certain commands). Second, the value of this parameter is passed to the runnables (and the runnables' plugins) on startup, making it no longer necessary to manually set the `zookeeper.connect.string` configuration parameter on the runnables.

i Some background info about passing the `zookeeper.connect.string` to the runnables on startup

The agent will pass its own `zookeeper.connect.string` to the runnables when it starts them, by setting a Java system property, which obviously only works for Java runnables. All Java-based runnables with ZKC inside will then be able to use this `zookeeper.connect.string` value to establish their Zookeeper connection, but will only do so, if the `zookeeper.connect.string` hasn't been set as a regular config parameter for this runnable (or has merely been set to the empty string or has been set to the default "`localhost:14281`"). For the non-Java-based runnables (i.e., the Apache-HTTPD-based load balancer, Postgres, and - for ARIS 9.x only - CouchDB (for ARIS 9.x)) and for Java-based 3rd-party runnables that do not have ZKC inside (Elasticsearch), the registration of the runnable in Zookeeper is done by the runnables' associated agent plugin. These plugins are also given the agent's `zookeeper.connect.string` to open their Zookeeper connection.

Further, under certain conditions, the ACC can now automatically set the `zookeeper.connect.string` parameter of all agents that are registered with it, as illustrated in the figure below:

FIGURE 28: AUTOMATIC SETTING OF THE ZOOKEEPER.CONNECT.STRING BY THE ACC



In particular, there must be only exactly one Zookeeper ensemble across all nodes, and this ensemble must be valid and consistent. Further, the nodes on which the Zookeeper instances are running must be registered with a "real" hostname or IP address, not with a loopback address (e.g., "localhost", "127.0.0.1" etc.).

If these conditions hold, the ACC will update the `zookeeper.connect.string` of the agents

- whenever a new node is added (and the conditions still hold)
- whenever an ensemble change is performed via Zookeeper manager

The automatic updating is enabled by default, but can be disabled by setting the ACC config parameter `zkmgr.maintain.agent.zookeeper.connect.string` to `false` (remember that like all ACC settings, this setting is not persisted and has to be re-done when restarting ACC).

4.7 Zookeeper in our example scenario

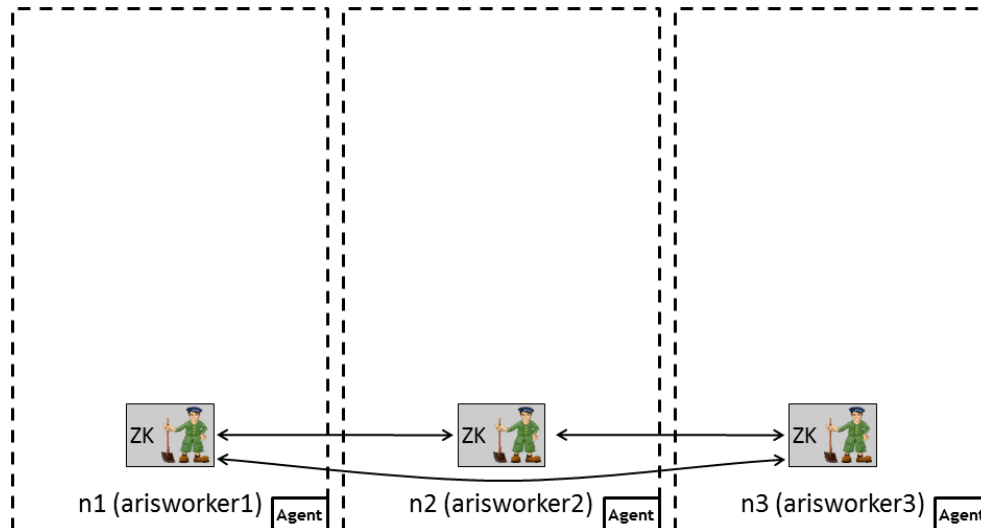
In our example scenario, where we have three nodes n1 to n3, we can now just use the basic Zookeeper manager "add zk" and "commit zk changes" commands introduced above to add one Zookeeper instance on each node:

```
on n1 add zk
on n2 add zk
on n3 add zk
commit zk changes
```

You can now run these commands directly in your ACC shell and/or you can add them to your template file.

Assuming that everything works OK, you will now have three already started Zookeeper instances that are correctly configured to form a working ensemble that can tolerate the failure of any single instance, as shown in the figure below:

FIGURE 29: THE ZOOKEEPER ENSEMBLE IN OUR EXAMPLE SCENARIO



5 Elasticsearch clustering and the Elasticsearch manager

5.1 Motivation - The role of Elasticsearch

In the ARIS 9/10 architecture Elasticsearch serves as the backend for several applications: It is used as search index for collaboration and document content. Further, it is the backend used to store all UMC data (users, user groups, licenses, license assignments,...) and also serves as backend for all currently valid user sessions. In particular the essential role it has for UMC makes Elasticsearch a service whose availability is pivotal for the availability of ARIS as a whole. So if one wants to set up a distributed, highly available ARIS installation, one has to set up Elasticsearch in a highly available way.

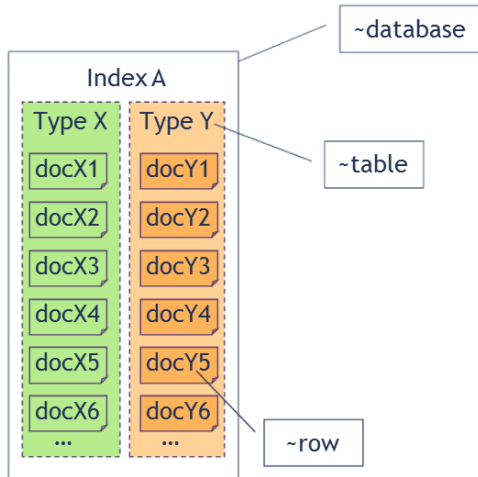
5.2 Elasticsearch terminology and clustering

Elasticsearch is first and foremost a scalable full-text search engine. It can store arbitrary documents in JSON format, which it can access in real-time, making it effectively also a NoSQL backend.

The basic unit of persistence in Elasticsearch is the **document**, which can be (very...) roughly compared to a tuple in a relational database system. A document consists of several **fields**, which can be roughly compared to a column in a relation database system. However, unlike a column in a relational database, a field can have complex, i.e., nested and repeating elements. A **type** is a category of documents which have similar properties and can be roughly compared to a relation or table in a relational database, however, different from an RDBMS, documents of one type need not have the same fields. An Elasticsearch **index** is now a logical collection of documents (of possibly different types), the closest analogy being a logical database in an RDBMS.

Figure 30 below gives an overview of these basic concepts: Index A consists of multiple documents assigned to two different types X and Y.

FIGURE 30: ELASTICSEARCH TERMINOLOGY: DOCUMENTS, TYPES, INDEXES, AND SHARDS

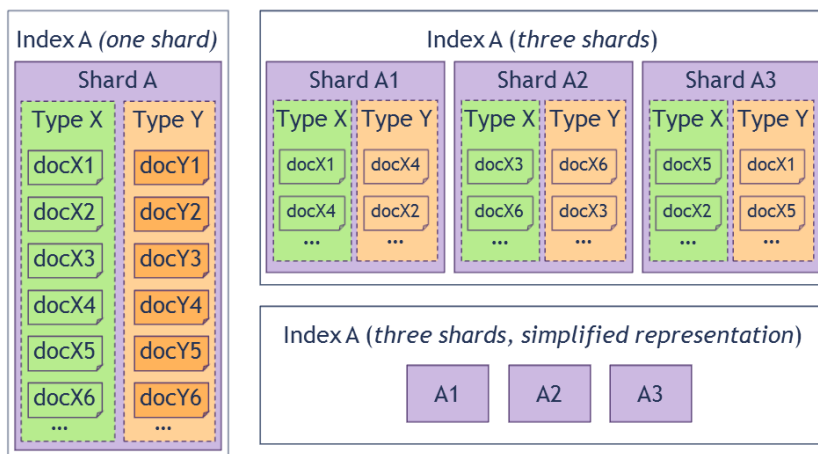


Aside from these logical data structures, Elasticsearch also gives some control over the physical storage of data. In particular, it allows to split a single index into several physical segments, the so-called **shards**. A shard is the unit of distribution and replication in an Elasticsearch cluster, i.e., the shards of an index can possibly be located on different Elasticsearch instances.

! Please note that in official Elasticsearch terminology what we call an "instance" is actually called a "node". We deliberately chose to not use Elasticsearch terminology here to avoid confusion between our meaning of nodes, i.e., as a machine on which an agent is running and which can therefore hold ARIS runnables.

Figure 31 below shows how our index from figure 1, and how it could be split up into shards.

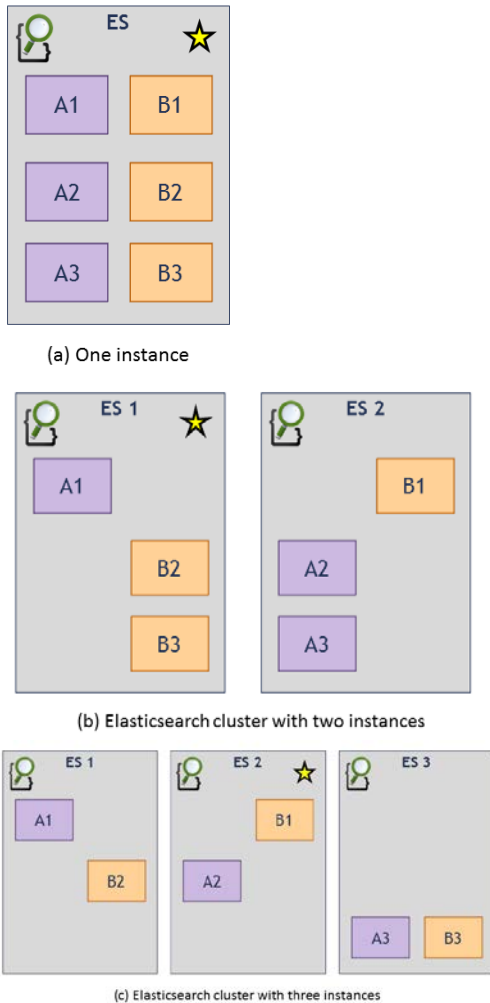
FIGURE 31: AN ELASTICSEARCH INDEX CONSISTING OF THREE SHARDS



To the left, the index is kept in a single shard. The upper right shows the same index, but this time split into three shards, each containing some of the documents. The lower right finally shows the same index, again with three shards, but in a simplified presentation, where we do not show the types and documents, but only the shards themselves.

The shards of an index can be distributed among the Elasticsearch instances of a cluster, as it is illustrated in figure 32 below.

FIGURE 32: DISTRIBUTION OF THE SHARDS OF AN ELASTICSEARCH INDEX ACROSS A CLUSTER



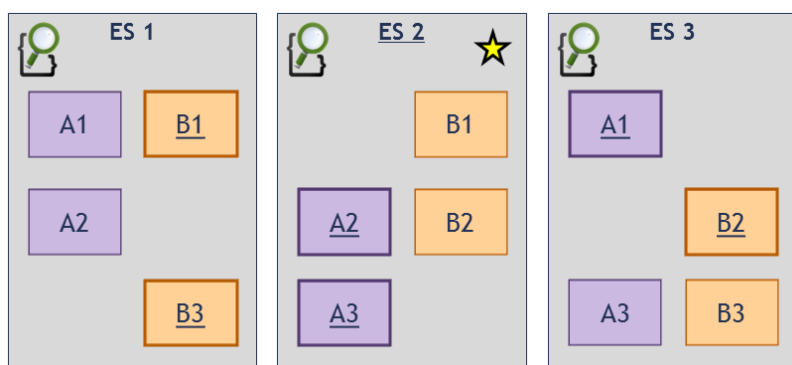
In figure 32 (a), we have a single Elasticsearch instance, with two indexes A and B with three shards each. With only one instance available, naturally all shards of all indexes are held by that instance. figure 32 (b) shows the same logical indexes, but now in an Elasticsearch cluster consisting of two instances, and a possible distribution of the shards among these instances. figure 32 (c) finally shows the same two indexes, but this time in a cluster consisting of three instances, and a possible distribution of the shards. Note that Elasticsearch will make sure that the shards are distributed more or less equally across all instances of the cluster. That way, not only the amount of data held by the cluster but also the load of handling requests can be distributed more evenly. Note that clients do not need to know which documents of a cluster reside in which shards - clients can perform requests on any instance and Elasticsearch will then transparently retrieve the data needed to answer the request from the other instances.

Note that in any of the scenarios in figure 32 above, we only have a single copy of each shard (the so-called "primary shard", so our Elasticsearch cluster is not highly available - once an instance holding a shard becomes unavailable, the index will not be available until the situation is resolved).

One important role we have not yet mentioned is that of the master instance. In each cluster (even in one consisting of only a single instance) one instance is elected to be the master, which is responsible for allocating the shards to the instances. In our figures, we indicate the master instance with a star symbol.

Just as we discussed in chapter 2 when we looked at the basics of high availability, we need data redundancy. For each shard, Elasticsearch can manage one or more redundant copies of shards, the so-called "replica shards". To control replica consistency, any changes are handled by the respective primary shard that is then responsible for updating the replica shards (by default with *eventual consistency* only). The replica shards can be used to serve read requests and will allow Elasticsearch to recover the data if the instance holding the primary shard of an index fails. figure 33 shows our two indexes A and B from the example above, but this time configured to hold an additional replica for each shard, and how these replicas could be distributed across the instances of a three node cluster. In the figure, we indicate the shards that are primary shards by underlining their name and using a bold border.

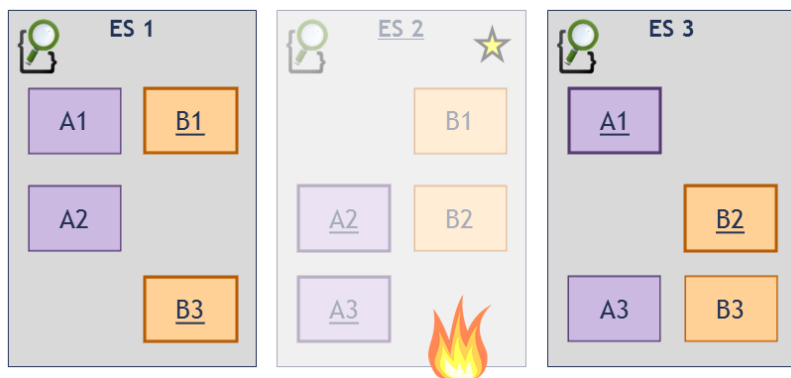
FIGURE 33: DISTRIBUTION OF PRIMARY AND REPLICA SHARDS IN AN ELASTICSEARCH CLUSTER



In the state shown in this figure, instance ES 2 has been elected to be the master instance. Observe that we indeed have two copies of each shard, and that these copies are always on different nodes. So if a single instance of our cluster fails, the data is still available.

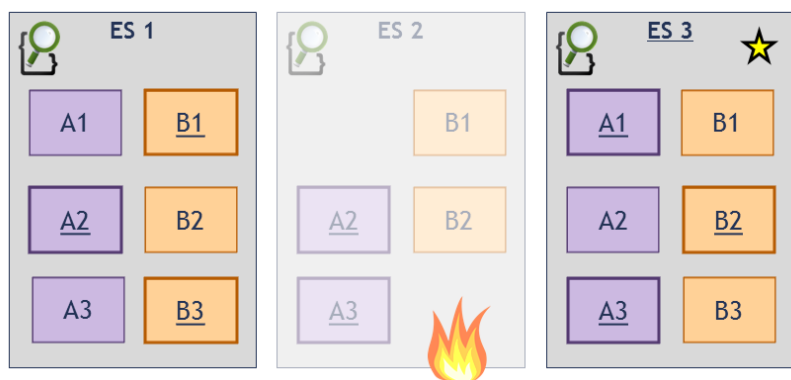
For example, assume that instance ES 2 (which happens to be the master) fails, as shown in figure 34 below.

FIGURE 34: FAILURE OF ONE INSTANCE IN A THREE NODE CLUSTER



Since the failed instance happens to be the master instance, the very first step is the election of a new master instance. Let's assume that instance ES 3 becomes the new master. Since the failed instance also contained the primary shards of A2 and A3, the new master will first make the former replica shards of these shards that exist in instances ES 1 and ES 3, respectively, the new primary shards. Next it will create new replica shards, which it will allocate on the remaining instances. After a short amount of time, all shards will again have a primary and a replica shard. A possible outcome of the automatic handling of the failed instance is shown in figure 35 below.

FIGURE 35: HANDLING OF A FAILED INSTANCE



If in the state shown in figure 35 another instance fails, we would still have one copy of each shard in the system and thus not lose any data. However, to avoid the dreaded split-brain problem that we discussed in chapter 2, an Elasticsearch cluster is in general configured to require a quorum of more than half of all instances to elect a new master. With only one out of three instances available, the cluster would therefore become unavailable, if for example, instance ES 1 would fail as well.

Of in the state shown in figure 35 the failed instance ES 2 comes back, it will rejoin the cluster and will be assigned new (primary and replica shards), to again balance the load across all instances.

The key to making Elasticsearch highly available is now to configure it properly, in order to obtain the behavior that we just sketched. This is the topic of the remainder of this chapter.

5.3 Why clustering Elasticsearch became more difficult with ARIS 10

With ARIS 10, the version of Elasticsearch shipped with ARIS has been updated to 2.4 which does no longer support the automatic cluster discovery via Zookeeper anymore. Further, the alternative, network-broadcast-based approach (a.k.a. multicast discovery) is also no longer standard functionality (and will not work for some environments anyway, in particular on Amazon EC2...). The only remaining discovery mechanism is unicast discovery, which requires that each Elasticsearch instance that should become part of a cluster is given a list of the hostnames and transport ports of at least a subset (but preferably all) other instances of the cluster, so that it can link up with them. Having to manually maintaining such a list, while technically easily possible with the configure parameter

"ELASTICSEARCH.discovery.zen.ping.unicast.hosts", has the significant drawback that a provisioning template would have to specify installation specific parameters (the hostnames!) and would no longer be portable. Even if one accepts this limitation, setting this parameter it is still easily done wrong. In addition, experience from Elasticsearch encountered in customer installations has shown that it is in general difficult to set up a proper Elasticsearch cluster, in particular setting the proper node names, cluster names, the proper quorum etc.

With ARIS 10 we provide a new functionality as part of ACC, the Elasticsearch manager. Similar to the Zookeeper manager that already comes with ARIS 9.x and which we used in the previous chapter to conveniently set up a Zookeeper ensemble, it considerably simplifies the creation of a working Elasticsearch cluster. But before we discuss the Elasticsearch manager, let's first have a look at how to manually configure an Elasticsearch cluster properly.

Making Elasticsearch run on Linux

With ARIS version 10 SR4, ARIS is now using Elasticsearch version 5.x instead of 2.x. With this new version, Elasticsearch is performing a number of initial checks of various system settings (the so-called **bootstrap checks**). If one of these checks fails, Elasticsearch will not start up. These problems can be recognized by checking the runnable's **elasticsearch.log** log file, located in the runnable's log directory (<installDir>/cloudagent/work/work_<ELASTICSEARCH_INSTANCEID>/logs).

One setting that is particularly troublesome on Linux systems is the **vm.max_map_count** setting, which in almost all environments will be set too low by default. Elasticsearch requires this value to be set to at least **262144**. The corresponding error message in elasticsearch.log will be:

"max virtual memory areas vm.max_map_count [xyz] is too low, increase to at least [262144]"

(xyz corresponds to the current value set on your node)

You can also check the current value by running the Linux shell command (e.g., via an SSH shell on the node) directly on the node(s) on which the Elasticsearch runnable(s) is/are located:

```
sysctl -n vm.max_map_count
```

You can change it transiently (i.e., only until next reboot) by running the Linux shell command directly on the node(s) on which the Elasticsearch runnable(s) is/are located:

```
sudo sysctl -w vm.max_map_count=262144
```

To change it persistently, as a root user, edit the **/etc/sysctl.conf** file on the node(s) on which the Elasticsearch runnable(s) is/are located and add the following line to the file:

```
vm.max_map_count=262144.
```

5.3.1 Manually setting up an Elasticsearch 2.x cluster

5.3.1.1 Configuring an Elasticsearch instance to find the other cluster members

Elasticsearch's unicast discovery is based on the gossip concept: a node joining the cluster only needs access to at least one running member of the cluster from which it can then retrieve the information about the other members.

With the configure parameter `ELASTICSEARCH.discovery.zen.ping.unicast.hosts` set for an Elasticsearch runnable, one can provide such a list of nodes. The structure of this parameter is a comma-separated list of the hostnames and ports (with the port separated from the hostname by a colon) of the other Elasticsearch instances in the cluster:

```
ELASTICSEARCH.discovery.zen.ping.unicast.hosts="<esInstance1IpOrHostName>:<esInstance1TransportTcpPort>,<esInstance2IpOrHostName>:<esInstance2TransportTcpPort>,<esInstance3IpOrHostName>:<esInstance3TransportTcpPort>"
```

For example, assuming that we want to tell an Elasticsearch instance on the first worker node of our example scenario about two other instances on the second and third worker node, we would set the parameter like this:

```
ELASTICSEARCH.discovery.zen.ping.unicast.hosts="arisworker2:14230,arisworker3:14230"
```

Note that the port in the list is **not** the ES client port (set with the parameter "`ELASTICSEARCH.http.port`"), but the so-called transport port for the communication between instances, which is set with the configure parameter "`ELASTICSEARCH.transport.tcp.port`").

i Due to the gossip protocol, it would be sufficient to have only one existing member listed here, but since any single existing member might not be available at any given time, one should probably aim to always provide a full list of all other cluster members to each cluster member.

i This is very similar to how the members of a Zookeeper ensemble find each other (via the indexed "server" configuration parameter). However, for Zookeeper, each ensemble member **has** to have the full list of members. Another difference is that each Zookeeper ensemble member has two server ports which have to be added to the list, while Elasticsearch only has one port for communication among the cluster instances, the transport port.

5.3.1.2 Additional Elasticsearch cluster configure parameters

`ELASTICSEARCH.discovery.zen.minimum_master_nodes`

The parameter "`ELASTICSEARCH.discovery.zen.minimum_master_nodes`" defines the size of the quorum, i.e., the number of possible master-eligible Elasticsearch instances that need to be available for a master election to proceed.

To avoid the split-brain problem the parameter always has to be set to values greater than (i.e., ">", **not** ">=") half of the total number of possible master nodes, (i.e., $\text{round}((n + 1)/2)$).

