

Institute of Architecture of Application Systems

University of Stuttgart
Universitätsstraße 38
D-70569 Stuttgart

Bachelorarbeit Nr. 12345

Self-Containment Packager Framework for TOSCA Cloud Service Archives

Yaroslav Nalivayko

Course of Study: Softwaretechnik

Examiner: Prof. Dr. Dr. h. c. Frank Leymann

Supervisor: M.Sc. Michael Zimmermann

Commenced: 10. February 2017

Completed: 10. August 2017

CR-Classification: I.7.2

Abstract

In recent years, cloud computing is gaining more and more popularity. **Topology and Orchestration Specification for Cloud Application (TOSCA)** provides a whole range of tools for managing and automation of cloud application. The TOSCA standard adds an additional level of abstraction to the cloud applications, in other words, a layer between the external interfaces of cloud application and cloud service provider's API. With the help of TOSCA, it's possible to describe several models of interaction with many different APIs, what allows to automate the rapid redeployment between providers, which are using completely different API. The University of Stuttgart implemented this specification in the runtime environment named **OpenTOSCA**. Description of a cloud application is stored in the **Cloud Service ARchive (CSAR)**, which contains all components necessary for a cloud application life-cycle.

Cloud systems are often described in such way that during they deployment, additional packages and programs need to be downloaded via the Internet. Even with a single server, this can slow down the deployment of cloud application. And if cloud application consists of a large number of servers, each of them downloading a large amount of data during the deployment, this can significantly increase both time and money consumption. During this work a software solution which will eliminate external dependencies in CSAR, resupply them with all packages necessary for deployment and also change the internal structure to display the achieved self-containment will be developed and implemented. For example, all commonly used "apt-get install" commands, which download and install packages, must be removed. Appropriate package must be downloaded and integrated into CSAR structure. Furthermore, all depended packages needed for new packages must also be added recursively.

This Document considers the concept and architecture for mentioned framework. In addition some aspects of implementation will be described.

Contents

1	Introduction	13
2	Basis	17
2.1	Cloud computing and cloud application	17
2.2	Topology and Orchestration Specification for Cloud Applications	19
2.3	OpenTOSCA	21
2.4	Package management	23
2.5	Configuration management tools	24
2.6	Encapsulation of CSAR	26
3	Requirements	29
4	Concept and Architecture	33
4.1	Concept	33
4.2	Architecture	39
5	Implementation	43
5.1	Global elements	43
5.2	References resolver	44
5.3	Language modules	46
5.4	Package manager modules	47
5.5	Package Handler	48
5.6	Topology handling	49
6	Add new package manager module	51
6.1	Aptitude	51
6.2	Implementing the new package manager module	52
6.3	Integrating Aptitude into the Bash module	53
7	Validation	55
7.1	Input CSAR	55
7.2	Processing	55
7.3	Displaying with Winery	56
7.4	Check artifacts	58

8 Summary	61
Bibliography	67

List of Figures

2.1	OpenTOSCA runtime environment control flow	22
2.2	TOSCA topology presented by <i>Winery</i>	23
2.3	Package management	25
4.1	General description of the software's work flow	33
4.2	An example tree describing how to find Service Templates and Node Templates for a given script	35
4.3	An example scheme representing several language modules containing package manager modules	36
4.4	Preprocessing: decompression, adding files and generating dependencies	40
4.5	The data flow scheme between language modules, package manager modules and the package handler.	42
6.1	A command line visual interface for the <i>aptitude</i> package manager. . . .	51
7.1	Processing by the framework.	56
7.2	Source CSAR represented by <i>Winery</i>	56
7.3	The output CSAR represented by <i>Winery</i>	57
7.4	The output CSAR represented by <i>Winery</i> , some nodes moved manually.	58
7.5	An ansible playbook's execution process	59

List of Listings

4.1	Unreadable bash script	36
5.1	Creating of a new Node Template	50
6.1	The <i>aptitude</i> inherited from the <i>PackageManager</i> abstract class	52
6.2	The <i>aptitude</i> module with some common elements	52
6.3	The <i>aptitude</i> <i>proceed</i> function	53
6.4	The <i>aptitude</i> line parser	54
7.1	Check bash installation script	58
8.1	Description for the script Artifact Type definition	63
8.2	Abstract language model	64
8.3	Abstract package manager model	65
8.4	Ansible proceeding	66

List of Abbreviations

API Application **P**rogramming Interface. 13

CSAR Cloud Service **AR**chive. 20

TOSCA Topology and **O**rchestration **S**pecification for **C**loud **A**pplication. 13

1 Introduction

Cloud applications market is increasing with great speed. Global annual growth is about 15%. [Sta16] Furthermore one can observe the growth in the number of firms which are using cloud applications. And it concerns not only some big corporations but also many small companies. [Bun14; Bun17]

One of the most important reasons for the development of cloud applications is the economy of resources. It is much easier and often cheaper to rent a part of another's big mainframe than to maintain one's own server. It is also easier and cheaper to send a small package by mail than to keep your own car (server) and driver (administrator) for rare traffic as well.

The growing popularity of cloud applications makes the automation and the ease of management increasingly important. Management is understood as deployment, administration, maintenance and the final roll-off of cloud applications.

The common problem of cloud applications is *affection*. The transfer of a cloud application configured to interact with the **A**pplication **P**rogramming **I**nterface (API) of one provider to work with another provider and another API is a difficult but important task. The ability to move a cloud application to the more suitable provider quickly is a key to the development of competition and reducing the cost of maintenance.

Topology and **O**rchestration **S**pecification for **C**loud **A**pplication (TOSCA) [13] provides an opportunity to solve this problem. TOSCA defines the language which makes cloud application's description and management portable and interoperable. The use of TOSCA allows to simplify and automate the management of cloud applications by different providers. According to TOSCA standard a cloud application is stored in a **C**loud **S**ervice **A**Rchive (CSAR). This archive contains the description of a cloud application, its external functions and internal dependencies and the data for the deployment and operation.

OpenTOSCA [Stu13] is an open source constantly improving and expanding ecosystem (runtime environment) for TOSCA standard developed by the University of Stuttgart. OpenTOSCA processes data in CSAR format and performs specified actions. These actions often contain links to external packages and programs which will be subsequently downloaded over the Internet for deployment of a cloud application,. These downloads can add expenses to the time required to download packages, money spent on rent of an idle server and Internet traffic for megabytes of pre-known data. If a cloud application consists of the only deployed server, this may mean a few seconds of delay. But when an application deploys a large number of linked servers (a cloud system), the costs can

increase significantly.

The other problems of external dependencies are security and stability. To ensure the security of information some firms restrict the Internet access. In other networks the Internet access is extremely limited. (For example, there can be no broadband access, slow communication only over a satellite at certain hours, etc) An attempt to deploy a cloud application with external dependencies in such networks may not succeed.

To solve these problems a software solution for removing external dependencies in CSARs will be developed and implemented during this work. This software will analyze a CSAR, identify dependencies to external packages and resolve them by downloading the necessary data to install the package (as well as data for all depended packages) and adding them to the CSAR's structure. The simplest example is to find in given CSAR all the commands like "apt-get install package", delete these commands, download the package and all depended packages and add them to the CSAR.

This software must easily be expanded (in other words - to be a framework) since it is impossible to predict and describe all possible types of external dependencies. The output of the framework is a CSAR which contains additions to the original structure, like all the packages necessary for the deployment of the cloud application, with the minimum possible level of access to the Internet during operation.

Structure

The work structure is as follows:

Chapter 2 – Basis. This chapter explains the basic terms of this work, which include definitions and descriptions of cloud applications (section 2.1), TOSCA standard (section 2.2), OpenTOSCA environment (section 2.3) and packet management (section 2.4).

Chapter 3 – Requirements. It clarifies requirements for the framework.

Chapter 4 – Concept and Architecture. The main concepts as well as architecture of the framework are explained and illustrated in chapter 4.

Chapter 5 – Implementation. This chapter contains the description of the implementation. It explains the design and development of individual components of the framework.

Chapter 6 – Add new package manager module. New package manager will be added in this chapter to proof ease of extensibility.

Chapter 7 – Validation. Output of the framework will be checked here.

Chapter 8 – Summary recaps the results of the work.

2 Basis

In this chapter, the terms used in this work will be explained. These include definitions for a cloud computing and cloud applications, description of a TOSCA standard and its implementation an OpenTOSCA. At the end, a package management and languages for its automation are described.

2.1 Cloud computing and cloud application

In everyday life, you can often hear the phrase "cloud computing", but what is it? Unfortunately, generally accepted definition of cloud computing that describes all possible situations doesn't exist. But in the scientific community, the definition put forward by National Institute of Standards and Technology (NIST) is commonly used. This definition appropriately describes the concept of cloud computing used in this paper, and therefore this definition will be used.

Cloud computing is a model for enabling ubiquitous, convenient, on-demand network access to a shared pool of configurable computing resources (e.g., networks, servers, storage, applications, and services) that can be rapidly provisioned and released with minimal management effort or service provider interaction. [Pet11]

But the computing is to abstract term, for our purpose we need something more practical, like an application. Also, there are no generally accepted definitions of cloud application, but it can be obtained from the definition of cloud computing.

A **Cloud application** is an application that is executed according to a cloud computing model. [tec]

In addition, a short definition of the cloud system, a provider, and a user will be provided. Composite cloud applications which consist of multiple small application will be called a **cloud system**. An owner of the physical platform, where cloud computing takes place is called a **provider**. An owner of the cloud application, renting a provider's platform is called a **user**.

Service models

NIST distinguishes between three types of service models.

- **Software as a Service (SaaS).** The capability provided to the consumer is to use the provider's applications running on a cloud infrastructure. The applications are accessible from various client devices through either a thin client interface, such as a web browser (e.g., web-based email), or a program interface. The consumer does not manage or control the underlying cloud infrastructure including network, servers, operating systems, storage, or even individual application capabilities, with the possible exception of limited userspecific application configuration settings.
- **Platform as a Service (PaaS).** The capability provided to the consumer is to deploy onto the cloud infrastructure consumer-created or acquired applications created using programming languages, libraries, services, and tools supported by the provider. The consumer does not manage or control the underlying cloud infrastructure including network, servers, operating systems, or storage, but has control over the deployed applications and possibly configuration settings for the application-hosting environment.
- **Infrastructure as a Service (IaaS).** The capability provided to the consumer is to provision processing, storage, networks, and other fundamental computing resources where the consumer is able to deploy and run arbitrary software, which can include operating systems and applications. The consumer does not manage or control the underlying cloud infrastructure but has control over operating systems, storage, and deployed applications; and possibly limited control of select networking components (e.g., host firewalls).

Usage

Now cloud computing and application can be found everywhere, and they number constantly grows. [Laz16] They are used for test and development, big data analyses, file storage and so on. Cloud computing allows using resources effectively, to distribute the load to a system from several physical servers and to shift the maintenance to the providers. If service uses a single one physical server and this server will be disabled, then entire service will be completely unavailable too. But if cloud application uses a hundred of physical servers, then disabling of one will not carry such serious consequences. In addition, a user doesn't need to maintain a team of administrators for the event of various problems.

Usually, a user doesn't have a direct access to the infrastructure (servers and operating systems), he uses only provided Application Programming Interface (API). An API

provides a set of methods to communicate with provider's infrastructure. Each provider defines his own set of methods, depending on his area of specialization. On the one hand, this specialization makes easier to work with the provider, but on the other hand, it becomes more difficult to redeploy an application to another provider.

2.2 Topology and Orchestration Specification for Cloud Applications

The OASIS [OAS] Topology and Orchestration Specification for Cloud Applications (TOSCA) standard provides a new way to enable portable automated deployment and management of cloud applications. TOSCA describes the structure of an application as a topology containing components and relationships between them. TOSCA application is a cloud application described according to the TOSCA standard. This standard can be used not only to describe all stages of a cloud application life-cycle but also to serve as a layer between the cloud application and provider's API, allowing to implement a single application suitable for working with different providers.

Structure

TOSCA specification provides a language to define services (described in section Service models) and relationships between them using *Service Templates*. In addition it describes the management procedures which create or modify services using orchestration processes. The description of the TOSCA components used in this work is provided.

A *Service Template* is the main component in a TOSCA structure. It contains an information about the structure (Topology Template) and interfaces (Plans) of the service. There can be many Service Templates within one TOSCA application. *Plans* provide interfaces to manage cloud applications. These components combine management capabilities to create high-level management tasks which can then be executed for fully automated deployment, configuration, management, and operation of the application. Plans are started by external messages and call management operations of the nodes in the topology. A *Topology Template* describes the topology of a cloud application, defining nodes (Node Templates) and relations between them (Relationship Templates). A *Node Template* instantiates a Node Type as a component of a service. A *Node Type* defines the properties of such a component and the operations available to manipulate the component. A *Relationship Template* instantiates a Relationship Type as a relationship between Node Templates in a Topology Template. The Relationship Template indicates that two nodes are connected and defines the direction of the connection. A

Relationship Type defines semantics and properties of the relationship. Any Node Types and Relationship Types can be instantiated multiple times. Those types are like abstract classes in high-level programming languages and Templates are objects of those classes. *Artifact* represents the content necessary for a management such as executables (e.g. a script, an executable program, an image), a configuration file, a data file, or something that might be needed for other executables (e.g. libraries). The TOSCA distinguishes two kinds of artifacts: Implementation Artifacts and Deployment Artifacts. An *Implementation Artifact* represents the executable of an operation described by a Node Type. An *Deployment Artifact* represents the executable for materializing instances of a node. An *Artifact Type* describes a common type of an artifact: python script, installation package and so on. An *Artifact Template* represents information about the artifact. The location of the artifact and other attendant data are stored here. Again, types are classes, templates are objects, and artifacts represent content or a value of an object, but these values of objects can not be changed. *Node Type Implementation* defines the artifacts needed for implementing the corresponding Node Type. For example, if Node Type contains *deploy* and *shutdown* operations, then Node Type Implementation can contain two Implementation Artifacts with scripts for these operations and one Deployment Artifact with data necessary for the deployment. Implementations are like final classes between Node Types and Node Templates, but in the TOSCA standard, the Implementation will be chosen only during execution. Types, Templates, and Implementations defining a TOSCA application are stored in definition document which have the XML format.

Usage

The combination of topology and orchestration in a Service Template defines what is needed to be preserved across deployments in different environments to enable interoperable deployment of cloud services and their management throughout the complete lifecycle (e.g. scaling, patching, monitoring, etc.). This is useful when an application is ported to alternative cloud environments. [OAS13]

To store a TOSCA application a **Cloud Service ARchive** (CSAR) is used. This is a ZIP-file with ".csar" extension that contains all the data needed for instantiation and management of TOSCA application. They include definition documents, artifacts and so on. In this form, a TOSCA application can be processed by a TOSCA runtime environment.

The root folder of any CSAR must contain the "Definitions" and "TOSCA-Metadata" folders. The "Definitions" folder contains definition documents and one of them must define a Service Template. The "TOSCA-Metadata" folder must contain TOSCA metadata in the form of a file with the "TOSCA.meta" name. This metafile consists of name/value pairs, one line for each pair. The first set of pairs describes CSAR itself (TOSCA version, CSAR version, creator and so on). All other pairs represent metadata of files from the

CSAR. The metadata is used by TOSCA runtime environment to process given files correctly.

During this work, such terms as input CSAR and output CSAR will be used. The input CSAR is the CSAR, which can contain external references and will be processed by the framework. The output CSAR is the CSAR, which was processed by the framework and doesn't contain external references.

2.3 OpenTOSCA

OpenTOSCA provides an open source web-based ecosystem for TOSCA applications. It was developed at the University of Stuttgart in october 2012. This ecosystem consists of three parts: a TOSCA **runtime environment**, a graphical modeling TOSCA tool **Winery**, and a self-service portal for the applications available in the container **Vinothek**. [Stu13] Descriptions of the runtime environment and Winery will be provided in more detail.

Runtime environment

The runtime environment enables a fully automated plan-based deployment and management of cloud applications in the CSAR container. First, the CSAR is unpacked and the files are put into the Files store. Then, the TOSCA definitions documents are loaded, resolved, validated, and processed by the Control component, which calls the Implementation Artifact Engine and the Plan Engine. The Implementation Artifact Engine deploys the referenced Implementation Artifacts and stores their endpoints in the Endpoints database. Finally, the Plan Engine binds and deploys the application's management plans. The endpoints of the management plans are stored in the Plans database. [BBH+ 13] This system can be visualized by figure 2.1

Winery

Winery provides a complete set of functions for graphically create, edit and delete elements in the TOSCA topology presented by a CSAR. It consists of four parts: the type and template management, the topology modeler, the BPMN4TOSCA plan modeler, and the repository.

The type, template and artifact management enables managing all TOSCA types, templates and related artifacts. This includes node types, relationship types, policy types, artifact types, artifact templates, and artifacts such as virtual machine images.

The topology modeler enables to create service templates. Service templates consist of

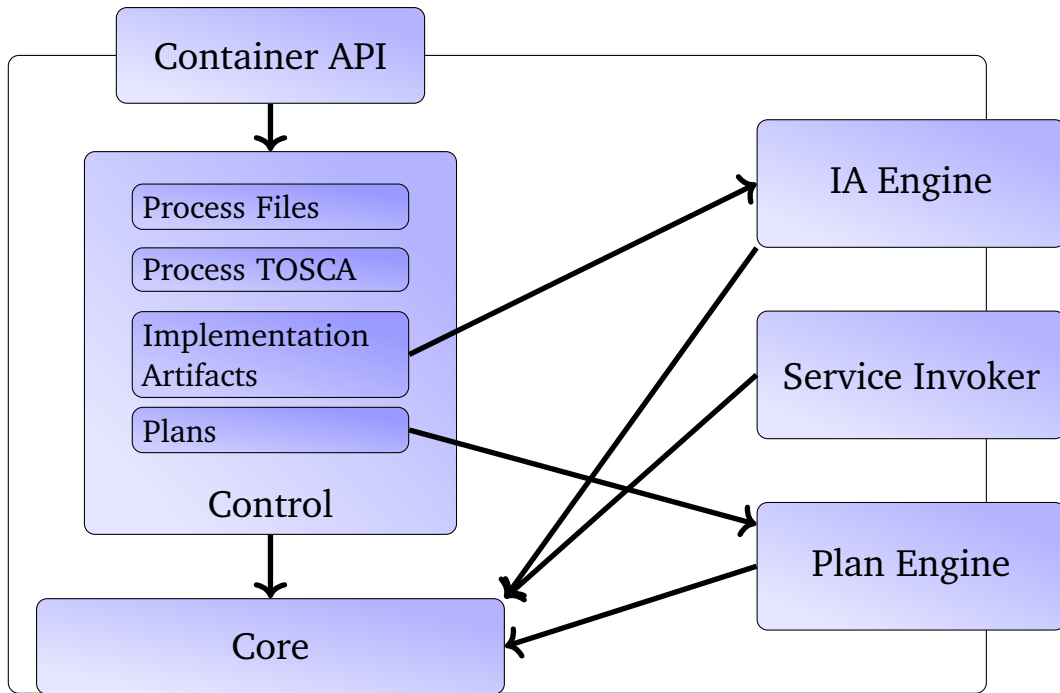


Figure 2.1: OpenTOSCA runtime environment control flow

instances of node types (node templates) and instances of relationship types (relationship templates). They can be annotated with requirements and capabilities, properties, and policies.

The BPMN4TOSCA plan modeler offers web-based creation of BPMN models with the TOSCA extension BPMN4TOSCA. That means the modeler supports the BPMN elements and structures required by TOSCA plans and not the full set of BPMN. The Stardust project offers a Browser Modeler, which supports BPMN. In the context of Winery, this modeler will be extended to support BPMN4TOSCA.

The repository stores TOSCA models and allows managing their content. For instance, node types, policy types, and artifact templates are managed by the repository. The repository is also responsible for importing and exporting CSARs, the exchange format of TOSCA files and related artifacts. [Win]

Winery works under a Tomcat server and therefore a visual interface is available in a browser. An example of the TOSCA topology visualization is presented in the figure 2.2.

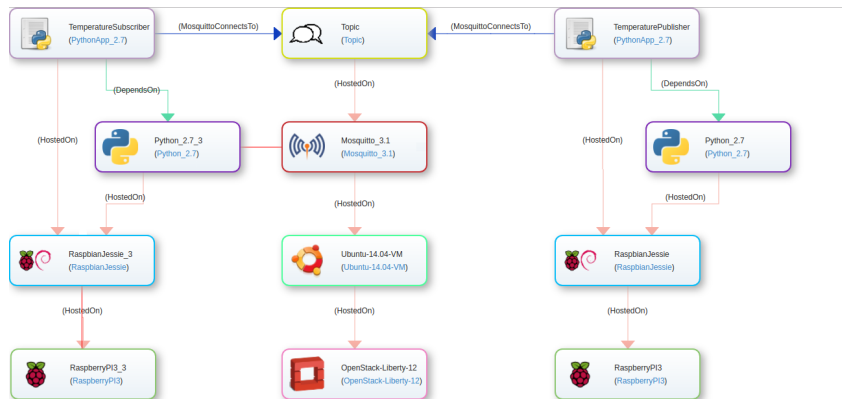


Figure 2.2: TOSCA topology presented by *Winery*.

2.4 Package management

Package is an archive file containing both data for installation of the program component and a set of metadata like name, function, version, producer, and a list of dependencies to other packages. Those program components can present not only a complete program but also a certain component of a large application. For a user, a **package manager** is a set of software tools that automate the process of installing, updating, configuring and removing packages. But from the operating system side, a package manager is used for managing the database of packages, their dependencies, and versions, to prevent erroneous installation of programs and missing dependencies. This task is especially complex in computer systems relying on dynamic library linking. Those systems share executable libraries of machine instructions across packages and applications. In these systems, complex relationships between different packages requiring different versions of libraries result in a challenge colloquially known as "dependency hell". Good package management is vital to these systems.

To give users more control over the kinds of programs that they allow to install on their systems, packages are often downloaded only from a number of software repositories. In Unix systems, a package manager uses official repositories appropriate for the operating system and the device architecture by default, but it's possible to use additional repositories, like third-party repositories or repositories for another architecture.

Package managers distinguish between two types of dependencies: *required* and *preRequired*. Dependency *package1 required package2* indicates that the *package2* must be installed for a proper **operation** of the *package1*. Dependency *package1 preRequired package2* indicates that the *package2* must be installed for a proper **installation** of the *package1*. In these examples, the *package2* is needed for the *package1*, but the *package2* itself can require additional packages. A structure describing all necessary packages and dependencies between them for the given root-package is called a dependency tree. The

dependency type *required* can lead to cycles in a dependency tree what differs it from the normal tree graph structures.

Example dependencies handling

The *apt-get* package manager will be considered to provide an example of a dependencies handling. This application is part of the **Advanced Packaging Tool (APT)** which uses *dpkg* to communicate with an operating system. The system keeps a database of packages and their condition. These relations are presented in figure 2.3.

apt-get has many functions: install, remove, update, autoremove, download and so on. We will consider the install, remove and autoremove operations to present the common algorithm of processing. When a package manager becomes a *package* installation command, it builds a dependencies tree for the *package* and checks the possibility to install these depended packages. For example it must check the compatibility with installed packages. If the check was successful, the *apt-get* downloads and installs the packages beginning with the bottom of the tree. The *package* is marked in database as manually installed and all other packages are marked as automatically installed. That will be helpful during the autoremove operation, when all automatically installed packages will be checked whether they are still needed. After installation of packages from the dependencies tree, the *package* will be ready to work.

A *package* can be deleted during the *remove package* command. That happens only if there are no other packages depending on the *package*. If the deletion is very important, then these packages can be removed too to keep the consistence of the database. The packages needed for the *package* itself will be deleted only by the autoremove command.