`ELASTICSEARCH.index.number_of_replicas`

The parameter "ELASTICSEARCH.index.number_of_replicas" controls the level of redundancy of your data, i.e., how many replica of a shard should be held in the cluster. Note that this setting refers to the copies of a shard in addition to the primary copy. The parameter should be set to 0 for single node "clusters" (as otherwise, a lone Elasticsearch instance will not have a second cluster instance available to hold the additional copy, leading to the cluster going into the degraded "yellow" state), and to 1 for actual clusters of more than one instance which are meant to be highly available. In theory, you could set this to even higher values (but never set it higher than the actual number of Elasticsearch instances in your cluster minus 1) for additional redundancy. For example, if you set this value to 2 (giving you a total of 3 copies of each shard including the primary), Elasticsearch will distribute these shards in such a way that even if any two instances of the cluster fail, there is still at least one copy of each shard available in one of the remaining instances. However, having an more replica will of course increase the storage overhead and the overhead to keep them in sync. Further, due to the quorum-based consistency control, it in general makes no sense to have significantly more than half as many replica as you have Elasticsearch instances - if you lose more than half your instances, your cluster will become inoperative anyway.

Note that it is technically possible to change the number of replica for an index, however, you cannot do that using this configure parameter, as it only affects the number of replica for newly created shards. In general, it is highly recommended to decide about the initial number of replica in your cluster at the beginning and keep it at that.

ELASTICSEARCH.index.number_of_shards

The parameter "ELASTICSEARCH.index.number_of_shards" controls into how many shards (i.e., physical parts which can be distributed across nodes) each index is divided. A change cannot be applied retroactively on existing indexes and will thus only take effect for indexes that are created afterwards. It should in general be left at its default value 3.

Note: Changes to this parameter will only take effect for indexes that are created afterwards, so make sure you set it properly before starting the system for the first time.

ELASTICSEARCH.gateway.expected_nodes

The parameter "ELASTICSEARCH.gateway.expected_nodes" has to be set to the total number of nodes in the cluster.

ELASTICSEARCH.gateway.recover_after_nodes

The parameter "ELASTICSEARCH.gateway.recover_after_nodes" controls how many nodes need to be available for the cluster to start shard recovery. It should be set to be a quorum of all master nodes, (i.e., $\text{round}((n + 1)/2)$).

5.3.2 Example of a manually configured three node Elasticsearch cluster

Continuing with our three node example scenario, the following configure commands will create a highly-available three-node Elasticsearch cluster, with one Elasticsearch instance on each node n1, n2 and n3.

Example of a manually configured Elasticsearch cluster, here all instances are located on one machine:

```

on n1 configure elastic_m es1 \
    ELASTICSEARCH.discovery.zen.ping.unicast.hosts =
"arisworker2:14230,arisworker3:14230" \
    ELASTICSEARCH.http.port = 14220 \
    ELASTICSEARCH.transport.tcp.port = 14230 \
    ELASTICSEARCH.node.local = false \
    ELASTICSEARCH.node.name = "node1" \
    ELASTICSEARCH.discovery.zen.minimum_master_nodes = 2 \
    ELASTICSEARCH.index.number_of_replicas = 1 \
    ELASTICSEARCH.index.number_of_shards = 3 \
    ELASTICSEARCH.gateway.expected_nodes = 3 \
    ELASTICSEARCH.gateway.recover_after_nodes = 2

on n2 configure elastic_m es2 \
    ELASTICSEARCH.discovery.zen.ping.unicast.hosts =
"arisworker1:14230,arisworker3:14230" \
    ELASTICSEARCH.http.port = 14220 \
    ELASTICSEARCH.transport.tcp.port = 14230 \
    ELASTICSEARCH.node.local = false \
    ELASTICSEARCH.node.name = "node2" \
    ELASTICSEARCH.discovery.zen.minimum_master_nodes = 2 \
    ELASTICSEARCH.index.number_of_replicas = 1 \
    ELASTICSEARCH.index.number_of_shards = 3 \
    ELASTICSEARCH.gateway.expected_nodes = 3 \
    ELASTICSEARCH.gateway.recover_after_nodes = 2

on n3 configure elastic_m es3 \
    ELASTICSEARCH.discovery.zen.ping.unicast.hosts =
"arisworker1:14230,arisworker2:14230" \
    ELASTICSEARCH.http.port = 14220 \
    ELASTICSEARCH.transport.tcp.port = 14230 \
    ELASTICSEARCH.node.local = false \
    ELASTICSEARCH.node.name = "node3" \
    ELASTICSEARCH.discovery.zen.minimum_master_nodes = 2 \
    ELASTICSEARCH.index.number_of_replicas = 1 \
    ELASTICSEARCH.index.number_of_shards = 3 \
    ELASTICSEARCH.gateway.expected_nodes = 3 \
    ELASTICSEARCH.gateway.recover_after_nodes = 2

```

i Please note that the parameter for the Elasticsearch client and transport ports (ELASTICSEARCH.http.port and ELASTICSEARCH.transport.tcp.port, respectively) have been added explicitly to the configure commands even though we actually left them at their respective default values. This is meant to help illustrate the point that the port listed in the "unicast host" list is the transport port.

As you can see, it is easy to get things wrong. While some parameters will be the same for any three-node cluster, others are not - in particular, the manual maintenance of the ELASTICSEARCH.discovery.zen.ping.unicast.hosts is rather painful, in particular since you have to adapt it based on the actual hostnames (or IP addresses) of the machines you want to deploy your cluster to - this also means that your template (i.e., the file containing all your configure commands) is not easily portable anymore, i.e., you have to adapt it to the hostnames resp. IP addresses of the concrete environment.

5.4 Setting up an Elasticsearch cluster with the Elasticsearch manager

To make setting up an Elasticsearch cluster more convenient and less error prone, we now provide tooling inside ACC that supports this task. This functionality is summarized as "Elasticsearch manager".

Conceptually similar to the already existing "Zookeeper manager", which solves basically the same problem for setting up a Zookeeper ensemble, it reduces the effort of setting up a cluster to just deciding about on which nodes to put an elasticsearch instance (and deciding how to name the cluster).

5.4.1 Configuration

The only ACC configuration parameter added for the Elasticsearch manager is "esmgr.elasticsearch.app.type", which has to be set to the name of the app type that will be used for configuring ES instances. If unset, or if no app type with the name given as value for this parameter is found, the Elasticsearch manager will not work.

5.4.2 Commands

5.4.2.1 Validate elasticsearch cluster

The "validate elasticsearch cluster" command performs a validation of the current state of all Elasticsearch instances found across all nodes currently registered in ACC. In particular, it will check if there is more than one cluster (or unclustered instance) across the nodes. Note that having more than one Elasticsearch cluster across the nodes of a single ARIS installation is not supported and will lead to a system that will not function properly.

Syntax:

```
"validate" ("es" | "elasticsearch") "cluster"
```

Example:

```
validate elasticsearch cluster
```

5.4.2.2 List elasticsearch instances

The "list elasticsearch instances" command displays an overview of all Elasticsearch instances found across all nodes currently registered in ACC. For each instance it will show the name of the cluster to which it belongs and the ports occupied.

Syntax:

```
"list" ("es" | "elasticsearch") "instances"
```

Example:

```
list es instances
```

5.4.2.3 Add elasticsearch instance

The "add elasticsearch instance" command allows "staging" of the addition of a new Elasticsearch instance on the specified node (or the current node, if no node is specified) to the cluster with the specified cluster name. Optionally, an instance ID that should be used for the instance can be specified.

Note that this command does not yet perform an actual change of the Elasticsearch cluster.

Syntax:

```
<nodeSpec>? "add" ("es" | "elasticsearch") "instance"? "to" "cluster"  
<identifier> ("using" "instanceId"? <instanceId>)?
```

Examples:

```
add es instance to cluster mycluster
```

Stages the addition of a new elasticsearch instance on the current node that will be part of the cluster "mycluster".

```
on n1 add elasticsearch instance to cluster myescluster using elastic1
```

Stages the addition of a new elasticsearch instance on node n1 that will be part of the cluster "myescluster" and will use the instance ID "elastic1".

5.4.2.4 Remove elasticsearch instance

Currently (and for the foreseeable future), the Elasticsearch manager does not support the removal of Elasticsearch instances due to the inherent danger to customer data imposed by such an operation.

The reason is that we do not know if all indexes have been created with replication. If we deconfigure the only instance holding the sole copy of a shard, the data will be gone.

5.4.2.5 Validate elasticsearch changes

The "validate elasticsearch changes" command performs a validation of the currently staged changes of the Elasticsearch cluster. In particular, it will check if the resulting cluster is highly available.

Syntax:

```
validate "pending"? ("es" | "elasticsearch") "cluster"? "changes"
```

Example:

```
validate es changes
```


5.4.2.6 Commit elasticsearch changes

The "commit elasticsearch changes" command applies the currently staged changes to the Elasticsearch cluster by running the required configure (for new instances) and reconfigure (for existing instances) commands (and in the future perhaps also deconfigure commands for instances staged for removal, if we ever support this).

Note that since a reconfigure of existing instances requires a restart, the cluster might be temporarily unavailable during this operation. Further, we also observed that despite proper configuration of the cluster instances after this operation, the Elasticsearch cluster would sometimes not return to a working state by itself, but only after a shutdown and restart of the entire cluster. It is therefore currently not recommended to try and change a running Elasticsearch cluster.

Syntax:

```
commit "pending"? ("es" | "elasticsearch") "cluster"? "changes"
```

Example:

```
commit es changes
```

5.4.2.7 Reset elasticsearch changes

The "reset elasticsearch changes" command discards all currently staged Elasticsearch cluster changes.

Syntax:

```
("discard" | "reset") "pending"? ("es" | "elasticsearch") "cluster"?
```

Example:

```
discard es changes
```

5.4.3 Usage examples

A common use case is to create a highly available Elasticsearch cluster consisting of three instances which should be spread across three nodes. Assuming your nodes are registered as n1, n2 and n3, respectively, you would need to run the following commands to create this cluster:

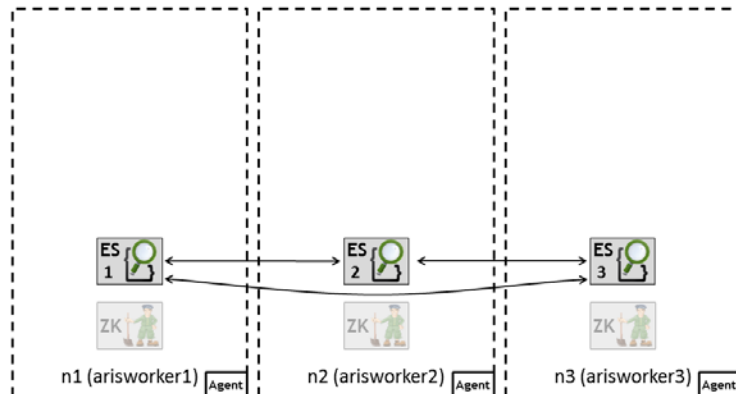
```
on n1 add es to cluster myelasticsearchcluster
on n2 add es to cluster myelasticsearchcluster
on n3 add es to cluster myelasticsearchcluster
commit es changes
```

Afterwards, start your newly created instances and you will have a working cluster, which can tolerate the outage of one of its instances.

5.5 Elasticsearch in our example scenario

You can use the commands from our Usage examples section above to set up a three-node Elasticsearch cluster. You can now run these commands directly in your ACC shell and/or you can add them to your template file. Assuming the commands complete successfully, your distributed example system will now look as it is shown in figure 36 below.

FIGURE 36: THE ELASTICSEARCH CLUSTER IN OUR EXAMPLE SCENARIO



6 Relational database backends

6.1 Motivation - The role of relational databases in the ARIS architecture

Despite the adoption of some "modern" backend systems in the ARIS 9/10 architecture (like Elasticsearch or Zookeeper), the primary storage location for most of the user-generated data (ARIS databases, models, definitions; ARIS method customizations; ARIS reports; collaboration content; document metadata and (optionally) content, process governance data) is still a "traditional" relational database system (RDBMS). RDBMSs have been around for decades now and are a mature technology that offers transaction semantics to ensure that any change to the data is done in an atomic, consistent, isolated and durable way (the ACID paradigm that we briefly introduced in chapter 2).

Due to the pivotal role of the relational backend in ARIS, when aiming for a highly available (and therefore inevitably distributed) ARIS installation, one has to obviously make the relational backend highly available as well.

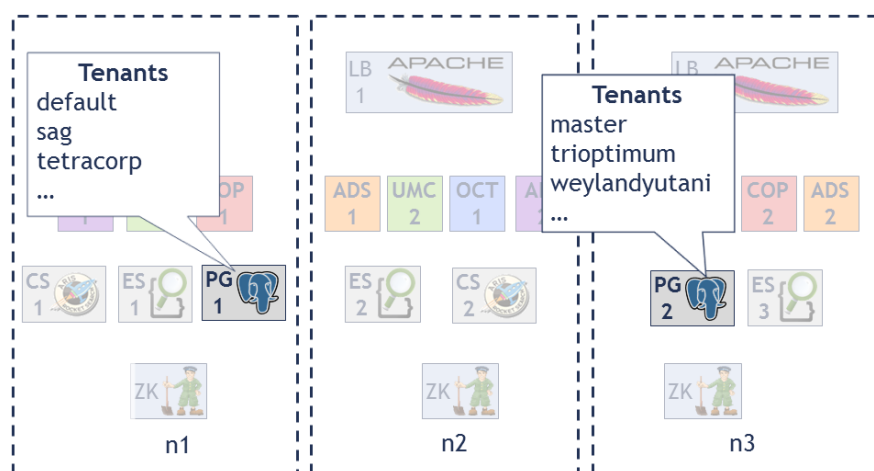
6.2 Scaling & high availability with relational databases

The standard RDBMS shipped with ARIS is PostgreSQL, which we provide packaged as a runnable, so that it can be installed and configured using the ARIS provisioning technology. Alas, while for services like Zookeeper or Elasticsearch it is possible to configure them to form a cluster as described in the previous chapters, we do not offer an out-of-the-box clustering solution for PostgreSQL. (There is a plethora of approaches to cluster our standard DBMS PostgreSQL, but none of them has achieved wide-spread acceptance in the industry as a foundation for mission-critical applications, so we currently do not support any of them officially for use with ARIS.)

6.2.1 Distributing the load of multiple tenants across several DB instances (no high availability)

The only kind of scaling that is possible with PostgreSQL at the moment is that of scaling "across" tenants: Currently, the data of a single tenant will always reside in a single database instance, so the load created by a single tenant on the database instance needs to be handleable by that instance. At the same time, a tenant's database instance will also be a single point of failure for that tenant - if the tenant's DB goes offline, the tenant will become unusable until the DB is available again. However, if high availability is not relevant for an ARIS usage scenario, but if a large number of tenants are to be supported, it can make sense to have multiple database instances available in the system. The data of the tenants will then be distributed (more or less) evenly across all available databases services, as illustrated in figure 37 below.

FIGURE 37: A POSSIBLE DISTRIBUTION OF TENANTS ACROSS MULTIPLE POSTGRESQL INSTANCES



Here, we have two instances of our PostgreSQL runnable, one on node n1, another on node n3. All tenants created in this ARIS installation will be distributed evenly across these two runnables, as indicated by the callouts.

While this offers a way to spread the load of multiple tenants across several instances, it is rarely used in practice. In particular for production use on mission critical systems, you are usually interested in high availability, which this approach will *not* provide.

6.2.2 High availability for the relational database backend by using an external, highly available database cluster (Oracle or SQL Server)

Since we do not support a highly-available configuration using our PostgreSQL runnable, you will have to use an external DBMS (where "external" means "software not provided together with ARIS") like Oracle or MS SQL, which themselves offer mechanisms for clustering and high availability. The details of setting up such a DBMS cluster are obviously beyond the scope of this text (or any other ARIS manual) but have to be handled by the database administrator.

From the perspective of ARIS, a clustered external DBMS behaves just like its respective non-clustered counterpart, i.e., clustering is transparent to the application. Consequently, also the steps needed to make ARIS use a clustered external are the same as for setting up ARIS with a NON-clustered external DBMS (the only difference might be the way ARIS connects to the cluster, in particular the database URL (see next section) can look a bit different from that of a non-distributed external DBMS - but this information will be provided by the database administrator). For the sake of completeness, the following sections will detail the steps needed for setting up ARIS with an external DBMS (or refer to other documentation for some aspects).

6.2.2.1 Registering the external DBMS as an external service

Since our runnables use Zookeeper to find services, the information needed for ARIS to find and use an external DBMS has to be added to Zookeeper as well. For We will be using the ACC commands for registering and updating external services. With the *register external service* command you perform the initial registration of the service. With the *list external services* command you can get an overview over all external services that are currently registered. With the *show external service* command you can review the properties of individual registered services. With the *update external service* command you can change individual properties of the registered service.

Note that for these commands to work, you need to have started your Zookeeper runnable (or Zookeeper runnables, if you have a highly available Zookeeper ensemble).

Registering an external Oracle DBMS

To register an external Oracle DBMS, you need to run the "register external service" command structured like this (note that the command below has been broken up across several lines by adding "\" character in front of every line break that we introduced for legibility).

```
register external service db \  
  host=<OracleHostnameOrIP> \  
  url="jdbc:oracle:thin:@<OracleHostnameOrIP>:<OraclePort>:<OracleServiceNa  
me>" \  
  driverClassName=oracle.jdbc.OracleDriver \  
  jmxEnabled=true \  
  username="<AppUserName>" \  
  password="<AppUserPwd>" \  
  maxIdle=15 \  
  maxActive=300 \  
  maxWait=10000 \  
  removeAbandoned=false \  
  removeAbandonedTimeout=600 \  
  logAbandoned=false \  
  initialSize=0 \  
  testOnBorrow=true \  
  validationQuery="select 1 from dual" \  
  defaultAutoCommit=false \  
  rollbackOnReturn=true \  
  jmxEnabled=true
```

Replace the placeholders (the parts in angled brackets), i.e., <OracleHostnameOrIP>, <OraclePort>, <OracleServiceName>, <AppUserName> and <AppUserPwd>, with the hostname (or the IP address), the port and the service name of your Oracle DBMS, the name of the application user and the password of this user (as you specified them in the envset.bat file that comes as a part of the schema creation scripts for the variables CIP_APP_USER and CIP_APP_PWD, respectively), respectively. Note that the <OracleHostnameOrIP> placeholder is used twice, once for the value of the "host" parameter, and a second time as part of the value of the "url" parameter.

The values for <OracleHostnameOrIP>, <OraclePort>, <OracleServiceName> (or the entire value of the "url" parameter) will have to be provided by your database administrator.

i You will observe that a lot of these parameters hold values that do not pertain in any way to the actual details of the database system being used. These parameters control the usage of the DBMS by the application, in particular the handling of the JDBC connection pool. These parameters can be used in demanding scenarios to tune the database usage to specific use case. For most scenarios, however, the default values you can find in the command above will serve nicely. Therefore, unless explicitly told by ARIS R&D to do so (or if you really know what you are doing), you shouldn't change these values from the defaults shown above.

For example, assuming that your Oracle DBMS is running on a machine myoracle.example.com on port 1521, uses the service name XE, and you used the application user name "arisuser" and password "TopS3cret", the actual command would be

```
register external service db \
  host=myoracle.example.com \
  url="jdbc:oracle:thin:@myoracle.example.com:1521:xe" \
  driverClassName=oracle.jdbc.OracleDriver \
  jmxEnabled=true \
  username="arisuser" \
  password="TopS3cret" \
  maxIdle=15 \
  maxActive=300 \
  maxWait=10000 \
  removeAbandoned=false \
  removeAbandonedTimeout=600 \
  logAbandoned=false \
  initialSize=0 \
  testOnBorrow=true \
  validationQuery="select 1 from dual" \
  defaultAutoCommit=false \
  rollbackOnReturn=true \
  jmxEnabled=true
```

Assuming that the command was executed successfully, it will return the technical ID given to the newly registered external service in Zookeeper:

```
ACC+ >register external service db host=myoracle.example.com [...]
New external service of type db registered with serviceId db0000000000.
```

Notice the service ID (here "db0000000000").

You will later need this ID when we later assign the ARIS tenants to the database system, as detailed in the section below. You can note it down now, or you can look it up later using the aforementioned "list external services" command.

Registering an external MS SQL DBMS

The procedure to register an external Microsoft SQL Server DBMS is pretty much the same as for Oracle, the only difference being the default (and fixed) value of the "driverClassName" parameter and the structure of the value of the "url" parameter. The basic structure of the "register external service" command for MS SQL is the following:

```
register external service db \  
  host=<SqlServerHostnameOrIP> \  
  url="jdbc:sqlserver://<SqlServerHostnameOrIP>:<SqlServerPort>;DatabaseNam  
e=<DatabaseName>" \  
  driverClassName=com.microsoft.sqlserver.jdbc.SQLServerDriver \  
  jmxEnabled=true \  
  username="<AppUserName>" \  
  password="<AppUserPwd>" \  
  maxIdle=15 \  
  maxActive=300 \  
  maxWait=10000 \  
  removeAbandoned=false \  
  removeAbandonedTimeout=600 \  
  logAbandoned=false \  
  initialSize=0 \  
  testOnBorrow=true \  
  validationQuery="select 1 from dual" \  
  defaultAutoCommit=false \  
  rollbackOnReturn=true \  
  jmxEnabled=true
```

Replace the parts in angled brackets, i.e., <SqlServerHostnameOrIP>, <SqlServerPort>, <DatabaseName>, <AppUserName> and <AppUserPwd> with the hostname (or IP address), the port, the database name and the database instance name of your SQLServer DBMS, and the application user name and password as you set it in the envset.bat before running the schema creation scripts. Note that the <SqlServerHostnameOrIP> placeholder is used twice, once for the value of the "host" parameter, and a second time as part of the value of the "url" parameter.

For example, assuming that your MS SQL Server DBMS is running on a machine `mysqlserver.example.com` on port 1433 and the database name is "arisdb", and assuming that the username and password of the application user are "arisuser" and "TopS3cret", respectively, the actual command would be

```
register external service db \  
  host=mysqlserver.example.com \  
  url="jdbc:sqlserver://mysqlserver.example.com:1433;DatabaseName=arisdb" \  
  driverClassName=com.microsoft.sqlserver.jdbc.SQLServerDriver \  
  jmxEnabled=true \  
  username="arisuser" \  
  password="TopS3cret" \  
  maxIdle=15 \  
  maxActive=300 \  
  maxWait=10000 \  
  removeAbandoned=false \  
  removeAbandonedTimeout=600 \  
  logAbandoned=false \  
  initialSize=0 \  
  testOnBorrow=true \  
  validationQuery="select 1 from dual" \  
  defaultAutoCommit=false \  
  rollbackOnReturn=true \  
  jmxEnabled=true
```

As was described in the section before about the registration of an Oracle DBMS, please note the service Id given to your newly registered database service.

6.2.2.2 Creating the database schemas for each tenant

ARIS is a multi-tenant-enabled application, i.e., a single physical ARIS installation (whether it is a single-node or multi-node, i.e., distributed installation) can host several "logical ARIS installations", i.e., "tenants". The different tenants share the infrastructure of their installation, but have complete separated content (i.e., users, ARIS databases, collaboration content, documents, process governance processes etc.). Inside the DBMS the data of different tenants is separated by holding it in different database schemas (or their respective equivalent, Oracle DBMS does not have the concept of a schema as a "container for database tables", but uses databases users for that purpose). Even if you do not plan to use more than one tenant for actual ARIS project work, each ARIS installation has **at minimum not one but two tenants**, named "default" and "master". "default" is the tenant that by default (pun intended) is used for ARIS project work. "master" is the tenant used for managing other tenants, in particular for creating new and deleting existing tenants, or for access to the ARIS health check UI.

Usually, the administrators of an Oracle or MS SQL database server will not provide an application with the necessary access rights so that it can create its own database schemas. Therefore, the administrator will need to run a set of scripts provided by ARIS (at the time of this writing these scripts are located on the ARIS DVD at the location `Add-ons\DatabaseScripts`) that create the required schema objects with the proper permissions. The preparation and usage of these scripts is described for the different database systems and installation platforms in the "ARIS Server Installation guides" found on the ARIS DVD (at the time of this writing located at `Documents\English\3 Installation\31 Initial installation\312 Server`, files "ARIS Server Installation - Linux.pdf" and "ARIS Server Installation - Windows.pdf"). In general, a dedicated schema needs to be created for each tenant.

Since as we mentioned above each ARIS installation has at minimum not one but two tenants ("default" and "master"), you need to create at least two schemas, plus additional schemas for each additional tenant you want use on the installation (additional tenants and their associated schemas can of course be added at a later stage, as described in section 1.18 of the document "ARIS Cloud Controller (ACC) Command-line Tool".

i In ARIS 9.x the database scripts for Oracle not only created the actual database schema object (i.e., the "container for tables") and the required roles and permissions, but also created the actual database tables in that schema. With ARIS 10, the database schema is created in a way such that the application has the necessary permissions to create the database tables itself automatically. While this makes little difference for a fresh installation, this also makes updating an ARIS 10 installation more convenient, as it also makes it possible to handle any schema update operations by the application (i.e., adding new tables or adding new columns to existing tables).

So, different from ARIS 9.x, where schema update scripts had to be executed when updating an ARIS installation that was using Oracle as external DBMS from, e.g., ARIS 9.7.x to ARIS 9.8.y, this is no longer necessary when updating from ARIS 10.x to ARIS 10.y.

6.2.2.3 Assigning the tenant to the database service

After you registered your Oracle or MS SQL Server database system with ARIS as described above and created the schemas for each tenant, you now need to tell ARIS two things

1. the database service that should be used for each tenant
2. the name of the schema in the respective database that should be used for each tenant

Point 1) is needed, because when an external DB service is used, ARIS will not automatically assign tenants to it (as it would do when our PostgreSQL runnable is used), and because there can be multiple database services, as explained in the section about horizontal scaling across tenants above. Point 2) is needed because the choice of naming the schemas for the different ARIS tenants is at the discretion of the database administrators and will in general be different from the default schema naming rules used when ARIS itself can create the schemas (i.e., when our PostgreSQL runnable is used).

Fortunately, both aspects can be handled with a single ACC command that needs to be run once for each tenant, the "assign tenant to service" command. The command creates a logical connection between a tenant and a concrete service instance (in this case an instance of the service type DB, as we created it in Zookeeper when we registered the external DB above). In addition, parameters can be specified for this logical connection. The only mandatory parameter is "com.aris.cip.db.schema", which is used for specifying the schema name.

The basic grammar of the "assign tenant to service" command is

```
"assign" "tenant" <tenantId> "to" "service" <serviceId> (<key> "=" <value>)*
```

Here, <tenantId> is the ID of the tenant ("default" and "master", respectively, for the two tenants that implicitly exist in every ARIS installation, or any string of a maximum length of 30 characters, that starts with a lower-case latin letter (i.e., "a" - "z") and consists of only lower-case latin letters and digits), <serviceId> the ID of the service (as it was returned by the "register external service command" above). <key> and <value> are the name and the value of various properties that can be stored with the tenant-2-service assignment. Multiple such key-value pairs can be specified.

For example, assuming your registered external DB service was given the service ID db0000000000 and you named the schema for the master tenant "aris_master" and the schema for the default tenant "aris_default", you would run the following two commands to make ARIS use this DB and these schemas for these tenants:

```
assign tenant master to service db0000000000
com.aris.cip.db.schema="aris_master"
assign tenant default to service db0000000000
com.aris.cip.db.schema="aris_default"
```

i Note that when creating additional tenants for an ARIS installation that uses an external DBMS service, you need to run the "assign tenant to service" command *and* create the corresponding schema **before** running the "create tenant" command.

This is detailed in section 1.12 of the "ARIS Cloud Controller (ACC) Command-line Tool" manual.

6.2.2.4 Enhancing the runnables with the JDBC driver for the external database system

The external database systems supported are commercial products. Licensing restrictions usually not only apply for the DBMS server itself, but also extend to the JDBC drivers needed for the ARIS to access the database. We can therefore not ship the JDBC drivers for Oracle and MS SQL Server directly with the ARIS runnables. Instead, the drivers have to be added to all runnables that access the database (ADSAdmin, ABS, APG, Octopus, Simulation, Dashboarding, ECP, Copernicus), by using the Provisioning enhancement functionality, using the commonsClasspath enhancement point. The basic structure of the required enhancement command is

```
on <nodeName> enhance <instanceId> with commonsClasspath path
<pathToJdbcDriverFileInRepo>
```

By default, the file with which you enhance the runnable is taken from the same repository as the runnables themselves. So the simplest way to make the JDBC driver available for the enhancement commands is to copy the file into a subdirectory of the remote repository that you are using to configure your runnables (as described in chapter 3, where we presented the preparational steps required before one can start setting up a distributed system).

Let's assume your JDBC driver JAR file was named driver.jar and you copied it into the subdirectory "jdbc" of your repository server's root directory. To enhance a runnable named "abs1" on node n1 with this file, you would run the command

```
on n1 enhance abs1 with commonsClasspath path "jdbc/driver.jar"
```

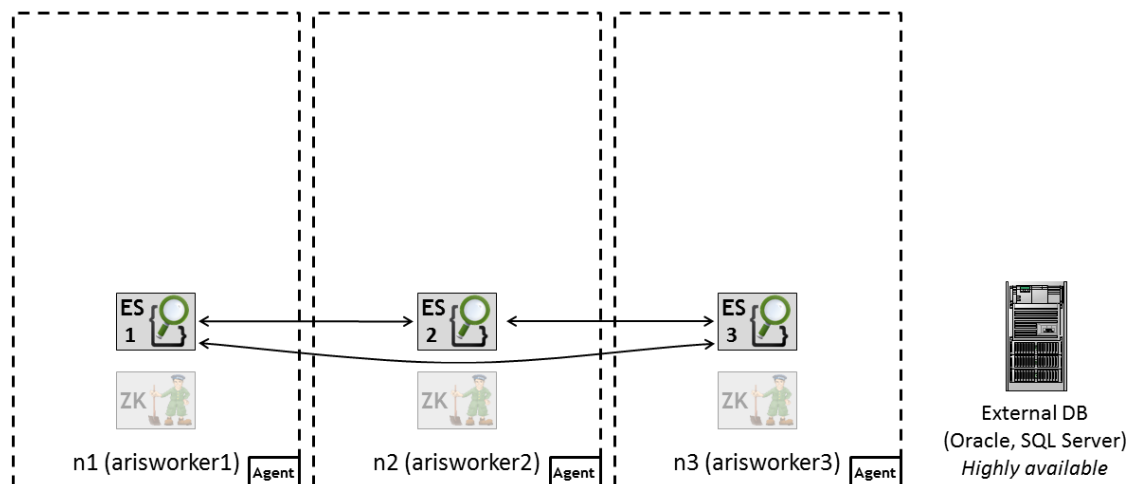
i While it would not be necessary to put the path to the driver JAR file in quotes in this example as it doesn't contain any blanks or special characters, it is in general a good habit to always use quotes here.

Once we configured all application runnables of our distributed system as it will be described over the course of the next sections, we would run this enhancement on all runnable instances of the aforementioned types over all nodes.

6.3 Relational databases in our example scenario

After configuring the external DB for use with ARIS as detailed above, our distributed example scenario will now look as it is shown in figure 38 below.

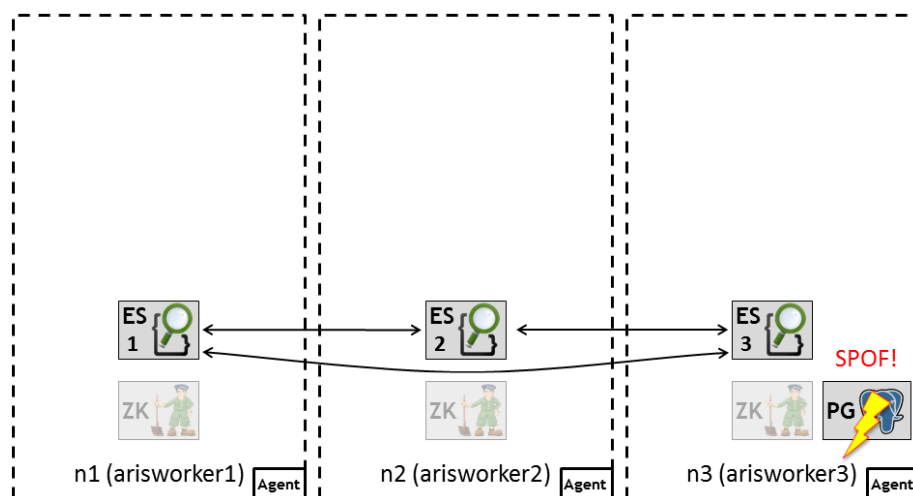
FIGURE 38: OUR DISTRIBUTED EXAMPLE SCENARIO AFTER REGISTERING A HIGHLY-AVAILABLE EXTERNAL DATABASE SYSTEM



Note that you cannot yet run the above commands at this point (and you also should not add these commands to the same template file that hold the "configure" and related commands). Since the "register external service" and "assign tenant" commands require a running Zookeeper cluster, you should run them after all runnables are configured (or add them to a dedicated "template" file).

If your goal is to make yourself familiar with distributed installations in general, but if you want to avoid the hassle of setting up an external DBMS for use with ARIS for now, you can also just configure a single instance of the PostgreSQL runnable on one of your nodes (or an additional node), as indicated in figure 39 below. This means however that this runnable and this node will consequently be a single point of failure for your installation, so this should be used for "exercise only":

FIGURE 39: OUR DISTRIBUTED EXAMPLE SCENARIO WHEN USING THE POSTGRESQL RUNNABLE INSTEAD (I.E., NO HIGH AVAILABILITY, STRICTLY EXERCISE ONLY)



To do this, you just need to run a configure command with a postgres app type, for example, to configure your PostgreSQL runnable instance on node n3 as shown above and name it "pg", you would run the command

```
on n3 configure postgres_m pg
```

(and/or add it to your template file).

7 Cloudsearch clustering

7.1 Motivation - The role of Cloudsearch in the ARIS architecture

Disclaimer

The following (admittedly rather high-level) description of the motivation for introducing the Cloudsearch component to the ARIS architecture requires a bit of understanding of relational databases. If you are not a database guy, don't bother.

One major change in the use of relational databases from ARIS 7.x to ARIS 9/10 was that the underlying relational schema has changed significantly. In ARIS 7.x, complex ARIS entities (like models or object definitions) were stored in a traditional "third normal form" (3NF) schema, where basically all properties of an entity are split up into individual columns and (for multi-valued properties) into dedicated tables connected with foreign-key relationships (just imagine that your complex object is first disassembled and then all parts are stored). While this traditional mapping of complex entities to the relational datamodel has its benefits (in particular, you can access any property of an entity directly with SQL, allowing you to express virtually any kind of query), the drawback is that in order to load such an entity from the database, it need to be reassembled by many, often complex queries across many database tables, which is often slow and puts a significant load on the DBMS. And since for most use cases, the complex entity was anyway loaded as a whole and rarely only individual parts, this structuring of the data made the most common use cases slow.

Starting with ARIS 9, complex entities are stored in an efficient, compressed binary form (as "binary large object", a.k.a. "BLOB" columns). That way, loading a complete entity is a single key-based database access, which is very fast. The drawback is however, that it is now no longer possible to perform complex queries on ARIS content via SQL. This is where Cloudsearch comes in: Cloudsearch holds indexes over the ARIS data, both indexes for attribute-based (including full-text) search and structural indexes that describe the relationships among ARIS entities. If you are a database guy you can now argue that relational databases also hold indexes over their data to speed up access. However, Cloudsearch indexes are held completely in memory, making it possible to answer even complex queries (that would correspond to large joins across many tables) without any disk access. With Cloudsearch indexes it is now feasible to perform even complex structural queries over ARIS content with good response times and with high throughput, which is needed in particular for the ARIS Connect portal, which is built to allow a large number of concurrent users a (primarily read-only) access to the ARIS repository.

So basically, the idea is to leave the aspect of securely storing the valuable customer data where it is done best (i.e., in the mature RDBMS), while the handling of complex queries specific to ARIS is done in a dedicated, highly optimized component.

Your next question should now be: "How are the database content and the Cloudsearch indexes kept in sync with each other when I change something?"

Good question! This is where the CIP module comes in. CIP is the component through which all database access of applications using Cloudsearch is handled. Whenever a change is made, the change is written to the database (e.g., if you save a model, that model's BLOB is written to the DB). Further, all updates to the Cloudsearch indexes that result from this change are sent synchronously(!) to Cloudsearch (or rather to all Cloudsearch instances holding the relevant index, more on that later), which will then apply the changes to its *in-memory* index structures (i.e., the change is not necessarily persisted). While the synchronous communication to the Cloudsearch instances of course incurs a bit of an overhead and also means that the slowest Cloudsearch instance will be the limiting factor for the response time of a change operation, this approach ensures that all Cloudsearch instances will immediately have all changes, so that a query sent directly after a change will directly see the change (i.e., from the perspective of the application Cloudsearch content is always fully consistent with the database content). In addition, all index changes are stored in the database as so-called "pending packets". This is needed to make sure that changes to the Cloudsearch indexes are not lost if a Cloudsearch instances goes down: Since in order to obtain a good performance for updates, the index changes in Cloudsearch are *not* persisted to disk immediately (as otherwise any change operation would synchronously write to the database and to the disk of all Cloudsearch instances, which is obviously *way* too slow) but *only* updated in memory. A background job in Cloudsearch will persist the changes from time to time in a batch operation to the Cloudsearch index files on disk. Only once a change has been written to disk by all Cloudsearch instances holding the relevant index(es), the "pending packet" stored in the database can be safely deleted. So this mechanism ensures that the Cloudsearch indexes are always in sync or in the case of a crashed Cloudsearch instance can easily be made to be in sync by "replaying" the remaining "pending packets" to the restarting Cloudsearch instance.

It is important to point out that the database still holds all information there is about the ARIS content - the data held in the index structures in Cloudsearch is just derived from the database content. This has the advantage that if for some reason the Cloudsearch indexes and the database content become out of sync (despite the mechanisms described above) or the index is corrupted in some way, it can be rebuilt from the database content, a process known as reindexing. Depending on the amount of ARIS content, this process can take a while, though, but your data is always safely stored in the database.

i Which applications actually use Cloudsearch?

Note that not all ARIS components make use of Cloudsearch and the CIP module. In ARIS 9.8.x, the business server (ABS), the portal (Copernicus), analysis (Octopus), and process governance (APG) are using CIP. In versions before 9.7.x, also user management (UMC) data was stored via CIP. The Collaboration component (ECP) is using a traditional 3NF schema and for that reason had its own database instance, "PostgreSQL-ECP" in ARIS 9.x.

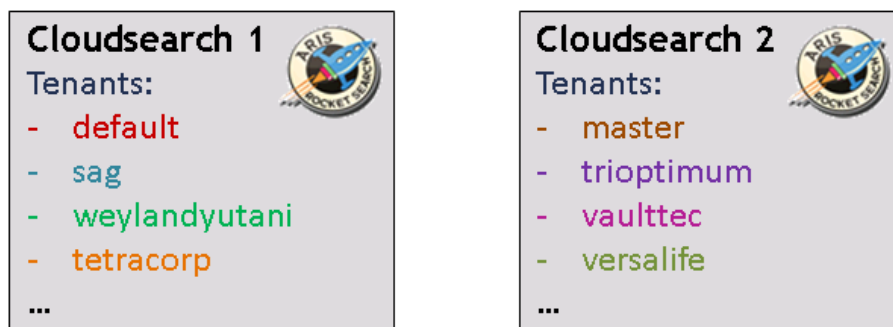
In ARIS 10, only the business server (ABS), portal (Copernicus) and analysis (Octopus) components make use of CIP and Cloudsearch, which have been further optimized to cover their specific use cases. APG (like ECP before) now also uses a more traditional 3NF schema, as does the document storage component (ADS), whose backend was changed from 9.x to 10 to use the database for the storage of document metadata (and optionally also document content). Despite the fundamental difference in how CIP-based and non-CIP-based applications represent their data in the database, with ARIS 10 all applications can share a single database system (and even share a single schema for each tenant), which is why the dedicated database system for Collaboration has been removed for ARIS 10.

7.2 Cloudsearch clustering, variant 1: data load distribution of multiple tenants across several instances, no high availability, no query load distribution

In a standard single-node installation, a single Cloudsearch exists which will hold the index data of all tenants in that installation.

When simply adding several Cloudsearch instances to the system, a tenant's Cloudsearch index data is handled similarly to how it is handled when having multiple database systems in an installation (as we discussed it in chapter 6): all data of an individual tenant will be stored in a single Cloudsearch instance. If multiple Cloudsearch instances are available, the data of different tenants will be distributed across the available instances, as illustrated in figure 40 below:

FIGURE 40: TWO CLOUDSEARCH INSTANCES WITHOUT REDUNDANCY



Observe that in this configuration, using two Cloudsearch instances in your system **only allow the distribution of the data and load caused by different tenants** across different instances - it **does not allow spreading the query load of a single tenant** across multiple instances, **nor does it offer high availability**. When a Cloudsearch instance becomes unavailable, the tenants whose data resided on that instance will not be able to use any ARIS components that require Cloudsearch access (in particular the business server and the portal) until the issue is resolved.

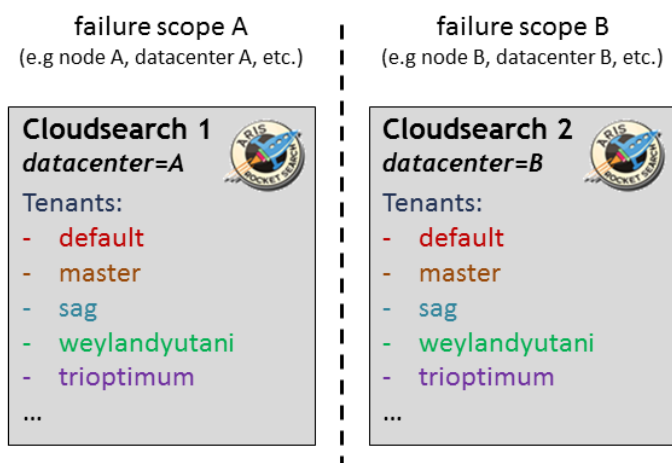
7.3 About failure scopes

To make a set of Cloudsearch instances highly available, it has to be configured appropriately by setting the data center parameter ("zookeeper.application.instance.data center", abbreviated with just "data center" in all figures below) on all instances. The idea is that all Cloudsearch instances having the same value for their data center parameter are considered to be located in the same "failure scope" (a term that we introduced in chapter 2 of this document), i.e. an environment which is likely to fail entirely. In the simplest cases, a "failure scope" would be an individual worker node - if that node fails, then all instances on that node will fail at once. In a larger deployment, the failure scope could be an individual host or hypervisor (i.e. rather the physical server machine running a bunch of VMs, some of which serving as ARIS workers), or it could be all physical servers located in the same server rack. If the hypervisor or the hardware it runs on fails, all VMs will fail at once. Similar, if the server rack has, e.g. a power supply or network issue, all machines in that rack will fail. In the most ambitious deployments, one might want to spread ARIS across different data centers, so here an entire data center would be a "failure scope" and this is also where the configure parameter got its name from.

7.4 Cloudsearch clustering, variant 2: high availability and query load distribution, no data load distribution

In order to make Cloudsearch highly available, there needs to be more than one Cloudsearch instance in the system and there need to be instances with different values for the "zookeeper.application.instance.data center" parameter, i.e., instances allocated to different data centers = failure scopes. In general, Cloudsearch and the CIP module will allocate the data of each tenant to exactly one instance per data center (i.e., you will have as many identical copies of each tenant's index data in the system as you have different values for the data center parameter). If for example you have two Cloudsearch instances and declare one to be in data center A and the other to be in data center B, as shown in figure 41 below, both instances will hold the full index data of all tenants in the system:

FIGURE 41: TWO CLOUDSEARCH INSTANCES ASSIGNED TO TWO DIFFERENT DATA CENTERS



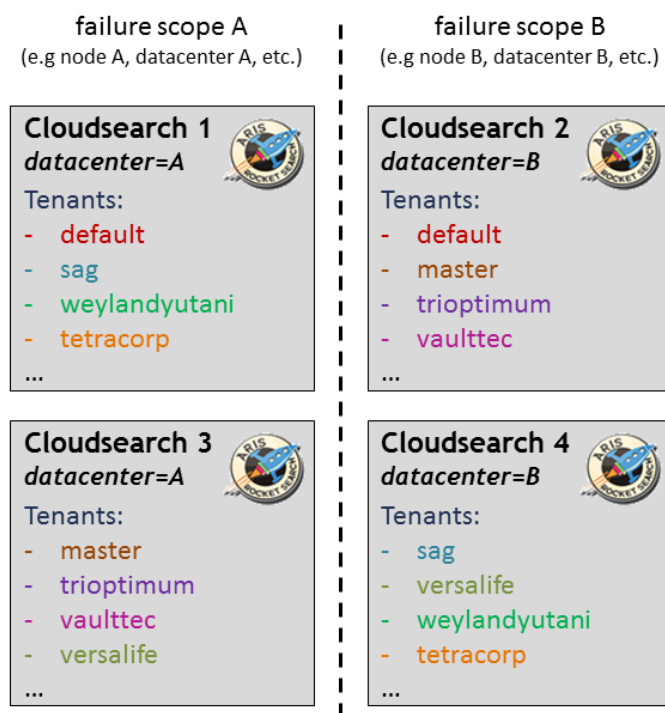
This redundancy will allow the system to continue working if one Cloudsearch instance (or the entire failure scope in which it is running, i.e., the node, host, rack or data center) fails, i.e., it **offers high availability**. At the same time, as long as both instances are available, **query load (but not update load) of individual tenants can also be distributed** across the two instances.