2.5 Configuration management tools

To ease the package management, various tools can be used. One can write an executable file, which checks an environment and installs the necessary packages using a package manager. Such files are called scripts and are commonly used. Bash and Ansible used in the framework will be described. Furthermore Chef and CFEngine will be presented too.

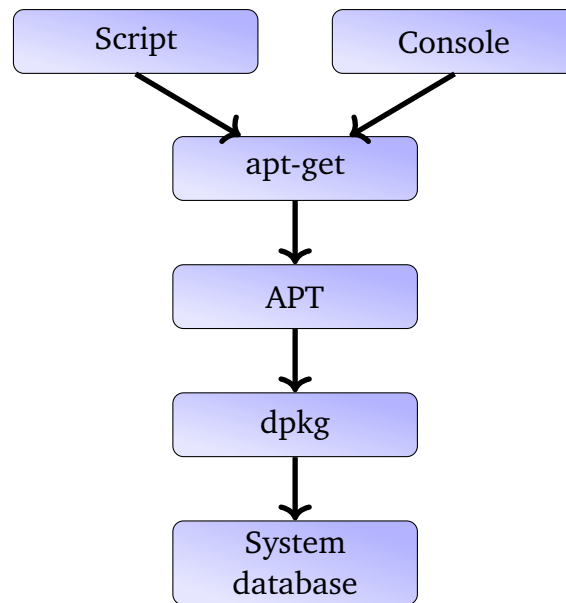


Figure 2.3: Package management

Bash

Bash is a Unix command language written as a free software. It provides enough capabilities to be used as a management tool. In addition Bash denotes a command processor that typically runs in a text window, where a user types commands that cause actions. Instead of typing commands direct into a command line, a script can be executed directly. [Ham11] These scripts can be used to configure a system, install package, create files, check environment and so on. Bash is a very popular and ease language, therefore a huge number of problems have solutions in Bash scripts already.

Ansible

Ansible is an open-source automation engine that automates software provisioning, configuration management, and application deployment. As with most configuration management software, Ansible has two types of servers: controlling machines and nodes. First, there is a single controlling machine which is where orchestration begins. Nodes are managed by a controlling machine over SSH. The controlling machine describes the location of nodes through its inventory. Ansible playbooks express configurations, deployment, and orchestration in Ansible. The playbook format is YAML. Each playbook maps a group of hosts to a set of roles. Each role is represented by calls to Ansible tasks. [com16a]

Chef

Chef is a configuration management tool, which uses Ruby for writing system configuration files called "recipes". They describe how Chef manages applications and utilities and how they are to be configured. These recipes (which can be grouped together as a "cookbook" for easier management) define a series of resources that should be in a particular state: packages that should be installed, services that should be running, or files that should be written. Chef can run in client/server mode, or in a standalone configuration named "chef-solo". In client/server mode, the Chef client sends various attributes about the node to the Chef server. In solo mode the local system will be configured. [Met15]

CFEngine

CFEngine is an open source configuration management system. Its primary function is to provide automated configuration and maintenance of large-scale computer systems, including the unified management of servers, desktops, consumer and industrial devices, embedded networked devices, mobile smartphones, and tablet computers. [Bur13] Configurations are described by "policy" files, which are plain text-files with .cf extension. These files define the necessary state of files, packages, users, processes, services and so on. [CFE]

2.6 Encapsulation of CSAR

According to requirements from the Institute of Architecture of Application Systems (IAAS), the encapsulation must be achieved through the download of external packages and generation of a new TOSCA node for each of them. But it can be interesting to analyze other techniques to encapsulate a CSAR. At first the methods not representing packages in a TOSCA topology and then the methods mirroring packages in a topology we will be described.

Generate custom repositories

It's possible to download all necessary packages and create an own custom package repository for each device used in application. Then one must rework any package installation commands or exchange system preferences to setup an access to the custom

repository. This method introduces minimal changes in a TOSCA structure. The main problem is the creation of the custom repositories. When a TOSCA application consists of many small devices with limited capabilities it can be difficult to start many big custom repositories.

Generate shared repository

Another opportunity is to create a single repository for all devices in a TOSCA application. It can be difficult to choose the right location for such server, but since an application represent the connected system this can redistribute the load to a powerful server. It is difficult to estimate the changes which will occur in a TOSCA topology during the applying of such method.

One node for one package

This method was suggested by IAAS. A new TOSCA node will be created for each downloaded package. All dependencies between packages will be mirrored to a TOSCA topology. This is a very visual method, facilitating the understanding of a TOSCA application and dependencies between packages.

Sets of packages

A set of depended packages from a dependencies tree related to an external reference can be archived and represented in a TOSCA topology as a single node. An installation of such node will lead to an installation of all needed packages. Of course some packages can be saved in different archives redundant, but a small size of a TOSCA application's structure will be achieved. It will be impossible to trace dependencies (since all packages are represented by one node), but it can help to avoid a difficult structure which consists of hundreds of nodes.

3 Requirements

It's necessary to develop new software, which will resolve external references and complement the TOSCA topology to represent the changes in a given CSAR file. Since the main purpose of the software is to Resolve References, further the *RR* can be used as an abbreviation. A primary type of external references is an artifact represented by a bash script which uses an apt-get package manager to install new packages. The exact command for this type of references is "*apt-get install package*". In this example, the *RR* should comment out such commands and add the installation of the **package** and installations of all packages from its dependencies tree as new separate nodes to the TOSCA topology.

A system package manager can be used for this purpose. An example for the *apt-get* package manager is provided. Using the "*apt-get download package*" command, the installation data for the **package** can be downloaded, and using the "*apt-cache depends package*" command, the list of dependencies for the **package** can be obtained. To install a package from the installation **data** the simple command "*dpkg -i data*" can be used. An installation data should be integrated into the TOSCA topology. A new node should be created for each package. During this, step a set of definitions will be created. A common description of a new node is provided via a Node Type, which contains the "install" operation. This operation will be implemented by a Node Type Implementation, which uses Artifact Templates. The created nodes should install packages using the same language, as the original node contained the external reference. In the example with the Bash language, the Artifact Templates will be represented by the deployment artifact with the installation data and by the implementation artifact with the script containing the installation command.

Then the nodes should be instantiated and referenced. For each Node Template which uses an artifact with external references to a **package**, a separate Node Template of the **package**'s Node Type should be defined and referenced. Then for each **package**'s Node Template additional Node Templates for packages from its dependencies tree should be defined and referenced. To define a reference between nodes a Relationship Template will be used and to find out which Node Template uses which artifact the CSAR will be preprocessed. During the preprocessing the topology of the SCAR will be analyzed and internal references found.

To distinguish between languages, package managers and they software handlers, a program component handling a language will be called the language module, and the

one handling a package manager - the package manager module. The Bash module with the apt-get module will be implemented first. The goal is to develop extendable software where new modules can easily be added later. This type of a software is called a framework.

After the minimal configuration with the *Bash* and *apt-get* modules is developed, an *Ansible* module with *apt* package manager module can be added. Ansible scripts are called playbooks. Ansible playbooks and related data are often packed to a zip archive for encapsulation. That makes the ansible module harder to implement since it should not only parse playbooks but additionally unpack archives. Ansible Node Type Implementations will contain only one artifact. This artifact will be an archive containing both the playbook and the installation data.

Handling

There is an example, representing how the framework should process a CSAR.

At start an input CSAR will be extracted to a temporary folder to handle its content. Then the internal structure will be analyzed during the preprocessing. In addition, the common TOSCA definitions not belonging to a specific node (like Artifact Types or Relationship Types) will be added. After that the processing starts. Each file from the input CSAR will be processed by each language module. If a language module accepts a file, then the file is transferred to the package manager modules belonging to the language module. A package manager module will resolve reference by commenting out the installation command and extracting package names from the command. Using the package name the package installation data will be downloaded, the installation script and the TOSCA node created. To create a new TOSCA node for the given package, the definitions for Node Type, Node Type Implementation, and Artifact Templates should be added to the CSAR. A separate Node Template, as well as Relationship Template, will be defined for each depending node. These actions will be recursively repeated for all depending packages from the package's dependencies tree which mirrors the tree to the TOSCA topology. At the end, the meta-file should be updated and the data packed back to the SCAR. This behavior will be described in the chapter 4 and implemented in the chapter 5 in more detail.

Result

As a result of the program's work, an output CSAR will be received. This CSAR must have the same functionality as the input CSAR, but all external references to additional

packages must be resolved. The output CSAR must be able to be deployed properly without downloading these packages over the Internet. In addition, the dependencies trees for packages from new nodes should be represented in the TOSCA topology. In order to validate the output CSAR, the TOSCA topology can be checked and printed by Winery and the defined artifacts can be validated through those installations on a test machine.

Summarize

The developed framework must:

- delete references to external packages from a CSAR.
- add the internal installation of the packages to the CSAR.
- handle different configuration management tools and package managers.
- have an easy to expendable structure.
- represent the packages into the TOSCA topology.
- generate a new CSAR with the same functionality, but without any external references.

4 Concept and Architecture

In this chapter, the concept and the architecture of the framework which can satisfy the requirements will be explained and substantiated. Solutions to some additional problems will be presented.

4.1 Concept

In this section, the main concept of this work are described. The general structure of the framework is represented in the block diagram 4.1. In section 4.1.1, it will become clear how to determine during the preprocessing stage the Node Templates, which use the given artifact. Then language modules and package manager modules functionality is going to be described. In section 4.1.3, it will be expressed how to create a new node for a TOSCA topology. After that, a problem of the determination the architecture of the final platform will be explained and a solution presented. In addition, it will be described, how the results can be validated.

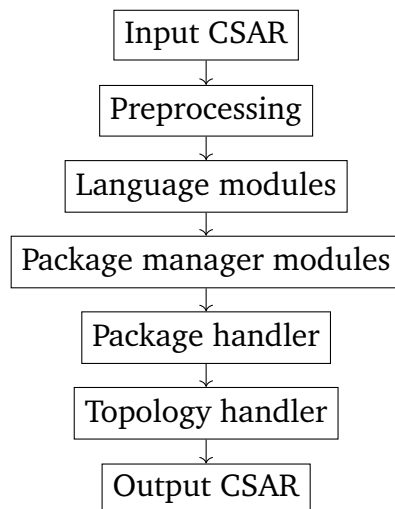


Figure 4.1: General description of the software's work flow

4.1.1 Analysis of existing TOSCA-Topology

To update the TOSCA topology properly, it is necessary to add references from the nodes where external references were to the newly created nodes which resolve the external references. According to the TOSCA standard, references between Node Templates can only be created in the same Service Template. That means that each Node Template which uses artifacts with external references must be found. Furthermore, Service Template where these Node Templates are instantiated must be determined to create there a Node Template for the new nodes and reference them to the Node Templates with external references. The Pointers to Artifacts are contained by Artifact Templates which are used by Node Type Implementations. By composing all the information a simple references chain can be built:

Artifact \rightarrow *Artifact Template* \rightarrow *Node Type Implementation* \rightarrow *Node Type* \rightarrow *Node Template* \rightarrow *Service Template*

Now consider the references in more detail.

- *Artifact* \rightarrow *Artifact Template*
An Artifact can be referenced by several Artifact Templates. (Despite the fact that this is a bad practice.)
- *Artifact Template* \rightarrow *Node Type Implementation*
The same way an Artifact Template can be used by several Node Type Implementations.
- *Node Type Implementation* \rightarrow *Node Type*
A Node Type Implementation can describe an implementation of only one Node Type.
- *Node Type* \rightarrow *Node Template*
Each Node Type can have any number of Node Templates.
- *Node Template* \rightarrow *Service Template*
But each Node Template is instantiated only once.

This structure can be described as a tree with an Artifact as a root, and Service Templates as leaves (The example is on figure 4.2) and will be called the internal dependencies tree.

There is an additional problem in the reference between a Node Type and a Node Type Implementation. A Node Type can have several implementations, but which one will be used is determined only during the deployment. The chosen solution to this problem is to use each Node Type Implementation in the hope, that they will not conflict.

The following steps to build the internal dependencies tree can be executed during the preprocessing.

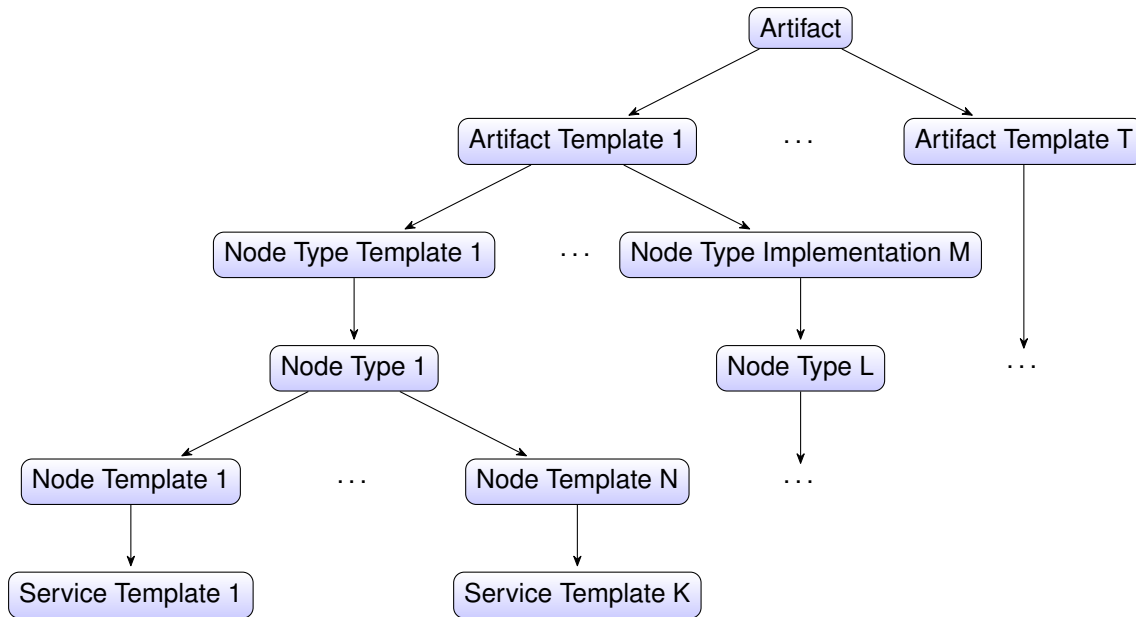


Figure 4.2: An example tree describing how to find Service Templates and Node Templates for a given script

- Find all Artifact Templates to build references from Artifacts to Artifact Templates.
- Find all Node Type Implementations. Since they contain references both to the Node Type and to the Artifact Templates, so the dependency from Artifact to Node Types can be built.
- Find all Service Templates and all contained Node Templates they contain. Each Node Template contains a reference to Node Type what is useful for building a dependency from Artifact to Node Template.

In this way the required internal dependencies tree can be built (with references *Artifact* \rightarrow *Node Template* and *Artifact* \rightarrow *Service Template*).

4.1.2 Modules and extensibility

Unfortunately it is impossible to identify all types of external references, even when only one language and one package manager are used (see an example in the listing 4.1). Since this work is aimed at creating the easily expanded and supplemented tool, initially only basic usage of package managers will be considered.

The framework should handle different languages, each of which can support various package managers. A language module should filter files not belonging to the language

Listing 4.1 Unreadable bash script

```
#!/bin/bash
set line = abcdefghijklmnoprst
set word1 = ${line:0:1}${line:14:1}${line:17:1}
set word2 = ${line:6:1}${line:4:1}${line:17:1}
$word1-$word2 install package
```

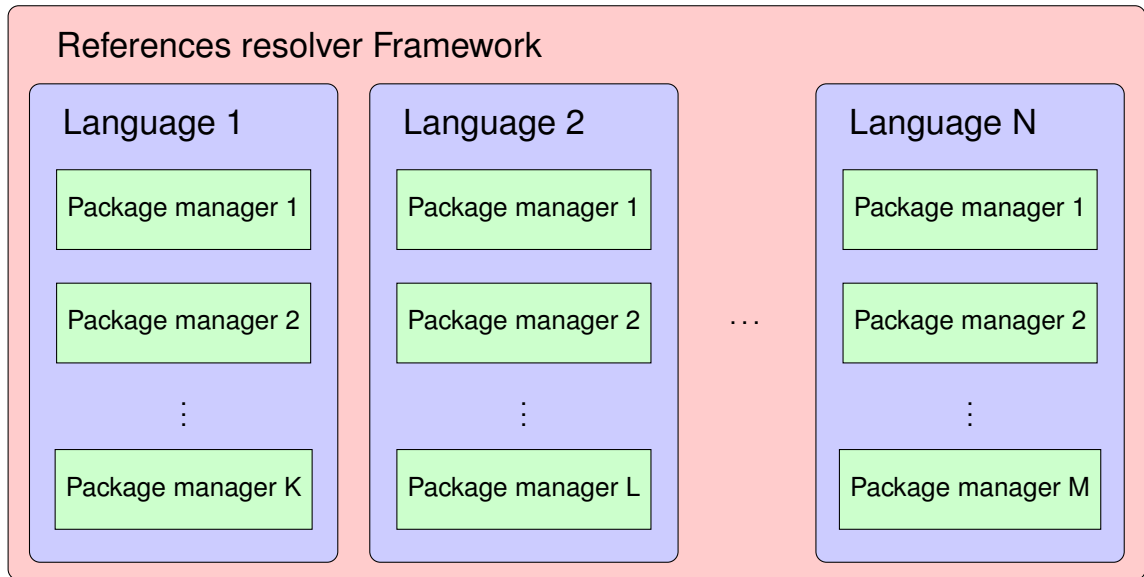


Figure 4.3: An example scheme representing several language modules containing package manager modules

and the accepted files will be transmitted to the corresponding package manager modules. This principle can be illustrated by figure 4.3.

A package manager module resolves an external reference and transmits the package name to a package handler described in section 4.2.5.

The framework will contain a list of all supported language modules and each language module will contain a list of supported package managers modules. Ease of adding new modules to the framework will prove the correctness of the architecture.

At the beginning the most popular combination will be implemented: the *bash* language with the *apt – get* package manager. This simple and powerful tool allows to install, delete or update the set of packages in one line of code. A line-by-line parser which analyses scripts and finds the installation commands should be developed. After the modules for this combination will be implemented, new language and package manager modules should be added.

4.1.3 Representing downloaded packages in TOSCA-Topology

A package node denotes the defined and instantiated element of TOSCA topology, the purpose of which is to install the package. The addition of new package nodes to TOSCA topology can be divided into several steps.

- Add definitions for common elements, like Artifact Types or Relationship Types. This can be done once at the preprocessing stage.
- The package node main definition will be represented by a Node Type. It must contain an *install* operation, which represents the capability to install the node.
- Artifacts (the downloaded data and the installation script) will be referenced by Artifact Templates.
- A Node Type Implementation will combine the artifacts to implement the *install* operation.
- A Node Template will instantiate the package node in the corresponding Service Templates. To determine the corresponding Service Template the autor will use the preprocessing described in the section Analysis of existing TOSCA-Topology.
- A Reference Template will provide topology information, allowing the observer (a user or a runtime environment) to determine which nodes the package must be installed for. References will be created from the Node Template which needs the package to the Node Template of the created package nodes.

After an execution of those steps, a definition of a package node will be finished and this node can be used.

4.1.4 Determining architecture of the final platform

Another problem happens while choosing the architecture of the device where packages will be installed. Unfortunately, it is impossible to analyze the structure of any CSAR and give an unambiguous answer to the question which architecture which node will be deployed on. There are many pitfalls here.

A single Service Template can use several physical devices with different architectures. Many Node Types and Node Templates instantiated on different platforms can refer to the same Implementation Artifact. This way one simple Implementation Artifact with a bash script containing "*apt-get install python*" command can be deployed on different devices within one Service Template (for example with the arm, amd64 and i386 architectures) and will result in the loading and installation of three different packages. For an end user, the ability to use such a simple command is a huge advantage,

but for the framework, it can greatly complicate the analysis. The following methods of architecture selection were designed.

- *Deployment environment analysis*
The script can analyze the system where it was started (for example using the "uname -a" command) and depending on the result, it will install the package corresponding to the system's architecture.
- *Unified architecture*
The architecture will be defined by the user for the whole CSAR.
- *Artifact specific architecture*
The architecture will be defined for each artifact separately.

Analysis of methods

The *deployment environment analysis*, which at first sight seems to be the most reliable solution, brings many additional problems. Packages for different platforms can differ not only by architecture but also by the version and the list of dependencies. As a consequence, chaos may occur while mirroring these different packages with different versions to the TOSCA-topology. The only found robust solution is to create a set of archives (one archive for one architecture), containing the entire dependency tree for the given package for each installed package. But this approach contradicts one of the main ideas of this work: the dependencies trees should be mapped to the topology.

The *artifact specific architecture* method carries an additional complexity to the user of the framework. It will make a user analyze each artifact and decide which architecture it will be executed on. All of this can be complicated by the fact that the same artifact can be executed on different architectures.

The method of the *unified architecture* was chosen as the simplest and easiest to implement. If it will be necessary, this method can easily be expanded to the *artifact specific architectures* method (by removing the user input at start, and choosing an architecture for each artifact separately) or to *deployment environment analysis* (by downloading packages for all available architectures and adding the architecture determining algorithm to the installation scripts).

4.1.5 Validation

Checking the output of the framework is an important stage in the development of the program. It is necessary to verify both the overall correctness of the output CSAR and the possibility to deploy generated package nodes. To the correctness it is possible to

use *Winery* tool from OpenTOSCA. This tool for creating and editing CSAR archives is also great for visualizing the results. Checking the deployment of the generated package nodes can be done manually by entering commands which start the artifact's execution.

4.2 Architecture

This section will present the architecture of the framework and the detailed description of its elements. The main elements are a *CSAR handler*, a *references resolver*, *language modules*, *package manager modules*, a *package handler*, and a *topology handler*.

4.2.1 CSAR handler

The CSAR handler provides an access to a CSAR and maintains its consistency. It describes the processes of adding the new files (to handle the metadata), archiving/unarchiving, and choosing the final platform architecture.

A new name/value pair must be added to the metadata for each new file integrated into the CSAR. The name represents the path and the name of the file. The value contains a type of the file. This type will be used by a runtime environment during the deployment to choose the right behavior. The input CSAR is initially archived and must be decompressed in order to handle the content first. When all external references will be resolved, the content will be archived to an output CSAR. As already was mentioned in section 4.1.4, an architecture of the final platform will be chosen for the entire CSAR. A command line interface must be provided for a user, to allow him to choose the architecture. The chosen architecture must be saved to the CSAR for the case of future processing by the framework to avoid the collisions between architectures of packages.

4.2.2 References resolver

This is the main element, whose execution can be divided into three stages: *preprocessing*, *processing*, *finishing*.