A disadvantage of this variant, as usual with high availability, is that it comes at **the price of increased resource usage** as we need double the amount of CPUs and memory for the additional Cloudsearch instances, as now again **each Cloudsearch instance has to hold the full amount of data** of all tenants. Remember that a key element of Cloudsearch query performance is that it can hold the data in memory. If you have many tenants and/or lots of data in your system (e.g., large ARIS databases), you will soon have to increase the amount of RAM available to your Cloudsearch instances (by increasing the Java maximum heap size using the "JAVA-Xmx" parameter), up to the point that you will not have enough RAM on a single node, or to a point that the heap size becomes so large that Java's garbage collection mechanisms become inefficient, leading to degrading performance due to increasingly long pauses for a full garbage collection.

7.5 Cloudsearch clustering, variant 3: high availability, query and data load distribution

To combine the advantages of clustering variants 1 and 2, one simply needs to have both: instances in different data centers and more than one instance per data center. figure 42 below shows a minimal configuration of this variant, with two instances in data center A, and another two in data center B:

FIGURE 42: FOUR CLOUDSEARCH INSTANCES ASSIGNED TO TWO DIFFERENT DATA CENTERS



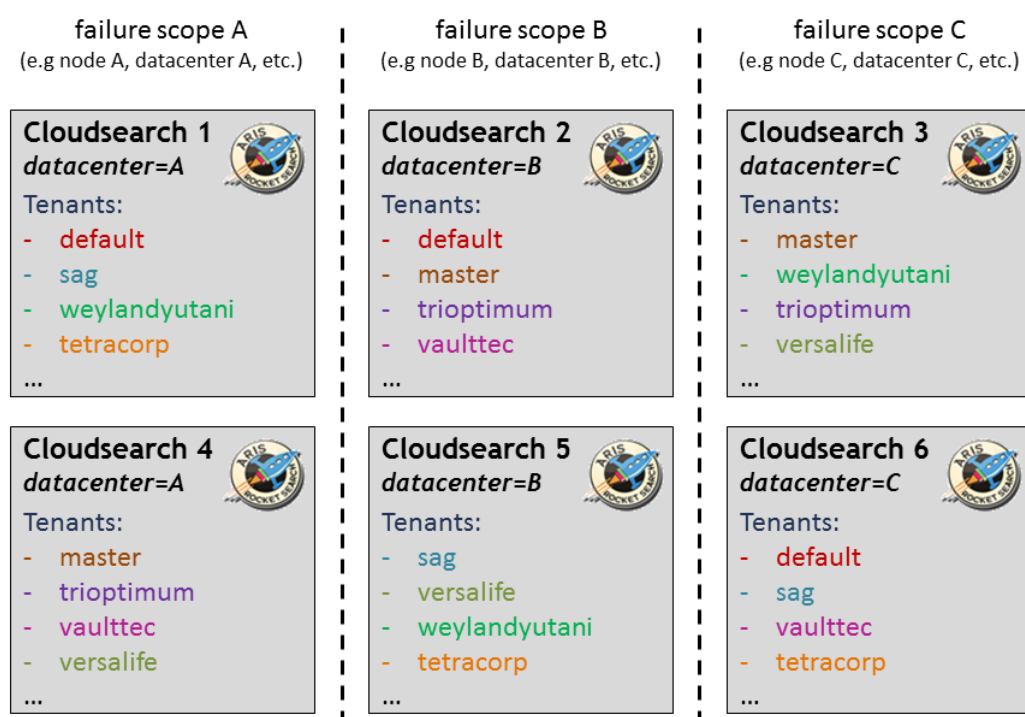
Observe how the index data of each tenant is found twice in the system, and once in each data center: For example, the data of the tenant "sag" is found on Cloudsearch instance 1 in data center A, and instance 4 in data center B. This redundancy again **offers high availability**, in that the system **can tolerate the outage of an entire failure scope**, as all data will still be available in the surviving failure scope.

Note that while in this configuration the system can tolerate the outage of the two Cloudsearch instances that are in the same failure scope, this system **cannot** tolerate the outage of any two instances! If for example Cloudsearch instances 1 and 2 fail, there is no more surviving instance that is allocated to handle the data of the default tenant.

So while the guaranteed maximum number of failing instances that the system can survive has not increased compared to variant 2, we now **distribute the query load of individual tenants across two instances** (one instance in each data center) and we **distribute the total amount of data (i.e. of all tenants) across two instances** (in the same data center). As with variant 2 before, variant 3 comes at the **price of increased resource usage** as we need double the amount of CPUs and memory for the additional Cloudsearch instances.

Depending on the number of tenants and the total amount of data in your system, you can now scale variant 3 in two ways: You can increase the number of failure scopes, i.e., have more data centers. If you added another two instances in a data center C, these two instances together would again hold replica of the data of all tenants, as shown in figure 43 below.

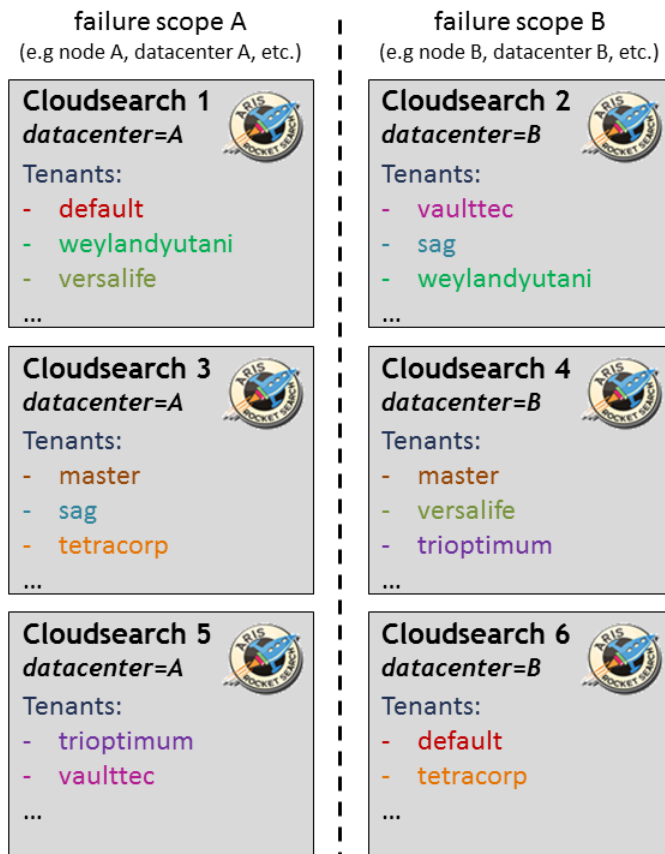
FIGURE 43: THREE DATA CENTERS WITH TWO CLOUDSEARCH INSTANCES IN EACH



Consequently, if set up this way, your **system could now tolerate the outage of two entire failure scopes** (or the failure of any two instances, regardless of which failure scope they are in). Of course since we now have three replica of each tenant's data in the system, we now tripled the resource requirements.

Alternatively, instead of having six instances spread over three data centers, you could also have only two data centers with three instances each, as illustrated in figure 44. Here, the load of the tenants is spread across more instances, but you as you now have only two replica of each tenant's data, your system can again only tolerate the outage of one failure scope.

FIGURE 44: THREE DATA CENTERS WITH THREE CLOUDSEARCH INSTANCES IN EACH

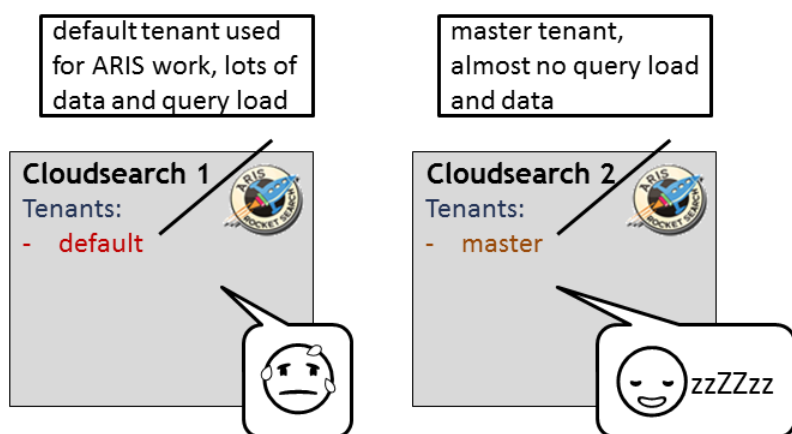


7.6 Some general guidelines & recommendations

Of course the different Cloudsearch deployment scenarios sketched above are only meant as examples to give you an understanding of how a different number of instances distributed across a different number of data centers give you varying levels of redundancy (and thus failure tolerance), query and data load distribution. You can always adjust these scenarios according to your requirements. In this section, we will give a few general guidelines to aid your Cloudsearch clustering decision.

7.6.1 Multiple instances in one data center only makes sense when using more than one tenant

If you are effectively using only one tenant in your system (as is the case in many on-premise scenarios, where many customers will just work with the default tenant), it makes little sense to have more than one Cloudsearch instance per data center. Of course there are always at least two tenants in the system (master and default), and if there are two instances per data center, the CIP module would put the default tenant on one instance and the master tenant on the other in each data center. However, since the master tenant does not contain any modeling data anyway, it makes little sense to dedicate a complete Cloudsearch instance for it, which would be basically idle, as illustrated in figure 45 below.

FIGURE 45: ONE ACTIVE TENANT AND TWO CLOUDSEARCH INSTANCES

Here the basic scenario described in variant 2 above is completely sufficient (of course you can add additional data centers to increase redundancy or distribute query (but not data) load) and dedicate the resources to increase the sizing of the single Cloudsearch instance in each data center (in particular increasing its Java heap).

7.6.2 Size your Cloudsearch instances equally

A general rule is that all Cloudsearch instances should be sized equally (i.e., be given the same amount of memory using the Java maximum heap size setting, "JAVA-Xmx") and be running on (more or less) identical hardware.

This is of particular importance as this is what the load distribution mechanisms assume. For example, if you have two instances in one data center and one has less resources than the other, CIP will still assume them to be sized equally and will thus put the same number of tenants (and thus the same amount of data, assuming that tenants are sized equally) on each. Having the same amount of (equally performing!) CPU cores for each instance is also important, as all change operations have to be done synchronously in all instances assigned to the respective tenant. If one instance is significantly slower, it will be the limiting factor.

7.6.3 Have the same number of Cloudsearch instances in each data center

Not only should your instances be of the same "size", but you also should make sure that you have the same number of instances in each data center, as otherwise the instances in the data center having fewer total instances will on average carry more load, leading to on average slower response time, again making them the limiting factor in particular for the synchronous index update operations.

7.6.4 Adding Cloudsearch instances retroactively will not help existing tenants

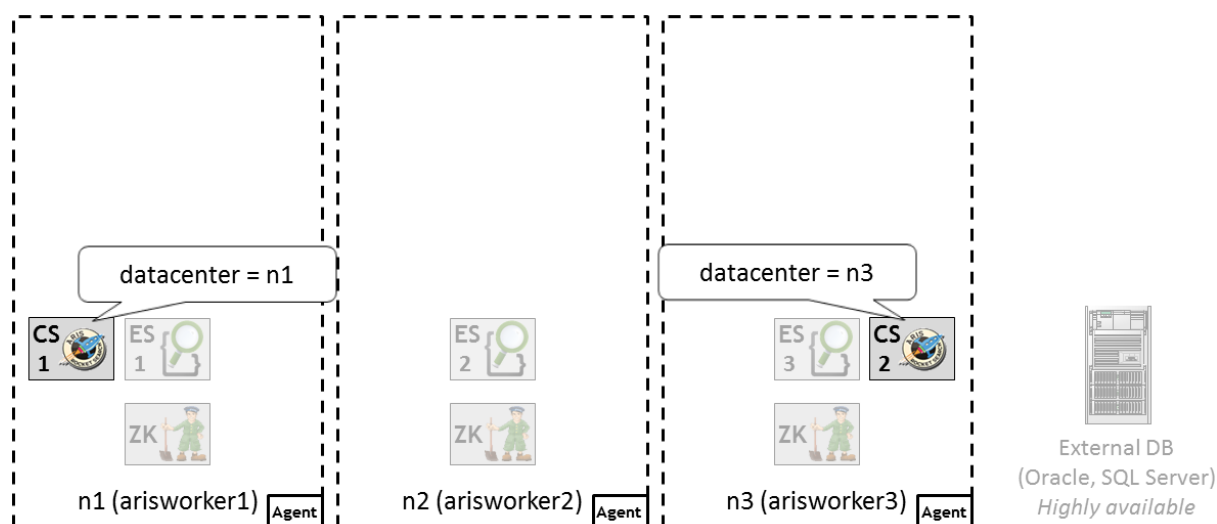
Adding new instances in the same data center will not make CIP module automatically redistribute the indexes of existing tenants to the new instances to balance the load due to the high cost of such an operation. Only new tenants created after the new instance(s) has/have been added will end up in these new instance(s), up to the point that the number of tenants is again roughly equal among all instances.

So if you have a bunch of tenants and only one Cloudsearch instance, and Cloudsearch performance or memory usage is becoming the bottleneck it is not enough to just add another instance (of course to the same data center). If you *really* want to redistribute the tenants among the Cloudsearch instances, the simplest approach is to stop the system, deconfigure *all* Cloudsearch instances, and configure them freshly in the new configuration. That way, CIP can make the tenant-to-cloudsearch assignment from scratch and then distribute the tenants equally among all instances of each data center. Of course this will require a reindexing of all Cloudsearch data, which can take its time.

7.7 Cloudsearch in our example scenario

For our example scenario, let's assume that we have only one or a few tenants (so that the data of all tenants can still be held in a single Cloudsearch instance), but we are of course interested in making Cloudsearch highly available. So we will need to have at least two instances of Cloudsearch in different failure scopes. In our small three-node scenario, we will consider each node as its own failure scope (or "mini data center"). Since we configured our other services (Zookeeper and Elasticsearch) so that we can survive the outage of at most one instance (or one node, i.e. one "mini data center"), it makes little sense to have three Cloudsearch instances distributed across three data centers (which would give the Cloudsearch service as a whole the ability to survive two outages), so we stick with two Cloudsearch instances, one on node n1, the other on node n3, as shown in figure 46 below. We consequently use the node name also as the name of the data center. When we later distribute the micro-service/application runnables to our example scenario, we are going to add a few more runnables to node n2 to make up for the imbalance we now created.

FIGURE 46: CLOUDSEARCH IN OUR EXAMPLE SCENARIO



To configure these two new instances, run the commands

```
on n1 configure cloudsearch_m cs1 zookeeper.application.instance.datacenter = n1
on n3 configure cloudsearch_m cs2 zookeeper.application.instance.datacenter = n3
```

(and/or add them to your template file).

8 Clustering ARIS micro-service-based applications

8.1 Motivation - How distributing micro-service-based applications works

As we learned in the chapter on the principles of high availability, the fundamental challenge for building highly available systems is the handling of persistent data, which we covered for all ARIS backend services in the previous chapters. So the hard part is already done!

With a few minor exceptions, the ARIS applications or microservices are fully stateless. So, making them highly available boils down to merely configuring more than one instance of each type of application required for the actual product. Depending on how you distribute the application across nodes or even data centers, you can protect yourself against the outage of a single runnable (if all instances are on one node), the outage of a node (if the instances are not all on the same node) or even against the outage of a data center (if the instances are distributed on nodes in different data centers).

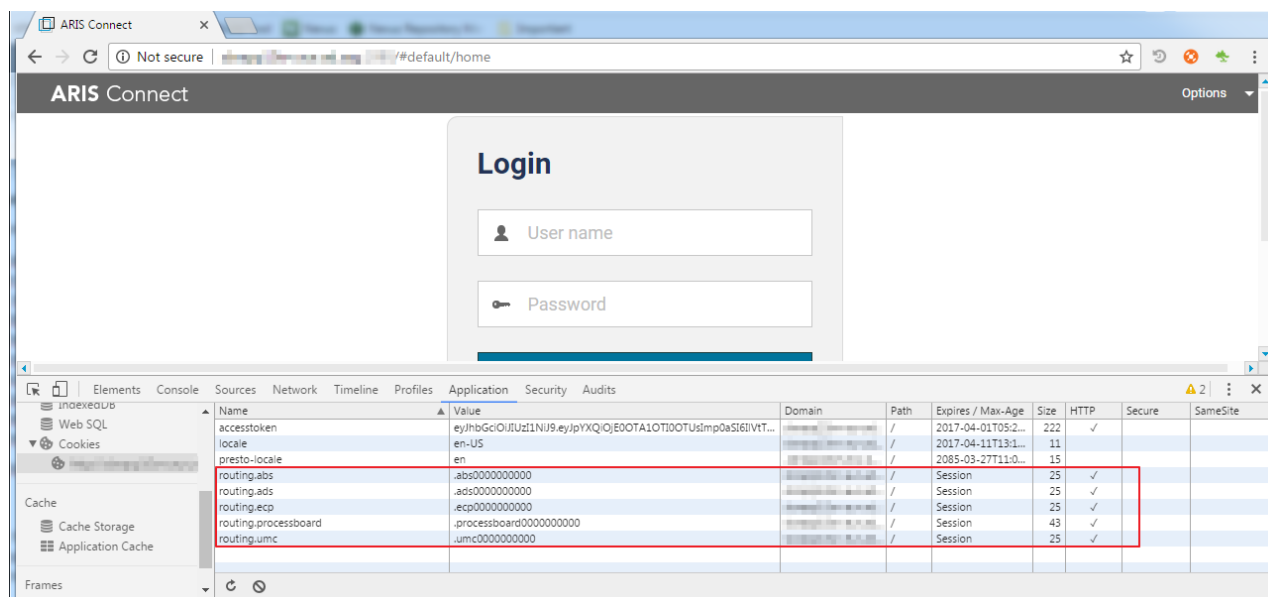
8.1.1 Load balancing and preferred routing

The ARIS LB will distribute the client requests for a specific type of service among all currently available instances, that way spreading the load and allowing the system to handle more requests.

To improve performance, for many applications we have a **preferred routing** mechanism in-place, where the load balancer takes care that requests done by the same client are preferably sent to the same application instances. This allows caching done in the applications to work more efficiently: To illustrate this, assume you have three instance of an application, and there was no preferred routing in place. Further assume that requests would be distributed round-robin, so every n -th request would go to instance $n \bmod 3$. (i.e., the 1st request would go to the 1st instance, the 2nd to the 2nd, the 3rd to the 3rd, the 4rd again to the 1st and so on). The first request would require the client's data to be loaded from the backends into the caches of the first instance. The next request would probably require most of the same data, but the second instance doesn't have it in the cache yet, and would therefore again have to go to the backends to get it and store it in its own cache, which is of course slower than if the request could be served directly from the cache (as instance 1 could do). The third request, going to the third instance, would again find the caches not containing the needed data, and again the backend would have to be accessed. So this obviously increases the load on the backends and slows more initial requests (until all instances have the required data in their caches). So instead of having one cache miss on the initial access to a certain data item, you will have n cache misses for n application instances. But not only that: the effective size of the caches over all application instances is reduced, as in general each instance will have to hold all data of all currently active clients. So the individual applications' caches will fill more quickly, and thus cached item will be evicted from the cache more often to free space, often only to be loaded again moments later when the next request for that client comes in. This will further reduce the cache hit ratio and increase the load on the backends, possibly massively reducing the system's throughput (i.e., the total number of requests that can be handled in a given time interval) and increasing response times. With preferred routing, the load balancer will, if possible, use the same instance to answer the request of one client. These requests are thus much more likely to find the required data already in the cache (this is usually referred to as the "the cache is hot") and can be answered much faster and with less load for the backends.

i How preferred routing works

The preferred routing mechanism is using cookies, as illustrated in the screenshot below. The cookie names all start with "routing.", followed by the application type (e.g., "abs", "ecp", etc.). The value of the cookie is a ".", followed by the internal ID (the Zookeeper instance ID) of the application instance to which requests should be routed, e.g., "abs0000000000").



When requests reach the load balancer, it will first decide which application type is responsible (based on the application's context, e.g., "/abs", "/collaboration" as usual). Then it will check if a cookie for that application is set and if it is pointing to a currently registered instance. If so, the load balancer will send the request to that instance. If no cookie is set for the application or if the cookie value is pointing to a non-existing or currently not serving instance, the normal round robin approach for request distribution will commence.

The cookies are set by the applications themselves: whenever a request reaches an application for which preferred routing is enabled, a servlet filter in that applications filter chain will check if the routing cookie for the application's type is set and that its value is the internal ID of the current application. If this is not the case, the response to the request will be enriched with a set-cookie header that sets such a cookie in the client.

8.1.2 Failure handling

If an instance of an application fails, the load balancer will notice this and redirect requests to the remaining ones. Naturally, any requests that were handled by the failing instance the moment it crashes will themselves fail, but the clients are usually built in a way to tolerate this situation, by simply retrying the failed request. Worst case, your users will observe a little hiccup and might have to, e.g., reload the page in the browser. Further, a user that is redirected to a new instance will of course not benefit from the preferred routing mechanism. The new instance will not yet have the data it is working with in the cache (for that user "the cache is cold"), so the first few requests will again have to go to the backends and will therefore be a bit slower momentarily.

Such minor glitches or slowdowns in the case of a failure aside, a failed instance will of course leave you with **reduced capacity** until it becomes available again, but the automatic restart should make sure that this doesn't take too long. If you want to make sure that even in the face of a failed instance the performance of the system does not suffer too badly, you can do **overprovisioning**, i.e., have more instances of the applications running than would be needed for high availability or for the regular system load by their own. For example, assume that your usual system load can be handled with good performance (i.e., acceptable response times) with three instances of an application running. While three instances are of course redundant and allow the system to operate if one or even two of them fail, to avoid degraded performance during such an outage, you might want to add a fourth (or even fifth) instance to your system.

i In the case of an instance failure, the preferred routing approach described in the previous section can be a bit problematic: During the outage of an instance, all clients that were assigned to it will be routed to the remaining instance(s), and will then have the preferred routing cookie set to point to these instances. Once the failed instance comes back, the clients will however continue to be routed to the surviving instances, i.e., there is no mechanism in place to redistribute clients to recovered (or to freshly started) instances. The recovered instance will only receive requests from "new" clients that do not yet have a preferred routing cookie. Since the preferred routing cookie's expiry is set to "session", they will be deleted if the browser is closed, which is a way to reset the preferred routing mechanism for an individual client (or by deleting the browser's cookies). This issue will probably be addressed in one of the next releases.