During the *preprocessing* stage, the CSAR will be unarchived, common files added, and internal dependencies trees generated. Figure 4.4 illustrates those steps. During the *processing* stage, all *language modules* will be activated, the operation is described in

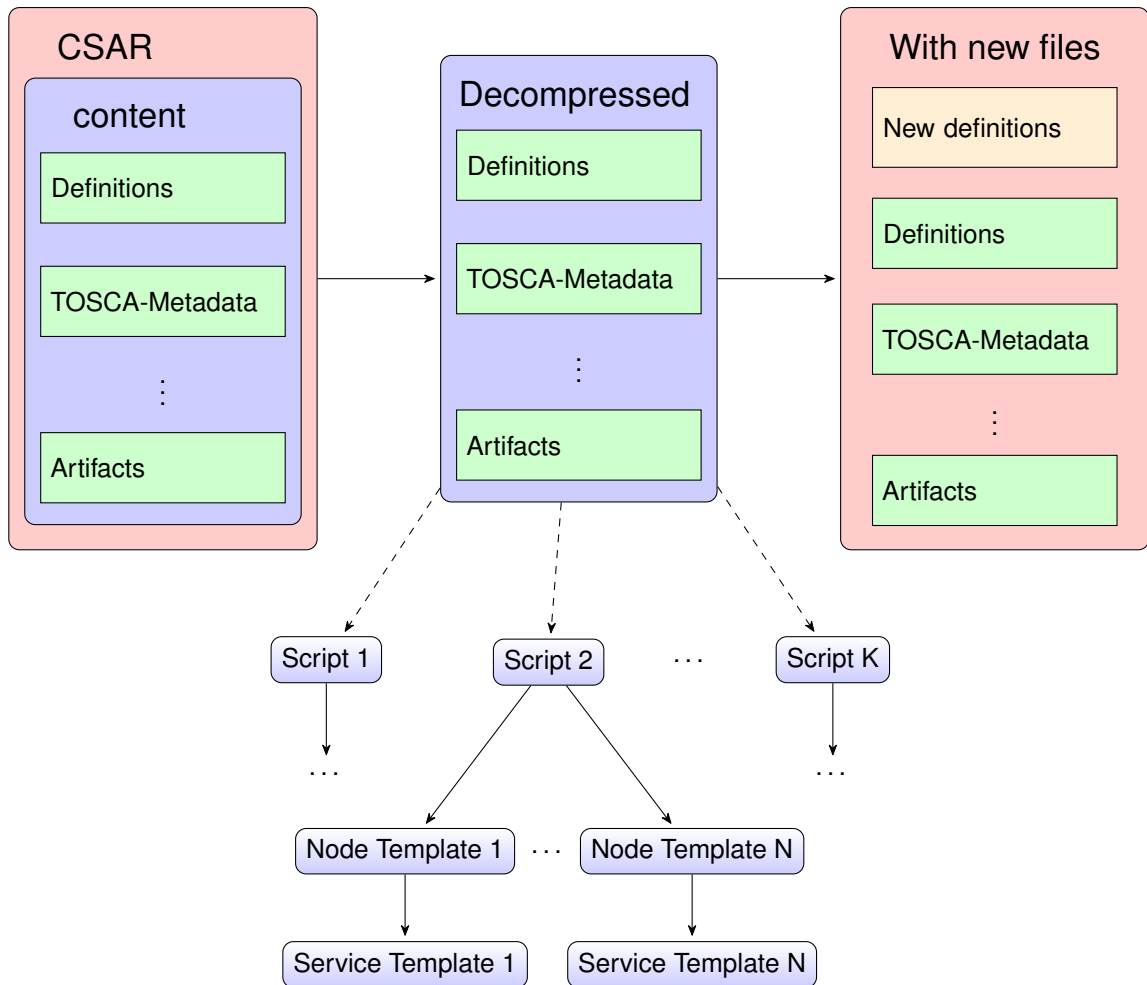


Figure 4.4: Preprocessing: decompression, adding files and generating dependencies

more detail in the next section. To finish the work all results will be packed into the output CSAR during the *finishing* stage.

4.2.3 Language modules

Each *language module* describes a handling of one language and chooses files written in the language. It also contains a list of supported package manager modules. Each language module must provide the capability to generate a TOSCA node for the given package and this node must use the same language to install the package. This means that a script and definitions for Artifact Templates, a Node Type, and a Node Type Implementation should be created by a language module.

As it is already mentioned above, during the *processing* stage a *language module* an-

alyzes all files one by one and checks their belonging to the language. Any files not belonging to the described language are filtered out. The remaining files are transferred to the *language module's package manager modules*. For example, a *bash* module will pass only files with ".sh" extension which start with the "#!/bin/bash" line. An *ansible* module should have an additional functionality to unpack zip archives where ansible playbooks can be stored. Since ansible playbooks don't contain a specific header or marker, the single sign of ansible files is the ".yaml" extension.

4.2.4 Package manager modules

A *package managers module* finds external references, resolves them and transmits the package name to the *package handler* described in the next section. Figure 4.5 illustrates data flow between language modules, package manager modules and the package handler.

To resolve an external reference a package manager module will parse the given file. In the case of the apt-get module for bash, the module will read a file line-by-line searching for the strings starting with "apt-get install". Such strings must be commented out and their ends should be divided into separate package names which will be transferred to the package handler.

4.2.5 Package Handler

The *package handler* gets a package name, downloads an installation data for an architecture specified by the CSAR handler, transfers the package name to the *topology handler* and repeats the actions for all depended packages recursively. To download an installation data the command "apt-get download **package**" can be used. The architecture can be specified by a ":*architecture*" suffix, for example, a "package:*arm*" mean the package of the *arm* architecture. The list of dependencies will be obtained using the "apt-cache depends **package**" command. The output of such command should be parsed in order to extract names of depended packages. Of course, in case of a fault during a download of a package, a user interface should be provided to find a solution. It can be: retry the download, ignore the package, rename the package or even break the framework's execution.

4.2.6 Topology Handler

This element should handle the TOSCA topology and it has two main tasks: to analyze the TOSCA topology during the preprocessing stage to create internal dependencies trees

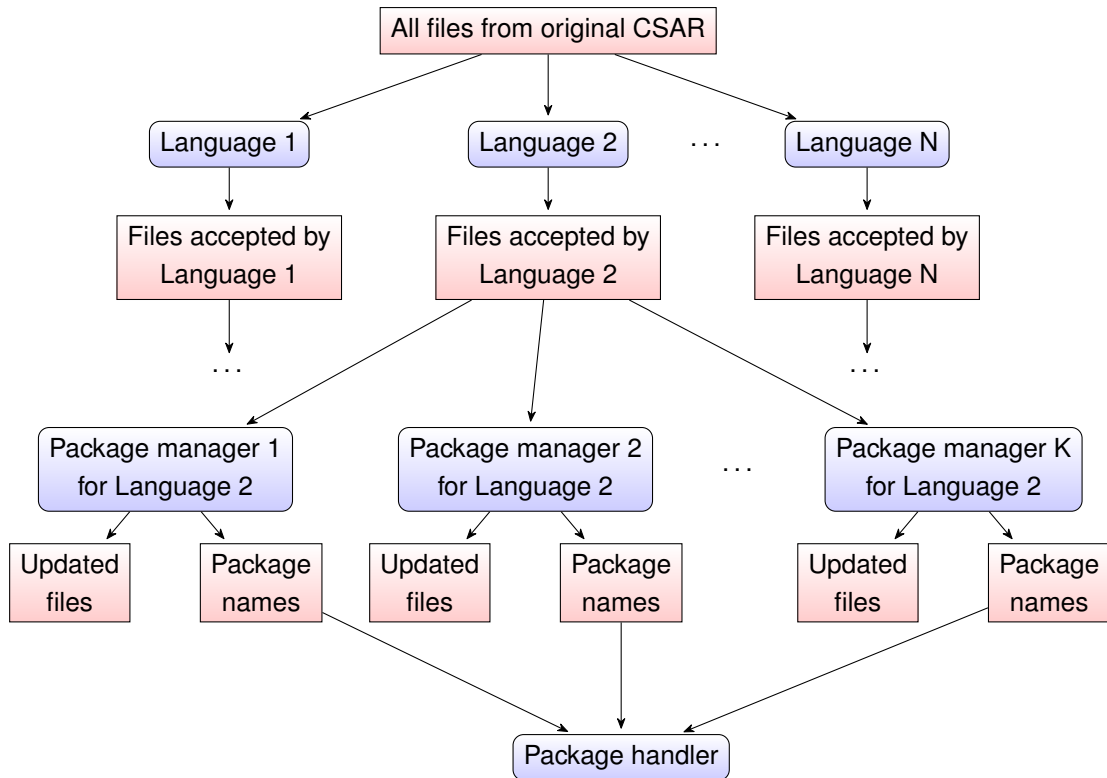


Figure 4.5: The data flow scheme between language modules, package manager modules and the package handler.

and to use those trees to create TOSCA definitions for Node Templates and Relationship Templates in the right places for the packages provided by the package handler. The analysis of the TOSCA topology was described in section 4.1.1 and the definition of Node Templates and Relationship Templates was given in section 4.1.3. The Internal dependencies trees must be updated to represent changes after addition of new Node Templates and Relationship Templates.

5 Implementation

This chapter provides an information about the implementation of the framework and its elements, whose behavior was described in chapter 4. Java language was chosen, because of its simplicity and strength. In this language, the elements are represented by classes. Java uses additional kind of packages which describe third-party modules and make programming easier. The used Java packages will be mentioned here and the necessary license will be listed in the "NOTICE.txt" file in the source code's root folder.

5.1 Global elements

This section describes the elements used throughout the whole framework's execution. A ZIP handler provides a functionality to operate ZIP archives, a CSAR handler keeps an interface to interact with a CSAR and a Utils helps to solve problems common to other elements.

Zip handler

This is a small element with straight functionality. It serves to pack and unpack ZIP archives which are used by the CSAR standard. It was decided to use the *java.util.zip* package for this task. The functions of archiving and unarchiving are called *zipIt* and *unZipIt* respectively.

CSAR handler

This element provides an interface to access the CSAR content and stores information about files associated with it. The most valuable data are the name of a temporary extraction folder, the list of files from the input CSAR, the meta-file entry, and the architecture of the target platform. All this data is encapsulated into the CSAR handler. The set of public functions allowing to operate with this element is available.

5 Implementation

- *unpack* and *pack* functions are used to extract the CSAR into the temporary folder and pack the folder to the output CSAR. These functions use the *ZIP handler*.
- *getFiles* returns the list with files presented by the input CSAR.
- *getFolder* returns the path to the folder where the CSAR was extracted.
- *getArchitecture* returns the chosen architecture of the target platform.
- *addFileToMeta* adds information about the new file to the meta-data.

Here is an example usage of the element. When the CSAR handler extracts the input CSAR to the temporary extraction folder during the *unpack*'s call, it saves the folder's name. Then other elements can use the *getFolder* function to get this name and access the data.

Utils

This class provides the *createFile*, *getPathLength*, and *correctName* methods used by many other elements. The main purpose of these functions is to make the code cleaner. Using the *createFile* an element can create a file with the given content. The *getPathLength* function returns the deep of the given file's path and it is very useful for creating references between files.

The OpenTOSCA uses some limitations to names of TOSCA nodes. Those names can't contain slashes, dots and so on. To obtain an acceptable name from a given name the function *correctName* can be used.

5.2 References resolver

This is the main module which starts by framework startup and is executed into three stages: preprocessing, processing and finishing.

Preprocessing

At the preprocessing stage, the CSAR is unpacked, common TOSCA definitions are generated and internal dependencies trees are built. As the first step, a user interface is provided to get the names of the input CSAR, output CSAR and the architecture of the final platform. To unpack the CSAR the function *unpack* from the CSAR handler is used. The *javax.xml.bind* package was chosen for creating the common TOSCA definition.

This Java package allows to generate a description - Java class describing an XML document. Those documents contain the following definitions.

- *DependsOn* and *PreDependsOn* describe Relationship Types between packages.
- *Package Artifact* defines a deployment Artifact Type for a package installation data.
- *Script Artifact* specify an implementation Artifact Type for a script installing a package.
- *Ansible Playbook* represent a deployment Artifact Type for a package installation via an ansible playbook.

An example description of the *Script Artifact* can be found in listing 8.1. Each description is represented by a separate Java class.

To build internal dependencies trees the topology handler described in section 5.6 is used.

Processing

During this stage, all language modules listed in the framework are started. For the references resolver element that is only two strings of code, but they start the main functionality of the framework. The languages modules check all files presented in the input CSAR. The list of these files is stored in the CSAR handler, a pointer to which the modules became during their instantiation and translate to the corresponding package manager modules. This system allows the modules to access the CSAR's content at any time.

Finishing

When all external references will be resolved, the framework can enter its last stage. At this stage, the changed data should be packed into an output CSAR, whose name was entered during the preprocessing stage. The function *pack* from the CSAR handler is used. After this operation, we become a more encapsulated CSAR. Its level of Internet access during a deployment will be significantly lower.

5.3 Language modules

This section will describe the language modules. Since the framework is initially oriented to easy extensibility, an abstract model for the modules will be defined, so that new modules can be added by implementing this model. The implementation of the bash and ansible modules will be provided at the end of the section.

Language model

To specify the common functionality and behavior of different language modules, the language model is used. In Java, this model is described by an abstract class. The abstract class *Language* is presented in listing 8.2. The common variables for all language modules are the name of the language, the list with package manager modules, and the extensions of files. And here are the common functions presented:

- *getName* returns the name of this language.
- *getExtensions* returns the list of extensions for this language.
- *proceed* checks all original files. Files written in the language should be transferred to every supported package manager module.
- *getNodeName* uses a package name to generate the name for a Node Type, which will install the package using this language.
- *createTOSCA_Node* creates the definitions for a TOSCA node. Since the created TOSCA nodes must install packages using the same language as the original node, all languages must provide the method for creating such definitions.

New language modules must be inherited from the language model and then can be added to the framework.

Bash module implementation

The processing of the popular language was implemented. The bash module should accept only files written on the bash language. To chose such files some signs inherent to all bash scripts can be used. These signs can be the file extensions (".sh" or ".bash") and the first line ("#!/bin/bash"). Each file which contains those signs will be passed to supported package managers modules, in our case to the *apt-get* module described later. The bash module must provide a capability for the given package to create a definition of a TOSCA node which uses the bash language to install the package. Such a bash

TOSCA node is defined by the Package Type, the Implementation, the Artifact, and the Script Artifact. The Package Type is a Node Type with an *"install"* operation and a name from the *getNodeName* function. The Package Implementation is a Node Type Implementation which refers the Package Artifact and the Script Artifact to implement the Package Type's *"install"* operation. The Package Artifact and the Script Artifact are Artifact Templates referencing the installation data and a bash installation script respectively. The installation script contains the bash header and an installation command, like *"dpkg -i installation_data"*. The topology handler will instantiate the package node later by defining a Node Template. Those definitions and the installation script are created by the *createTOSCA_Node* function.

Ansible implementation

To test the extensibility of the framework, the ansible language was added. Since ansible playbooks are often packed into archives, it may be necessary to unpack them first and then to analyze the content. Thus, the files are either immediately transferred to the package manager modules, or they are unzipped first. Listing 8.4 represents those operations. As a sign of the ansible language, the *".yaml"* extension is used, since its playbooks don't contain any specific header.

Creating an ansible TOSCA node for a package is a complicated operation. As the first step, the original files should be analyzed to determine the configuration (the set of options like a user name or a proxy server). If the implemented analyzer is unable to find all necessary options, a user interface will be provided to fulfill any missing parameters. Having the configuration a playbook and a configuration file will be created in a temporary folder. After the installation data has been added to the folder, it can be packed to a zip archive. This archive is an implementation artifact, which the Artifact Template should be created for. A Node Type with an *"install"* operation should be defined. And finally, a Node Type Implementation linking the operation and the Artifact Templates should be defined. A Node Template will be added later by the topology handler.

5.4 Package manager modules

In this section, package manager modules will be specified. Like languages, an abstract model will be defined to make the extensibility easier. The apt-get module for bash and an apt module for ansible will be implemented.

Package manager model

The model is described by an abstract class. Its description contains only one function *proceed* (in listing 8.3), that finds and eliminates external references, as well as passes the found package names to the package handler.

Apt-get for bash

The apt-get package manager module is a simple line-by-line file parser which searches for the lines starting with the "*apt-get install*" string, comments them out and passes this command's arguments to the package handler's public function *getPackage*.

Apt for ansible

Since the package installation written in the *ansible* language with the *apt* package manager can be described in many different ways, then the processing will be a complicated task too. It's worth mentioning that the processing uses a simple state machine and regular expression from the *java.util.regex* package.

5.5 Package Handler

Package handler provides an interface for interaction with the package manager of the operating system. It allows to download packages and to determine the type of dependencies between them.

Package downloading

The download operation is performed using one recursive function *getPacket*. The Arguments of the function will be defined shortly.

- *language* is a pointer to the language module which has accepted the original artifact.
- *packet* is a name of the package.
- *listed* holds a list with already downloaded packages. It is not necessary to download them again, but new dependencies will be created.

- *source* defines the parent element of the package. It will be the original artifact file for the root package, and the depending package for other packages.
- *sourcefile* is a name of the original artifact.

This function downloads packages, calls the language's function *createTOSCA_Node* to create the TOSCA node for the package and the topology handler's functions *addDependencyToPacket* or *addDependencyToArtifact* to update the topology. Then this function calls itself recursively for all depended packages. After those operations, a dependencies tree for the *packet* will be built.

The command *apt-get download package* is used for downloading. If the process fails, a user input is provided to solve the problem. The user will be able to rename the package, ignore it or even break the processing.

Dependencies

To determine the dependency type the *getDependencies* function was developed. It becomes a *package* as an argument and uses the command *apt-cache depends package* to build a list with dependencies for the *package*. The *apt-cache* command is a part of the *apt-get* package manager and uses a packages database to print the dependencies. The output is parsed to find strings like "Depends: *dependent_package*". These dependent packages are combined to a list and returned back.

5.6 Topology handling

The topology handler serves to update the TOSCA topology. It builds the internal dependencies trees during the preprocessing stage. Later the trees are used to find the right places for definitions of Node Templates and Dependency Templates.

Building internal dependencies trees

At the preprocessing stage, this element analyzes all original definitions and constructs internal dependencies trees. To read those definitions from the XML files the package *org.w3c.dom* was taken.

As the first step, all definitions of Artifact Templates are analyzed and pairs consist of an Artifact Template's ID and an artifact itself are built. Then each Node Type Implementation will be read and Node Types and Artifact Template's IDs found. Now each artifact has a set with Node Types where it is used. After the analysis of Service

Listing 5.1 Creating of a new Node Template

```
Element template = document.createElement("tosca_ns:NodeTemplate");
template.setAttribute("xmlns:RRnt",
    RR_NodeType.Definitions.NodeType.targetNamespace);
template.setAttribute("id", getID(package));
template.setAttribute("name", package);
template.setAttribute("type", "RRnt:" + RR_NodeType.getTypeName(package));
topology.appendChild(template);
```

Templates, analog sets of Node Templates for each artifact will be created. In addition, for each Node Template one should keep a Service Templates, where this Node Template was defined.

Updating Service Templates

To update Service Templates two functions are provided.

- *addDependencyToPacket(sourcePacket, targetPacket, dependencyType)* generates a dependency between two package nodes.
- *addDependencyToArtifact(sourceArtifact, targetPacket)* generates a dependency between the original node and a package node.

Both functions find all Node Templates which instantiate the given *sourcePacket* or *sourceArtifact*. Besides, they find Service Templates where the Node Templates are defined. The search is done with the help of the internal dependencies trees. For each found Node Template a package node for the *targetPacket* package should be instantiated by creating a new Node Template. Then the dependency between the found Node Template and the new Node Template is created by defining a Relationship Template. The Relationship Template references both Node Templates. The type of dependency is the value of the *dependencyType* for the *addDependencyToPacket* function and the *preDependsOn* for the *addDependencyToArtifact*.

To update the existing TOSCA definition the *org.w3c.dom* and *org.xml.sax* packages are used. The definition of a new Node Template for the given *topology* and *package* is presented in listing 5.1.

6 Add new package manager module

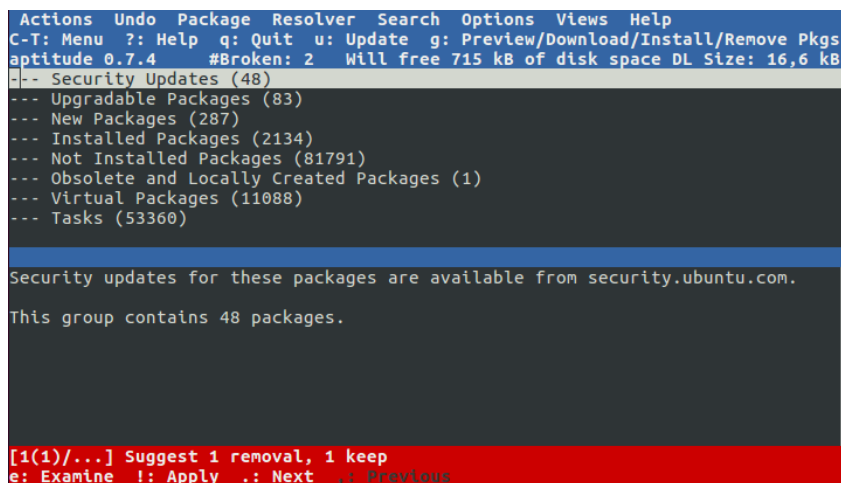
This chapter shows the extensibility of the framework by adding a new *aptitude* package manager module for the *bash* language.

Section 6.1 provides a common information about the *aptitude*.

In section 6.2 the module is implemented and in section 6.3 is integrated into the *bash* language module.

6.1 Aptitude

This section describes the *aptitude* package manager. Like to the *apt-get*, the *aptitude* is a command line program, where a package can be installed using the *aptitude install package* command. In additional, it can be started in a pseudo-graphical mode to provide a visual interface (An example in figure 6.1). An another advantage compared to the *apt-get* is the capability to search for packages by a part of the name (or by any other attributes) using the *aptitude search text* command.



```
Actions Undo Package Resolver Search Options Views Help
C-T: Menu ? : Help q : Quit u : Update g : Preview/Download/Install/Remove Pkgs
aptitude 0.7.4 #Broken: 2 Will free 715 kB of disk space DL Size: 16,6 kB
--- Security Updates (48)
--- Upgradable Packages (83)
--- New Packages (287)
--- Installed Packages (2134)
--- Not Installed Packages (81791)
--- Obsolete and Locally Created Packages (1)
--- Virtual Packages (11088)
--- Tasks (53360)

Security updates for these packages are available from security.ubuntu.com.
This group contains 48 packages.

[1(1)/...] Suggest 1 removal, 1 keep
e: Examine !: Apply .: Next <: Previous
```

Figure 6.1: A command line visual interface for the *aptitude* package manager.

6 Add new package manager module

Listing 6.1 The *aptitude* inherited from the *PackageManager* abstract class

```
public final class PM_aptitude extends PackageManager {
    @Override
    public void proceed(String filename, String source)
        throws FileNotFoundException, IOException, JAXBException {
        // TODO Auto-generated method stub
    }
}
```

Listing 6.2 The *aptitude* module with some common elements

```
public final class PM_aptitude extends PackageManager {

    // name of the package manager
    static public final String Name = "aptitude";

    /**
     * Constructor
     */
    public PM_aptitude(Language language, CSAR_handler ch) {
        this.language = language;
        this.ch = ch;
    }

    @Override
    public void proceed(String filename, String source)
        throws FileNotFoundException, IOException, JAXBException {
        // TODO Auto-generated method stub
    }
}
```

6.2 Implementing the new package manager module

The implementation of the *aptitude* module will be described here. At first, the *aptitude* class will be inherited from the abstract class *PackageManager*. This is presented in the listing 6.1.

After that, the *aptitude* class can be used as a regular package manager module (but it lacking functionality). It is necessary to add the common code, like the constructor and the manager's name. After these operations, the *ansible* module can be presented in the listing 6.2.

Since the package manager will read files from an input CSAR, the CSAR handler is stored by the constructor to the *ch* variable for a further use. In addition, the language

Listing 6.3 The aptitude *proceed* function

```
@Override
public void proceed(String filename, String source)
    throws FileNotFoundException, IOException, JAXBException {
    if (ch == null)
        throw new NullPointerException();
    System.out.println(Name + " proceed " + filename);
    BufferedReader br = new BufferedReader(new FileReader(filename));
    boolean isChanged = false;
    String line = null;
    String newFile = "";
    while ((line = br.readLine()) != null) {
        // TODO parsing will be done here
    }
    br.close();

    if (isChanged)
        Utils.createFile(filename, newFile);
}
```

(*bash* in this case) is stored too, to be propagated later to the package handler.