8.2 Distributing stateless applications

As was said above, for most applications, getting them highly available is as simple as configuring additional instances. This works for the business server (ABS), user management (UMCAdmin), document storage (ADSAdmin), Collaboration (ECP), portal (Copernicus), analysis (Octopus) and process governance (APG) components. For the ADS component, the configuration has to be changed from its default settings for horizontal scaling and high availability to work. For the simulation, dashboarding and ARCM components, the situation is slightly different, as these components are either not fully stateless or have other technical restrictions, that will be explained below.

8.2.1 Configuring document storage for high availability

With ARIS 10, for the document storage application, a few special configuration settings have to be done in order to make it highly available and horizontally scalable:

While in ARIS 9 documents were held in the CouchDB runnable which we could make highly available using a master-slave replication mechanism, in ARIS 10 we removed the CouchDB runnable. Document metadata is now stored in the relational database backend. Document content, however, will by default, be stored in the file system of the respective node, in the working directory of the ADSAdmin runnable. If you now added multiple ADSAdmin runnables, each would have its own storage, and they would not be able to see each other's documents, which would of course lead to weird errors.

There are two options to remedy this situation:


1. It is possible to also **store the document content in the database**. To enable this, set the configure parameter "JAVA-Dcom.aris.ads.filesystem.active" to "false" *on all(!) ADSAdmin runnables* in your installation, *before* starting any ADSAdmin instance. You can either set this parameter directly when configuring the instances, or later using reconfigure, e.g.

```
configure adsadmin_m ads1 JAVA-Dcom.aris.ads.filesystem.active =
false
```

or

```
reconfigure ads1 JAVA-Dcom.aris.ads.filesystem.active = false
```

respectively.

 If you want to use the database as the backend for document content then perform the configuration change above before starting to store documents, as otherwise already existing documents stored in the default location will no longer be accessible.

The disadvantage of this approach is of course that storing document content in the DB adds additional load on this backend system. Depending on how intensively you intend to use Document storage, this might or might not be a problem. Another problem is that often when using ARIS with an external database system (Oracle or MS SQL Server), the storage space available for the DBMS (i.e., the size of the tablespace) is rather restricted. So in case you want to use the DBMS as storage location for document content, you should check with the administrator of your external database system.

2. You can **use a shared network folder** to hold the documents. This shared folder has to be made available on each node that has ADSAdmin instances running on it. Of course, the degree of high availability that can be achieved for document storage will now depend directly on the degree of high availability of the shared network folder - if this a simple single-node file server, this server will be a single point of failure for the document storage part of ARIS.

The approach to make a shared folder available to ADS will differ depending on whether you are using Windows or Linux:

On **Linux** systems, you can mount your shared folder to a local file system location (e.g., /mnt/arisdocs) using the usual operating system mechanism. The details of how to do this are not in the scope of this document. In any case, you need to make sure that the user under which ARIS is running (by default this is the user "aris", as created by the agent installer package) has permissions to read and write data and create subdirectories in this directory.

All you need to do is tell each ADS instance to use the mount point of the shared folder as its document content backend, by setting the JAVA-Dcom.aris.ads.filesystem.path configure parameter appropriately, either directly when configuring the runnable on configure, e.g., with a command like this

```
configure adsadmin_<SIZING> <adsInstanceID> JAVA-
Dcom.aris.ads.filesystem.path="<adsDocumentFolderMountPoint>"
```

for example

```
configure adsadmin_m ads JAVA-Dcom.aris.ads.filesystem.path="/mnt/arisdocs"
```


or later with a reconfigure command like this

```
reconfigure <adsInstanceID> JAVA-
Dcom.aris.ads.filesystem.path="<adsDocumentFolderMountPoint>"
```

for example

```
reconfigure ads JAVA-Dcom.aris.ads.filesystem.path="/mnt/arisdocs"
```

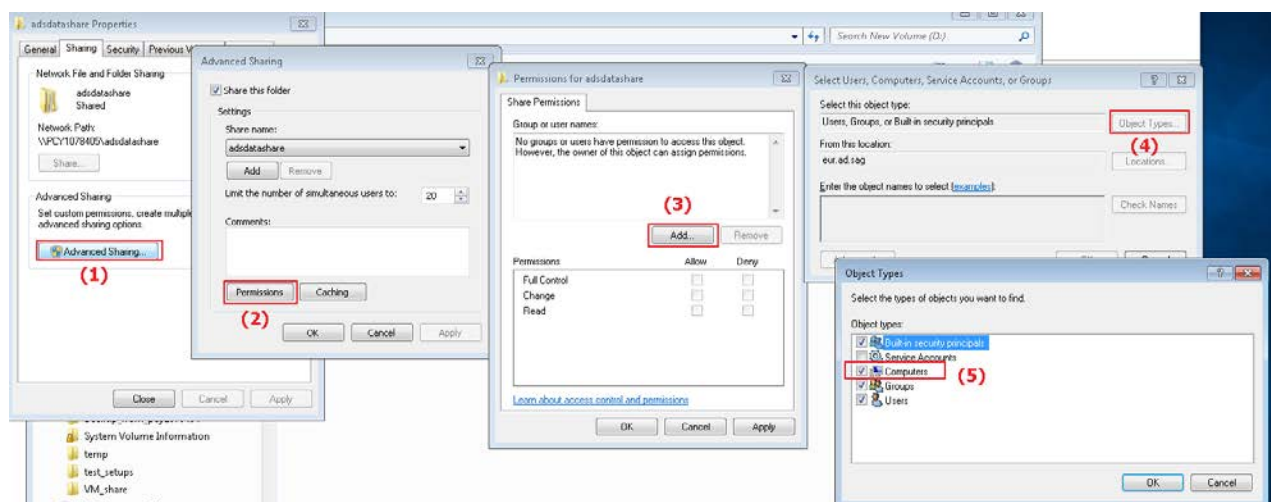
Naturally, **this has to be done for all ADS instances in the installation.**

! If you want to change the folder used as the backend for document content then perform the configuration change above before starting to store documents, as otherwise already existing documents stored in the default location will no longer be accessible.

On **Windows**, accessing a shared network folder unfortunately isn't as straightforward. By default, the ARIS agent will be running with the local SYSTEM user (a.k.a. "LocalSystem" or "NT AUTHORITY\SYSTEM"), which in most environments doesn't have access to network shares. There are various workarounds for this:

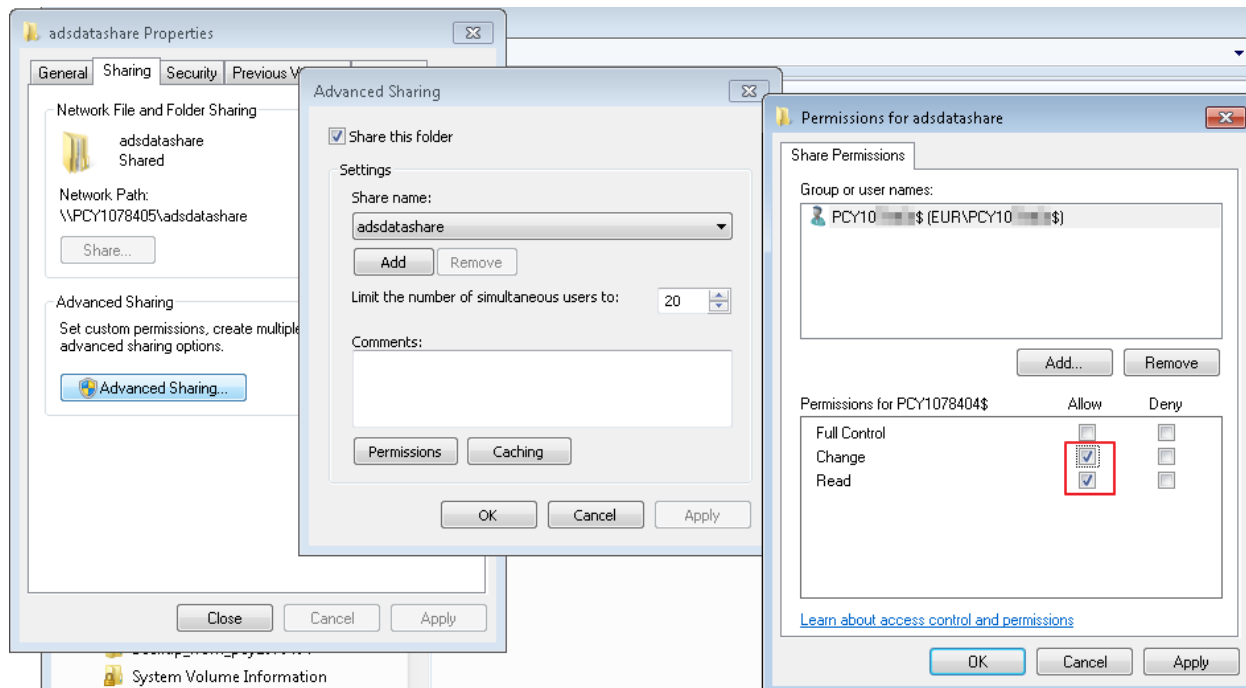
- You **can change the user under which the ARIS agent is running** to a user who has access to the network share (e.g., a user in your Windows domain). This approach is **not recommended** and will therefore not be described in detail here.
- Assuming that the machines on which the ADS runnable instances are running are part of a Windows domain, you can **grant each machine's system user access to the share**. To actually be able to add a permission for a computer, go to the "Advanced sharing" dialog for the shared folder, then click "Permissions", then "Add", then "Object Types", then check the checkbox "Computer", as illustrated in the figure below.

FIGURE 47



You can then search for the computer name of the machine(s) on which ADSAdmin will be running, add them to the list of permitted users and assign both "Read" and "Change" to each of them, as shown in the figure 48:

FIGURE 48



To make ADS use the share, you need to set the configure parameter `JAVA-Dcom.aris.ads.filesystem.path` to point to the share's location, again either directly when configuring the runnable with a command like this

```
configure adsadmin_<SIZING> <adsInstanceID> JAVA-
Dcom.aris.ads.filesystem.path="//<HOSTNAME>/<SHARENAME>"
```

for example

```
configure adsadmin_m ads JAVA-
Dcom.aris.ads.filesystem.path="//myfileserver.company.com/arisdocuments"
```

or later with a reconfigure command like this

```
reconfigure <adsInstanceID> JAVA-
Dcom.aris.ads.filesystem.path="//<HOSTNAME>/<SHARENAME>"
```

for example

```
reconfigure ads JAVA-
Dcom.aris.ads.filesystem.path="//myfileserver.company.com/arisdocuments"
```

Naturally, **this has to be done for all ADS instances in the installation.**

i Note that we used forward slashes ("/") in the above commands to specify the path to the shared folder, instead of the Windows standard backslash ("\") character. This is because the backslash character has a special role of an escape character in the ACC grammar and needs to be escaped by replacing each backslash with four (!) backslashes. Since in Java applications, backward and forward slashes can often be used interchangeably in paths, it is much easier to just use forward slashes.

w If you want to change the folder used as the backend for document content then perform the configuration change above before starting to store documents, as otherwise already existing documents stored in the default location will no longer be accessible.

c) You can also **mount the shared folder as a drive letter in a way that it is available to the system user**. This, however, requires external tools (psexec by Sysinternals) to open a shell in the context of the system user and is therefore not recommended.

8.2.2 Restrictions of the tenant management application (inside UMCAdmin) regarding high availability

The UMCAdmin runnable does not only contain the UMC web application, but also the ACC server and tenant management web applications. The ACC server application provides both the healthcheck UI (which you can access with the URL `http://<arisServer>:<arisHttpPort>/acc/ui` or `https://<arisServer>:<arisHttpsPort>/acc/ui`, respectively) and also the create, backup, restore and delete tenant functionality, on top of which the tenant management web application provides a user interface. The ACC server itself is stateless and can be scaled just like UMC. However, the tenant management application is not fully stateless, as it will store any tenant backups in a local folder (by default the folder "backup" inside the runnable's working directory). Consequently, the backup files are always only available on the instance that created them. So if you have multiple instances of the UMCAdmin runnable (and thus also the tenant management application) in your installation, and if you try to access a backup file, but your request is directed to an instance that does not have this particular file in the local directory, you will encounter an error. There are two solutions to this problem:

Similar to document storage, you can configure the tenant management web application to use a shared network folder. Here, the path to the folder can be specified by setting the configure parameter "JAVA-Dtm.backup.folder", by setting it directly during the configure command

```
configure umcadmin_m <instanceId> JAVA-Dtm.backup.folder =
"<pathToSharedFolder>"
```

resp.

```
reconfigure umcadmin_m <instanceId> JAVA-Dtm.backup.folder =
"<pathToSharedFolder>"
```

Similar to the use of shared network folders as backend for document storage that we described in the previous section, the shared folder can be given as the name of a network share (e.g., "\\server\share"). Alternatively, a local drive letter to which the network share is mounted can be specified. The same considerations regarding the accessibility of the share for the application as discussed above apply here as well.

Examples

```
configure umcadmin_m umcadmin_m JAVA-Dtm.backup.folder = "n:/tenantBackups"
```


resp.


```
reconfigure umcadmin_m <instanceId> JAVA-Dtm.backup.folder =
"<pathToSharedFolder>"
```

will make TM place the backups to the folder "tenantBackups" on drive n, while

```
reconfigure umcadmin_m JAVA-Dtm.backup.folder =
"//fileserver.company.com/arisTenantBackups"
```

will make TM use the share \\fileserver.company.com/arisTenantBackups.

 Note that just as with ADS we used forward slashes ("/") in the above commands to specify the path to the shared folder, instead of the Windows standard backslash ("\") character. This is because the backslash character has a special role of an escape character in the ACC grammar and needs to be escaped by replacing each backslash with four (!) backslashes. Since in Java applications, backward and forward slashes can often be used interchangeably in paths, it is much easier to just use forward slashes.

 If you want to change the folder used as the storage location for tenant backups, then perform the above configuration change before starting to do backups, as otherwise already existing backups stored in the default location will no longer be accessible.

While using a network share is the recommended approach, in particular as it will also make sure that your tenant backups are located "outside" the ARIS installation, for small scale scenarios where you do not need the tenant management (and the ACC server) webapps to be highly available, you can also configure **all but one** UMCAdmin instance such that only the UMC web application inside the runnable will be active, by setting the configure parameter "umcOnly" to true, e.g.

```
reconfigure umcadmin_m umcOnly=true
```

That way, all your backup files will be located on the UMCAdmin instance where this parameter is left unset (resp. is explicitly set to "false").

8.2.3 Restrictions of the simulation component regarding high availability

A simulation will always run in one simulation runnable instance, which will hold the current state of the simulation in memory. So if the simulation runnable crashes, any simulation runs that were just happening on that instance will be aborted and have to be repeated. Still it makes sense to have multiple simulation runnables in the system, first to distribute the load of several parallel simulation runs across multiple instances (which should of course preferably run on different node), but it will also allow the users to restart their simulation run immediately, instead of having to wait for the crashed instance to be restarted.


8.2.4 Restrictions regarding the Dashboarding ("ARIS Aware") component regarding high availability

Due to technical restrictions, with ARIS 10.0 (and 10.0 SR1), the dashboarding runnable cannot be scaled horizontally as straightforwardly as it is possible with most of the other runnables. While it is possible to have several instances for horizontal scaling for load distribution, one selected instance has to be configured to be the master node. The master instance is the only instance that contains the Universal Messaging (UM) component, which is an essential part of the architecture of the underlying MashZone product. The other instances need to have UM removed and have to be set to act as slaves only. If the master instance fails, dashboarding will be unavailable until it is restarted or until a slave instance has been manually configured to take over the master role.

The procedure to configure multiple dashboarding instances in a master-slave setting differs between ARIS version 10.0.0 (GA) and version 10.0.1 (SR1).

For versions 10.0.2 (SR2) and later, Dashboarding will no longer need special settings to work with multiple instances in a single installation, but can be scaled horizontally without further precautions like most of our other applications.¹¹

Configuration for ARIS 10.0 GA

 Note that currently, most of these configuration steps will not be preserved during an update or a reconfigure command and will have to be re-applied.

To configure one instance of the Dashboarding runnable to be the **master**, apply the following steps

1. Stop the runnable if it isn't already stopped.
2. In the file `<runnableWorkingDir>\um\um.properties` you need to add the property "um-server.slaveNodes". The value of this property is the comma-separated list of the hostnames of all slave nodes. For example:

```
um-server.slaveNodes=arisworker1.example.com,arisworker3.example.com
```
3. Delete the folder `<runnableWorkingDir>\um\presto`
4. Start the runnable

To configure an instance of the Dashboarding runnable to be a **slave**, apply the following steps:

1. Stop the runnable if it isn't already stopped.
2. Delete the entire universal-messaging web application by deleting the folder `<runnableWorkinddir>/base/webapps/0um-launcher`
3. In the file `<runnableWorkingDir>\runnable.properties`, set the value of the property "zookeeper.bundle.parts" to be only "dashboarding", i.e., find a line that starts with `zookeeper.bundle.parts` and usually looks like this

```
zookeeper.bundle.parts=dashboarding,universal-messaging
```

and change it by removing the ",universal-messaging" part so that it looks like this

```
zookeeper.bundle.parts=dashboarding
```
4. In the file `<runnableWorkingDir>\um\um.properties` set the value of the property ensure that property `um-core.sessionFactory.realmUrl.1` so that it points to the hostname and the Universal messaging port (parameter `um.port`) of the Dashboarding master instance, e.g. `nsp://arisworker1:14492`
5. Delete the folder `<runnableWorkingDir>\um\presto`
6. Start the runnable

As was said above, in case of a failover event, a new master node must be chosen and configured manually and the configuration of all remaining slave nodes must be adjusted accordingly to point to the new master. Due to the manual effort involved (in particular, you have to restore the folder `<runnableWorkinddir>/base/webapps/0um-launcher` in a former master instance), unless the outage of the node holding the master instance will foreseeable not be resolved quickly, it is in general more advisable to focus one's efforts on bringing the existing master back online.

¹¹ New content for ARIS 10.0 SR2 release is marked with |

Configuration for ARIS 10.0 SR1

Master instance

To configure the dashboarding master instance, you need to add the configure parameter `um-server.slaveNodes` (or you can set these configure parameters later with a reconfigure command). The parameter's value needs to be set to a single string (i.e., surrounded by double-quotes) with the comma-separated list of the full-qualified host names (FQN) and the IP addresses of all slave nodes. For example, if you have two slave instances, the configure parameter would take the form

```
um-server.slaveNodes="<slave1FQN>,<slave2FQN>,<slave1IPAddress>,<slave2IPAddress>"
```

Assuming that your slave nodes are located on the machines `arisworker1.example.com` and `arisworker3.example.com` with IP addresses `10.23.42.1` and `10.23.42.3`, respectively, the configure parameter would be

```
um-server.slaveNodes="arisworker1.example.com,arisworker3.example.com,10.23.42.1,10.23.42.3"
```

Slave instances

To configure a dashboarding instance to act as a slave, you need to set two configure parameters: `"um.deactivated"` needs to be set to `"true"` to - you guessed it - deactivate the Universal Messaging component. Further, the configure parameter `"um-core.sessionFactory.realmUrl.1"` needs to be set to a URL that points to the FQN of the master instance:

```
um-core.sessionFactory.realmUrl.1="nsp://<masterFQN>:${um-server.port}"
```

So assuming that your master instance is located on machine `arisworker2.example.com`, we would set the following two configure parameters on each dashboarding slave instance:

```
um.deactivated = true
```

```
um-core.sessionFactory.realmUrl.1="nsp://arisworker2.example.com:${um-server.port}"
```

Configuration for ARIS 10.0 SR2¹²

For versions 10.0.2 (SR2) and later, Dashboarding will no longer need special settings to work with multiple instances in a single installation, but can be scaled horizontally without further precautions like most of our other applications.

8.2.5 Restrictions regarding the ARCM component regarding high availability

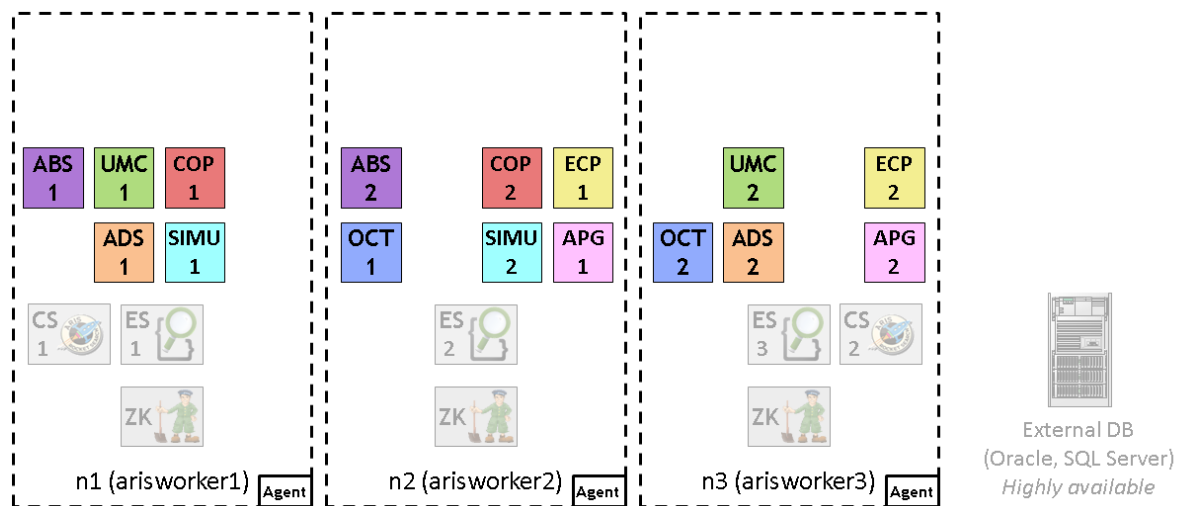
With ARCM, all requests of an active client session are routed to the same ARCM instance (using the preferred routing mechanism explained above). However, here this is not only used for performance optimization, but it is necessary for correct functioning of the application, as the ARCM runnable needs to hold the current client state in memory. This means that if an ARCM runnable crashes, all clients that were handled by this instance will lose any unsaved data and will be returned to the home screen. As long as other ARCM instances are available, the affected users will be rerouted to one of the remaining instances and can resume their work from the last saved state. Users on other ARCM instances will not be affected.

¹² New content for ARIS 10.0 SR2 release is marked with |

8.3 Distributing applications in our example scenario

Continuing with our example scenario from the previous chapters, we will now add redundant application instances for all application types needed for a minimal ARIS Connect installation (i.e., business server (ABS), user management (UMCAdmin), document storage (ADSAdmin), Collaboration (ECP), portal (Copernicus), process governance (APG), Analysis (Octopus) and Simulation.). figure 49 below shows a possible resulting configuration. Note that this is meant only as an example. Depending on the amount of RAM needed and load observed by the different runnables for your specific usage scenarios, you might notice an imbalance in the memory and CPU usage across these nodes, which you can easily adjust for simply by deconfiguring the application runnable on one and configuring it freshly on another node that is currently underutilized (just make sure that the two instances of the same application type are always on different nodes).