Now focus on the *proceed* function. A line-by-line file analyzer is needed, which can modify the data and in the case of a modification, the entire file should be rewritten.

The *isChanged* variable indicates that the file must be rewritten with a new content from the *newFile* variable. Now an *aptitude* line parser will be implemented, which reads a line from the *line* variable and stores it or its changed version to the *newFile* variable. If the data is changed, then the *isChanged* variable must be set to true. Any *ansible* package installation calls should be detected, commented out and its arguments (package names) should be propagated one by one to the package handler's function *getPackage*.

During the parsing which is described in the listing 6.4, the line is divided into words. Each found package name is transmitted to the packet handler as an argument of its public function *getPackage*. In additional, this function must take the language and the source artifact's name as the arguments.

6.3 Integrating Aptitude into the Bash module

Now the *aptitude* module can be added to the bash module. The only thing to do is to add the *aptitude* to the bash's list of package manager modules (the list is stored in the *packetManagers* variable). This is done by the bash's constructor with the command:

6 Add new package manager module

Listing 6.4 The aptitude line parser

```
String[] words = line.replaceAll("[:&]", "").split("\\s+");
// skip spaces at the beginning of string
int i = 0;
if (words[i].equals(""))
    i = 1;
// looking for aptitude
if (words.length >= 1 + i && words[i].equals("aptitude")) {
    // aptitude found
    if (words.length >= 3 + i && words[1 + i].equals("install")) {
        System.out.println("aptitude found:" + line);
        isChanged = true;
        for (int packet = 2 + i; packet < words.length; packet++) {
            System.out.println("package: " + words[packet]);
            ch.getPackage(language, words[packet], source);
        }
    }
    newFile += "##References resolver//" + line + '\n';
}
else
    newFile += line + '\n';
```

"packetManagers.add(new PM_aptitude(this, ch));".

The new package manager module is ready to work.

7 Validation

In this chapter, the developed framework will be tested. An input CSAR will be described in section 7.1. The processing by the framework is described in section 7.2. The output CSAR will be added to and displayed by Winery in section 7.3. Generated Artifacts will be checked in section 7.4.

7.1 Input CSAR

In this test, an CSAR from the OpenTOSCA Demos is used. The CSAR provides a service for Automating the Provisioning of Analytics Tools based on Apache Flink. [com16b] The structure of the service is provided in figure 7.2. The service uses a server virtualization environment named *vSphere* (The *VSphere_5.5* node from the structure.). In the environment works the *Ubuntu* virtual server (The *Ubuntu-14.04-VM* node). The *Ubuntu* hosts two applications: the *Python* (*Python_2.7*) and the *Flink Simple* (*Flink_Simple_1.0.3*). An analyze shows two external references. The *Python* node installs the python package and the *Flink Simple* node - the Java package. The service has two submodules: a Data Prediction and a Data Delivery, both a hosted on the *Flink Simple* node and require the Python node.

7.2 Processing

Since the framework is written in the Java, to start it a JDK (version 1.8 or above) is necessary. Additionally, the apt-get package manager must be installed. To start the framework an Java environment is used. After the start, a user should enter the input CSAR name, the output CSAR name, and the architecture. After that, the framework should work fully automatically, analyzing the artifacts and resolving any external references. Figure 7.1 provides an example.

```

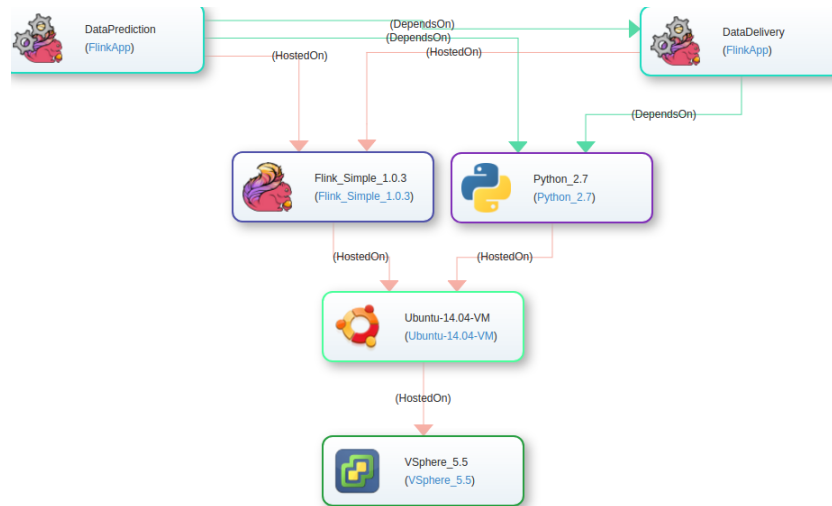
jerry@jerry-note:~/TOSCA$ java -jar RR.jar
enter the input CSAR name: FlinkApp_Demo_Small_On_VSphere.csar
enter the output CSAR name: output.csar
source: FlinkApp_Demo_Small_On_VSphere.csar
target: output.csar
Proceeding file FlinkApp_Demo_Small_On_VSphere.csar
Please enter the architecture.
Example: i386, amd64, arm, noarch.
architecture: amd64
Parse Artifacts
Parse Implementations
Parse ServiceTemplates
RefToNodeType
artifacttemplates_httpP3AP2FP2Fopentosca_orgP2Fartifacttemplates_VSphere_5_5_Clo
udProviderInterface_IA_files_org_opentosca_nodetypes_VMWare5_5__CloudProviderInt
erface_war : [VSphere_5.5]
artifacttemplates_httpP3AP2FP2Fopentosca_orgP2Fartifacttemplates_FlinkApp_IA_fil
es_start_sh : [FlinkApp]
artifacttemplates_httpP3AP2FP2Fopentosca_orgP2Fartifacttemplates_Python_2_7_Impl
_InstallIA_files_install_sh : [Python_2.7]
artifacttemplates_httpP3AP2FP2Fopentosca_orgP2Fnodetypes_Ubuntu-14_04-VM_Operati
ngSystemInterface_IA_files_org_opentosca_NodeTypes_Ubuntu-14_04-VM__OperatingSys
temInterface_war : [Ubuntu-14.04-VM]
artifacttemplates_httpP3AP2FP2Fopentosca_orgP2Fartifacttemplates_Flink_Simple_1

```

Figure 7.1: Processing by the framework.

7.3 Displaying with Winery

Winery was installed to test the correctness of the output CSAR. This is an environment for the development of TOSCA systems and is useful for checking the results. The input CSAR's representation by Winery is displayed in figure 7.2. Those external references will be resolved by the framework and exchanged by new nodes in output CSAR.

Figure 7.2: Source CSAR represented by *Winery*.

Add to Winery

The output CSAR is added to Winery. Due to a significant increase in size, this can be a fairly lengthy procedure. It was only six nodes in the input CSAR, but after the processing, the output CSAR contains more than 100 of nodes. During the addition to Winery, the CSAR's syntax is tested. In a case of errors, messages will be displayed.

Display by Winery

The output CSAR will be displayed. Again, due to the high number of nodes, the processing can take a long time. At the time, the correctness of the internal references will be checked. If something was defined not properly, these erroneous nodes or links between them will not be displayed. The representation of the output CSAR by Winery is shown on figure 7.3 (Only the part of the CSAR is visible).

It seems pretty beloved. To verify the TOSCA's structure some nodes was moved manually (figure 7.4). The correctness of dependencies was verified by checking several nodes with the *apt-cache depends* command. By opening the content of the new nodes, it was verified, that there are right artifacts.

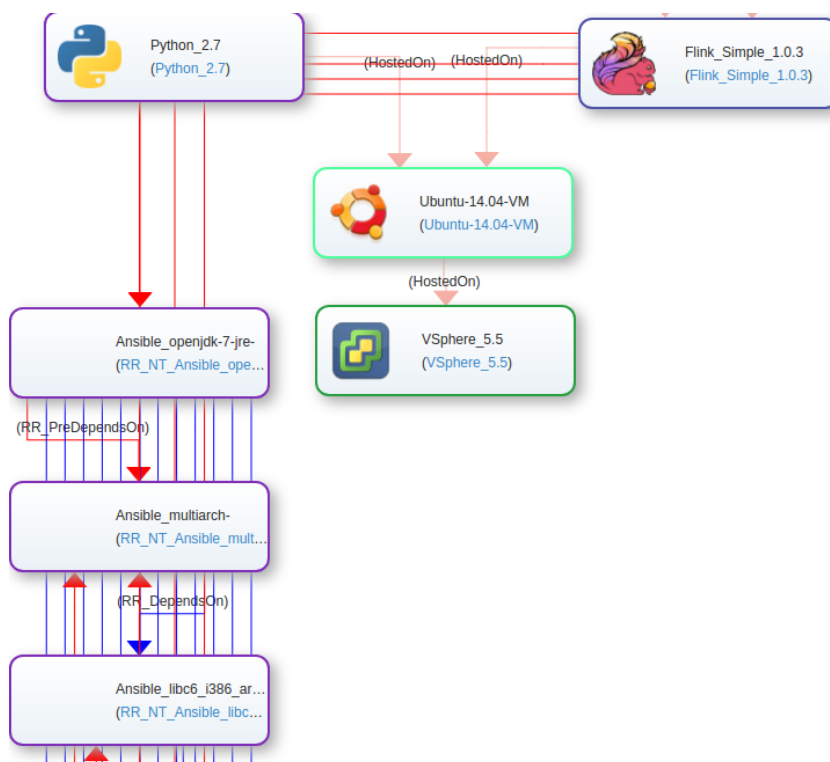


Figure 7.3: The output CSAR represented by *Winery*.

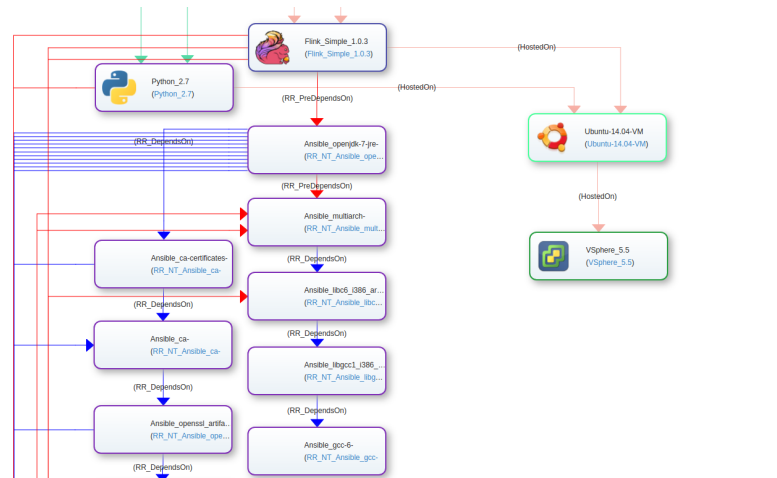


Figure 7.4: The output CSAR represented by *Winery*, some nodes moved manually.

Listing 7.1 Check bash installation script

```
user@user:~$ sudo RR_python2_7-minimal.sh
(Reading database ... 286091 files and directories currently installed.)
Preparing to unpack python2_7-minimal.deb ...
Unpacking python2.7-minimal (2.7.12-1ubuntu0~16.04.1) over
(2.7.12-1ubuntu0~16.04.1) ...
Setting up python2.7-minimal (2.7.12-1ubuntu0~16.04.1) ...
Processing triggers for man-db (2.7.5-1) ...
```

7.4 Check artifacts

Also, it is necessary to check whether it is possible to install new packages using the generated artifacts. At first bash scripts will be tested, then ansible playbooks.