FIGURE 49: POSSIBLE HIGHLY AVAILABLE DISTRIBUTION OF THE APPLICATION RUNNABLES FOR AN ARIS CONNECT INSTALLATION



The actual configure commands required to set up this scenario are the following:

```
on n1 configure abs_m abs1
on n2 configure abs_m abs2

on n1 configure umcadmin_m umc1
on n3 configure umcadmin_m umc2

on n1 configure copernicus_m cop1
on n2 configure copernicus_m cop2

on n2 configure ecp_m ecp1
on n3 configure ecp_m ecp2

on n2 configure octopus_m oct1
on n3 configure octopus_m oct2

on n1 configure adsadmin_m ads1 "JAVA-Dcom.aris.ads.filesystem.active" = "false"
on n3 configure adsadmin_m ads2 "JAVA-Dcom.aris.ads.filesystem.active" = "false"

on n1 configure simulation_m simu1
on n2 configure simulation_m simu2

on n2 configure apg_m apg1
on n3 configure apg_m apg2
```

Observe that the only configure parameters that we had to set were the ones to make ADS use the database as backend for document content storage, all other runnables were left at the default configuration as defined in the "m" sizing of their respective app type.

Depending on how much resources you have available on your three worker nodes, they might or not be able to actually hold all these runnables once the system is actively used and the runnables actually start consuming memory and CPU. In that case you have several options (aside from increasing the amount of hardware resources available):

- In the case you are only interested in familiarizing yourself with a distributed installation for evaluation purposes and do not plan to actually use the system for production use (which this scenario isn't meant for anyway), you can always decide to live with the resulting performance loss (in particular if memory run out, the operating system will start swapping, which will significantly decrease performance).
- Alternatively, you can configure the runnables in small ("s") sizing instead of the medium ("m") sizing we were using. That way, the individual runnables' Java heaps cannot grow as much, however, certain uses cases might now fail (in particular bulk operations with large amounts of data, running large and complex reports that handle lots of data etc.). If your interest is focused on one or a few applications (e.g., business server), you can also choose to use "s" sizing for only the less relevant application types and leave the relevant ones at "m". However, be careful **not to have differently sized runnables of the same type** in the system (e.g., one ABS instance in sizing "m", the other in sizing "s"). Having different sizings for different app types is okay (e.g., all UMC, Octopus, and Simulation runnables configured with "s" sizing, and leaving the others at "m" sizing).
- Another alternative is that you might actively decide to **leave some applications without high availability**. For example, if collaboration or document storage are not essential for your scenario's main use cases, you can configure only a single instance of each and risk an outage of the corresponding part of the application if the single instance crashes or the entire node it resides on goes down. In general, the rest of the application will continue to work fine in such a situation (users might have to reload the current portal page, though).
- Leave out (or deactivate) some applications entirely. If you are not using, e.g., simulation, you can just leave out those runnables altogether. The application will in general handle this situation gracefully. However, since the functionalities of some applications depend on each other, it is not often easy to decide, which ones you can do without. Further, certain UI elements of these applications will remain visible even if they are deactivated both in the thin and rich clients, but using these elements will either not give any feedback at all, yield error messages, or give you empty dialogues. It is therefore not recommended to remove any of the above applications completely from a production system. The notable exception are of course the optional app types that we didn't even put into the example configuration above, i.e., the HD server (HDS), the CDF, and the Dashboarding runnables.

8.4 Finishing touches

In addition to the basic configure command to get the runnables onto the nodes as desired, you need to add a few finishing touches.

8.4.1 Adding JDBC drivers

If you are using an external database system (and are not using a single Postgres runnable to keep things simple, albeit not highly available), before you can start your application runnables, you still need to first add the JDBC drivers to each instance (except UMCAdmin, which doesn't use the DB backend service), as we explained in the chapter about relational database backends.

Assuming you used the configure commands as they were given in our example above, and assuming that the JDBC driver JAR was named "driver.jar" and was copied into a subdirectory "jdbc" in your remote repository, you would need to run the following commands:

```
on n1 enhance abs1 with commonsClasspath path "jdbc/driver.jar"
on n2 enhance abs2 with commonsClasspath path "jdbc/driver.jar"

on n1 enhance cop1 with commonsClasspath path "jdbc/driver.jar"
on n2 enhance cop2 with commonsClasspath path "jdbc/driver.jar"

on n2 enhance ecp1 with commonsClasspath path "jdbc/driver.jar"
on n3 enhance ecp2 with commonsClasspath path "jdbc/driver.jar"

on n2 enhance oct1 with commonsClasspath path "jdbc/driver.jar"
on n3 enhance oct2 with commonsClasspath path "jdbc/driver.jar"

on n1 enhance ads1 with commonsClasspath path "jdbc/driver.jar"
on n3 enhance ads2 with commonsClasspath path "jdbc/driver.jar"

on n1 enhance simul1 with commonsClasspath path "jdbc/driver.jar"
on n2 enhance simu2 with commonsClasspath path "jdbc/driver.jar"

on n2 enhance apg1 with commonsClasspath path "jdbc/driver.jar"
on n3 enhance apg2 with commonsClasspath path "jdbc/driver.jar"
```


8.4.2 Adding help files

By default, the application runnables do not contain the files of the online help. These files have to be added using the help enhancement. Different from the JDBC driver enhancement in the previous section, we do not specify the files to add to the runnable using the path, but by using the "artifact coordinates", i.e., the technical name and version of the file.

The basic structure of an enhance command for adding the help for most applications is

```
on <nodeName> enhance <instanceId> with help <helpArtifact> <version> type war
```

where the placeholders in angled brackets have the following meaning


<nodeName> is the logical name of the node, e.g. n1

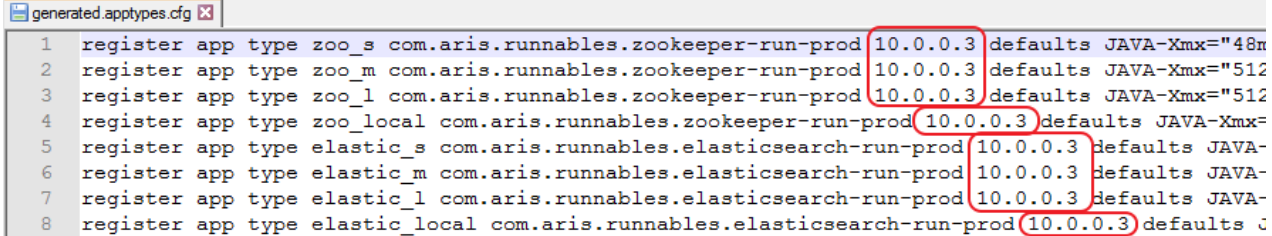
<instanceId> is the instance ID of the runnable you want to add the help to, e.g. "abs1"

<helpArtifact> is the technical name of the actual set of help files, e.g. com.aris.documentation.architect.

The technical name will be different for the help of different runnables.

<version> is the version you are installing, e.g., "10.0.0".

 To find out the actual version of the runnables and help files on your DVD, just have a look at generated.apptypes.cfg. It will contain multiple "register app type" commands, which contain the versions of the runnables - the versions used here are also the versions of the help files. For example, if your DVD contains an ARIS version 10.0.0.3, the generated.apptypes.cfg file would look like this:



```
generated.apptypes.cfg
1 register app type zoo_s com.aris.runnables.zookeeper-run-prod 10.0.0.3 defaults JAVA-Xmx="48m
2 register app type zoo_m com.aris.runnables.zookeeper-run-prod 10.0.0.3 defaults JAVA-Xmx="512
3 register app type zoo_l com.aris.runnables.zookeeper-run-prod 10.0.0.3 defaults JAVA-Xmx="512
4 register app type zoo_local com.aris.runnables.zookeeper-run-prod 10.0.0.3 defaults JAVA-Xmx=
5 register app type elastic_s com.aris.runnables.elasticsearch-run-prod 10.0.0.3 defaults JAVA-
6 register app type elastic_m com.aris.runnables.elasticsearch-run-prod 10.0.0.3 defaults JAVA-
7 register app type elastic_l com.aris.runnables.elasticsearch-run-prod 10.0.0.3 defaults JAVA-
8 register app type elastic_local com.aris.runnables.elasticsearch-run-prod 10.0.0.3 defaults c
```

For the business server runnable (ABS), some parts of the online help have been separated into dedicated files for different languages. The language can be specified with an additional option, the "classifier", in the enhance command. So for adding help to ABS, the structure of the enhance command is

```
on <nodeName> enhance <instanceId> with help <helpArtifact> <version>
classifier <lang> type war
```

The placeholders <nodeName>, <instanceId>, <helpArtifact>, and <version> have the same meaning as above.

<lang> is the language of the help artifact. Currently, there are ABS help files for five different languages: en (English), de (German), fr (French), es (Spanish), ja (Japanese)

You can simply add all languages using several enhance commands, or you can pick only those languages that you need, depending on your user base.

The names (and classifiers) of the help artifacts of the different runnables are the following:

Runnable	Help artifacts (since version)	Languages
ABS	com.aris.documentation.architect	de, en, fr, es, ja pt (since 9.8.2) ru, zh, nl (since 9.8.6)
	com.aris.documentation.architect-scripthelp (since 9.8.2)	-
	com.aris.documentation.architect-methodhelp (since 9.8.2)	- ¹³ , de, en, fr, es, ja pt, ru, zh, nl (since 9.8.6)
ADSAdmin	com.aris.documentation.administration.ads	-
UMCAdmin	com.aris.documentation.administration.administration	-
	com.aris.documentation.administration.tm (for versions 9.8.6 (=9.8 SR6) or later)	-
Copernicus	com.aris.documentation.connect	-
APG	com.aris.documentation.processboard	-
Businesspublisher	com.aris.documentation.publisher.publisher	-
ARCM	com.aris.documentation.arcm.arcm	-

Again assuming that you are following our example exactly, and assuming that you are using repository of a version 10.0.0, the commands to add all available help files to your application runnables are the following:

Help enhancements for the ABS runnables

```
on n1 enhance abs1 with help com.aris.documentation.architect 10.0.0classifier de type war
on n1 enhance abs1 with help com.aris.documentation.architect 10.0.0classifier en type war
on n1 enhance abs1 with help com.aris.documentation.architect 10.0.0classifier fr type war
on n1 enhance abs1 with help com.aris.documentation.architect 10.0.0classifier es type war
on n1 enhance abs1 with help com.aris.documentation.architect 10.0.0classifier ja type war
on n1 enhance abs1 with help com.aris.documentation.architect 10.0.0classifier pt type war
on n1 enhance abs1 with help com.aris.documentation.architect 10.0.0classifier ru type war
on n1 enhance abs1 with help com.aris.documentation.architect 10.0.0classifier zh type war
on n1 enhance abs1 with help com.aris.documentation.architect 10.0.0classifier nl type war
on n1 enhance abs1 with help com.aris.documentation.architect-scripthelp 10.0.0 type war
on n1 enhance abs1 with help com.aris.documentation.architect-methodhelp 10.0.0 type war
on n1 enhance abs1 with help com.aris.documentation.architect-methodhelp 10.0.0 classifier de type
war
on n1 enhance abs1 with help com.aris.documentation.architect-methodhelp 10.0.0 classifier en type
war
on n1 enhance abs1 with help com.aris.documentation.architect-methodhelp 10.0.0 classifier es type
war
on n1 enhance abs1 with help com.aris.documentation.architect-methodhelp
10.0.0classifier fr type war
on n1 enhance abs1 with help com.aris.documentation.architect-methodhelp 10.0.0 classifier ja type
war
on n1 enhance abs1 with help com.aris.documentation.architect-methodhelp 10.0.0 classifier pt type
war
on n1 enhance abs1 with help com.aris.documentation.architect-methodhelp 10.0.0 classifier ru type
war
on n1 enhance abs1 with help com.aris.documentation.architect-methodhelp 10.0.0 classifier zh type
war
on n1 enhance abs1 with help com.aris.documentation.architect-methodhelp 10.0.0 classifier nl type
war
```

¹³ The methodhelp consists of a common part (an artifact without language classifier) and parts for the different languages. The common part should always be enhanced, while the language specific parts can be chosen depending on the languages required by one's user base.

```

on n2 enhance abs2 with help com.aris.documentation.architect 10.0.0 classifier de type war
on n2 enhance abs2 with help com.aris.documentation.architect 10.0.0 classifier en type war
on n2 enhance abs2 with help com.aris.documentation.architect 10.0.0 classifier fr type war
on n2 enhance abs2 with help com.aris.documentation.architect 10.0.0 classifier es type war
on n2 enhance abs2 with help com.aris.documentation.architect 10.0.0 classifier ja type war
on n2 enhance abs2 with help com.aris.documentation.architect 10.0.0 classifier pt type war
on n2 enhance abs2 with help com.aris.documentation.architect 10.0.0 classifier ru type war
on n2 enhance abs2 with help com.aris.documentation.architect 10.0.0 classifier zh type war
on n2 enhance abs2 with help com.aris.documentation.architect 10.0.0 classifier nl type war
on n2 enhance abs2 with help com.aris.documentation.architect-scripthelp 10.0.0 type war
on n2 enhance abs2 with help com.aris.documentation.architect-methodhelp 10.0.0 type war
on n2 enhance abs2 with help com.aris.documentation.architect-methodhelp 10.0.0 classifier de type
war
on n2 enhance abs2 with help com.aris.documentation.architect-methodhelp 10.0.0 classifier en type
war
on n2 enhance abs2 with help com.aris.documentation.architect-methodhelp 10.0.0 classifier es type
war
on n2 enhance abs2 with help com.aris.documentation.architect-methodhelp 10.0.0 classifier fr type
war
on n2 enhance abs2 with help com.aris.documentation.architect-methodhelp 10.0.0 classifier ja type
war
on n2 enhance abs2 with help com.aris.documentation.architect-methodhelp 10.0.0 classifier pt type
war
on n2 enhance abs2 with help com.aris.documentation.architect-methodhelp 10.0.0 classifier ru type
war
on n2 enhance abs2 with help com.aris.documentation.architect-methodhelp 10.0.0 classifier zh type
war
on n2 enhance abs2 with help com.aris.documentation.architect-methodhelp
10.0.0 classifier nl type war

```

Help enhancements for the UMCAdmin runnables

```

on n1 enhance umc1 with help com.aris.documentation.administration.administration 10.0.0 type war
on n1 enhance umc1 with help com.aris.documentation.administration.tm 10.0.0 type war

on n3 enhance umc2 with help com.aris.documentation.administration.administration 10.0.0 type war
on n3 enhance umc2 with help com.aris.documentation.administration.tm 10.0.0 type war

```

Help enhancements for the Copernicus runnables

```

on n1 enhance cop1 with help com.aris.documentation.connect 10.0.0 type war
on n2 enhance cop2 with help com.aris.documentation.connect 10.0.0 type war

```

Help enhancements for the ADSAdmin runnables

```

on n1 enhance ads1 with help com.aris.documentation.administration.ads 10.0.0 type war
on n3 enhance ads2 with help com.aris.documentation.administration.ads 10.0.0 type war

```

Help enhancements for the APG runnables

```

on n2 enhance apg1 with help com.aris.documentation.processboardhelp 10.0.0 type war
on n3 enhance apg2 with help com.aris.documentation.processboardhelp 10.0.0 type war

```

9 Highly available client access to a distributed installation

9.1 The roles of the ARIS load balancer (ARIS LB)

In our chapter about the distribution of microservice applications, we already briefly explained the role of the load balancer both as single point of access for all kinds of clients and its role in the preferred routing mechanism that we use to improve performance. However, we did not really explain the rationale behind routing all clients' requests through the load balancer. In particular in a single-node installation without redundant applications, it might seem that since there is no load distribution to be done among different instances of the same application type, we could do entirely without a load balancer. However, there is an important restriction that makes it necessary to use a load balancer also in this case.

9.1.1 The same origin policy

The ARIS connect portal is based on a framework for the creation of modern interactive web applications that contain content and scripts that can come from different applications, which are then integrated to offer a single, seamless user experience. However, the same origin policy, an essential security model of all browser-based applications, only allows scripts in a page to execute if they were loaded from the same origin as the main page, where "origin" here means the combination of URI scheme (i.e., "http" or "https"), hostname (or IP address) *and port*.

So only if from the perspective of the browser all scripts seem to be coming from the same origin, the application can function properly. If we were to do without a loadbalancer and load content and scripts directly from the different applications via their respective HTTP connector port, the same origin policy would not allow them to be executed (as even if the scheme and port are the same, the ports, and thus the "origin" would be different).

If however all applications are accessed through the load balancer, from the perspective of the browser it will appear as if all scripts are coming from the same server (same scheme, host *and* port), thus the same origin policy is fulfilled and the scripts can be executed to offer the connect portal functionality.

9.1.2 TLS termination proxy

Another role of the load balancer is that of handling de- and encryption of the traffic between the clients and the server installation, i.e., serve as a TLS termination proxy. Instead of having a full end-to-end encryption of client requests and responses from the client down right to the individual application servers and back, the ARIS security model assumes that communication among the runnables is happening in a secure environment and thus need not be encrypted. This not only makes it much easier to use encrypted communication, as you only need to provide an SSL certificate for the load balancer(s), but it also takes the load of en- and decrypting from the applications and puts it to the load balancer(s), which are much better suited to the tasks and can be scaled independently.

9.1.3 Load balancing and high availability

Aside from allowing us to work around the same origin policy and for handling encryption, in a distributed, highly available installation, the load balancer of course has to fulfill its eponymous role - i.e., balancing the load of all client requests more or less equally across all instances of our applications.

However, if all requests have to go through the load balancer, that quite obviously makes it a single point of failure. Fortunately, it is easily possible to simply have multiple load balancers in a distributed installation, and each one will allow access to the complete application. If for example we added a load balancer runnable to each of our three nodes of our example scenario, clients could access the application via the hostnames or IP addresses of any of the worker nodes.

This approach, while possibly acceptable for small scale scenarios, has a few limitations: First, the users themselves have to be made aware of the fact that they can access the same application via multiple hostnames, which can be confusing for them. More importantly, if one of the load balancers (or the node it is running on) fails, all users that chose this machine for accessing ARIS will lose the connection. While most ARIS clients can deal with connection losses between client and server in that any unsaved changes are stored on client side. If the outage persists, the users will have to manually switch over to one of the remaining nodes.

There are basically two solutions to this problem: one is the use of a round-robin DNS. Here a single hostname is chosen for the ARIS application. In the DNS server, not a single IP address is registered, but the IP addresses of all ARIS workers with load balancers running on them. Clients will now use the single hostname to connect to the application, upon which the DNS will return an arbitrary worker node address.

Round robin DNS has several drawbacks. While it in theory allows to hide the fact that the application is distributed to the users to some degree, it will not help in the case of an outage, as clients that were connected to the failed load balancer instance will lose the connection. In addition, a vanilla DNS server doesn't know which of the servers are actually currently available and will happily return the IP addresses of failed instances, too. Further, the caching of DNS lookups on client side is hard to influence. Even if the DNS server would check the workers for availability and take unavailable ones out of the rotation, clients could still try to access the unavailable instances. Since in practice these issues are hard to overcome, using round robin DNS is not recommended.

A more suitable approach is to put a highly available load balancer (HA LB) in front of the ARIS LBs. Such a HA LB is usually a hardware appliance (like "F5") or specialized software on highly available hardware and serves as single point of entry for their users. The details of how to configure such a device are obviously out of the scope of this manual. Basically, all worker nodes holding ARIS LBs have to be added as workers to the HA LB's configuration.

Usually, the HA LB also takes over the role of the TLS termination proxy, i.e., it decrypts the requests from clients that are communicating with the ARIS server in encrypted form via TLS, then forwards the decrypted requests to one of the ARIS LBs, receives the unencrypted response from the application via the ARIS LB, then encrypts this response and sends it back to the clients. If the HA LB takes over this responsibility, you do not need to worry about adding SSL certificates to the ARIS LBs, you can even disable HTTPS entirely on them by setting the configure parameter `HTTPD.ssl.port` to "0".

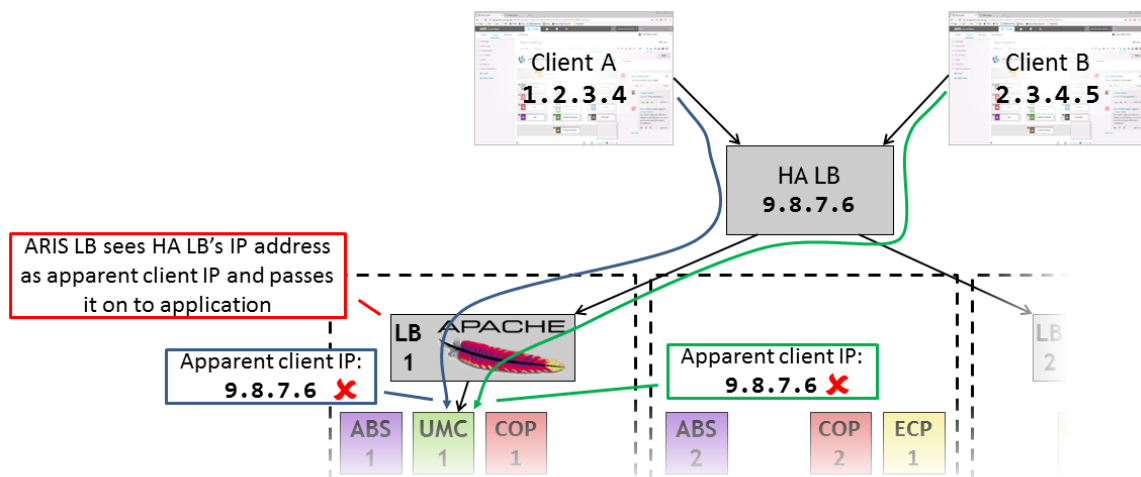
9.1.4 The X-Forwarded-* headers

When using a HA LB in front of the ARIS LBs, additional considerations have to be made, in particular regarding the correct handling of the different **X-Forwarded-*** headers, in particular X-Forwarded-For.

The **X-Forwarded-For (XFF) HTTP** header is a de-facto standard method for identifying the originating IP address of a client accessing a web application (like ARIS) through an HTTP proxy or load balancer.

To understand how this header works, consider the situation shown in Figure 50:

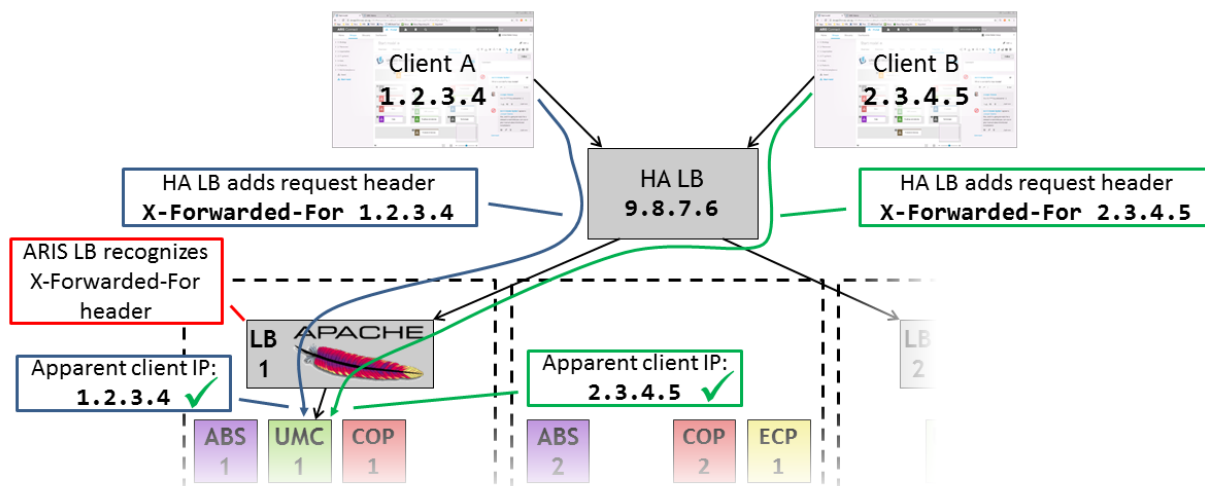
FIGURE 50: WITHOUT X-FORWARDED-FOR, APPLICATIONS SEE THE HA LB'S IP ADDRESS AS APPARENT CLIENT IP ADDRESS



Here, two clients are accessing ARIS through a HA LB. From the perspective of the ARIS LBs, the request is coming from the IP address **9.8.7.6** of the HA LB, and this IP address is passed on to the applications as alleged client IP address. The result is that from the perspective of the application, all requests seem to come from a single client IP address. Amongst other negative effects, this makes troubleshooting more difficult.

When the HA LB is configured to send an **X-Forwarded-For** header whose value is the actual client IP address with each request, the ARIS LB can pass this information on to the applications, resolving the problem, as shown in Figure 51:

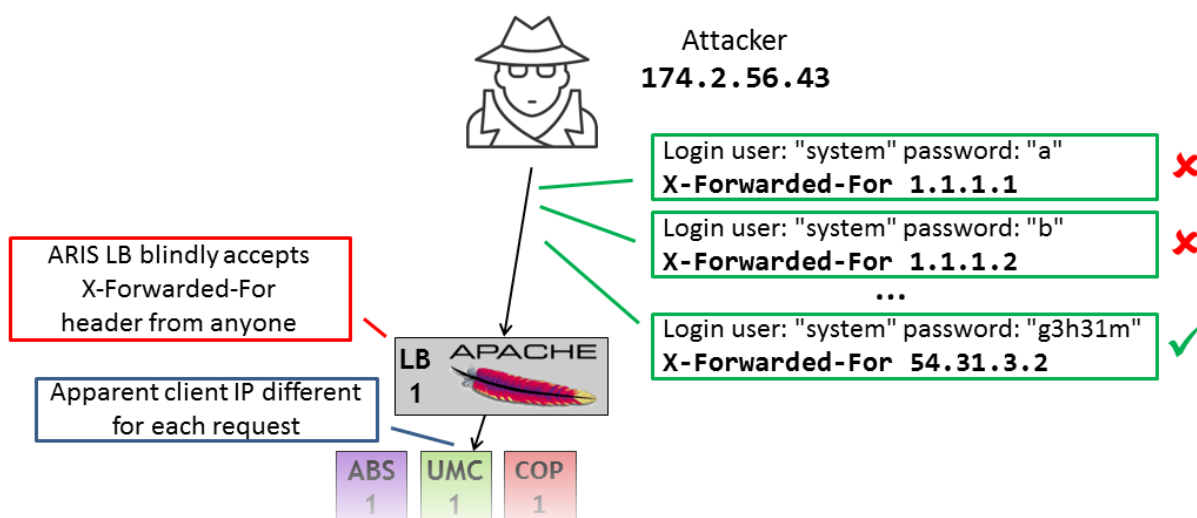
FIGURE 51: HA LB SENDS X-FORWARDED-FOR, APPLICATIONS SEE ACTUAL CLIENT IP ADDRESS



Alas, if the ARIS LB simply accepts **X-Forwarded-For** headers from any source, this would allow an attacker to fake his IP address (as it is seen by the ARIS applications) by itself adding a faux **X-Forwarded-For** header. UMC has an (optional) security mechanism in place to prevent attackers from performing arbitrarily many login attempts to “try out” or “guess” a user’s password (a form of brute-force attack): if too many failed login attempts for a certain user account coming from a certain IP address are observed, that combination of user account and IP address is locked for any further login attempts for a configurable amount of time.

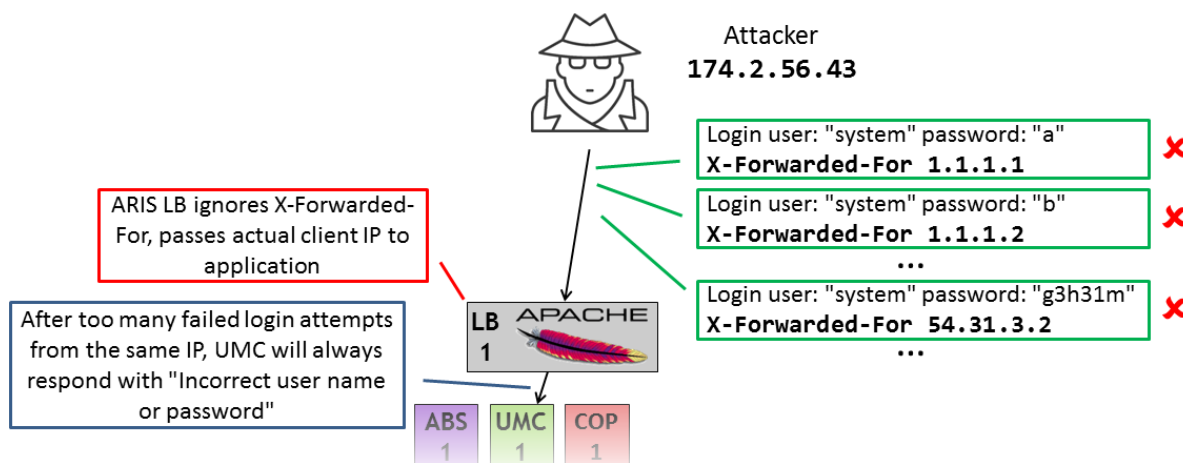
Up to ARIS version 10.0.1, the ARIS LB would accept **XFF** headers in any request. As a consequence, an attacker could set this header to make ARIS believe that his requests are coming from a different than his actual IP address. That way, the protection mechanism against brute force login attacks could be circumvented, allowing the attacker to try out arbitrary many passwords using different “fake” IP addresses in the **XFF** header, a situation is shown in Figure 52:

FIGURE 52: ATTACKER ABUSING X-FORWARDED-FOR TO FAKE IP ADDRESSES



To prevent this problem, starting with ARIS 10.0.2, the ARIS LB, will by default discard any **XFF** headers received from clients except requests coming from localhost (127.0.0.1), as shown in Figure 53:

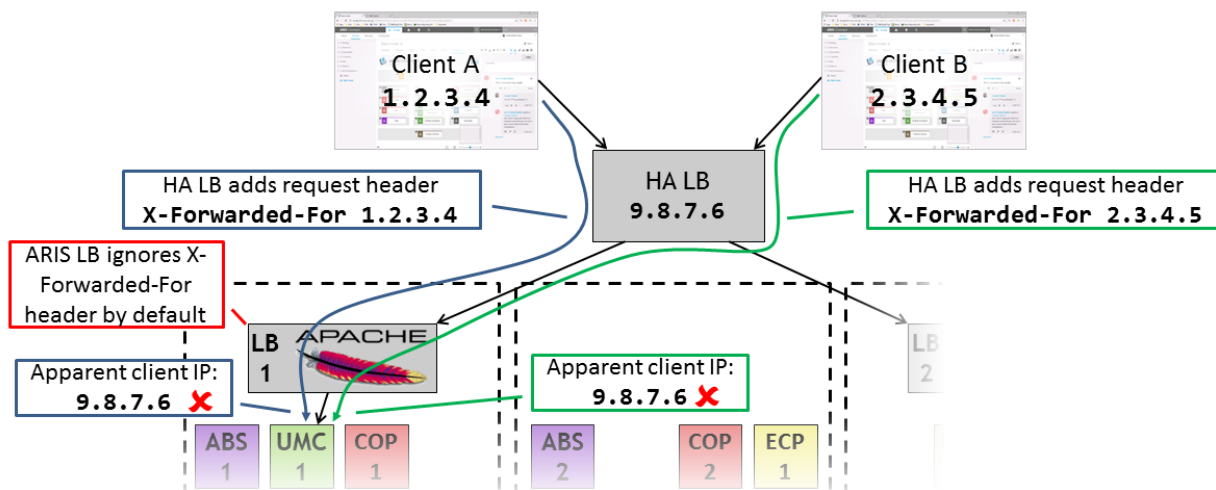
FIGURE 53: ARIS LB IGNORES X-FORWARDED-FOR HEADERS BY DEFAULT, PREVENTING BRUTE-FORCE PASSWORD ATTACKS



Since the applications (in particular UMC) now see the actual client IP, the attack prevention mechanism works as intended: after too many failed login attempts coming from the same IP address, UMC will always respond with "Incorrect user name or password", even if the attacker did find the correct username/password combination.

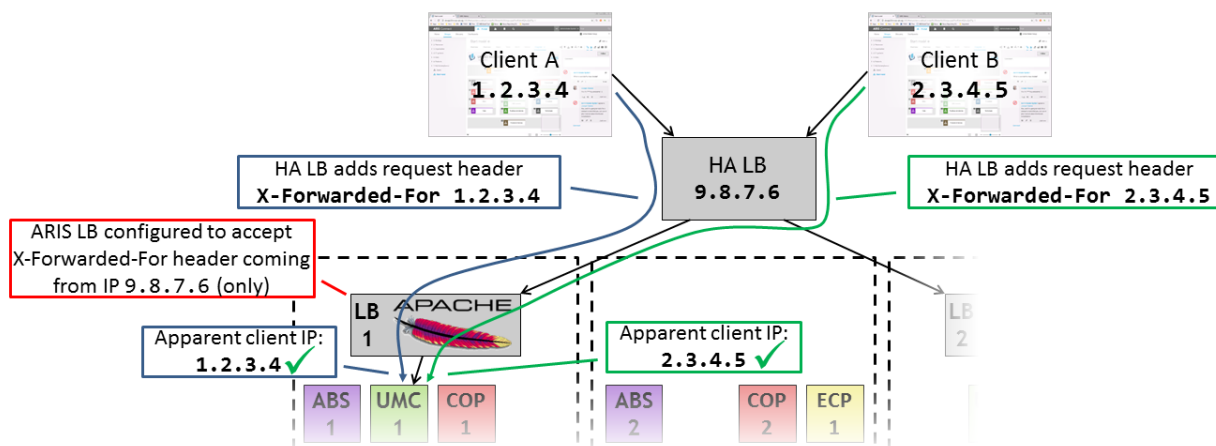
While this change will foil such brute-force password guessing attacks (in particular in simple setups, where the ARIS LB is directly exposed to clients, and thus also to possible attackers), it disables the **X-Forwarded-For** mechanism also for the legitimate use case of putting ARIS behind a HA LB, as illustrated in Figure 54:

FIGURE 54: PROBLEM WITH IGNORED X-FORWARDED-FOR HEADERS IN HA LB SCENARIO



For these scenarios, it is possible to configure the IP address(es) from which the ARIS LB will accept **XFF** headers. To do this, you can use the **HTTPD.X-Forwarded-For.trusted.proxy.regex** load balancer configure parameter, with which you can provide a regular expression (regex, see <https://httpd.apache.org/docs/2.4/glossary.html#regex>) that will make the ARIS LB accept **XFF** headers from any IP address that matches this regex. Note that any backslash characters that are to appear in your regex need to be escaped by replacing them with four(!) backslash characters. For example, to accept **XFF** headers from the IP address **9.8.7.6** of the HA LB in the example above, you would set the **HTTPD.X-Forwarded-For.trusted.proxy.regex** to **"9\\8\\8\\8\\7\\7\\7\\6"**. Observe the use of the quadruple backslashes to escape the dot (".") character that has a special meaning in regular expressions. The resulting scenario is shown in Figure 55:

FIGURE 55: ARIS LB ACCEPTS X-FORWARDED-FOR HEADERS ONLY FROM THE HA LB



To obtain the behavior as it was in previous versions (i.e., the scenario shown in Figure 51, where the ARIS LB would accept the **XFF** header from any client), you can set the **HTTPD.X-Forwarded-For.trusted.proxy.regex** parameter to **"*"** (a regular expression matching every IP address). This is an acceptable setting if you ensure that clients (and thus possible attackers) cannot access the ARIS LBs directly and that your HA LB itself properly handles **XFF** headers coming from clients (i.e., in general by discarding them).

9.2 Configuring load balancer in a distributed, highly available installation

Now that we learned about the roles of the ARIS LBs, we need to make sure that all LBs are configured correctly. First, you obviously need more than one ARIS LB in your installation. Further, when following the recommended approach of putting a HA LB in front of your ARIS LBs, some configuration changes have to be made.

9.2.1 Make the ARIS LBs correctly to play well with a HA loadbalancer

To make our ARIS LBs play together nicely with a HA LB, two things have to be changed in their configuration:

1) Each ARIS LB has to register itself in Zookeeper with the scheme, host and port of the HA LB instead of the using the default (The default is to use the hostname or IP of the machine on which the ARIS LB is running, the ARIS LB's HTTPS port and the scheme "https" if HTTPS is enabled on the ARIS LB; or the HTTP port and scheme "http" if HTTPS is disabled). To make this happen, we need to set a few configure parameters properly on *all* ARIS LBs in the installation. First, you need to set

- the parameter "HTTPD.servername" to the hostname or IP address of the HA LB

If you want to make the URLs sent out by the application use https (which of course requires that you have https enabled on your HA LB), you need to set

- the parameter "HTTPD.zookeeper.application.instance.port" to the port on which our "HA" LB is receiving HTTPS requests (by default that is usually 443) and set
- the parameter "zookeeper.application.instance.scheme" to "https" (**no**, there is no "HTTPD." missing at the beginning of this parameter... just don't ask...)

Otherwise, If you prefer the URLs to point to the unencrypted HTTP endpoint of the HA LB (or if you simply do not support HTTPS on the HA LB), you need to set

- the parameter "HTTPD.zookeeper.application.instance.port" to the port on which our "HA" LB is receiving HTTP requests (by default that is usually 80) and set
- the parameter "zookeeper.application.instance.scheme" to "http"

So why are we actually doing this? The Zookeeper registration of our load balancers is needed, as applications that send out links to the application (in particular links generated into notification emails) need to know through which host the users access the applications, so that they can use this machine's host and the proper scheme and port of the load balancer when building the links. For this, the applications will simply do a lookup of the load balancer's information in Zookeeper. By setting the above configuration parameters, we make our ARIS LBs basically "fake" their entries so that any generated links point to the HA LB instead of pointing directly to any of the ARIS LBs.

2) We also have to disable the FQN redirect that could get in the way by setting the configure parameter "HTTPD.EnforceFQN" to "false".

3) We need to make sure that each ARIS LB will accept the **X-Forwarded-For** header that will be send by the HA LB in front of ARIS. For that, you need to know the IP address of the HA LB which will be used for communication with the ARIS LBs. You then need to create a regular expression that matches only this IP address and escape it as described in the section about the handling of the **X-Forwarded-For** header above, by replacing each backslash character that is to be in your regular expression by four (!) backslashes. For example, if the HA LB has the IP address 9.8.7.6 and you want to make sure that the X-

Forwarded-For header is only ever accepted from this IP address by the ARIS LBs, you would set the ARIS LB's configure parameter **HTTPD.X-Forwarded-For.trusted.proxy.regex** to "9\\.\.8\\.\.7\\.\.6".

To make 1), 2), and 3) happen, use the following configure parameters command on each of your ARIS LBs (either directly when configuring them, or afterwards, within a reconfigure command):

```
HTTPD.servername=<hostnameOfFakeHaLb>
HTTPD.zookeeper.application.instance.port=<httpPortOfHaLb>
zookeeper.application.instance.scheme="http" HTTPD.X-Forwarded-
```

if you want to use HTTP, or, if you want to use HTTPS, to

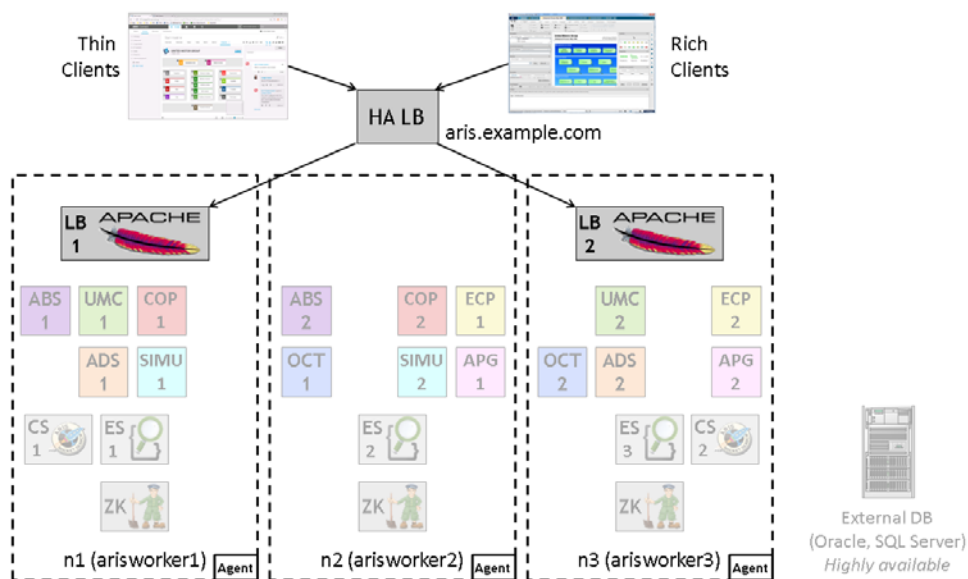
```
HTTPD.servername=<hostnameOfFakeHaLb>
HTTPD.zookeeper.application.instance.port=<httpsPortOfHaLb>
zookeeper.application.instance.scheme="https" HTTPD.X-Forwarded-
```

if you want to use HTTPS, respectively.

9.3 The load balancers in our example scenario

Returning to our example scenario, we will now add two load balancers to our nodes n1 and n3 and configure them as described above so that they work properly with a highly available load balancer (HA LB), as shown in figure 56.

FIGURE 56: REDUNDANT LOADBALANCERS IN OUR EXAMPLE SCENARIO



Assuming that the HA load balancer can be reached by the clients through the hostname "aris.company.com", further assuming that it has the IP address **9.8.7.6**, and assuming that we want any generated links to point to this machine are using HTTPS protocol (i.e., communication is encrypted with SSL/TLS) and the HTTPS standard port, we can configure the load balancers with the following commands (adjust the value given for `HTTPD.zookeeper.application.instance.port` in case you are using a non-standard port for HTTPS on your HA load balancer):

Configuring the ARIS loadbalancers to work with a HA LB with HTTPS support

```
on n1 configure loadbalancer_m lb1 HTTPD.servername="aris.example.com" \
HTTPD.zookeeper.application.instance.port=443 \
zookeeper.application.instance.scheme="https"
on n3 configure loadbalancer_m lb2 HTTPD.servername="aris.example.com" \
HTTPD.zookeeper.application.instance.port=443 \
zookeeper.application.instance.scheme="https"
```

If you prefer that any links generated by the application use HTTP protocol or in case you do not have HTTPS configured on your HA load balancer, you should configure your ARIS LBs with the following commands (adjust the value given for HTTPD.zookeeper.application.instance.port in case you are using a non-standard port for HTTP):

Configuring the ARIS loadbalancers to work with a HA LB without HTTPS support

```
on n1 configure loadbalancer_m lb1 HTTPD.servername="aris.example.com" \
HTTPD.zookeeper.application.instance.port=80 \
zookeeper.application.instance.scheme="http"
on n3 configure loadbalancer_m lb2 HTTPD.servername="aris.example.com" \
HTTPD.zookeeper.application.instance.port=80 \
zookeeper.application.instance.scheme="http"
```

If you do not have access to a HA loadbalancer, the next section describes how to simulate such a device with software (purely for educational purposes).

In case you do not want to use a HA LB at all (neither a real one nor the "fake" software solution), you should instead leave the ARIS LBs at their default configuration, i.e., instead of the configure commands above, you would simply use the following commands:

Configuring the ARIS loadbalancers to work without a HA LB

```
on n1 configure loadbalancer_m lb1
on n3 configure loadbalancer_m lb2
```

9.4 Optional - Simulating a highly available hardware load balancer with software (for educational purposes only)

In case you want to experiment with the behavior and configuration of ARIS when put behind a highly-available (HA) load balancer LB, but do not have access to such a device, this section will provide a description of how to set up a "fake" HA LB by using software (by installing an Apache HTTPD that is configured to act as a reverse proxy in front of the two ARIS LBs in our example scenario). This will give you single point of entry for clients to the application that should behave like a HA LB, but **will of course not be truly highly available** and therefore **must not be used in production environments**, but purely for testing and educational purposes. Consequently, we will also not worry about how to secure the Apache installation, so **do not put this on a machine that is publicly accessible on the internet**. Further, we will also not discuss how to enable support for HTTPS on your "fake" HA LB. Therefore, for this example to work, you will have to use the configure commands for the ARIS LBs found in the section "Configuring the ARIS LBs to work with a HA LB without HTTPS support" in the previous section.

i If you are interested in HTTPS support for your "fake" HA LB, you can just follow one of the many tutorials about this topic that are available, for example the one provided by the Apache documentation¹⁴. In this case you will have to add the configuration settings we describe below to the VirtualHost section for the HTTPS endpoint.

Throughout this section, we will assume that in addition to the three worker machines (arisworker1, arisworker2, arisworker3) that we used so far, you have another machine available on which you can install our "fake" HA LB. We will assume that this machine's hostname is "aris.company.com". Naturally, you will have to replace all hostnames with those in your environment. Further, the descriptions below are for installing Apache on a Windows machine, but the changes to the actual Apache configuration file itself will be the same regardless of the platform it is running on. If you do not have another machine available for running the "fake" HA LB on, you can also install it on one of your three worker machines (or directly on your workstation).