Check bash scripts

Since the bash is used in the Linux's command line, it will be pretty easy to check bash installation scripts by starting them (of course that must be done with the necessary privileges). An example of the *python2.7* installation is presented in the listing 7.1. The process ended without any warnings or errors, which means that it was completed successfully. This way any bash installation script can be checked.

Check ansible playbooks

To check an ansible playbook manually we need to extract the zip file containing the playbook. During the regular execution, this work will be done by the runtime environment. The call of the ansible runtime which proceeds the playbook is a simple procedure too. An example is provided in figure 7.5. *Ok* signals that the installation was completed successfully.

```
jery@jery-note:~/TOSCA/temp$ sudo ansible-playbook ./main.yml
[WARNING]: provided hosts list is empty, only localhost is available

PLAY [install package] *****
TASK [setup] *****
ok: [localhost]
TASK [install task] *****
changed: [localhost]
PLAY RECAP *****
localhost                : ok=2    changed=1    unreachable=0    failed=0
```

Figure 7.5: An ansible playbook's execution process

8 Summary

!!!Entwurf!!!

External references can negatively affect the performance of cloud applications. Unfortunately, very many TOSCA applications access external sources for installing various packages at the stage of deployment. If we have high level of an information security or an limited access to the Internet, these external dependencies lead to a lot of problems which encapsulation can solve. The purpose of this work was to develop a software for the elimination of such dependencies.

The software was developed in the Java language, which ensures its portability, ease of maintenance and expanding. To enable the ability to add processing of new types of external dependencies, the software was implemented in the form of a modular framework, with separate modules for languages and package managers. This allows the user to easily add their new handlers to package managers and languages. In the first version, the framework handled only *Bash* scripts which use the *apt-get* package manager. To check the simplicity of the extensibility, the processing of *Bash* scripts with the *aptitude* package manager and the Ansible playbooks with the *apt* package manager. The framework handles the CSAR as follows. The structure of a CSAR is analyzed to determine the internal dependencies between artifacts and Node Templates. Each language module selects artifacts for which the language was used. All such artifacts are transferred for processing to the package managers of the language. They find external dependencies, remove them and pass the names of the required packages to the package handler. It loads each package along with all its dependencies and sends information about them to the topology handler, which creates TOSCA nodes for these packages and dependencies from the original nodes to the new ones. In the future, a runtime environment could analyze these dependencies and install the necessary packages.

In order to show the extensibility of the framework, the process of adding the *aptitude* package manager module into *Bash* module was described in detail. It was shown how to create a module which can be added to the framework, how to implement its basic functions, pass data to the packet handler, and connect the module to *Bash* language. In the end, the results of the framework were validated. The output SCARs were visualized and analyzed with the help of Winery. Generated artifacts were tested and executed.

As a result, the framework was obtained that eliminates external dependencies in a CSAR. It handles *Bash* language with *aptitude* and *apt-get* package managers and

8 Summary

Ansible language with *apt* package managers and can be easily expanded to handle additional types of external references.

Listings

Listing 8.1 Description for the script Artifact Type definition

```
public class RR_ScriptArtifactType {

    @XmlElement(name = "tosca:Definitions")
    @XmlAccessorType(XmlAccessType.PUBLIC_MEMBER)
    public static class Definitions {

        @XmlElement(name = "tosca:ArtifactType", required = true)
        public ArtifactType artifactType;

        @XmlAttribute(name = "xmlns:tosca", required = true)
        public static final String tosca="http://docs.oasis-open.org/tosca/ns/2011/12";
        @XmlAttribute(name = "xmlns:winery", required = true)
        public static final String
            winery="http://www.opentosca.org/winery/extensions/tosca/2013/02/12";
        @XmlAttribute(name = "xmlns:ns0", required = true)
        public static final String ns0="http://www.eclipse.org/winery/model/selfservice";
        @XmlAttribute(name = "id", required = true)
        public static final String id="winery-defs-for_tbt-RR_ScriptArtifact";
        @XmlAttribute(name = "targetNamespace", required = true)
        public static final String
            targetNamespace="http://docs.oasis-open.org/tosca/ns/2011/12/ToscaBaseTypes";

        public Definitions() {
            artifactType = new ArtifactType();
        }

        public static class ArtifactType {
            @XmlAttribute(name = "name", required = true)
            public static final String name = "RR_ScriptArtifact";
            @XmlAttribute(name = "targetNamespace", required = true)
            public static final String
                targetNamespace="http://docs.oasis-open.org/tosca/ns/2011/12/ToscaBaseTypes";
            ArtifactType() {}
        }
    }
}
```

Listing 8.2 Abstract language model

```
public abstract class Language {

    // List of package managers supported by language
    protected List<PackageManager> packetManagers;

    // Extensions for this language
    protected List<String> extensions;

    // Language Name
    protected String Name;

    // To access package topology
    protected Control_references cr;

    // List with already created packages
    protected List <String> created_packages;

    /**   Generate node name for specific packages
     * @param packet
     * @param source
     * @return
     */
    public abstract String getNodeName(String packet, String source);

    /**   Generate Node for TOSCA Topology
     * @param packet
     * @param source
     * @return
     * @throws IOException
     * @throws JAXBException
     */
    public abstract String createTOSCA_Node(String packet, String source) throws
        IOException, JAXBException;
}
```

Listing 8.3 Abstract package manager model

```
public abstract class PackageManager {

    // Name of manager
    static public String Name;

    protected Language language;

    protected Control_references cr;

    /**
     * Proceed given file with different source (like archive)
     *
     * @param filename
     * @param cr
     * @param source
     * @throws FileNotFoundException
     * @throws IOException
     * @throws JAXBException
     */
    public abstract void proceed(String filename, String source) throws
        FileNotFoundException, IOException,
        JAXBException;
}
```

Listing 8.4 Ansible proceeding

```
public void proceed() throws FileNotFoundException,
    IOException, JAXBException {
    if (ch == null)
        throw new NullPointerException();
    for (String f : cr.GetFiles())
        for (String suf : extensions)
            if (f.toLowerCase().endsWith(suf.toLowerCase())) {
                if (suf.equals(".zip")) {
                    proceedZIP(f);
                } else
                    proceed(f, f);
            }
}

public void proceed(String filename, String source)
    throws FileNotFoundException, IOException, JAXBException {
    for (PacketManager pm : packetManagers)
        pm.proceed(filename, source);
}

private void proceedZIP(String zipfile) throws FileNotFoundException,
    IOException, JAXBException {
    boolean isChanged = false;
    String folder = new File(cr.getFolder() + zipfile).getParent()
        + File.separator + "temp_RR_ansible_folder" + File.separator;
    List<String> files = zip.unZipIt(cr.getFolder() + zipfile, folder);
    for (String file : files)
        if (file.toLowerCase().endsWith(".yaml"))
            proceed(folder + file, zipfile);
    if (isChanged) {
        new File(cr.getFolder() + zipfile).delete();
        zip.zipIt(cr.getFolder() + zipfile, folder);
    }
    zip.delete(new File(folder));
}
```

Bibliography

- [13] OASIS Standard. *Topology and Orchestration Specification for Cloud Applications Version 1.0*. Nov. 25, 2013. URL: <http://docs.oasis-open.org/tosca/TOSCA/v1.0/TOSCA-v1.0.html> (cit. on p. 13).
- [BBH+13] T. Binz, U. Breitenbücher, F. Haupt, O. Kopp, F. Leymann, A. Nowak, S. Wagner. “OpenTOSCA - A Runtime for TOSCA-based Cloud Applications.” English. In: *Proceedings of 11th International Conference on Service-Oriented Computing (ICSOC’13)*. Vol. 8274. LNCS. Springer Berlin Heidelberg, Dec. 2013, pp. 692–695. DOI: [10.1007/978-3-642-45005-1_62](https://doi.org/10.1007/978-3-642-45005-1_62). URL: http://www2.informatik.uni-stuttgart.de/cgi-bin/NCSTRL/NCSTRL_view.pl?id=INPROC-2013-45&engl=1 (cit. on p. 21).
- [Bun14] S. Bundesamt. *12 % der Unternehmen setzen auf Cloud Computing*. Dec. 19, 2014. URL: https://www.destatis.de/DE/PresseService/Presse/Pressemitteilungen/2014/12/PD14_467_52911.html (cit. on p. 13).
- [Bun17] S. Bundesamt. *17 % der Unternehmen nutzten 2016 Cloud Computing*. Mar. 20, 2017. URL: https://www.destatis.de/DE/PresseService/Presse/Pressemitteilungen/2017/03/PD17_102_52911.html (cit. on p. 13).
- [Bur13] M. Burgess. *Cfengine V2.0 : A network configuration tool*. 2013. URL: http://www.iu.hio.no/~mark/papers/cfengine_history.pdf (cit. on p. 26).
- [CFE] CFEngine. URL: <https://cfengine.com/learn/learnhow-to-write-cfengine-policy/> (cit. on p. 26).
- [com16a] A. community. *Ansible (software)*. 2016. URL: http://docs.ansible.com/ansible/playbooks_best_practices.html#task-and-handler-organization-for-a-role (cit. on p. 25).
- [com16b] O. community. *OpenTOSCA for the 4th Industrial Revolution*. University of Stuttgart. 2016. URL: <http://www.opentosca.org/demos/smart-prediction-demo/index.htm> (cit. on p. 55).
- [Ham11] N. Hamilton. *The A-Z of Programming Languages: BASH/Bourne-Again Shell*. Computerworld: 2. 2011. URL: https://www.computerworld.com.au/article/222764/a-z_programming_languages_bash_bourne-again_shell/?pp=2&fp=16&fpid=1 (cit. on p. 25).

- [Laz16] M. Lazar. *Current Cloud Computing Statistics Send Strong Signal of What's Ahead*. Nov. 3, 2016. URL: https://www.insight.com/en_US/learn/content/2016/11032016-current-cloud-computing-statistics.html (cit. on p. 18).
- [Met15] C. Metz. *The Chef, the Puppet, and the Sexy IT Admin*. Wired. 2015. URL: https://www.wired.com/2011/10/chef_and_puppet/ (cit. on p. 26).
- [OAS] OASIS. *Organization for the Advancement of Structured Information Standards*. URL: <https://www.oasis-open.org/> (cit. on p. 19).
- [OAS13] OASIS. "Topology and Orchestration Specification for Cloud Applications Version 1.0." In: *OASIS Committee Specification 01*. <http://docs.oasis-open.org/tosca/TOSCA/v1.0/cs01/TOSCA-v1.0-cs01.html>. Mar. 18, 2013 (cit. on p. 20).
- [Pet11] T. G. Peter Mell. *The NIST Definition of Cloud Computing*. Recommendations of the National Institute of Standards and Technology, Special Publication 800-145. National Institute of Standards and Technology, Sept. 2011 (cit. on p. 17).
- [Sta16] Statista. *Umsatz mit Cloud Computing** weltweit von 2009 bis 2016 und Prognose bis 2020 (in Milliarden US-Dollar)*. 2016. URL: <https://de.statista.com/statistik/daten/studie/195760/umfrage/umsatz-mit-cloud-computing-weltweit-seit-2009/> (cit. on p. 13).
- [Stu13] I. U. Stuttgart. *OpenTOSCA - Open Source TOSCA Ecosystem*. 2013. URL: <http://www.iaas.uni-stuttgart.de/OpenTOSCA/> (cit. on pp. 13, 21).
- [tec] techopedia. *Definition - What does Cloud App mean?* URL: <https://www.techopedia.com/definition/26517/cloud-app> (cit. on p. 17).
- [Win] Winery. *Eclipse Winery*. URL: <https://projects.eclipse.org/projects/soa.winery> (cit. on p. 22).

All links were last followed on July 30, 2017.

Declaration

I hereby declare that the work presented in this thesis is entirely my own and that I did not use any other sources and references than the listed ones. I have marked all direct or indirect statements from other sources contained therein as quotations. Neither this work nor significant parts of it were part of another examination procedure. I have not published this work in whole or in part before. The electronic copy is consistent with all submitted copies.

place, date, signature