9.4.1 Setting up Apache HTTPD

The first step is to install an Apache HTTPD, which will act as our simulated "hardware" load balancer on the "myaris" machine.

You can use any halfway recent 2.4.x version of Apache HTTPD for Windows.

Install Apache into a directory of your choice. We will refer to this directory as <apacheInstallDir>.

After installing, make sure that Apache itself works, by running the httpd.exe located in <apacheInstallDir>\bin from a shell and accessing the aris.company.com machine from a web browser (usually, the default config will configure Apache to run on port 80, so "localhost" will do).

You should see the familiar "It works!" message:

FIGURE 57



¹⁴Apache SSL/TLS Strong Encryption: How-To - https://httpd.apache.org/docs/2.4/ssl/ssl_howto.html

Now you need to modify the httpd.conf file, which you will find in <apacheInstallDir>\conf directory.

In the initial section of this file where the "LoadModule" statements are found, make sure that the lines for the modules

mod_proxy, mod_proxy_balancer, mod_proxy_html, mod_proxy_http and mod_xml2enc are NOT commented out (the usually are by default, so remove the "#" in front of the line):

```
[...]
LoadModule lbmethod_byrequests_module modules/mod_lbmethod_byrequests.so
[...]
LoadModule proxy_module modules/mod_proxy.so
[...]
LoadModule proxy_balancer_module modules/mod_proxy_balancer.so
[...]
LoadModule proxy_html_module modules/mod_proxy_html.so
LoadModule proxy_http_module modules/mod_proxy_http.so
[...]
LoadModule xml2enc_module modules/mod_xml2enc.so
[...]
```

Then make sure that "Listen" directive is set to the port you want ARIS to be available (this is for HTTP access, HTTPS access is not in the scope of this guide), e.g.,

```
Listen 80
```

Also make sure that the DocumentRoot directive points to an existing directory, from which you want to serve static content (we actually DON'T want our fake HA LB serve static content aside from what is inside ARIS; but if that directive points to a non-existing or otherwise inaccessible directory, HTTPD will not start). For example, if your <apacheInstallDir> is "C:\Program Files\Apache\Apache24_fakeHWLB", the line should look like this (make sure that the htdocs subdirectory actually exists).

```
DocumentRoot "C:\Program Files\Apache\Apache24_fakeHWLB\htdocs"
```

Finally, add the following magic lines to the end of httpd.conf, which will make HTTPD load-balance all requests to our two ARIS LBs:

```
1 ProxyPass "/" "balancer://myaris/"
2 ProxyPassReverse "/" "balancer://myaris/"
3 ProxyPreserveHost On
4
5 <Proxy "balancer://myaris">
6     BalancerMember "http://arisworker1:80"
7     BalancerMember "http://arisworker3:80"
8 </Proxy>
9
10 <Location "/balancer-manager">
11     SetHandler balancer-manager
12 </Location>
```

Lines 5-8 define the "balancer", i.e., a representation of a group of worker nodes that can handle the requests of the application to which our "fake" HA LB is providing access:

- Line 5 names our balancer (here we name it "myaris" - this is purely an internal name and does not correlate to any hostname etc.),
- lines 6-7 enumerate the "BalancerMembers", i.e., the actual workers, which in case of an ARIS installation are the nodes containing ARIS LBs. Based on our example scenario, where we have two ARIS LBs, one on machine arisworker1 and another on arisworker3, we have two such entries. Each "BalancerMember" entry contains a URL pointing to the hostname and HTTP port of one ARIS LB.

Of course you need to adjust these lines to the actual hostnames of your ARIS worker nodes (and you also need adjust the ports in case you configured the ARIS LBs to use a port other than port 80 for HTTP requests).

Line 1 tells Apache that all requests starting with "/" (i.e. ALL requests to our "fake" HA LB) should be handled by the balancer "myaris". (Do not forget the trailing slash after "myaris".) A more detailed documentation of the ProxyPass directive can be found on the Apache documentation page for `mod_proxy`¹⁵.

Line 2 tells Apache to rewrite certain parts of a response so that they point to our "fake" HA LB itself instead of pointing to the worker. A more detailed documentation of the ProxyPassReverse directive can be found on the Apache documentation page for `mod_proxy`¹⁶.

Line 3 make our "fake" HA LB preserve the host used in client requests and pass it to the downstream server (i.e., our ARIS LBs).

Lines 10-12 add the balancer-manager application, which is accessible on the HA LB at `<fakeHALBHostName>/balancer-manager`.

After making these changes, start your Apache (again using `httpd.exe` in `<apacheInstallDir>\bin`).

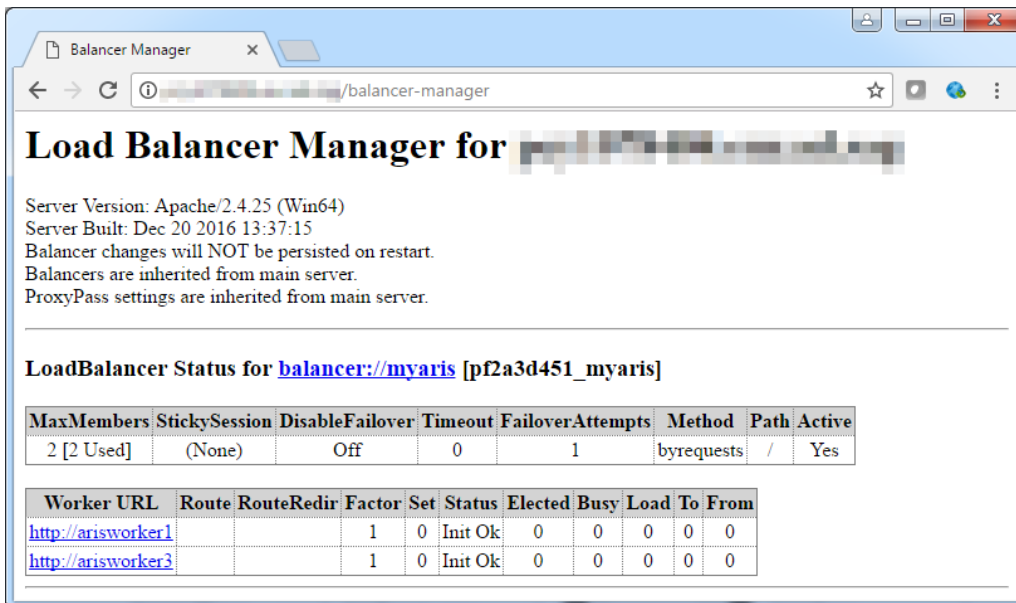
After starting your system (which we will describe in the next chapter), you will then be able to access your ARIS application through the "fake" HA LB, i.e., by using the URL `http://aris.example.com` (resp. the actual hostname of the machine on which you installed the "fake" HA LB) instead of using the hostnames of the worker machines containing ARIS LBs.

To check the state of your "fake" HA LB and the workers, access the balancer-manager application with a browser at `http://<fakeHALBHostName>/balancer-manager` (i.e., assuming your "fake" HA LB runs on `aris.example.com`, the URL would be `http://aris.example.com/balancer-manager`). You should see an output like the following:

¹⁵ Apache Module `mod_proxy` documentation version 2.4 -
ProxyPass https://httpd.apache.org/docs/2.4/mod/mod_proxy.html#proxypass

¹⁶ Apache Module `mod_proxy` documentation version 2.4 -
ProxyPassReverse https://httpd.apache.org/docs/2.4/mod/mod_proxy.html#proxypassreverse

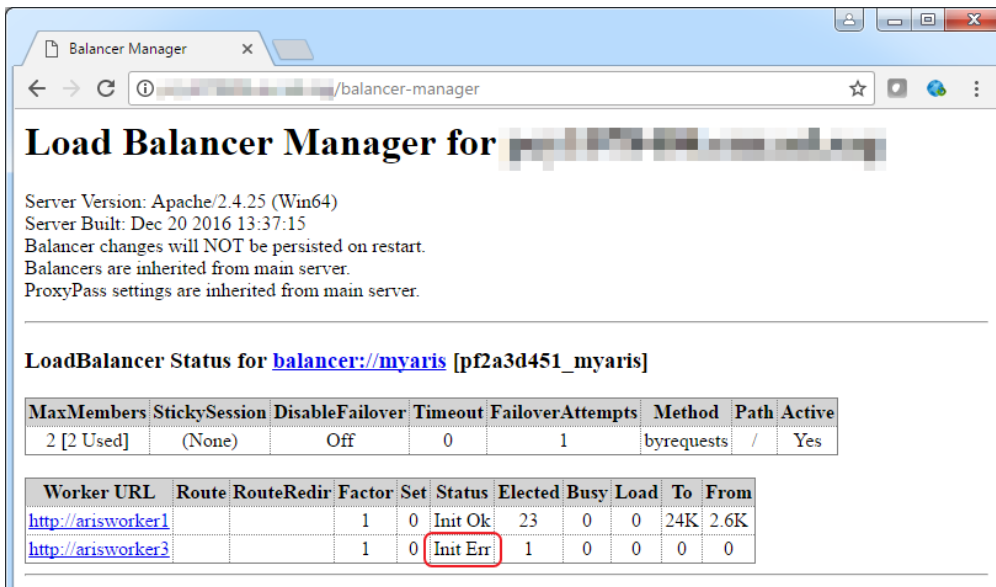
FIGURE 58



Observe that both workers are shown in Status "Init Ok".

You can then test the HA behavior by, e.g., shutting down one of the ARIS LBs. You will observe that clients will still be able to access ARIS without taking notice of the outage - only requests that were still active in the very moment the ARIS LB was shut down will of course be interrupted, but the clients (both rich clients and browser-based thin clients) are built to tolerate this, in the worst case requiring you to reload the page. The balancer-manager application should then show the worker on which you stopped the ARIS LB in Status "Init Err":

FIGURE 59



10 Administration of a distributed installation

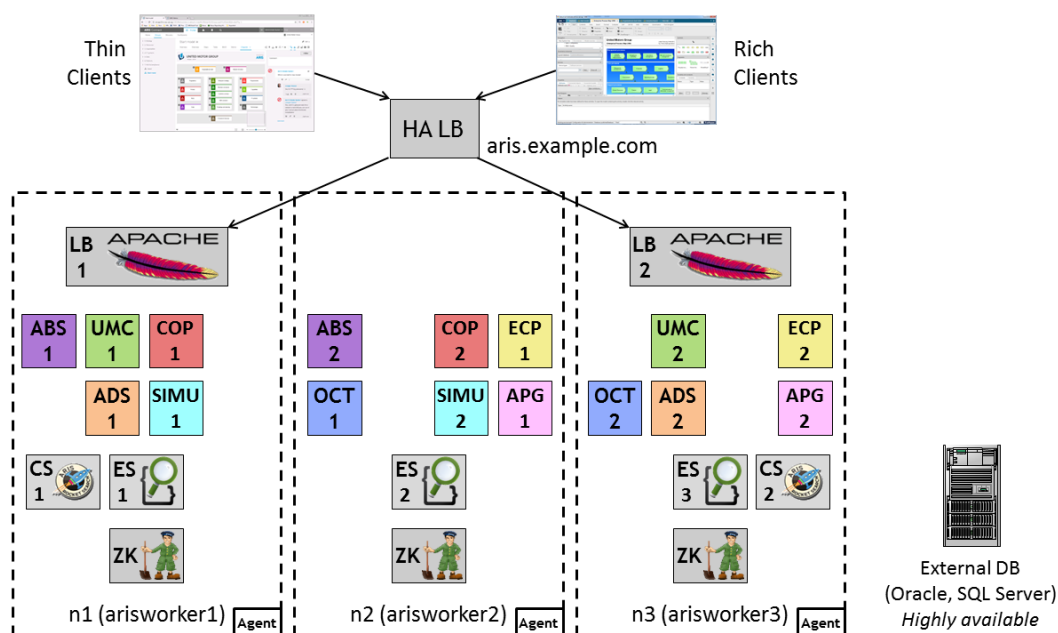
In many aspects, the administration of distributed installation does not differ too much from a single node installation, in particular all administrative steps that do not directly pertain to the runnables or the nodes they reside on. For example, tenants can be created, deleted, backed up and restored using ACC (or the tenant management UI) as before - one only has to make sure that all nodes that belong to the installation (and *only* nodes belonging to the installation, not those of other installations!) are registered in ACC or with the ACC server, as we described it in our chapter about the preparation steps for setting up a distributed system.

We also already explained the basics of how to use ACC to control multiple agents in another section of that same chapter.

10.1 Starting and stopping a distributed system

A single-node installation can be started and stopped using the startall and stopall commands on the node. One could now simply use startall/stopall commands individually for all nodes of a distributed installation, e.g., by running "on n1 startall", "on n2 startall" etc. This however, is not only rather inconvenient, but it also has a significant drawback: while a simple startall (or stopall) command executed on a single node will preserve a certain order in which the runnables are started (or stopped, respectively), it will only preserve this order, among the runnables of *one* node. While the order in which runnables are started does not have to be observed as strictly as it was necessary in early versions of ARIS 9, there are still situations where violating the proper start order can lead to problems. While runnables should no longer fail to start up eventually when a service or other application they depend on is not available, they might still be stuck in STARTING state. For example, consider our distributed example scenario, shown once more in figure 60.

FIGURE 60: OUR DISTRIBUTED EXAMPLE SCENARIO



Assuming that all runnables of the installation are completely stopped, if we were to execute a `startall` command on node n1 first, this would start the Zookeeper instance on that node. Since the other two instances of our Zookeeper ensemble are still stopped, this instance cannot achieve the required quorum (2 out of 3 instances) and thus, while it will reach `STARTED` state, it will not server Zookeeper requests. Next, the "on n1 startall" command would proceed with starting Cloudsearch and Elasticsearch, which at the time of this writing, reach `STARTED` state even without a working Zookeeper access. Next runnable to be started is UMC, which, however, will wait for Zookeeper to become available (since it needs to do a Zookeeper lookup for its backend, i.e., Elasticsearch) and will be stuck in `STARTING` state, preventing the "on n1 startall" command from completing. figure 55 below shows how this would show up in a second ACC console opened up parallel to the running "on n1 startall" command.

FIGURE 61: THE PROBLEM OF USING “STARTALL” ON INDIVIDUAL NODES OF A DISTRIBUTED SYSTEM

```
ACC+ >on all nodes list
Node n1 - 9 installed runnables.
ID      State   Version  Type                                     Extended state
zoo0    STARTED 10.0.0   com.aris.runnables.zookeeper-run-prod  NOT SERVING
es0     STARTED 10.0.0   com.aris.runnables.elasticsearch-run-prod -
cs1     STARTED 10.0.0   com.aris.cip.y-cloudsearch-run-prod    -
abs1    STOPPED 10.0.0   com.aris.modeling.components.y-server-run-prod -
umc1    STARTING 10.0.0   com.aris.umcadmin.y-umcadmin-run-prod  -
cop1    STOPPED 10.0.0   com.aris.copernicus.copernicus-portal-server-run-prod -
ads1    STOPPED 10.0.0   com.aris.adsadmin.y-adsadmin-run-prod  -
lb1     STOPPED 10.0.0   com.aris.runnables.httpd.httpd-run-prod -
simu1   STOPPED 10.0.0   com.aris.modeling.components.simulationserver.y-simuserver-run-prod -

Node n2 - 10 installed runnables.
ID      State   Version  Type                                     Extended state
zoo0    STOPPED 10.0.0   com.aris.runnables.zookeeper-run-prod  -
es0     STOPPED 10.0.0   com.aris.runnables.elasticsearch-run-prod -
pg      STOPPED 10.0.0   com.aris.runnables.PostgreSQL-run-prod -
abs2    STOPPED 10.0.0   com.aris.modeling.components.y-server-run-prod -
cop2    STOPPED 10.0.0   com.aris.copernicus.copernicus-portal-server-run-prod -
ecp1    STOPPED 10.0.0   com.aris.runnables.ecp-run-prod        -
lb2     STOPPED 10.0.0   com.aris.runnables.httpd.httpd-run-prod -
oct1    STOPPED 10.0.0   com.aris.octopus.y-octopus_server-run-prod -
simu2   STOPPED 10.0.0   com.aris.modeling.components.simulationserver.y-simuserver-run-prod -
app1    STOPPED 10.0.0   com.aris.age.age-run-prod              -

Node n3 - 9 installed runnables.
ID      State   Version  Type                                     Extended state
zoo0    STOPPED 10.0.0   com.aris.runnables.zookeeper-run-prod  -
es0     STOPPED 10.0.0   com.aris.runnables.elasticsearch-run-prod -
cs2     STOPPED 10.0.0   com.aris.cip.y-cloudsearch-run-prod    -
umc2    STOPPED 10.0.0   com.aris.umcadmin.y-umcadmin-run-prod  -
ecp2    STOPPED 10.0.0   com.aris.runnables.ecp-run-prod        -
ads2    STOPPED 10.0.0   com.aris.adsadmin.y-adsadmin-run-prod  -
```

Note that the single started Zookeeper instance is shown with a "NOT SERVING" extended state (since it cannot get a quorum) and that the UMCAdmin instance on node n1 is stuck in `STARTING` state.

Obviously, the proper start (stop) order needs to be maintained not merely among the runnables of the individual nodes, but across all runnables belonging to the installation. This is what the **on all nodes startall** and **on all nodes stopall** commands do: using information about the start order specified among all different runnable types (using the ACC config parameter "app.type.order"), these commands will start all instances of each app type across all nodes, before proceeding with instances of app types coming later in the app.type.order.

So in order to start (resp. stop) a distributed system (and of course assuming that the ACC was started with a node file, such that all nodes that belong to the distributed installation are known to ACC), you simply need to use the **on all nodes** variant of **startall** (resp. **stopall**):

Starting a distributed application

```
on all nodes startall
```

Stopping a distributed application

```
on all nodes stopall
```

Like their single node equivalents, these commands will indicate the current progress of the operation, i.e., which runnables are currently being started (resp. stopped).

10.2 Collecting and deleting log files

Another operation that you might need eventually, in particular when you have to involve ARIS support to solve a problem with the system is to collect all log files of your distributed installation.

For quite some time now, one can use the collectLogFiles script (a batch file on Windows, or a Bash script on Linux) that collects all log files on a single node. However, to run such a script, one needs to have access to each machine on the operating system level (i.e., log in remotely), then run the script, and then retrieve the resulting ZIP file. This can be getting quite tedious when working with a multinode installation.

A similar problem is that on a system that is used intensively, log files tend to grow to considerable size, eventually consuming large amounts of disk space. So far, you could either manually delete the biggest log files, or you could use the "deleteLogFiles" script (again, coming as a batch on Windows, and as a Bash script on Linux) (but only on a stopped system!). Again, this task requires you to log in remotely to each machine and is thus rather inconvenient.

Since ARIS 9.8.6, ACC offers two new commands, "collect log files" and "delete log files", respectively, that allows you to conveniently collect (resp. delete) the log files of the agent and all runnables on individual nodes or on all nodes registered with ACC.

10.2.1 The "collect log files" command

The "collect log files" command allows to collect the log files (and also certain configuration files etc. which might help R&D to analyze an issue) of all runnables and the agent (and on Windows also the log files of the setup) on the current node, an explicitly specified node, or on all nodes, and put them into a ZIP file that is then placed on the machine running ACC.

The basic syntax of the command is

```
("on all nodes" | "on" <nodeName>)? collect log files ("of
<agentOrInstanceId> ("," <agentOrInstanceId>*) ("to" <pathToResultFile>)?
```

You can run the command either on a specific node name (by prefixing it with the clause "on <nodeName>"), on all nodes (by prefixing it with "on all nodes", or on the current node (by omitting either clause).

Further, you can optionally specify for which runnables and/or the agent you want to collect the log files by using the optional "of" clause. If the clause is omitted, the logs of the agent and all runnables on the selected node(s) are collected.

Finally, you can optionally specify the name of the ZIP file into which the log files should be placed, by using the optional "to" clause. If the clause is omitted, the ZIP file will be placed in the current working directory of ACC.

Examples

```
collect log files
```

Collects the log files of all runnables and the agent (and on Windows also those of the setup) of the current node and puts them into a ZIP file located in the ACC's working directory.

```
on node n1 collect log files to c:\temp
```

Collects the log files of all runnables and the agent (and on Windows also those of the setup) of the node with the logical node name "n1" and puts them into a ZIP file located in the directory c:\temp.

```
on all nodes collect log files
```

Collects the log files of all runnables and the agent (and on Windows also those of the setup) of the all nodes currently registered with ACC and puts them into a ZIP file located in the ACC's working directory.

```
collect log files for agent, abs_m
```

Collects the log files of the agent (and on Windows also those of the setup) and the runnable named "abs_m" from the current node puts them into a ZIP file located in the ACC's working directory.

10.2.2 The "delete log files" command

The "delete log files" command allows to delete the log files of all runnables and the agent (and on Windows also the log files of the setup) on the current node, an explicitly specified node, or on all nodes.

Note that the log files of some runnables cannot be deleted while the runnable is not in STOPPED state, in particular the logs of the Postgres and the Loadbalancer runnable. Deleting of these log files will be skipped and a warning will be issued.

The basic syntax of the command is:

```
("on all nodes" | "on" <nodeName>)? "force"? delete log files ("of  
<agentOrInstanceId> ("," <agentOrInstanceId>)*)
```

You can run the command either on a specific node name (by prefixing it with the clause "on <nodeName>"), on all nodes (by prefixing it with "on all nodes", or on the current node (by omitting either clause).

Further, you can optionally specify for which runnables and/or the agent you want to delete the log files by using the optional "of" clause. If the clause is omitted, the logs of the agent and all runnables on the selected node(s) are deleted.

Finally, you can optionally specify the "force" keyword, to skip the security query that normally requires you to confirm the command. This is useful in particular if you want to run the command from within scripts.

10.3 Updating a distributed installation

Please refer to the ARIS Update Cookbook which has extensive coverage for manually updating single-node and distributed, multi-node installations.

ABOUT SOFTWARE AG

The digital transformation is changing enterprise IT landscapes from inflexible application silos to modern software platform-driven IT architectures which deliver the openness, speed and agility needed to enable the digital real-time enterprise. Software AG offers the first end-to-end Digital Business Platform, based on open standards, with integration, process management, in-memory data, adaptive application development, real-time analytics and enterprise architecture management as core building blocks. The modular platform allows users to develop the next generation of application systems to build their digital future, today. With over 45 years of customer-centric innovation, Software AG is ranked as a leader in many innovative and digital technology categories. Learn more at www.SoftwareAG.com.

© 2017 Software AG. All rights reserved. Software AG and all Software AG products are either trademarks or registered trademarks of Software AG. Other product and company names mentioned herein may be the trademarks of their respective owners

